

深度学习

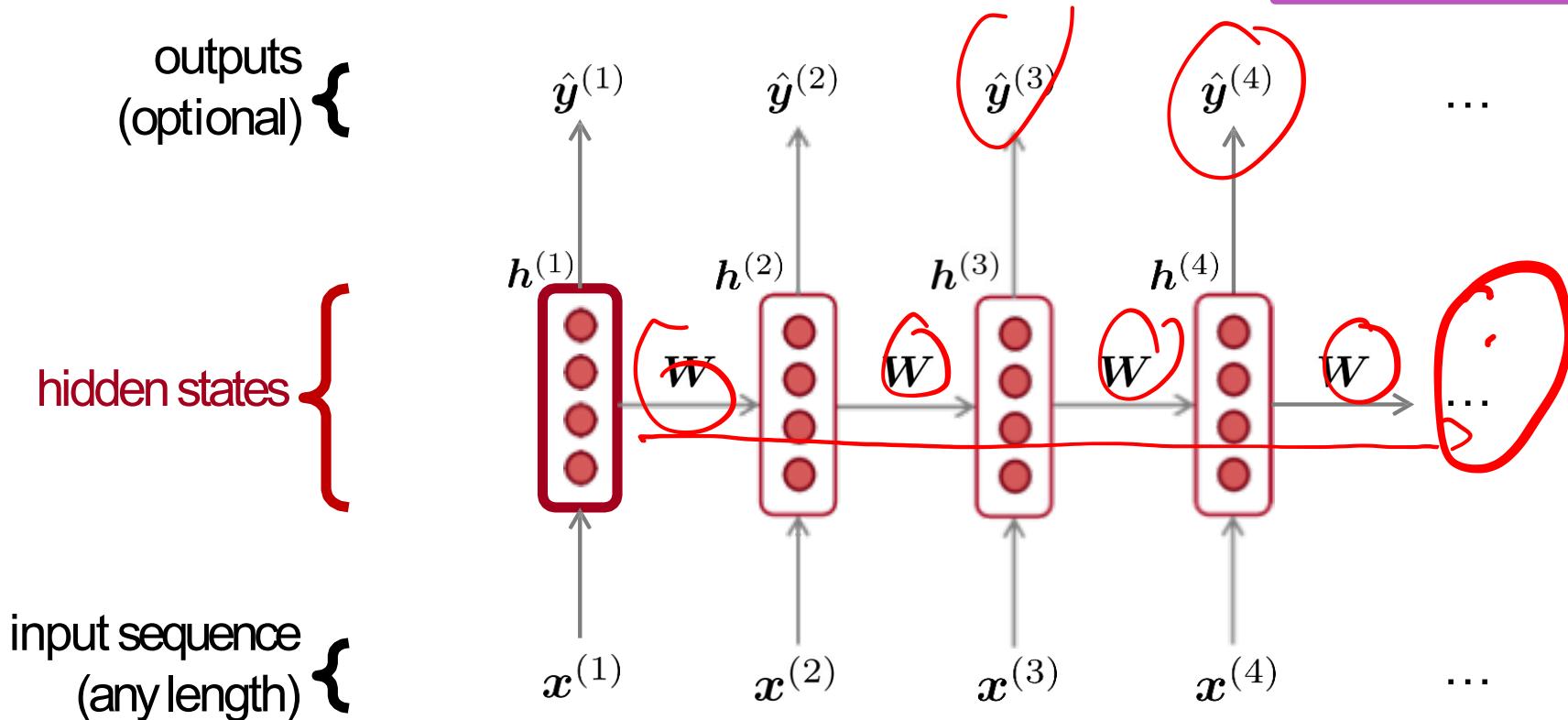
第十讲

王胤

Recurrent Neural Networks (RNN)

A family of neural architectures

Core idea: Apply the same weights W repeatedly



A RNN Language Model

$\hat{y}^{(4)} = P(\mathbf{x}^{(5)} | \text{the students opened their})$

output distribution

$$\hat{y}^{(t)} = \text{softmax} (\mathbf{U} \mathbf{h}^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

hidden states

$$\mathbf{h}^{(t)} = \sigma (\mathbf{W}_h \mathbf{h}^{(t-1)} + \mathbf{W}_e \mathbf{e}^{(t)} + \mathbf{b}_1)$$

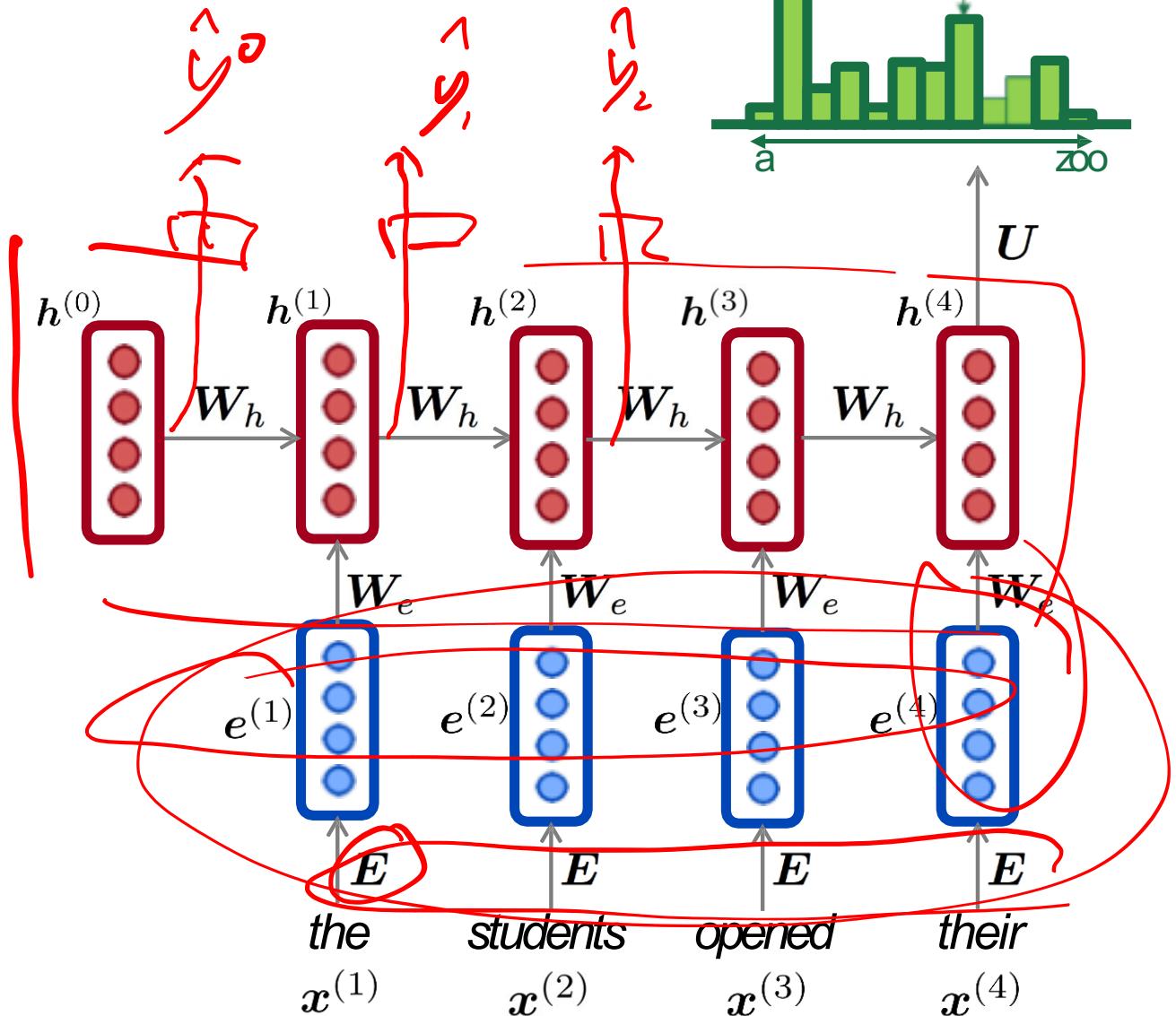
$\mathbf{h}^{(0)}$ is the initial hidden state

word embeddings

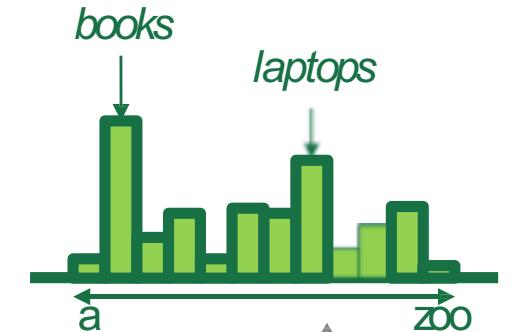
$$\mathbf{e}^{(t)} = \mathbf{E} \mathbf{x}^{(t)}$$

words / one-hot vectors

$$\mathbf{x}^{(t)} \in \mathbb{R}^{|V|}$$



Note: this input sequence could be much longer, but this slide doesn't have space!



RNN

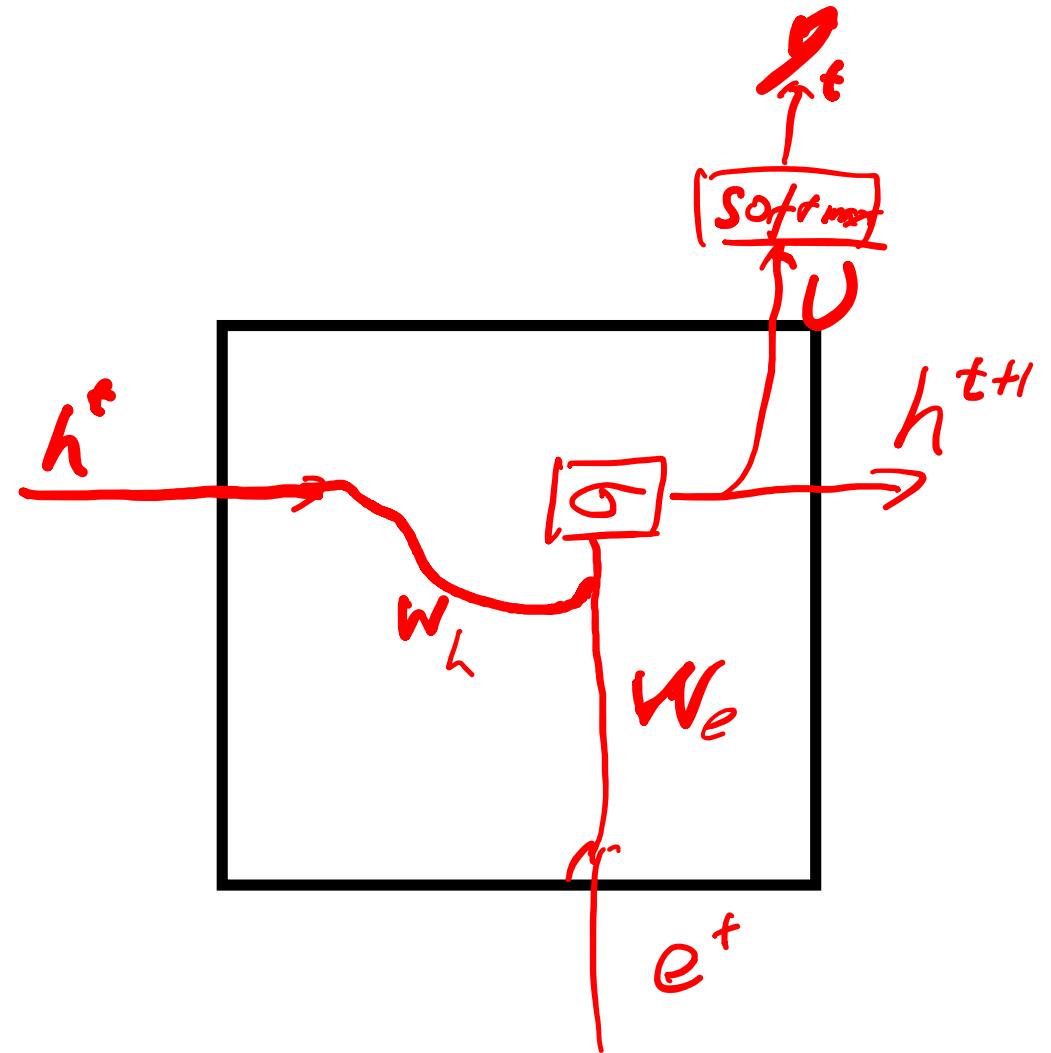
output distribution

$$\hat{y}^{(t)} = \text{softmax}(\mathbf{U}h^{(t)} + \mathbf{b}_2) \in \mathbb{R}^{|V|}$$

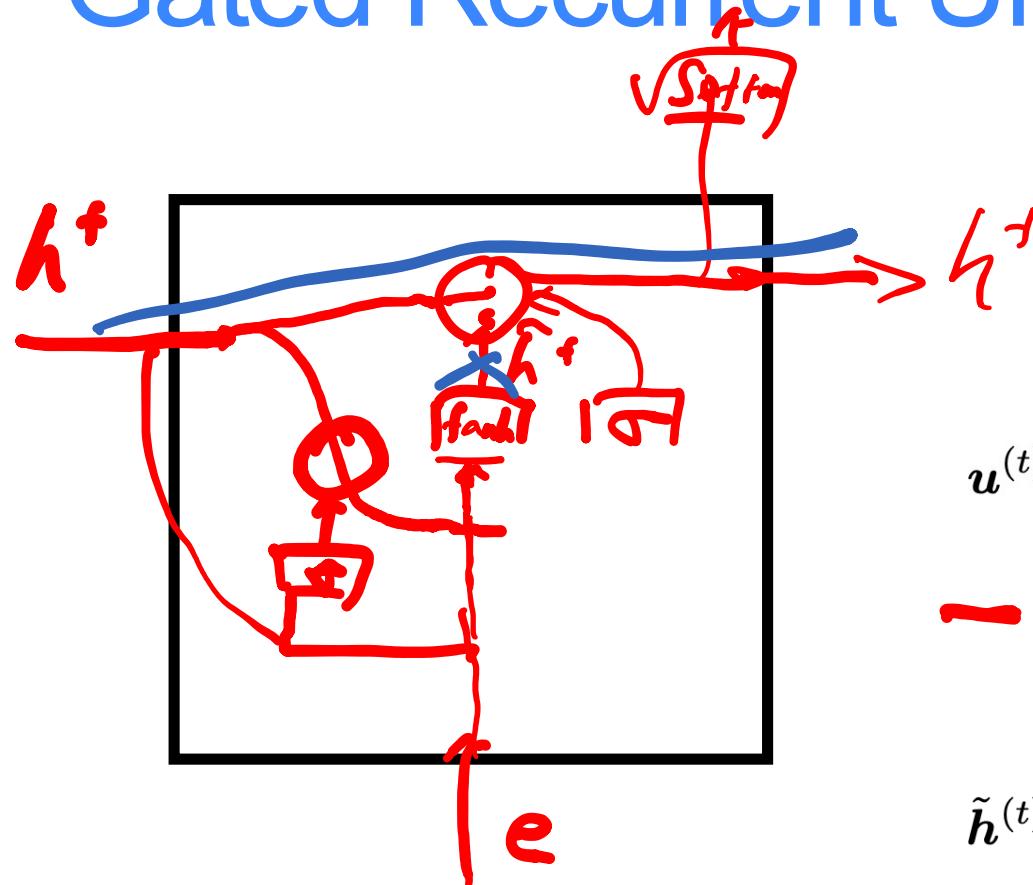
hidden states

$$h^{(t)} = \sigma(\mathbf{W}_h h^{(t-1)} + \mathbf{W}_e e^{(t)} + \mathbf{b}_1)$$

$h^{(0)}$ is the initial hidden state



Gated Recurrent Units (GRU)



$$u^{(t)} = \sigma \left(W_u h^{(t-1)} + U_u x^{(t)} + b_u \right)$$

—

$$\tilde{h}^{(t)} = \tanh \left(W_h (\quad) h^{(t-1)} + U_h x^{(t)} + b_h \right)$$

$$h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

How does this solve vanishing gradient?

GRU makes it easier to retain info long-term
(e.g. by setting update gate to 0)

GRU vs. LSTM

GRU

$$\begin{aligned} \underline{\mathbf{u}}^{(t)} &= \sigma \left(\mathbf{W}_u \mathbf{h}^{(t-1)} + \mathbf{U}_u \mathbf{x}^{(t)} + \mathbf{b}_u \right) \\ \underline{\mathbf{r}}^{(t)} &= \sigma \left(\mathbf{W}_r \mathbf{h}^{(t-1)} + \mathbf{U}_r \mathbf{x}^{(t)} + \mathbf{b}_r \right) \end{aligned}$$

$$\begin{aligned} \tilde{\mathbf{h}}^{(t)} &= \tanh \left(\mathbf{W}_h (\underline{\mathbf{r}}^{(t)} \circ \mathbf{h}^{(t-1)}) + \mathbf{U}_h \mathbf{x}^{(t)} + \mathbf{b}_h \right) \\ \mathbf{h}^{(t)} &= (1 - \underline{\mathbf{u}}^{(t)}) \circ \mathbf{h}^{(t-1)} + \underline{\mathbf{u}}^{(t)} \circ \tilde{\mathbf{h}}^{(t)} \end{aligned}$$

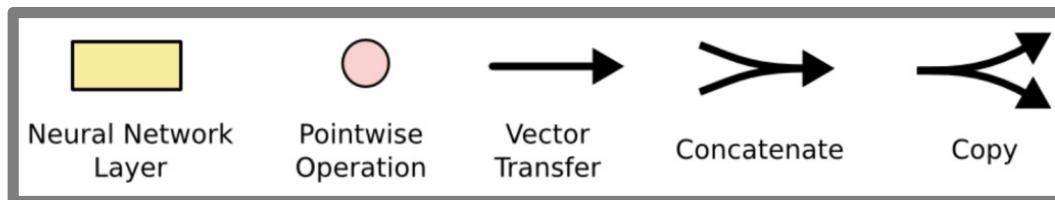
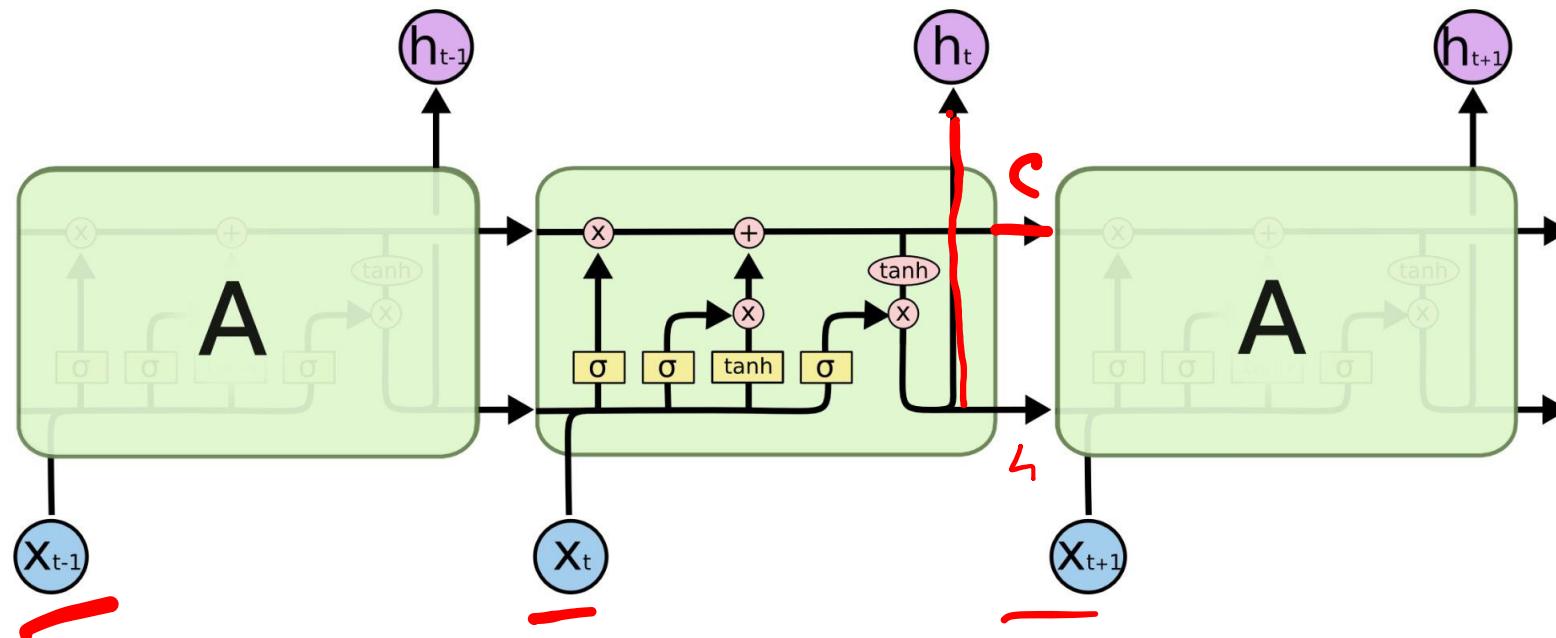
LSTM

$$\begin{aligned} \mathbf{f}^{(t)} &= \sigma \left(\mathbf{W}_f \mathbf{h}^{(t-1)} + \mathbf{U}_f \mathbf{x}^{(t)} + \mathbf{b}_f \right) \\ \mathbf{i}^{(t)} &= \sigma \left(\mathbf{W}_i \mathbf{h}^{(t-1)} + \mathbf{U}_i \mathbf{x}^{(t)} + \mathbf{b}_i \right) \\ \mathbf{o}^{(t)} &= \sigma \left(\mathbf{W}_o \mathbf{h}^{(t-1)} + \mathbf{U}_o \mathbf{x}^{(t)} + \mathbf{b}_o \right) \end{aligned}$$

$$\begin{aligned} \tilde{\mathbf{c}}^{(t)} &= \tanh \left(\mathbf{W}_c \mathbf{h}^{(t-1)} + \mathbf{U}_c \mathbf{x}^{(t)} + \mathbf{b}_c \right) \\ \mathbf{c}^{(t)} &= \mathbf{f}^{(t)} \circ \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \circ \tilde{\mathbf{c}}^{(t)} \\ \mathbf{h}^{(t)} &= \mathbf{o}^{(t)} \circ \tanh \mathbf{c}^{(t)} \end{aligned}$$

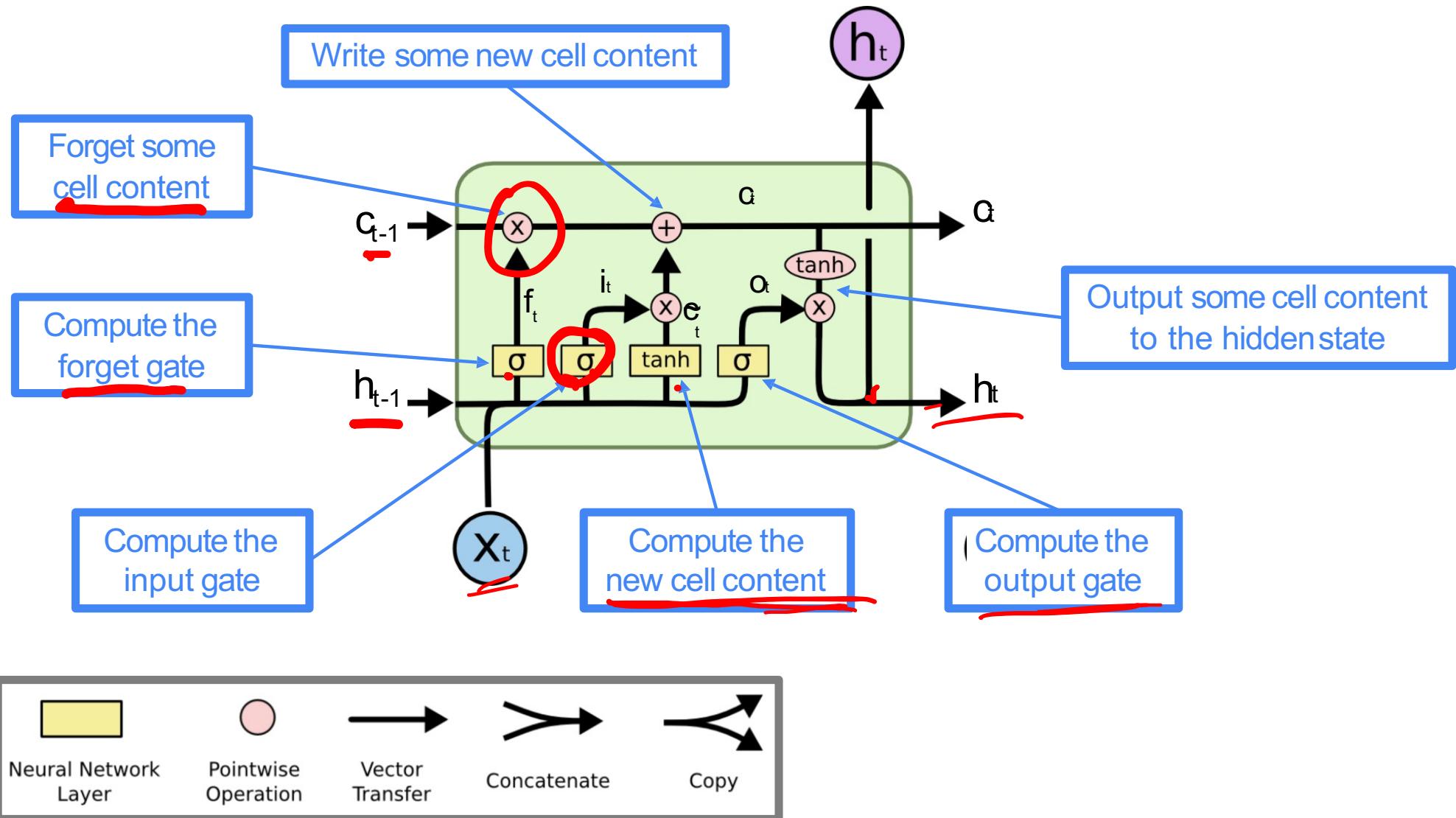
Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



Is vanishing/exploding gradient just a RNN problem?

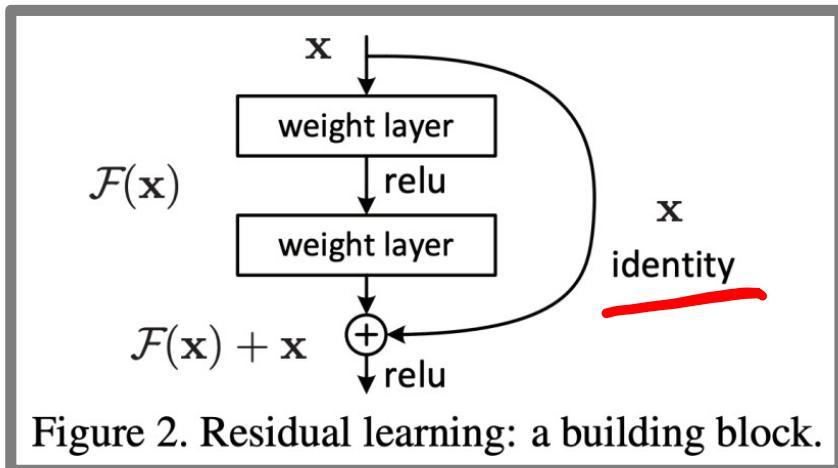


Figure 2. Residual learning: a building block.

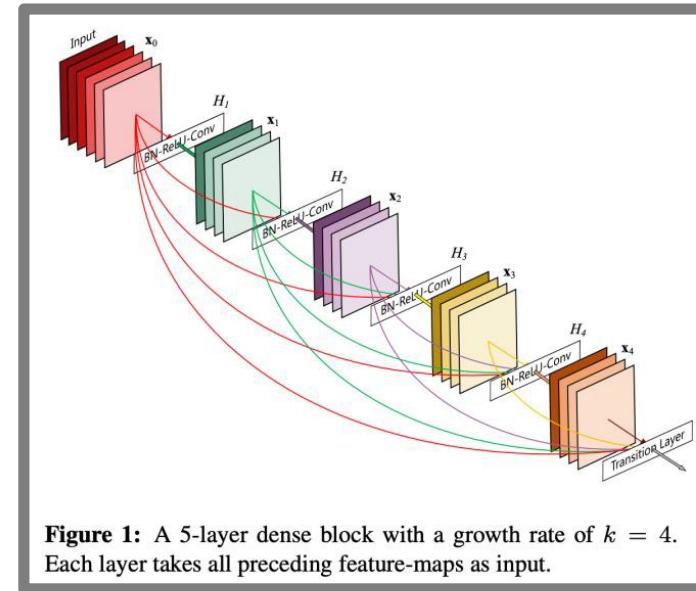
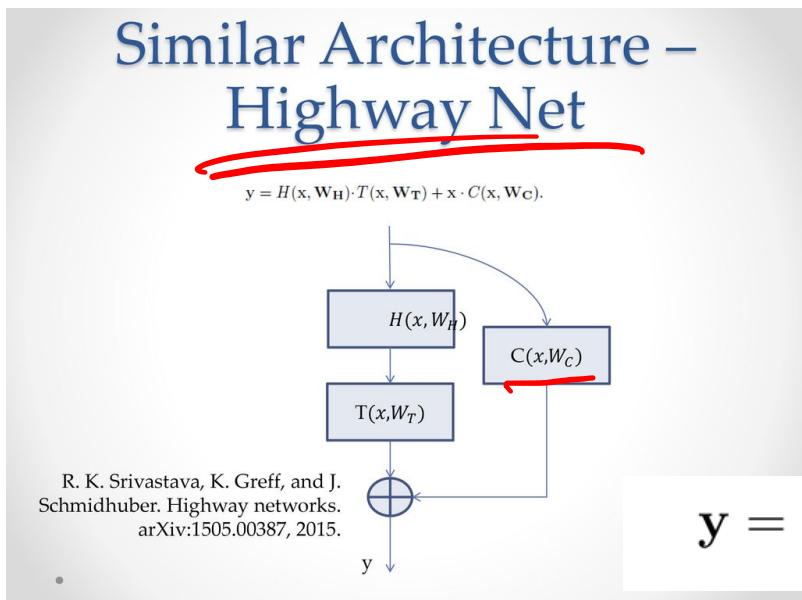


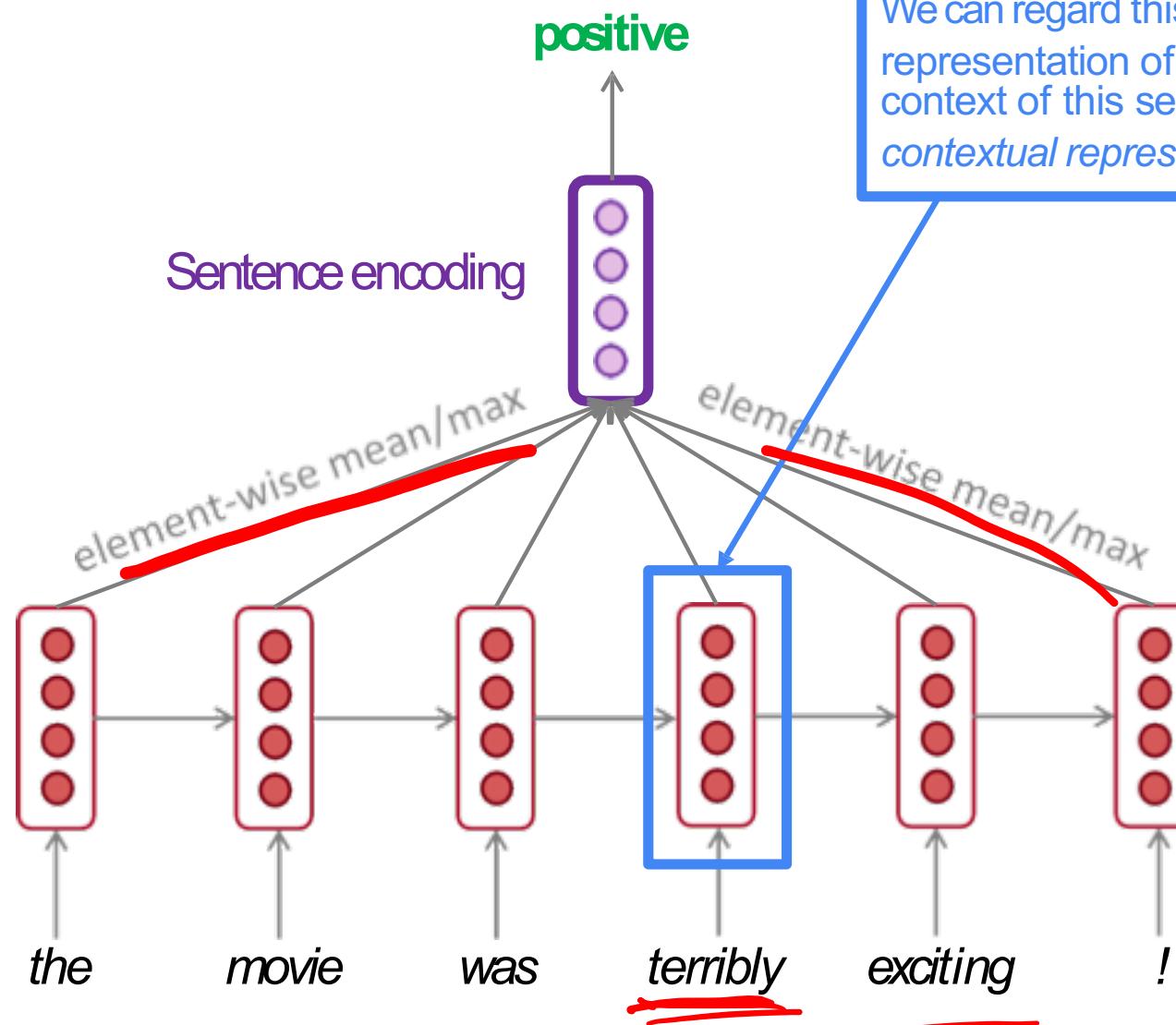
Figure 1: A 5-layer dense block with a growth rate of $k = 4$. Each layer takes all preceding feature-maps as input.



$$y = H(\mathbf{x}, \mathbf{W}_H) \cdot T(\mathbf{x}, \mathbf{W}_T) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W}_T)).$$

Bidirectional RNNs: motivation

Task: Sentiment Classification



We can regard this hidden state as a representation of the word “*terribly*” in the context of this sentence. We call this a *contextual representation*.

