

CS224n: NLP with Deep Learning

Lecture 8: Machine Translation

Statistical Machine Translation

Idea:

Learn a probabilistic model, using our data

Using Bayes Rule, it can be decomposed into:

- Translation Model:
Know about translation of local chunks of phrases
- Language Model:
Writing good English (target language)

1990s-2010s: Statistical Machine Translation

- Core idea: Learn a **probabilistic model** from **data**
- Suppose we're translating French \rightarrow English.
- We want to find **best English sentence y** , given French sentence x

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

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

$$= \operatorname{argmax}_y \underbrace{P(x|y)}_{\text{Translation Model}} \underbrace{P(y)}_{\text{Language Model}}$$

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

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

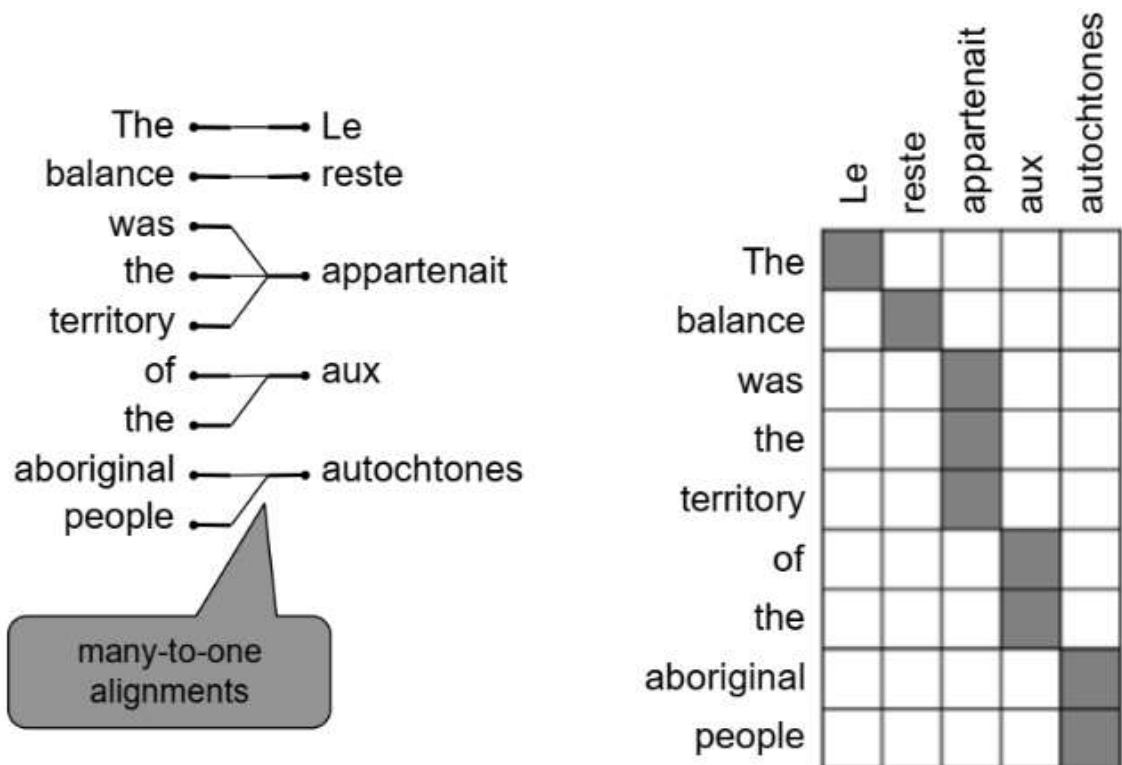
7

Alignment

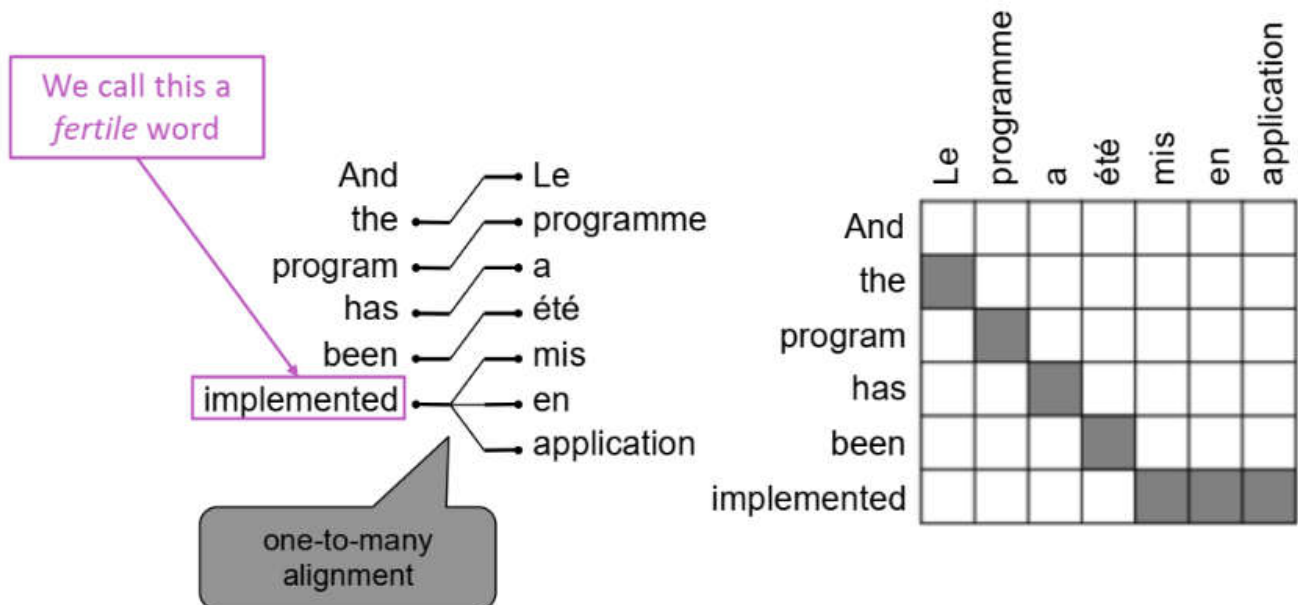
Alignment = correspondence between particular words

- Not necessarily bijective:
Can be:
 - Many-to-one
 - One-to-many: the **one** word is called a *fertile* word
 - Many-to-many

Alignment can be many-to-one

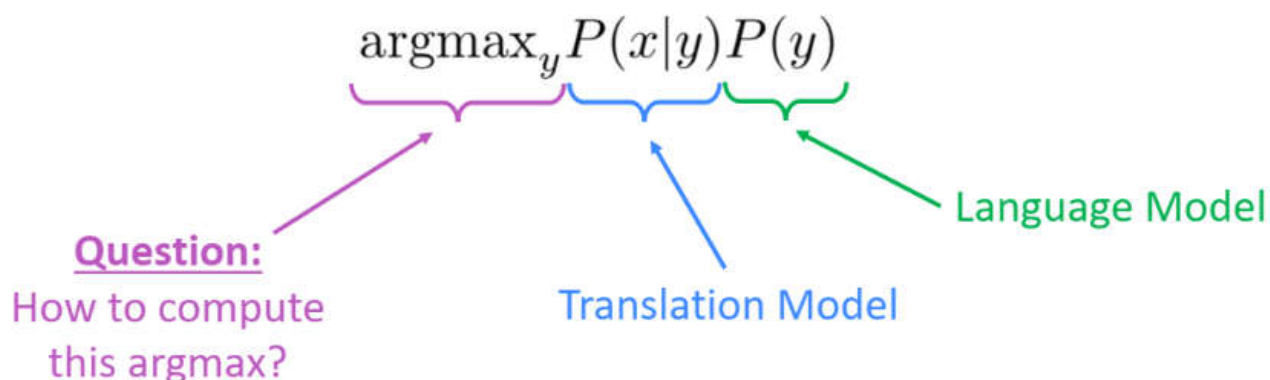


Alignment can be one-to-many



Back to SMT

Decoding for SMT



- We can't enumerate every y to calculate this probability
- We explore the different possibilities, and discard the low-probability ones as we go
→ Beam Search

Neural Machine Translation

- Neural Machine Translation is called **sequence-to-sequence** (seq2seq)
- Involves 2 RNN

