# Natural Language Processing with Deep Learning
# CS224N/Ling284

Lecture 8:
Machine Translation,
Sequence-to-sequence and Attention

**Abigail See**

# Announcements

- We are taking attendance today
  - Sign in with the TAs outside the auditorium
  - No need to get up now – there will be plenty of time to sign in after the lecture ends
  - For attendance policy special cases, see Piazza post for clarification

- Assignment 4 content covered today
  - Get started early! The model takes 4 hours to train!

- Mid-quarter feedback survey:
  - Will be sent out sometime in the next few days (watch Piazza).
  - Complete it for 0.5% credit

# Overview

Today we will:

- Introduce a <u>new task</u>: Machine Translation

  **is a major use-case of**

- Introduce a <u>new neural architecture</u>: sequence-to-sequence

  **is improved by**

- Introduce a <u>new neural technique</u>: attention

# Section 1: Pre-Neural Machine Translation

history

# 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).

*x:*        *L'homme est né libre, et partout il est dans les fers*

*y:*        *Man is born free, but everywhere he is in chains*

- Rousseau

# 1950s: Early Machine Translation

Machine Translation research began in the early 1950s.

- Russian → English (motivated by the Cold War!)

    AI hype



**1 minute video showing 1954 MT:**
https://youtu.be/K-HfpsHPmvw

- Systems were mostly rule-based, using a bilingual dictionary to map Russian words to their English counterparts

# 1990s-2010s: Statistical Machine Translation

- Core idea: Learn a probabilistic model from data

- Suppose we're translating French → English.

- We want to find best English sentence *y,* given French sentence *x*

target                                                    source

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

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

division of labor                    y|x

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

local

**Translation Model**

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

**Language Model**

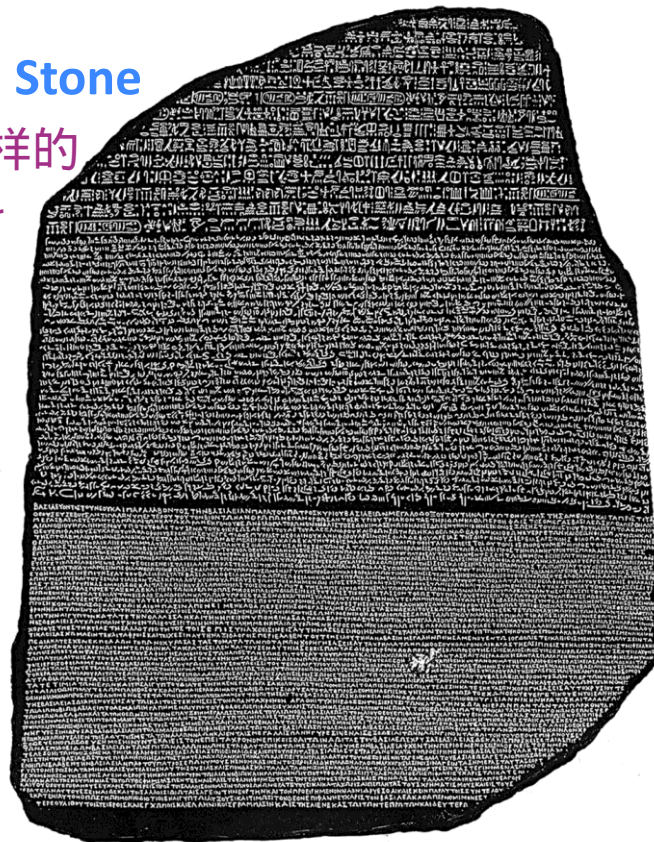Models how to write good English (*fluency*). Learnt from monolingual data.

# 1990s-2010s: Statistical Machine Translation

- Question: How to learn translation model $P(x|y)$ ?
- First, need large amount of parallel data
  (e.g. pairs of human-translated French/English sentences)

**The Rosetta Stone**

3
——

Ancient Egyptian

Demotic

Ancient Greek

# Learning alignment for SMT

- <u>Question:</u> How to learn translation model $P(x|y)$ from the parallel corpus?
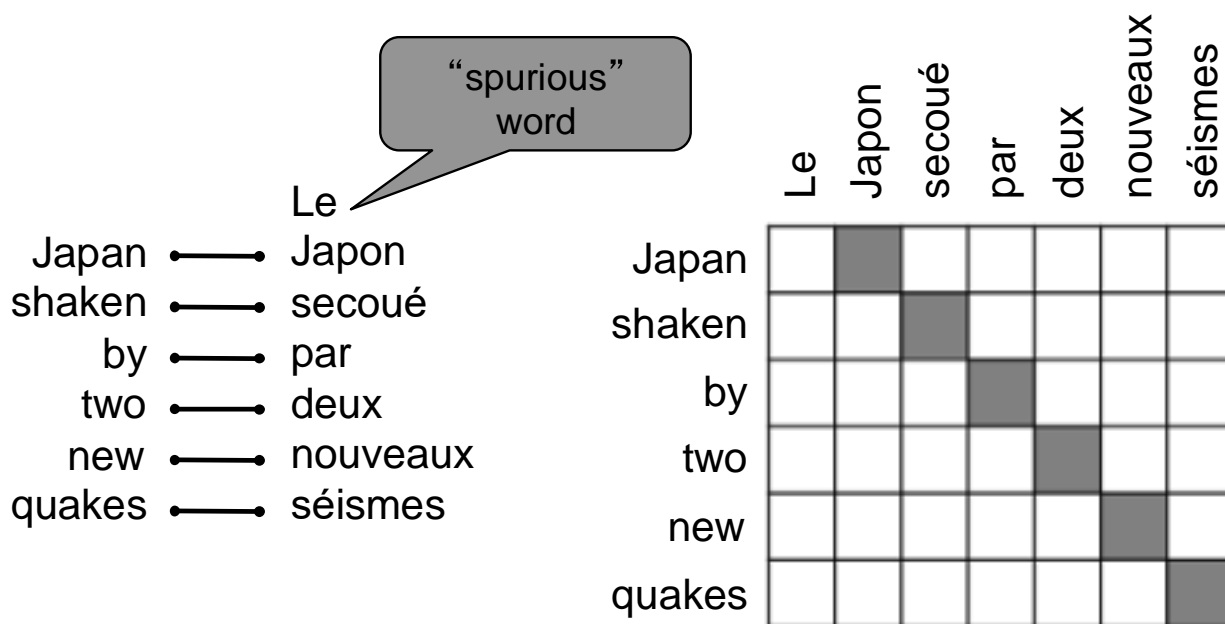
- Break it down further: we actually want to consider

$$P(x, a|y)$$

  where *a* is the alignment, i.e. word-level correspondence between French sentence *x* and English sentence *y*
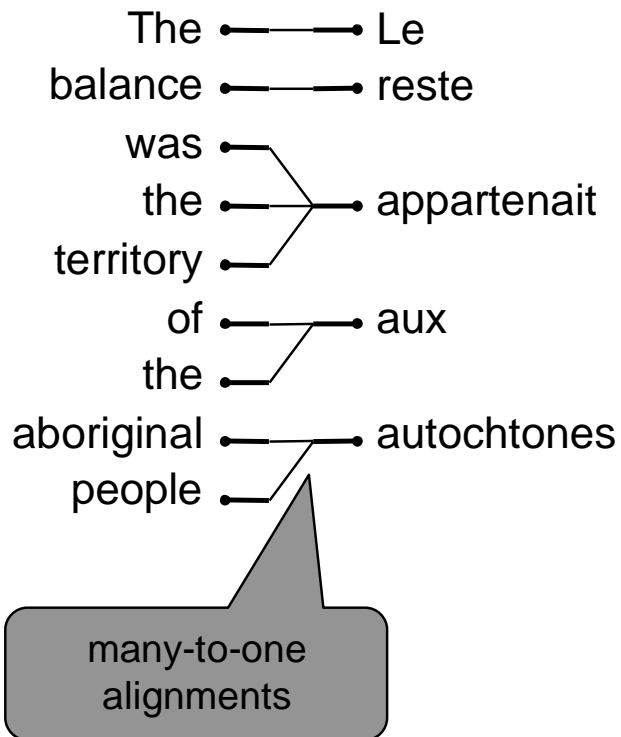
# What is alignment?

Alignment is the correspondence between particular words in the translated sentence pair.

- Note: Some words have no counterpart



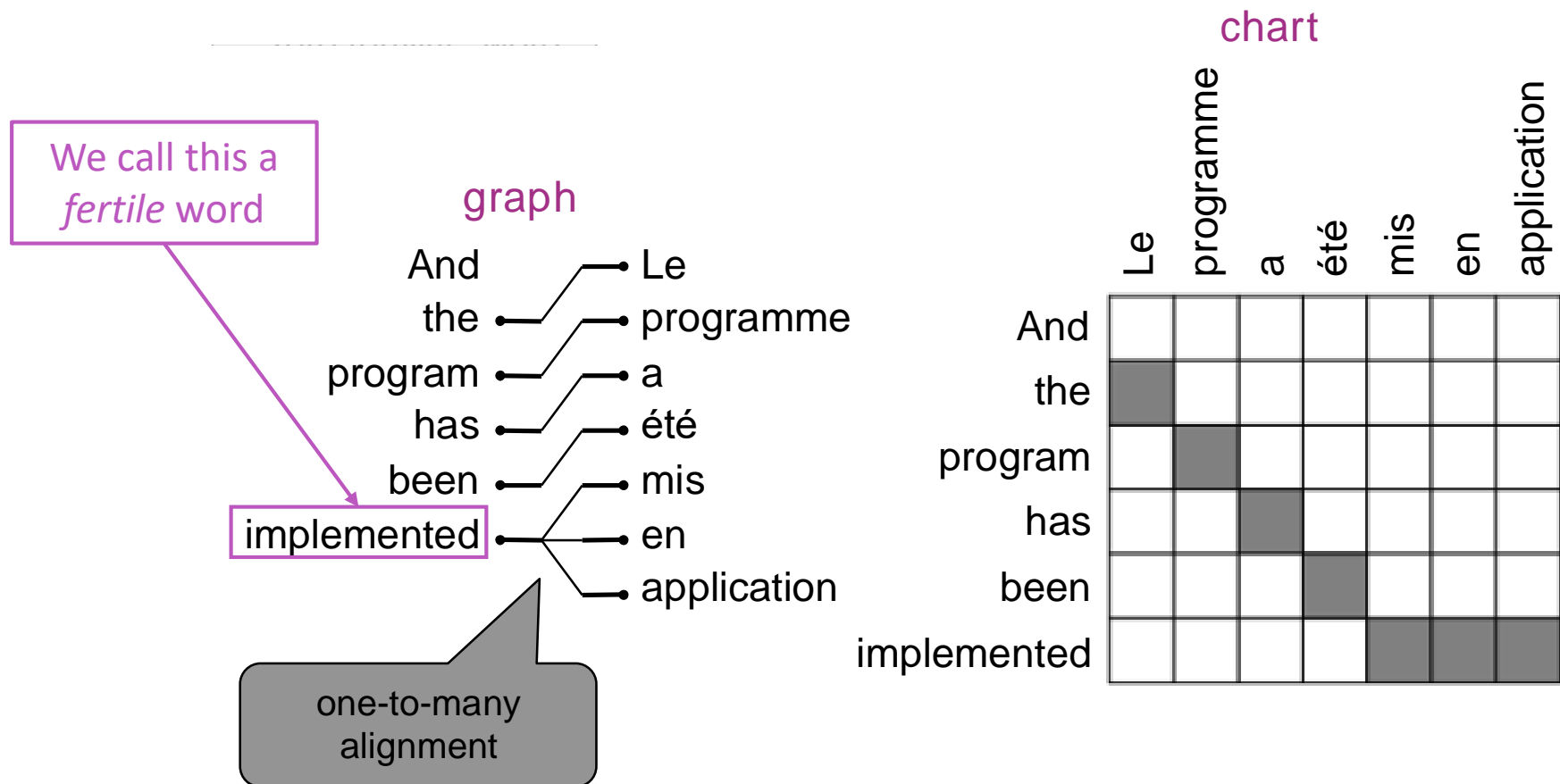**Examples from:** "The Mathematics of Statistical Machine Translation: Parameter Estimation", Brown et al, 1993. http://www.aclweb.org/anthology/J93-2003

# Alignment is complex

Alignment can be many-to-one



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

many-to-one
alignments

**Examples from:** "The Mathematics of Statistical Machine Translation: Parameter Estimation", Brown et al, 1993. http://www.aclweb.org/anthology/J93-2003

# Alignment is complex

Alignment can be one-to-many

We call this a *fertile* word

graph

And ——— Le
the ——— programme
program ——— a
has ——— été
been ——— mis
implemented ——— en
——— application

one-to-many alignment

chart

|  | Le | programme | a | été | mis | en | application |
|---|---|---|---|---|---|---|---|
| And |  |  |  |  |  |  |  |
| the | ▓ |  |  |  |  |  |  |
| program |  | ▓ |  |  |  |  |  |
| has |  |  | ▓ |  |  |  |  |
| been |  |  |  | ▓ |  |  |  |
| implemented |  |  |  |  | ▓ | ▓ | ▓ |

# Alignment is complex

Some words are very fertile!



il ●————————————● he

a ————————————● hit

m' ●————————————● me

entarté ●————————————● with

● a

● pie

This word has no single-word equivalent in English

|         | he | hit | me | with | a | pie |
|---------|----|-----|----|----|----|----|
| il      | ■  |     |    |    |   |    |
| a       |    |     |    |    |   |    |
| m'      |    |     | ■  |    |   |    |
| entarté |    | ■   |    | ■  | ■ | ■  |

# Alignment is complex

Alignment can be many-to-many (phrase-level)

The — Les
poor — pauvres
don't — sont
have — démunis
any
money

many-to-many alignment

|  | Les | pauvres | sont | démunis |
|---|---|---|---|---|
| The | ■ |  |  |  |
| poor |  | ■ |  |  |
| don't |  |  | ■ | ■ |
| have |  |  | ■ | ■ |
| any |  |  | ■ | ■ |
| money |  |  | ■ | ■ |

phrase alignment

**Examples from:** "The Mathematics of Statistical Machine Translation: Parameter Estimation", Brown et al, 1993. http://www.aclweb.org/anthology/J93-2003

