



COS 484

Natural Language Processing

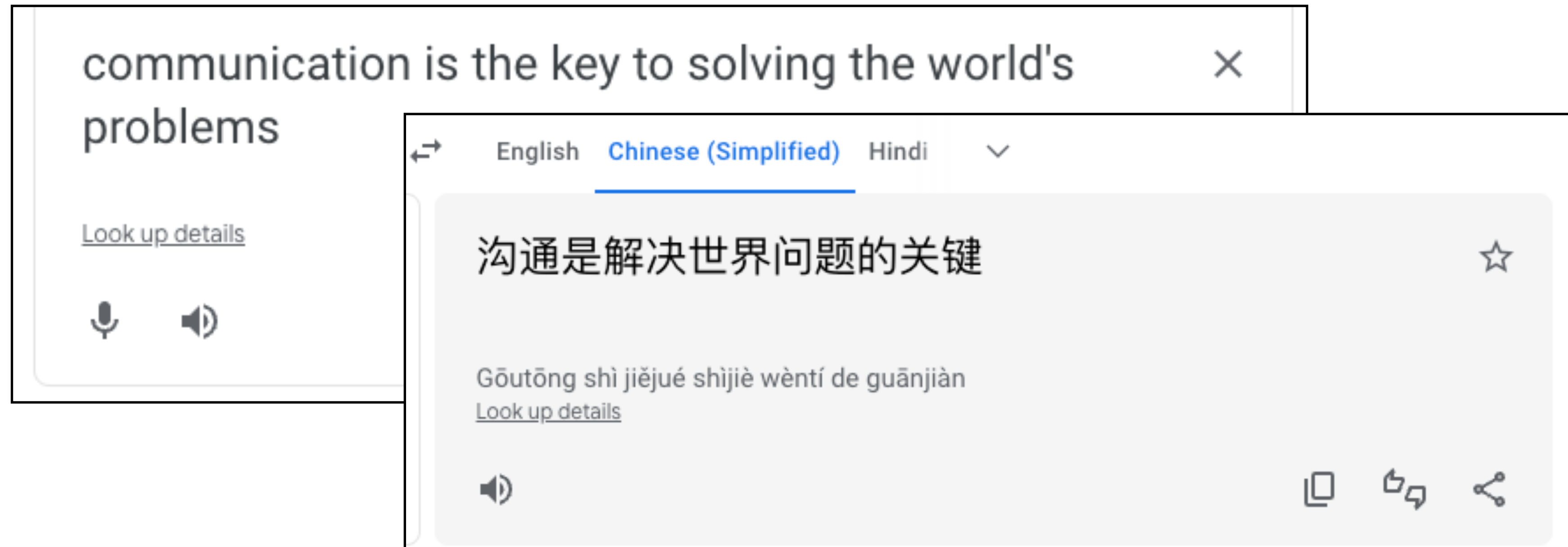
# L13: Machine translation +seq2seq models

Spring 2023

# Lecture plan (two lectures)

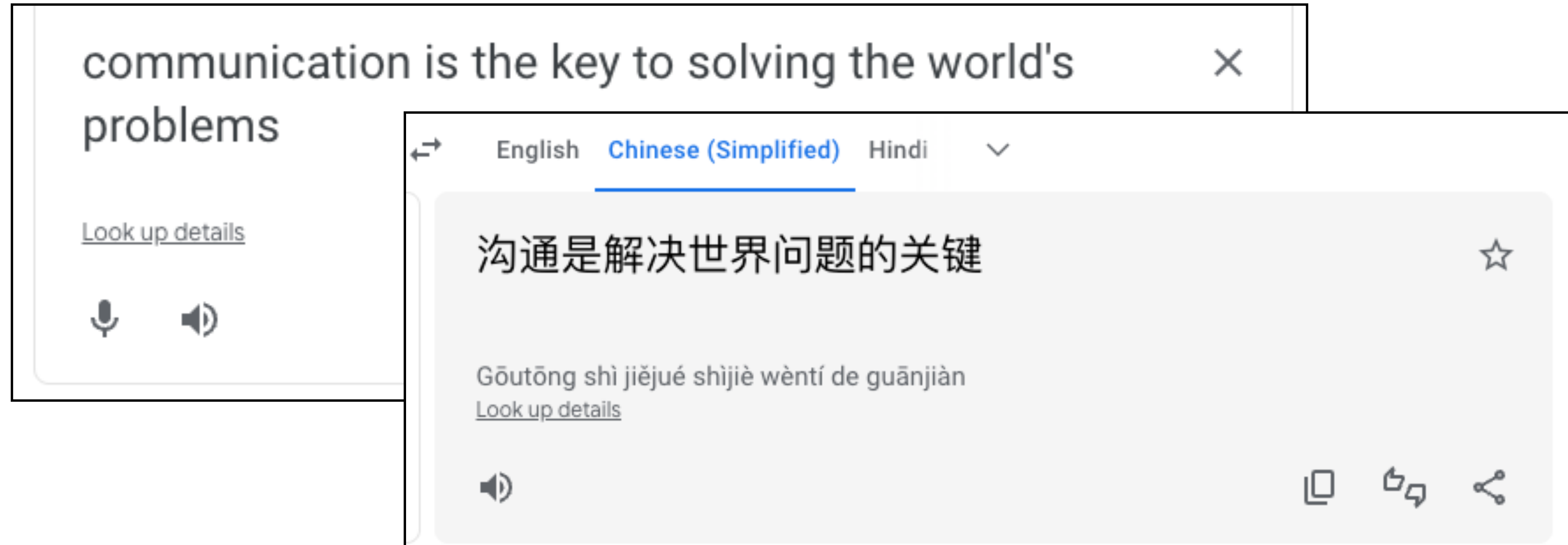
- An introduction to machine translation
  - Statistical machine translation (SMT) → **neural machine translation (NMT)**
- **Sequence to sequence models (seq2seq)**
  - Seq2seq is a general approach for many NLP tasks (summarization, dialogue, parsing, code generation)
- Subword tokenization: how to segment words into word pieces?
- Seq2seq + **Attention**: an effective mechanism to address the fixed representation problem

# Translation



- One of the “holy grail” problems in artificial intelligence
- Practical use case: Facilitate communication between people in the world
- Extremely challenging (especially for low-resource languages)

# Translation



How many languages do you speak?

- A) 1
- B) 2
- C) 3
- D) 4+

# Machine translation (MT)

- Goal: Translate a sentence  $\mathbf{w}^{(s)}$  in a source language (input) to a sentence  $\mathbf{w}^{(t)}$  in the target language (output)

I like apples  $\leftrightarrow$  ich mag Äpfel (German)

- Why is MT challenging?

- Single words may be replaced with multi-word phrases:

I like apples  $\leftrightarrow$  J'aime les pommes (French)

- Reordering of phrases:

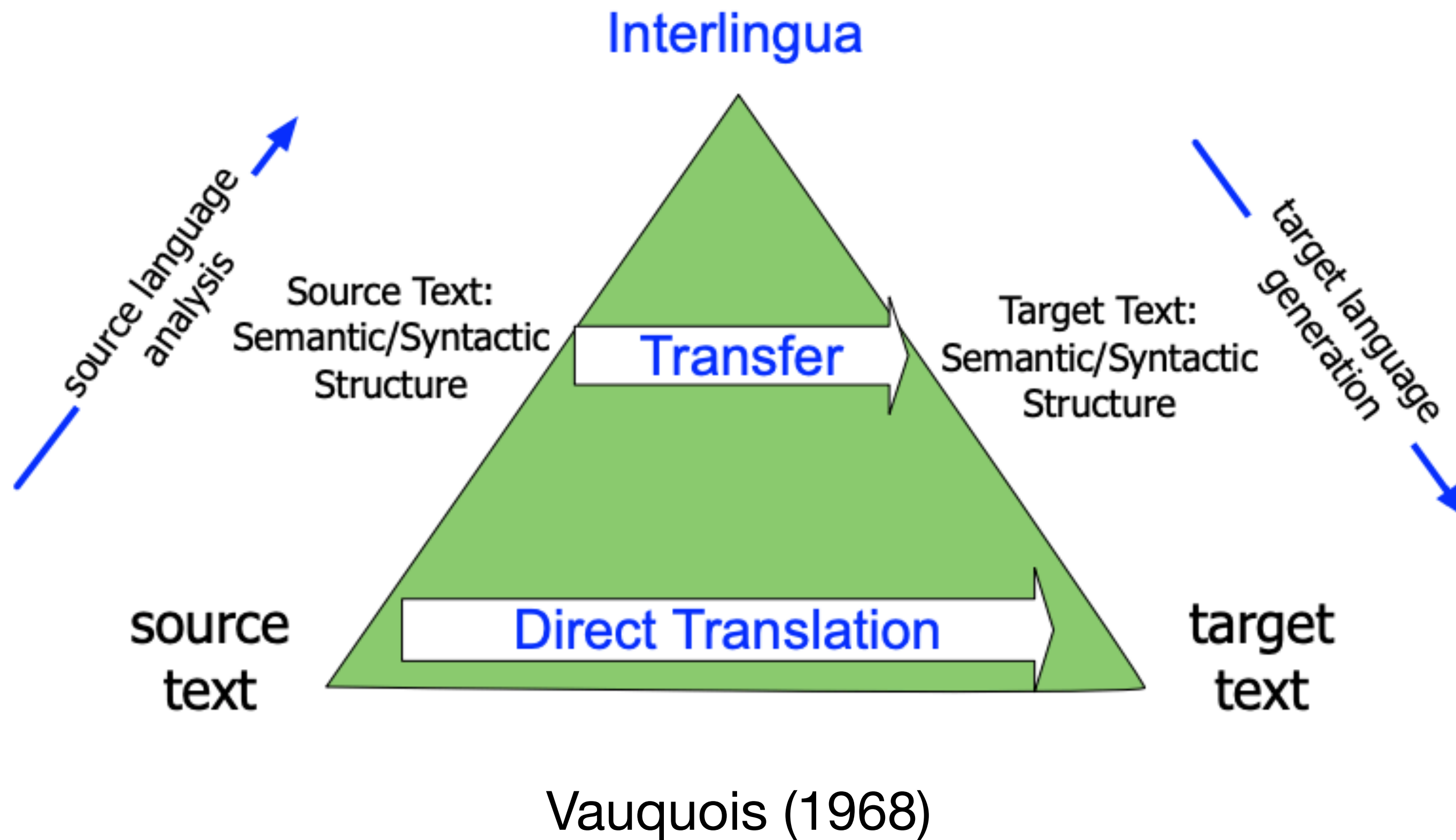
I like red apples  $\leftrightarrow$  J'aime les pommes rouges (French)