Other usages of Seq2seq

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)

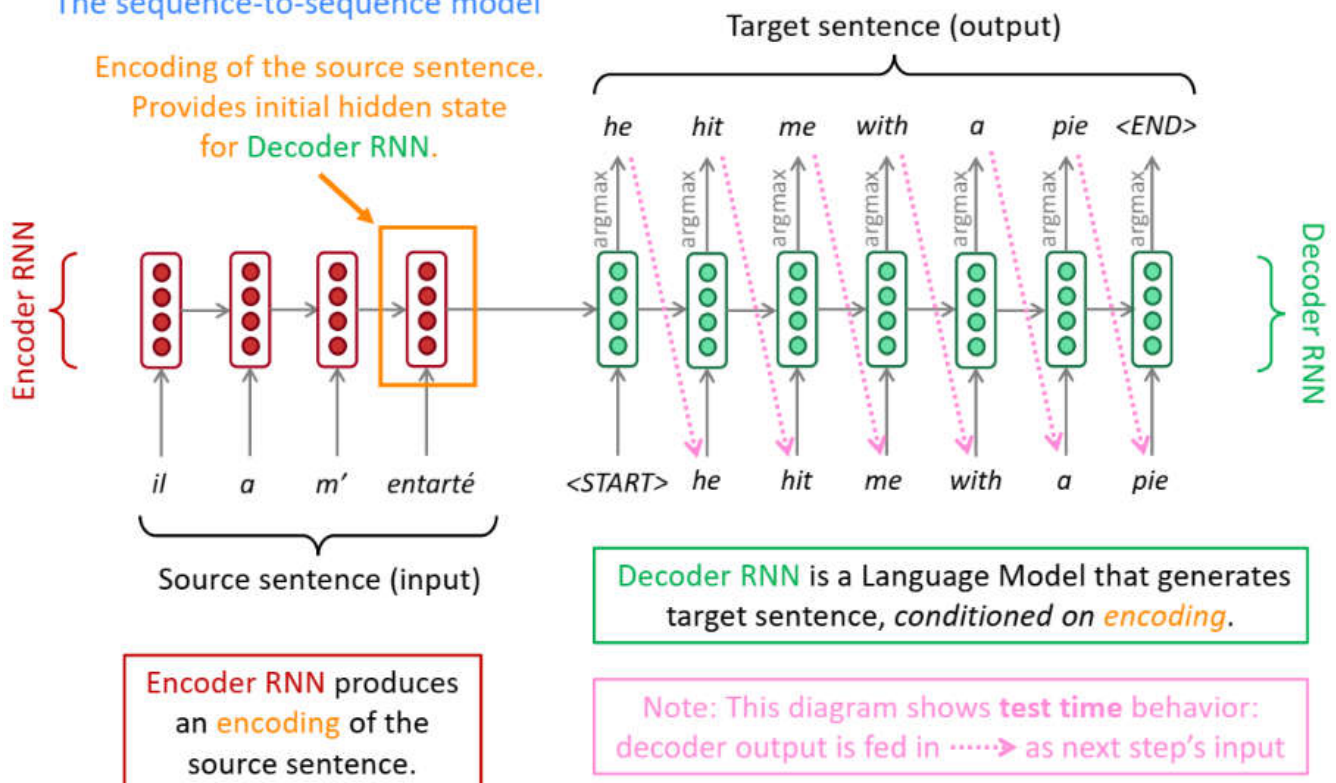
Encoding / Decoding

- Encoding of the source sentence = last hidden state of the Decoder RNN
- Decoder RNN is a **conditional** Language Model, as it is conditioned on this encoding

At Test time

Neural Machine Translation (NMT)

The sequence-to-sequence model



24

NMT Principle

Direct calculation of $P(y|x)$:

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

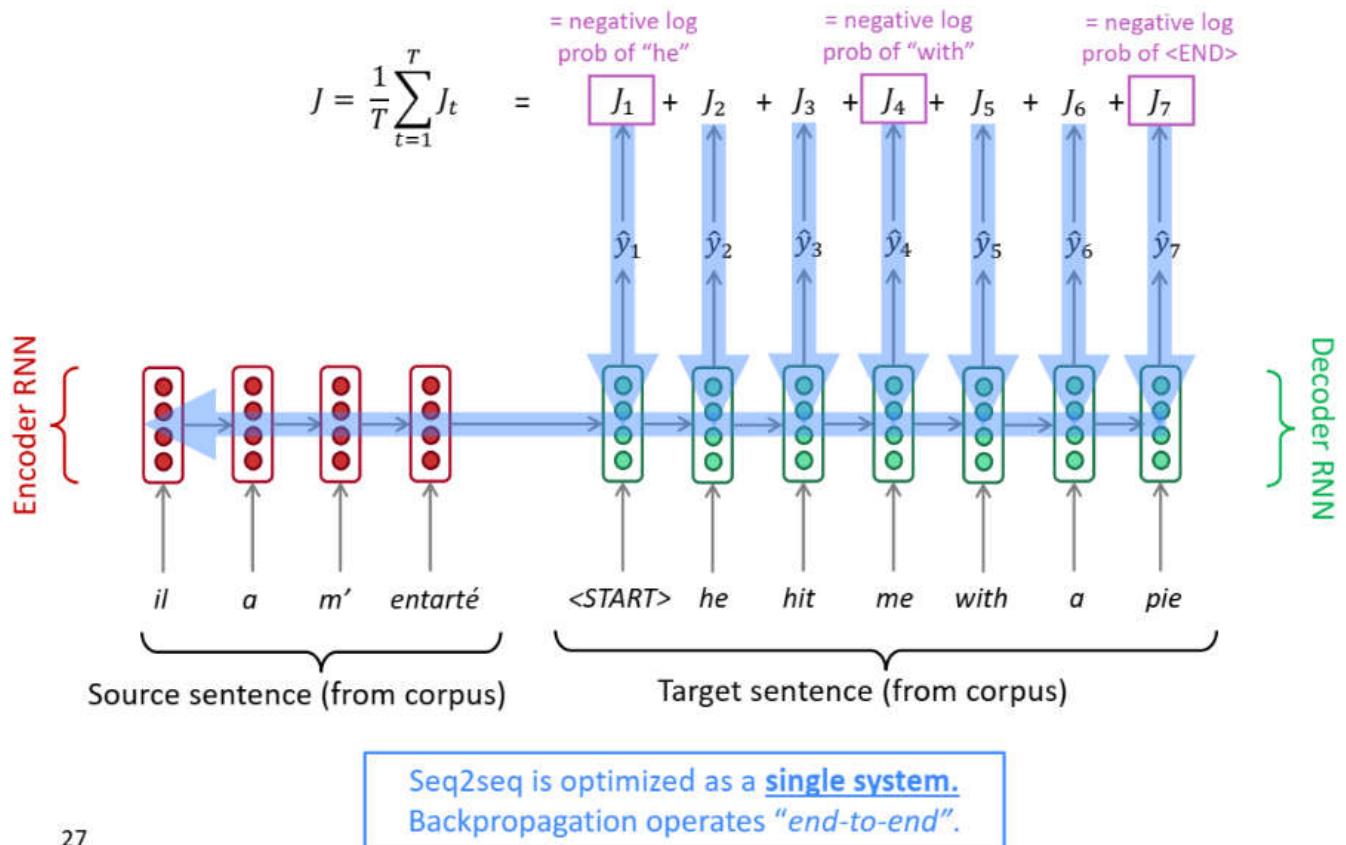
Training a Model

We need:

- A parallel corpus

- Words embeddings for words in both languages (source & target)
- Here, in NMT, we learn the translation objective directly
- While in SMT, we were doing it indirectly, by separating the task into different subtasks

Training a Neural Machine Translation system



27

The backprop is happening end-to-end:

- one end is the loss functions
- the other end is the beginning state of the encoder RNN

Backprop flows through the entire system

⚠ Difference between Training & Testing ⚠

During Training, in the Decoder RNN, we don't feed the previous prediction into the next step. Instead, we feed each state the correct previous word from the corpus.