These contextual representations only contain information about the *left context* (e.g. “*the movie was*”).

What about *right context*?

In this example, “*exciting*” is in the right context and this modifies the meaning of “*terribly*” (from negative to positive)

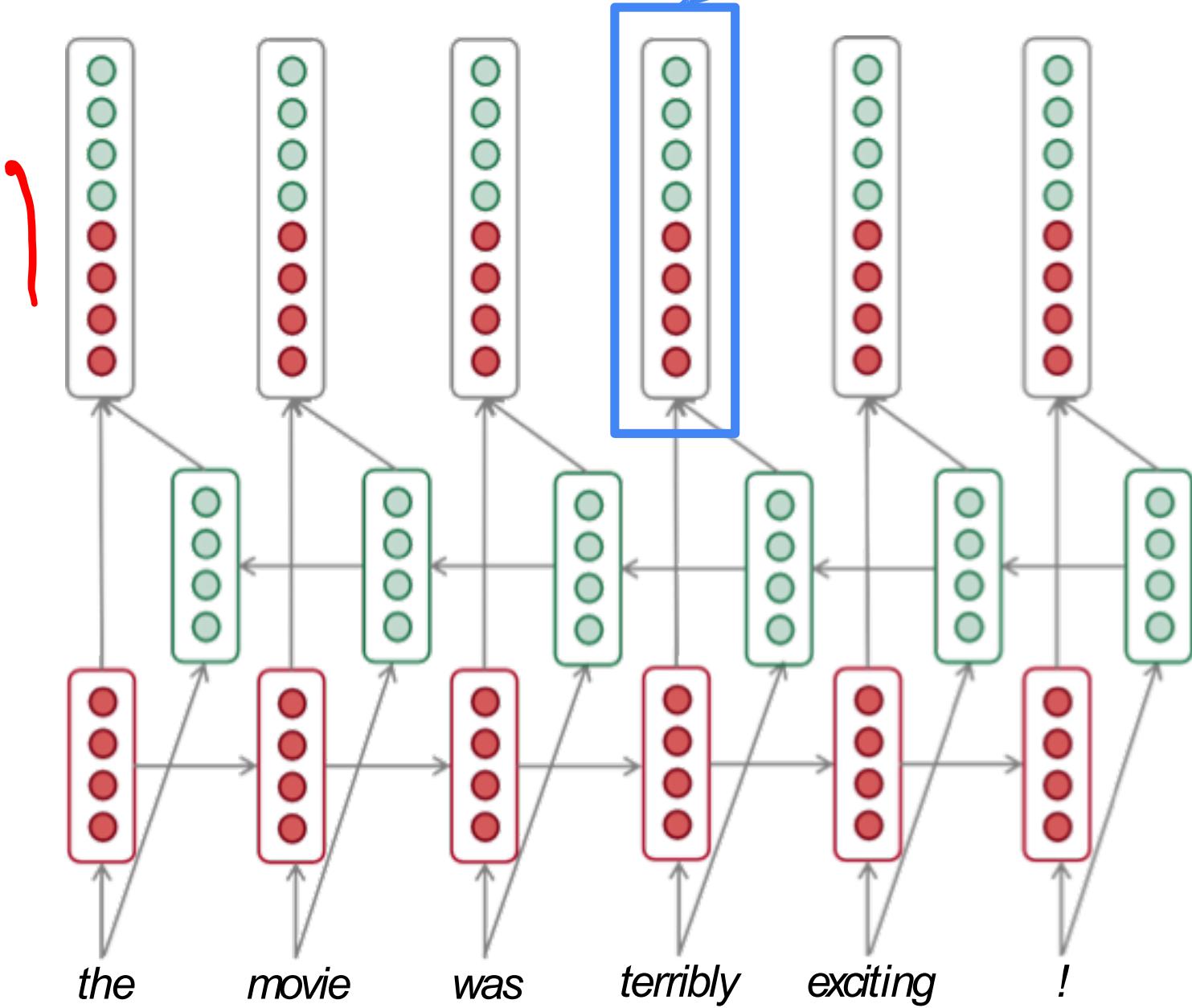
Bidirectional RNNs

This contextual representation of “terribly” has both left and right context!

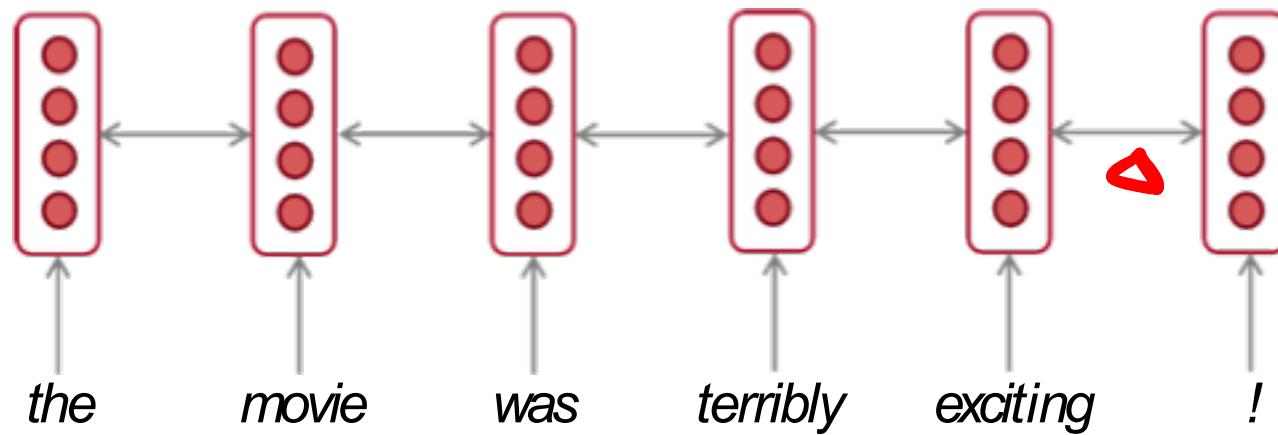
Concatenated hidden states

Backward RNN

Forward RNN



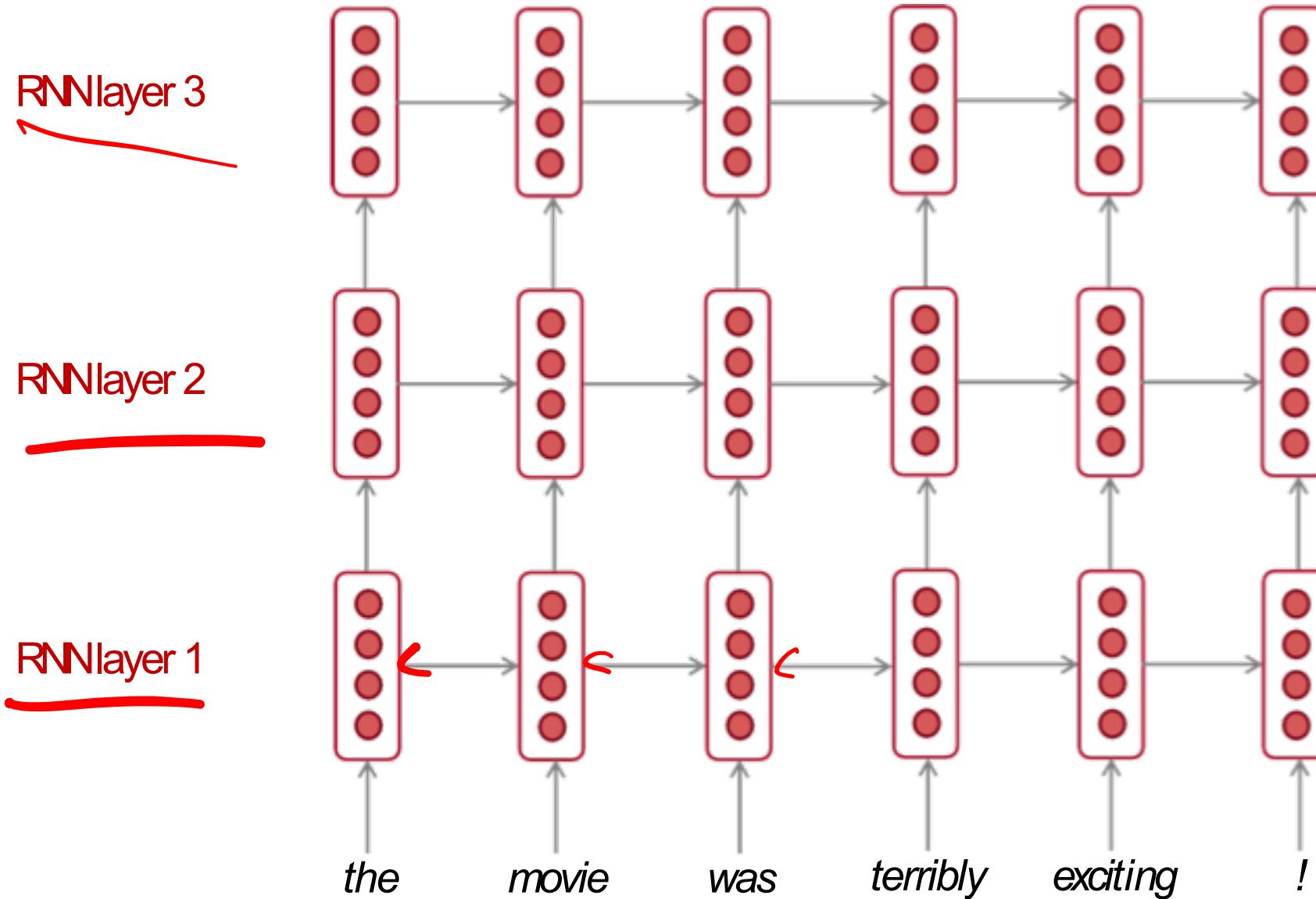
Bidirectional RNNs: simplified diagram



The two-way arrows indicate bidirectionality and the depicted hidden states are assumed to be the concatenated forwards+backwards states.

Multi-layer RNNs

The hidden states from RNNlayer_i are the inputs to RNNlayer_{i+1}



Neural Machine Translation (NMT)

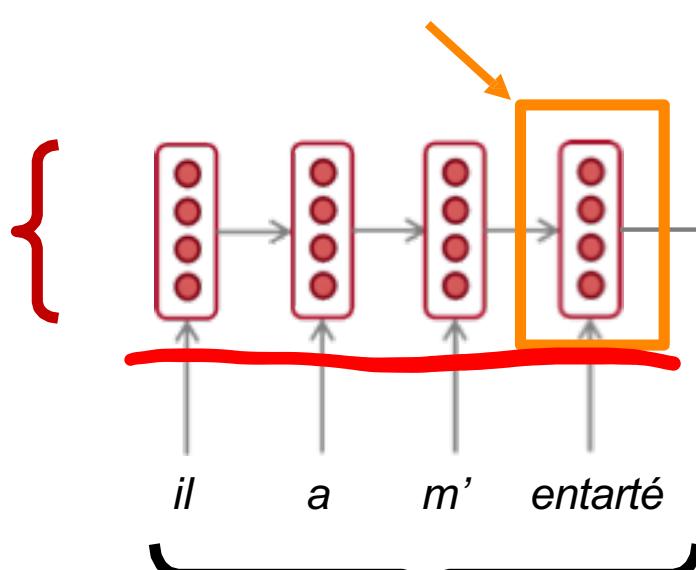
The sequence-to-sequence model

Encoding of the source sentence.

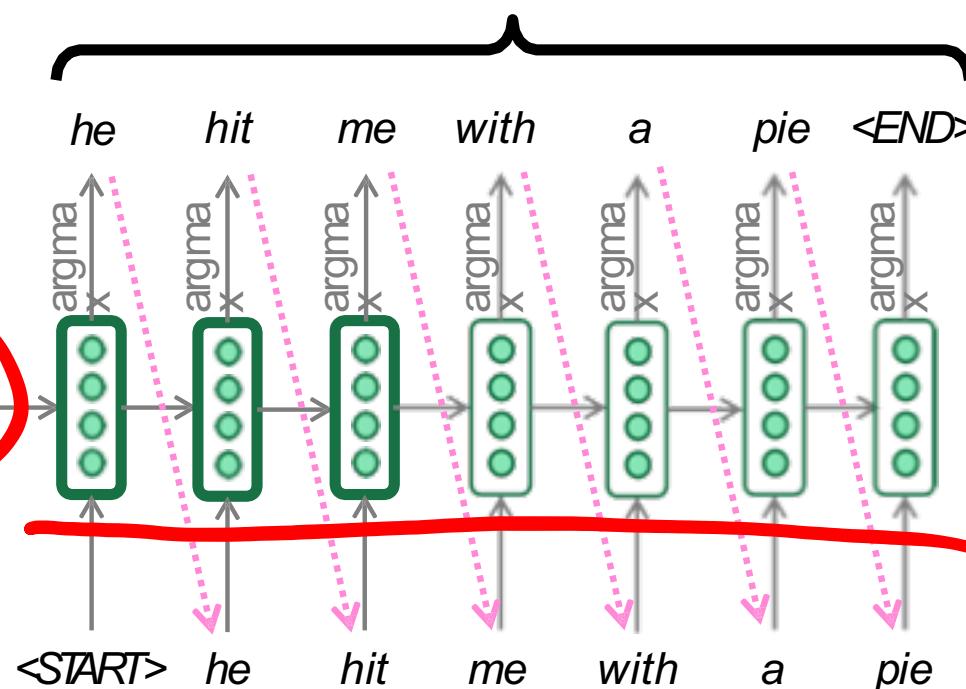
Provides initial hidden state

for Decoder RNN.

Encoder RNN



Target sentence (output)



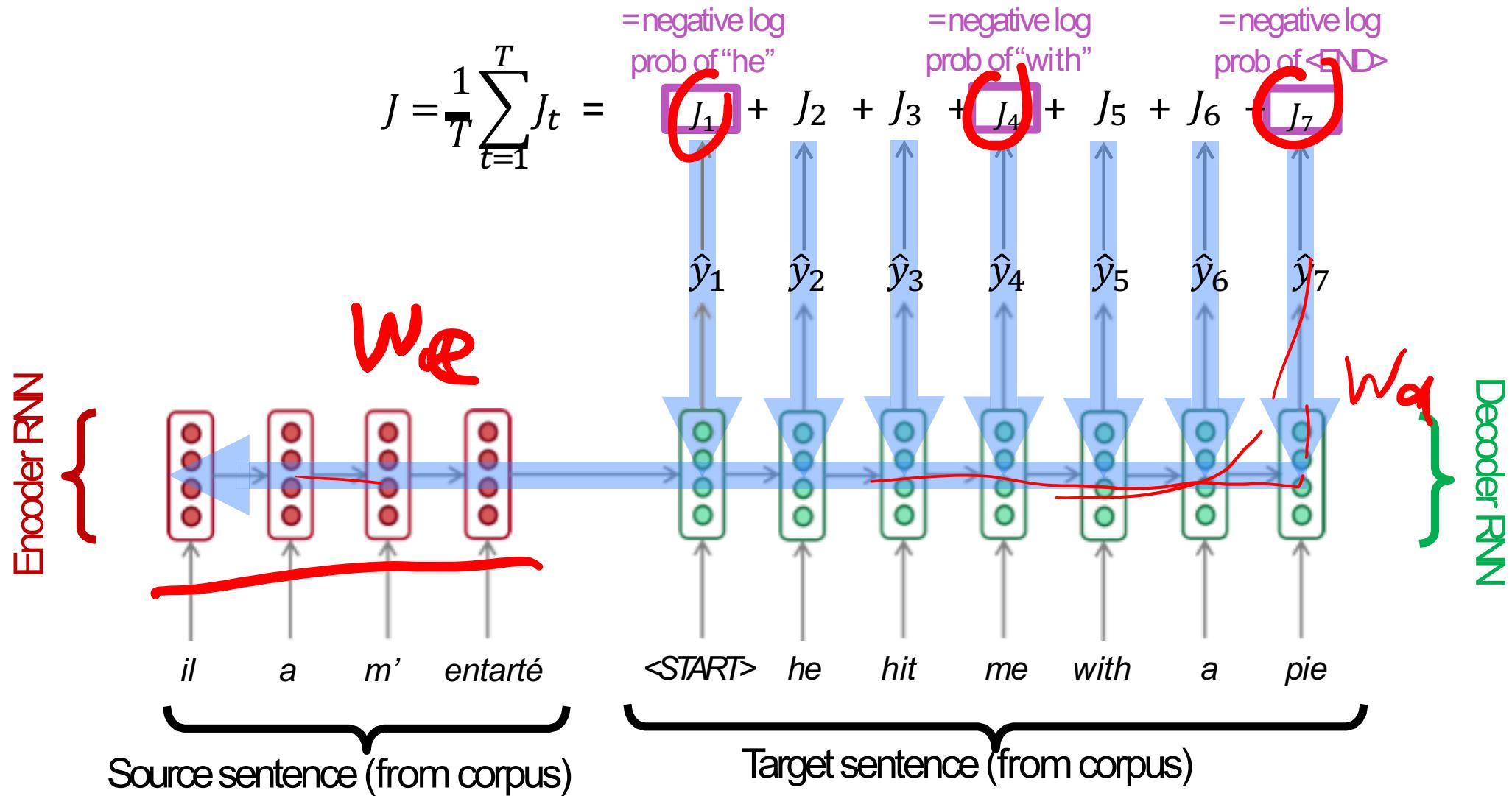
Decoder RNN

Encoder RNN produces
an encoding of the
source sentence.

Decoder RNN is a Language Model that generates
target sentence, conditioned on encoding.

Note: This diagram shows test time behavior:
decoder output is fed in as next step's input

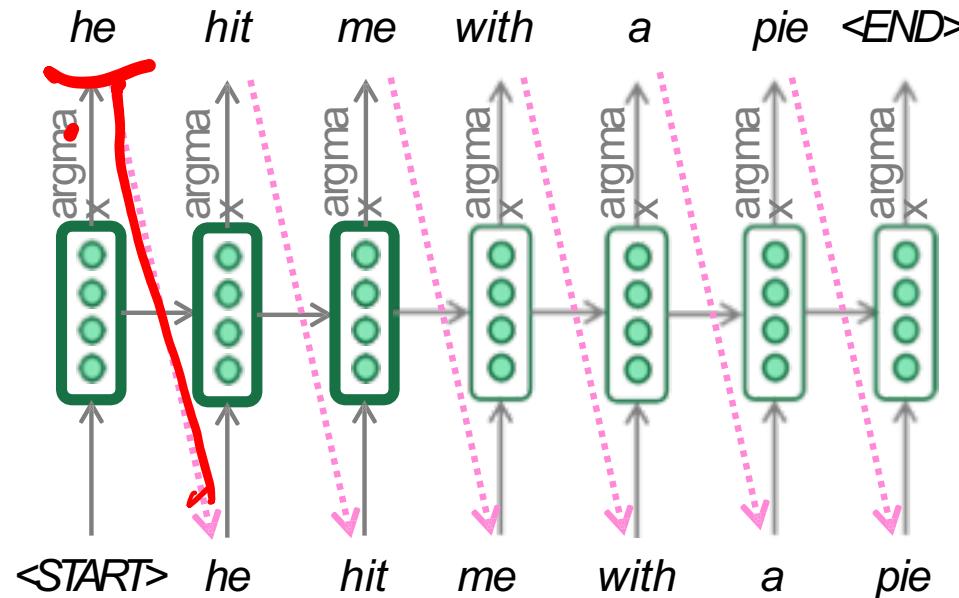
Training a Neural Machine Translation system



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

Greedy Decoding

- We saw how to generate (or “decode”) the target sentence by taking argmax on each step of the decoder



- This is **greedy decoding** (take most probable word on each step)
- **Problems with this method?**

Beam search decoding: example

Beam size = k = 2. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$

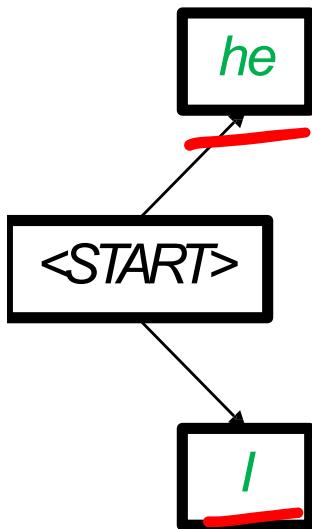
<START>

Calculate prob
dist of next word

Beam search decoding: example

Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$

$$-0.7 = \log P_{\text{LM}}(\text{he} | \text{<START>})$$

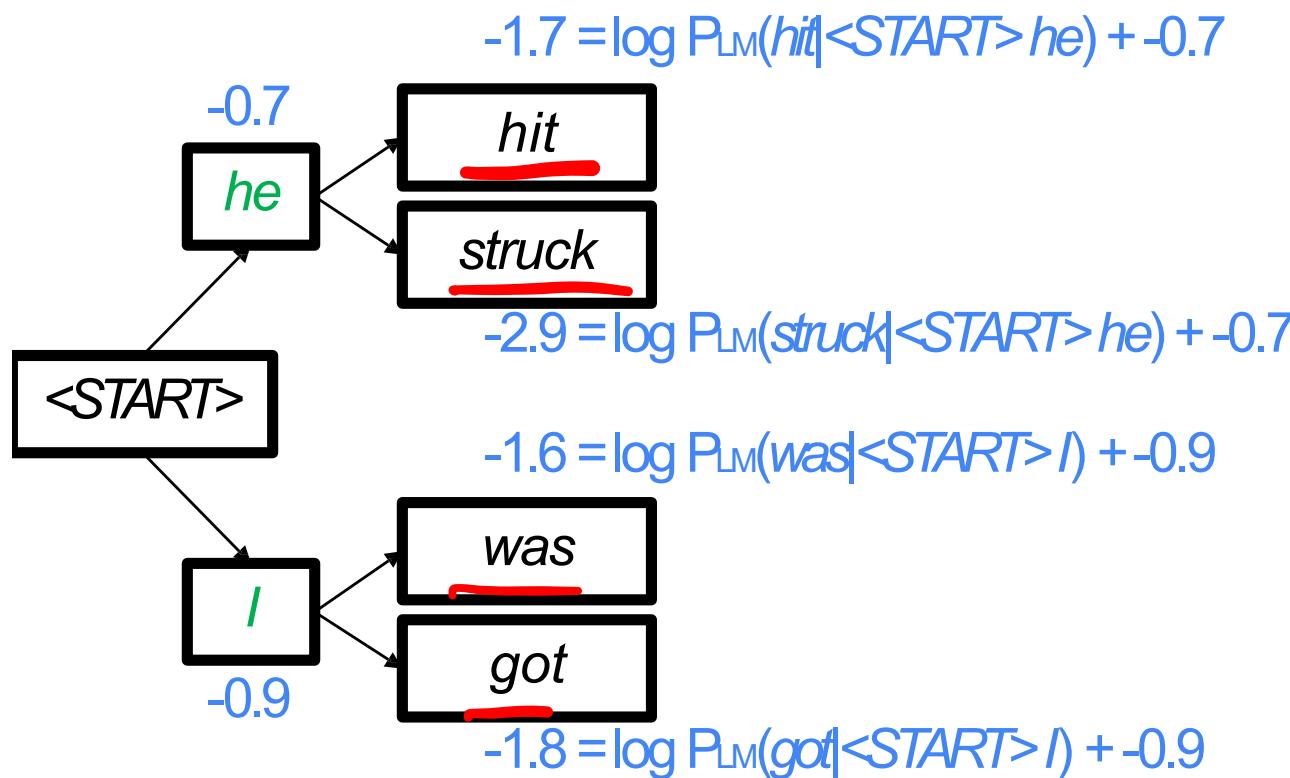


$$-0.9 = \log P_{\text{LM}}(\text{/} | \text{<START>})$$

Take top k words
and compute scores

Beam search decoding: example

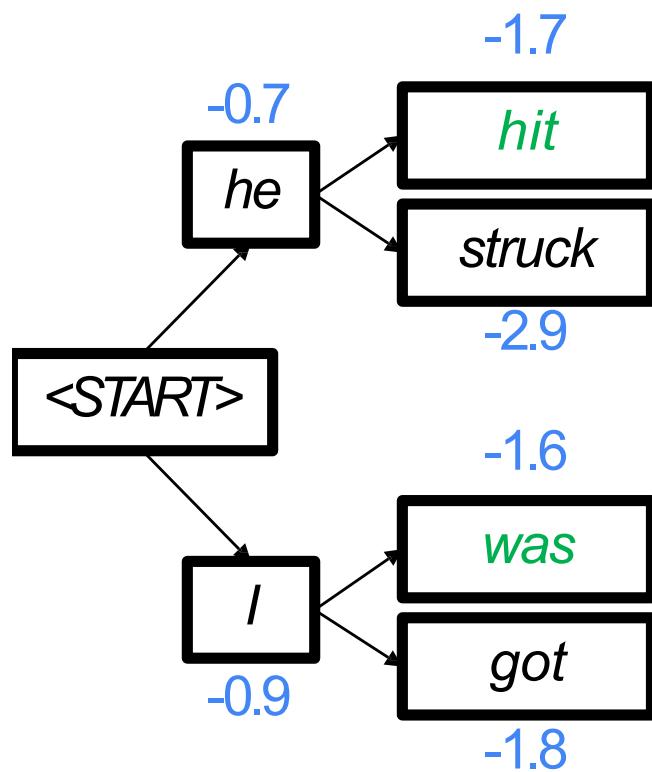
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



For each of the k hypotheses, find
top k next words and calculate scores

Beam search decoding: example

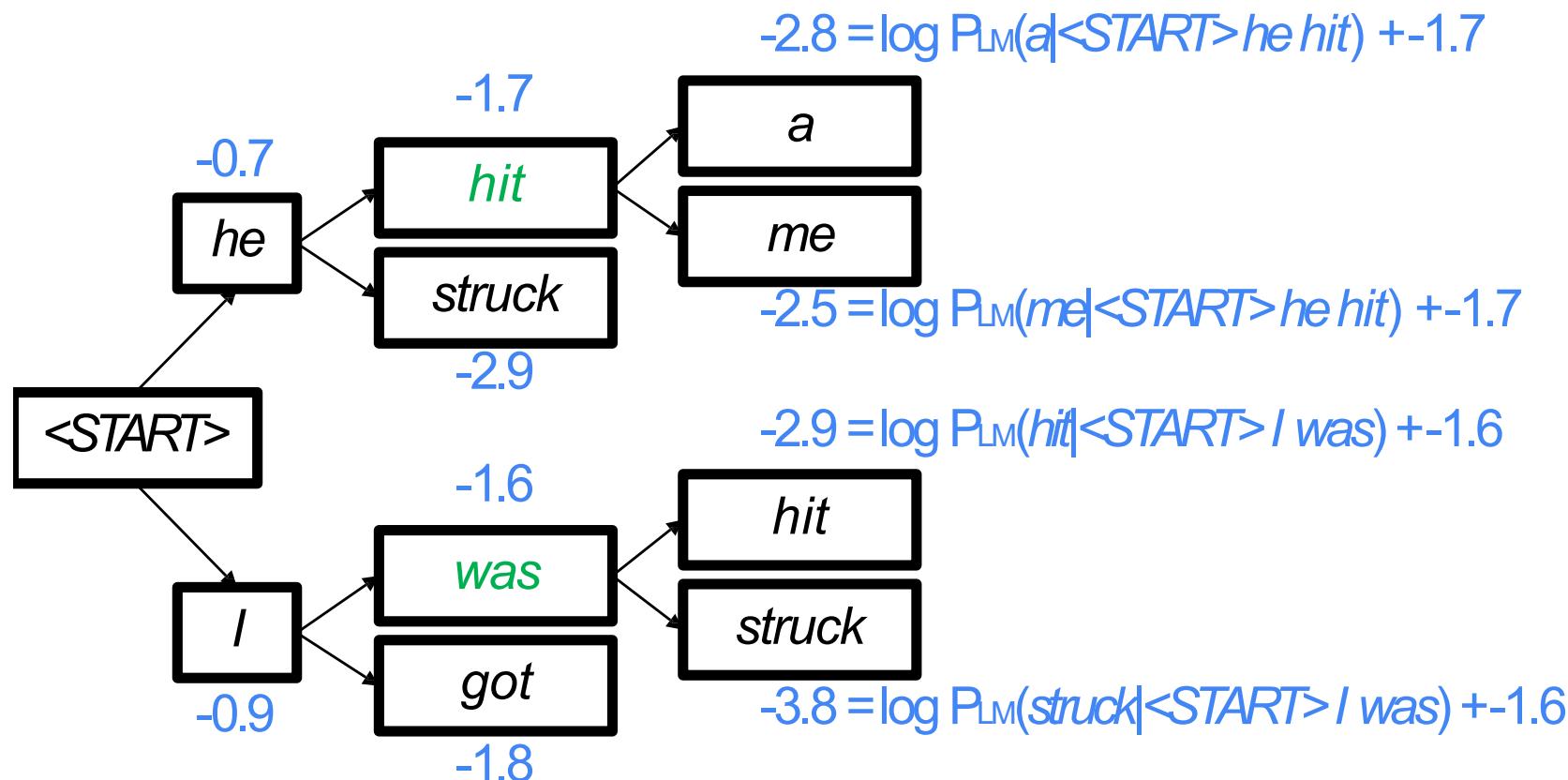
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



Of these k^2 hypotheses,
just keep k with highest scores

Beam search decoding: example

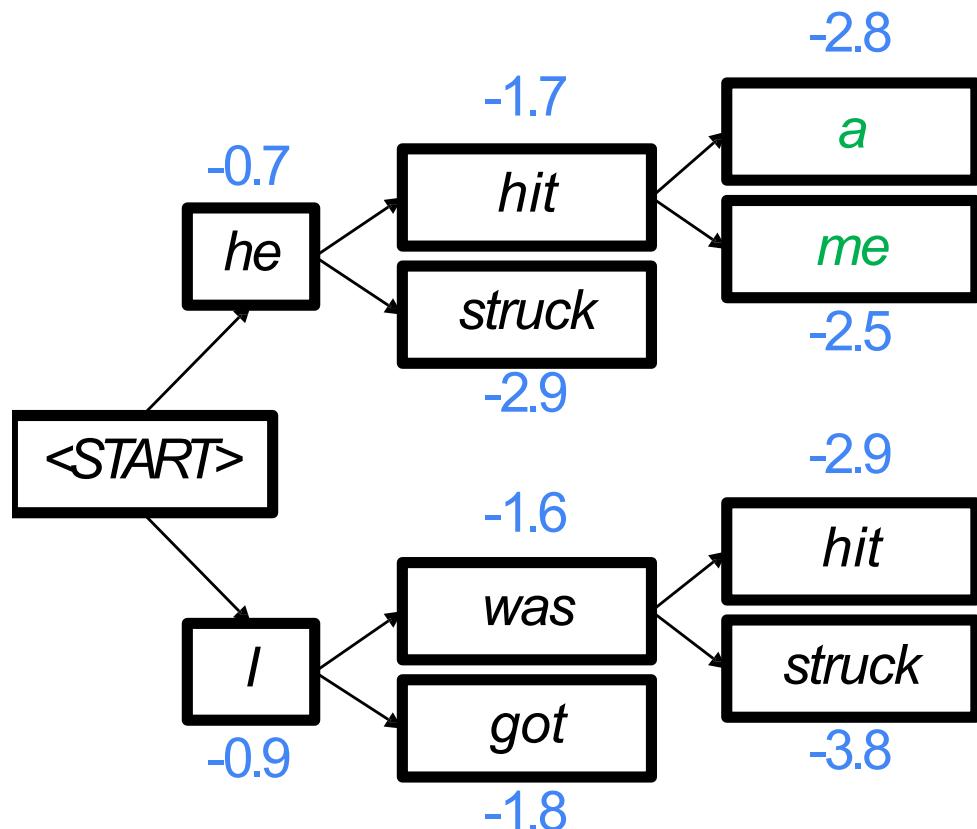
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