- Context-dependent translations:

*les*  $\leftrightarrow$  *the*    but    *les pommes*  $\leftrightarrow$  *apples*

Extremely large output space  $\implies$  Decoding is NP-hard

# Vauquois triangle



- **Direct translation:** word-by-word
- **Transfer approaches:** we first parse the input text and transform the source-language parse into a target-language parse
- **Interlingua approaches:** we first transform source text into an interlingua (generic language-agnostic representation of meaning) and then generate into target language

# Evaluating machine translation



Two main criteria:

- **Adequacy:** Translation  $\mathbf{w}^{(t)}$  should adequately reflect the linguistic content of  $\mathbf{w}^{(s)}$
- **Fluency:** Translation  $\mathbf{w}^{(t)}$  should be fluent text in the target language

---

---

*To Vinay it like Python*  
*Vinay debugs memory leaks*  
*Vinay likes Python*

---

Different translations of  
“A Vinay le gusta Python” (Spanish)

Which of these translations is both  
adequate and fluent?

- A) first
- B) second
- C) third
- D) none of them



# Evaluating machine translation



Two main criteria:

- **Adequacy:** Translation  $\mathbf{w}^{(t)}$  should adequately reflect the linguistic content of  $\mathbf{w}^{(s)}$
- **Fluency:** Translation  $\mathbf{w}^{(t)}$  should be fluent text in the target language

	Adequate?	Fluent?
<i>To Vinay it like Python</i>	yes	no
<i>Vinay debugs memory leaks</i>	no	yes
<i>Vinay likes Python</i>	yes	yes

Different translations of  
“A Vinay le gusta Python” (Spanish)

Which of these translations is both  
adequate and fluent?

- A) first
- B) second
- C) third
- D) none of them

The answer is (C).



# Evaluation metrics

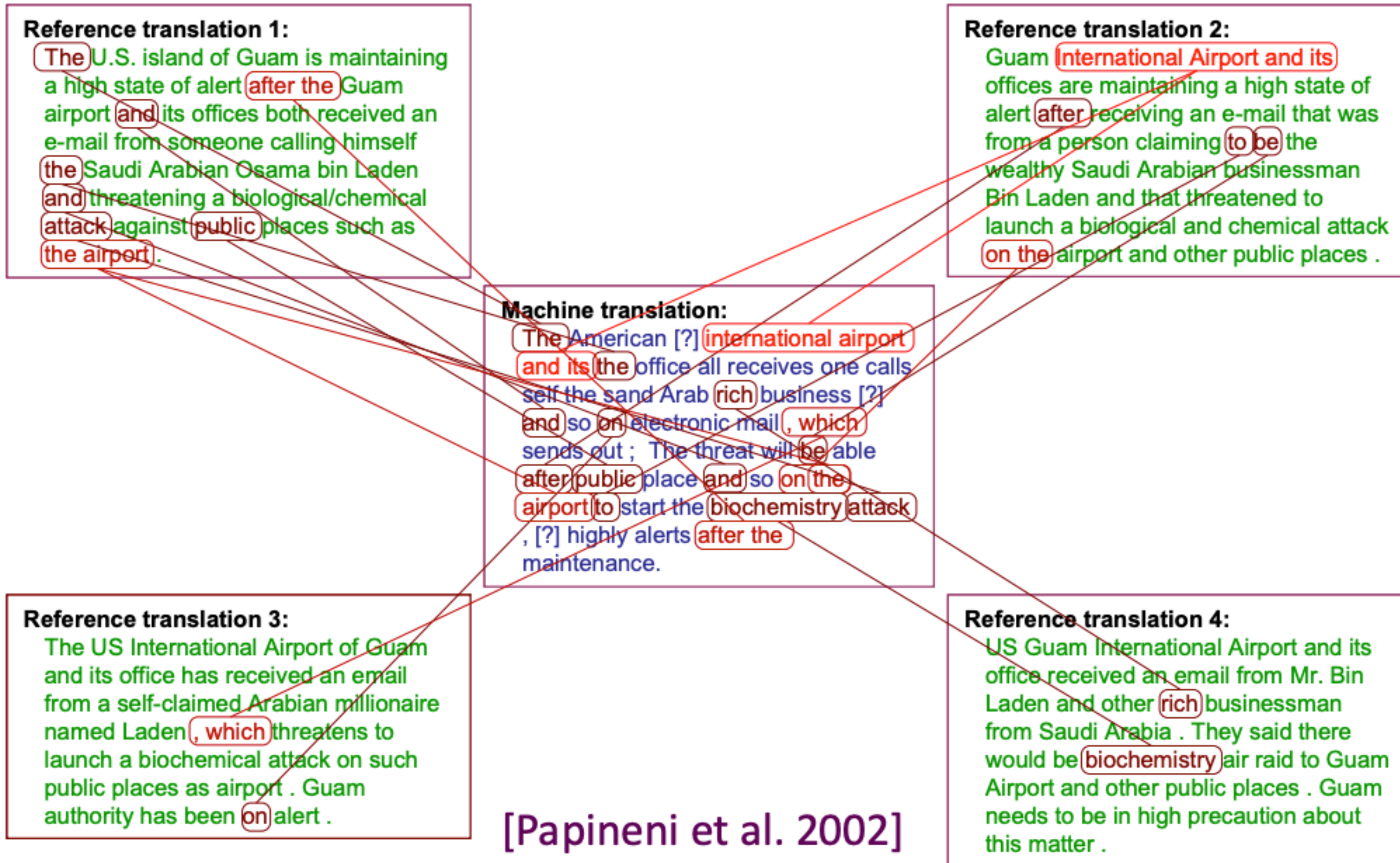
- **Manual evaluation:** ask a native speaker to verify the translation
  - Most accurate, but expensive
- **Automated evaluation metrics:**
  - Compare system hypothesis with reference translations
  - BiLingual Evaluation Understudy (BLEU) (Papineni et al., 2002):
    - Modified n-gram precision

$$p_n = \frac{\text{number of } n\text{-grams appearing in both reference and hypothesis translations}}{\text{number of } n\text{-grams appearing in the hypothesis translation}}$$

Reference translation

System predictions

# Evaluation metric: BLEU





# Evaluation metric: BLEU

- Calculate modified n-gram precision  $p_n$  (usually for 1, 2, 3 and 4-grams)
- Plus a “brevity penalty” for too-short system translations
- The final BLEU score takes the geometric mean of  $p_n$  (with smoothing)  $\times$  brevity penalty
- BLEU ranges between 0 and 1 and people usually express them in percentage

BP: brevity penalty

	Translation	$p_1$	$p_2$	$p_3$	$p_4$	BP	BLEU
Reference	<i>Vinay likes programming in Python</i>						
Sys1	<i>To Vinay it like to program Python</i>	$\frac{2}{7}$	0	0	0	1	.21
Sys2	<i>Vinay likes Python</i>	$\frac{3}{3}$	$\frac{1}{2}$	0	0	.51	.33
Sys3	<i>Vinay likes programming in his pajamas</i>	$\frac{4}{6}$	$\frac{3}{5}$	$\frac{2}{4}$	$\frac{1}{3}$	1	.76

BLEU is **useful (and widely used)** but **far from perfect**

A **good** translation can get a **poor** BLEU score because it has low n-gram overlap with human translation

Sample BLEU scores for various system outputs

# Machine translation: Data

- Statistical MT relies requires **parallel corpora (bilingual)**

1. Chapter 4, Koch (DE)	de	es
<p><b>context</b> We would like to ensure that there is a reference to this <b>as early as the recitals</b> and that the period within which the Council has to make a decision - which is not clearly worded - is set at a maximum of three months .</p>	<p>Wir möchten sicherstellen , daß hierauf bereits in den Erwägungsgründen hingewiesen wird und die uneindeutig formulierte Frist , innerhalb der der Rat eine Entscheidung treffen muß , auf maximal drei Monate fixiert wird .</p>	<p>Quisiéramos asegurar que se aluda ya a esto en los considerandos y que el plazo , imprecisamente formulado , dentro del cual el Consejo ha de adoptar una decisión , se fije en tres meses como máximo .</p>
2. Chapter 3, FÅarm (SV)	de	es
<p><b>context</b> Our experience of modern administration tells us that openness , decentralisation of responsibility and qualified evaluation are often <b>as effective as detailed bureaucratic supervision</b> .</p>	<p>Unsere Erfahrungen mit moderner Verwaltung besagen , daß Transparenz , Dezentralisation der Verantwortlichkeiten und eine qualifizierte Auswertung oft ebenso effektiv sind wie bürokratische Detailkontrolle .</p>	<p>Nuestras experiencias en materia de administración moderna nos señalan que la apertura , la descentralización de las responsabilidades y las evaluaciones bien hechas son a menudo tan eficaces como los controles burocráticos detallados .</p>

*(Europarl, Koehn, 2005)*

- And lots of it!
- Not easily available for many low-resource languages in the world

# Machine translation: Data

**21 European languages:** Romanic (French, Italian, Spanish, Portuguese, Romanian), Germanic (English, Dutch, German, Danish, Swedish), Slavik (Bulgarian, Czech, Polish, Slovak, Slovene), Finni-Ugric (Finnish, Hungarian, Estonian), Baltic (Latvian, Lithuanian), and Greek.