It is possible to train the 2 RNN separately:

for example, training a strong Language Model on its own, then initializing the Decoder RNN with it

Remarks

- In practice, we pad the short sentences up to a certain length

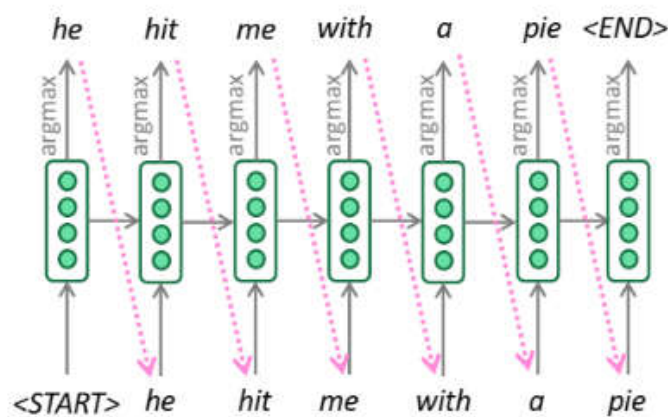
- We can use pre-trained word embeddings, or use word2vec or GloVe to train them
And then, do the training of the Encoder/Decoder

Towards Beam search Decoding

Problem of Greedy Decoding

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?**

⚠ Problem ⚠

The argmax at each step doesn't necessarily give us the global argmax for the whole sentence !

(ie, usually, the global optimum is not the combination of all the local optima)

- We can't search all the possibilities for the argmax, as it would be much too expensive to compute

→ Use a search algorithm: **Beam search**

Beam Search Decoding

- Beam size k :
how big our search space is at any time

Beam search decoding

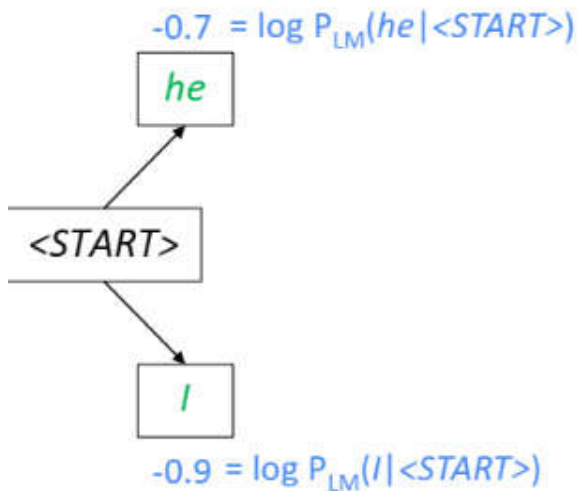
- Core idea: On each step of decoder, keep track of the *k* most probable partial translations (which we call *hypotheses*)
 - *k* is the *beam size* (in practice around 5 to 10)
- A hypothesis y_1, \dots, y_t has a *score* which is its log probability:
$$\text{score}(y_1, \dots, y_t) = \log P_{\text{LM}}(y_1, \dots, y_t | x) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$$
 - Scores are all negative, and higher score is better
 - We search for high-scoring hypotheses, tracking top *k* on each step
- Beam search is *not guaranteed* to find optimal solution
- But *much more efficient* than exhaustive search!