# Learning alignment for SMT

- We learn $P(x, a|y)$ as a combination of many factors, including:

  - Probability of particular words aligning (also depends on position in sent)                    align

  - Probability of particular words having particular fertility (number of corresponding words)

  - etc.
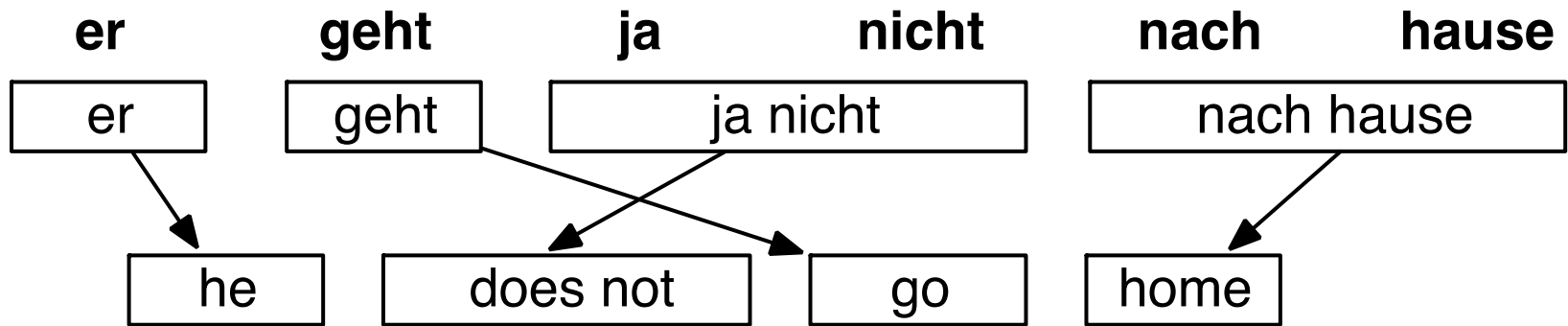
# Decoding for SMT

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

**Question:**
How to compute
this argmax?

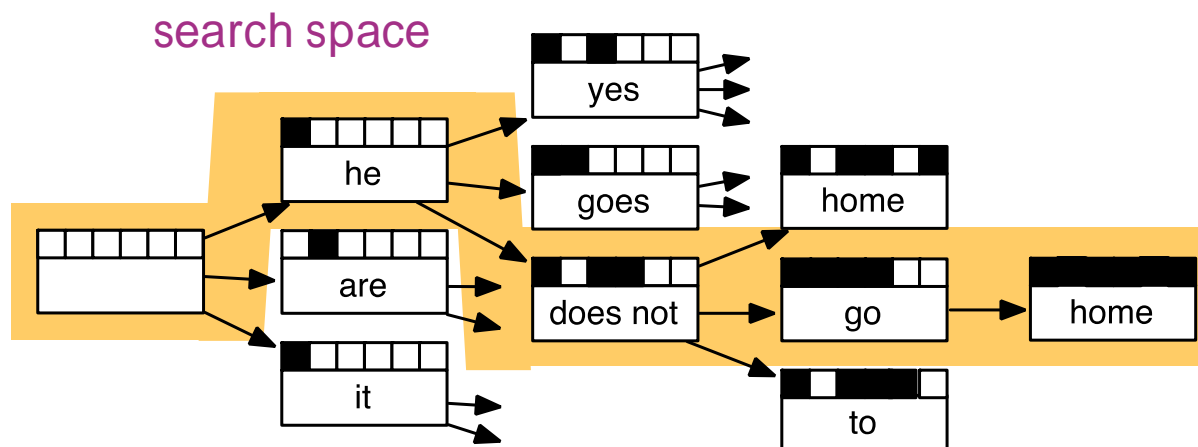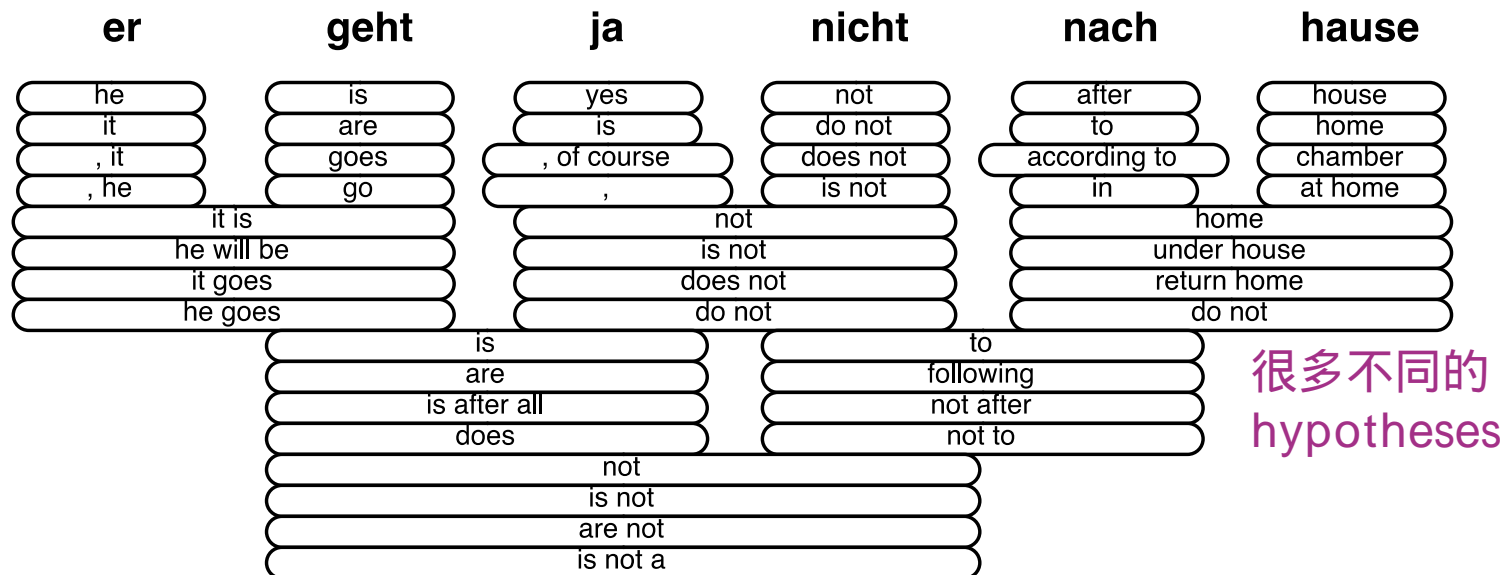Translation Model

Language Model

- We could enumerate every possible *y* and calculate the probability?  →  Too expensive!

- **Answer:** Use a heuristic search algorithm to search for the best translation, discarding hypotheses that are too low-probability

- This process is called *decoding*

# Decoding for SMT

# Decoding for SMT

| er | geht | ja | nicht | nach | hause |
|----|------|-----|-------|------|-------|
| he | is | yes | not | after | house |
| it | are | is | do not | to | home |
| , it | goes | , of course | does not | according to | chamber |
| , he | go | , | is not | in | at home |

| it is | | not | | home |
| he will be | | is not | | under house |
| it goes | | does not | | return home |
| he goes | | do not | | do not |

| is | to |
| are | following |
| is after all | not after |
| does | not to |

| not |
| is not |
| are not |
| is not a |

hypotheses

search space



**Source:** "Statistical Machine Translation", Chapter 6, Koehn, 2009.
https://www.cambridge.org/core/books/statistical-machine-translation/94EADF9F680558E13BE759997553CDE5

# 1990s-2010s: Statistical Machine Translation

- SMT was a huge research field
- The best systems were extremely complex
    - Hundreds of important details we haven't mentioned here
    - Systems had 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            phrase
    - Lots of human effort to maintain
        - Repeated effort for each language pair!

# Section 2: Neural Machine Translation

**2014**

**(dramatic reenactment)**

2014

Neural Machine Translation

MT research

(dramatic reenactment)

# What is 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.
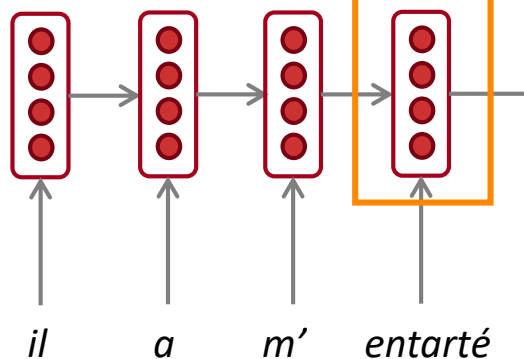
# Neural Machine Translation (NMT)

The sequence-to-sequence model

embedding matrix

Target sentence (output)

Encoding of the source sentence.
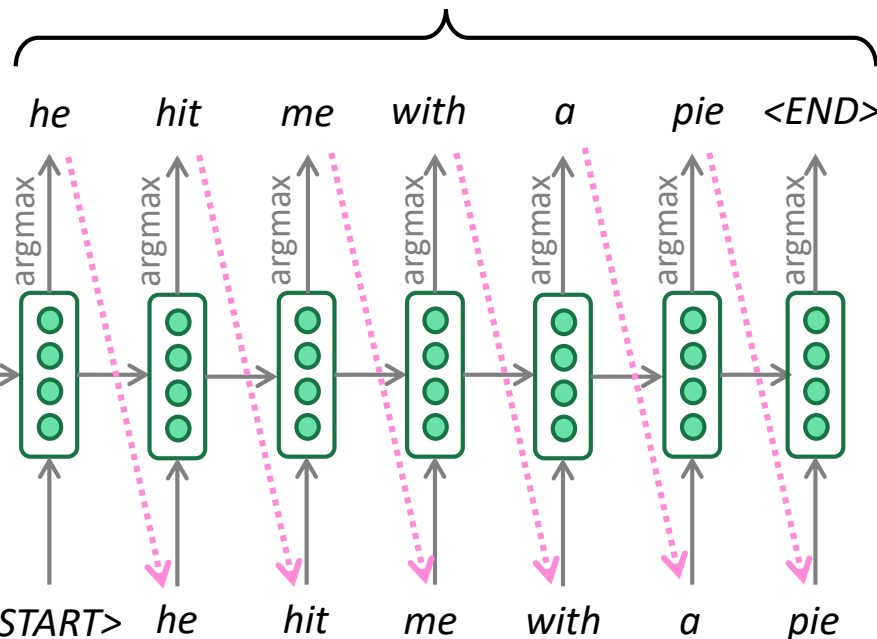Provides initial hidden state
for Decoder RNN.

RNN



Encoder RNN

Decoder RNN

he    hit    me    with    a    pie    <END>

argmax  argmax  argmax  argmax  argmax  argmax  argmax

il    a    m'    entarté

<START>  he    hit    me    with    a    pie

Source sentence (input)

conditional language model

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

Encoder RNN produces
an encoding of the
source sentence.

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

test          training

24

# 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)

# 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 $P(y|x)$ :   SMT

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

RNN      timestep

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

- **Question**: How to train a NMT system?

- **Answer**: Get a big parallel corpus…

26

# Training a Neural Machine Translation system

padding

pad

hidden states

= negative log prob of "he"

= negative log prob of "with"

= negative log prob of <END>

training

timestep

labeled word

test

<END>

state

$$J = \frac{1}{T}\sum_{t=1}^{T} J_t \quad = \quad \boxed{J_1} + J_2 + J_3 + \boxed{J_4} + J_5 + J_6 + \boxed{J_7}$$

output loss

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$   $\hat{y}_4$   $\hat{y}_5$   $\hat{y}_6$   $\hat{y}_7$

Encoder RNN

Decoder RNN

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

Source sentence (from corpus)

Target sentence (from corpus)

end-to-end   train

pre-train

loss

encoder

Seq2seq is optimized as a **single system.**
Backpropagation operates "*end-to-end*".

system

# 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?** argmax

# Problems with greedy decoding

- Greedy decoding has no way to undo decisions!
    - Input: *il a m'entarté*     *(he hit me with a pie)*
    - → *he ____*
    - → *he hit ____*
    - → *he hit a ____*     (whoops! no going back now…)

- How to fix this?       beam search

# 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 far too expensive!

# 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 $y_1, \cdots, y_t$ has a score which is its log probability:

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

  - Scores are all negative, and higher score is better  prob
  - 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

# Beam search decoding: example

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

*<START>*

Calculate prob
dist of next word

32

# Beam search decoding: example

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

hypotheses

-0.7 = log P$_{\mathrm{LM}}$(*he*|*<START>*)

| he |

| *<START>* |

| I |

-0.9 = log P$_{\mathrm{LM}}$(*I*|*<START>*)

