# Introduction to Natural Language Processing

## Goals, Main Problems

# Sequence to Sequence Modeling

# Sequence-to-Sequence Models

Example usages of Seq2Seq architecture:

- Machine Translation (sentence in A lang -> sentence in B lang)
- Text Summarization (long text -> short text)
- Dialogue (previous utterances -> next utterance)
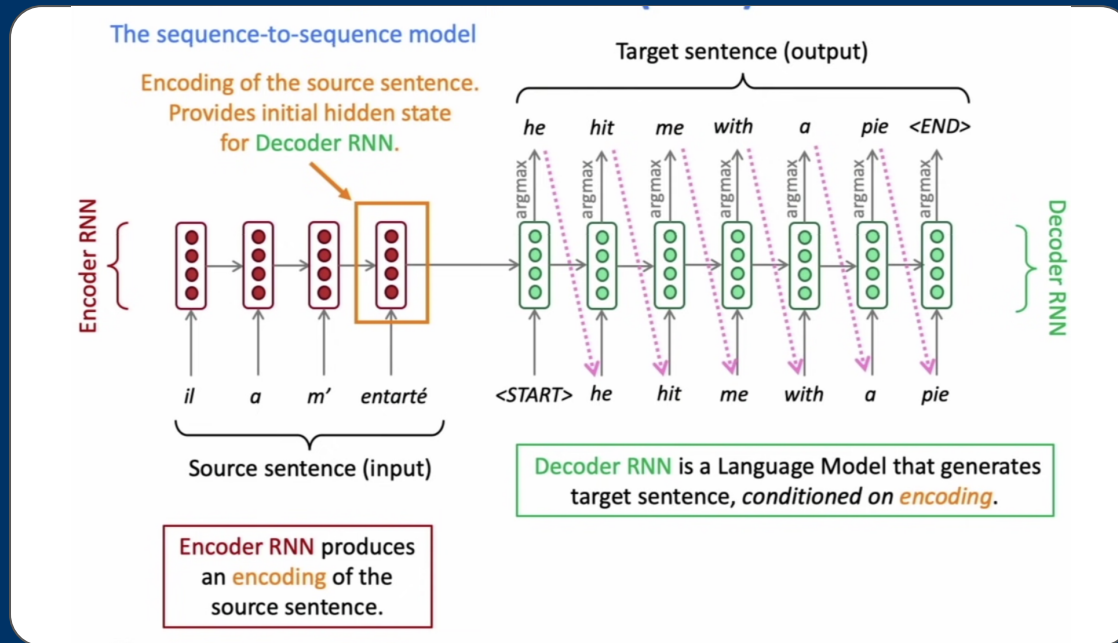- Code Generation (natural language -> Python code)

Seq2Seq architectures involve two RNNs. One is for encoding a sequence and the other one for decoding.

We will consider the Neural Machine Translation problem.

It is a way to do Machine Translation with a single NN.