31

Example

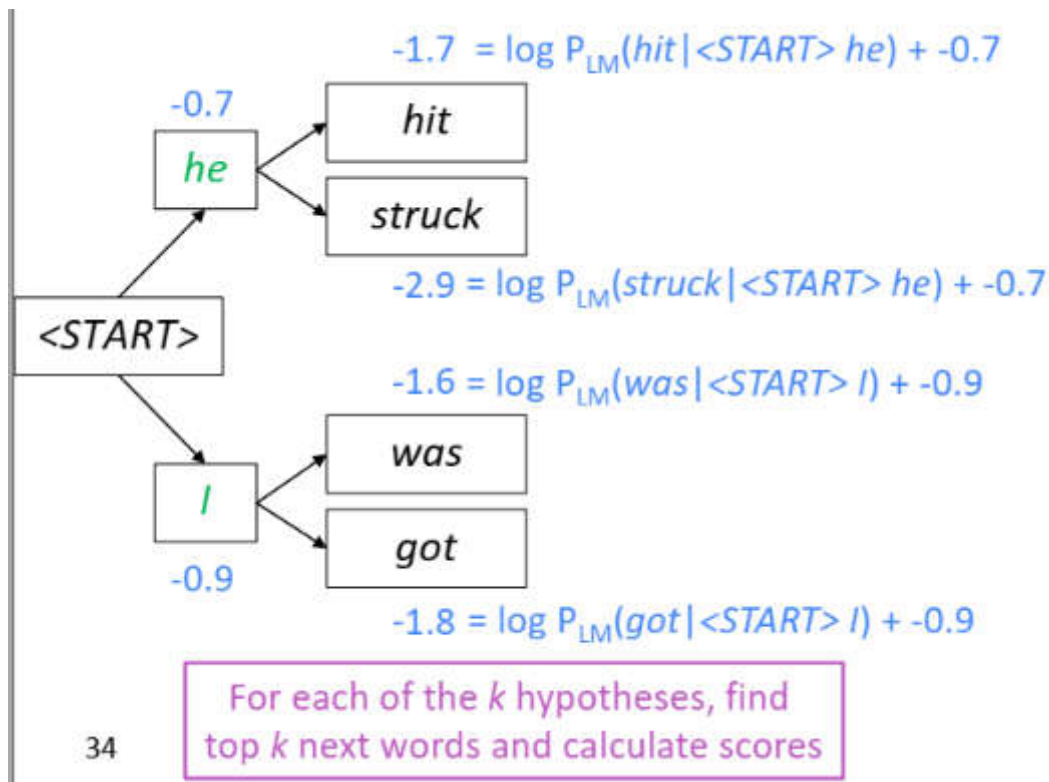
Beam search decoding: example

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



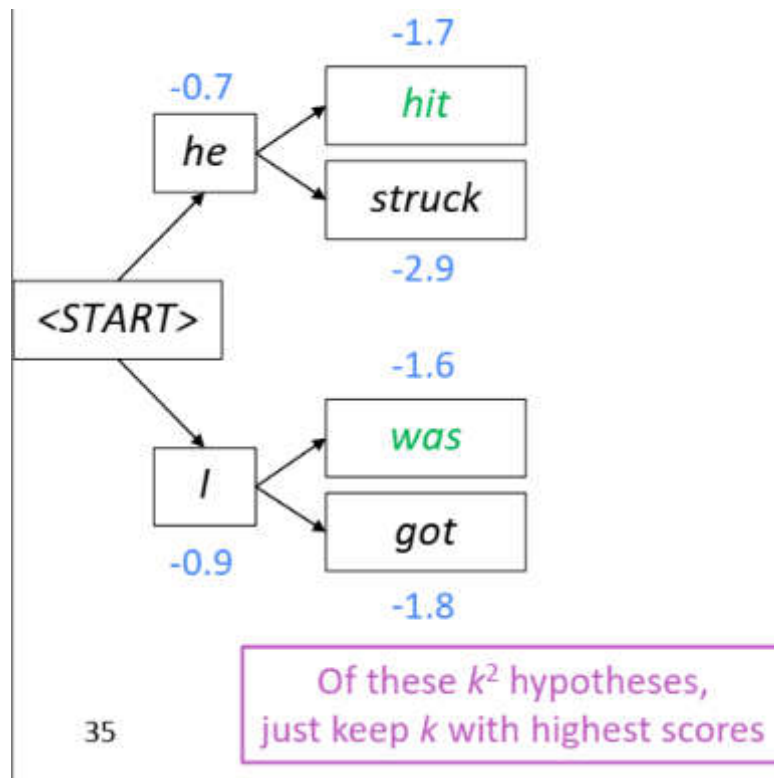
33

Take top k words
and compute scores

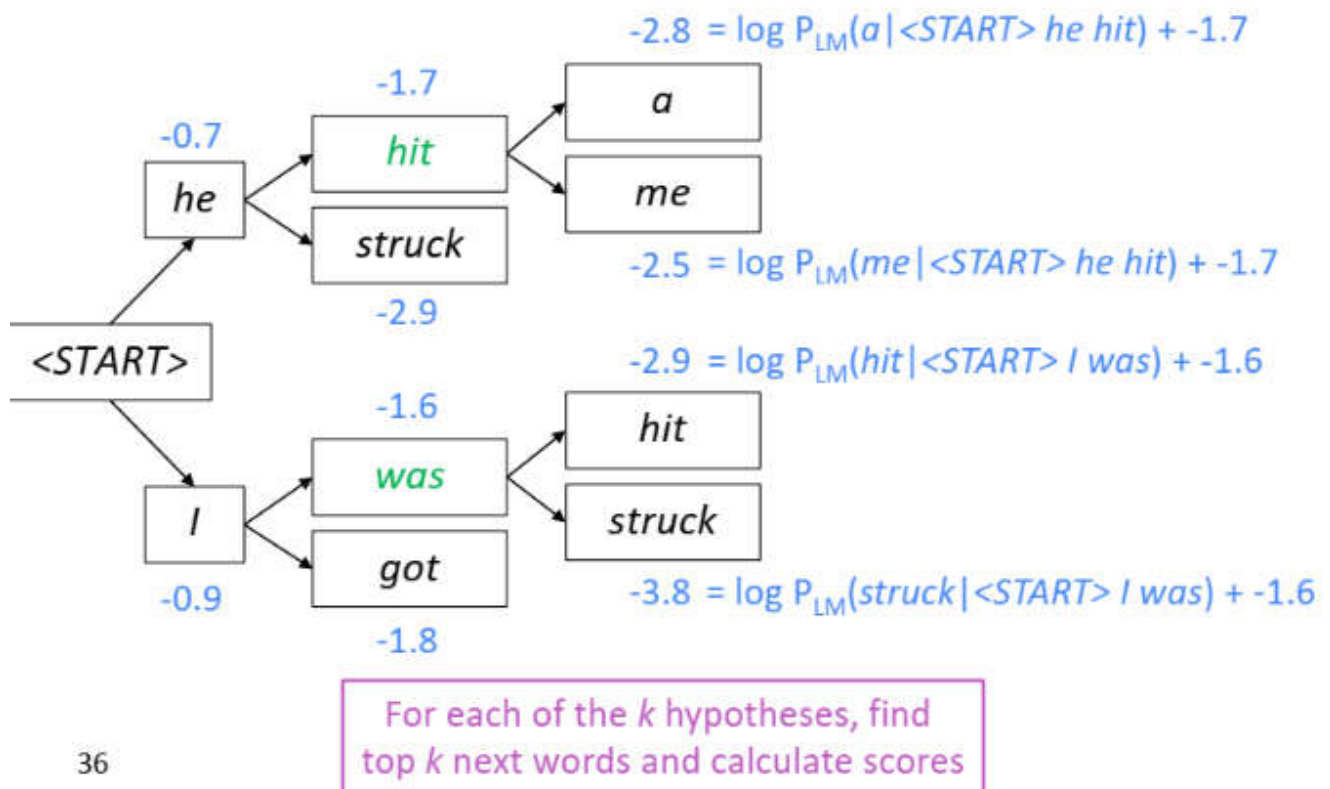


34

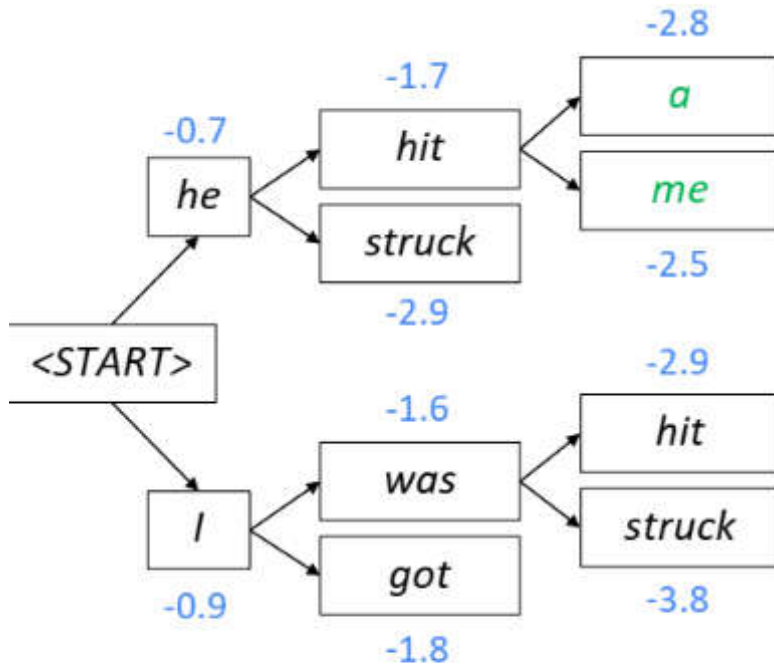
For each of the k hypotheses, find
top k next words and calculate scores