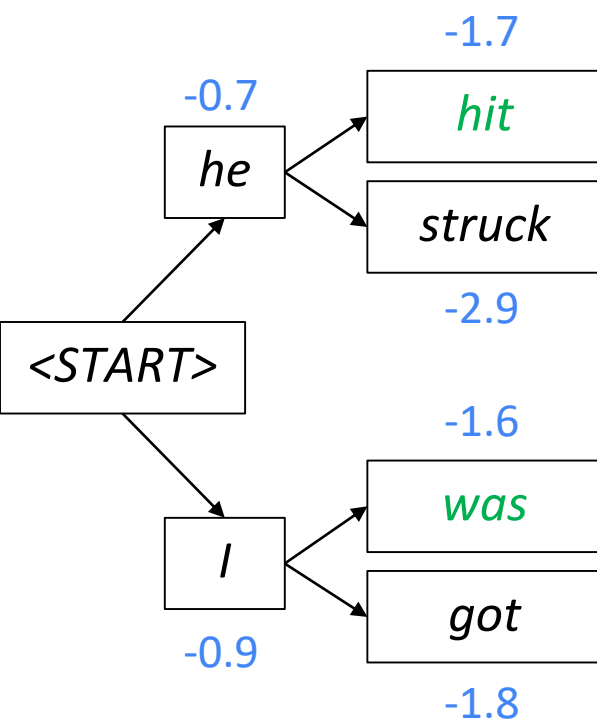Take top *k* words
and compute scores

# Beam search decoding: example

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

-1.7 = log P$_{LM}$(*hit*|*<START> he*) + -0.7

-0.7

*he*

*hit*

*struck*

-2.9 = log P$_{LM}$(*struck*|*<START> he*) + -0.7

*<START>*

-1.6 = log P$_{LM}$(*was*|*<START> I*) + -0.9

*was*

*I*

*got*

-0.9

-1.8 = log P$_{LM}$(*got*|*<START> I*) + -0.9

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

34

# Beam search decoding: example

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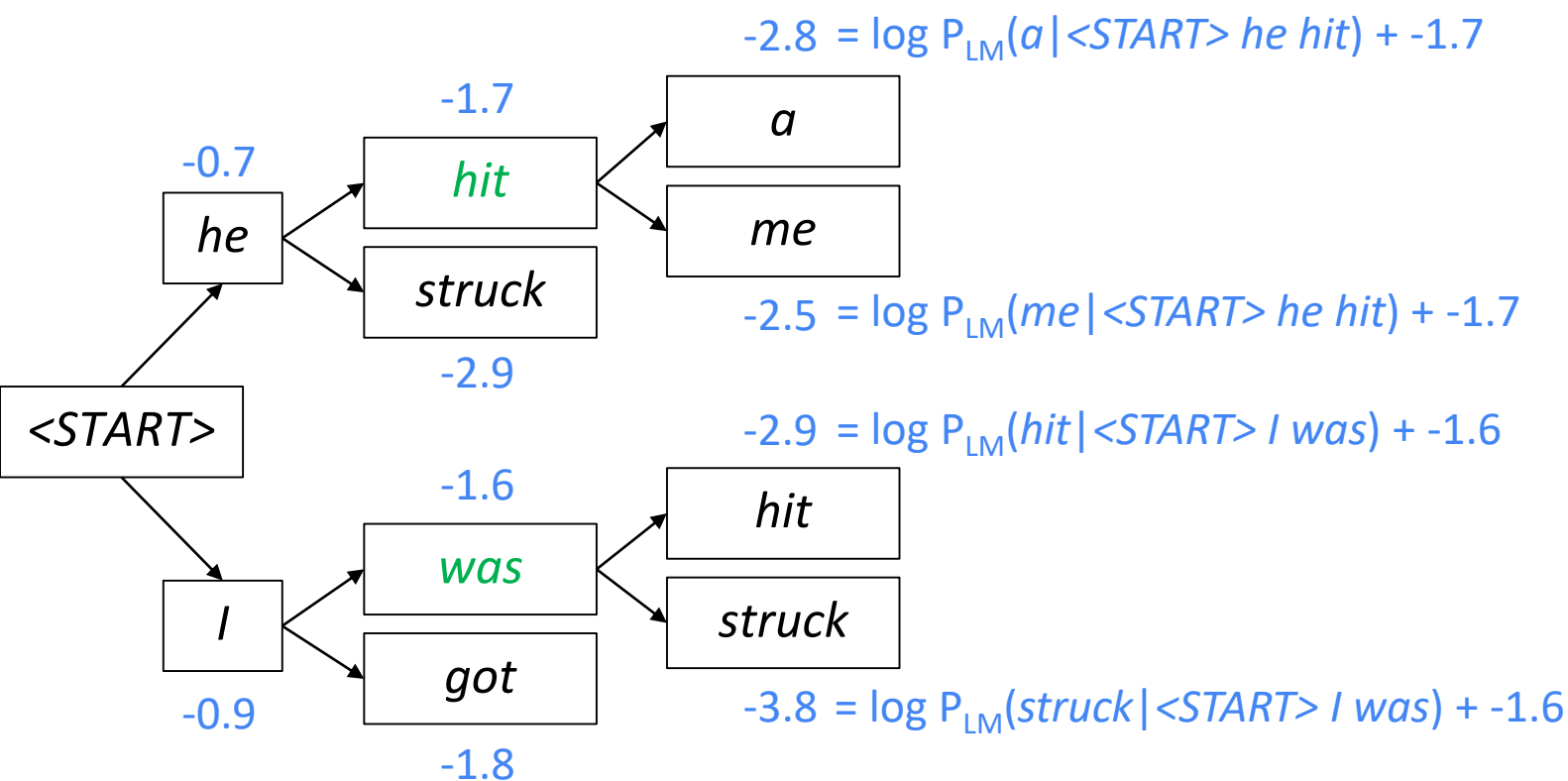


-1.7
-0.7
hit
he
struck
<START>
-2.9
-1.6
was
I
got
-0.9
-1.8

k

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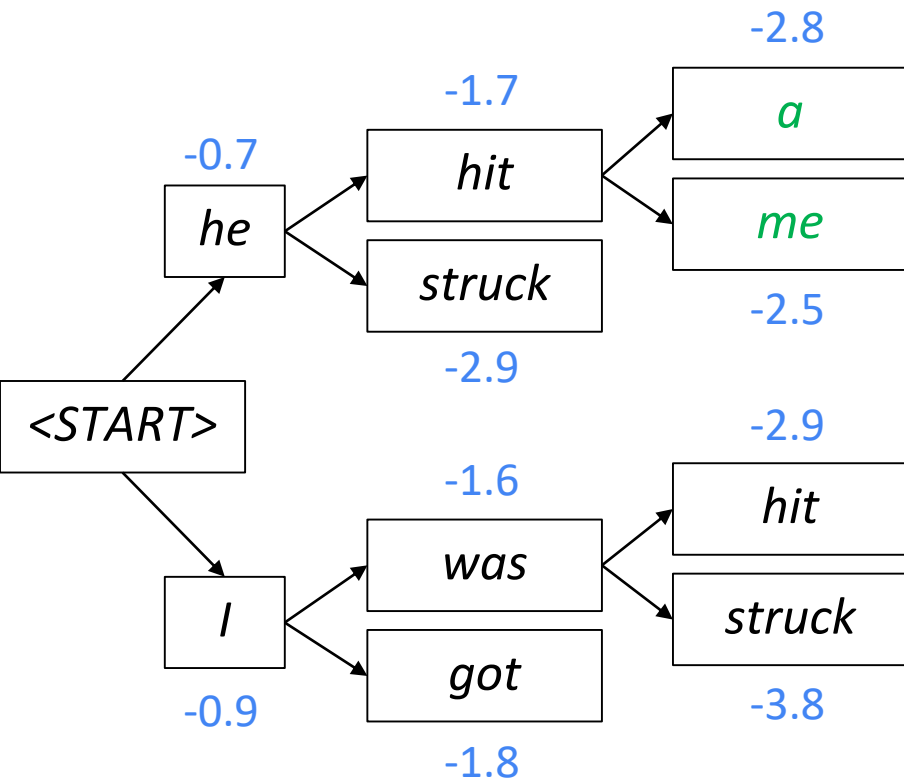just keep $k$ with highest scores

35

# Beam search decoding: example

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

-2.8 = log $P_{\text{LM}}$(*a*|*<START> he hit*) + -1.7

-1.7

*a*

-0.7

*hit*

*he*

*me*

*struck*

-2.5 = log $P_{\text{LM}}$(*me*|*<START> he hit*) + -1.7

-2.9

*<START>*

-2.9 = log $P_{\text{LM}}$(*hit*|*<START> I was*) + -1.6

-1.6

*hit*

*was*

*I*

*struck*

*got*

-3.8 = log $P_{\text{LM}}$(*struck*|*<START> I was*) + -1.6

-0.9

-1.8

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

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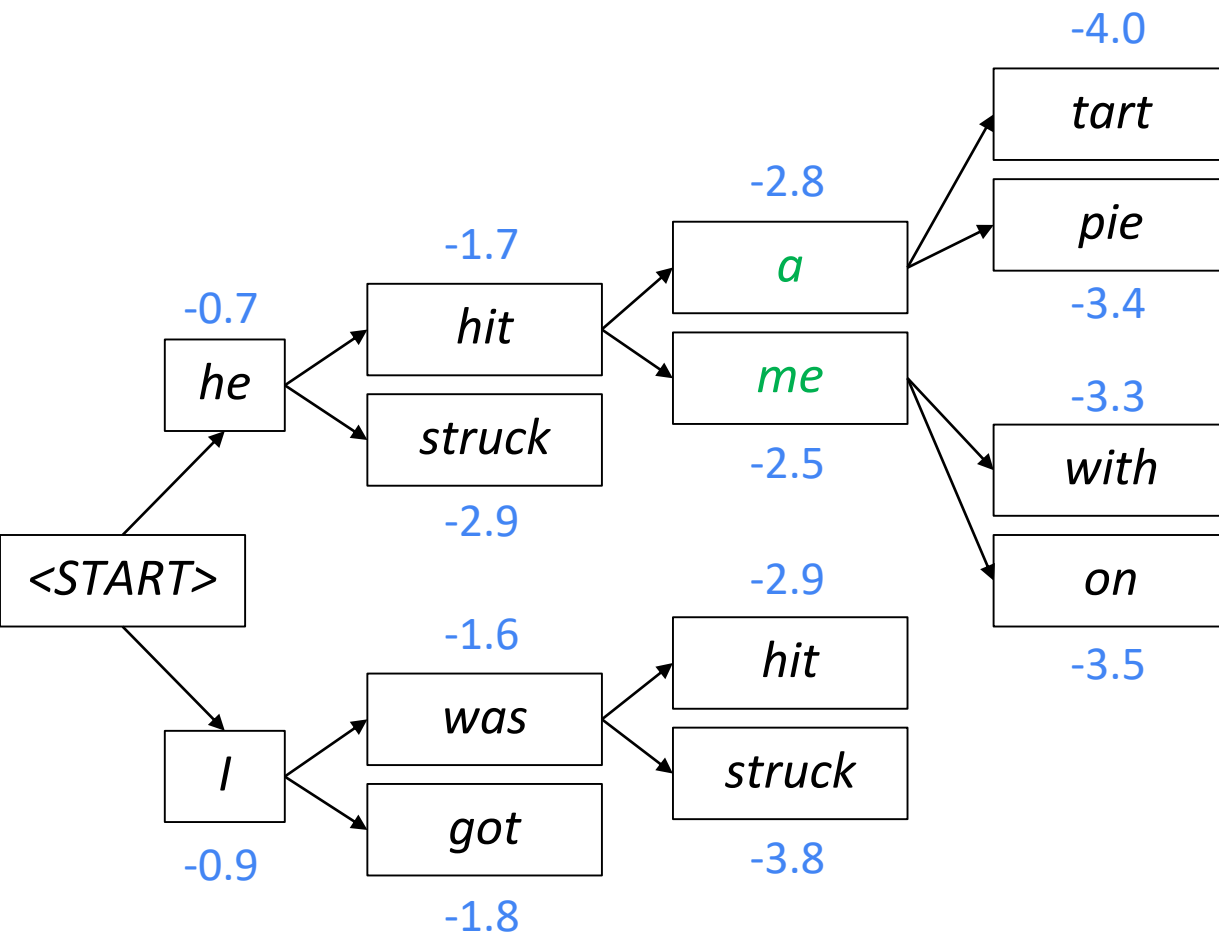
Beam size = k = 2. Blue numbers = $\text{score}(y_1, \ldots, y_t) = \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$

```
                                          -2.8
                            -1.7         ┌────────┐
                          ┌────────┐  ┌─▶│   a    │
              -0.7        │  hit   │──┤  └────────┘
            ┌──────┐   ┌─▶└────────┘  └─▶┌────────┐
            │  he  │───┤                 │   me   │
            └──────┘   └─▶┌────────┐     └────────┘
          ▲              │ struck │        -2.5
          │              └────────┘
┌──────────┐              -2.9
│ <START>  │
└──────────┘              -2.9
          │              ┌────────┐
          │     -1.6  ┌─▶│  hit   │
          ▼   ┌────────┐ └────────┘
            ┌──────┐ │  was  │─┤
            │  I   │─▶└────────┘ └─▶┌────────┐
            └──────┘ │              │ struck │
           -0.9      └─▶┌────────┐  └────────┘
                       │  got   │     -3.8
                       └────────┘
                        -1.8
```

Of these $k^2$ hypotheses,
just keep $k$ with highest scores
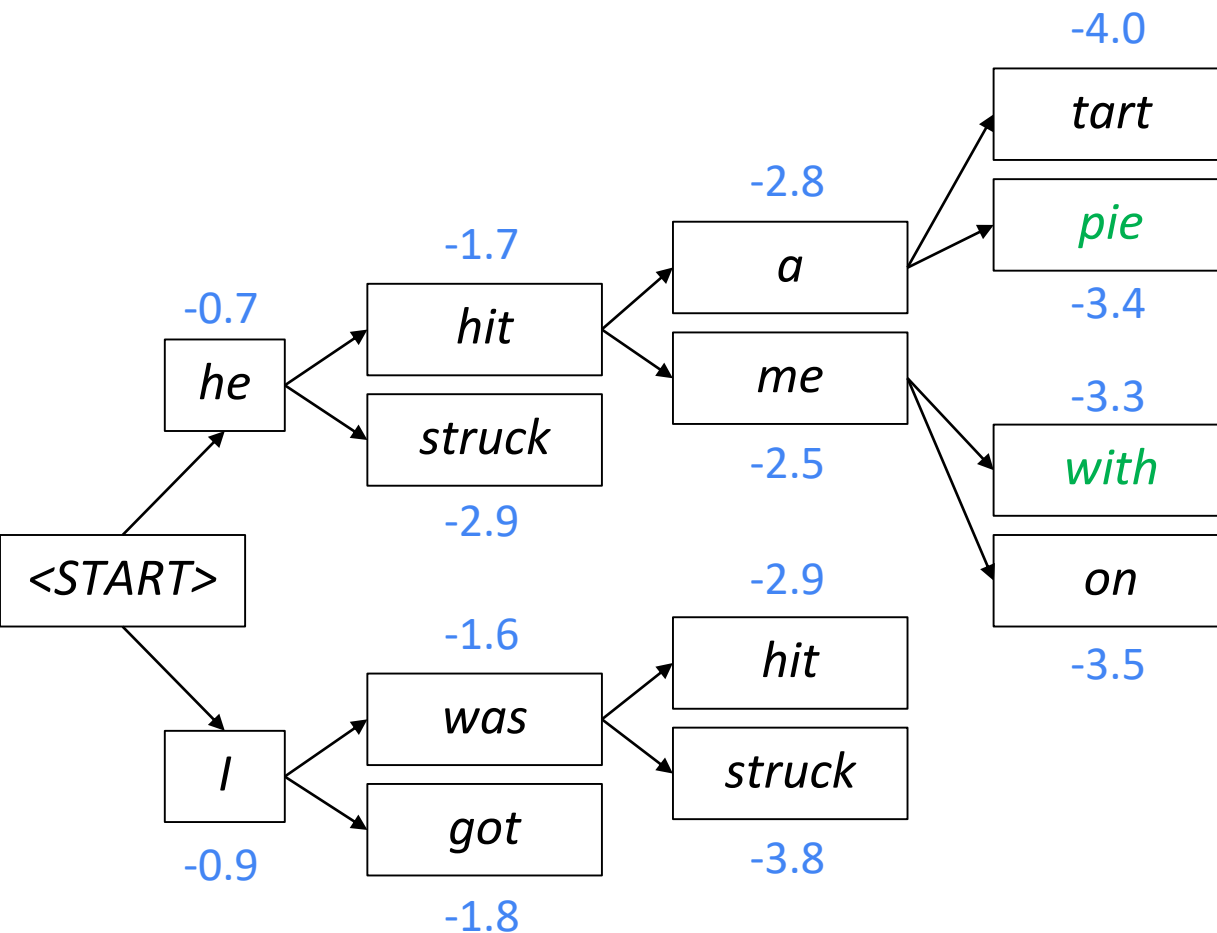
# Beam search decoding: example