# NMT

The Decoder is a conditional LM, which takes the encoded output from the Encoder and generates a new sequence.
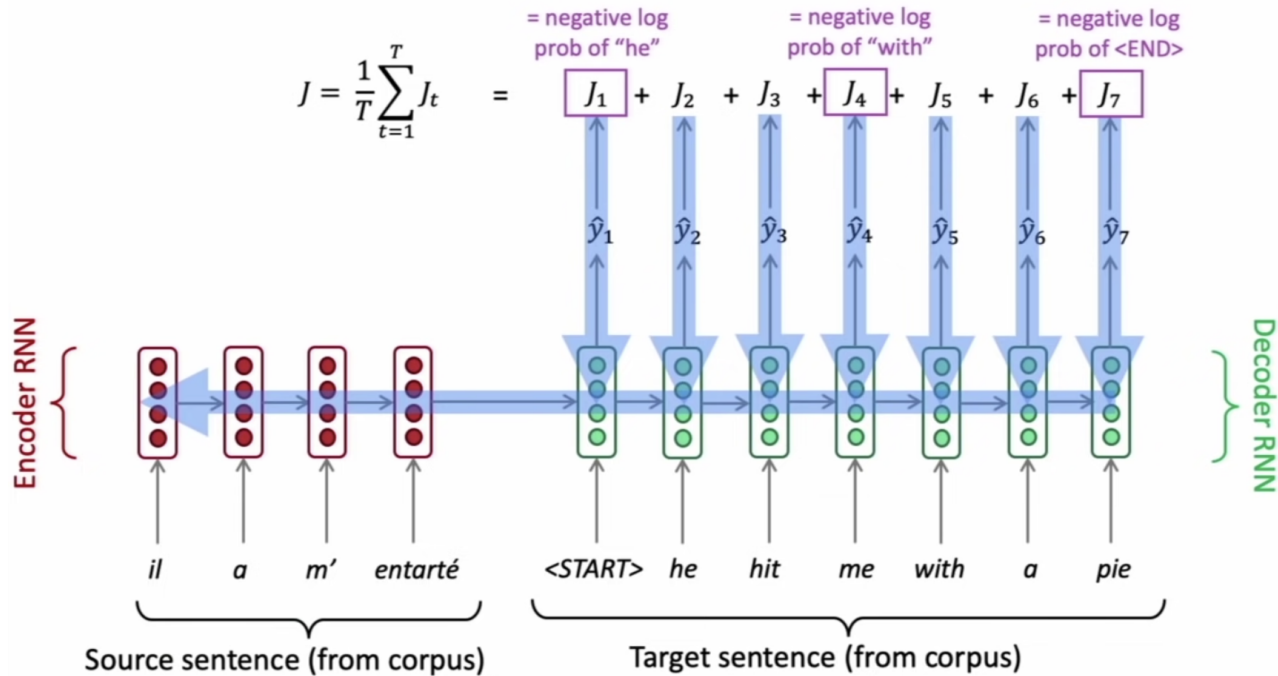
# NMT - task formulation

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

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

So how to train this system?

Get a Big parallel corpus and do the following:

# NMT - Training Architecture
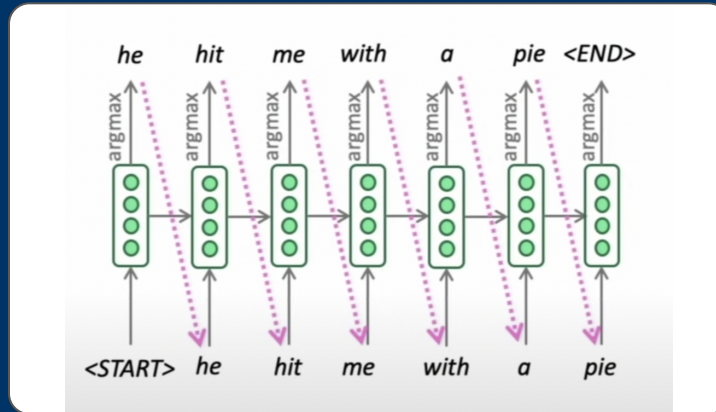
# Decoding Algorithms

How to generate/choose the most plausible output sentence?

Decoding strategies (most common):

- Greedy Decoding
- Beam Search

# Greedy Decoding

We saw how to generate (or "decode") the target sentence by taking argmax on each time step of the decoder.



This is **greedy decoding** - take most probable word on each step

# Exhaustive Search Decoding

Ideally we want to find a (length T) translation y that maximizes:

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

$$= \prod_{t=1}^{T} P(y_t|y_1,\ldots,y_{t-1},x)$$

We could try computing all possible sequences y

- This means that on each step t of the decoder, we're tracking $V^t$ possible partial translations, where V is vocab size.
- This $O(V^t)$ complexity is too expensive.

# Beam Search Decoding

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)

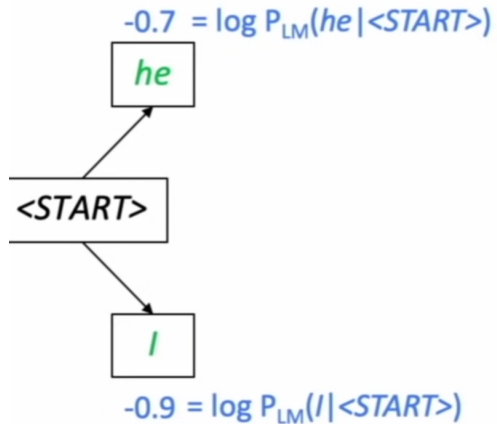A hypothesis $y_1,...,y_t$ has a score which is its log probability:

$$score(y_1, \ldots, y_t) = logP_{LM}(y_1, \ldots, y_t|x) = \sum_{i=1}^{t} logP_{LM}(y_i|y_1, \ldots, y_{i-1}, x)$$