For each of the k hypotheses, find
top k next words and calculate scores

Beam search decoding: example

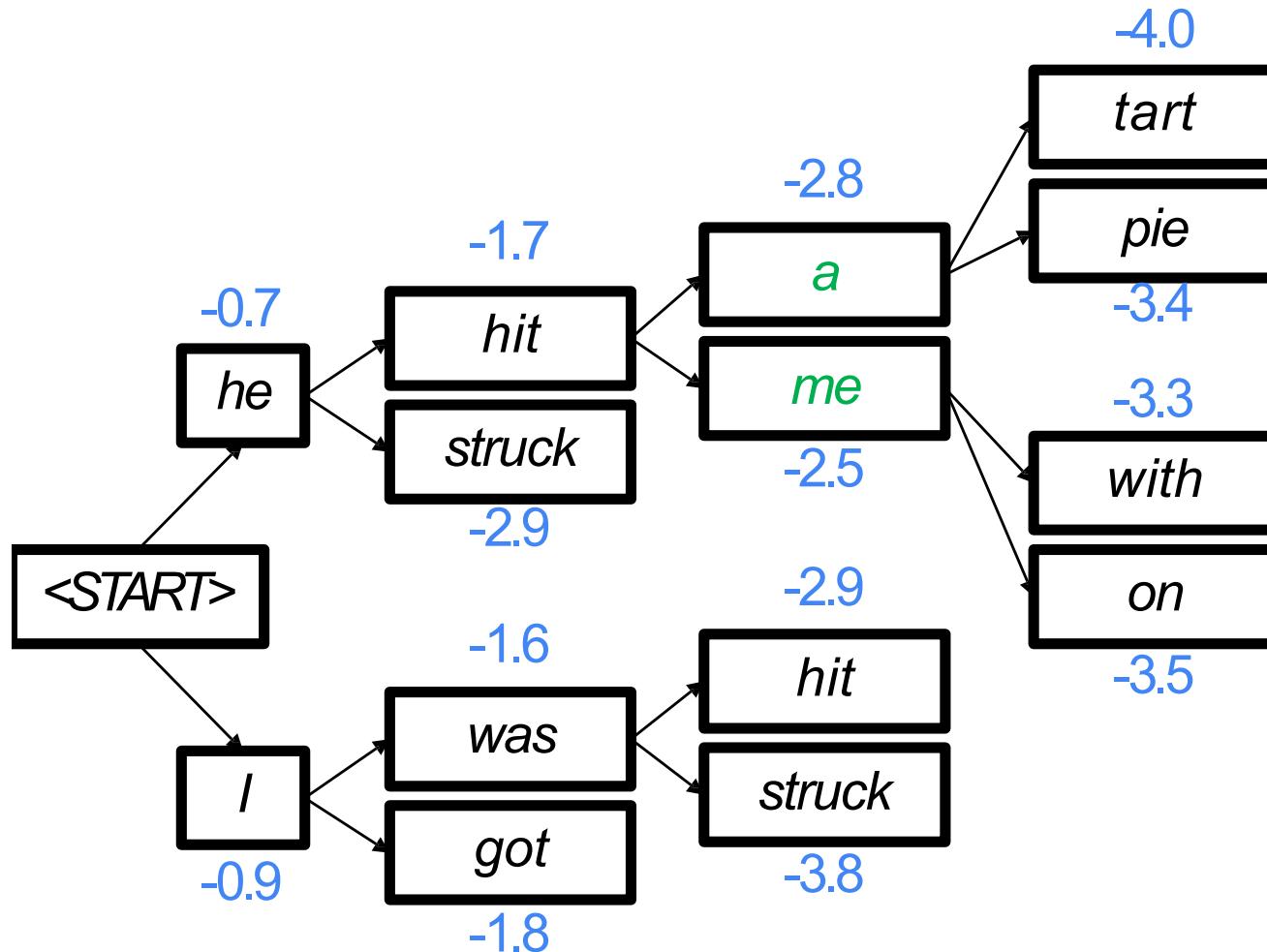
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



Of these k^2 hypotheses,
just keep k with highest scores

Beam search decoding: example

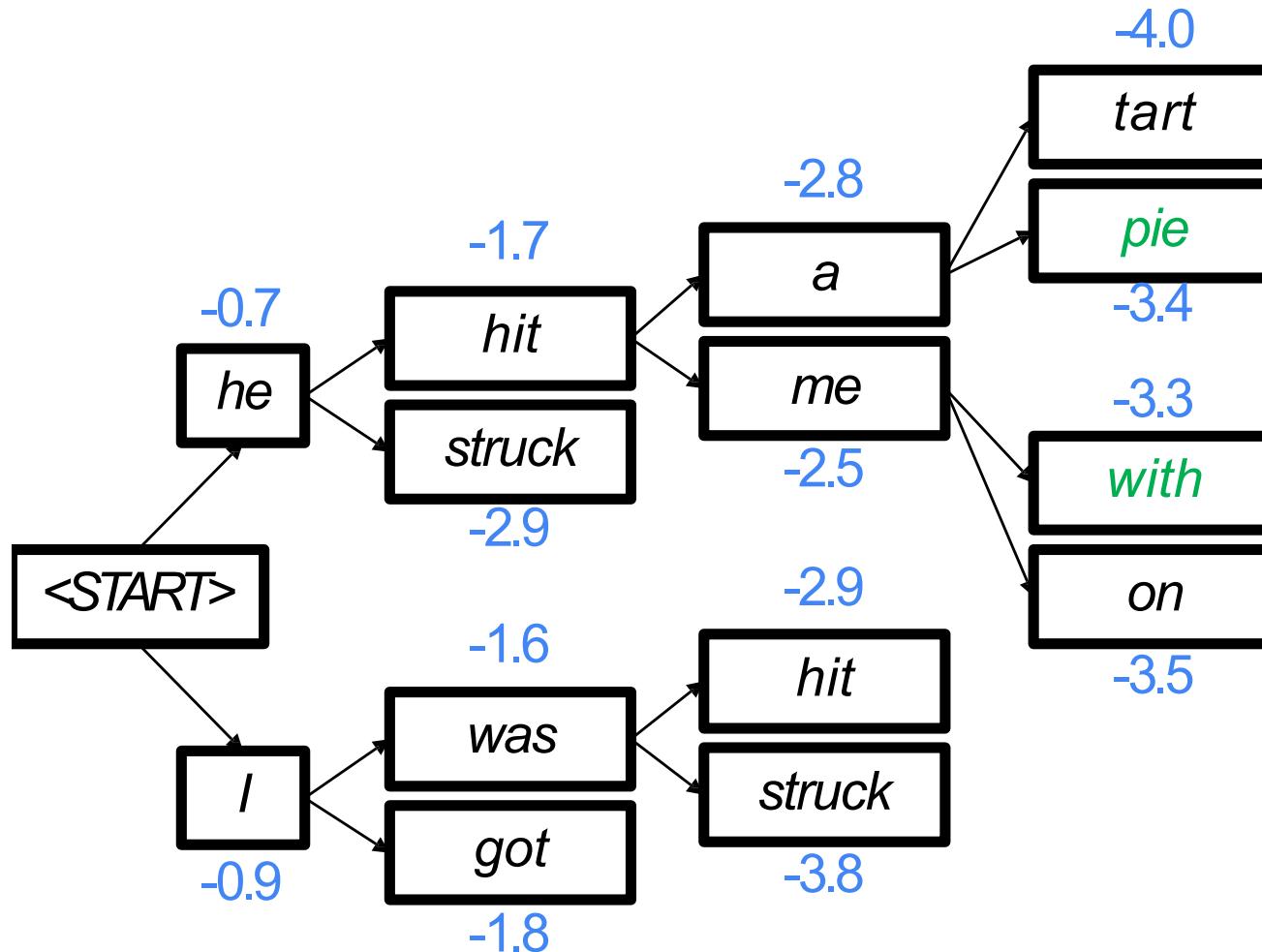
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



For each of the k hypotheses, find top k next words and calculate scores

Beam search decoding: example

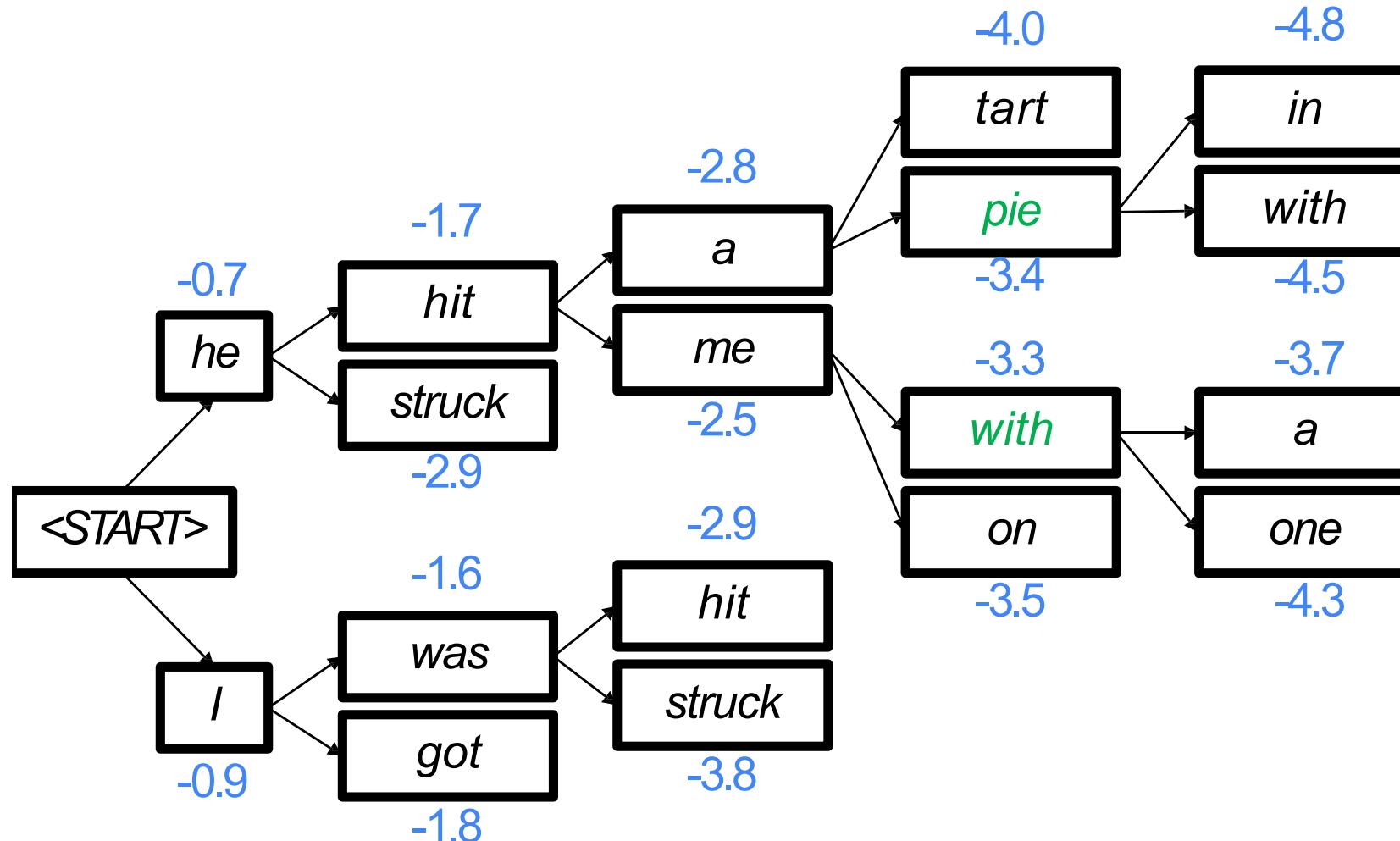
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



Of these k^2 hypotheses,
just keep k with highest scores

Beam search decoding: example

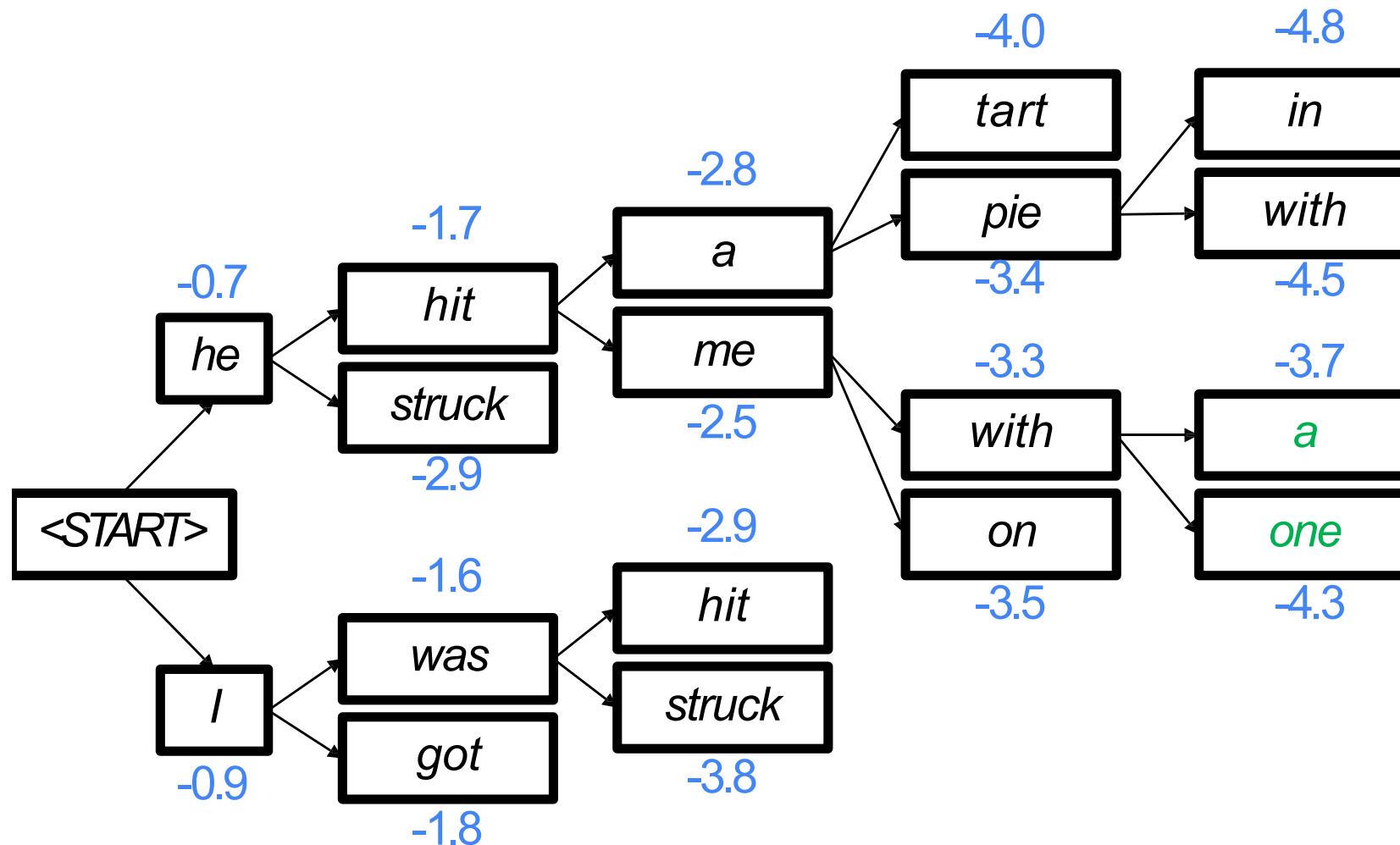
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



For each of the k hypotheses, find top k next words and calculate scores

Beam search decoding: example

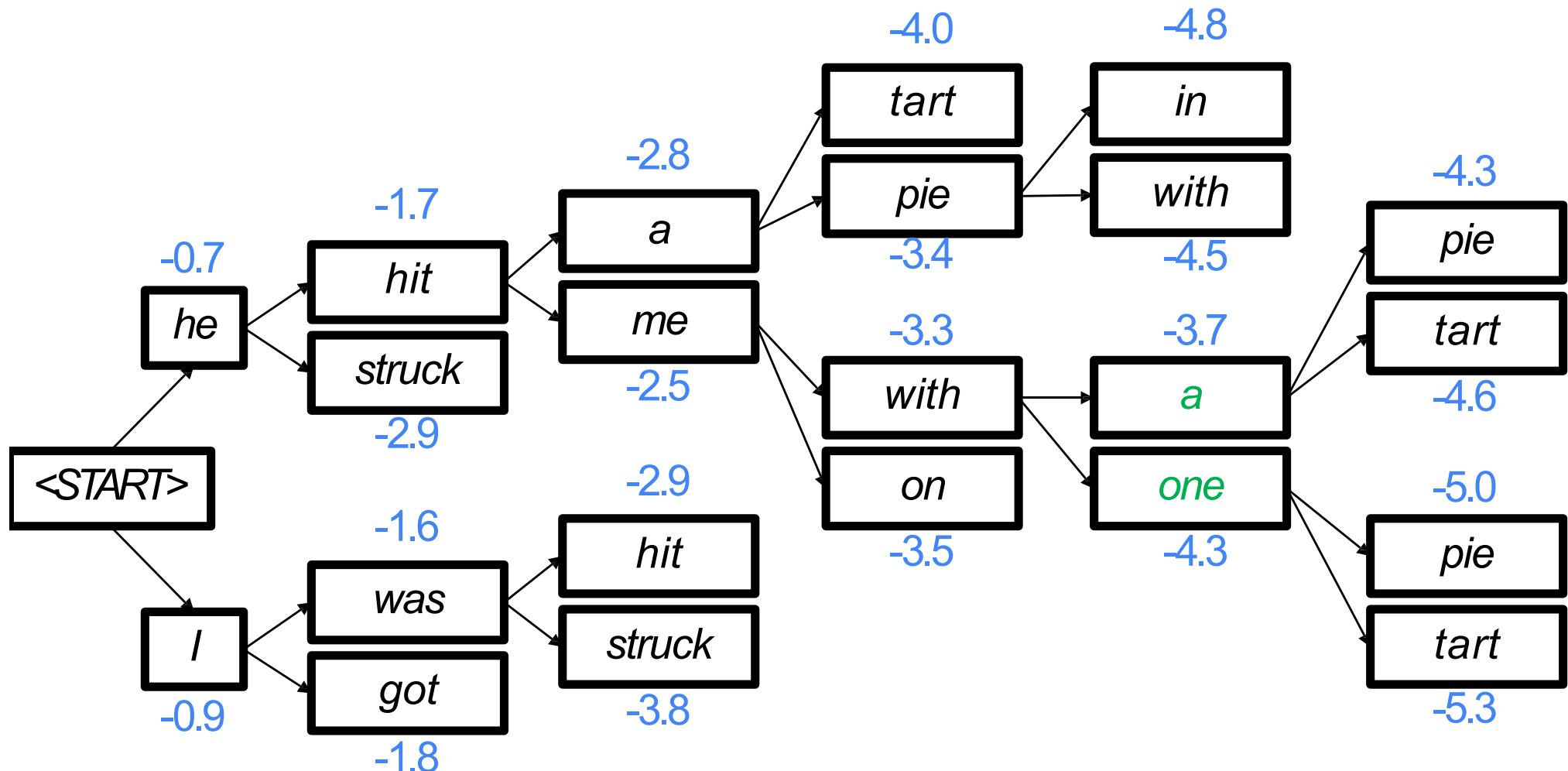
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



Of these k^2 hypotheses,
just keep k with highest scores

Beam search decoding: example

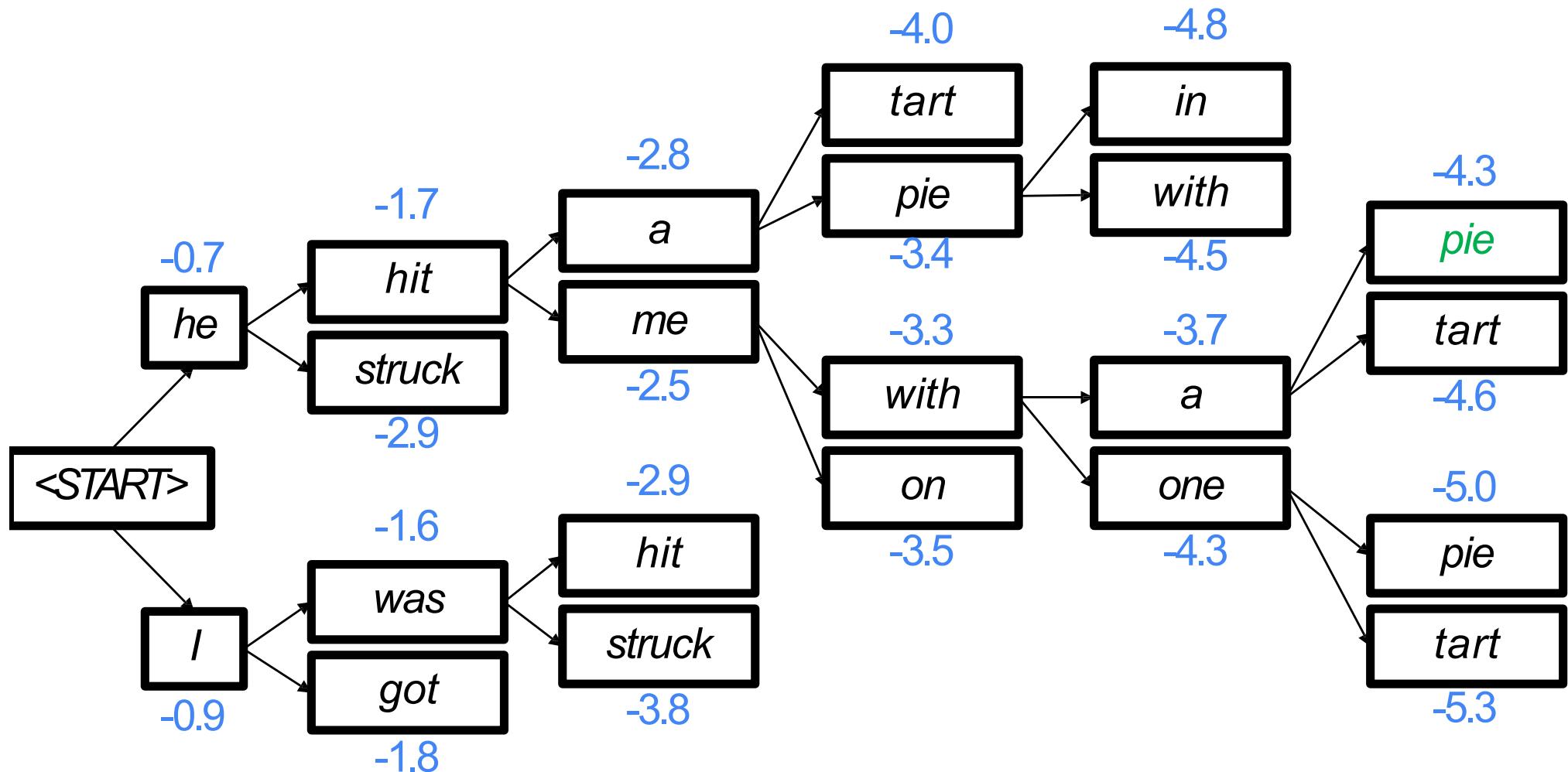
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



For each of the k hypotheses, find top k next words and calculate scores

Beam search decoding: example

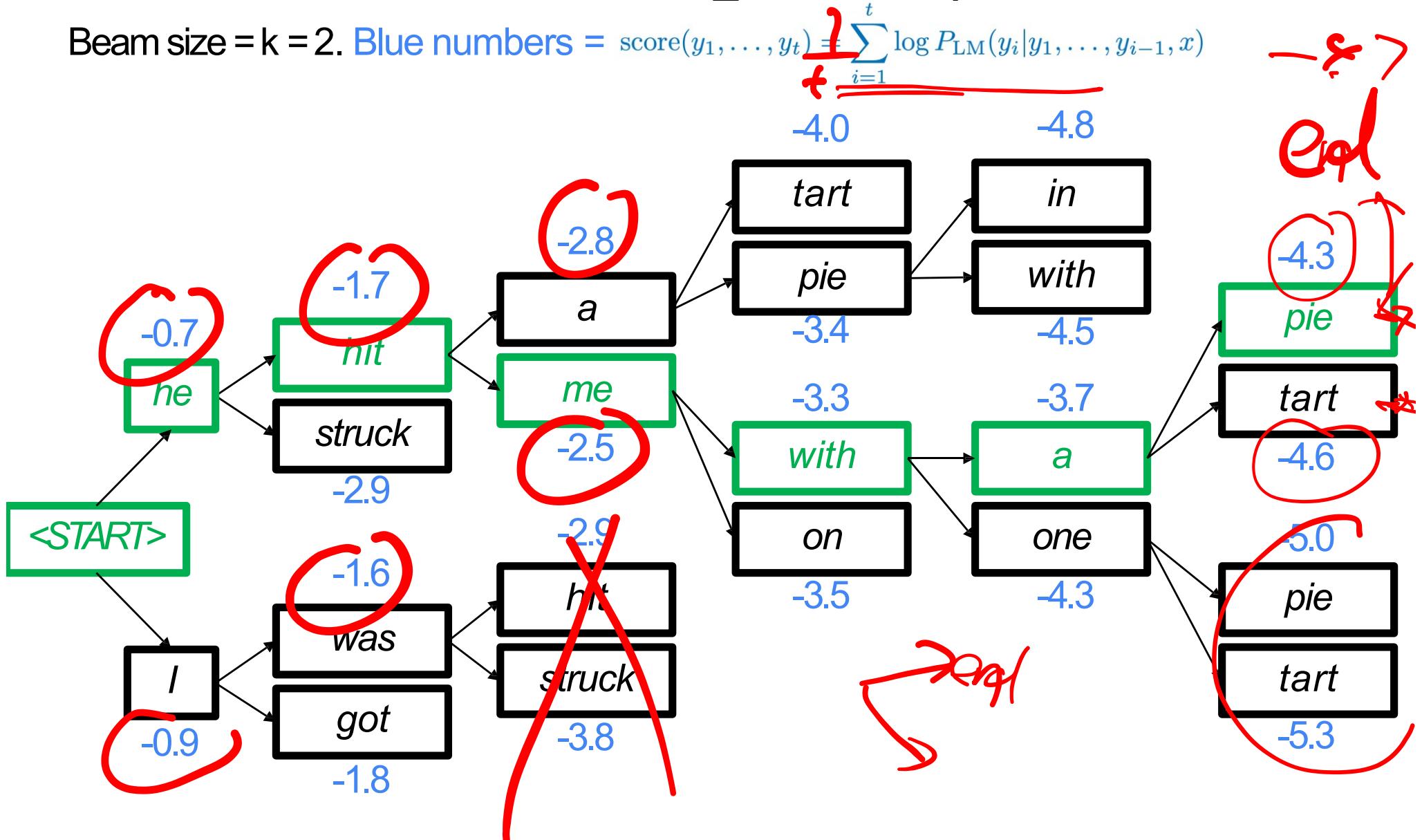
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



This is the top-scoring hypothesis!

Beam search decoding: example

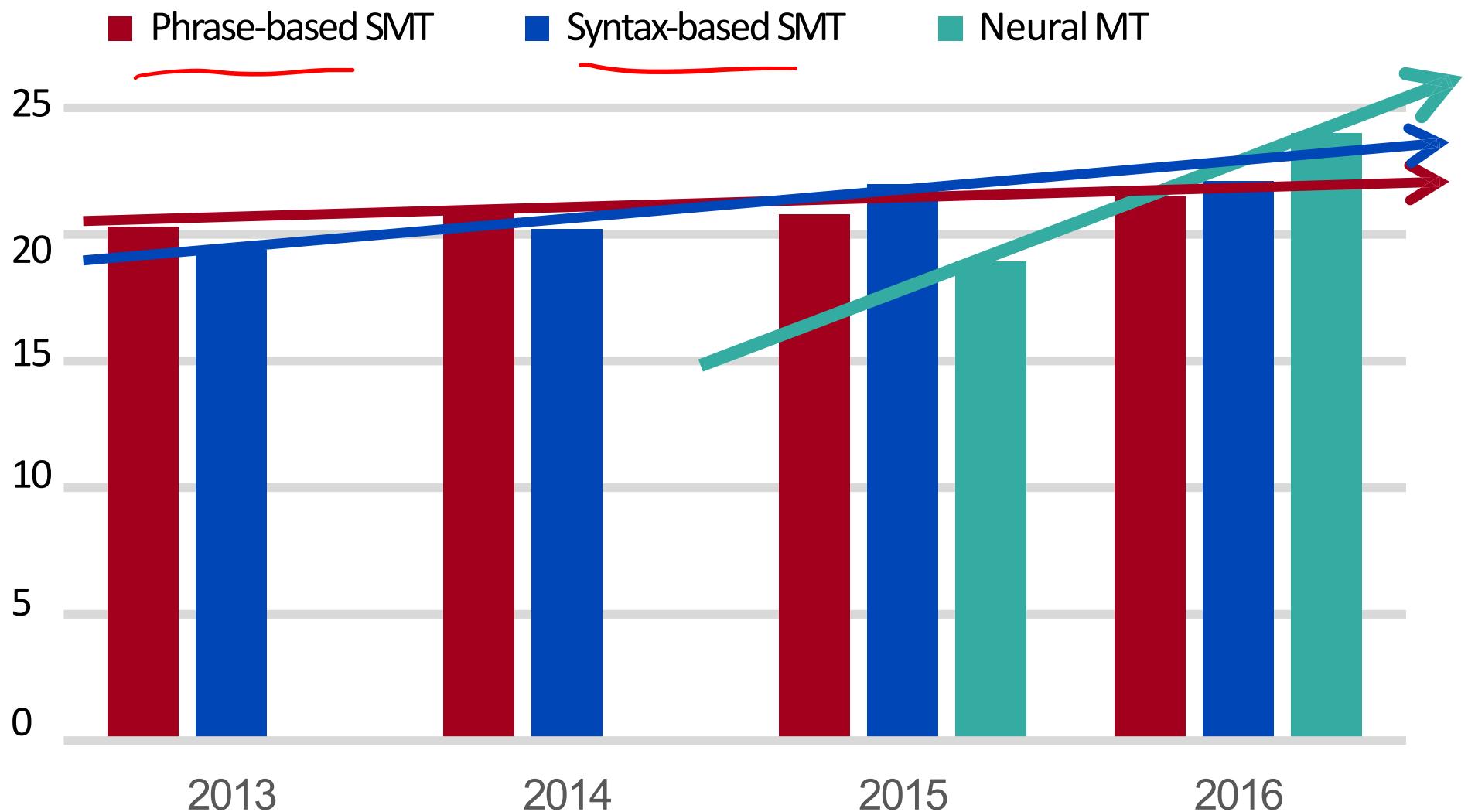
Beam size = $k = 2$. Blue numbers = score(y_1, \dots, y_t) = $\sum_{i=1}^t \log P_{LM}(y_i | y_1, \dots, y_{i-1}, x)$



Backtrack to obtain the full hypothesis

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

So is Machine Translation solved?

- **Nope!**
- Using common sense is still hard

The image shows a screenshot of the Google Translate interface. On the left, under 'English', the text 'paper jam' is displayed with a red underline. On the right, under 'Spanish', the translation 'Mermelada de papel' is shown. The interface includes language selection dropdowns, microphone and speaker icons, and a refresh button.

[Open in Google Translate](#)

[Feedback](#)



So is Machine Translation solved?

- **Nope!**
- NMT picks up **biases** in training data

Malay - detected ▾

English ▾

Dia bekerja sebagai jururawat.

Dia bekerja sebagai pengaturcara. Edit

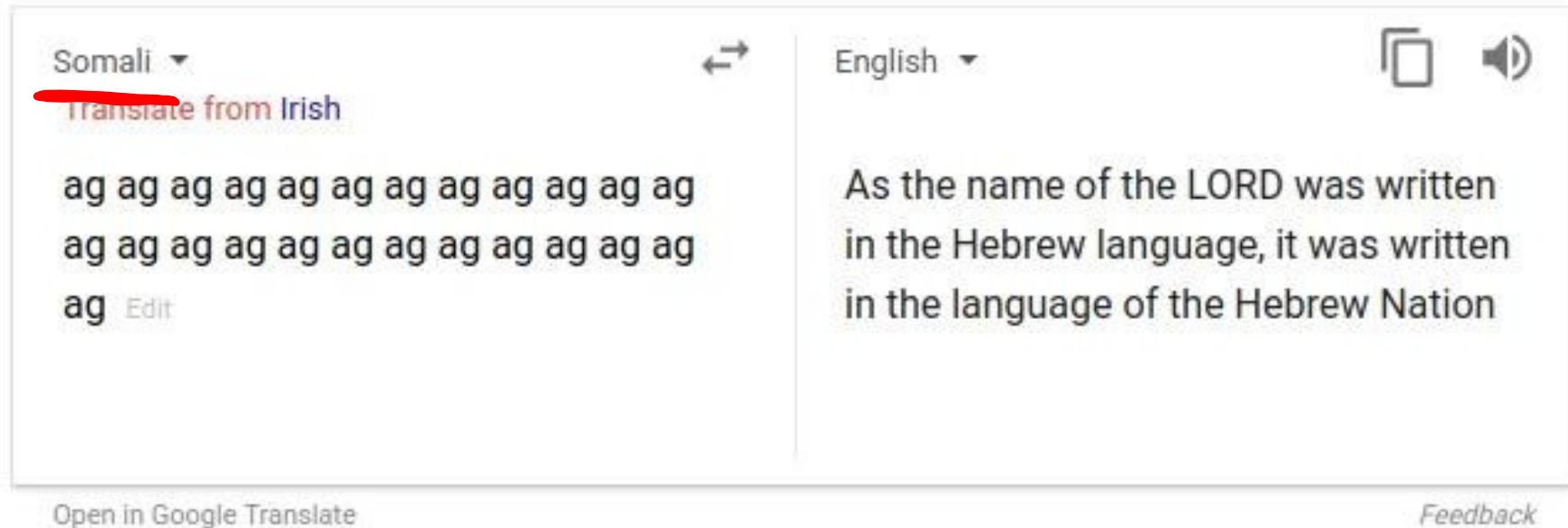
She works as a nurse.

He works as a programmer.

Didn't specify gender

So is Machine Translation solved?