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



For each of the *k* hypotheses, find top *k* next words and calculate scores
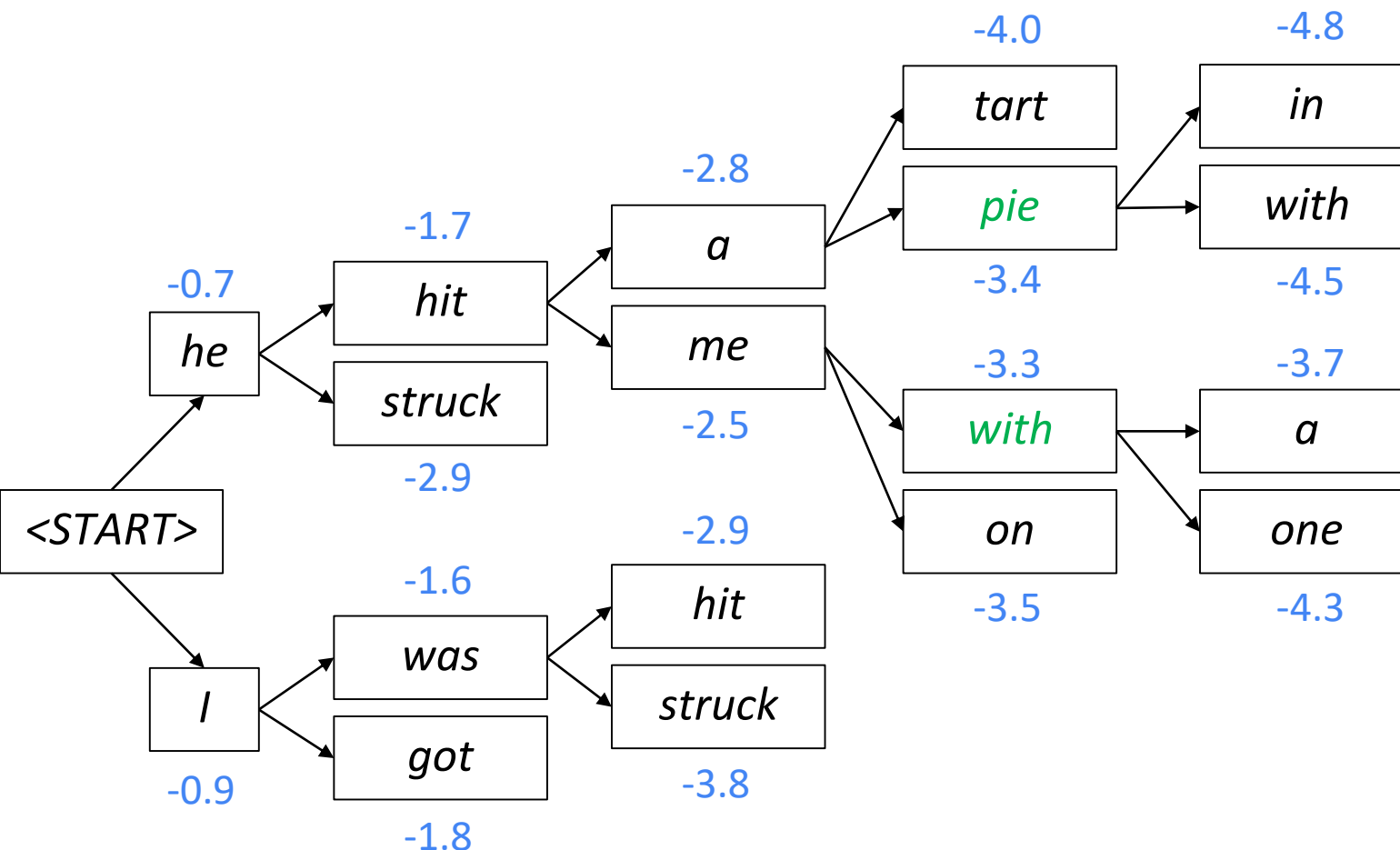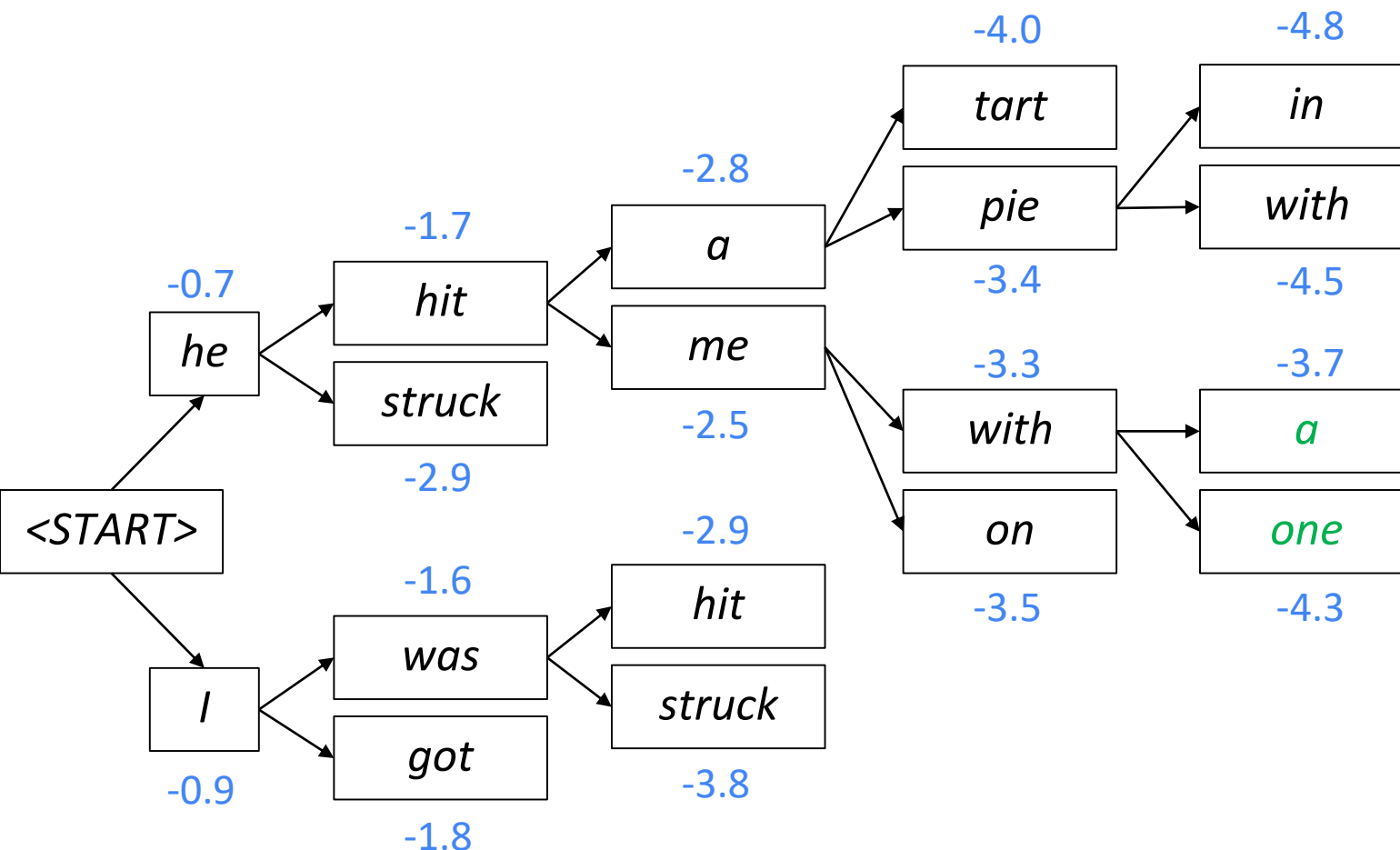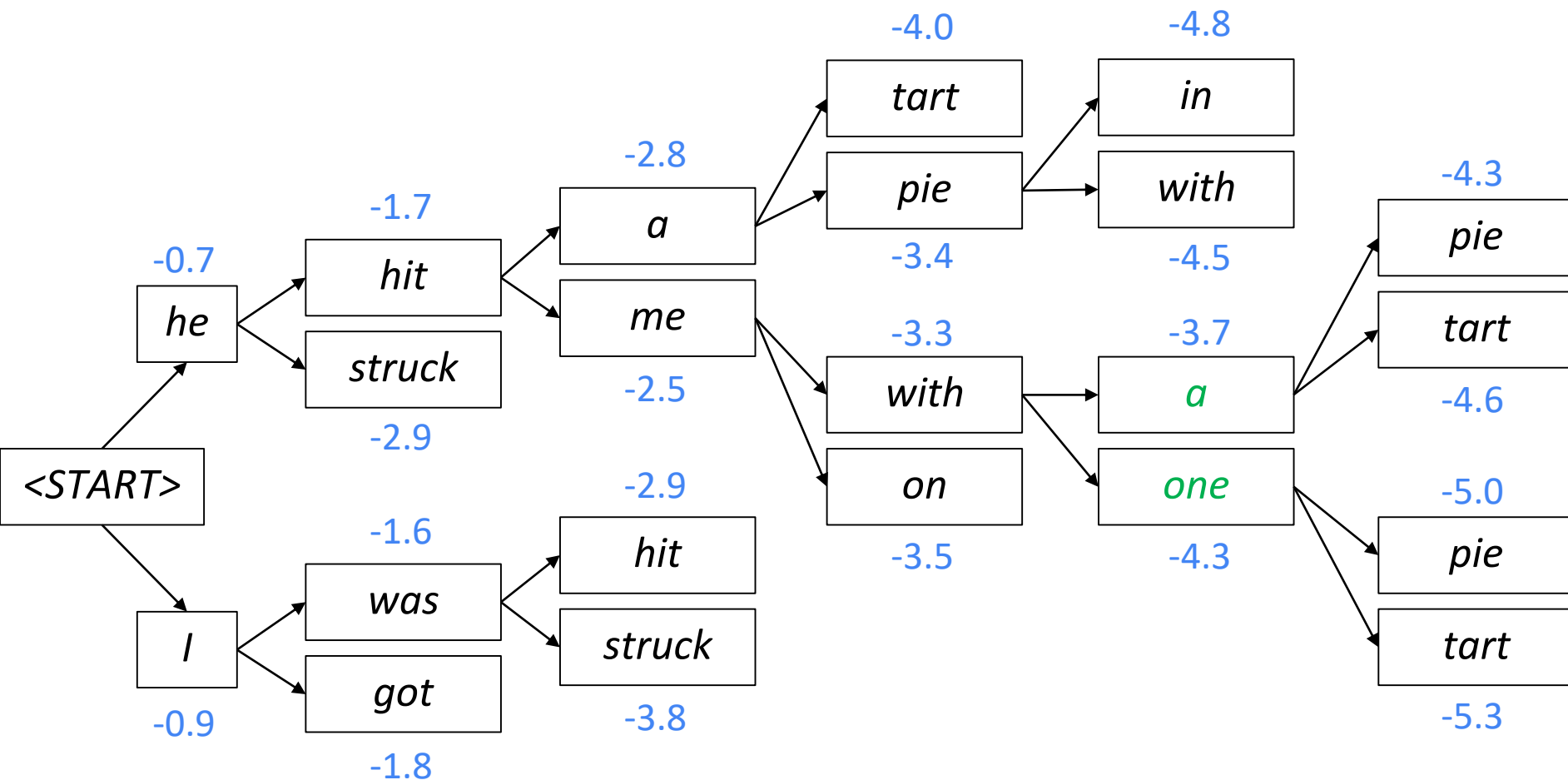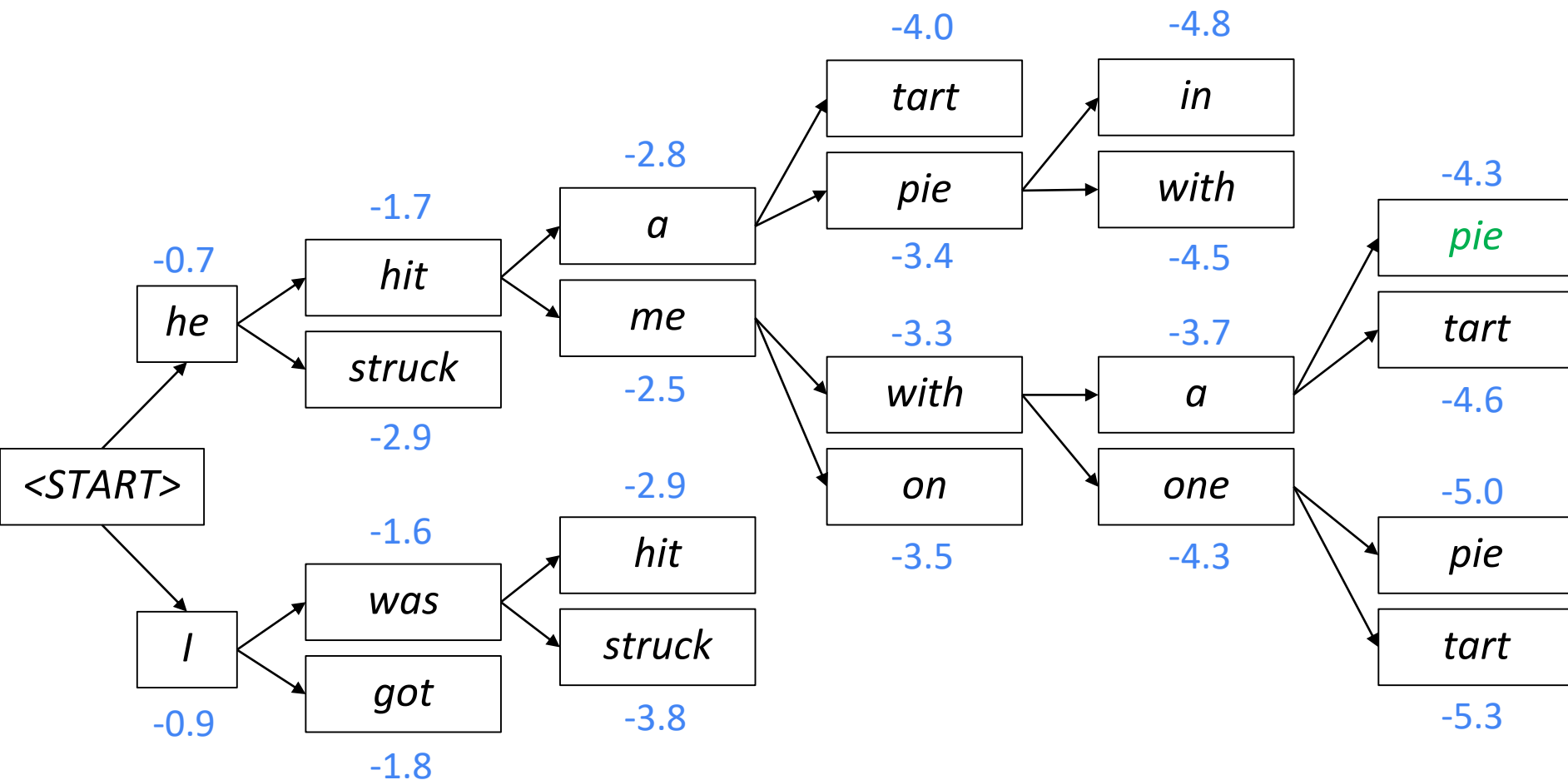
# Beam search decoding: example