The scores are all negative, and higher score is better

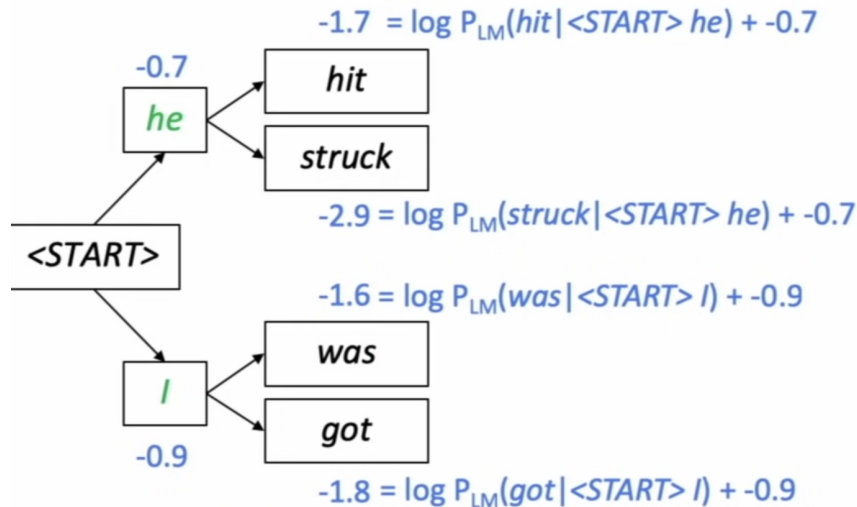We search for high-scoring hypotheses, tracking top k on each step.
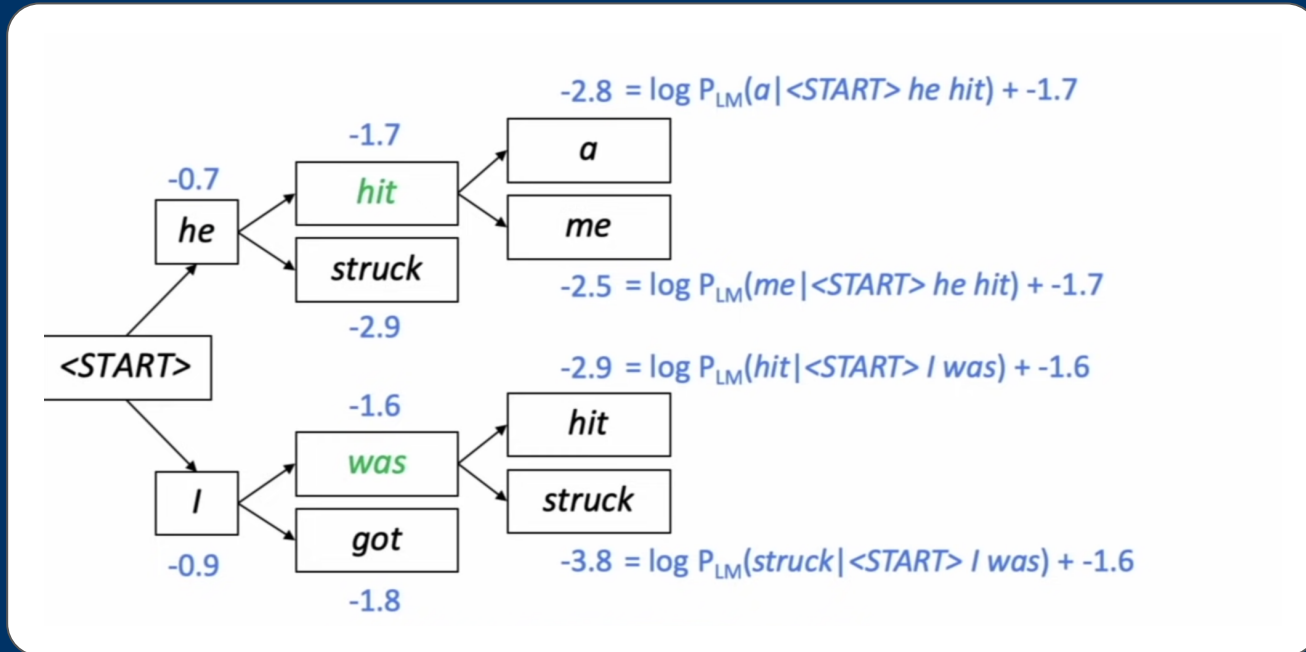
# Beam Search Decoding: Example