- Nope!
- Uninterpretable systems do strange things



Picture source: https://www.vice.com/en_uk/article/j5npeg/why-is-google-translate-spitting-out-sinister-religious-prophecies

Explanation: <https://www.skynettoday.com/briefs/google-nmt-prophecies>

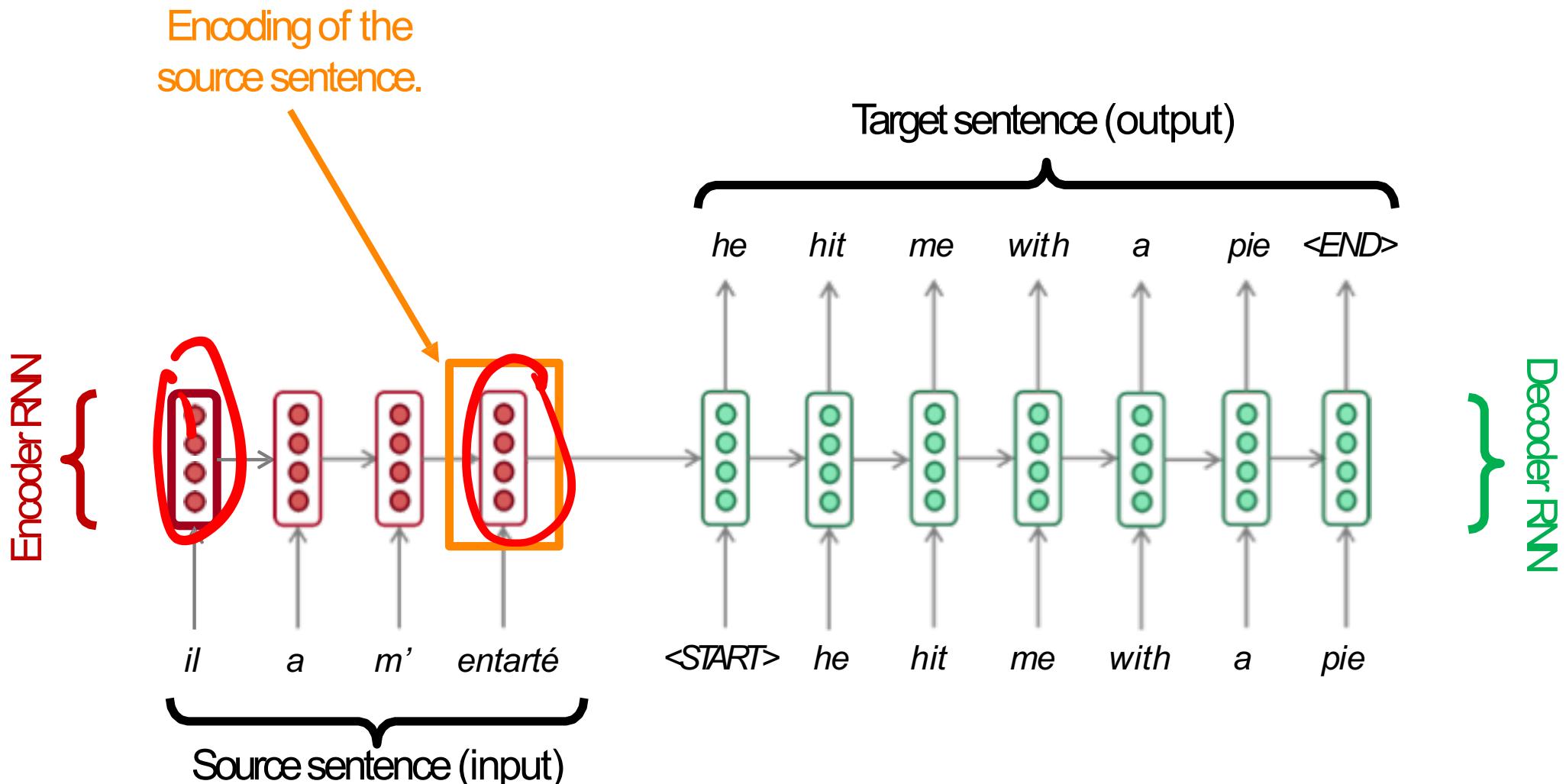
NMT research continues

NMT is the **flagship task** for NLP Deep Learning

- NMT research has **pioneered** many of the recent **innovations** of NLP Deep Learning
- In **2019**: NMT research continues to **thrive**
 - Researchers have found *many, many improvements* to the “vanilla” seq2seq NMT system we’ve presented today
 - But **one improvement** is so integral that it is the new vanilla...

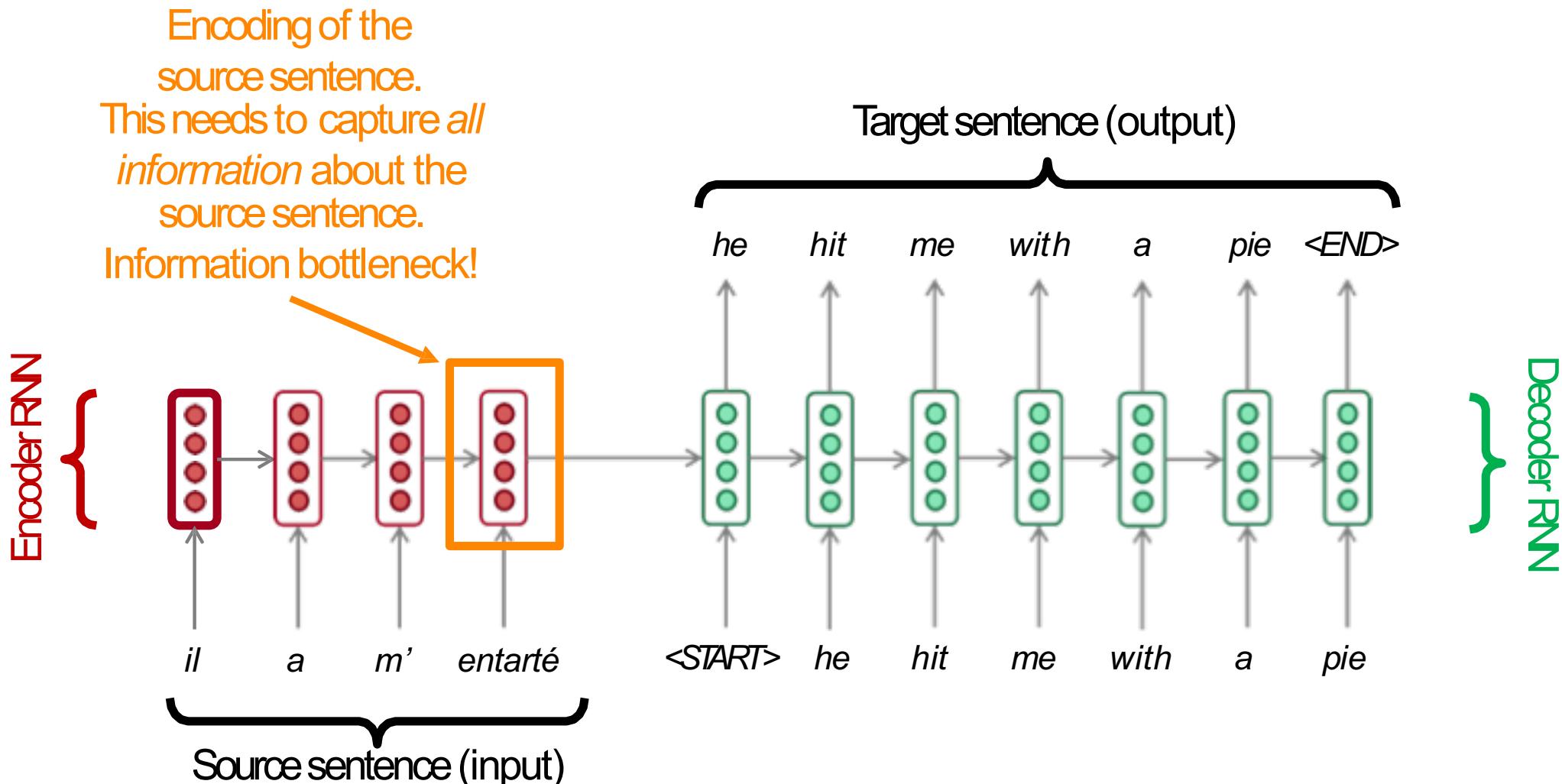
ATTENTION

Sequence-to-sequence



Problems with this architecture?

Sequence-to-sequence: the bottleneck



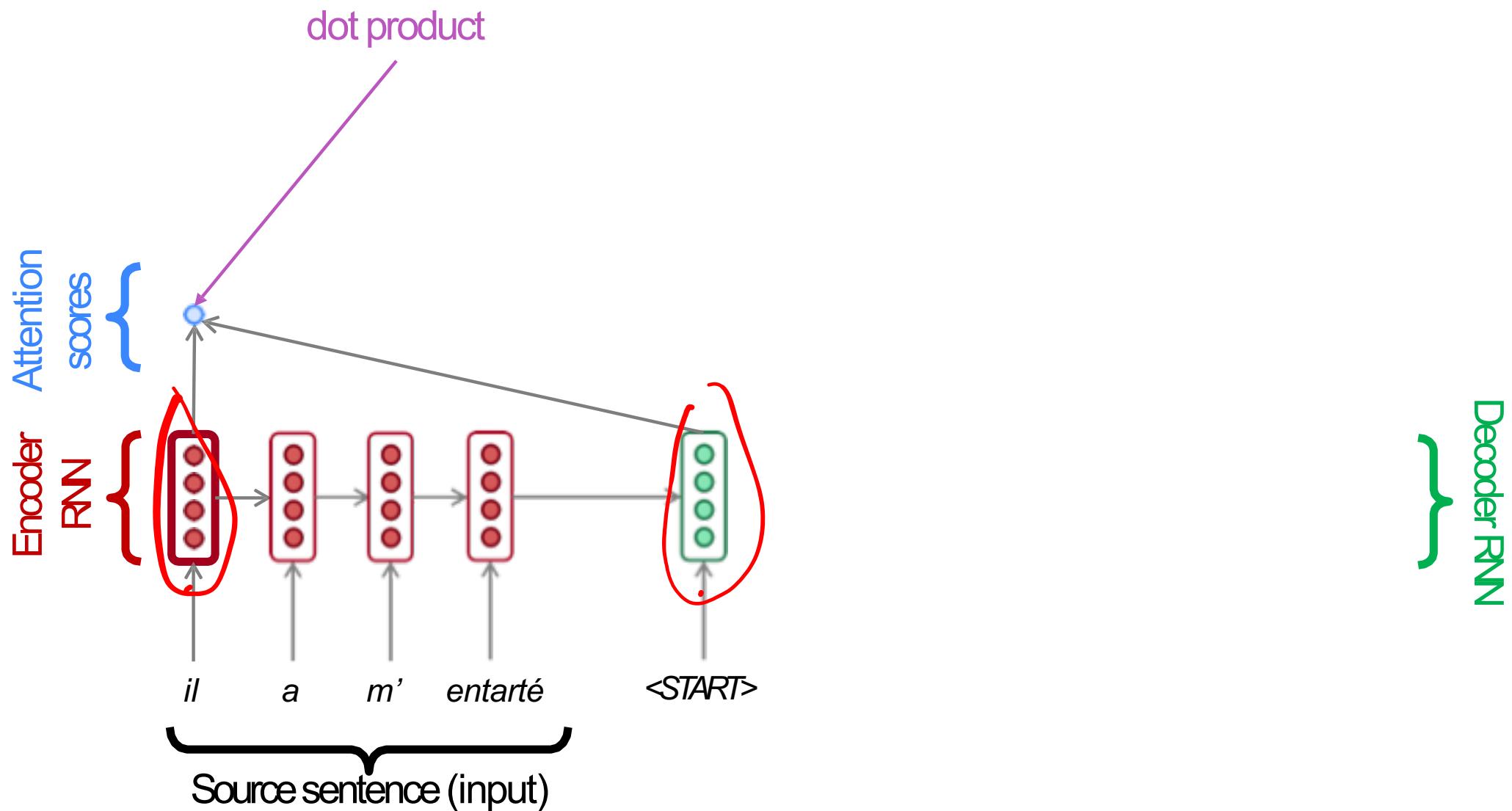
Attention

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

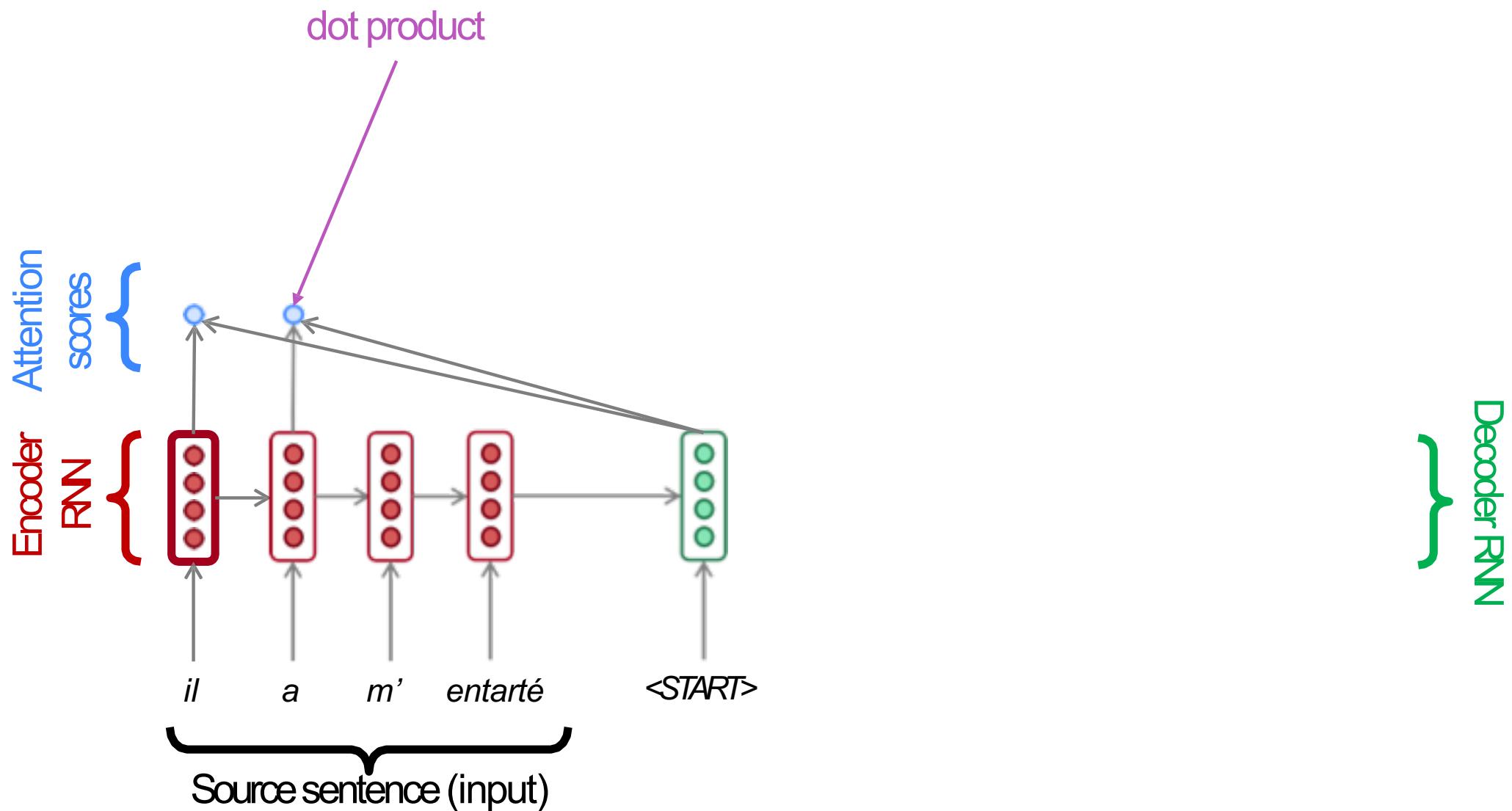


- First we will show via diagram (no equations), then we will show with equations