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



-0.7
he

-0.9
I

-1.7
hit

-2.9
struck

-1.6
was

-1.8
got

-2.8
a

-2.5
me

-2.9
hit

-3.8
struck

-4.0
tart

-3.4
pie

-3.3
with

-3.5
on

<START>

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

39

# Beam search decoding: example

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



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

# Beam search decoding: example

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



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

41

# Beam search decoding: example

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



For each of the *k* hypotheses, find top *k* next words and calculate scores
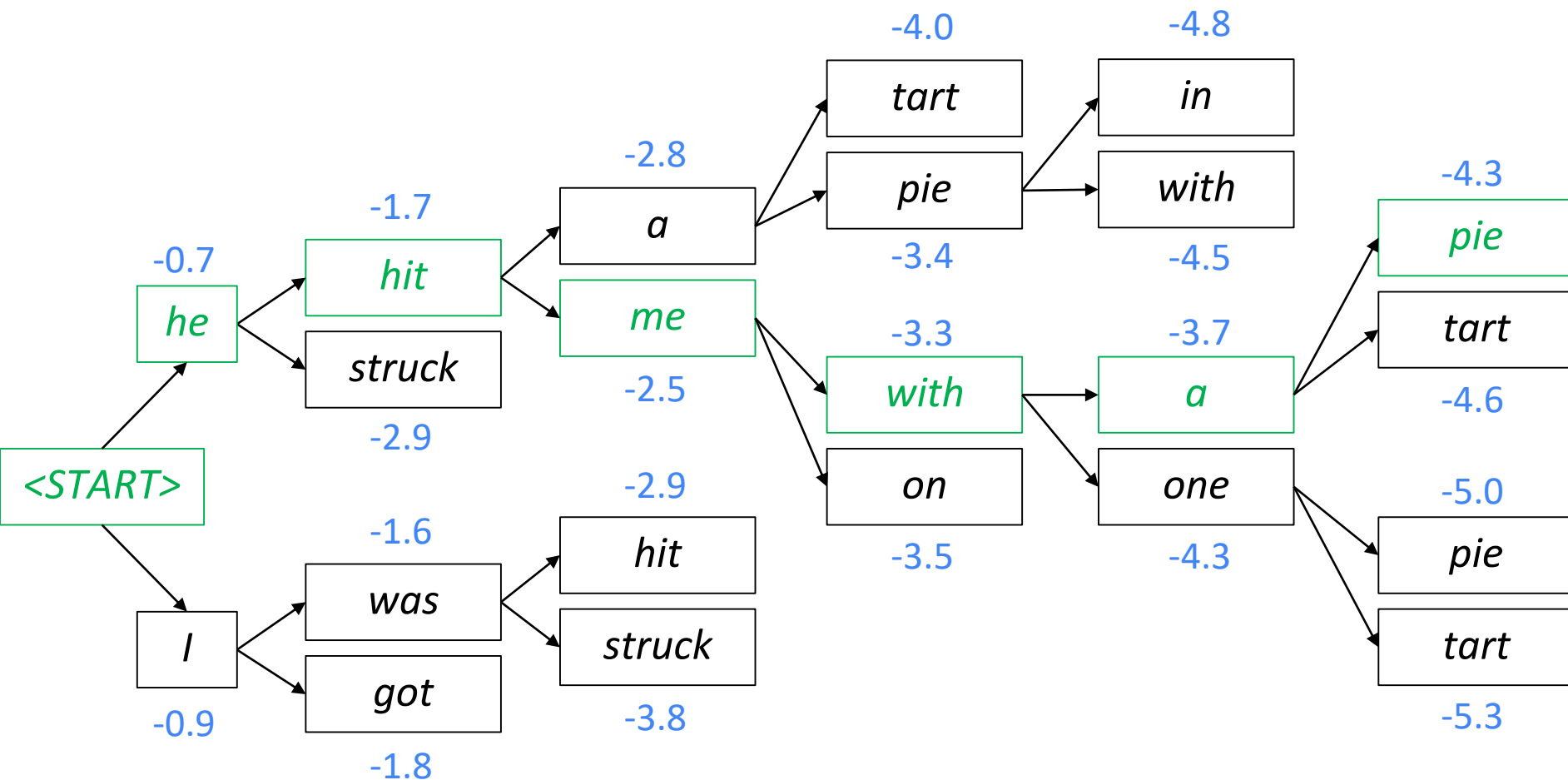
# Beam search decoding: example

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

This is the top-scoring hypothesis!

# Beam search decoding: example

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

Backtrack to obtain the full hypothesis

# Beam search decoding: stopping criterion

- In greedy decoding, usually we decode until the model produces a <END> token
  - For example: *<START> he hit me with a pie <END>*

- In beam search decoding, different hypotheses may produce <END> tokens 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 reach timestep $T$ (where $T$ is some pre-defined cutoff), or
  - We have at least $n$ completed hypotheses (where $n$ is pre-defined cutoff)

# Beam search decoding: finishing up

- We have our list of completed hypotheses.

- How to select top one with highest score?

- Each hypothesis $y_1, \ldots, y_t$ on our list has a score

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

- <u>Problem with this:</u> longer hypotheses have lower scores

  negative value     prob

- <u>Fix:</u> Normalize by length. Use this to select top one instead:

  top-score

$$\frac{1}{t} \sum_{i=1}^{t} \log P_{\text{LM}}(y_i | y_1, \ldots, y_{i-1}, x)$$

# Advantages of NMT

Compared to SMT, NMT has many advantages:

- Better performance
  - More fluent        RNN
  - Better use of context    condition on the source sentence
  - 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!

# How do we evaluate Machine Translation?

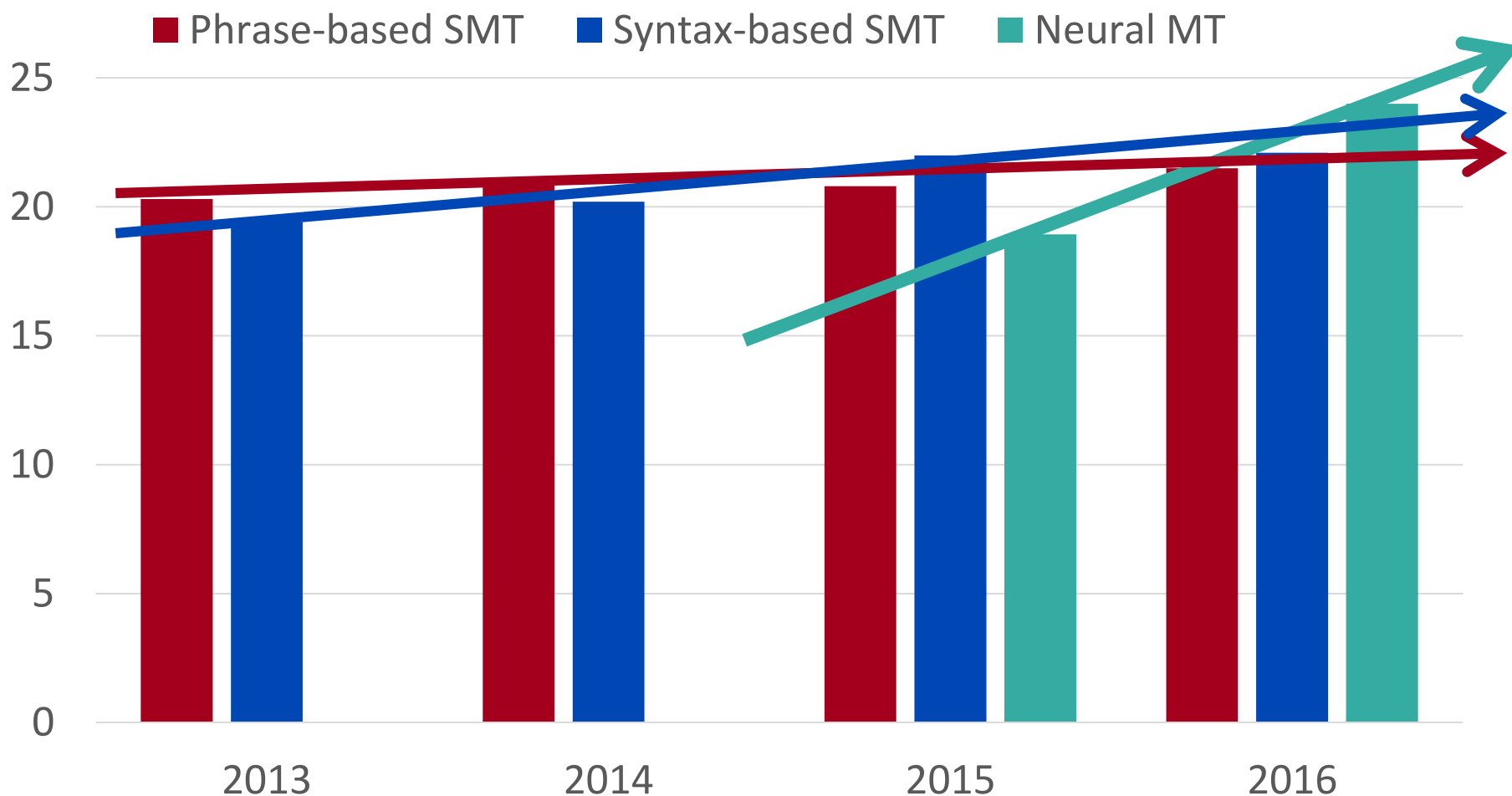**BLEU** (**Bil**ingual **E**valuation **U**nderstudy)

You'll see BLEU in detail in Assignment 4!

- BLEU compares the <u>machine-written translation</u> to one or several <u>human-written translation</u>(s), and computes a similarity score based on:
  - *n*-gram precision (usually for 1, 2, 3 and 4-grams)
  - Plus a penalty for too-short system translations

- BLEU is useful but imperfect
  - 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 ☹

49

**Source:** " BLEU: a Method for Automatic Evaluation of Machine Translation", Papineni et al, 2002.

# 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: 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
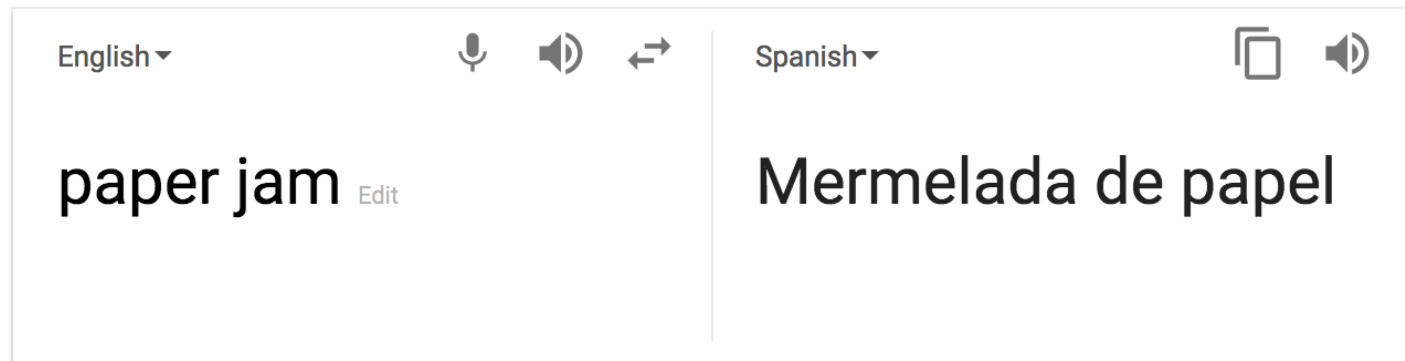
# So is Machine Translation solved?