Beam size: k=2

# Beam Search Decoding: Example

Beam size: k=2



$-1.7 = \log P_{LM}(hit | \text{<START>} he) + -0.7$

$-0.7$

**he**

**hit**

**struck**

$-2.9 = \log P_{LM}(struck | \text{<START>} he) + -0.7$

**<START>**

$-1.6 = \log P_{LM}(was | \text{<START>} I) + -0.9$

**I**

**was**

**got**

$-0.9$

$-1.8 = \log P_{LM}(got | \text{<START>} I) + -0.9$
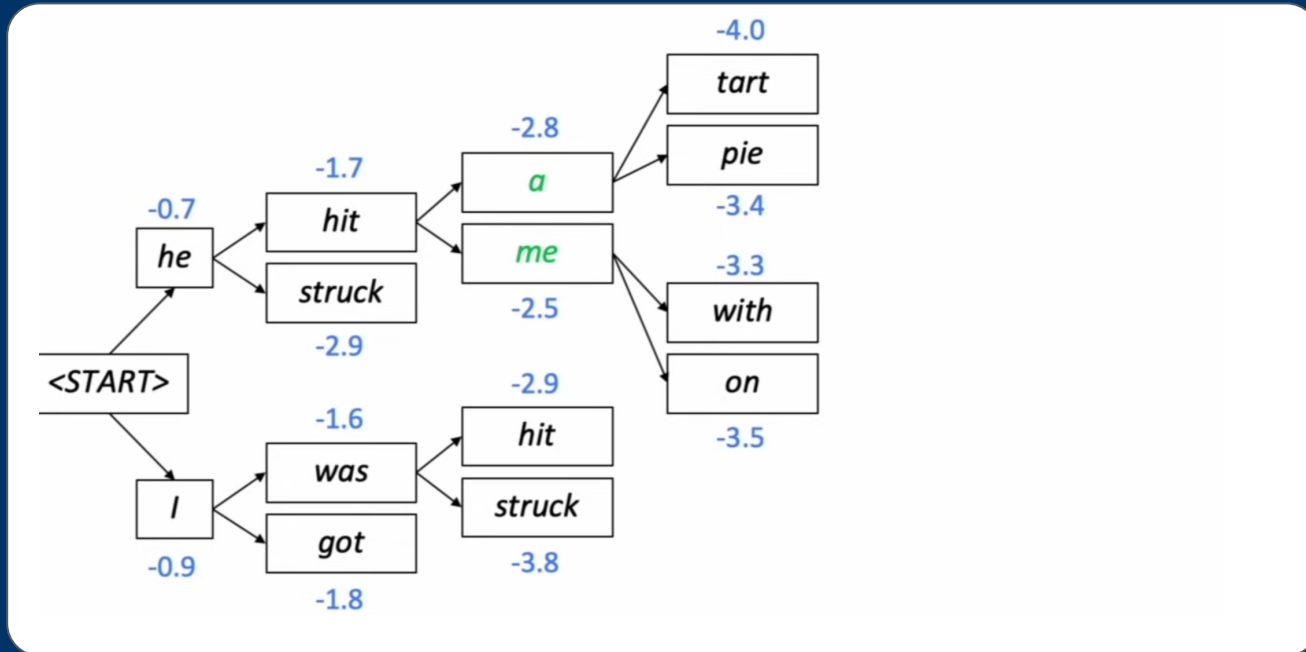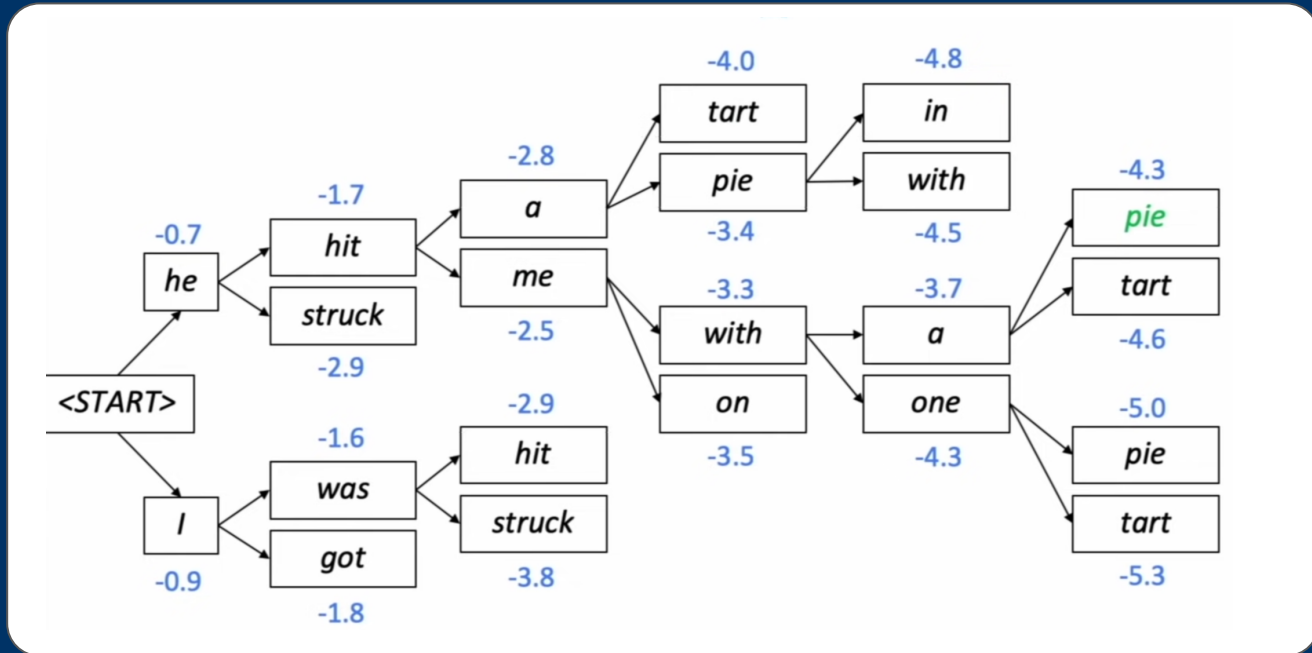
# Beam Search Decoding: Example