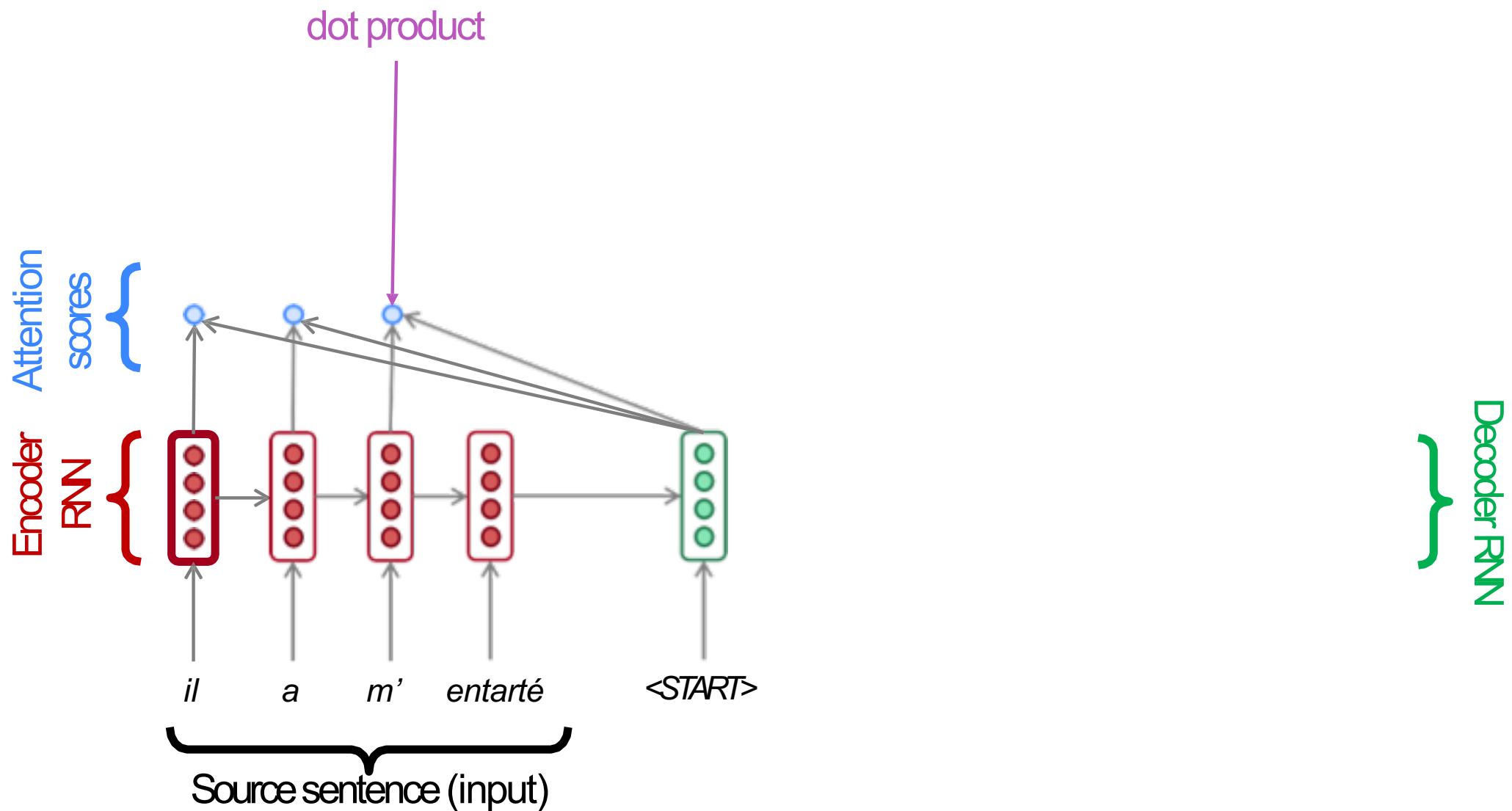
Sequence-to-sequence with attention



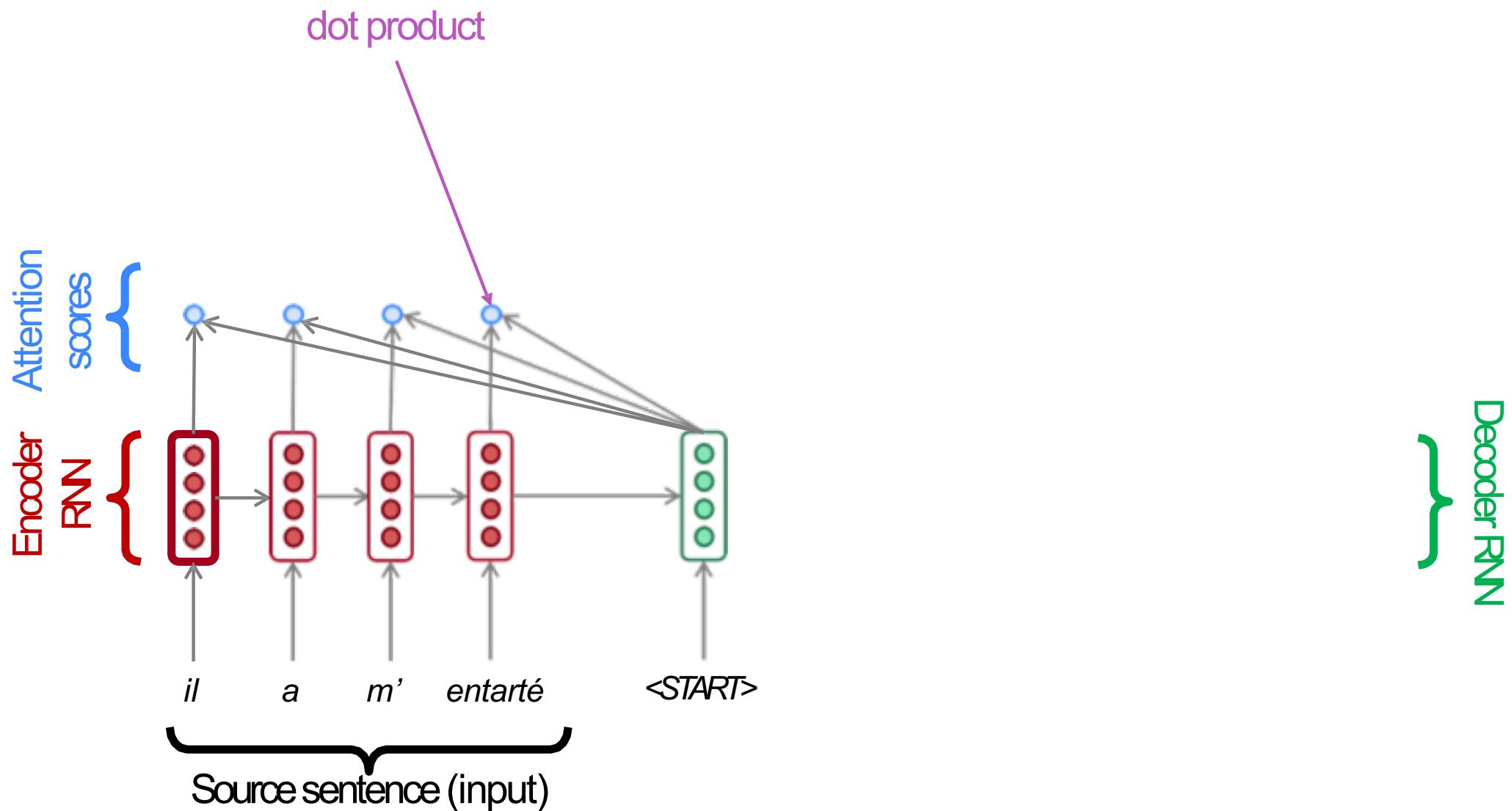
Sequence-to-sequence with attention



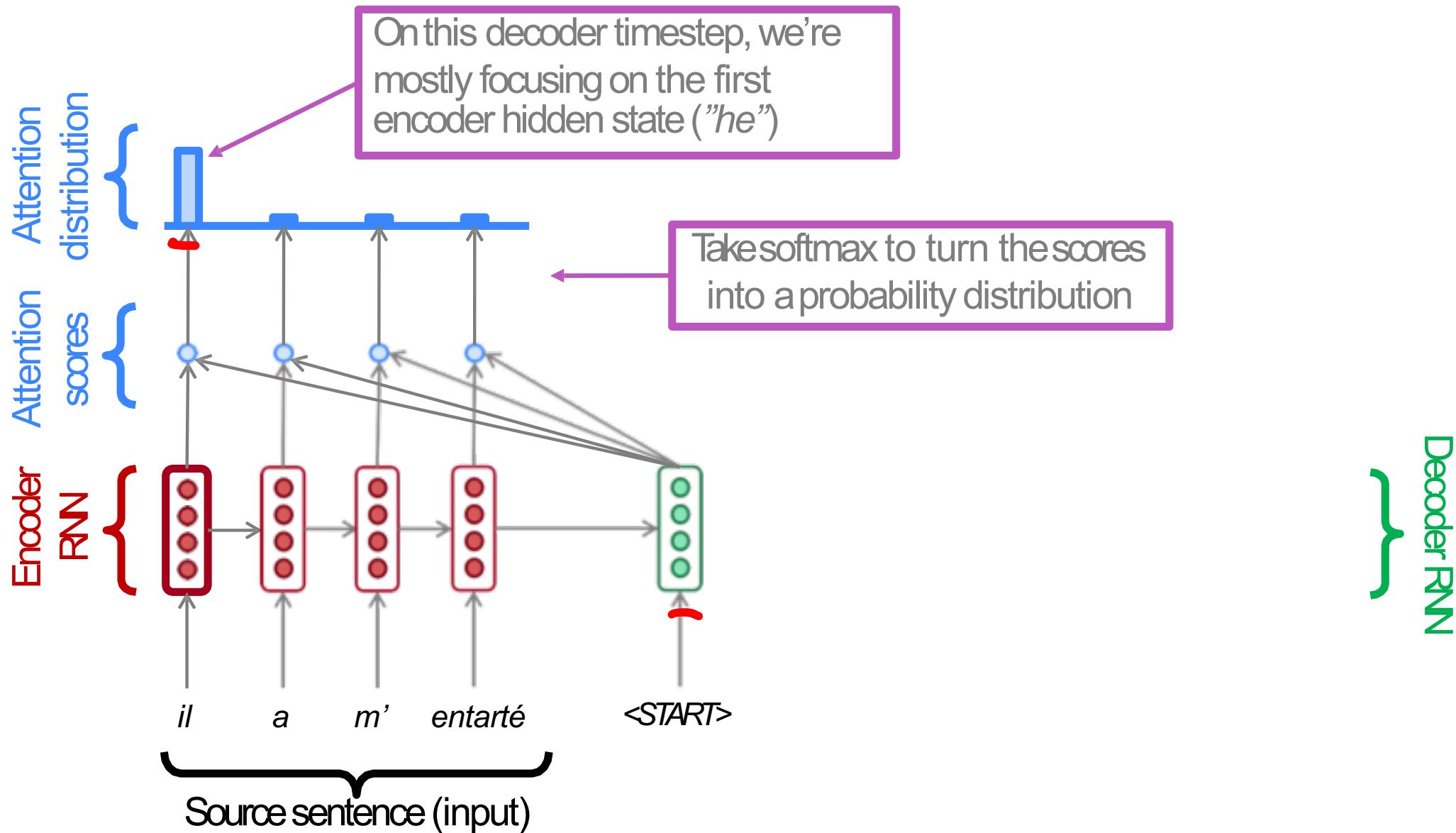
Sequence-to-sequence with attention



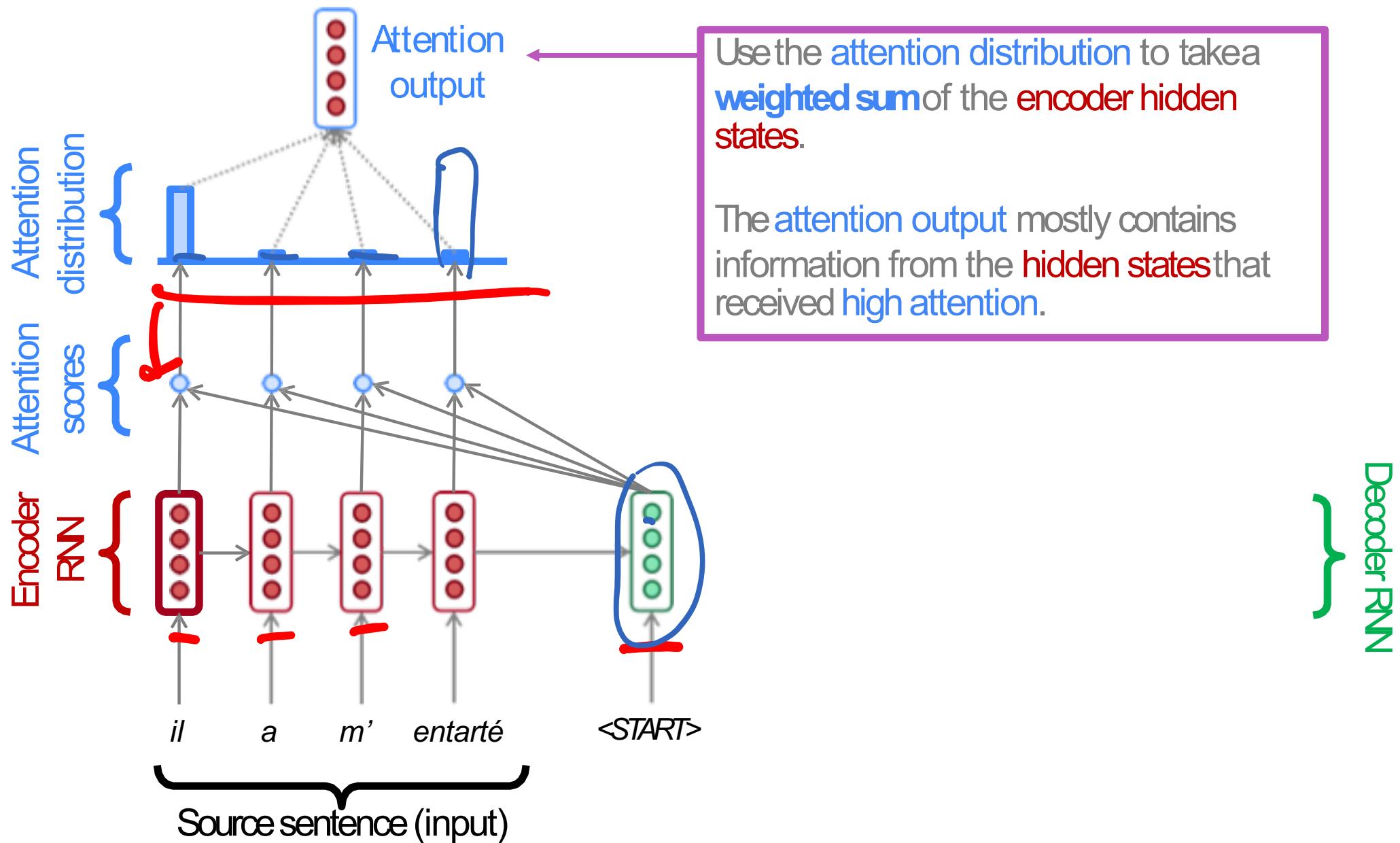
Sequence-to-sequence with attention



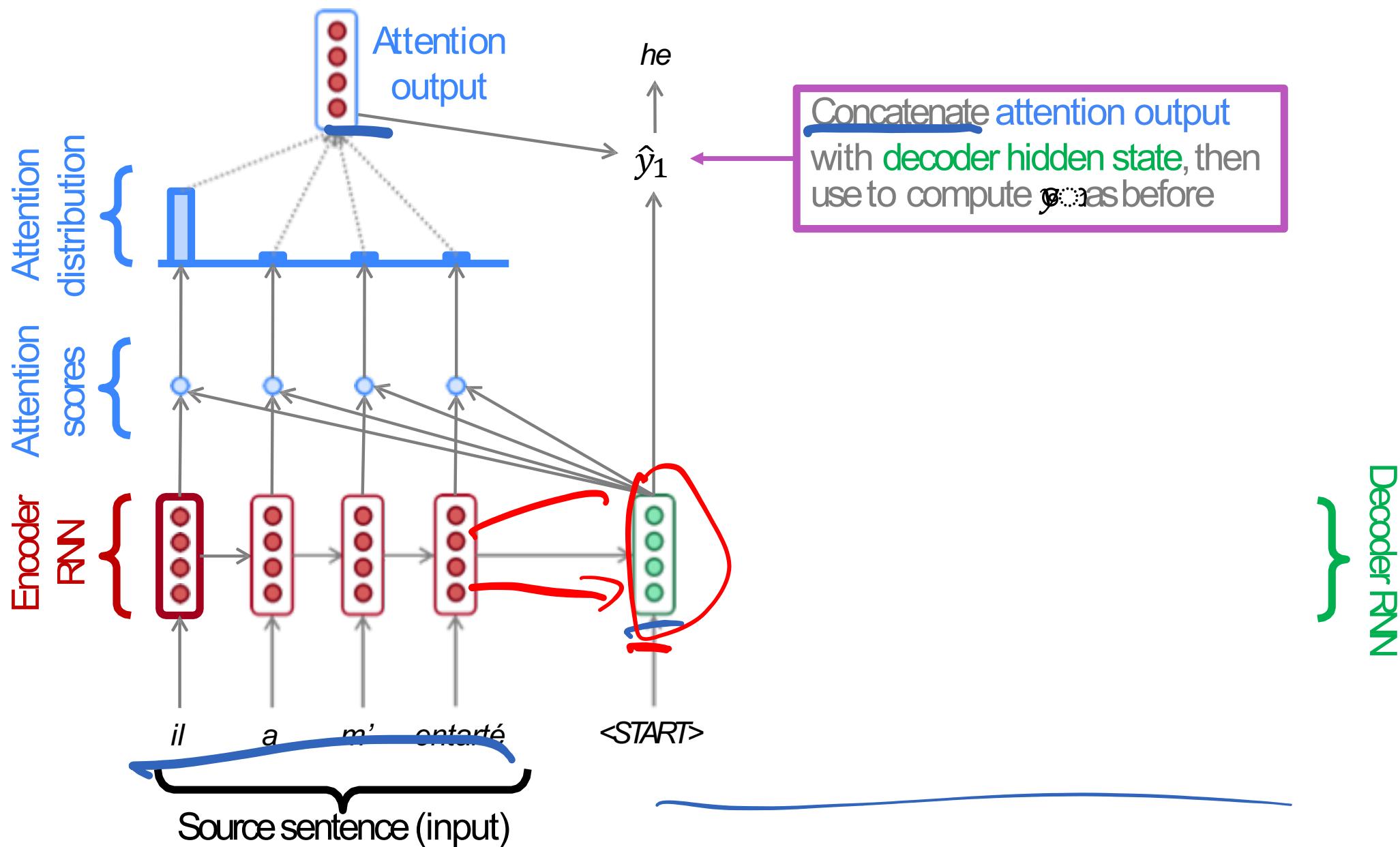
Sequence-to-sequence with attention



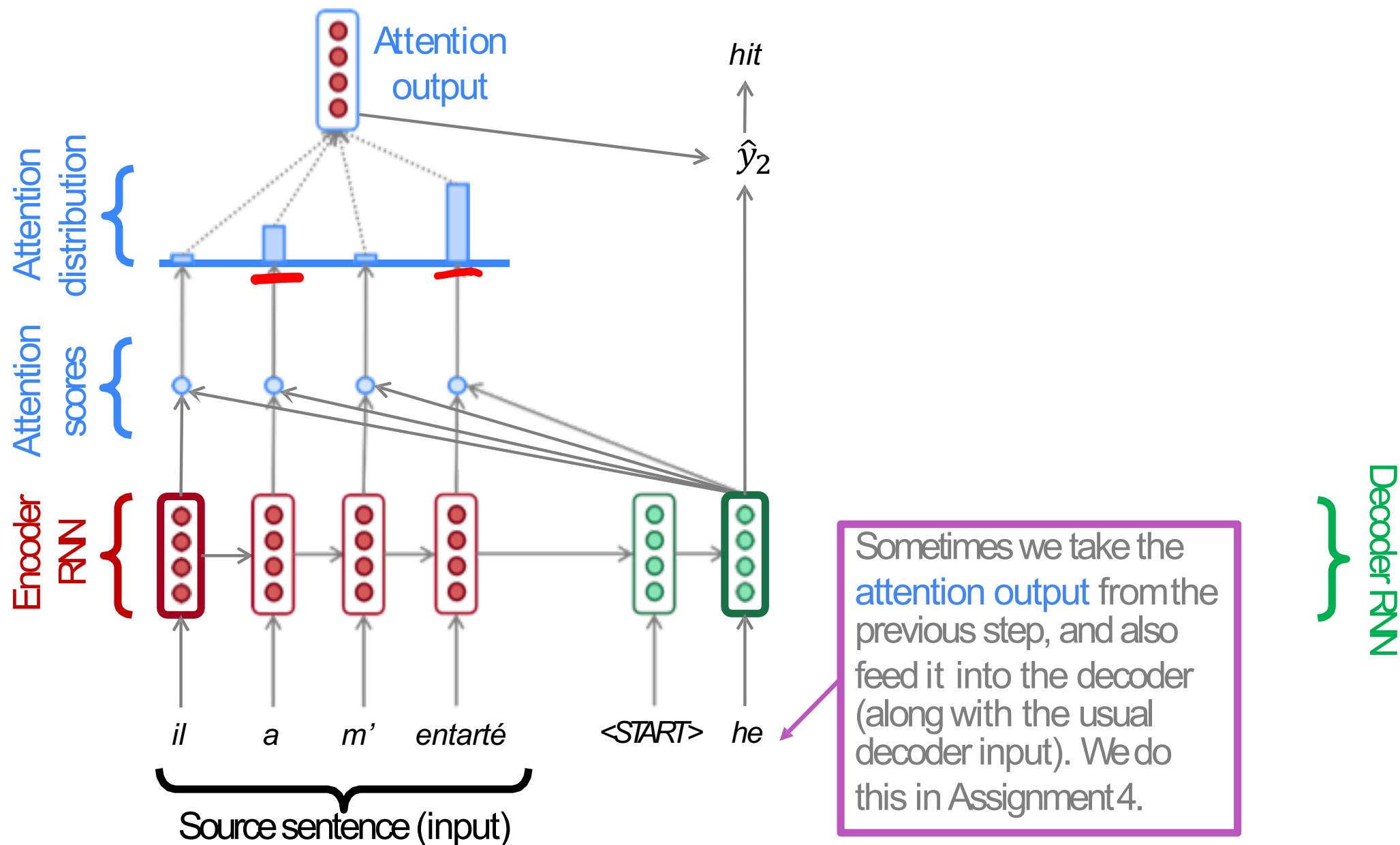
Sequence-to-sequence with attention



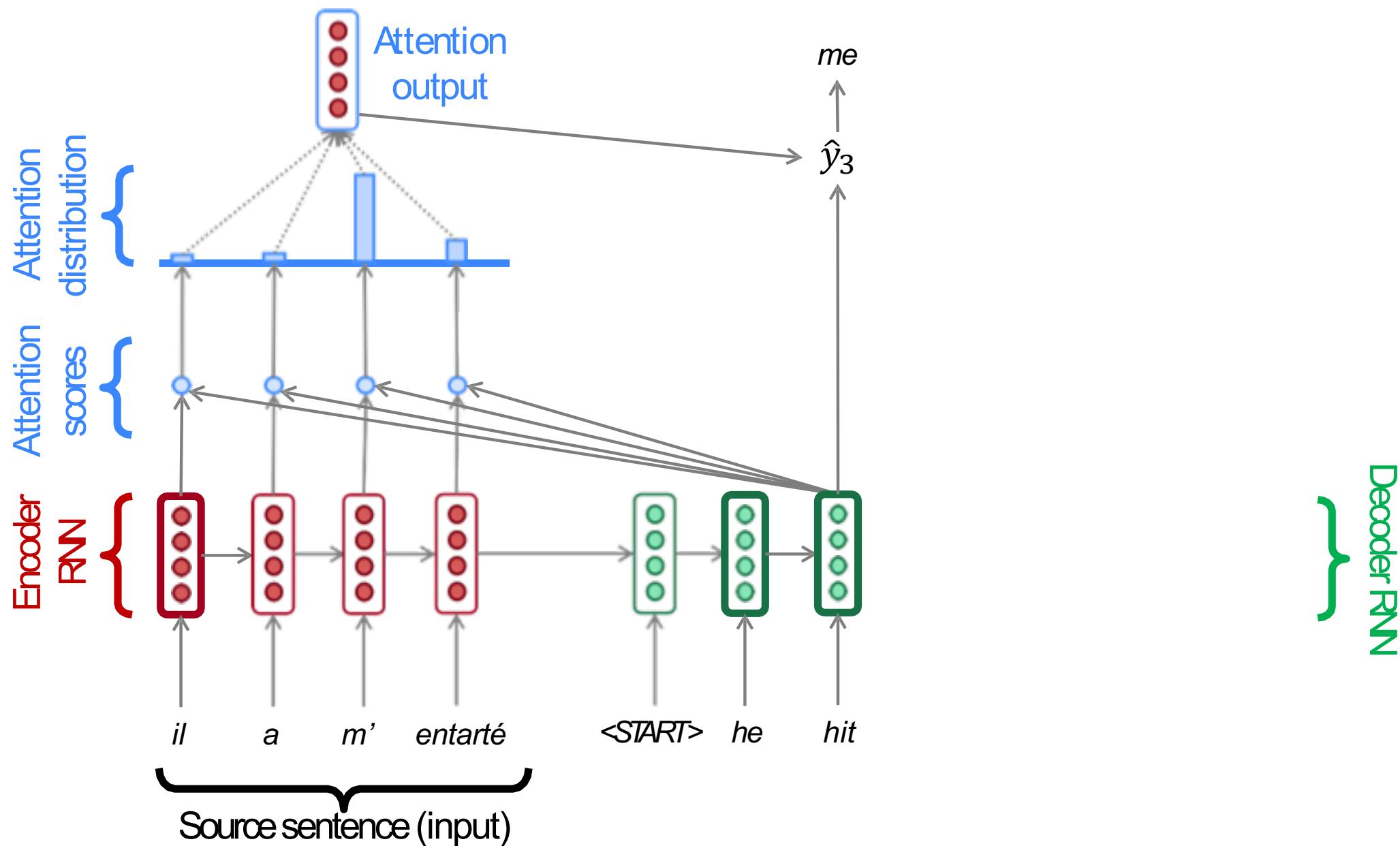
Sequence-to-sequence with attention



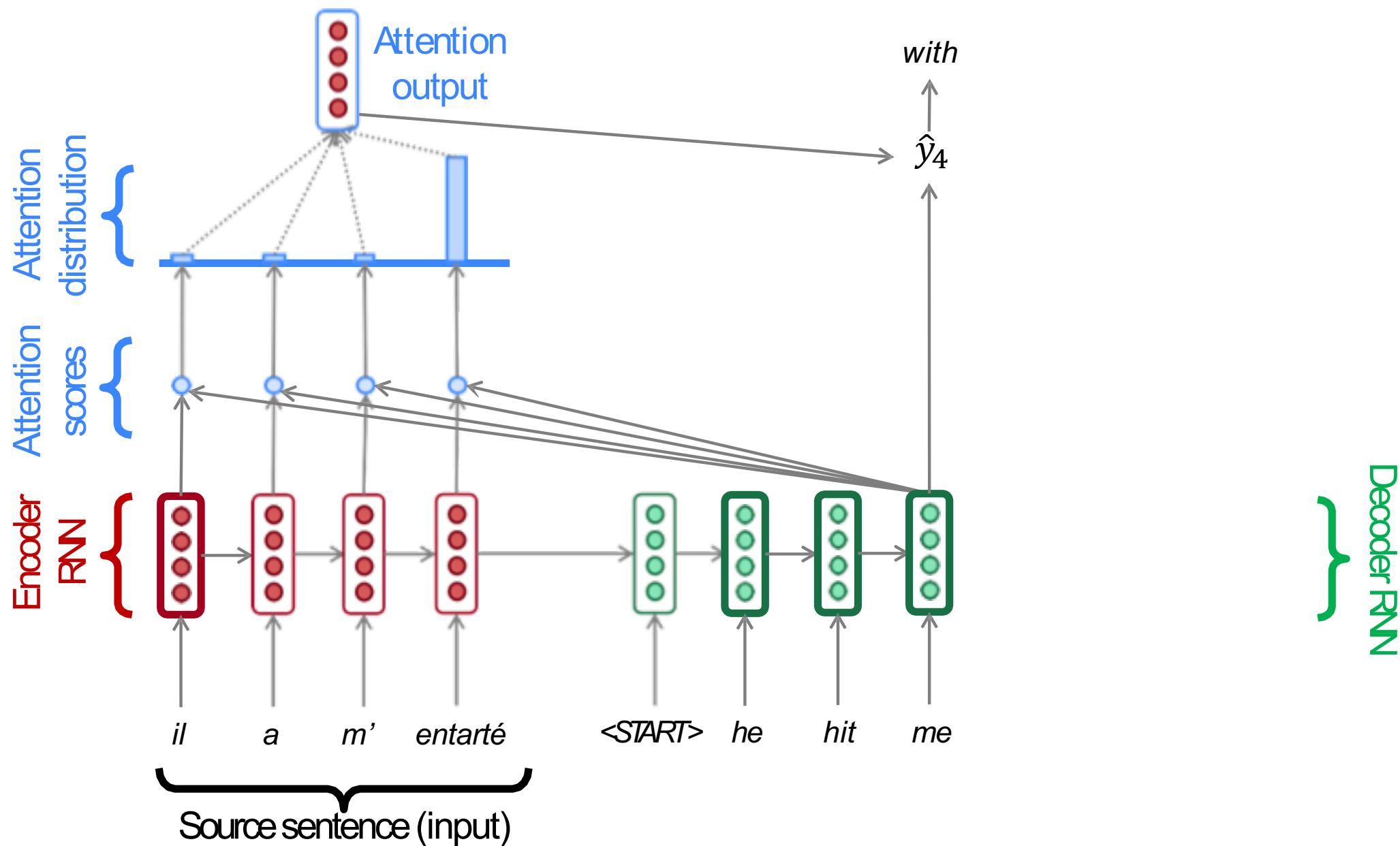
Sequence-to-sequence with attention



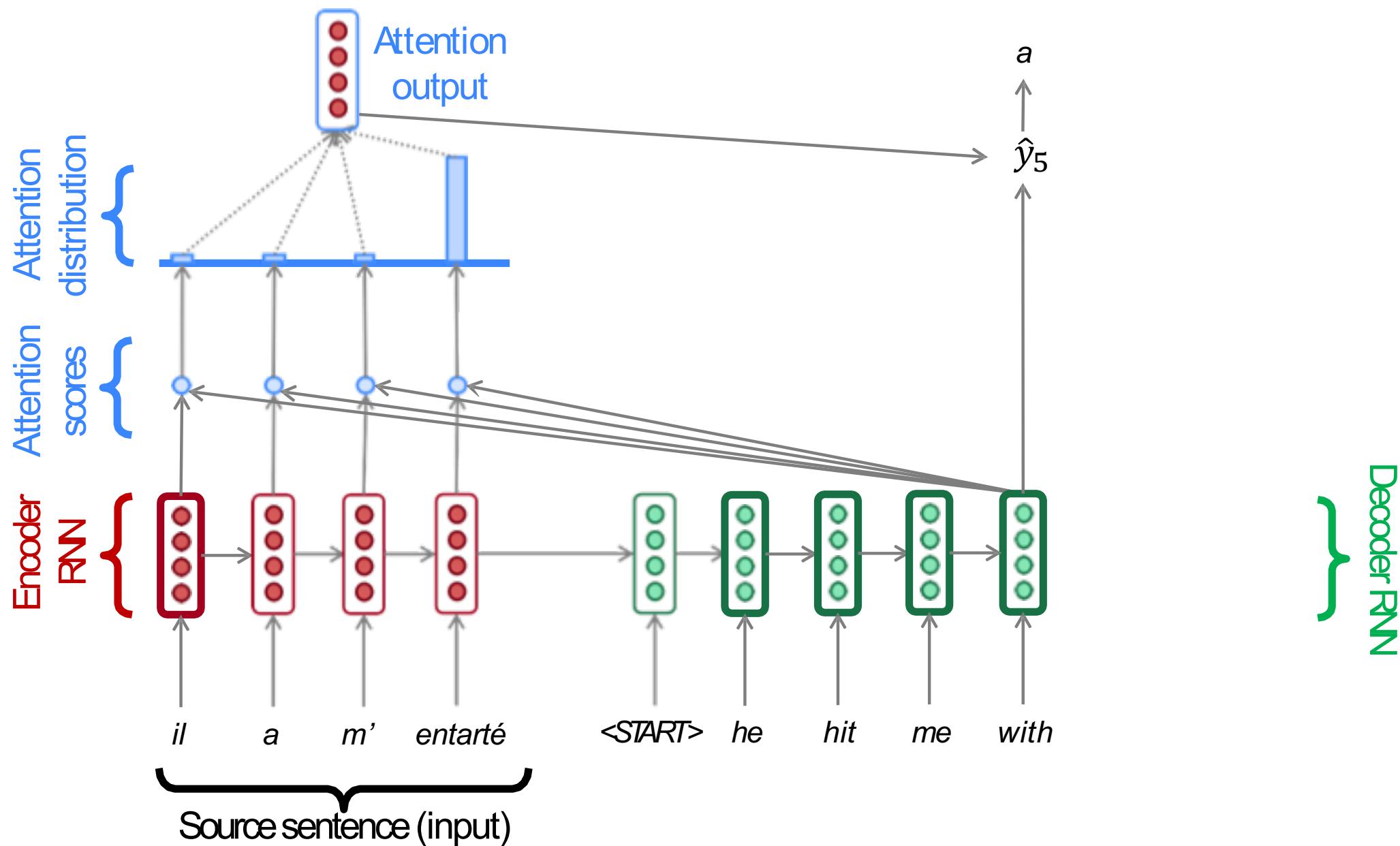
Sequence-to-sequence with attention



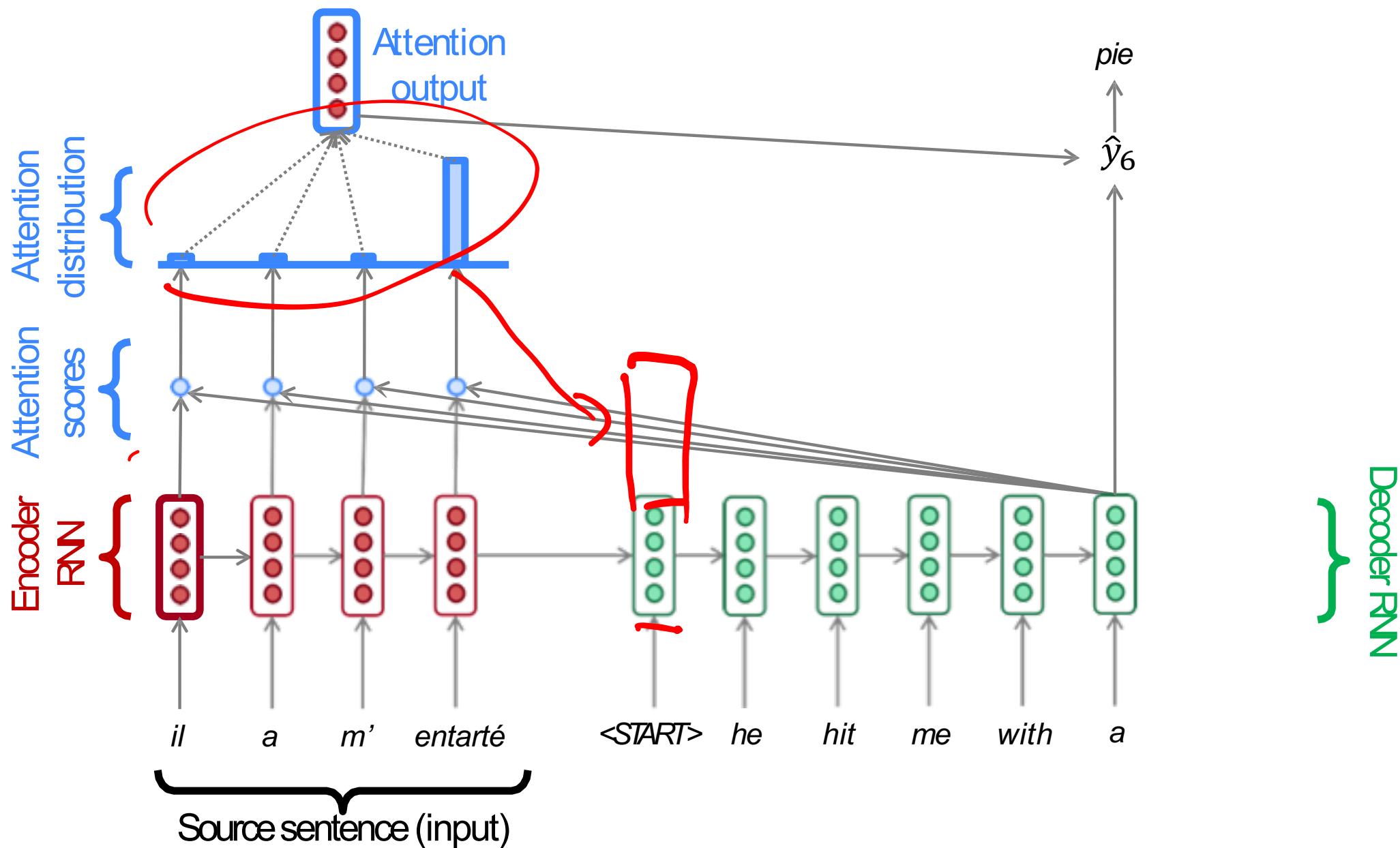
Sequence-to-sequence with attention



Sequence-to-sequence with attention



Sequence-to-sequence with attention



Attention: in equations

- We have encoder hidden states $\underline{h_1, \dots, h_N} \in \mathbb{R}^h$
- On timestep t , we have decoder hidden state $\underline{s_t} \in \mathbb{R}^h$
- We get the attention scores $\underline{e^t}$ for this step:

$$\underline{e^t} = [\underline{s_t^T h_1}, \dots, \underline{s_t^T h_N}] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution $\underline{\alpha^t}$ for this step (this is a probability distribution and sums to 1)

$$\underline{\alpha^t} = \underline{\text{softmax}(e^t)} \in \mathbb{R}^N$$

- We use $\underline{\alpha^t}$ to take a weighted sum of the encoder hidden states to get the attention output $\underline{a_t}$

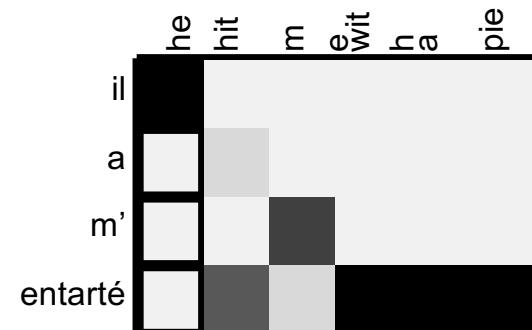
$$\underline{a_t} = \sum_{i=1}^N \underline{\alpha_i^t h_i} \in \mathbb{R}^h$$

- Finally we concatenate the attention output $\underline{a_t}$ with the decoder hidden state $\underline{s_t}$ and proceed as in the non-attention seq2seq model

$$[\underline{a_t; s_t}] \in \mathbb{R}^{2h}$$

Attention is great

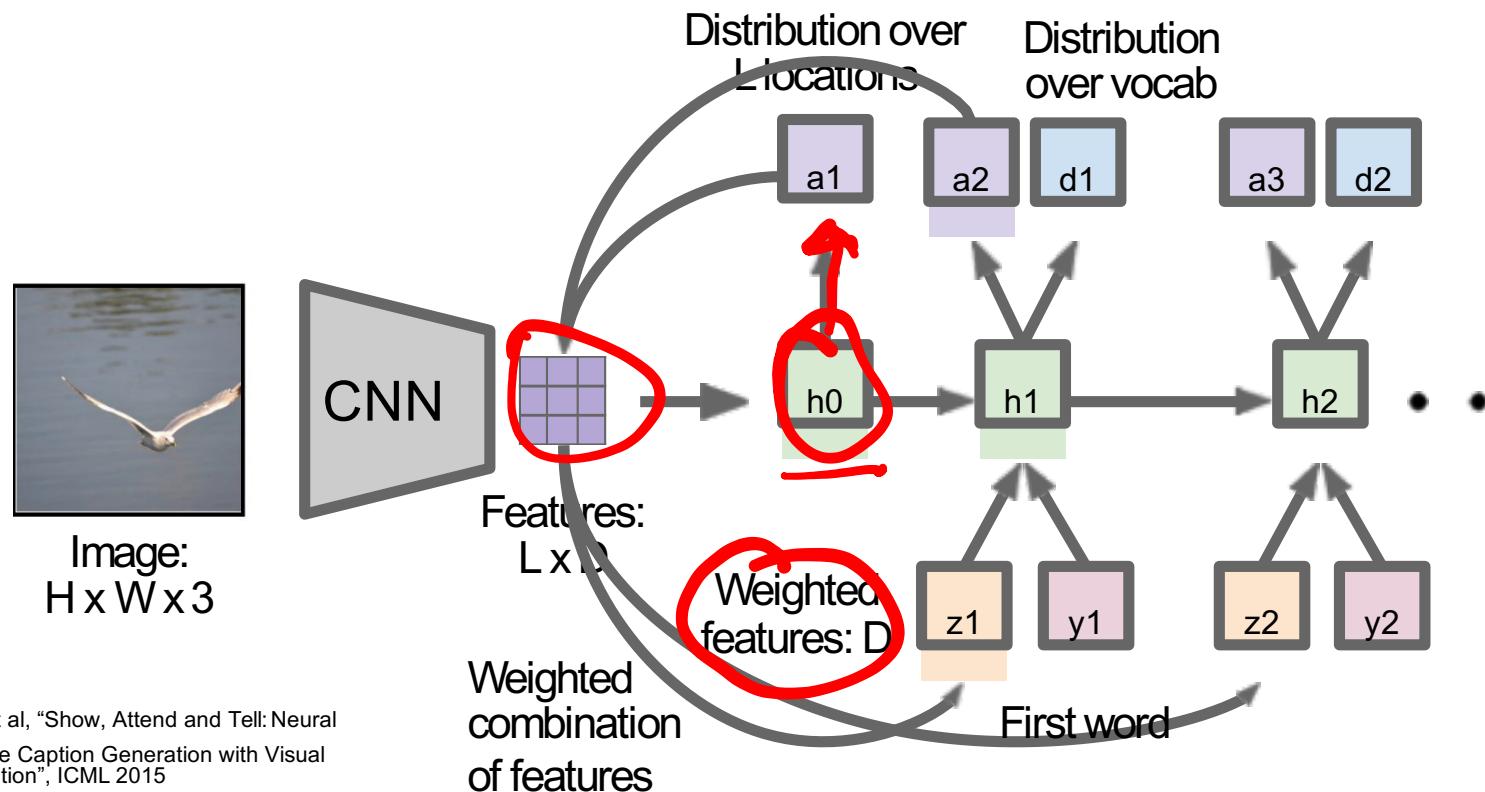
- Attention significantly improves NMT performance
 - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
 - Provides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



Attention is a *general* Deep Learning tech

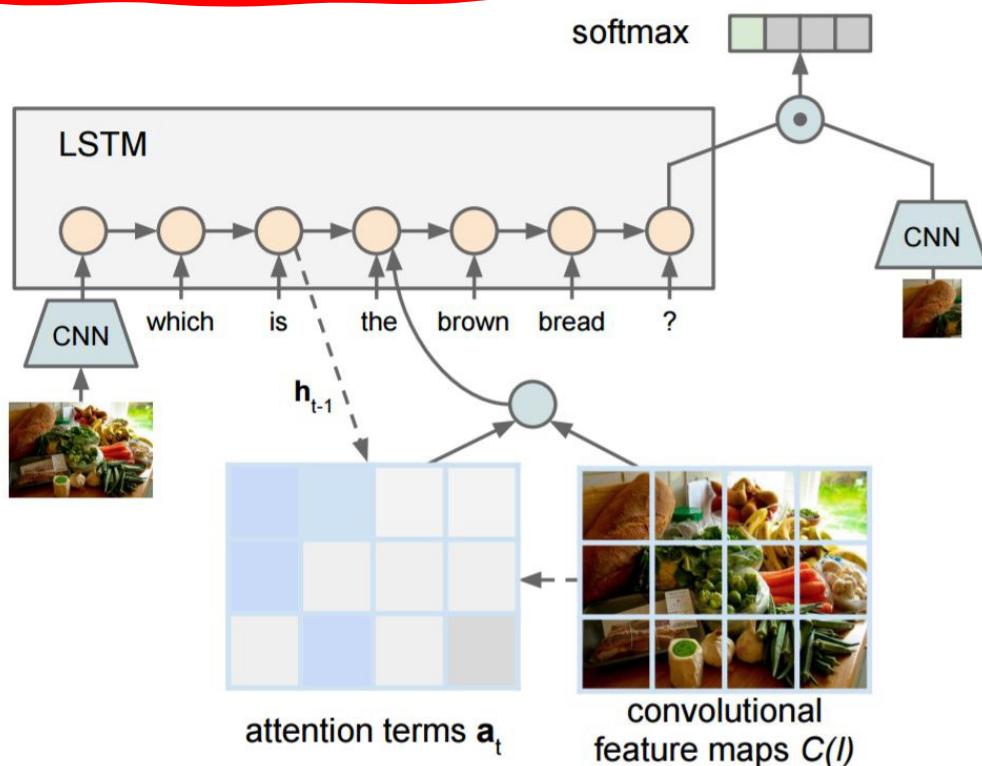
- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.
 - However: You can use attention in **many architectures** (not just seq2seq) and **many tasks** (not just MT)
- **More general definition of attention:**
 - Given a set of vector *values*, 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 (query) *attends to* all the encoder hidden states (values).

Image Captioning with Attention

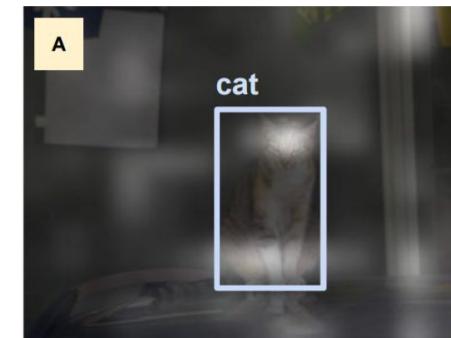


Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Visual Question Answering: RNNs with Attention



Zhu et al., "Visual 7W: Grounded Question Answering in Images", CVPR 2016
Figures from Zhu et al., copyright IEEE 2016. Reproduced for educational purposes.



What kind of animal is in the photo?

A **cat**.



Why is the person holding a knife?

To cut the **cake** with.

There are *several* attention variants

- We have some *values* $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$ and a *query* $\mathbf{s} \in \mathbb{R}^{d_2}$
- Attention always involves:
 1. Computing the *attention scores* $\mathbf{e} \in \mathbb{R}^N$
 2. Taking softmax to get *attention distribution* α :

$$\alpha = \text{softmax}(\mathbf{e}) \in \mathbb{R}^N$$

There are multiple ways to do this

3. Using attention distribution to take weighted sum of values:

$$\mathbf{a} = \sum_{i=1}^N \alpha_i \mathbf{h}_i \in \mathbb{R}^{d_1}$$

thus obtaining the *attention output* \mathbf{a} (sometimes called the *context vector*)

Attention variants

There are **several ways** you can compute $e \in \mathbb{R}^N$ from $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$ and $\mathbf{s} \in \mathbb{R}^{d_2}$:

- **Basic dot-product attention:** $e_i = \mathbf{s}^T \mathbf{h}_i \in \mathbb{R}$
 - Note: this assumes $d_1 = d_2$
 - This is the version we saw earlier
- **Multiplicative attention:** $e_i = \mathbf{s}^T \mathbf{W} \mathbf{h}_i \in \mathbb{R}$
 - Where $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$ is a weight matrix
- **Additive attention:** $e_i = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}) \in \mathbb{R}$
 - Where $\mathbf{W}_1 \in \mathbb{R}^{d_3 \times d_1}$, $\mathbf{W}_2 \in \mathbb{R}^{d_3 \times d_2}$ are weight matrices and $\mathbf{v} \in \mathbb{R}^{d_3}$ is a weight vector.
 - d_3 (the attention dimensionality) is a hyperparameter

More information:

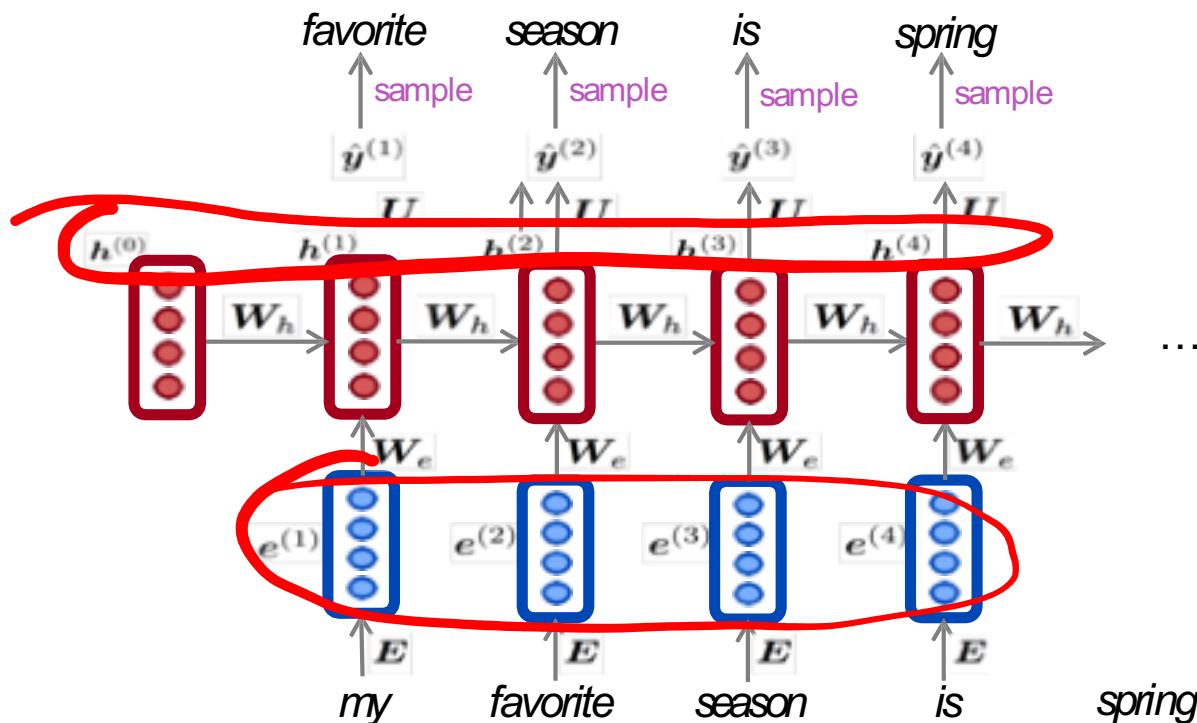
“Deep Learning for NLP Best Practices”, Ruder, 2017, <http://ruder.io/deep-learning-nlp-best-practices/index.html#attention>
“Massive Exploration of Neural Machine Translation Architectures”, Britz et al, 2017, <https://arxiv.org/pdf/1703.03906.pdf>

Contextual Word Representations

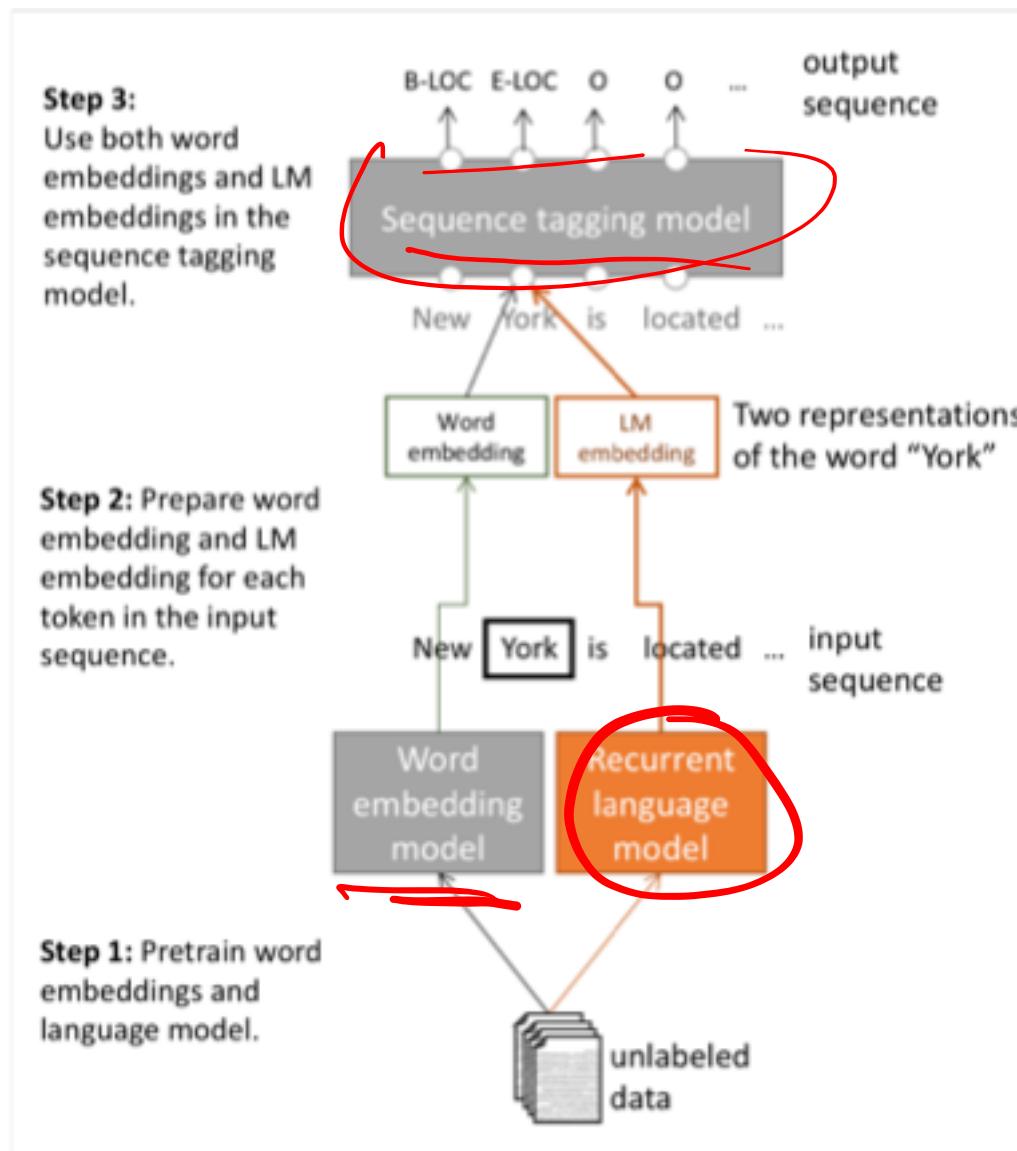
- Problem with word2vec?