- **Nope!**
- Many difficulties remain:
  - Out-of-vocabulary words
  - Domain mismatch between train and test data
  - Maintaining context over longer text
  - Low-resource language pairs

**Further reading:** *"Has AI surpassed humans at translation? Not even close!"*
https://www.skynettoday.com/editorials/state_of_nmt

# So is Machine Translation solved?

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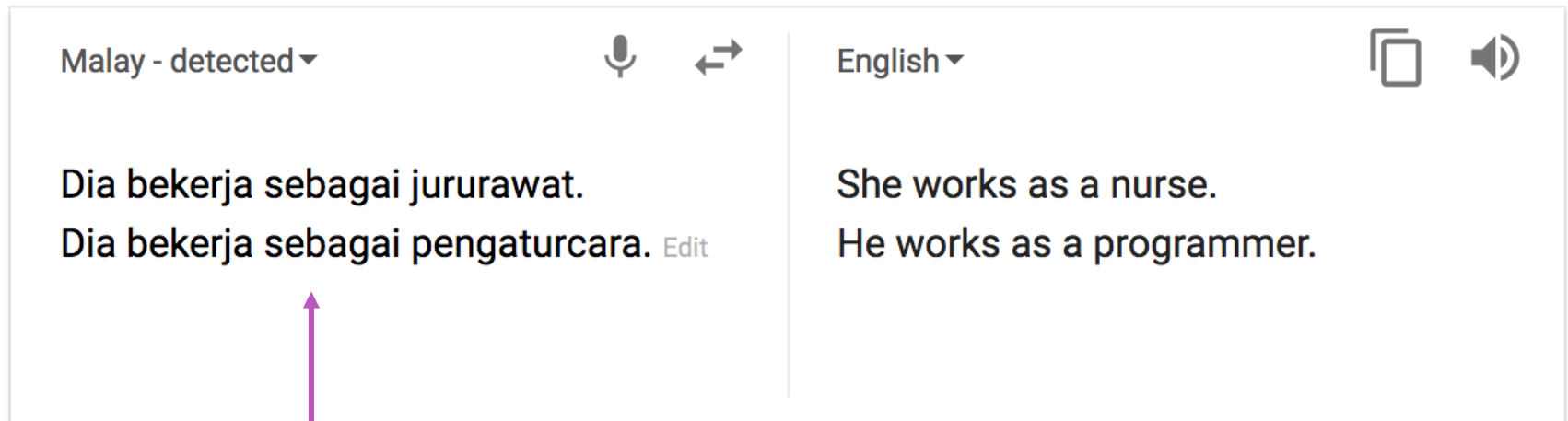- Using common sense is still hard

# So is Machine Translation solved?

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



| Malay - detected | English |
|---|---|
| Dia bekerja sebagai jururawat.<br>Dia bekerja sebagai pengaturcara. Edit | She works as a nurse.<br>He works as a programmer. |

Didn't specify gender

**Source:** https://hackernoon.com/bias-sexist-or-this-is-the-way-it-should-be-ce1f7c8c683c

# So is Machine Translation solved?

- Nope!

- Uninterpretable systems do strange things



Somali ▼
Translate from Irish

ag ag ag ag ag ag ag ag ag ag ag ag
ag ag ag ag ag ag ag ag ag ag ag ag
ag Edit

English ▼

As the name of the LORD was written in the Hebrew language, it was written in the language of the Hebrew Nation

Open in Google Translate                    Feedback

# NMT research continues

NMT is the **flagship task** for NLP Deep Learning
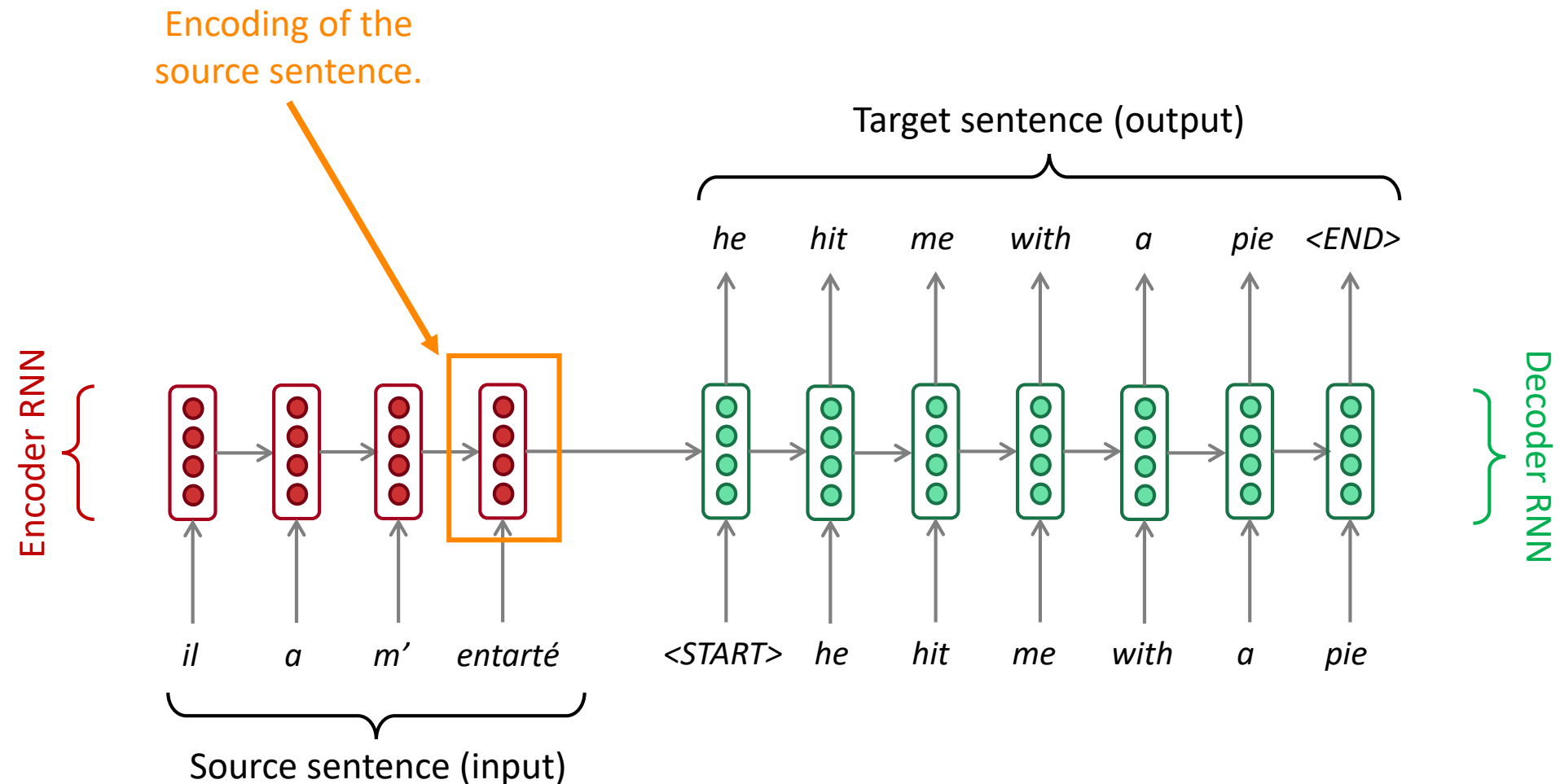
- NMT research has pioneered many of the recent innovations of NLP Deep Learning

- 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...
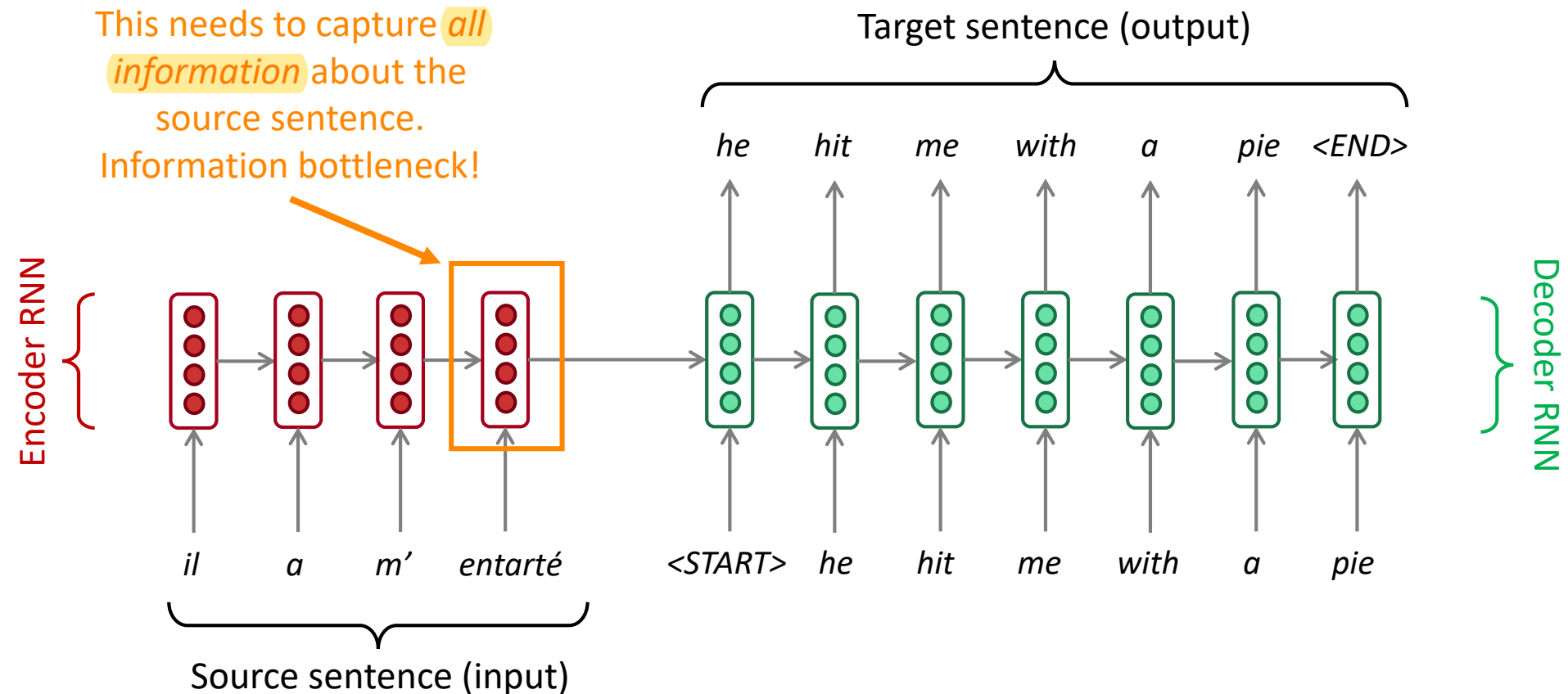
# ATTENTION

# Section 3: Attention

# Sequence-to-sequence: the bottleneck problem



Encoding of the source sentence.

Target sentence (output)

Encoder RNN

Decoder RNN

Source sentence (input)

he hit me with a pie <END>

<START> he hit me with a pie

il a m' entarté

**Problems with this architecture?**

# Sequence-to-sequence: the bottleneck problem

Encoding of the source sentence. This needs to capture *all* *information* about the source sentence. Information bottleneck!

Target sentence (output)

Encoder RNN

Decoder RNN

he    hit    me    with    a    pie    <END>

il    a    m'    entarté    <START>    he    hit    me    with    a    pie
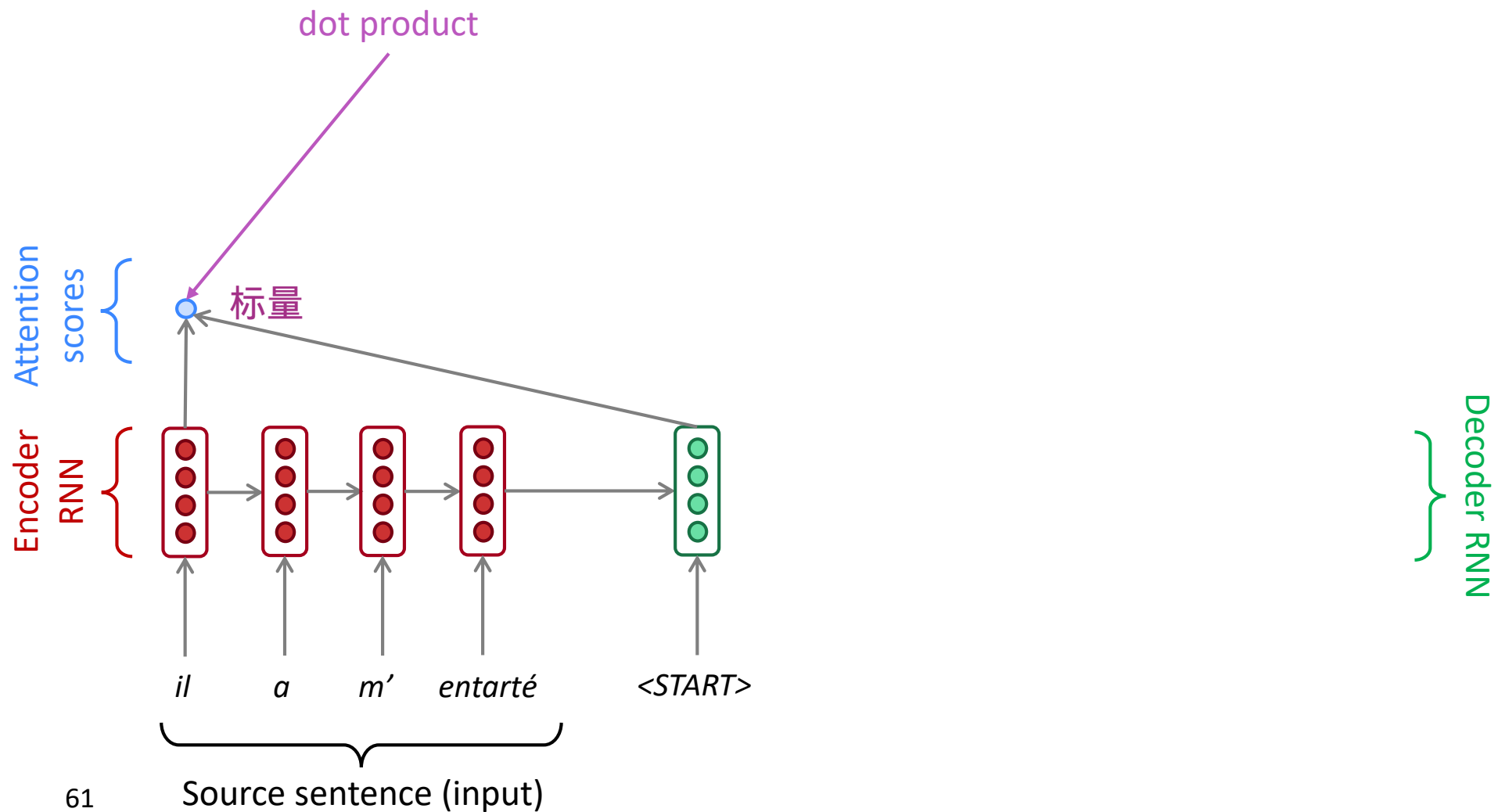
Source sentence (input)