Beam size: k=2

# Beam Search Decoding: Example

Beam size: k=2
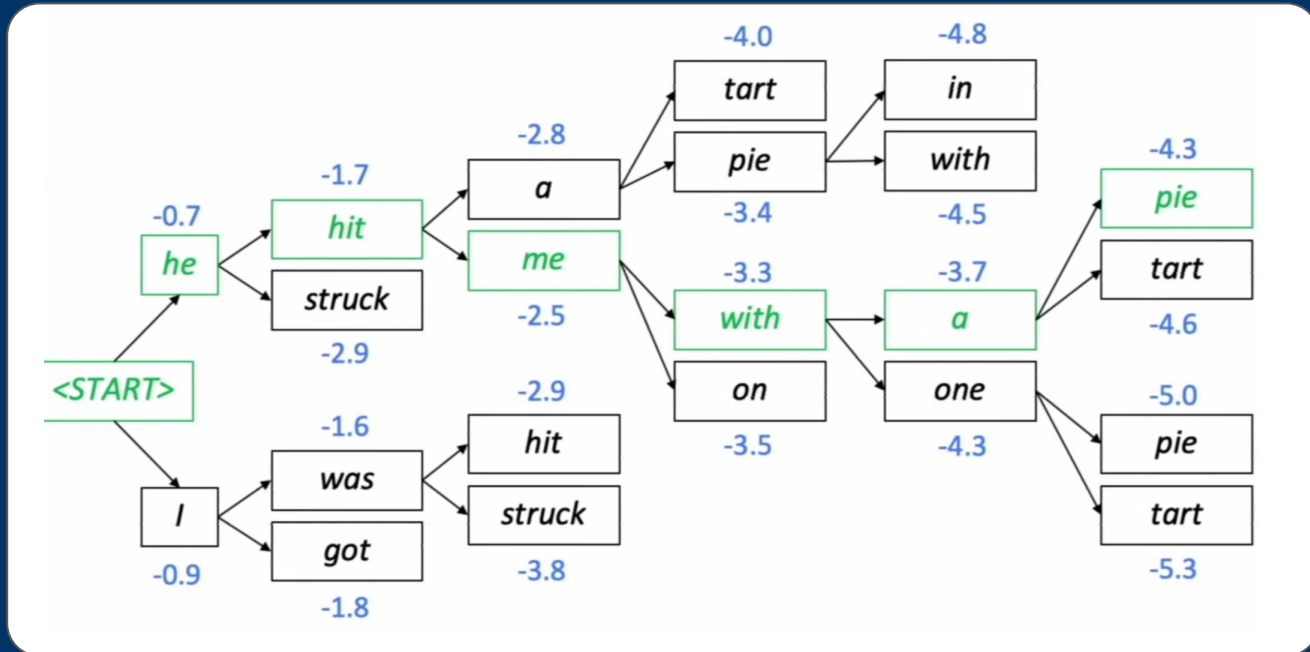
# Beam Search Decoding: Example