<b>Parallel Corpus (L1-L2)</b>	<b>Sentences</b>	<b>L1 Words</b>	<b>English Words</b>
Bulgarian-English	406,934	-	9,886,291
Czech-English	646,605	12,999,455	15,625,264
Danish-English	1,968,800	44,654,417	48,574,988
German-English	1,920,209	44,548,491	47,818,827
Greek-English	1,235,976	-	31,929,703
Spanish-English	1,965,734	51,575,748	49,093,806
Estonian-English	651,746	11,214,221	15,685,733
Finnish-English	1,924,942	32,266,343	47,460,063
French-English	2,007,723	51,388,643	50,196,035

<https://www.statmt.org/europarl/>



# Statistical machine translation (SMT)

- Core idea: Learn a probabilistic model from data
- Suppose we are translating French  $\rightarrow$  English
- We want to find **best target sentence**  $\mathbf{w}^{(t)}$ , given **source sentence**  $\mathbf{w}^{(s)}$

$$\arg \max_{\mathbf{w}^{(t)}} P(\mathbf{w}^{(t)} \mid \mathbf{w}^{(s)})$$

- According to Bayes' rule, we can break this down into two components:

$$= \arg \max_{\mathbf{w}^{(t)}} P(\mathbf{w}^{(s)} \mid \mathbf{w}^{(t)}) P(\mathbf{w}^{(t)})$$

**Translation model:** models whether the target sentence reflects the linguistic content of the source language (adequacy)  
Learned from **parallel** data

**Language model:** models how fluent the target sentence is (fluency)  
Can be learned from **monolingual** data



# Statistical machine translation (SMT)

$$= \arg \max_{\mathbf{w}^{(t)}} P(\mathbf{w}^{(s)} | \mathbf{w}^{(t)}) P(\mathbf{w}^{(t)})$$

**Translation model:** models whether the target sentence reflects the linguistic content of the source language (adequacy)  
Learned from **parallel** data

**Language model:** models how fluent the target sentence is (fluency)  
Can be learned from **monolingual** data

How should we align words in source to words in target?

	<i>A</i>	<i>Vinay</i>	<i>le</i>	<i>gusta</i>	<i>python</i>
<i>Vinay</i>					
<i>likes</i>					
<i>python</i>					

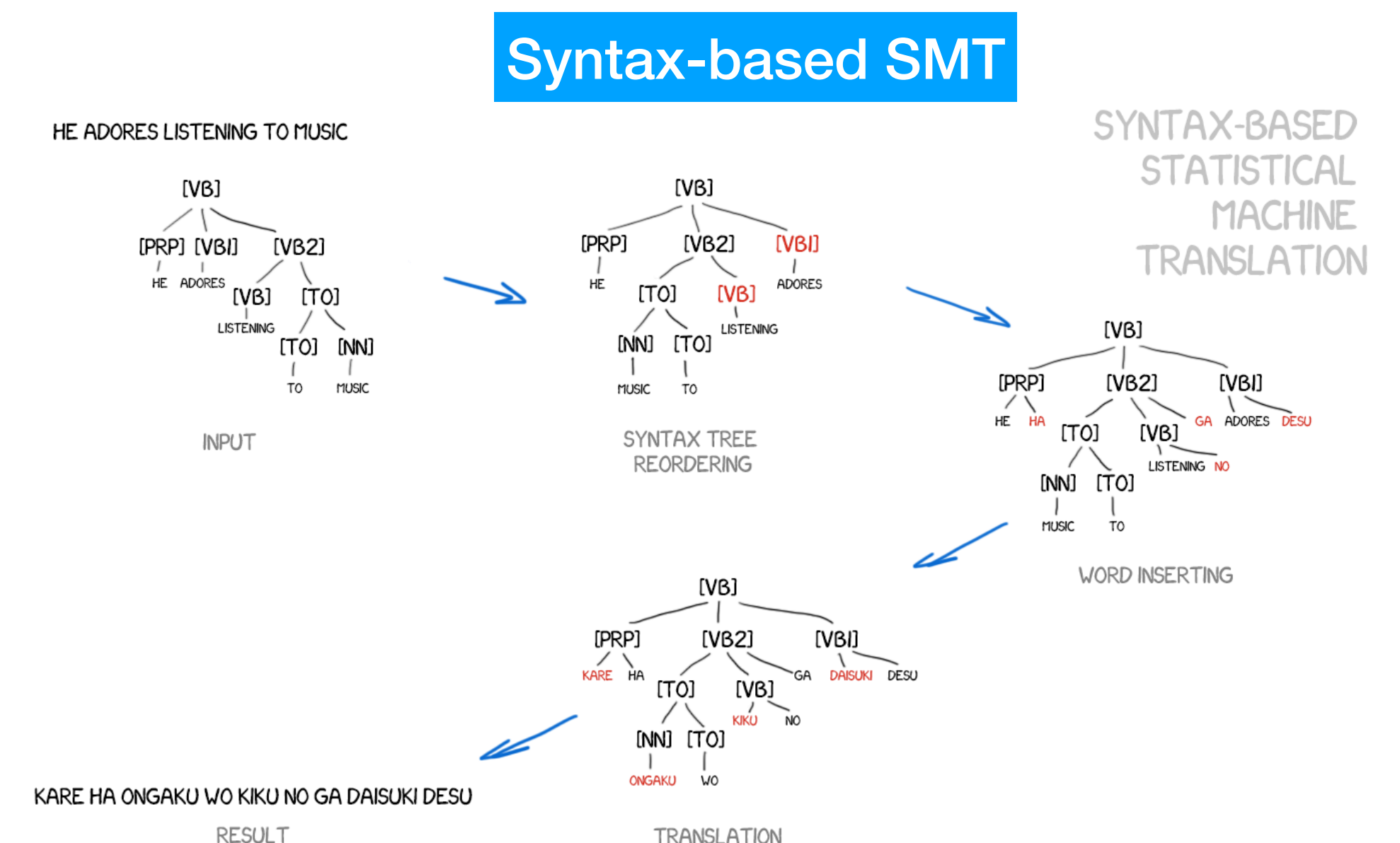
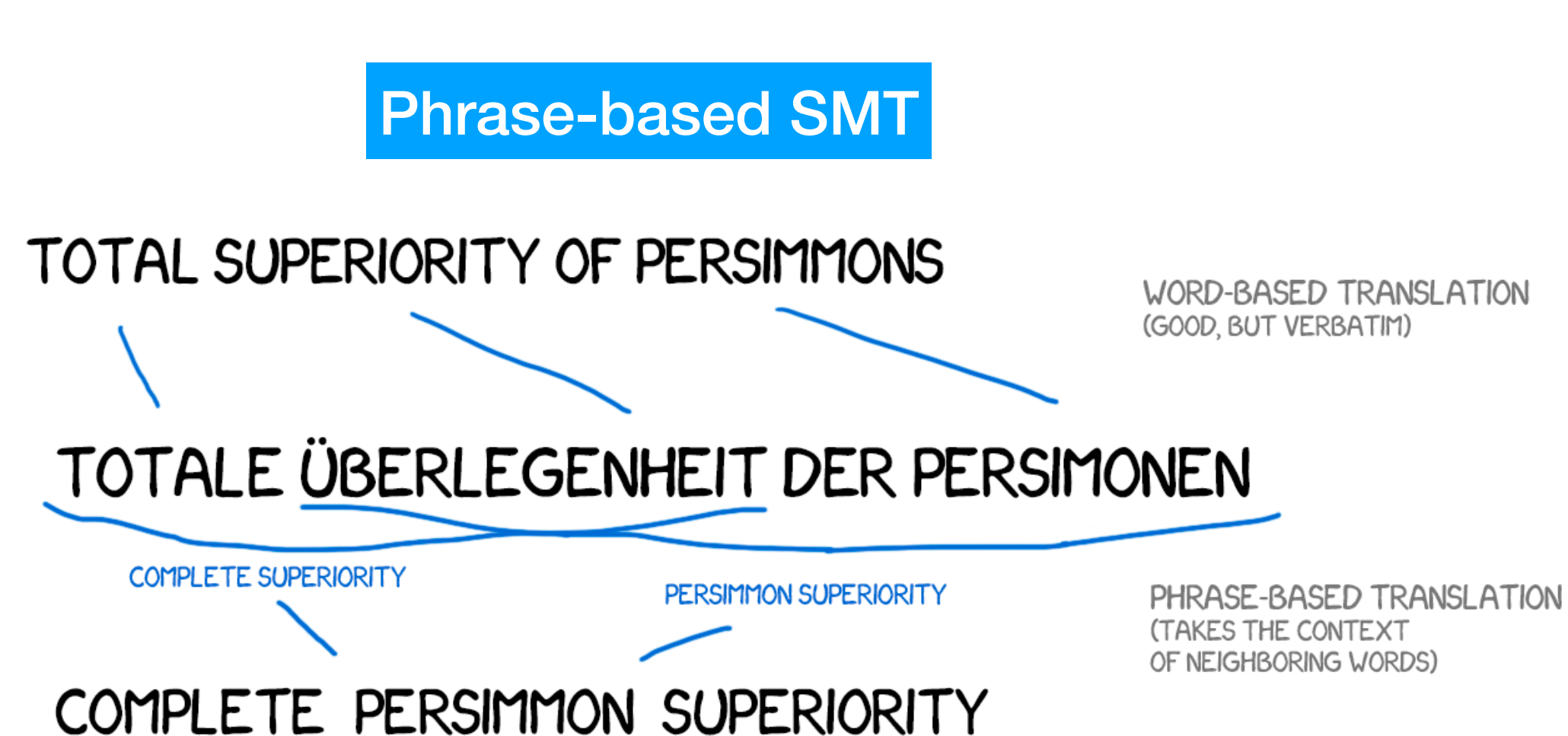
**good**  $\mathcal{A}(\mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \{(A, \emptyset), (Vinay, Vinay), (le, likes), (gusta, likes), (Python, Python)\}$ .