# 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
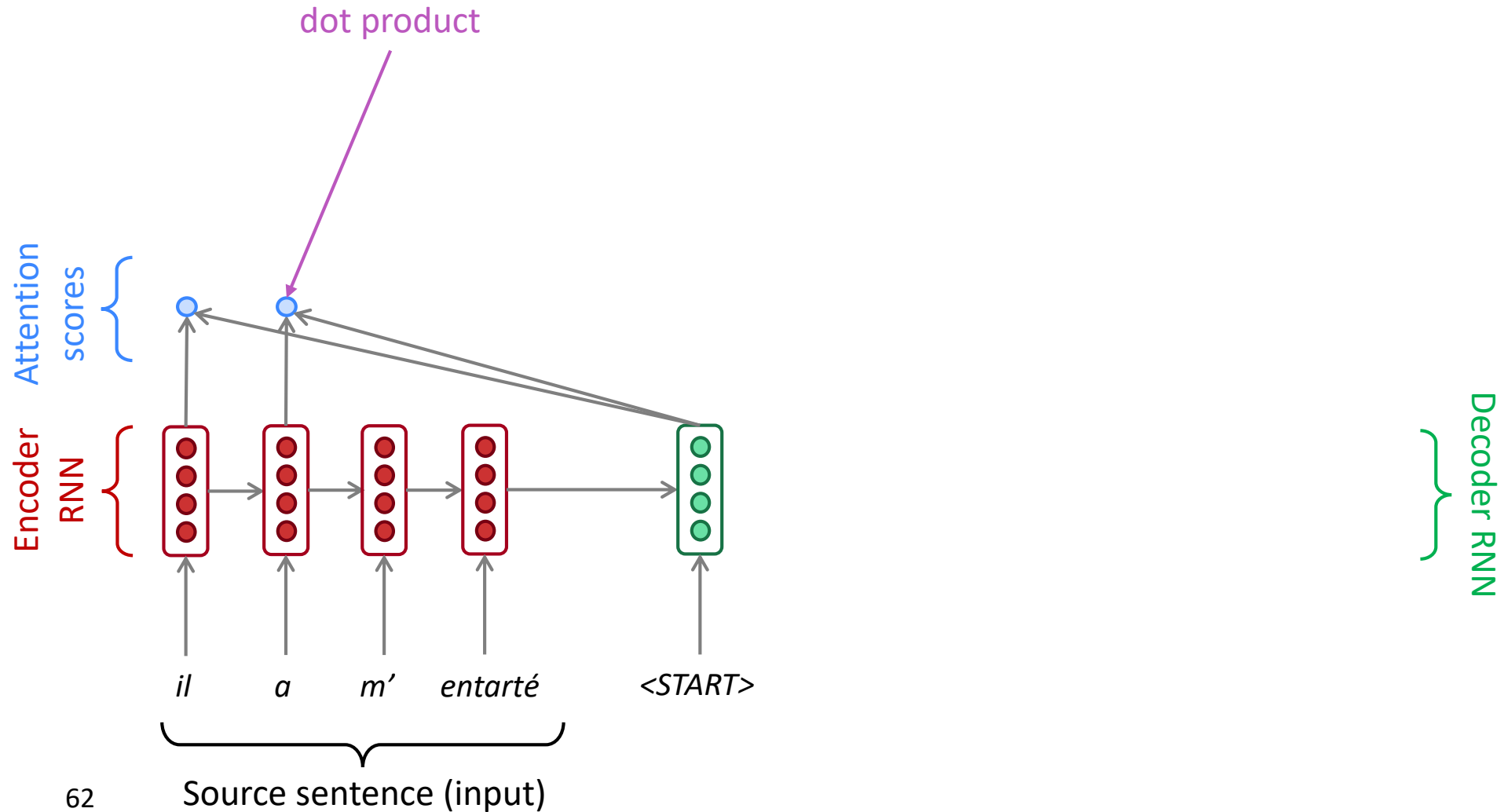
- First we will show via diagram (no equations), then we will show with equations
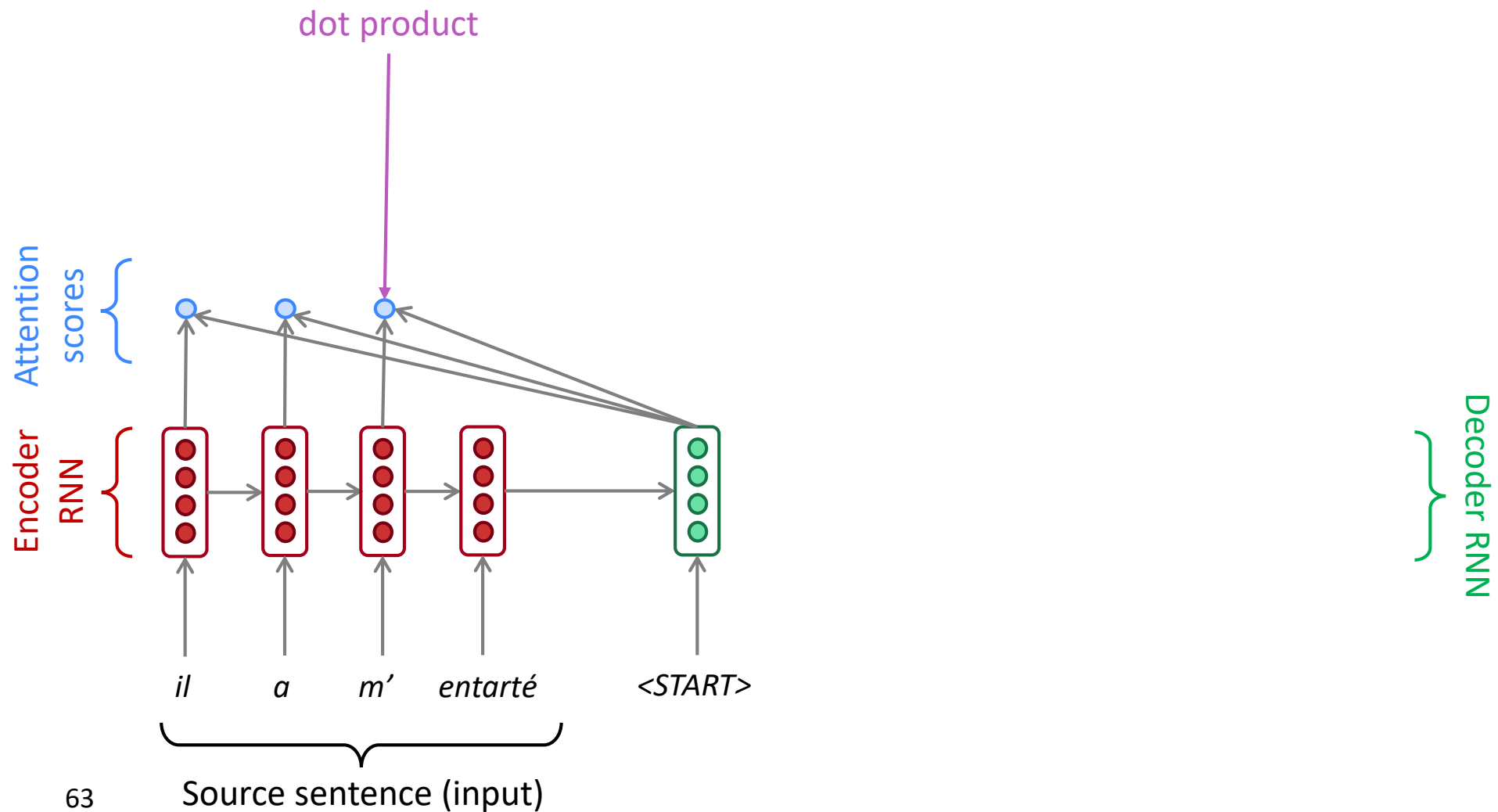
# Sequence-to-sequence with attention

dot product

Encoder RNN
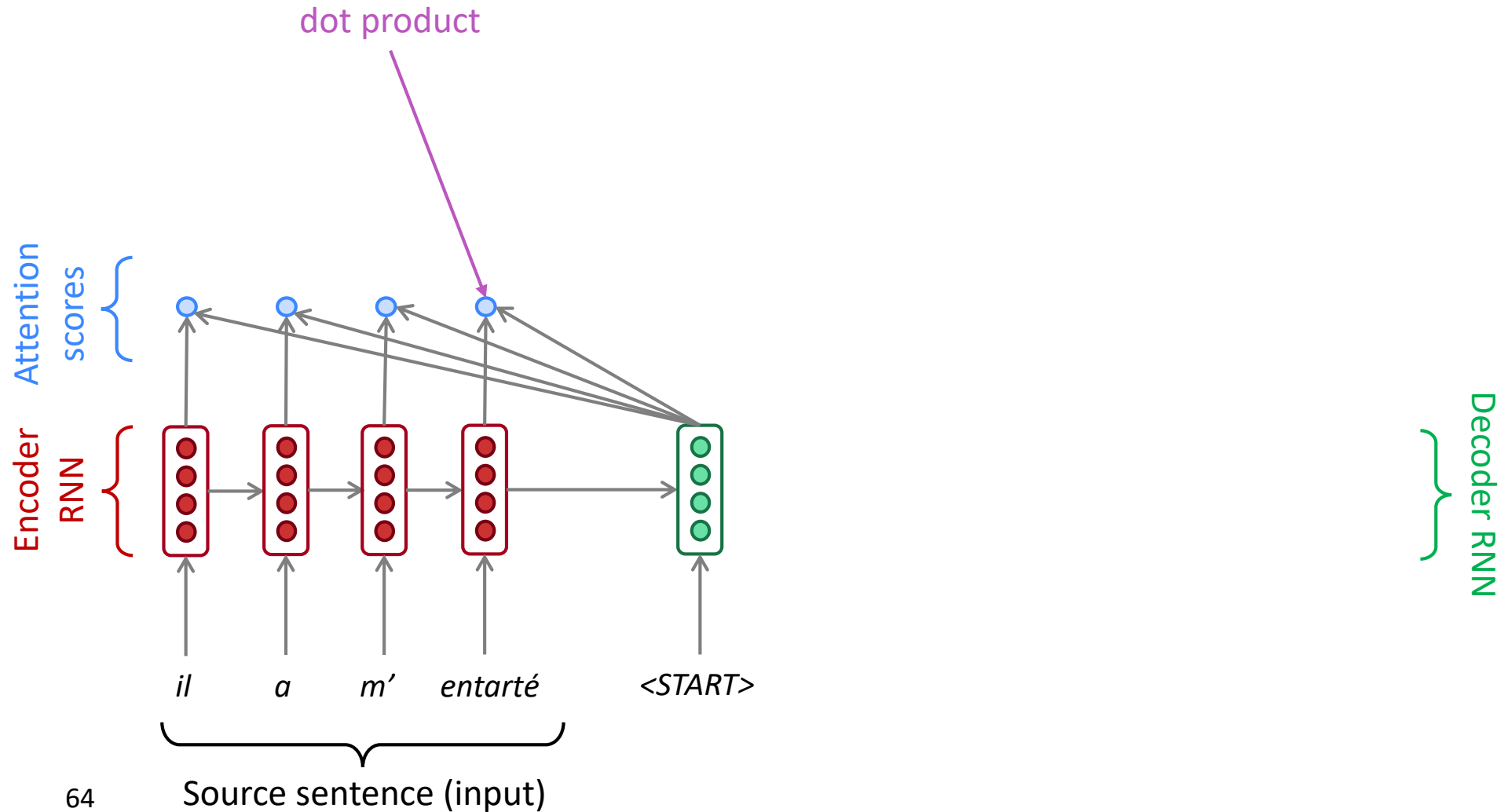
Decoder RNN

il    a    m'    entarté        <START>

Source sentence (input)

61

# Sequence-to-sequence with attention

dot product

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté        <START>

Source sentence (input)

# Sequence-to-sequence with attention

dot product

Attention scores

Encoder RNN

Decoder RNN

*il*    *a*    *m'*    *entarté*     *<START>*

Source sentence (input)

# Sequence-to-sequence with attention



dot product

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté    <START>

Source sentence (input)

64

# Sequence-to-sequence with attention



On this decoder timestep, we're mostly focusing on the first encoder hidden state (*"he"*)

Take softmax to turn the scores into a probability distribution

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

il    a    m'    entarté      &lt;START&gt;

Source sentence (input)

65

# Sequence-to-sequence with attention



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.