35

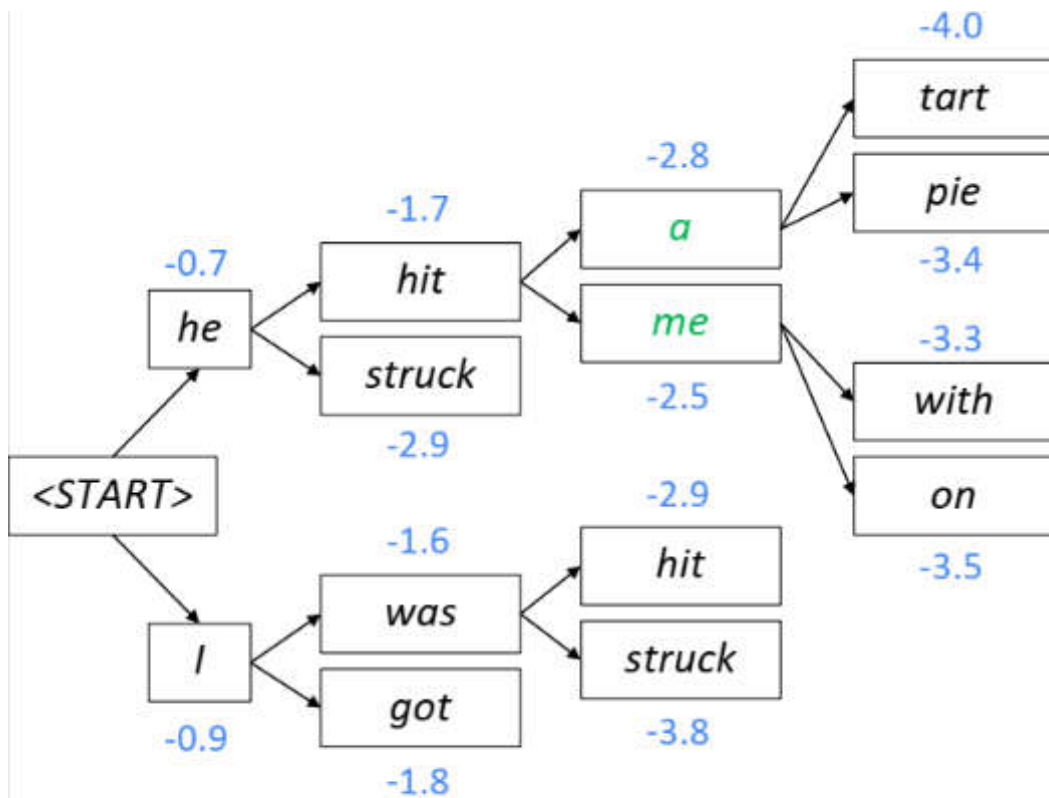


36



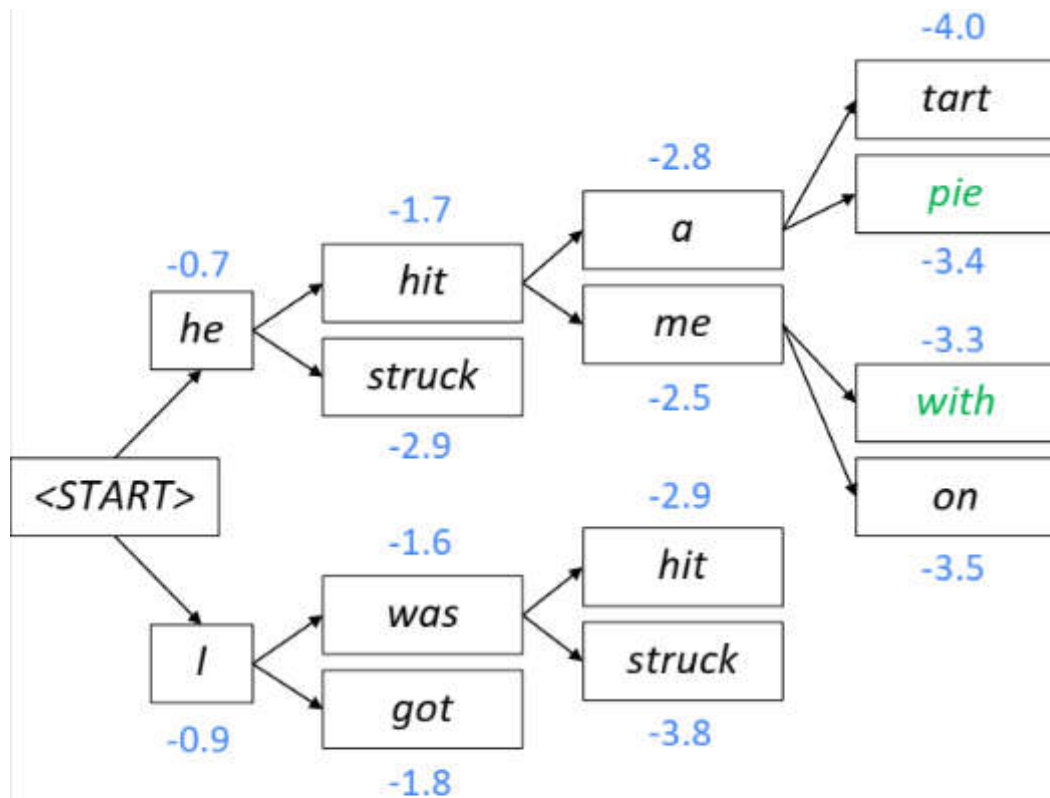
37

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



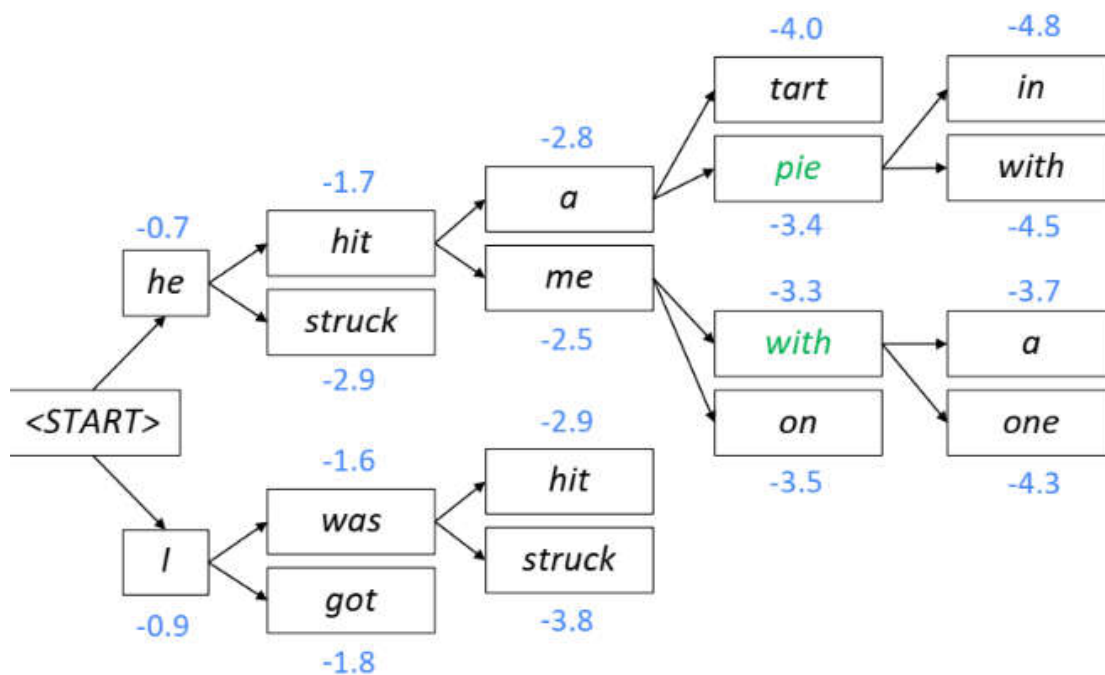
38

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



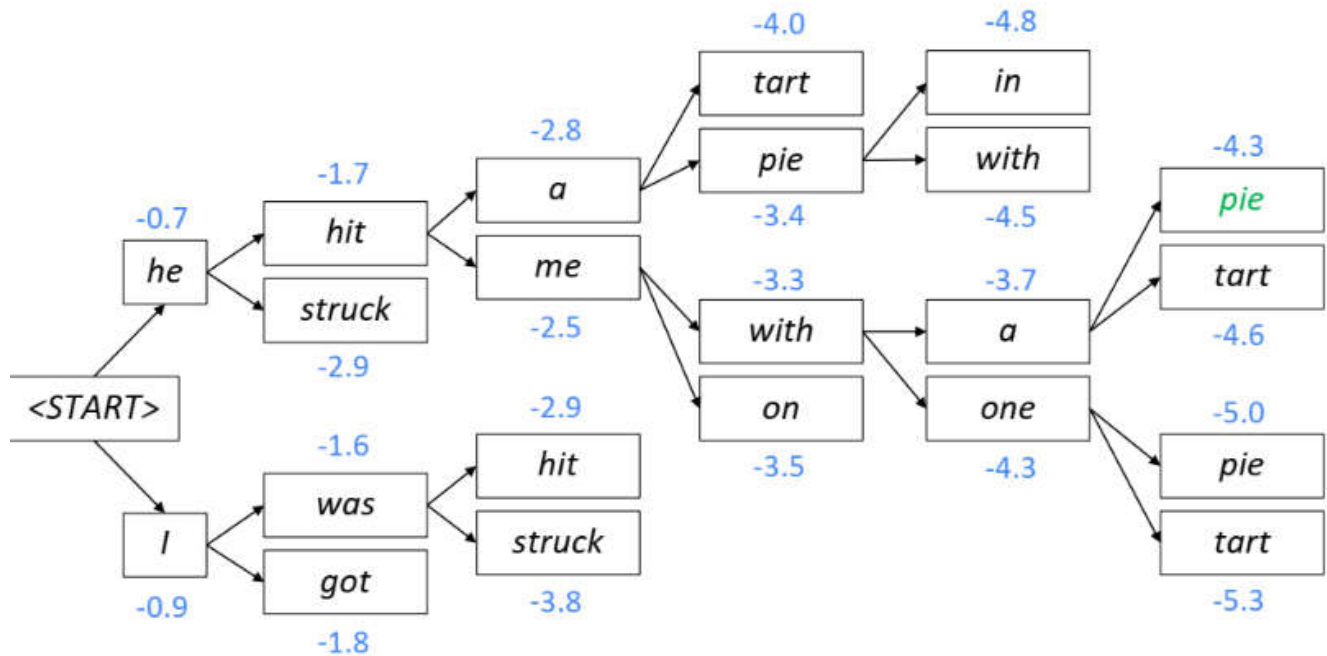
39

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



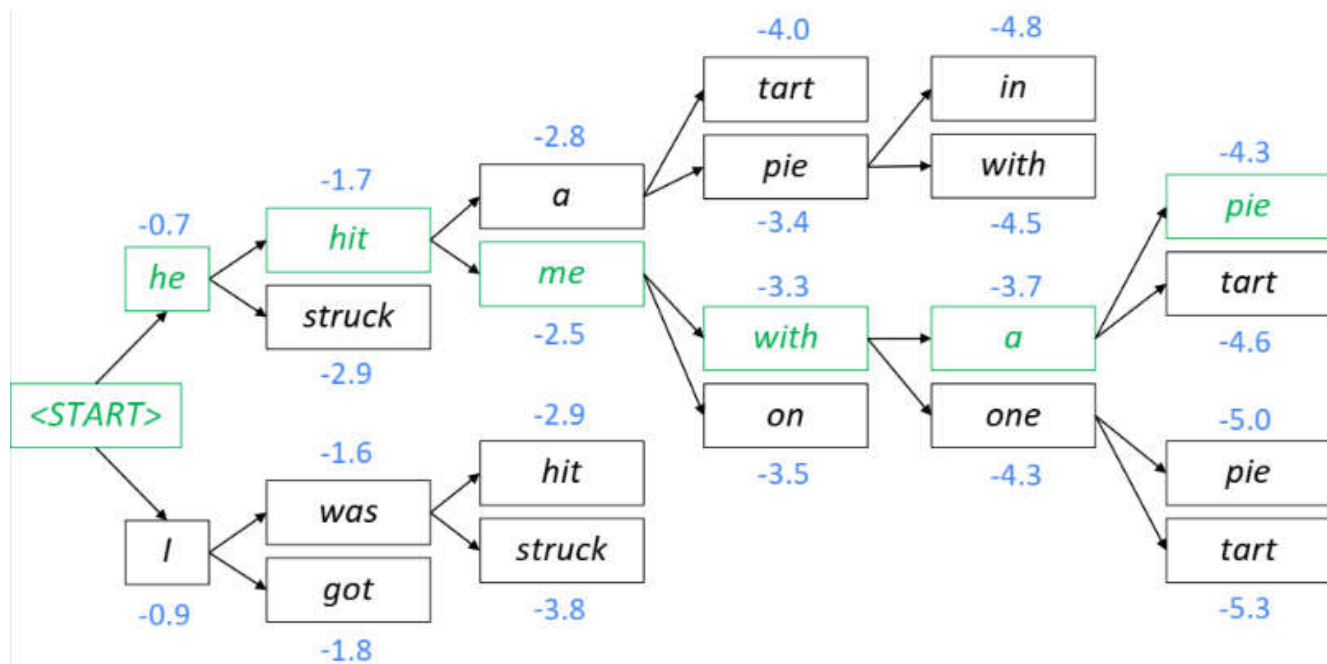
40

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



43

This is the top-scoring hypothesis!



44

Backtrack to obtain the full hypothesis

Stopping Criterion

Continue exploring our hypothesis until:

- We reach a certain predefined timestep T
- We have a certain number of completed hypothesis n (ie, where the hypothesis has produced)

⚠ Problem ⚠

- As our hypothesis get longer, their probabilities get smaller, and thus their scores get smaller

→ Need to normalize these log-scores by the length of the sequence

Advantages of NMT

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

47

Disadvantages of NMT

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!
 - Can't create rules such as:
"I always want to translate this word A by A' "
-

Evaluation of Machine Translation

BLEU (Bilingual Evaluation)

- Compare Machine Translation with several Human translations
- Compute a similarity score:
 - n-gram precision:
"For all the n-grams that appear in the Machine Translation, how many actually appear in one of the Human translations?"
 - Penalty score to penalize for too-short translations
(which could try to game the n-gram precision!)

Reward translations that have a high n-gram overlap with Human translations

⚠ Problem ⚠

There can be several right translations, we may be giving accurate ones bad scores !

Looking back at the past