**bad**  $\mathcal{A}(\mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \{(A, Vinay), (Vinay, likes), (le, Python), (gusta, \emptyset), (Python, \emptyset)\}$ .

Examples: IBM models 1, 2, 3, 4, 5

# Statistical machine translation (SMT)

- SMT was a huge field (1990s-2010s) - The best systems were **extremely complex**
- Systems had many separately-designed subcomponents
  - Need to **design features** to capture particular language phenomena
  - Required compiling and maintaining **extra resources**
  - Lots of **human effort** to maintain - repeated effort for each language pair!



# SMT → NMT

Q. Do you know when Google Translate was first launched?

Launched in April 2006 as a [statistical machine translation](#) service, it used [United Nations](#) and [European Parliament](#) documents and transcripts to gather linguistic data. Rather than translating languages directly, it first translates text to English and then pivots to the target language in most of the language combinations it posits in its grid,<sup>[7]</sup> with a few exceptions including Catalan-Spanish.<sup>[8]</sup> During a translation, it looks for patterns in millions of documents to help decide which words to choose and how to arrange them in the target language. Its accuracy, which has been criticized on several occasions,<sup>[9]</sup> has been measured to vary greatly across languages.<sup>[10]</sup> In November 2016, Google announced that Google Translate would switch to a [neural machine translation](#) engine – [Google Neural Machine Translation](#) (GNMT) – which translates "whole sentences at a time,



# Google's NMT system in 2016

RESEARCH > PUBLICATIONS >

Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Table 10: Mean of side-by-side scores on production data

	PBMT	GNMT	Human	Relative Improvement
English → Spanish	4.885	5.428	5.504	87%
English → French	4.932	5.295	5.496	64%
English → Chinese	4.035	4.594	4.987	58%
Spanish → English	4.872	5.187	5.372	63%
French → English	5.046	5.343	5.404	83%
Chinese → English	3.694	4.263	4.636	60%



# SMT → NMT

1519年600名西班牙人在墨西哥登陆，去征服**几百万人口**的**阿兹特克帝国**，初次交锋他们**损兵三分之二**。

In 1519, six hundred Spaniards landed in Mexico to conquer **the Aztec Empire with a population of a few million**. They lost two thirds of their soldiers in the first clash.

[translate.google.com](https://translate.google.com) (2009): 1519 600 Spaniards landed in Mexico, **millions of people to conquer the Aztec empire**, the first two-thirds of soldiers against their loss.

[translate.google.com](https://translate.google.com) (2013): 1519 600 Spaniards landed in Mexico **to conquer the Aztec empire, hundreds of millions of people**, the initial confrontation loss of soldiers two-thirds.

[translate.google.com](https://translate.google.com) (2015): 1519 600 Spaniards landed in Mexico, **millions of people to conquer the Aztec empire**, the first two-thirds of the loss of soldiers they clash.

The screenshot shows the Google Translate interface with the source language set to Chinese (Simplified) and the target language set to English. The input text is: "1519年600名西班牙人在墨西哥登陆，去征服几百万人口的阿兹特克帝国，初次交锋他们损兵三分之二。". The output text is: "In 1519, 600 Spaniards landed in Mexico to conquer the Aztec Empire with a population of several million. They lost two-thirds of their troops in the first confrontation." The interface also shows a "Look up details" link and a "拼" (Pinyin) button.

# Neural machine translation (NMT)

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

---

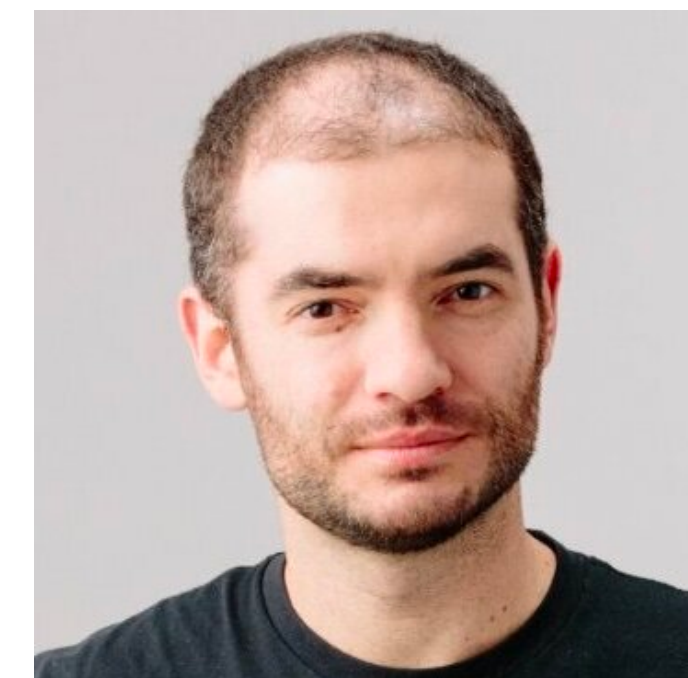
## Sequence to Sequence Learning with Neural Networks