Beam size: k=2

# Beam Search Decoding: Example

Beam size: k=2

# Beam Search: Stopping Criterion

In **greedy decoding**, usually we decode until the model produces a <END> token.

In beam **search decoding**, different hypotheses may produce <END> token on different timesteps.

- When a hypothesis produces <END>, that hypothesis is complete.
- Place it aside and continue exploring other hypotheses via beam search.

Usually we continue beam search until:

- We reached timestep T (where T is some predefined cutoff), or
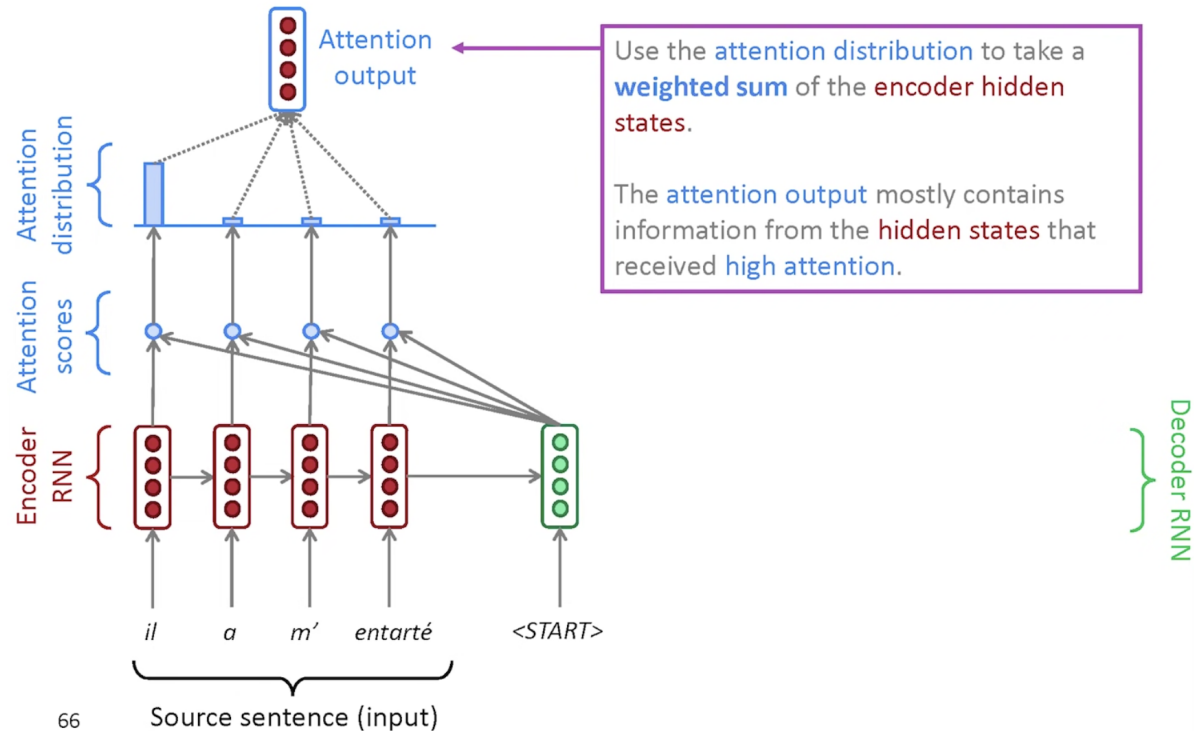- We have at least n completed hypothesis (where n i predefined cutoff)