- NMT managed to outperform SMT rapidly, with only a few engineers compared to SMT !!

Remaining Difficulties

- How to translate out-of-vocabulary words?
- Domain mismatch: formal \leftrightarrow informal
- Maintaining context over long text, without getting too expensive computationally
- Low-resource language pairs (languages with not much parallel data available online)
For low-resource languages, one of the best sources of parallel text is the Bible...
(ex of gibberish, translated as a Biblical sentence, from Somali to English)
- Still a lack of common sense:
There can't be "jam of paper" !!
- Keeps biases of the training data:
(Nurse \rightarrow She) & (Programmer \rightarrow He)
- Uninterpretability:
Really weird effects can happen, and it's difficult to understand why !

Attention

Why Attention?

- Our whole source sentence has been captured as the last hidden state of the Encoder RNN

Thus, we have lost all the information that is not stored there

Plus, it may suffer from a recency bias

\rightarrow Information **bottleneck** !

What is Attention?

Idea

On each step of the decoder, use direct connection to the Encoder to focus on a particular part of the source sequence

- We take the dot products between the decoder hidden state, and all the Encoder hidden states

→ We get **attention scores**

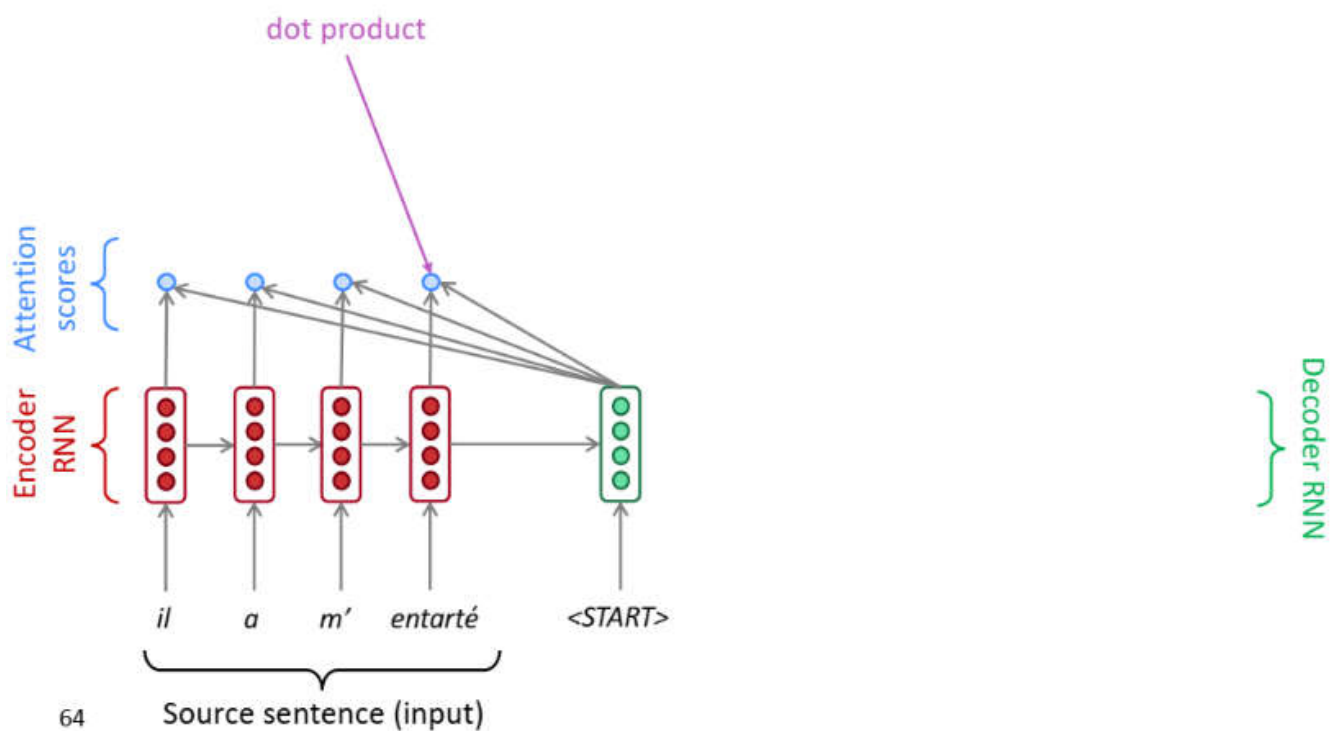
- We apply softmax to all these attention scores to get a probability distribution to know which parts of the encoded sentence are most relevant to our Decoder hidden state
→ **Attention Distribution**

- Use the Attention Distribution to make a weighted score of the Encoder hidden states
→ **Attention Output**
It mostly contains information from the Hidden states with high Attention !