---

**Ilya Sutskever**  
Google  
ilyasu@google.com

**Oriol Vinyals**  
Google  
vinyals@google.com

**Quoc V. Le**  
Google  
qvl@google.com

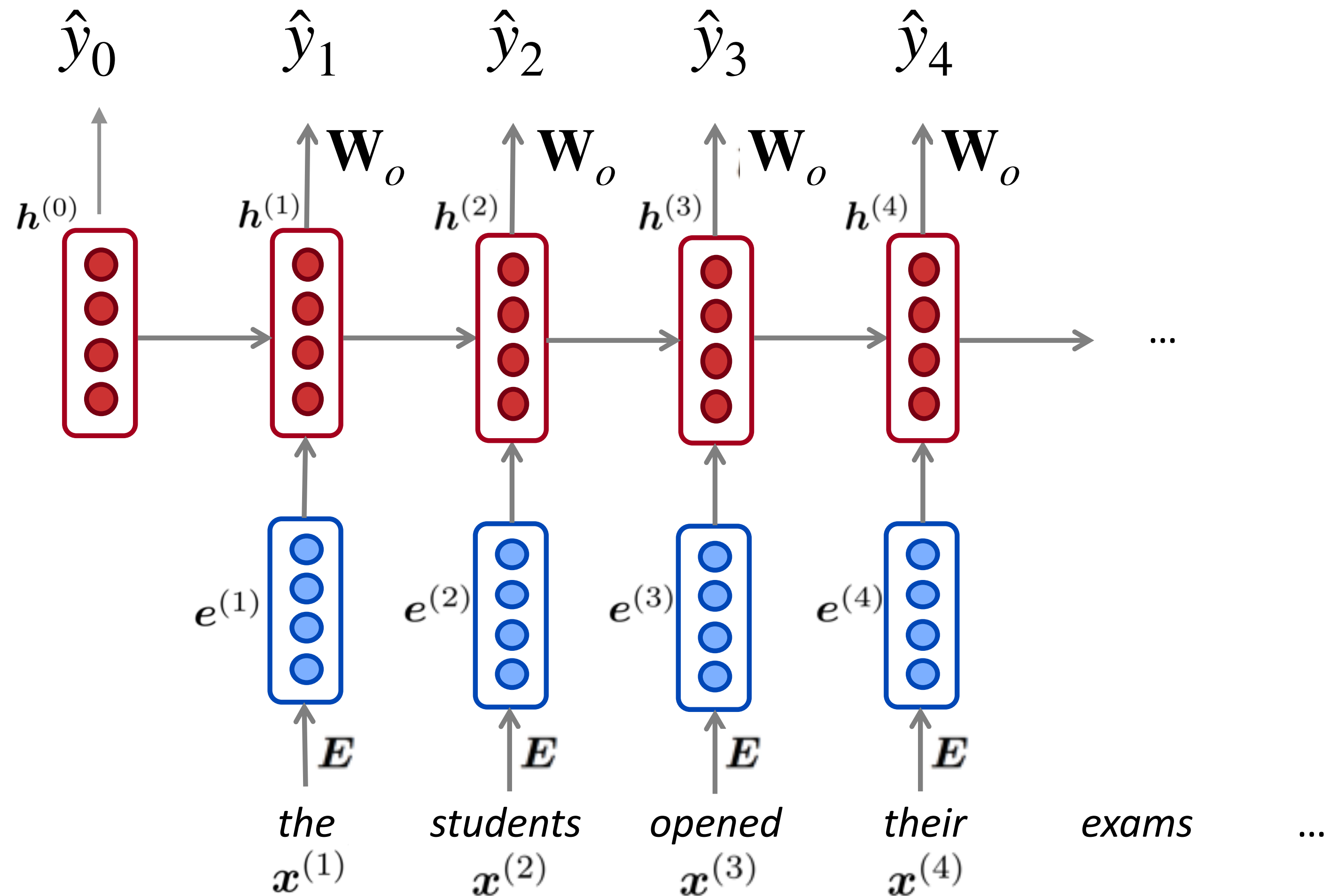


Ilya Sutskever

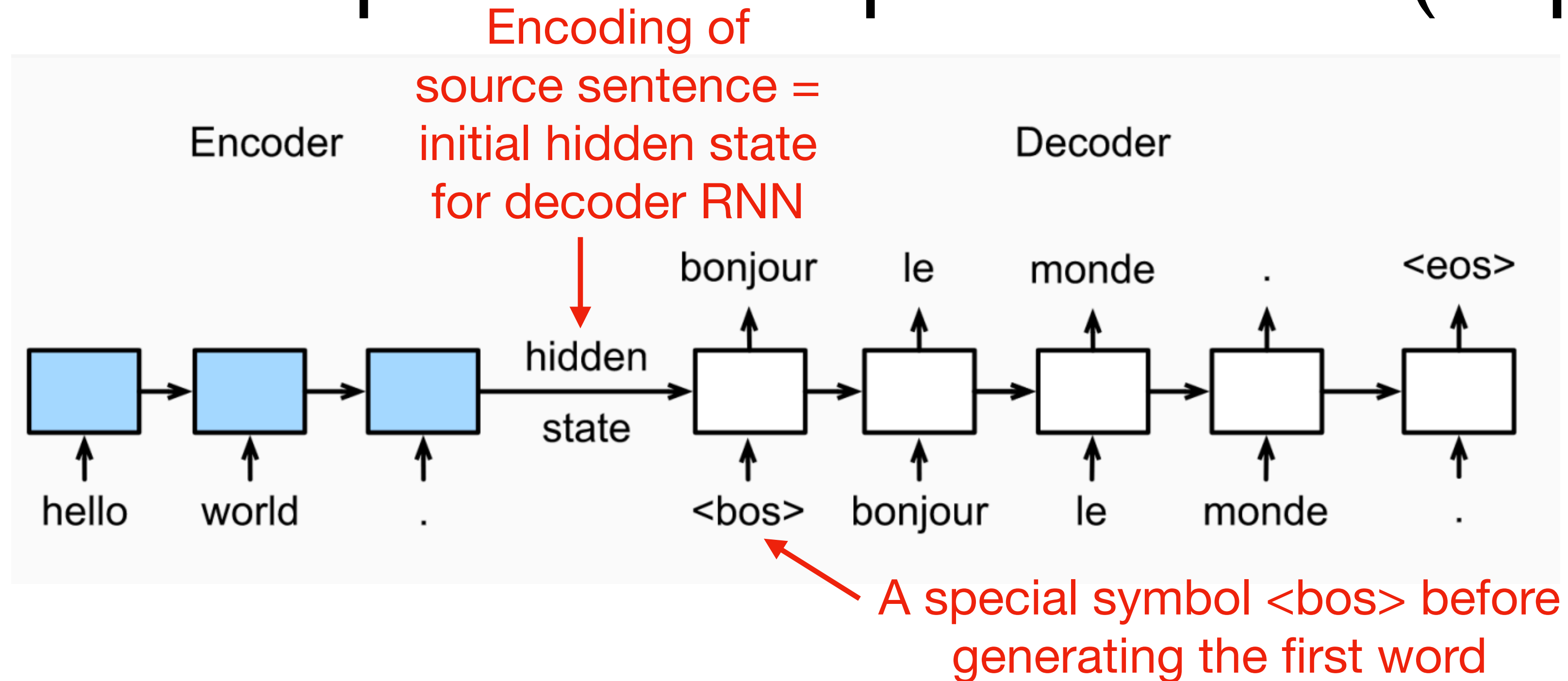
(Sutskever et al., 2014)



# Recall: RNNLMs



# The sequence-to-sequence model (seq2seq)

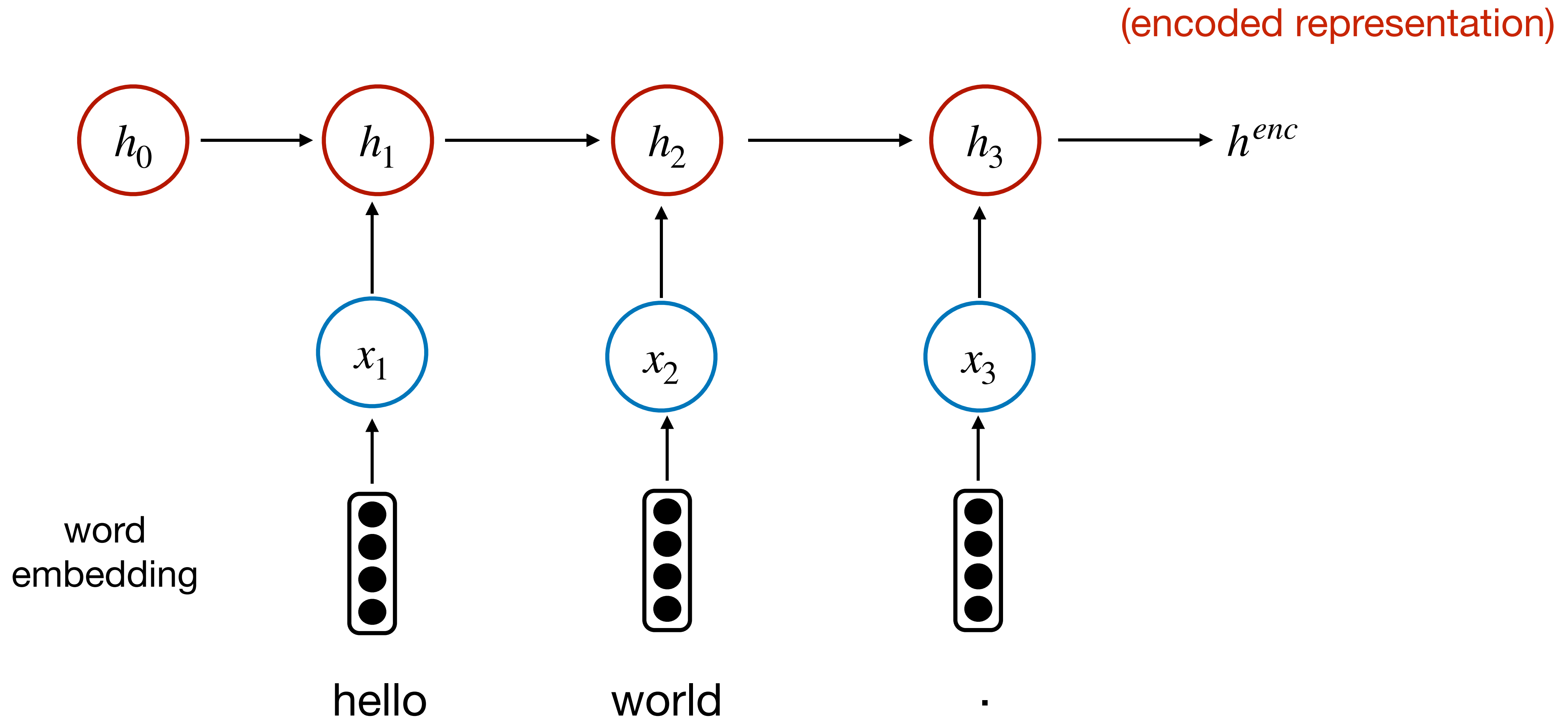


It is called an **encoder-decoder** architecture

- The encoder is an RNN to read the input sequence (source language)
- The decoder is another RNN to generate output word by word (target language)