# Sequence-to-sequence with attention



Attention output

*he*

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

$\hat{y}_1$

Attention distribution

Attention scores

Encoder RNN

Decoder RNN

*il*     *a*     *m'*     *entarté*     *<START>*

Source sentence (input)

# Sequence-to-sequence with attention



soft alignment        hard

Attention output

Attention distribution

hit

tense

Attention scores

Encoder RNN

y1 hat

Decoder RNN

hit

$\hat{y}_2$

*il      a      m'      entarté          <START>   he*

Source sentence (input)

Sometimes we take the attention output from the previous step, and also feed it into the decoder (along with the usual decoder input). We do this in Assignment 4.

# Sequence-to-sequence with attention

# Sequence-to-sequence with attention

# Sequence-to-sequence with attention

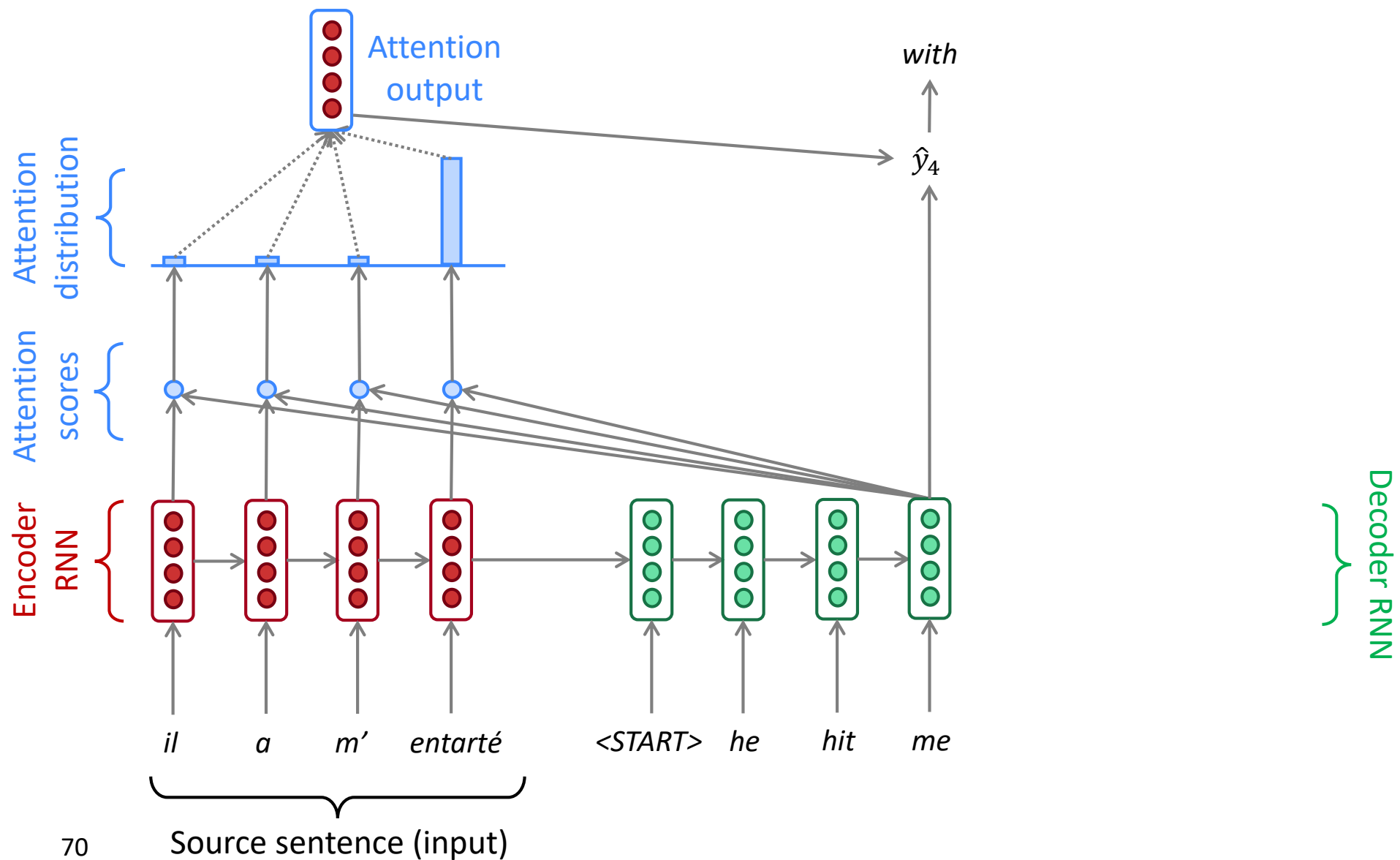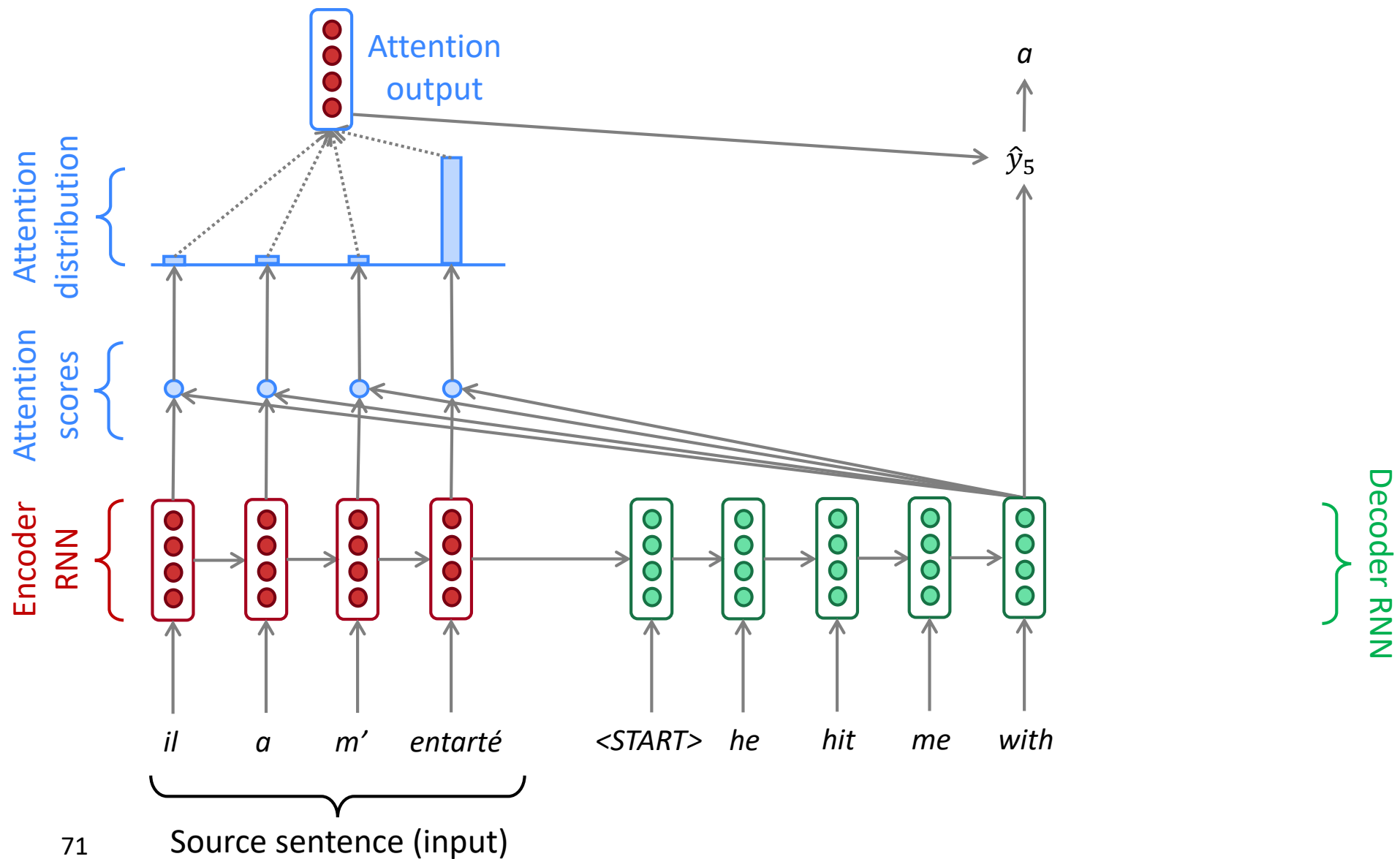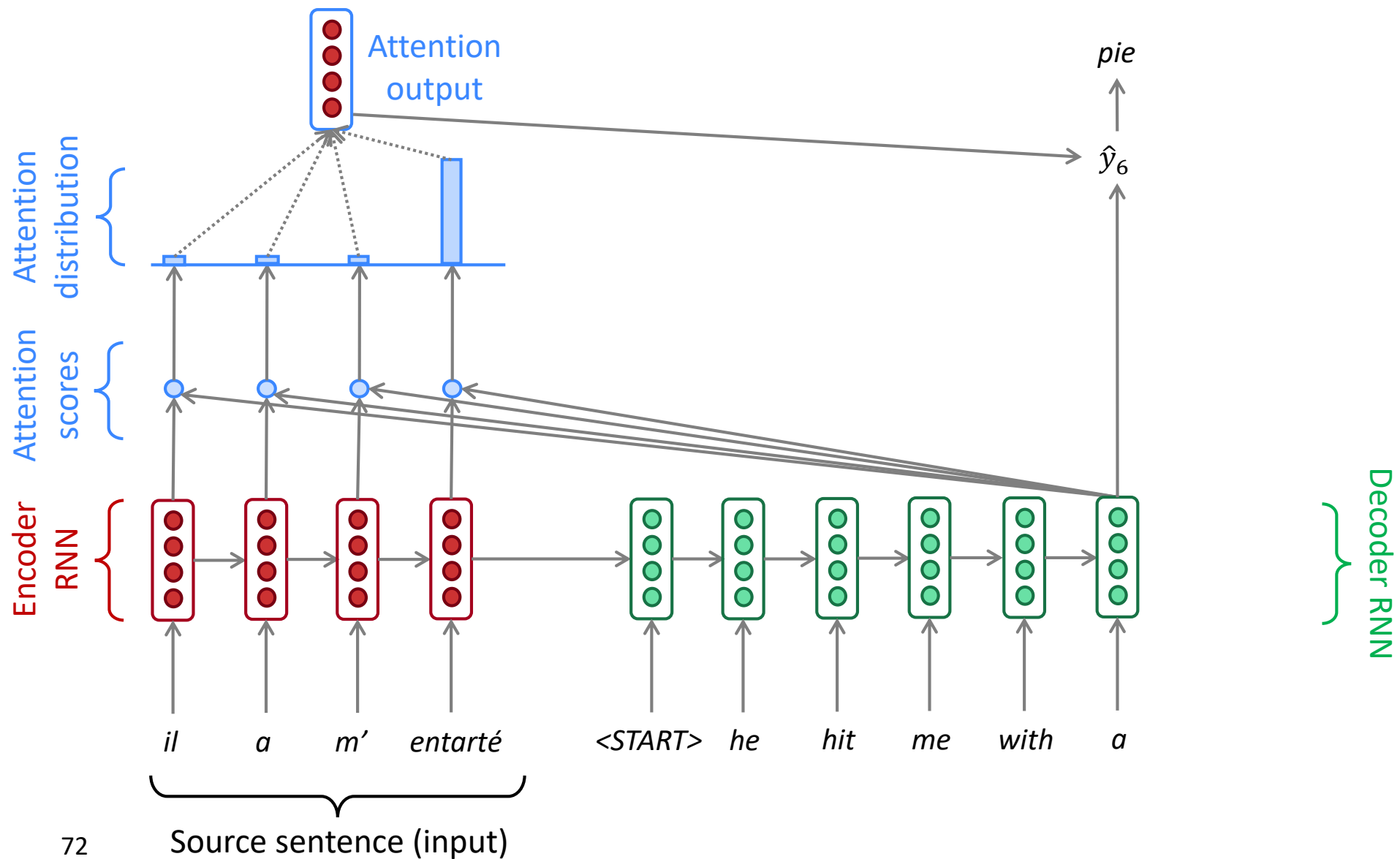# Sequence-to-sequence with attention

# Attention: in equations

- We have <u>encoder hidden states</u> $h_1, \ldots, h_N \in \mathbb{R}^h$

- On timestep $t$, we have <u>decoder hidden state</u> $s_t \in \mathbb{R}^h$

- We get the attention scores $\boldsymbol{e}^t$ for this step:

$$\boldsymbol{e}^t = [\boldsymbol{s}_t^T \boldsymbol{h}_1, \ldots, \boldsymbol{s}_t^T \boldsymbol{h}_N] \in \mathbb{R}^N$$

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

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

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

$$\boldsymbol{a}_t = \sum_{i=1}^{N} \alpha_i^t \boldsymbol{h}_i \in \mathbb{R}^h$$

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

$$[\boldsymbol{a}_t; \boldsymbol{s}_t] \in \mathbb{R}^{2h}$$

# 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

high attention

|        | he | hit | me | with | a | pie |
|--------|----|-----|----|------|---|-----|
| il     |    |     |    |      |   |     |
| a      |    |     |    |      |   |     |
| m'     |    |     |    |      |   |     |
| entarté|    |     |    |      |   |     |

# Attention is a *general* Deep Learning technique

- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.

- <u>However</u>: 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).

75

# 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.

**Intuition**:

- The weighted sum is a *selective summary* of the information contained in the values, where the query determines which values to focus on.
- Attention is a way to obtain a *fixed-size representation of an arbitrary set of representations* (the values), dependent on some other representation (the query).

LSTM     context     gate

# There are *several* attention variants

- We have some *values* $\boldsymbol{h}_1, \ldots, \boldsymbol{h}_N \in \mathbb{R}^{d_1}$ and a *query* $\boldsymbol{s} \in \mathbb{R}^{d_2}$

- Attention always involves:
  1. Computing the *attention scores* $\boldsymbol{e} \in \mathbb{R}^N$ ← There are multiple ways to do this
  2. Taking softmax to get *attention distribution* α:

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

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

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

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

# Attention variants

There are several ways you can compute $e \in \mathbb{R}^N$ from $h_1, \ldots, h_N \in \mathbb{R}^{d_1}$ 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

# 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!