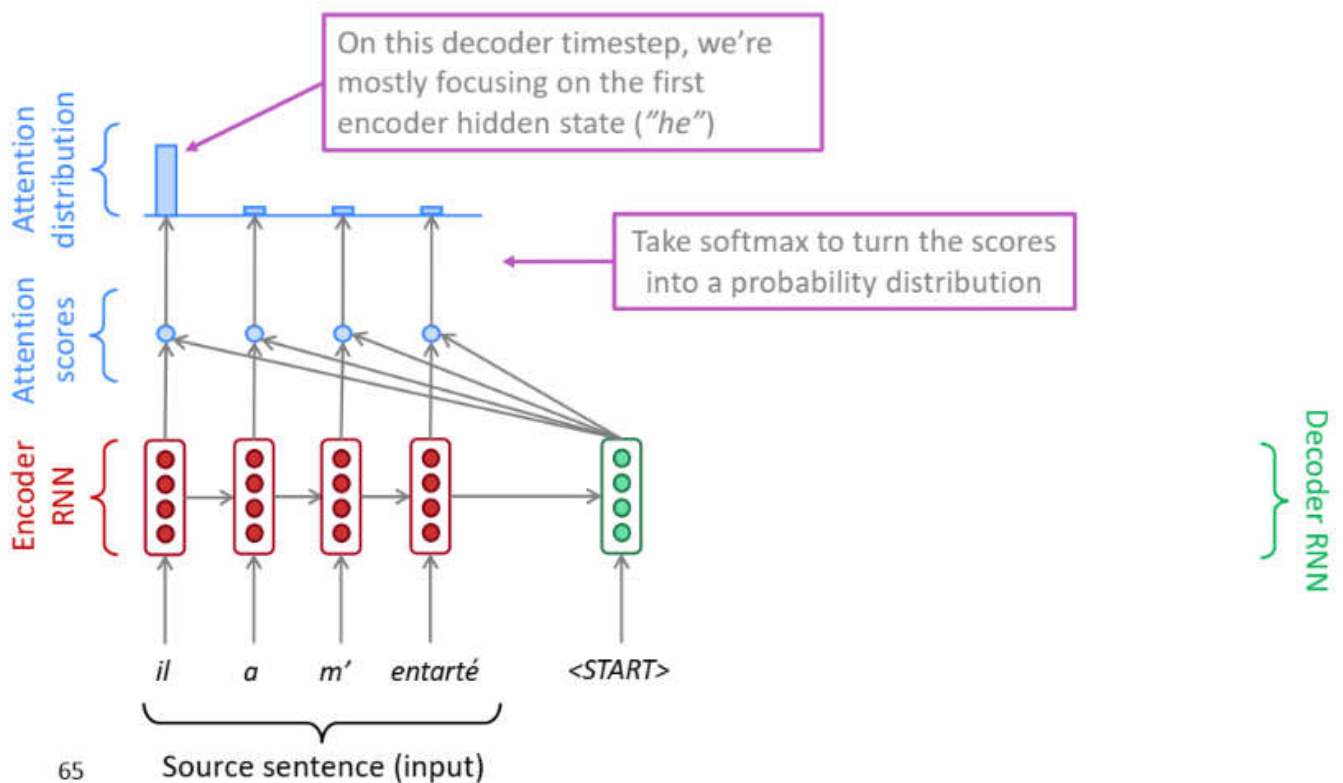
- Concatenate this Attention Output with Decoder Hidden State, to compute our output