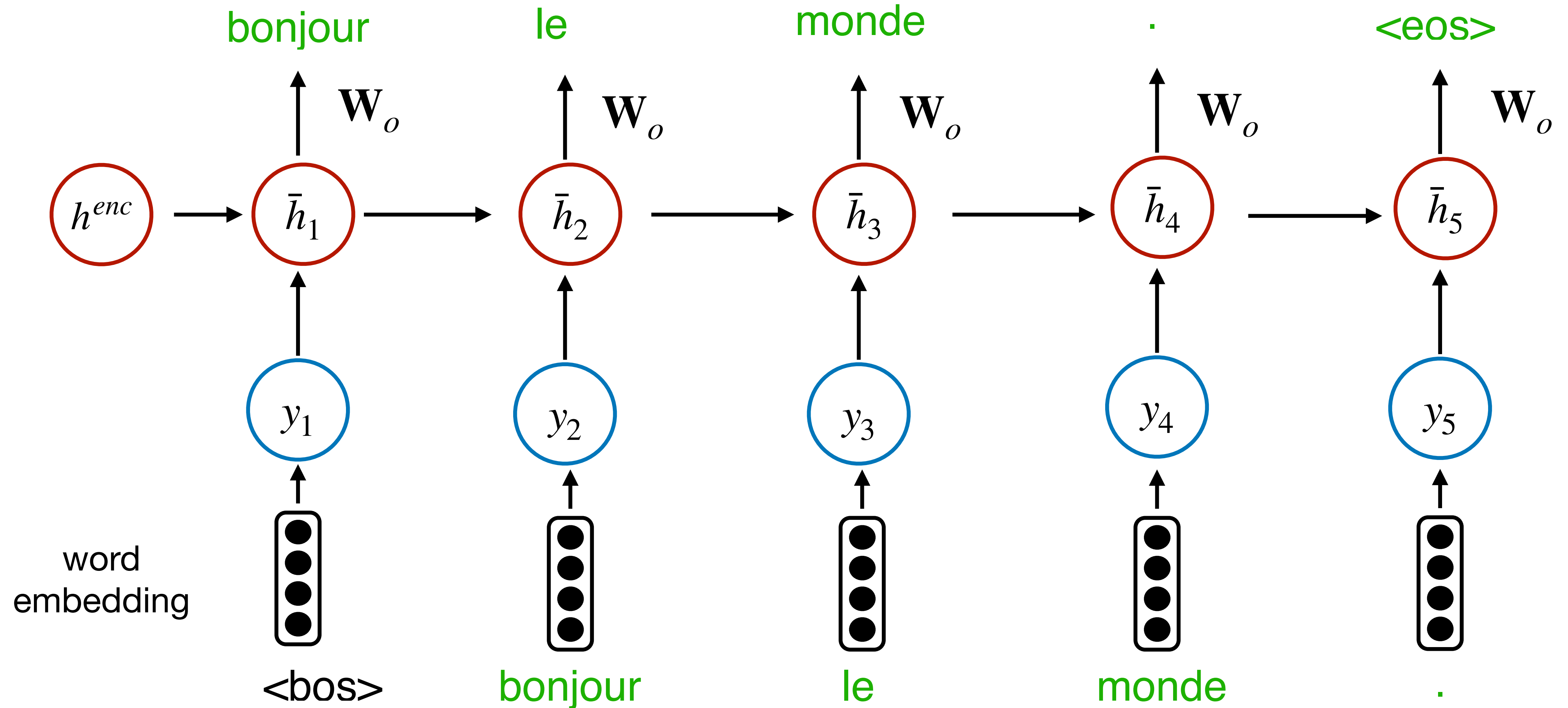
# Seq2seq: Encoder

*Sentence: hello world .*



# Seq2seq: Decoder

- A **conditional** language model

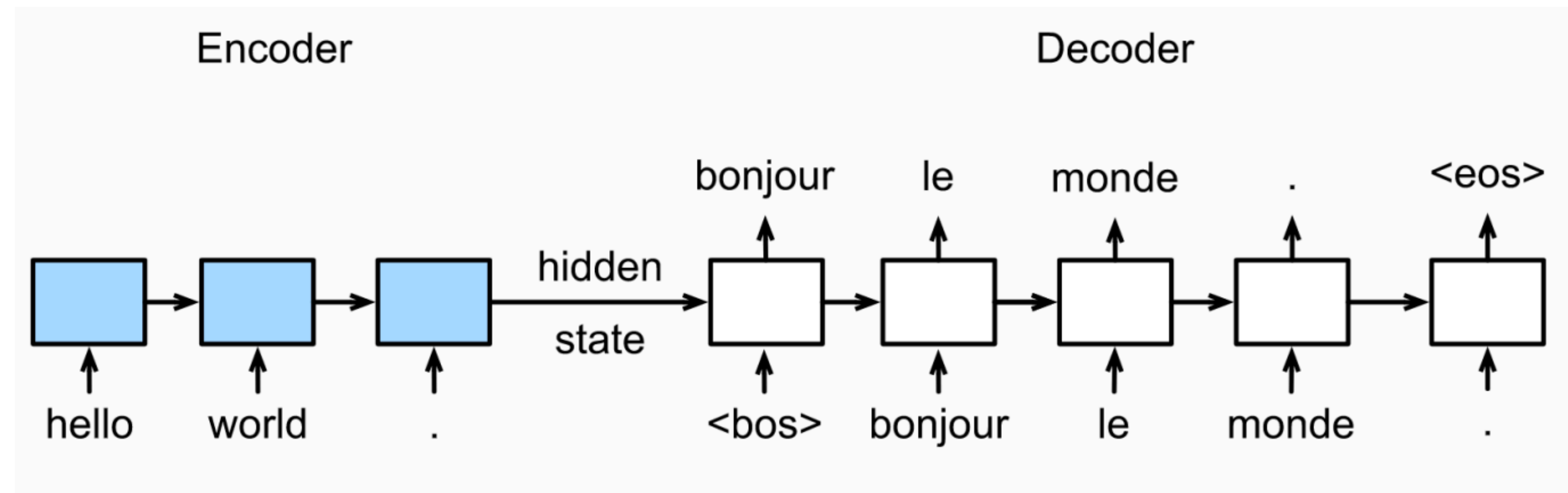


# Seq2seq: Decoder

- A **conditional** language model
  - It is a **language model** because the decoder is predicting the next word of the target sentence
  - **Conditional** because the predictions are also conditioned on the source sentence through  $h^{enc}$
- NMT directly calculates  $P(\mathbf{w}^{(t)} \mid \mathbf{w}^{(s)})$ 
  - Denote  $\mathbf{w}^{(t)} = y_1, \dots, y_T$

$$P(\mathbf{w}^{(t)} \mid \mathbf{w}^{(s)}) = P(y_1 \mid \mathbf{w}^{(s)})P(y_2 \mid y_1, \mathbf{w}^{(s)})P(y_3 \mid y_1, y_2, \mathbf{w}^{(s)}) \dots P(y_T \mid y_1, \dots, y_{T-1}, \mathbf{w}^{(s)})$$

# Understanding seq2seq



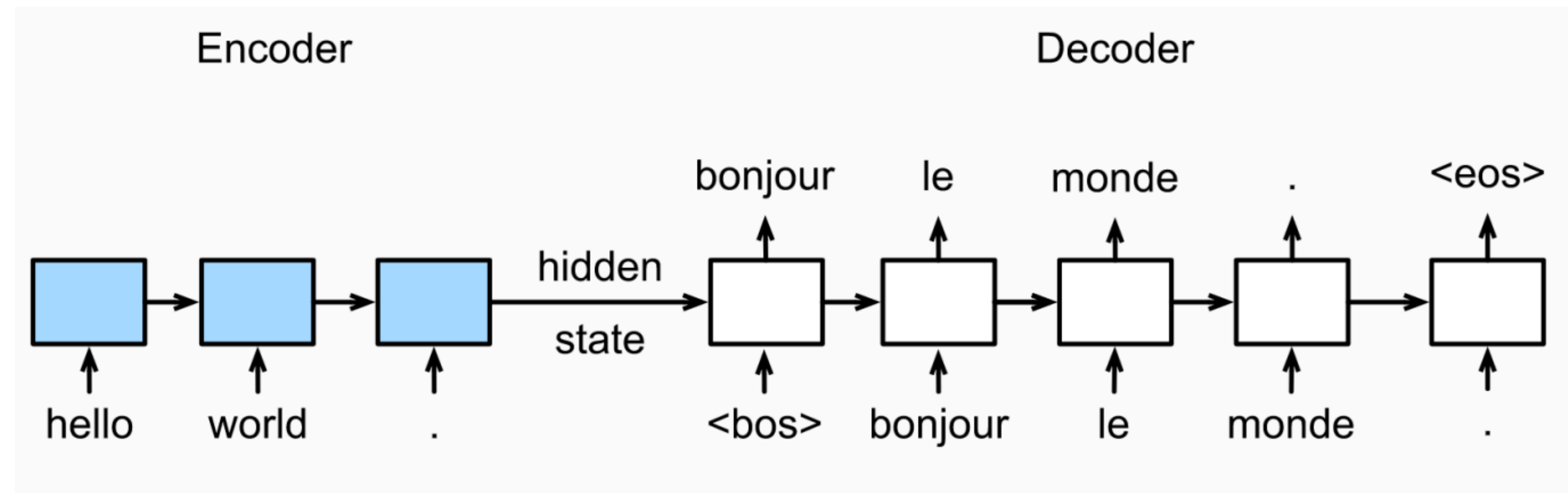
Which of the following is correct?

- (A) We can use bidirectional RNNs for both encoder and decoder
- (B) The decoder has more parameters because of the output matrix  $\mathbf{W}_o$
- (C) The encoder and decoder have separate word embeddings
- (D) The encoder and decoder's parameters are optimized together

Both (C) and (D) are correct.



# Understanding seq2seq



## Encoder RNN:

- word embeddings  $\mathbf{E}^{(s)}$  for source language
- RNN parameters, e.g.,  $\{\mathbf{W}, \mathbf{U}, \mathbf{b}\}$  for simple RNNs and 4x parameters for LSTMs
- Encoder RNN can be bidirectional!

## Decoder RNN:

- word embeddings  $\mathbf{E}^{(t)}$  for target language
- RNN parameters, e.g.,  $\{\mathbf{W}, \mathbf{U}, \mathbf{b}\}$  for simple RNNs and 4x parameters for LSTMs
- Output embedding matrix  $\mathbf{W}_o =$  can be tied with  $\mathbf{E}^{(t)}$
- **Decoder RNN has to be unidirectional (left to right)!**