# Attention

The NMT can be potentially improved with the mechanism of attention.
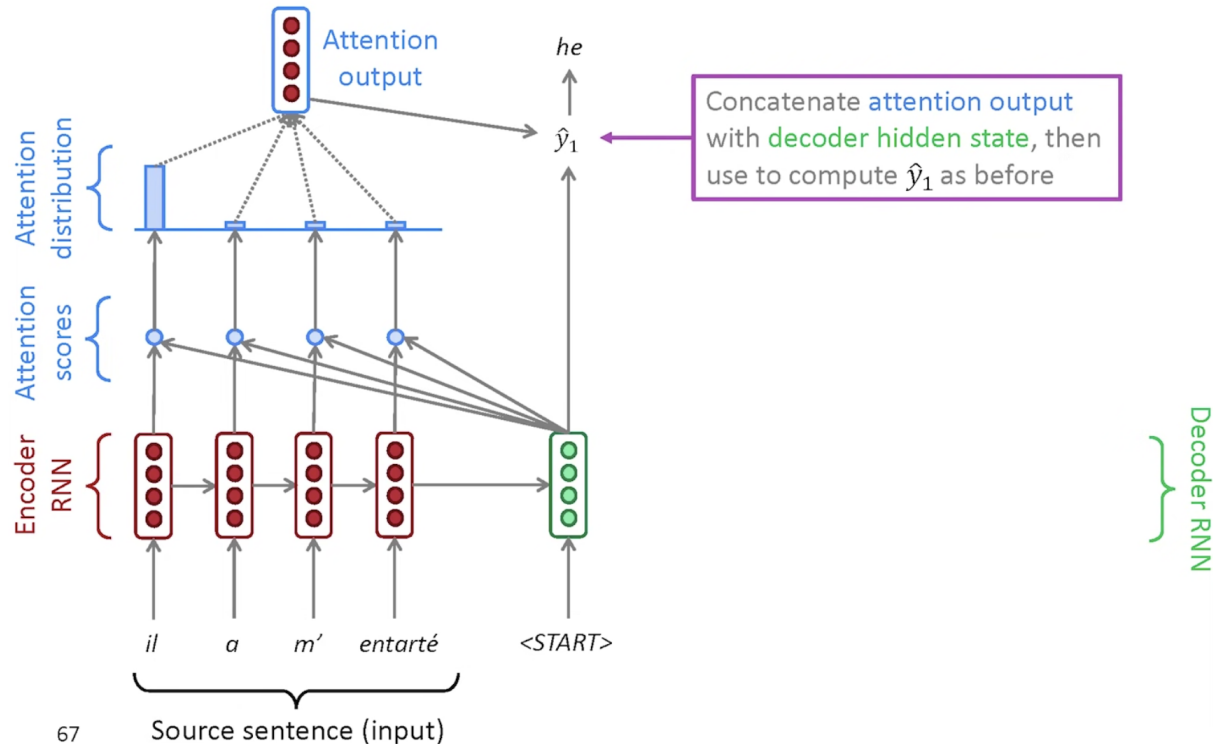
Attention provides a solution to the bottleneck problem, which is we are pushing all the information in the encoder state to a single vector.

Core idea: on each step of the decoder, use direct connection to the encoder to focus on a particular part of the source sentence.

# Attention: Example



Use the attention distribution to take a **weighted sum** of the encoder hidden states.

The attention output mostly contains information from the hidden states that received high attention.

Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>

Source sentence (input)

66

# Attention: Example



Concatenate attention output with decoder hidden state, then use to compute $\hat{y}_1$ as before

# Attention: Example



Attention output

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

$\hat{y}_5$

$a$

il    a    m'    entarté         <START>    he    hit    me    with

Source sentence (input)

71

# Attention

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

Thank You