Example

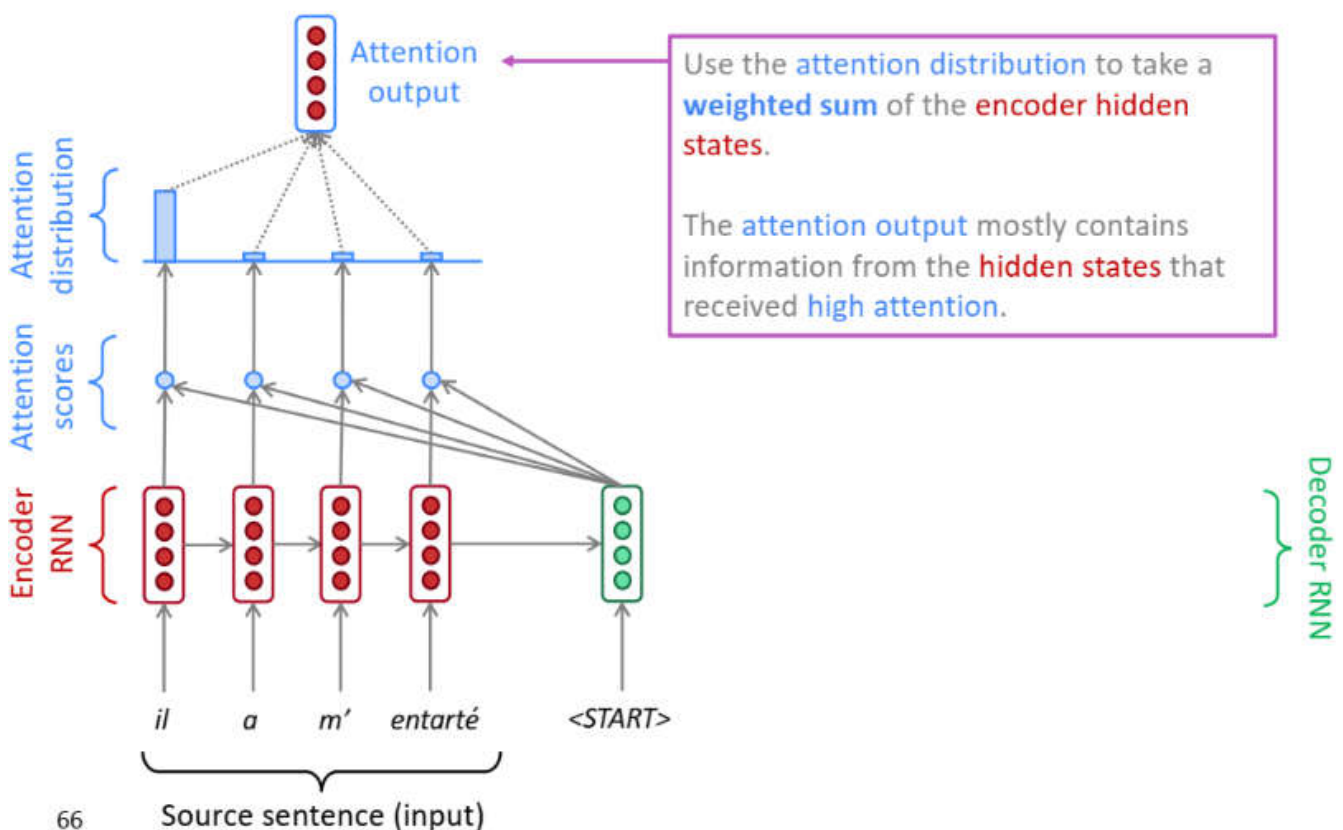
Sequence-to-sequence with attention



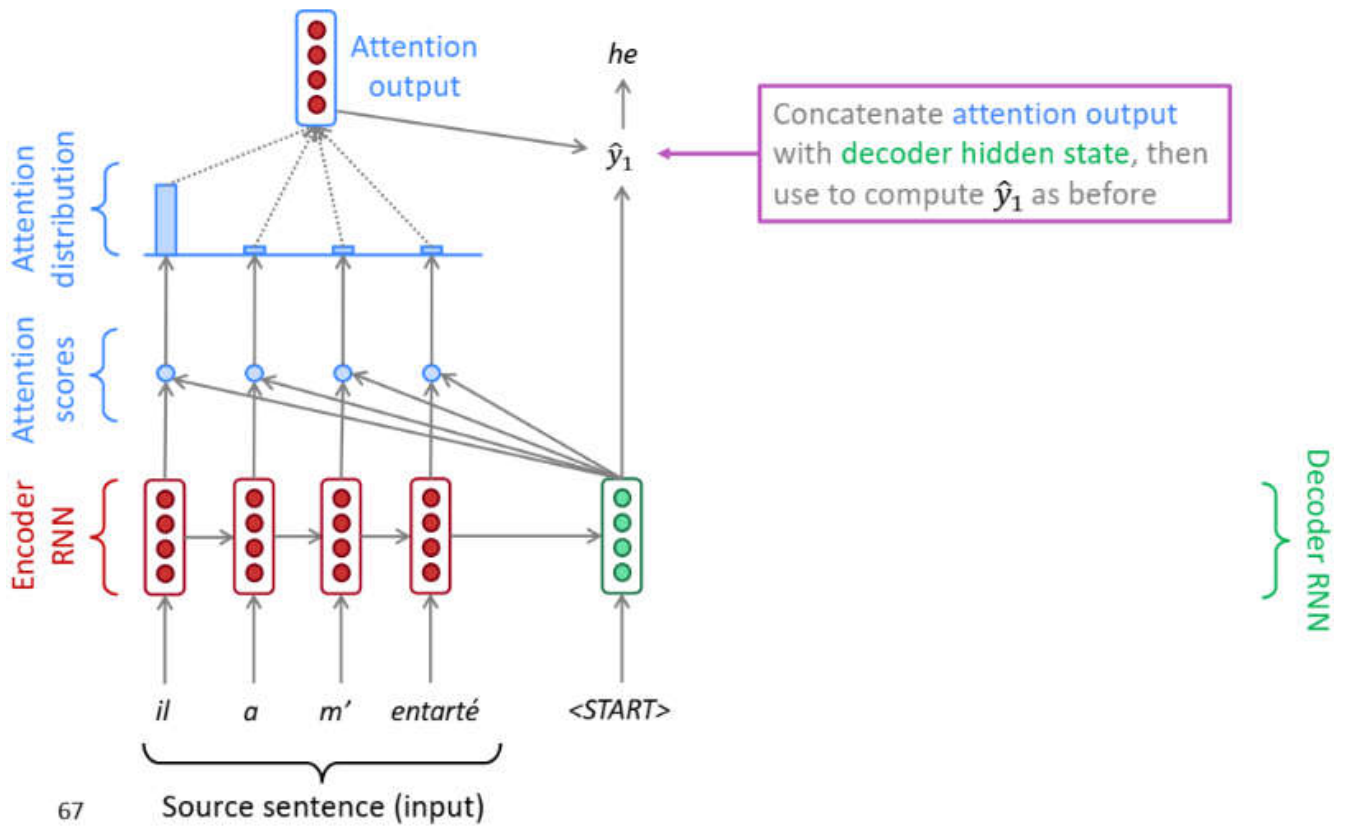
Sequence-to-sequence with attention



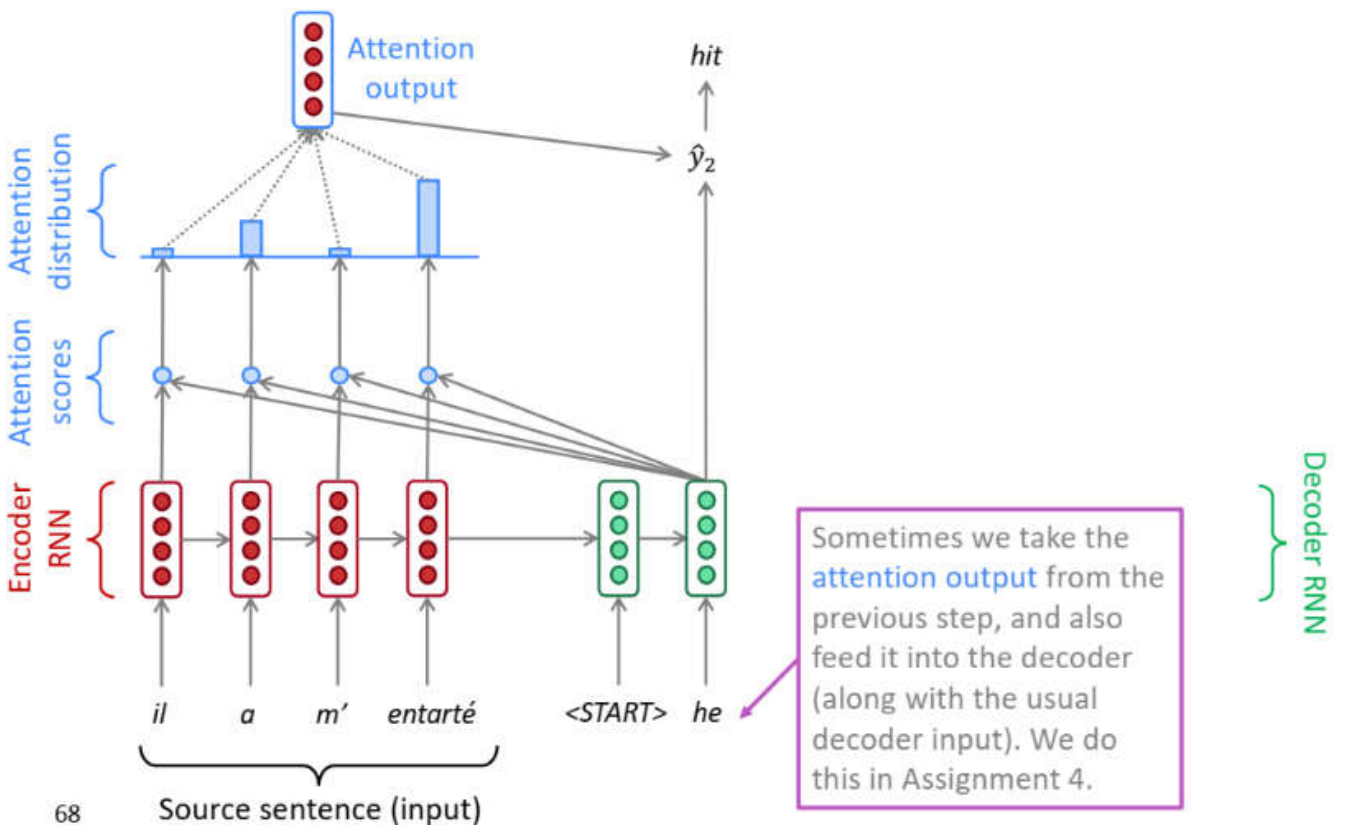
Sequence-to-sequence with attention



Sequence-to-sequence with attention



Sequence-to-sequence with attention



We can see here that Attention is a bit like a softer version of Alignment, which was much restrictive (binary)

Attention: in equations

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

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

- We take softmax to get the attention distribution α^t for this step (this is a probability distribution and sums to 1)

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

- We use α^t to take a weighted sum of the encoder hidden states to get the attention output a_t

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

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

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

73

Attention Advantages

- Attention allows Decoder to focus on only some parts of the source text
- Attention bypasses the bottleneck:
Direct access to the source !
- The shortcut (skipping connections) helps with the **vanishing gradient** problem !!
- Provides some **interpretability** (~soft alignment)

	he	hit	me	with	a	pie
il						
a						
m'						
entarté						

Other Attention uses

Attention can also be used for other Deep Learning tasks

Attention needs:

- Vector values
- Vector query

→ Compute a weighted sum of values, which is dependent on the query

- The weighted sum is a selective summary of the values information
- The query determines which values we focus on in our weighted sum

~ Similar to LSTM Gates, which depend on the context

- Attention is a way to get a **fixed-size** representation of any number of values

Variants

There are *several* attention variants

- We have some *values* $h_1, \dots, h_N \in \mathbb{R}^{d_1}$ and a *query* $s \in \mathbb{R}^{d_2}$

- Attention always involves:

1. Computing the *attention scores* $e \in \mathbb{R}^N$
2. Taking softmax to get *attention distribution* α :

There are multiple ways to do this

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

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

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

thus obtaining the *attention output* a (sometimes called the *context vector*)

Attention variants

You'll think about the relative advantages/disadvantages of these in Assignment 4!

There are **several ways** you can compute $e \in \mathbb{R}^N$ from $h_1, \dots, h_N \in \mathbb{R}^{d_1}$ and $s \in \mathbb{R}^{d_2}$:

- Basic dot-product attention: $e_i = s^T h_i \in \mathbb{R}$
 - Note: this assumes $d_1 = d_2$
 - This is the version we saw earlier
- Multiplicative attention: $e_i = s^T W h_i \in \mathbb{R}$
 - Where $W \in \mathbb{R}^{d_2 \times d_1}$ is a weight matrix
- Additive attention: $e_i = v^T \tanh(W_1 h_i + W_2 s) \in \mathbb{R}$
 - Where $W_1 \in \mathbb{R}^{d_3 \times d_1}$, $W_2 \in \mathbb{R}^{d_3 \times d_2}$ are weight matrices and $v \in \mathbb{R}^{d_3}$ 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>

78

Lecture Summary

Summary of today's lecture

- We learned some history of Machine Translation (MT)
- Since 2014, **Neural MT** rapidly replaced intricate Statistical MT
- **Sequence-to-sequence** is the architecture for NMT (uses 2 RNNs)
- **Attention** is a way to *focus on particular parts* of the input
 - Improves sequence-to-sequence a lot!