Did we all along have a solution to this problem?

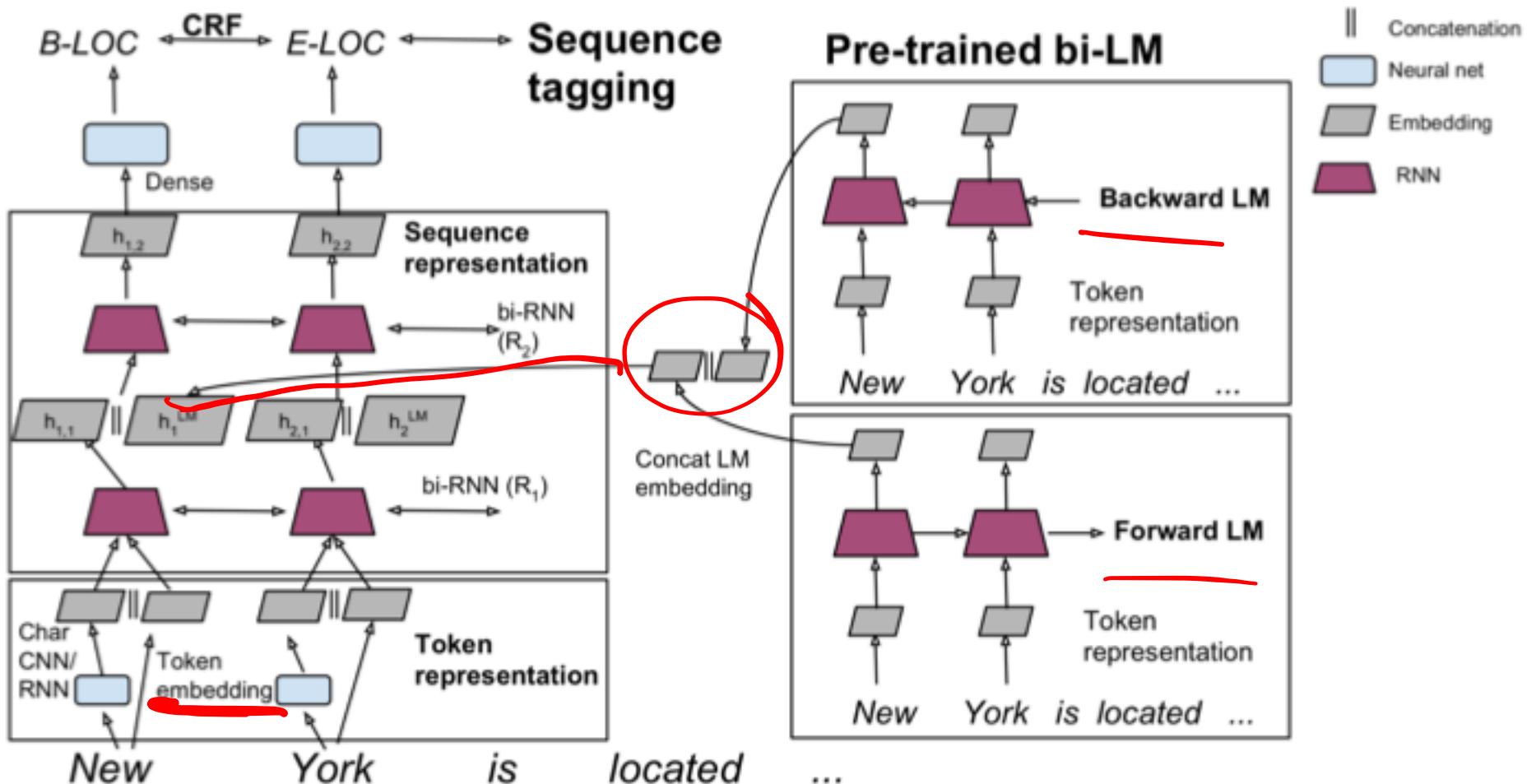
- In, an NLM, we immediately stuck word vectors (perhaps only trained on the corpus) through LSTM layers
- Those LSTM layers are trained to predict the next word
- But those language models are producing context-specific word representations at each position!



TagLM [Peters et al. 2017]



Tag LM



$$\mathbf{h}_{k,1} = [\overrightarrow{\mathbf{h}}_{k,1}; \overleftarrow{\mathbf{h}}_{k,1}; \mathbf{h}_k^{LM}].$$

Named Entity Recognition (NER)

- A very important NLP sub-task: find and classify names in text, for example:
 - The decision by the independent MP Andrew Wilkie to withdraw his support for the minority Labor government sounded dramatic but it should not further threaten its stability. When, after the 2010 election, Wilkie, Rob Oakeshott, Tony Windsor and the Greens agreed to support Labor, they gave just two guarantees: confidence and supply.

Person
Date
Location
Organization

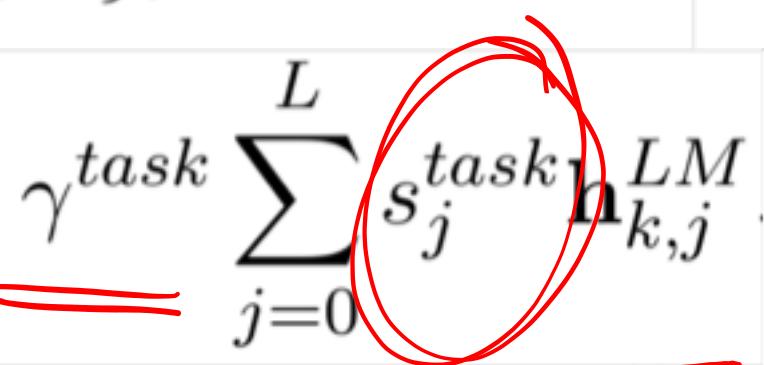
CoNLL2003 Named Entity Recognition

Name	Description	Year	F1
TagLM Peters	LSTM BiLM in BiLSTMtagger	2017	91.93
Ma + Hovy	BiLSTM + char CNN + CRFlayer	2016	91.21
Tagger Peters	BiLSTM + char CNN + CRFlayer	2017	90.87
Ratinov + Roth	Categorical CRF+Wikipedia+word cls	2009	90.80
Finkel et al.	Categorical feature CRF	2005	86.86
IBM Florian	Linear/softmax/TBL/HMM ensemble, gazettes++	2003	88.76
Stanford Klein	MEMM softmax markov model	2003	86.07

ELMo [Peters et al. 2018]

- ELMo learns task-specific combination of biLM representations
- This is an innovation that improves on just using top layer of LSTM stack

$$\begin{aligned} R_k &= \{\mathbf{x}_k^{LM}, \overrightarrow{\mathbf{h}}_{k,j}^{LM}, \overleftarrow{\mathbf{h}}_{k,j}^{LM} \mid j = 1, \dots, L\} \\ &= \{\mathbf{h}_{k,j}^{LM} \mid j = 0, \dots, L\}, \end{aligned}$$

$$\text{ELMo}_k^{\text{task}} = E(R_k; \Theta^{\text{task}}) = \gamma^{\text{task}} \sum_{j=0}^L s_j^{\text{task}} \mathbf{h}_{k,j}^{LM}$$


CoNLL2003 Named Entity Recognition

Name	Description	Year	F1
ELMo	ELMo in BiLSTM	2018	<u>92.22</u>
TagLM Peters	LSTM BiLM in BiLSTMtagger	2017	91.93
Ma + Hovy	BiLSTM + char CNN + CRFlayer	2016	91.21
Tagger Peters	BiLSTM + char CNN + CRFlayer	2017	90.87
Ratinov + Roth	Categorical CRF+Wikipedia+word cls	2009	90.80
Finkel et al.	Categorical feature CRF	2005	86.86
IBM Florian	Linear/softmax/TBL/HMM ensemble, gazettes++	2003	88.76
S ² tanford	MEMM softmax markov model	2003	86.07

ELMo results: Great for all tasks

TASK	PREVIOUS SOTA	OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17
SRL	He et al. (2017)	81.7	81.4	84.6
Coref	Lee et al. (2017)	67.2	67.2	70.4
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5

ELMo: Weighting of layers

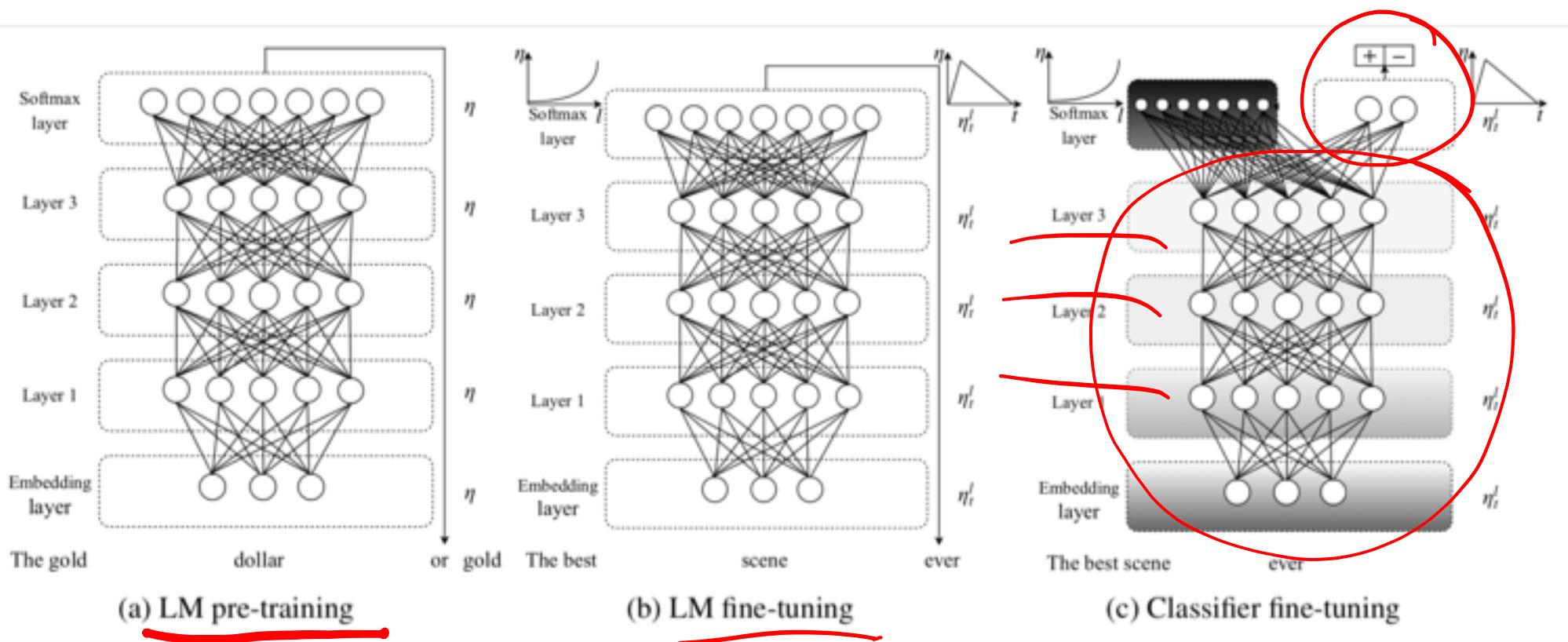
- The two biLSTM NLM layers have differentiated uses/meanings
 - Lower layer is better for lower-level syntax, etc.
 - Part-of-speech tagging, syntactic dependencies, NER
 - Higher layer is better for higher-level semantics
 - Sentiment, Semantic role labeling, question answering, SNLI
- This seems interesting, but it'd seem more interesting to see how it pans out with more than two layers of network

ULMfit [Howard and Ruder 2018]

Train LM on big general domain corpus (use biLM)

Tune LM on target task data

Fine-tune as classifier on target task

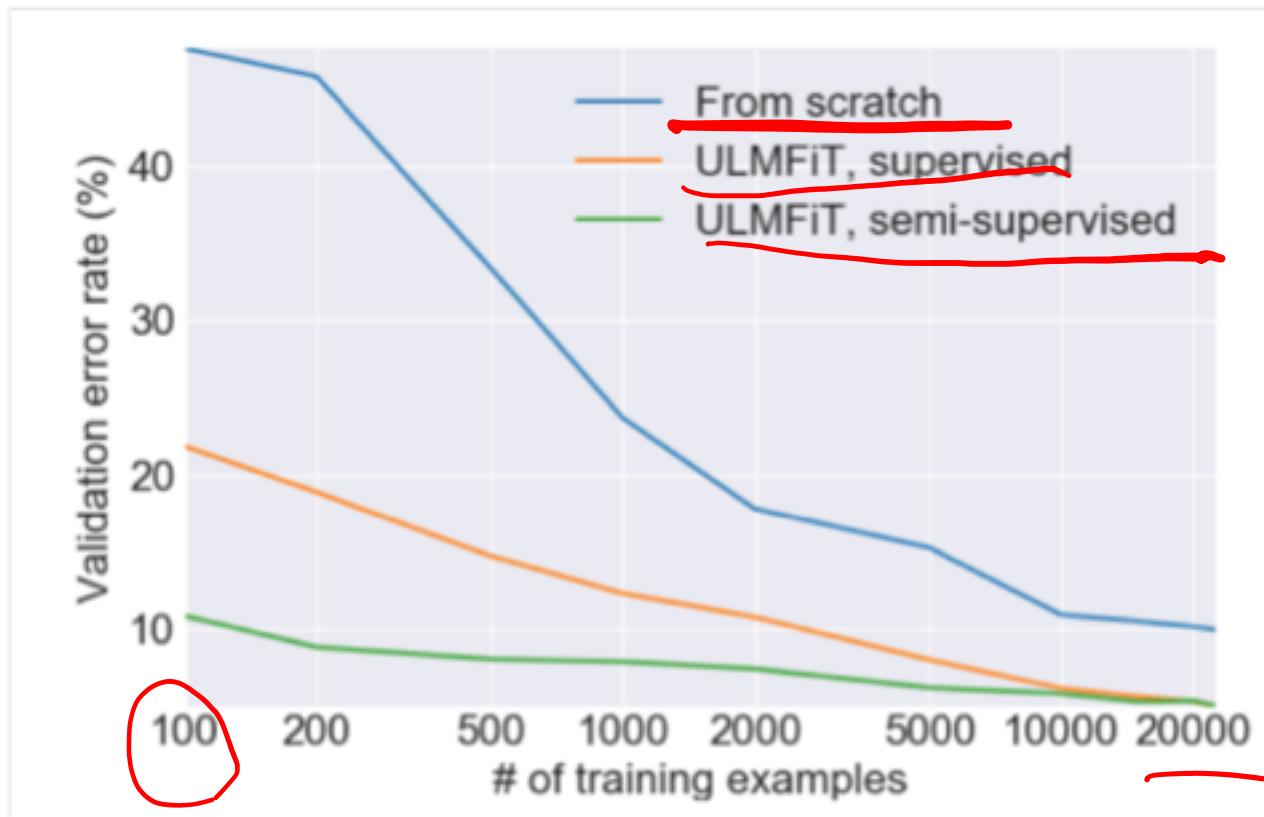


ULMfit performance

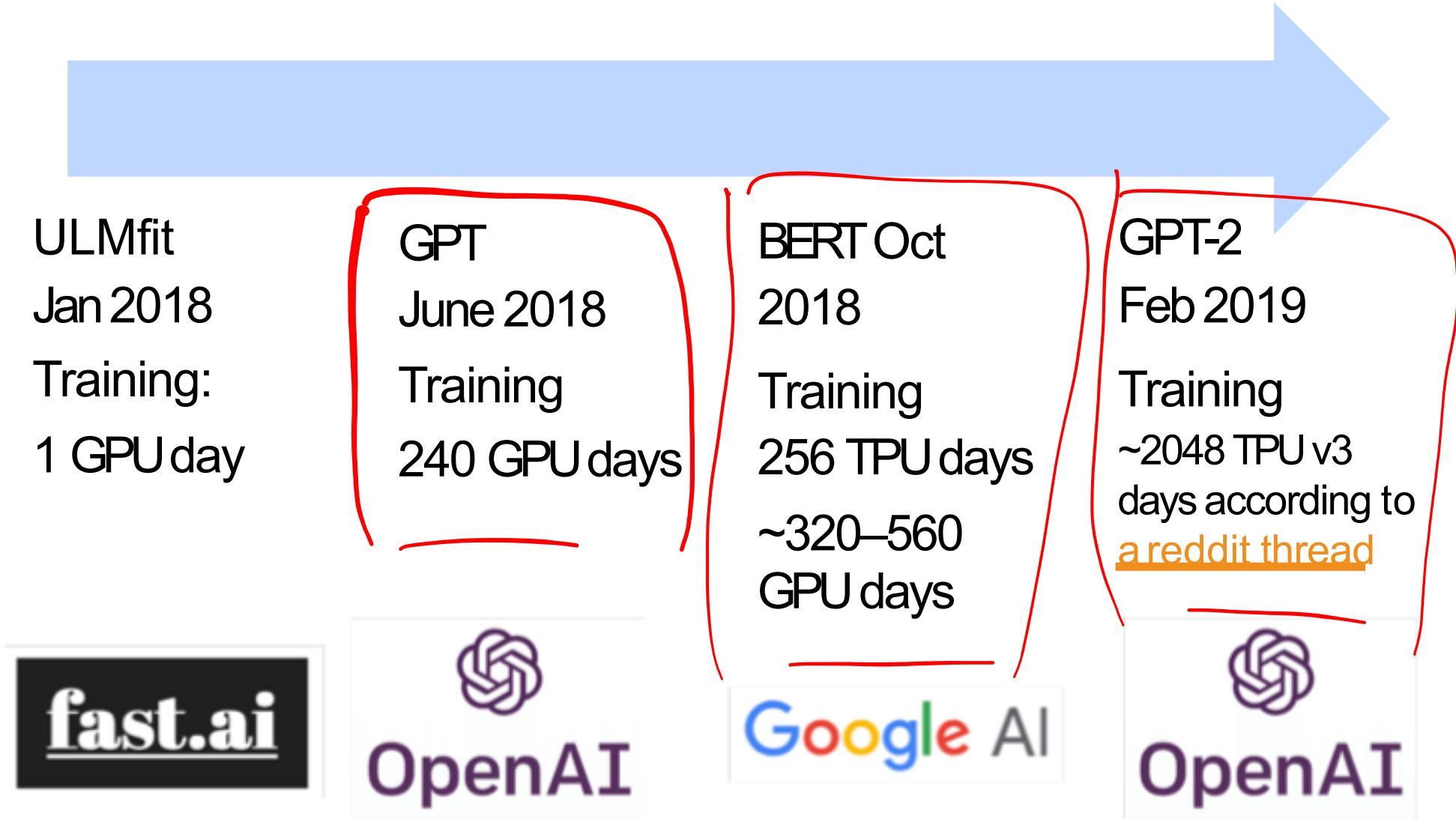
- Text classifier error rates

	Model	Test	Model	Test
IMDb	CoVe (McCann et al., 2017)	8.2	CoVe (McCann et al., 2017)	4.2
	oh-LSTM (Johnson and Zhang, 2016)	5.9	TBCNN (Mou et al., 2015)	4.0
	Virtual (Miyato et al., 2016)	5.9	LSTM-CNN (Zhou et al., 2016)	3.9
	ULMFiT (ours)	4.6	ULMFiT (ours)	3.6

ULMfit transfer learning



Let's scale it up!



GPT-2 Reaction

Due to concerns about large language models being used to generate deceptive, biased, or abusive language at scale, we are only releasing a much smaller version of GPT-2 along with sampling code. We are not releasing the dataset, training code, or GPT-2 model weights. Nearly a year ago we wrote in the OpenAI Charter:

GPT-2 Reaction



Elon Musk-founded OpenAI builds artificial intelligence so powerful it must be kept locked up for the good of humanity



Jasper Hamill Friday 15 Feb 2019 10:06 am

Machine-generated text is about to break the internet



Mark Rickerby | Guest writer

OpenAI built a text generator so good, it's considered too dangerous to release

Zack Whittaker @zackwhittaker 3 weeks ago

 Comment

Transformer models

All of these models are Transformer architecture models ...so maybe we had better learn about Transformers?

ULMfit
Jan 2018
Training:
1 GPU day

GPT
June 2018
Training
240 GPU days

BERT Oct
2018
Training
256 TPU days
~320–560
GPU days

GPT-2
Feb 2019
Training
~2048 TPU v3
days according to
[a reddit thread](#)



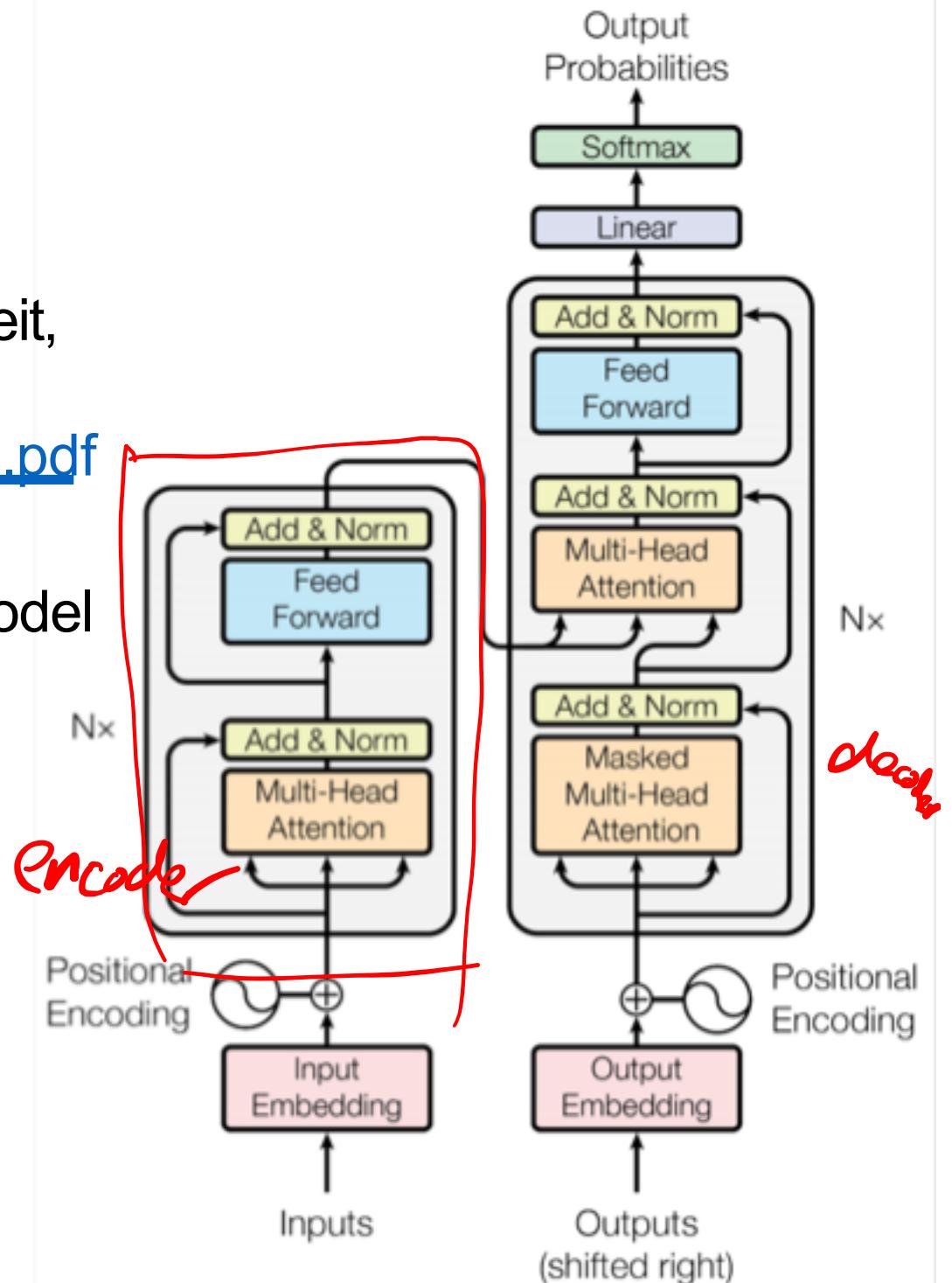
Transformer Overview

Attention is all you need. 2017.

Aswani, Shazeer, Parmar, Uszkoreit,
Jones, Gomez, Kaiser, Polosukhin

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

- Non-recurrent sequence-to-sequence encoder-decoder model
- 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

- Learning about transformers on your own?
 - Key recommended resource:
 - <http://nlp.seas.harvard.edu/2018/04/03/attention.html> 
 - The Annotated Transformer by Sasha Rush
 - An Jupyter Notebook using PyTorch that explains everything!
- For now: Let's define the basic building blocks of transformer networks: first, new attention layers!

Dot-Product Attention

- 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 d_k value have d_v

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

Dot-Product Attention

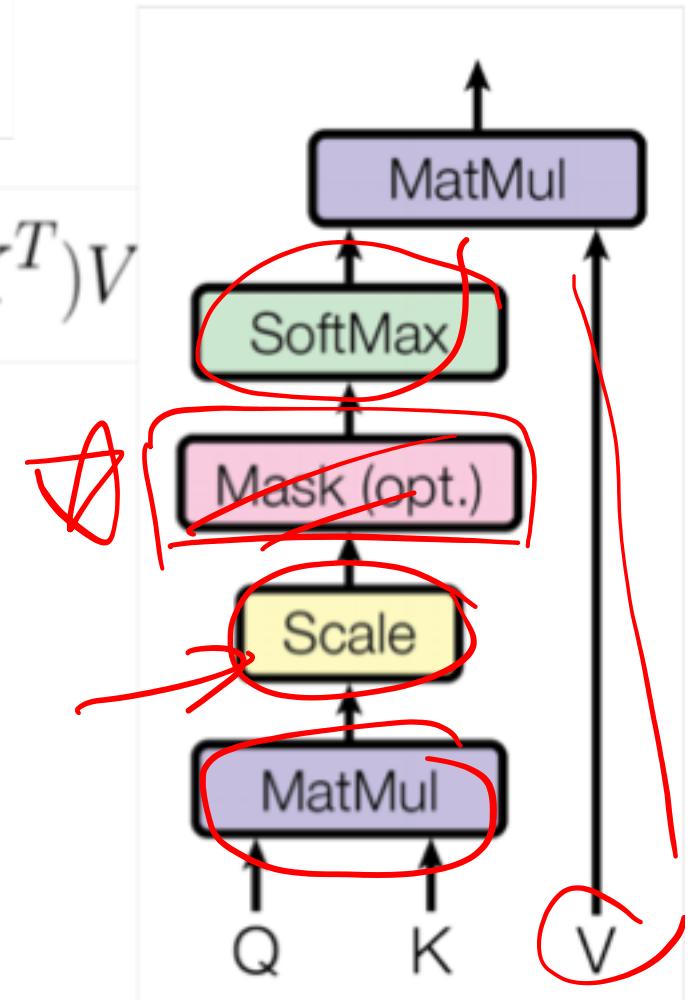
- When we have multiple queries q , we stack them in a matrix Q :

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

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

- Scale by length of query/key vectors:

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



Self-attention in the encoder

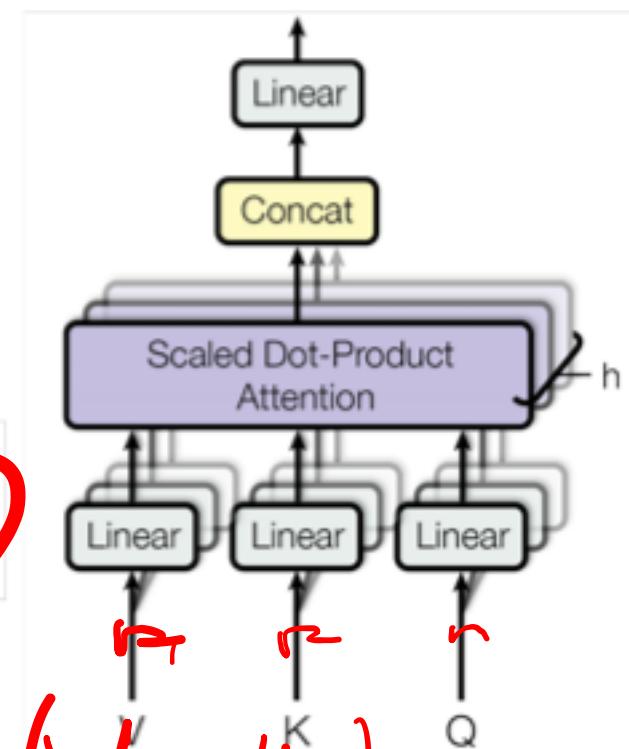
- The input word vectors are the queries, keys and values
- In other words: the word vectors themselves select each other
- Word vector stack = $\underline{Q = K = V}$
- We'll see in the decoder why we separate them in the definition

Multi-head attention

- Problem with simple self-attention:
- Only one way for words to interact with one-another
- Solution: Multi-head attention
- First map Q, K, V into h=8 many lower dimensional spaces via W matrices
- Then apply attention, then concatenate outputs and pipe through linear layer

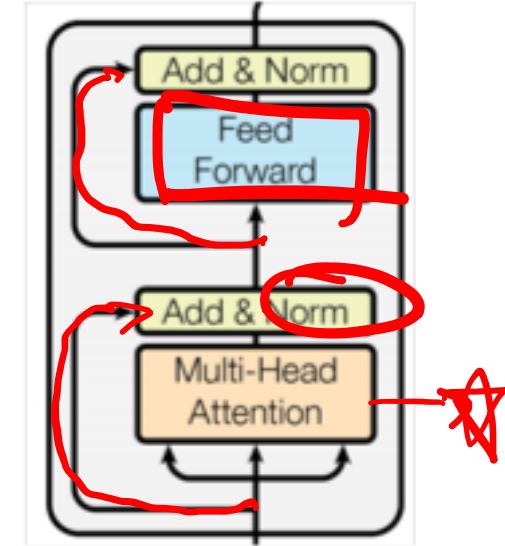
MultiHead(Q, K, V) = Concat(head₁, ..., head_h) W^O
where head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)

$$\sum \text{softmax}(QW_i^Q \cdot (KW_i^K)^T)_{ij} (VW_i^V)_j$$



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

LayerNorm(x + Sublayer(x))

LayerNorm changes input to have mean 0 and variance 1, per layer and per training point (and adds two more parameters)

$$\mu^l = \frac{1}{H} \sum_{i=1}^H a_i^l \quad \sigma^l = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i^l - \mu^l)^2}$$

$$h_i = f\left(\frac{g_i}{\sigma_i} (a_i - \mu_i) + b_i\right)$$

Layer Normalization by Ba, Kiros and Hinton, <https://arxiv.org/pdf/1607.06450.pdf>