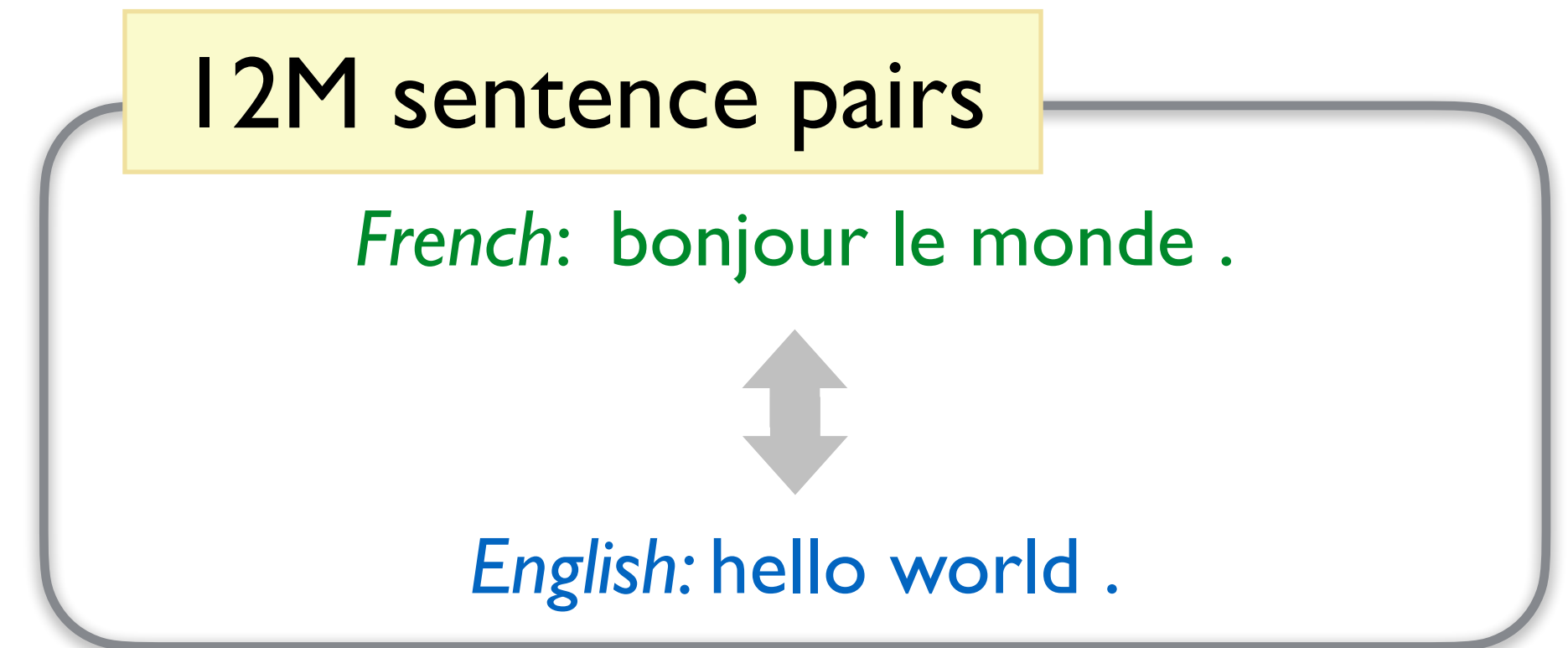
# Training seq2seq models

- Training data: parallel corpus  $\{(\mathbf{w}_i^{(s)}, \mathbf{w}_i^{(t)})\}$
- Minimize cross-entropy loss:

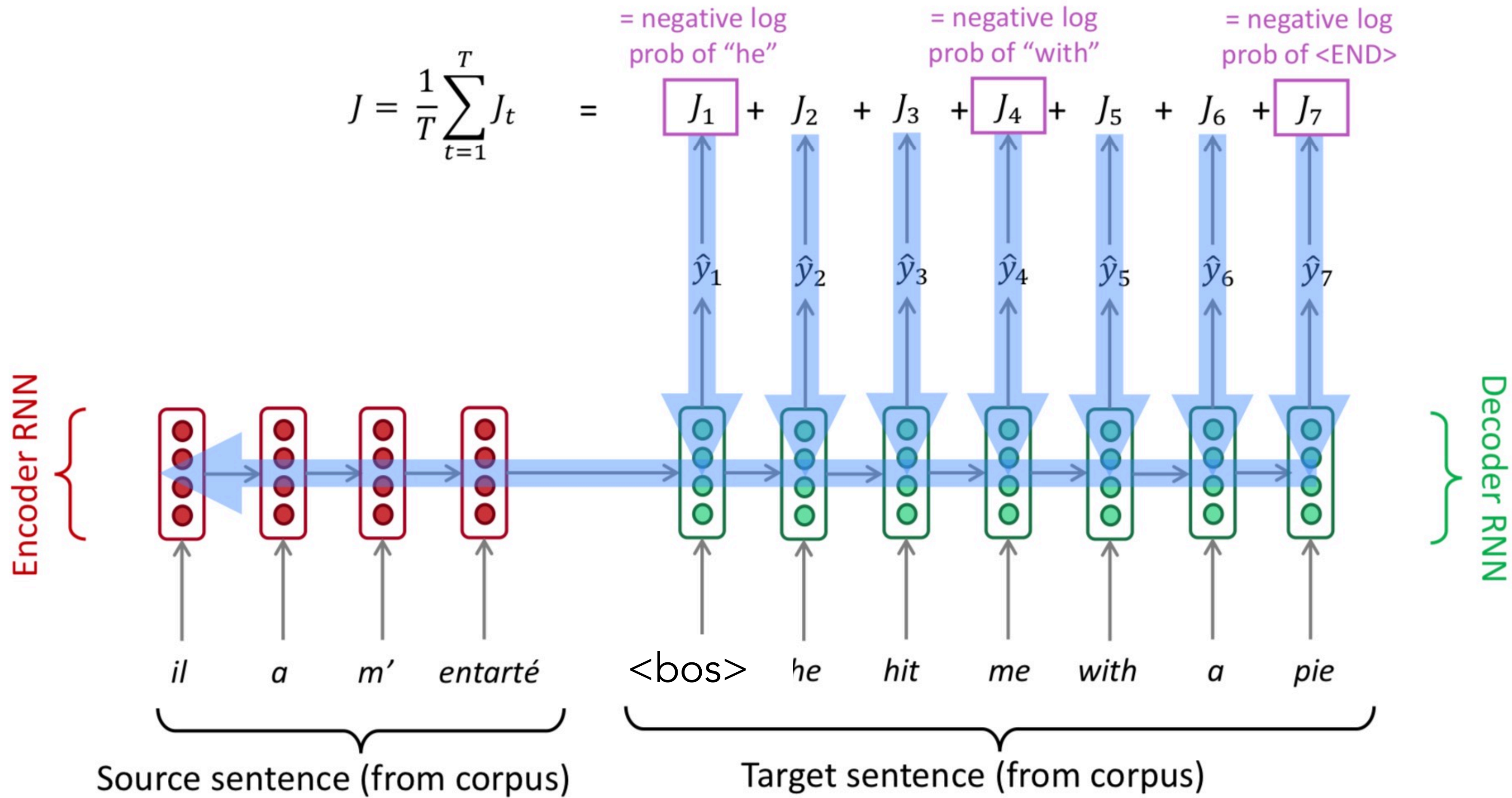
$$\sum_{t=1}^T -\log P(y_t | y_1, \dots, y_{t-1}, \mathbf{w}^{(s)})$$

(denote  $\mathbf{w}^{(t)} = y_1, \dots, y_T$ )

- Back-propagate gradients through both encoder and decoder



# Training seq2seq models

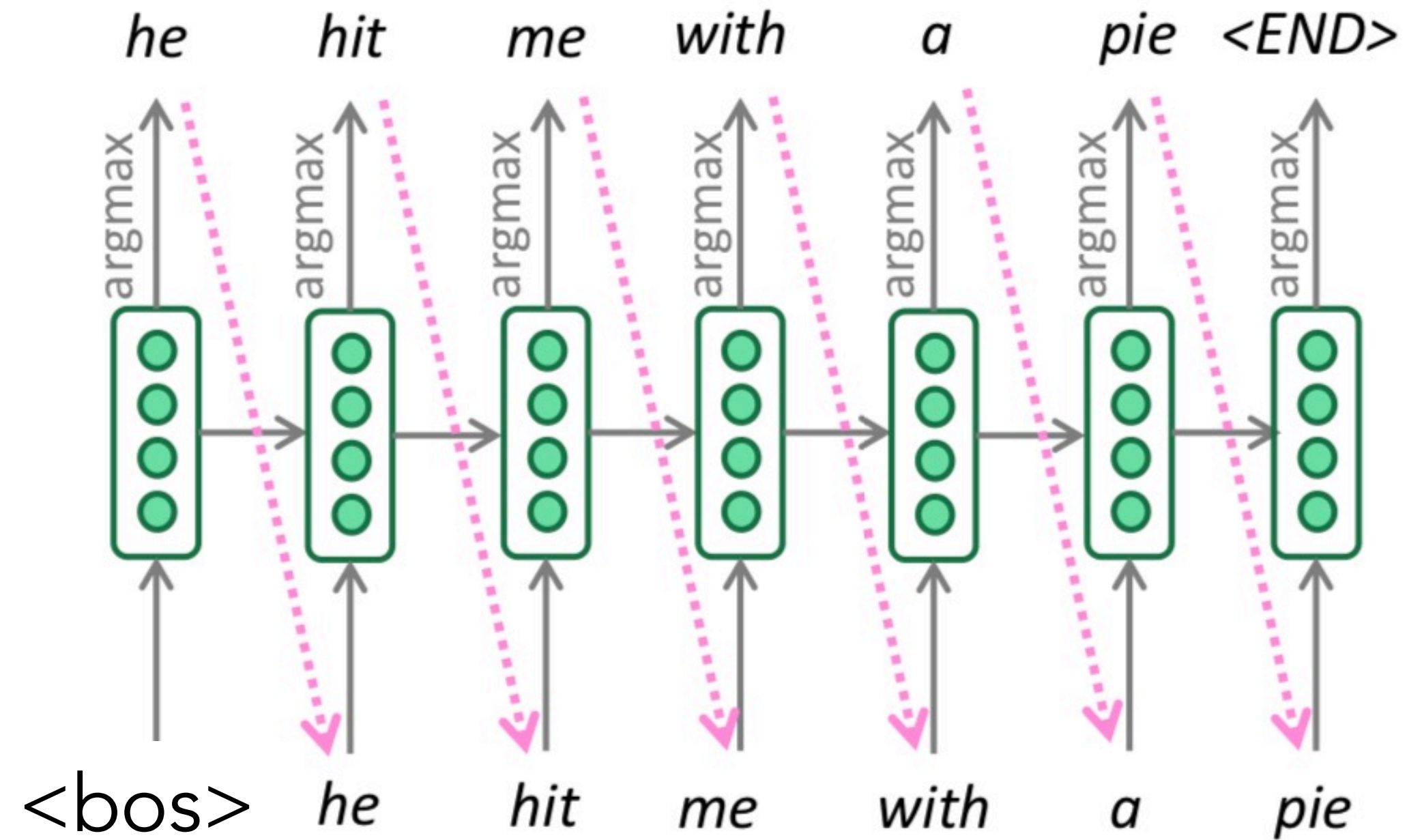


Seq2seq is optimized as a **single system**.  
Backpropagation operates "end-to-end".