Encoder Input

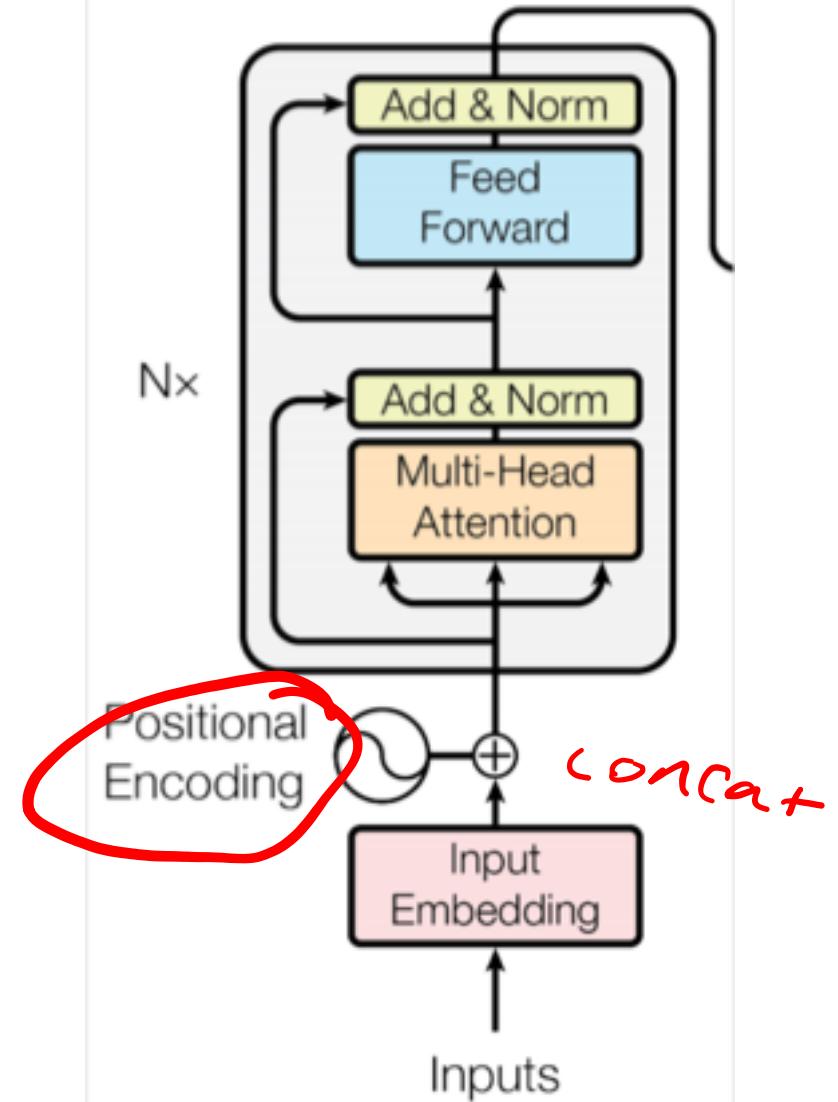
- Actual word representations are byte-pair encodings
- Also added is a **positional encoding** so same words at different locations have different overall representations:

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

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

Complete Encoder

- For encoder, at each block, we use the same Q, K and V from the previous layer
- Blocks are repeated 6 times
 - (in vertical stack)

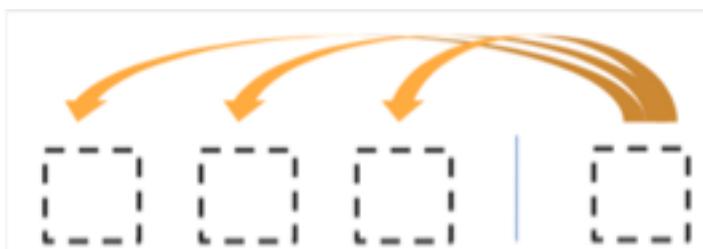


Transformer Decoder

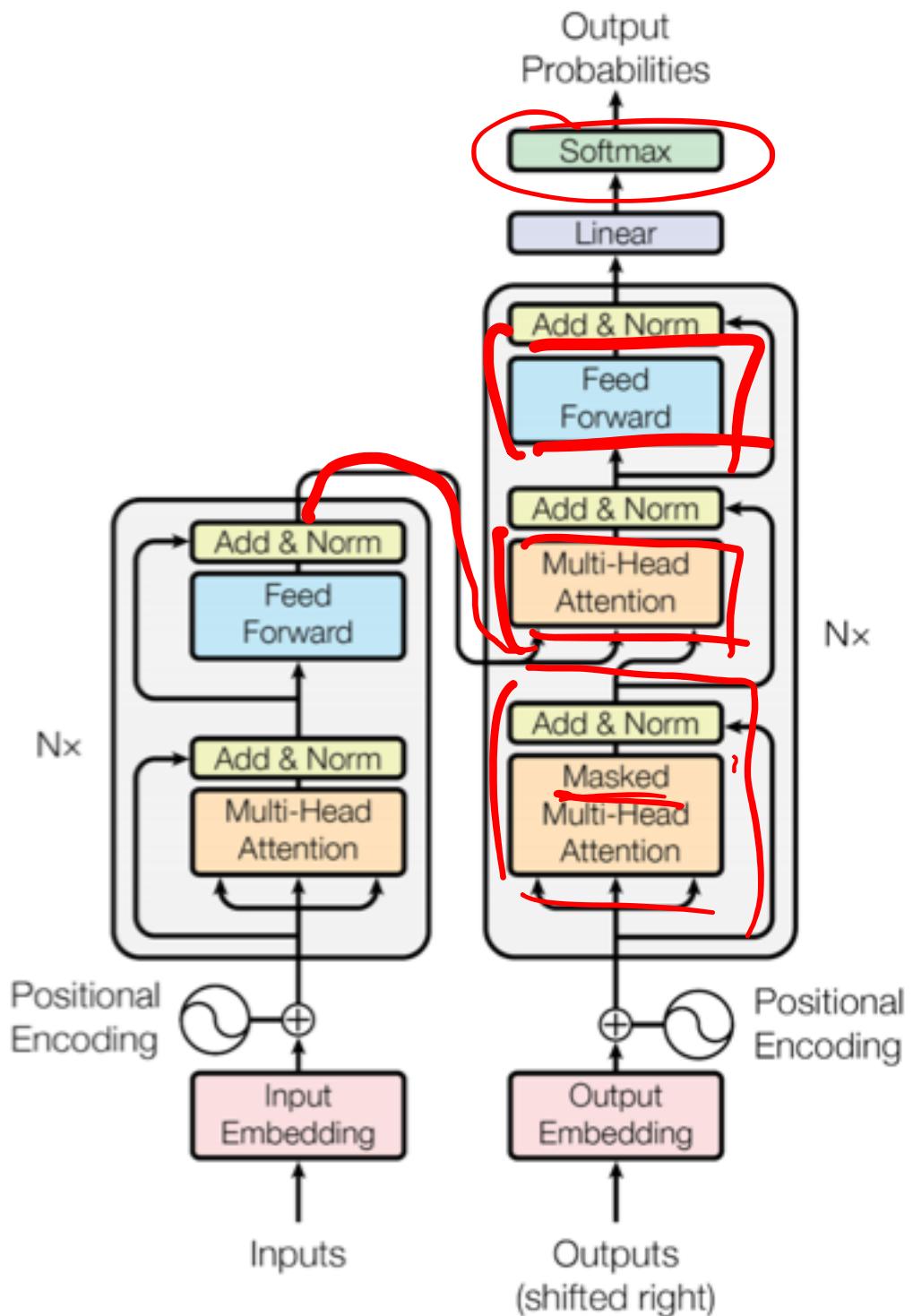
- 2 sublayer changes in decoder
- Masked decoder self-attention on previously generated outputs:



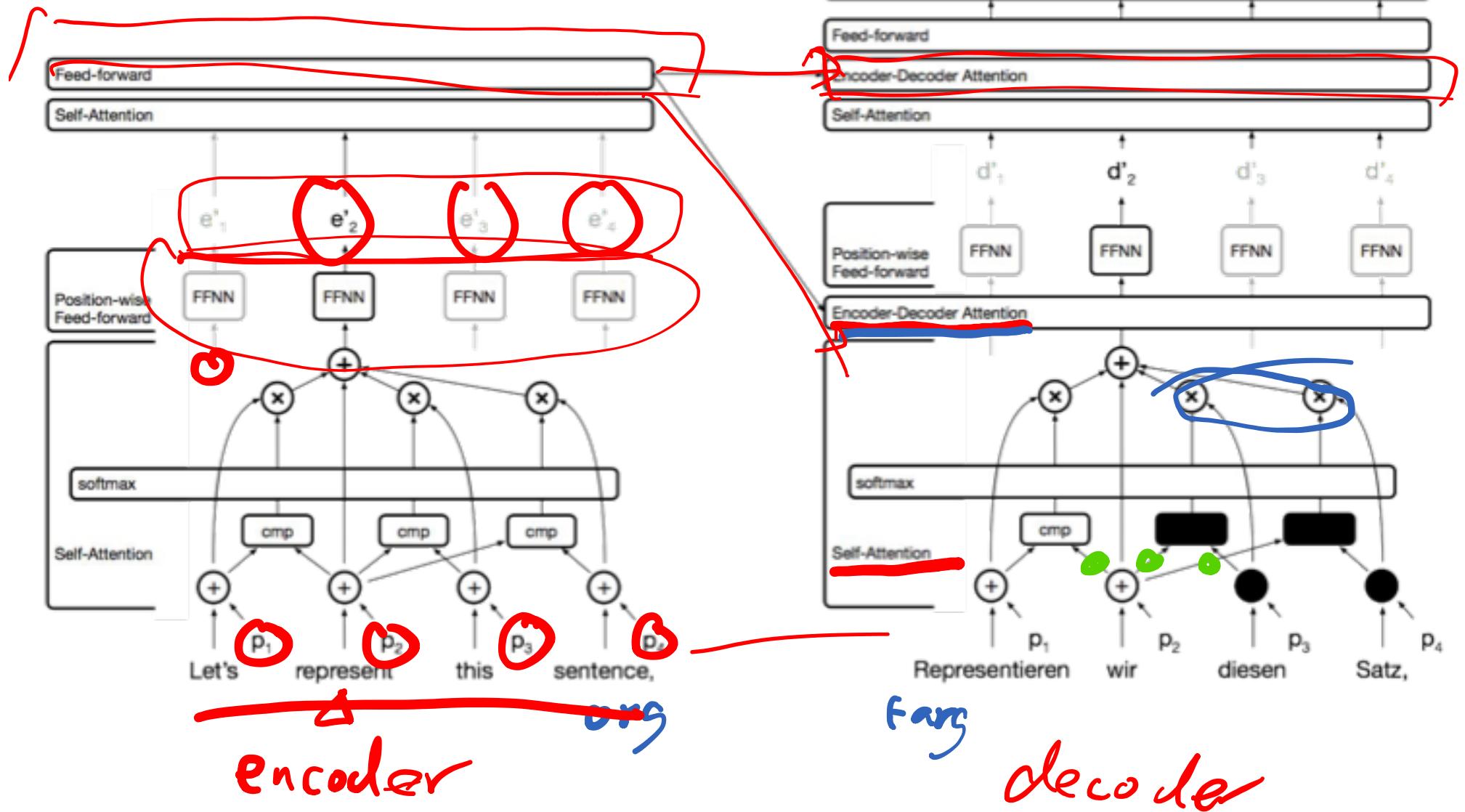
- Encoder-Decoder Attention, where queries come from previous decoder layer and keys and values come from output of encoder



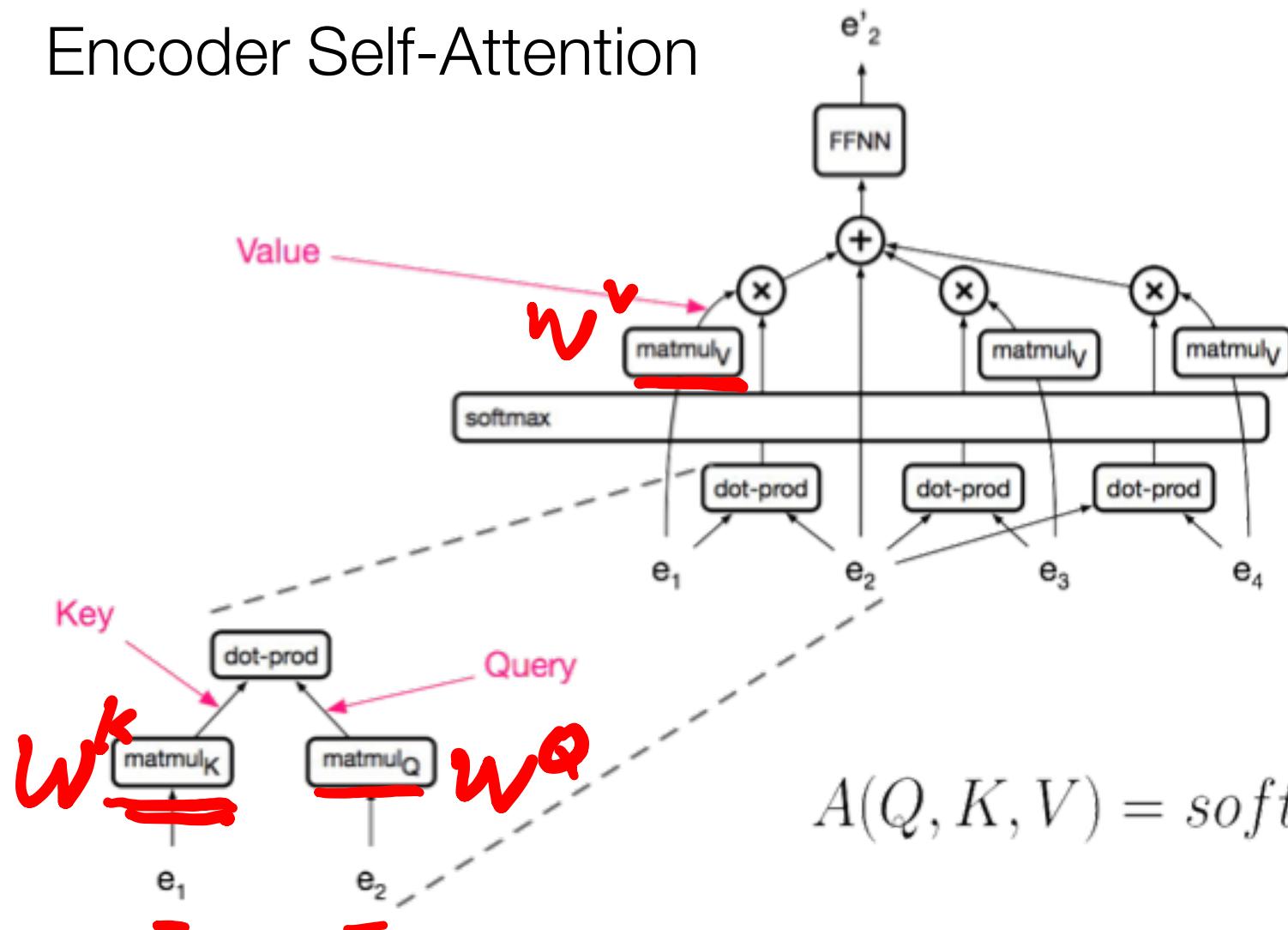
Blocks repeated 6 times also



The Transformer



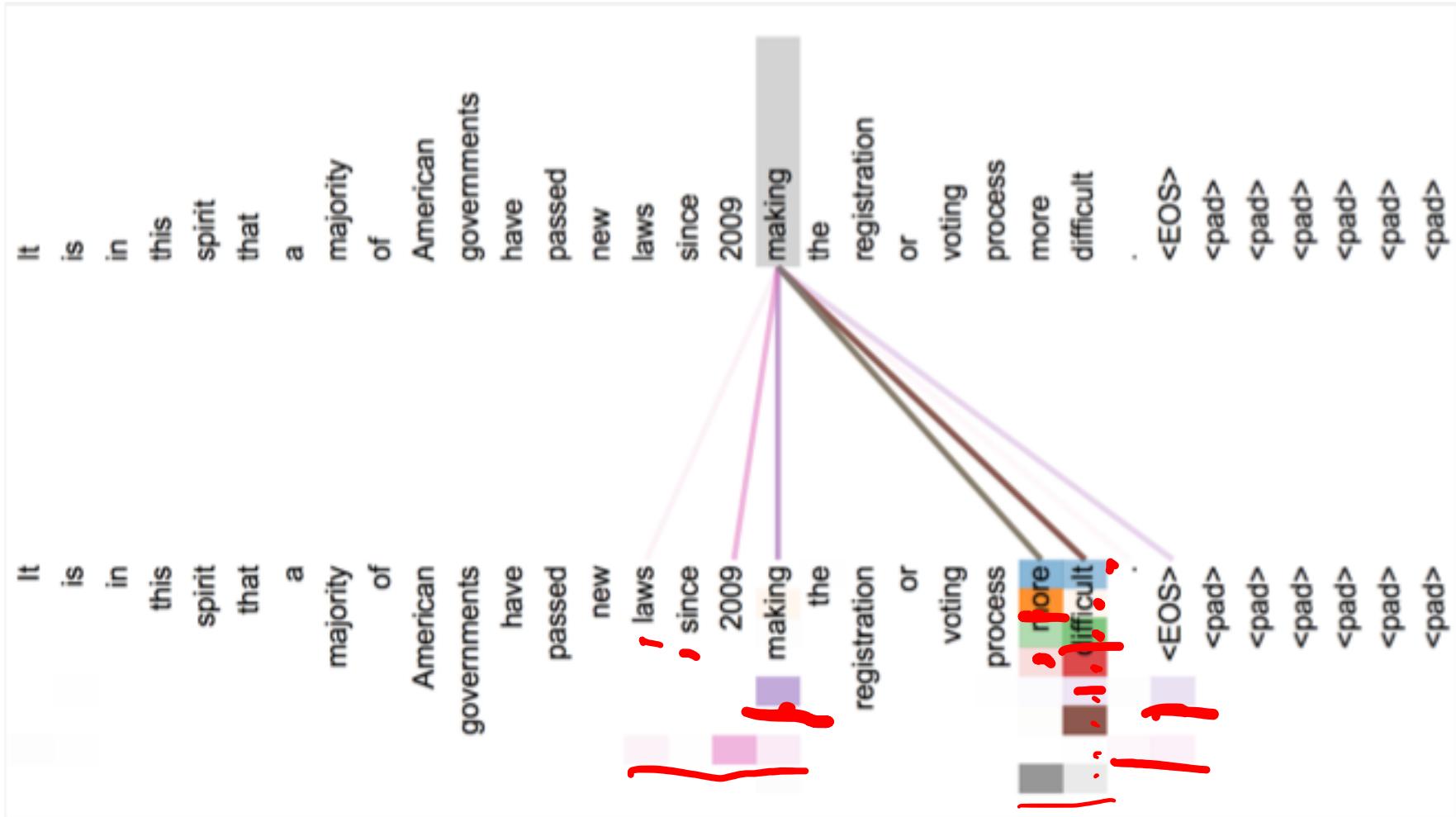
Encoder Self-Attention



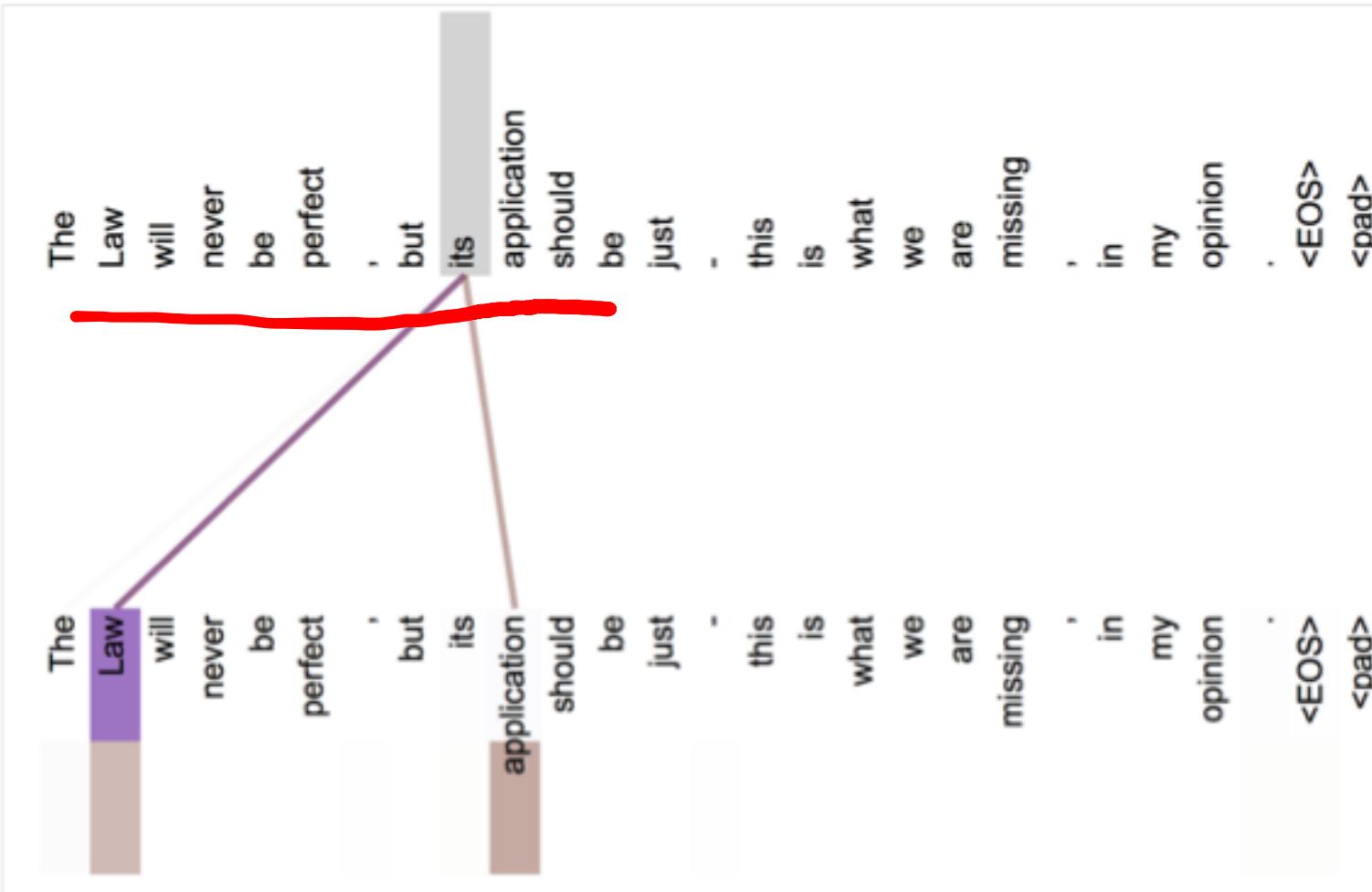
$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Attention visualization in layer 5

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



Attention visualization



In 5th layer. Isolated attentions from just the word 'its' for attention heads 5 and 6. Note that the attentions are very sharp for this word.

Tips and tricks of the Transformer

Details (in paper):

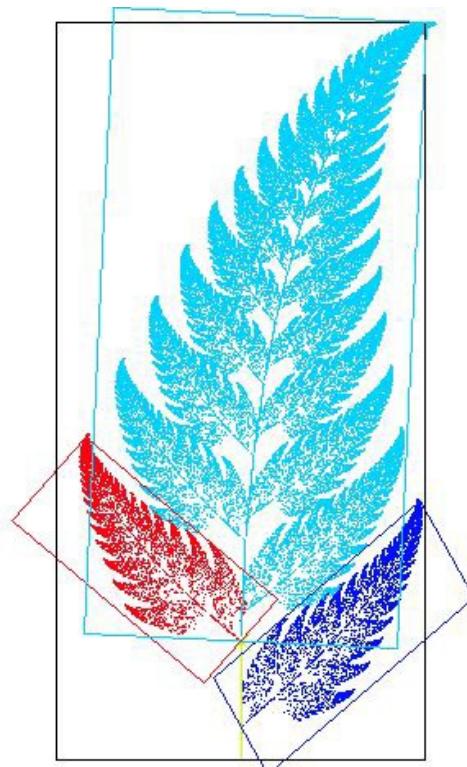
- Byte-pair encodings
 - Checkpoint averaging
 - ADAM optimizer with learning rate changes
 - Dropout during training at every layer just before adding residual
 - Label smoothing
 - Auto-regressive decoding with beam search and length penalties
- Use of transformers is spreading but they are hard to optimize
~~and unlike LSTMs don't usually just work out of the box and~~
they don't play well yet with other building blocks on tasks.

Experimental Results for MT

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1		$3.3 \cdot 10^{18}$
Transformer (big)	28.4	41.8		$2.3 \cdot 10^{19}$

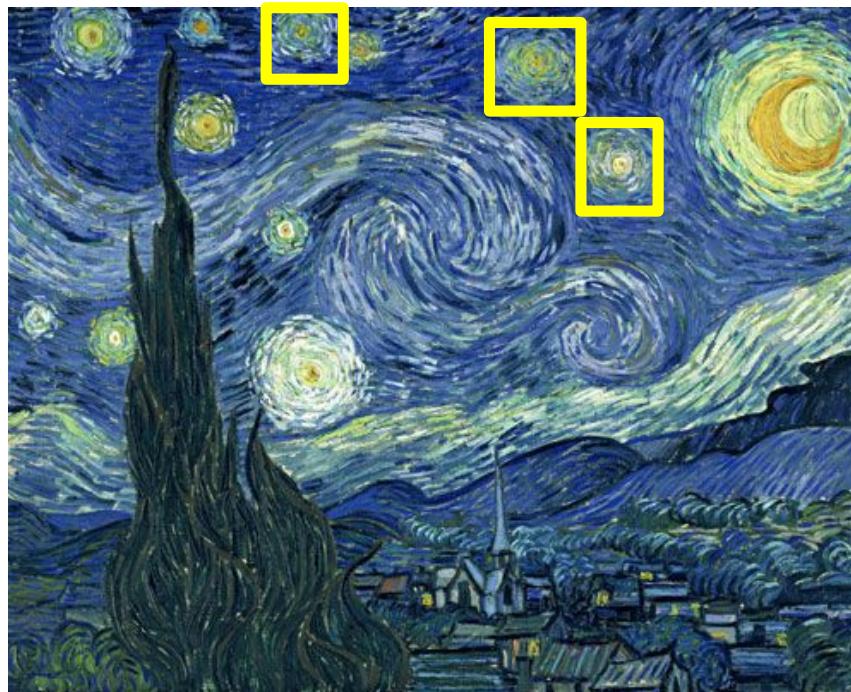
Self-Similarity, Image and Music Generation

Self-similarity in images



<https://en.wikipedia.org/wiki/Self-similarity>

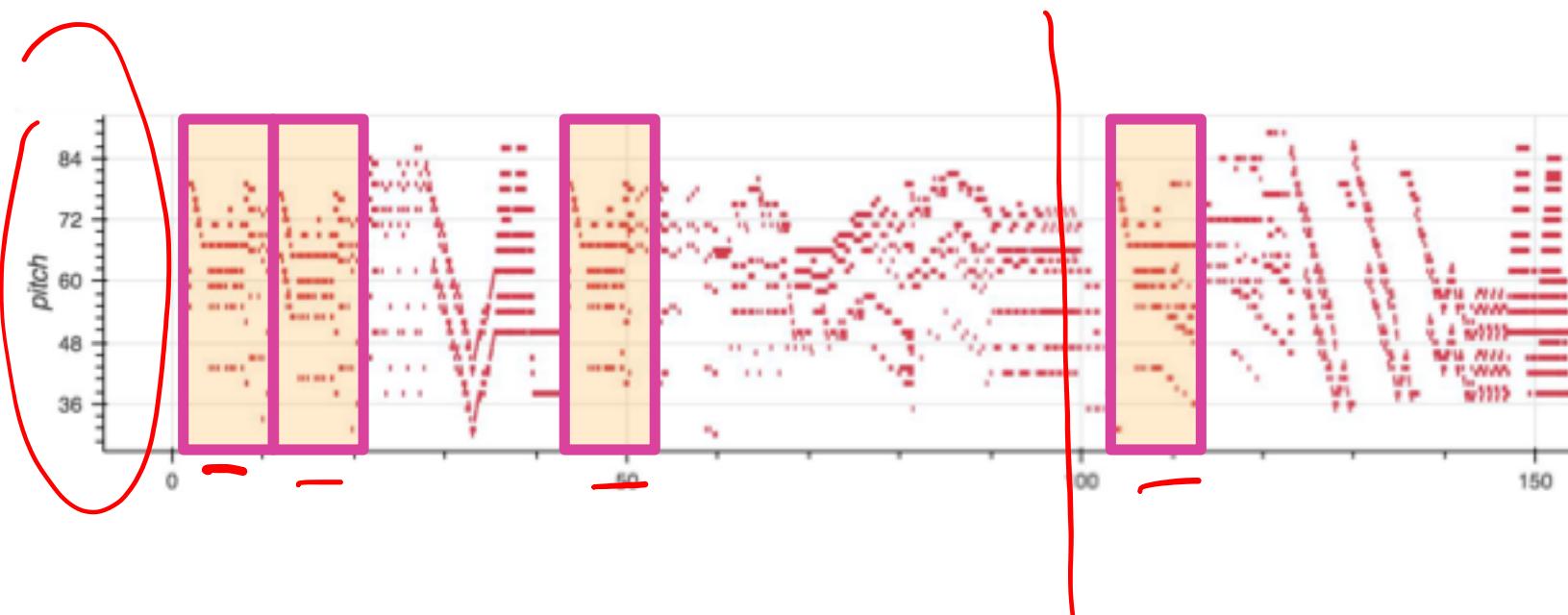
Self-Similarity in Images



Starry Night (Van Gogh, June 1889)

Self-similarity in music

Motifs repeat, immediately and also at a distance

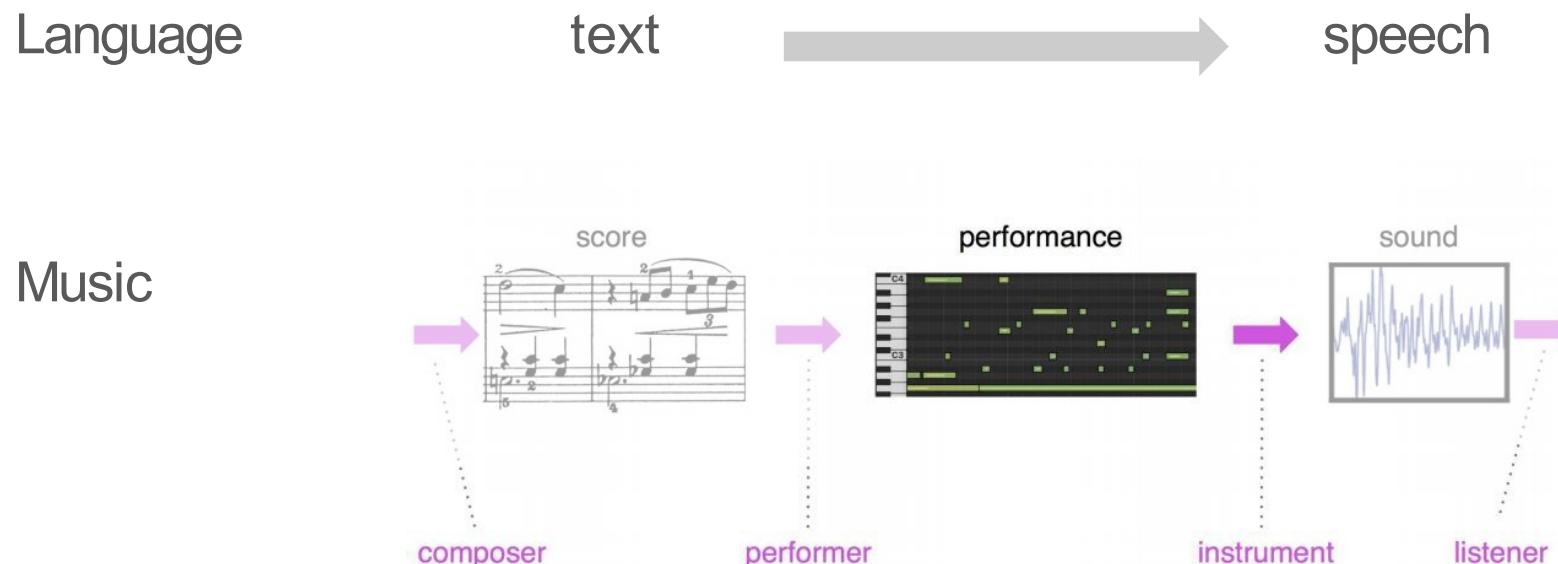


Music generation using relative self-attention

[Music Transformer](#) (ICLR 2019) by [Cheng-Zhi Anna Huang](#), [Ashish Vaswani](#), [Jakob Uszkoreit](#), [Noam Shazeer](#), [Ian Simon](#), [Curtis Hawthorne](#), [Andrew M. Dai](#), [Matthew D. Hoffman](#), [Monica Dinculescu](#) and [Douglas Eck](#).

Blog post: <https://magenta.tensorflow.org/music-transformer>

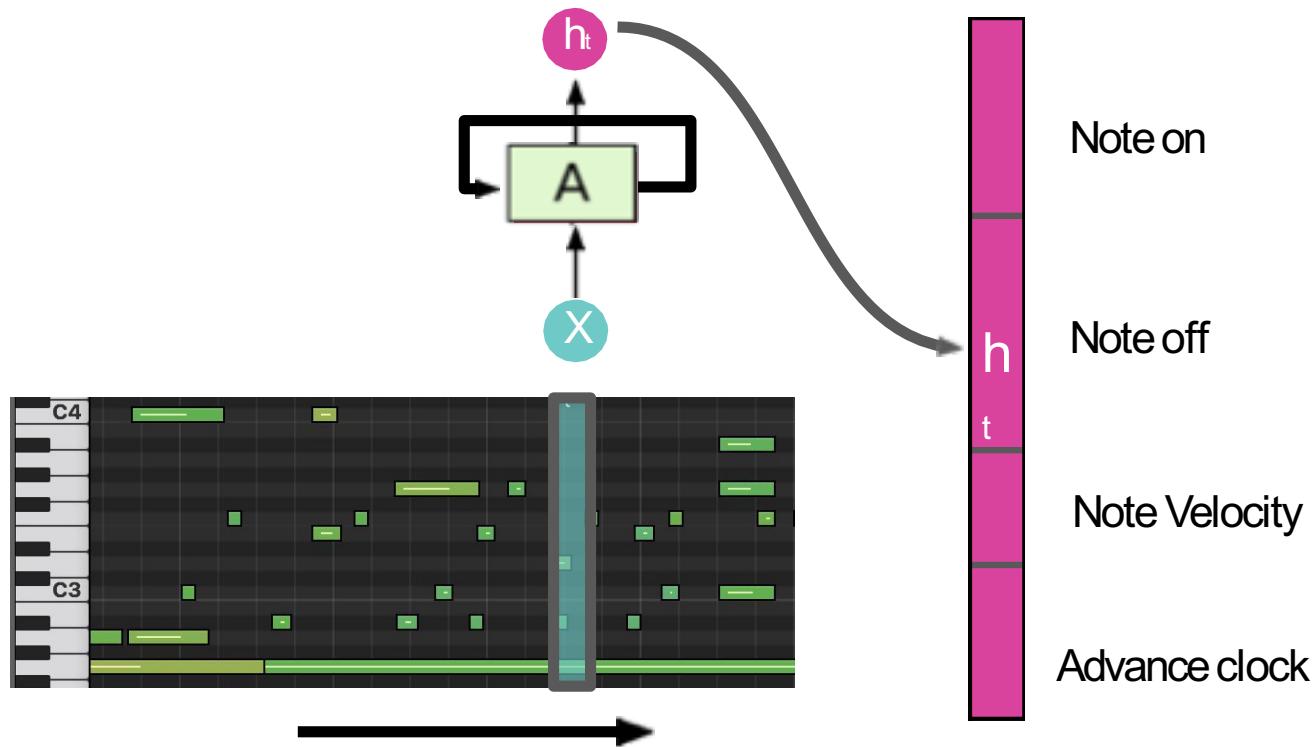
Raw representations in music and language



(Image from Simon & Oore, 2016)

Music Language model:

Prior work Performance RNN (Simon & Oore, 2016)

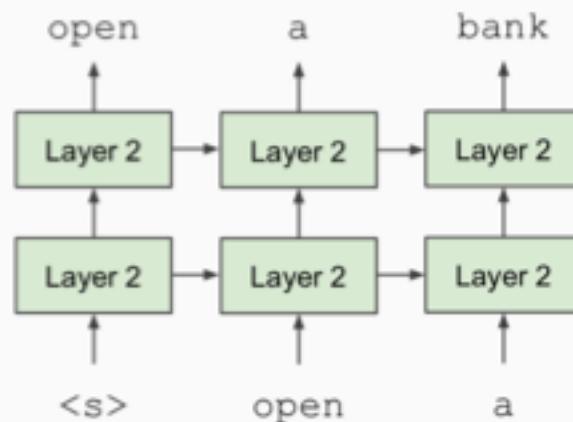


Sample from Music Transformer

BERT:Devlin, Chang, Lee, Toutanova (2018)

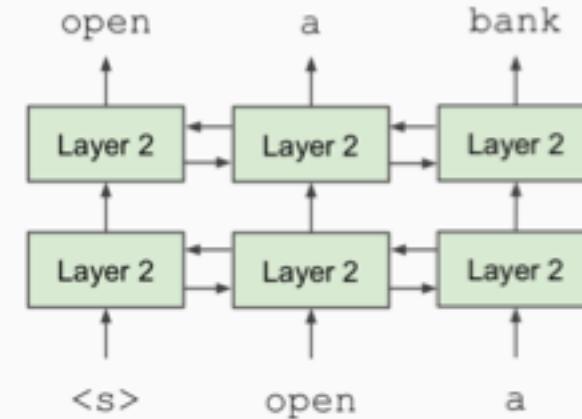
Unidirectional context

Build representation incrementally



Bidirectional context

Words can “see themselves”



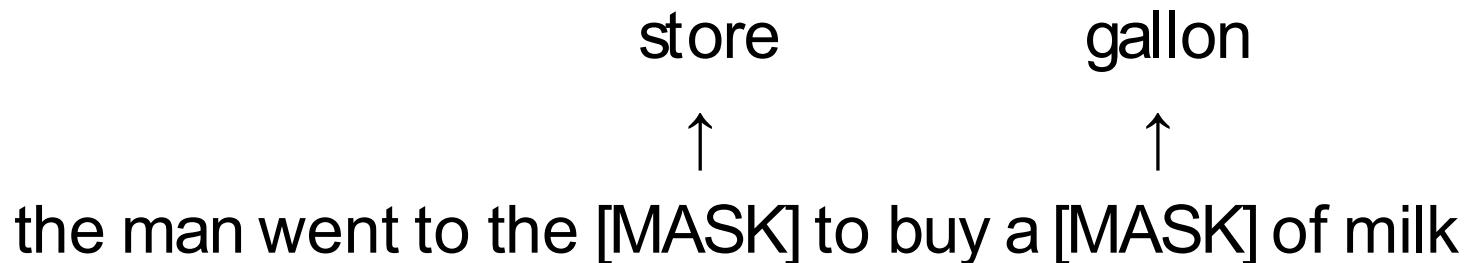
GPT

BERT

BERT:Devlin, Chang, Lee, Toutanova (2018)

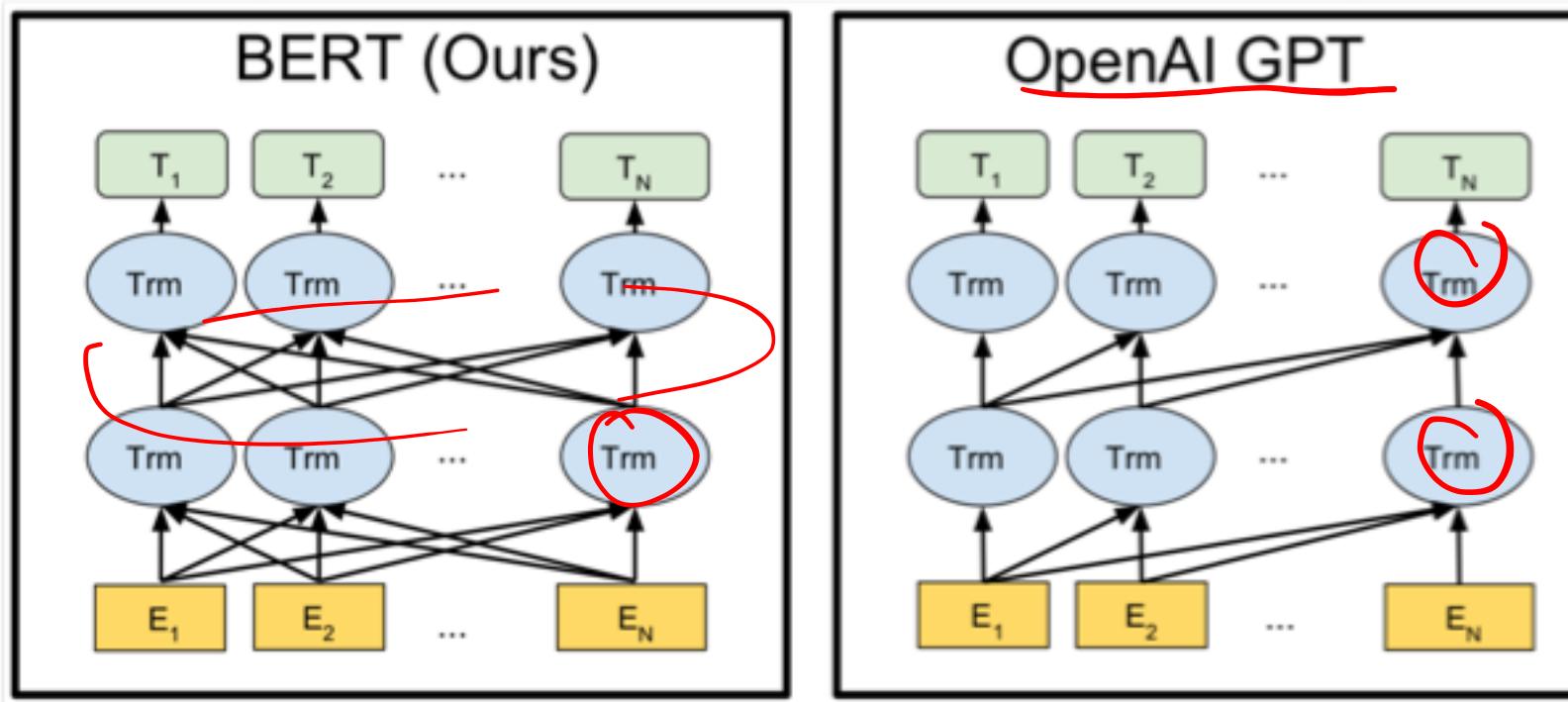
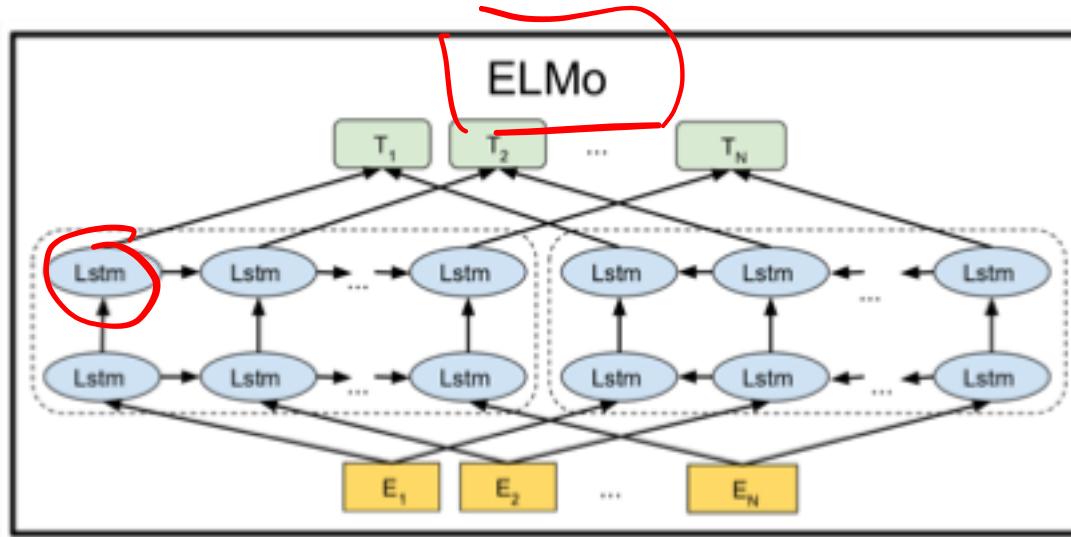
- **Solution:** Mask out $k\%$ of the input words, and then predict the masked words
 - They always use $k = 15\%$

M L M



- Too little masking: Too expensive to train
- Too much masking: Not enough context

BERT: Devlin, Chang, Lee, Toutanova (2018)



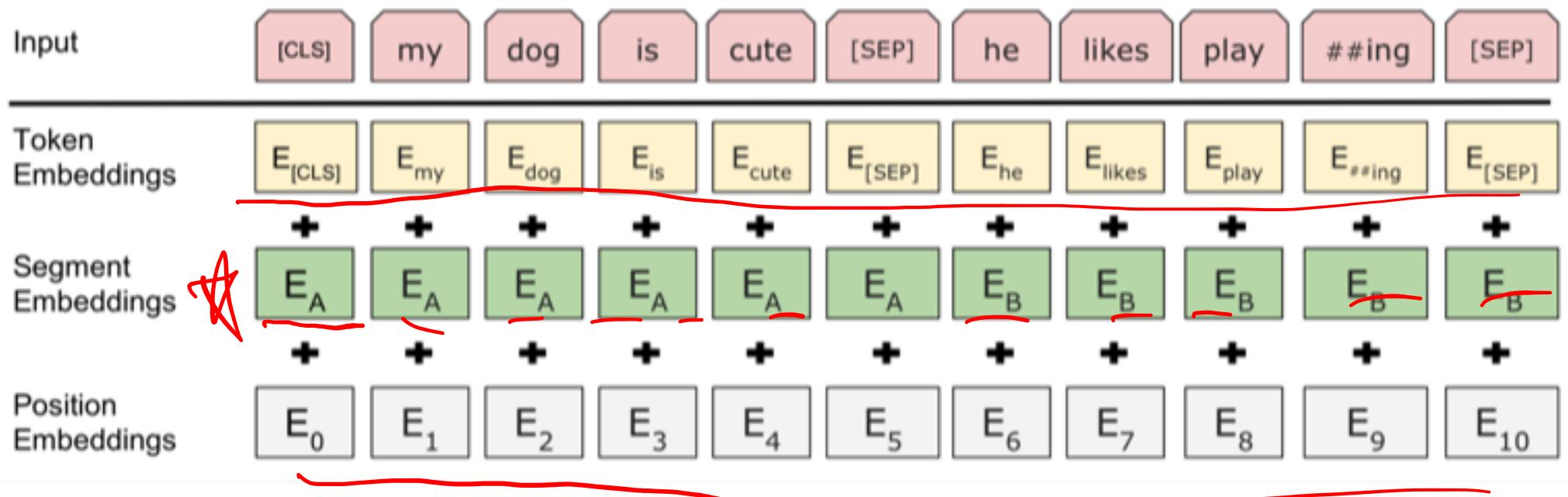
BERT complication: Next sentence prediction

- To learn *relationships* between sentences, predict whether Sentence B is actual sentence that proceeds Sentence A, or a random sentence

Sentence A = The man went to the store.
Sentence B = He bought a gallon of milk.
Label = IsNextSentence

Sentence A = The man went to the store.
Sentence B = Penguins are flightless.
Label = NotNextSentence

BERT sentence pair encoding



Token embeddings are word pieces

Learned segment embedding represents each sentence

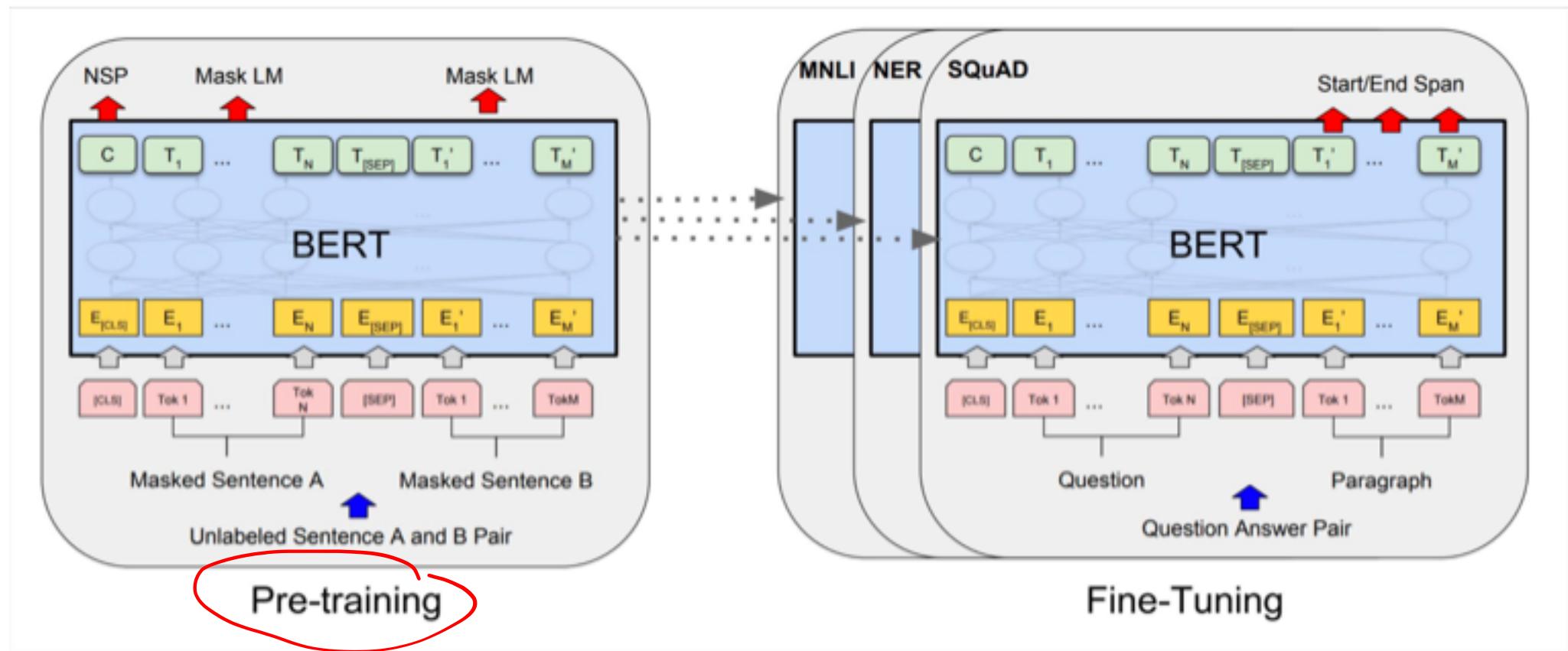
Positional embedding is as for other Transformer architectures

BERT model architecture and training

- Transformer encoder (as before)
 - Self-attention \Rightarrow no locality bias
 - Long-distance context has “equal opportunity”
 - Single multiplication per layer \Rightarrow efficiency on GPU/TPU
 - Train on Wikipedia + BookCorpus
 - Train 2 model sizes:
 - BERT-Base: 12-layer, 768-hidden, 12-head
 - BERT-Large: 24-layer, 1024-hidden, 16-head
 - Trained on 4x4 or 8x8 TPU slice for 4 days
-

BERT model fine tuning

- Simply learn a classifier built on the top layer for each task that you fine tune for



BERT results on GLUE tasks

- GLUE benchmark is dominated by natural language inference tasks, but also has sentence similarity and sentiment
- MultiNLI
 - Premise: Hills and mountains are especially sanctified in Jainism.
Hypothesis: Jainism hates nature.
Label: Contradiction
- CoLa
 - Sentence: The wagon rumbled down the road. Label: Acceptable
 - Sentence: The car honked down the road. Label: Unacceptable

BERT results on GLUE tasks

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

CoNLL2003 Named Entity Recognition

Name	Description	Year	F1
Elair (Zalando)	Character-level language model	2018	93.09
BERT Large	Transformer bidi LM + fine tune	2018	92.8
CVT Clark	Cross-view training + multitask learn	2018	92.61
BERT Base	Transformer bidi LM + fine tune	2018	92.4
ELMo	ELMo in BiLSTM	2018	92.22
TagLM Peters	LSTM BiLM in BiLSTMtagger	2017	91.93
Ma + Hovy	BiLSTM + char CNN + CRFlayer	2016	91.21
Tagger Peters	BiLSTM + char CNN + CRFlayer	2017	90.87
Ratinov + Roth	Categorical CRF+Wikipedia+word cls	2009	90.80
Finkel et al.	Categorical feature CRF	2005	86.86
IBM Florian	Linear/softmax/TBL/HMM ensemble, gazettes++	2003	88.76
S a tanford	MEMM softmax markov model	2003	86.07

BERT results on SQuAD 1.1

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar et al. '16)	82.304	91.221
1	BERT (ensemble) Google AI Language https://arxiv.org/abs/1810.04805	87.433	93.160
2	BERT (single model) Google AI Language https://arxiv.org/abs/1810.04805	85.083	91.835
2	nlnet (ensemble) Microsoft Research Asia	85.954	91.677
5	nlnet (single model) Microsoft Research Asia	83.468	90.133
3	QANet (ensemble) Google Brain & CMU	84.454	90.490

SQuAD 2.0 leaderboard, 2019-02-07

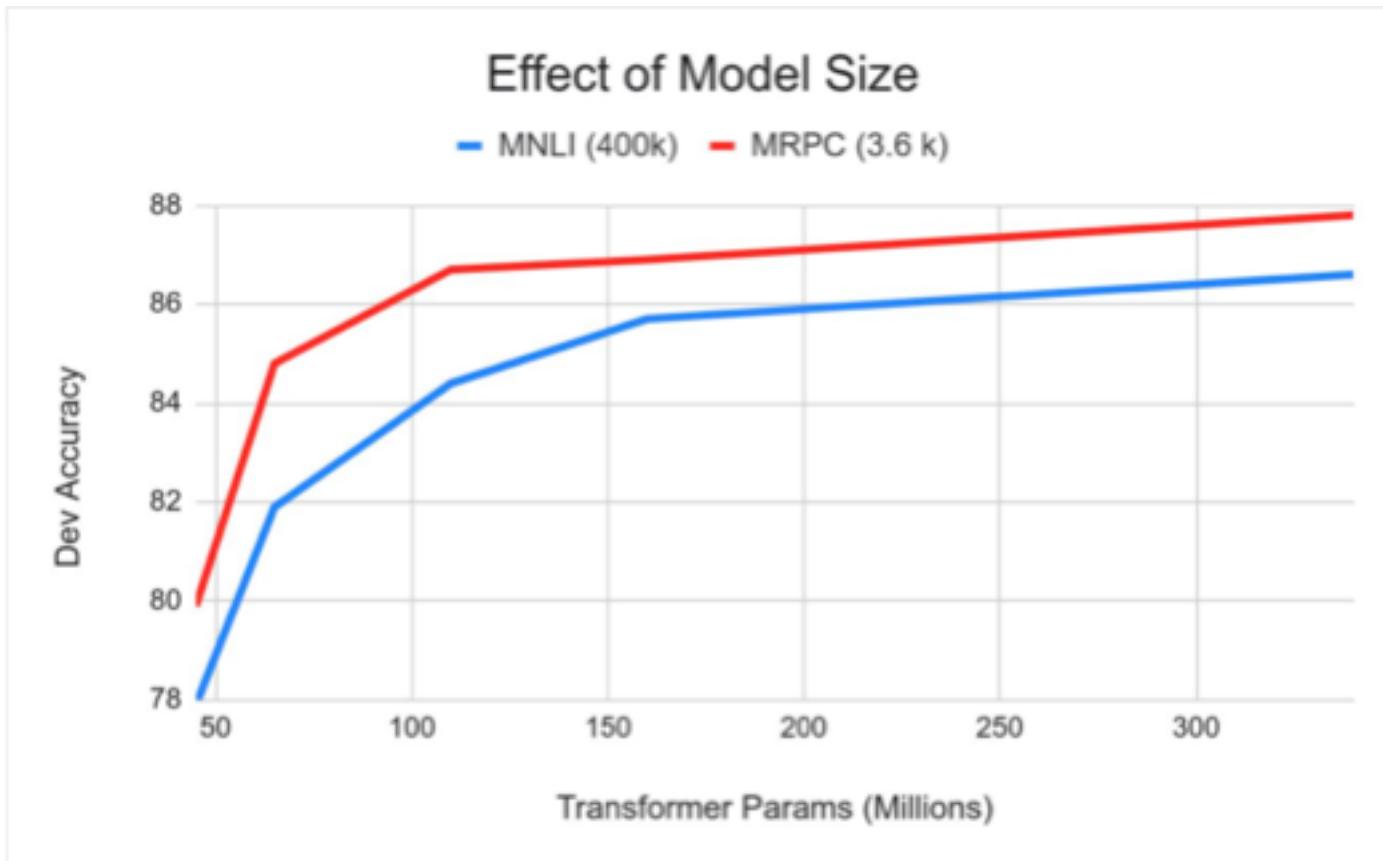
Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	BERT + MMFT + ADA (ensemble) Microsoft Research Asia <small>Jan 15, 2019</small>	85.082	87.615
2	BERT + Synthetic Self-Training (ensemble) Google AI Language https://github.com/google-research/bert <small>Jan 10, 2019</small>	84.292	86.967
3	BERT finetune baseline (ensemble) Anonymous <small>Dec 13, 2018</small>	83.536	86.096
4	Lunet + Verifier + BERT (ensemble) Layer 6 AI NLP Team <small>Dec 16, 2018</small>	83.469	86.043
4	PAML+BERT (ensemble model) PINGAN GammaLab <small>Dec 21, 2018</small>	83.457	86.122
5	Lunet + <u>Verifier</u> + BERT (single model) Layer 6 AI NLP Team <small>Dec 15, 2018</small>	82.995	86.035

SQuAD 2.0 leaderboard, 2019-09-04

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	XLNet + <u>DAAF + Verifier</u> (ensemble) <u>PINGAN Omni-Sinitic</u>	88.592	90.859
2	UPM (ensemble) Anonymous	88.231	90.713
3	XLNet + <u>SG-Net Verifier</u> (ensemble) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	88.174	90.702
4	XLNet + SG-Net Verifier++ (single model) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	87.238	90.071
5	UPM (single model) Anonymous	87.193	89.934

Size matters

- Going from 110M to 340M parameters helps a lot
- Improvements have not yet asymptoted



SQuAD 2.0 leaderboard, 2019-12-04

Rank	Model	EM	F1	Date	Team	Score	Score
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452				
1	ALBERT + DAAF + Verifier (ensemble) PINGAN Omni-Sinitic	90.002	92.425	Nov 06, 2019		88.174	90.702
2	ALBERT (ensemble model) Google Research & TTIC https://arxiv.org/abs/1909.11942	89.731	92.215	Sep 18, 2019		87.926	90.689
3	XLNet + DAAF + Verifier (ensemble) PINGAN Omni-Sinitic	88.592	90.859	Jul 22, 2019		87.238	90.071
3	albert+verifier (single model) Ping An Life Insurance Company AI Team	88.355	91.019	Nov 22, 2019		87.193	89.934
4	ALBERT (single model) Google Research & TTIC https://arxiv.org/abs/1909.11942	88.107	90.902	Sep 16, 2019		86.933	90.037
4	UPM (ensemble) Anonymous	88.231	90.713	Jul 26, 2019		87.147	89.474
5				Aug 04, 2019	XLNet + SG-Net Verifier (ensemble) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	86.820	89.795
6				Nov 15, 2019	XLNet (single model) Google Brain & CMU	86.651	89.595
7				Aug 04, 2019	XLNet + SG-Net Verifier++ (single model) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	86.448	89.586
8				Jul 26, 2019	UPM (single model) Anonymous	86.448	89.586
8				Nov 27, 2019	RoBERTa+Verify (ensemble) CW	86.238	90.071
8				Mar 20, 2019	BERT + DAE + AoA (ensemble) Joint Laboratory of HIT and iFLYTEK Research	86.000	89.934
9				Jul 20, 2019	RoBERTa (single model) Facebook AI	85.775	89.795
10				Sep 12, 2019	RoBERTa+Span (ensemble) CW	85.551	89.595
				Nov 12, 2019	RoBERTa+Verify (single model) CW	85.448	89.586

ALBERT: A LITE BERT FOR SELF-SUPERVISED LEARNING OF LANGUAGE REPRESENTATIONS

Zhenzhong Lan¹ Mingda Chen^{2*} Sebastian Goodman¹ Kevin Gimpel²

Piyush Sharma¹ Radu Soricut¹

¹Google Research

²Toyota Technological Institute at Chicago

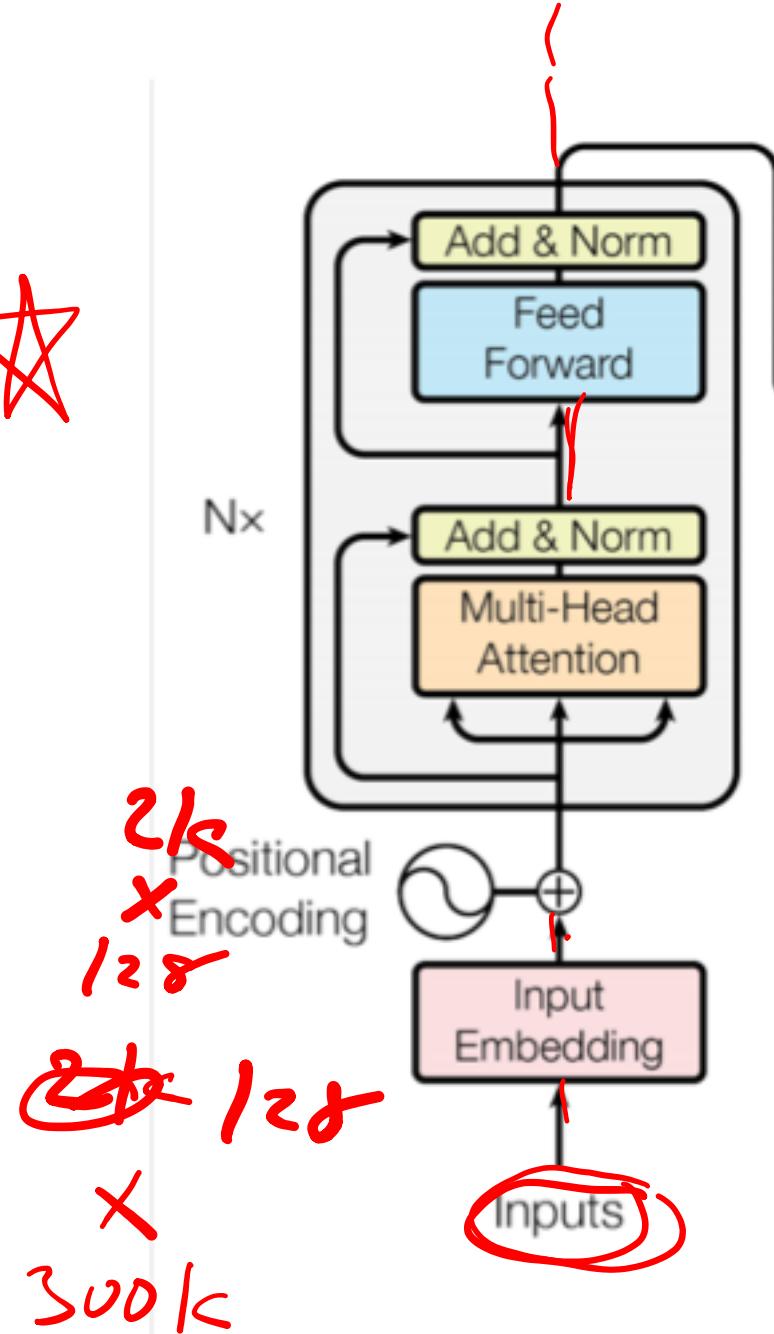
{lanzhzh, seabass, piyushsharma, rsoricut}@google.com
{mchen, kgimpel}@ttic.edu

ABSTRACT

Increasing model size when pretraining natural language representations often results in improved performance on downstream tasks. However, at some point further model increases become harder due to GPU/TPU memory limitations, longer training times, and unexpected model degradation. To address these problems, we present two parameter-reduction techniques to lower memory consumption and increase the training speed of BERT (Devlin et al., 2019). Comprehensive empirical evidence shows that our proposed methods lead to models that scale much better compared to the original BERT. We also use a self-supervised loss that focuses on modeling inter-sentence coherence, and show it consistently helps downstream tasks with multi-sentence inputs. As a result,

ALBERT vs. BERT

- Factorized embedding parameterization ~~☆~~
- Cross-layer parameter sharing
- Inter-sentence coherence loss ~~☆~~



ALBERT vs. BERT

	Model	Parameters	Layers	Hidden	Embedding	Parameter-sharing
BERT	base	108M	12	768	768	False
	large	334M	24	1024	1024	False
	xlarge	1270M	24	2048	2048	False
ALBERT	base	12M	12	768	128	True
	large	18M	24	1024	128	True
	xlarge	60M	24	2048	128	True
	xxlarge	235M	12	4096	128	True

	Model	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg	Speedup
BERT	base	108M	90.4/83.2	80.4/77.6	84.5	92.8	68.2	82.3	17.7x
	large	334M	92.2/85.5	85.0/82.2	86.6	93.0	73.9	85.2	3.8x
	xlarge	1270M	86.4/78.1	75.5/72.6	81.6	90.7	54.3	76.6	1.0
ALBERT	base	12M	89.3/82.3	80.0/77.1	81.6	90.3	64.0	80.1	21.1x
	large	18M	90.6/83.9	82.3/79.4	83.5	91.7	68.5	82.4	6.5x
	xlarge	60M	92.5/86.1	86.1/83.1	86.4	92.4	74.8	85.5	2.4x
	xxlarge	235M	94.1/88.3	88.1/85.1	88.0	95.2	82.3	88.7	1.2x

State-of-the-art (2019-9-16)

Models	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
<i>Single-task single models on dev</i>										
BERT-large	86.6	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet-large	89.8	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa-large	90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4	-	-
ALBERT (1M)	90.4	95.2	92.0	88.1	96.8	90.2	68.7	92.7	-	-
ALBERT (1.5M)	90.8	95.3	92.2	89.2	96.9	90.9	71.4	93.0	-	-
<i>Ensembles on test (from leaderboard as of Sept. 16, 2019)</i>										
ALICE	88.2	95.7	90.7	83.5	95.2	92.6	69.2	91.1	80.8	87.0
MT-DNN	87.9	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0	87.6
XLNet	90.2	98.6	90.3	86.3	96.8	93.0	67.8	91.6	90.4	88.4
RoBERTa	90.8	98.9	90.2	88.2	96.7	92.3	67.8	92.2	89.0	88.5
Adv-RoBERTa	91.1	98.8	90.3	88.7	96.8	93.1	68.0	92.4	89.0	88.8
ALBERT	91.3	99.2	90.5	89.2	97.1	93.4	69.1	92.5	91.8	89.4