# Decoding seq2seq models

- Greedy decoding
  - = Compute argmax at every step of decoder to generate word



- Exhaustive search is very expensive:  $\arg \max_{y_1, \dots, y_T} P(y_1, \dots, y_T | \mathbf{w}^{(s)})$  - we even don't know what T is

# Decoding with beam search

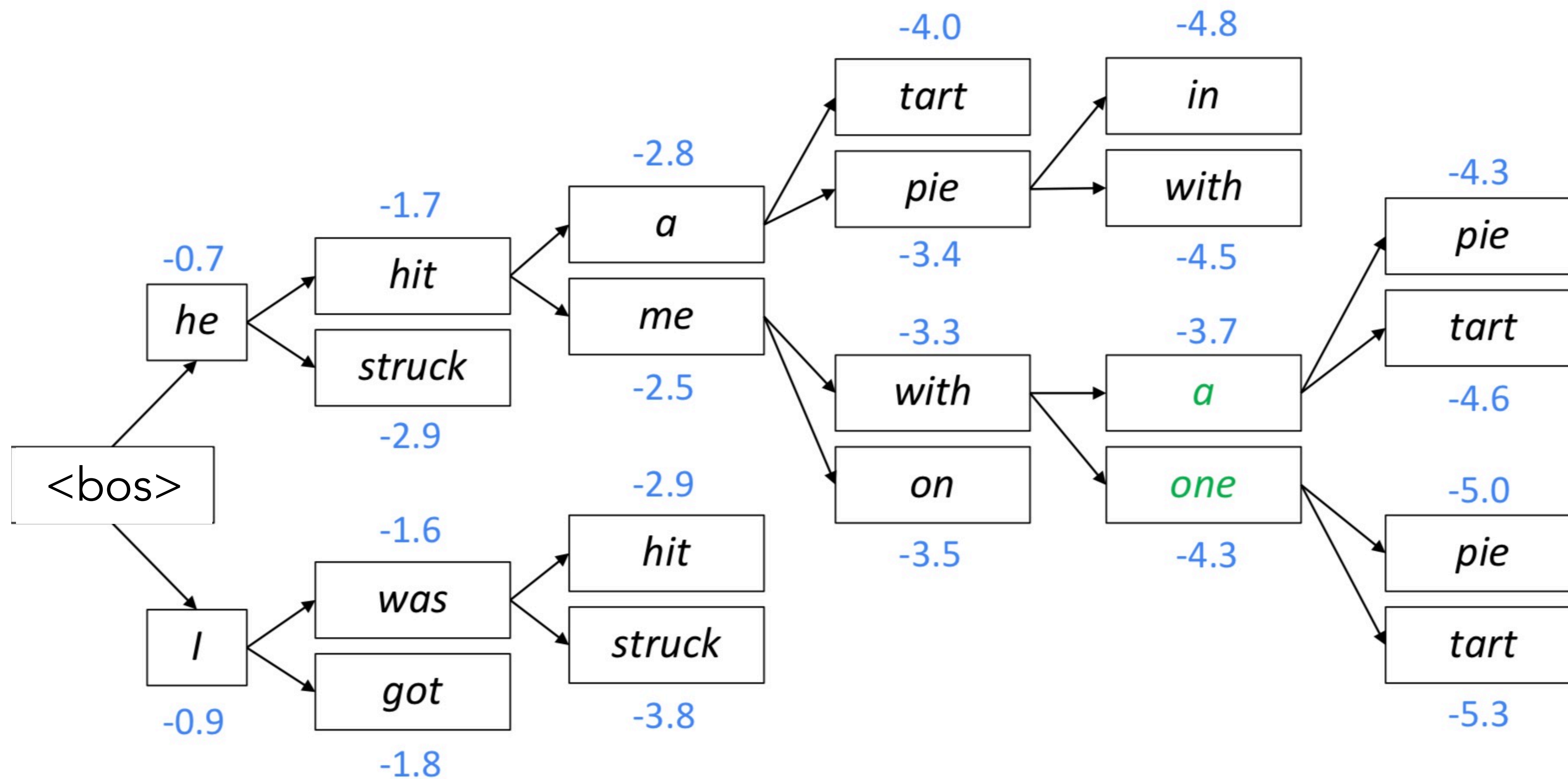
- At every step, keep track of the  $k$  most probable partial translations (hypotheses)
- Score of each hypothesis = log probability of sequence so far

$$\sum_{t=1}^j \log P(y_t | y_1, \dots, y_{t-1}, \mathbf{w}^{(s)})$$

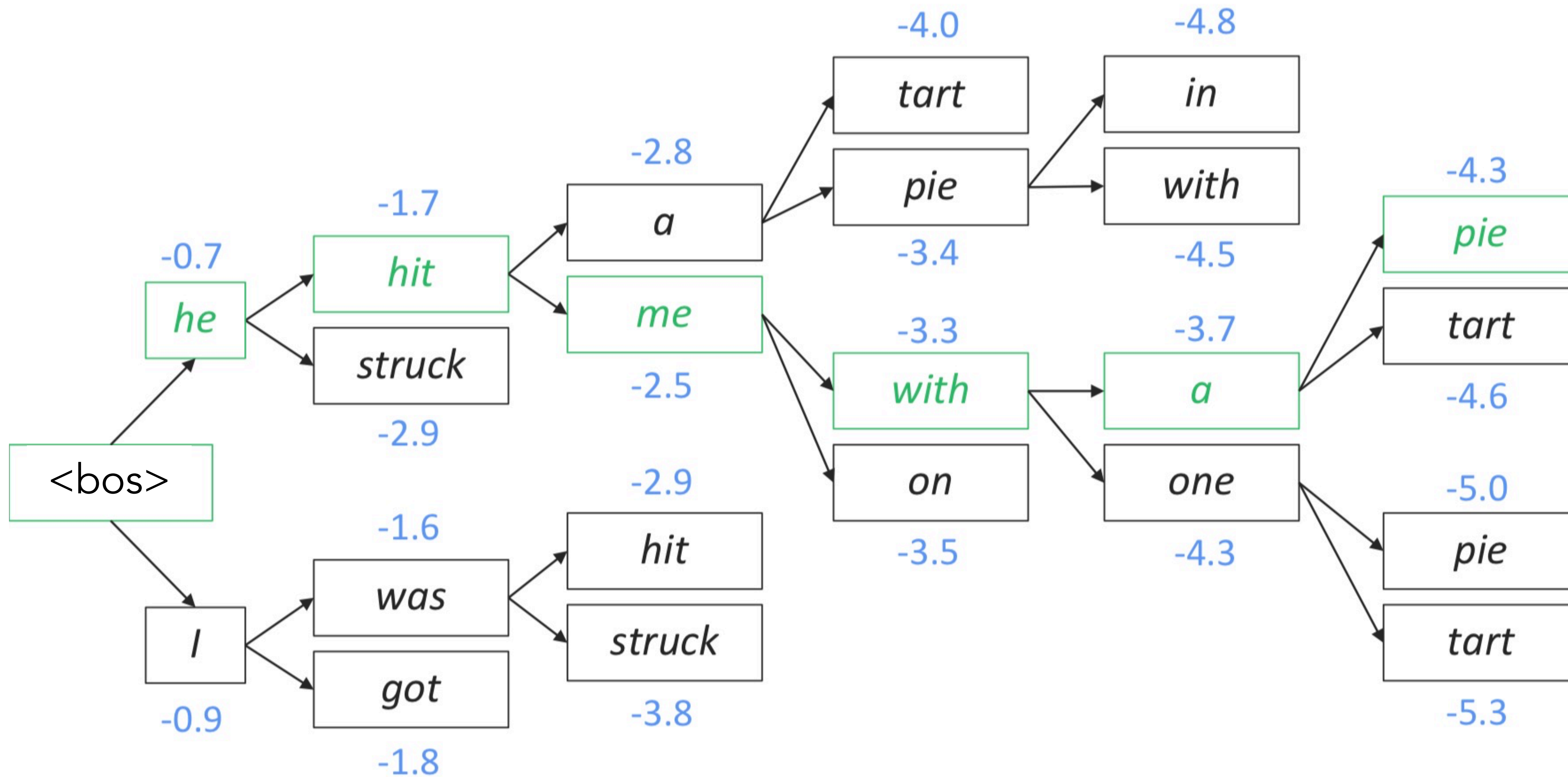
- Not guaranteed to be optimal
- Works better than greedy decoding in practice

# Beam search

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



# Beam search: Backtrack





# Beam search: details

- ▶ Different hypotheses may produce  $\langle eos \rangle$  token at different time steps
  - ▶ When a hypothesis produces  $\langle eos \rangle$ , stop expanding it and place it aside
- ▶ Continue beam search until:
  - ▶ All  $k$  hypotheses produce  $\langle eos \rangle$  OR
  - ▶ Hit max decoding limit  $T$
- ▶ Select top hypotheses using the *normalized* likelihood score

$$\frac{1}{T} \sum_{t=1}^T \log P(y_t | y_1, \dots, y_{t-1}, \mathbf{w}^{(s)})$$

- ▶ Otherwise shorter hypotheses have higher scores