



COS 484/584

LI 6: Neural Machine Translation - I

Spring 2021

Last time: IBM Model I

- Assume $p(a_m | m, M^{(s)}, M^{(t)}) = \frac{1}{M^{(t)}}$

- We then have:

$$p(w^{(s)}, w^{(t)}) = p(w^{(t)}) \sum_A \left(\frac{1}{M^{(t)}}\right)^{M^{(s)}} p(w^{(s)} | w^{(t)})$$

- How do we estimate $p(w^{(s)} = v | w^{(t)} = u)$?

IBM Model I

- If we have word-to-word alignments, we can compute the probabilities using the MLE:
- $$p(v | u) = \frac{\text{count}(u, v)}{\text{count}(u)}$$
- where $\text{count}(u, v) = \# \text{instances where target word } u \text{ was aligned to source word } v \text{ in the training set}$
- However, word-to-word alignments are often hard to come by

What can we do?

EM for Model I



- **(E-Step)** If we had an accurate translation model, we can estimate likelihood of each alignment as:

$$q_m(a_m | \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) \propto p(a_m | m, M^{(s)}, M^{(t)}) \times p(w_m^{(s)} | w_{a_m}^{(t)}),$$

Remember
these are
fixed

- **(M Step)** Use expected count to re-estimate translation parameters:

How would you compute the new probabilities $p(v | u)$?

A) $p(v | u) = \frac{E_q[\text{count}(u, v)]}{\text{count}(u)}$

B) $p(v | u) = \frac{E_q[\text{count}(u, v)]}{\text{count}(v)}$

C) $p(v | u) = E_q[\text{count}(u, v)]$

EM for Model I

- **(E-Step)** If we had an accurate translation model, we can estimate likelihood of each alignment as:

$$q_m(a_m | \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) \propto p(a_m | m, M^{(s)}, M^{(t)}) \times p(w_m^{(s)} | w_{a_m}^{(t)}),$$

Remember
these are
fixed

- **(M Step)** Use expected count to re-estimate translation parameters:

$$E_q [\text{count}(u, v)] = \sum_m q_m(a_m | \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) \times \delta(w_m^{(s)} = v) \times \delta(w_{a_m}^{(t)} = u).$$

$$p(v | u) = \frac{E_q[\text{count}(u, v)]}{\text{count}(u)}$$

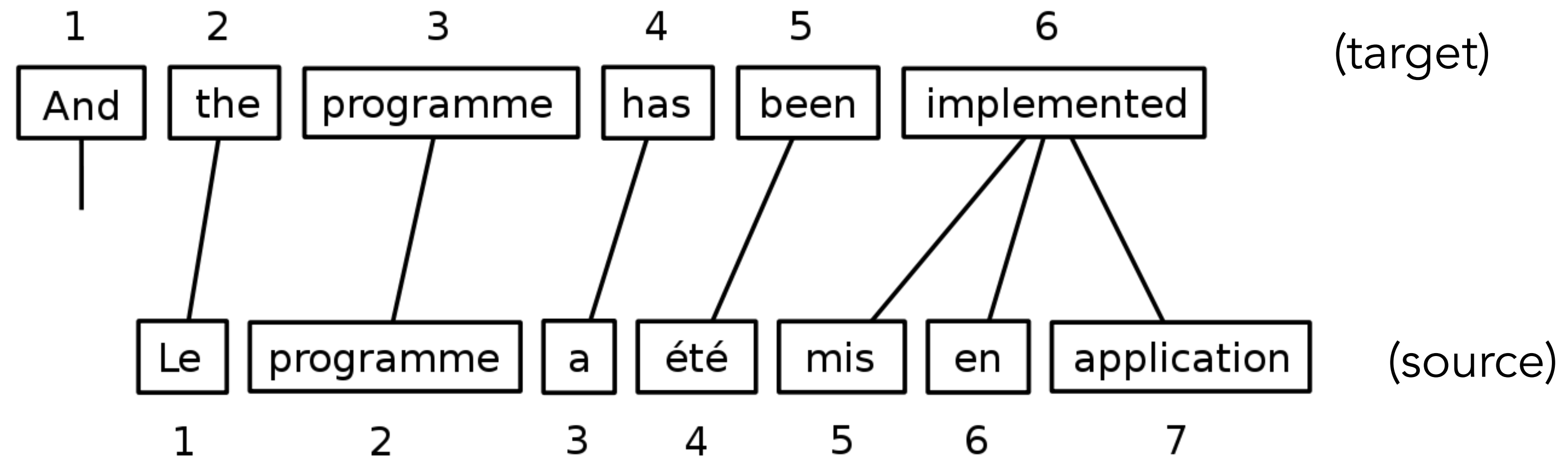
Decoding: How do we translate?

- We want: $\arg \max_{w^{(t)}} p(w^{(t)} | w^{(s)}) = \arg \max_{w^{(t)}} \frac{p(w^{(s)}, w^{(t)})}{p(w^{(s)})}$
- Sum over all possible alignments:

$$\begin{aligned} p(w^{(s)}, w^{(t)}) &= \sum_{\mathcal{A}} p(w^{(s)}, w^{(t)}, \mathcal{A}) \\ &= p(w^{(t)}) \sum_{\mathcal{A}} p(\mathcal{A}) \times p(w^{(s)} | w^{(t)}, \mathcal{A}) \end{aligned}$$

- Alternatively, take the max over alignments
- Decoding: Greedy/beam search

Model I: Decoding



At every step m , pick target word $w_m^{(t)}$ to maximize product of:

1. Language model: $p_{LM}(w_m^{(t)} | w_1^{(t)}, \dots, w_{m-1}^{(t)})$

2. Translation model: $p(w_{b_m}^{(s)} | w_m^{(t)})$

where b_m is the inverse alignment from target to source

IBM Model I

- Assume $p(a_m | m, M^{(s)}, M^{(t)}) = \frac{1}{M^{(t)}}$
- Each source word is aligned to at most one target word
- We then have:

$$p(w^{(s)}, w^{(t)}) = p(w^{(t)}) \sum_A \left(\frac{1}{M^{(t)}}\right)^{M^{(s)}} p(w^{(s)} | w^{(t)})$$

Restrictive assumptions

IBM Model 2

- Slightly relaxed assumption:
 - $p(a_m | m, M^{(s)}, M^{(t)})$ is also estimated/learned
- Some independence assumptions from Model 1 still required:
 - Alignment probability factors across tokens:

$$p(\mathcal{A} | \mathbf{w}^{(s)}, \mathbf{w}^{(t)}) = \prod_{m=1}^{M^{(s)}} p(a_m | m, M^{(s)}, M^{(t)}).$$

- Translation probability factors across tokens:

$$p(\mathbf{w}^{(s)} | \mathbf{w}^{(t)}, \mathcal{A}) = \prod_{m=1}^{M^{(s)}} p(w_m^{(s)} | w_{a_m}^{(t)}),$$

Other IBM models

Model 1: lexical translation

Model 2: additional absolute alignment model

Model 3: extra fertility model

Model 4: added relative alignment model

Model 5: fixed deficiency problem.

Model 6: Model 4 combined with a [HMM](#) alignment model in a log linear way

- Models 3 - 6 make successively weaker assumptions
 - But get progressively harder to optimize
- Simpler models are often used to 'initialize' complex ones
 - e.g train Model 1 and use it to initialize Model 2 translation parameters

Phrase-based MT

(literal)

Nous allons prendre un verre
We will take a glass

(actual)

We'll have a drink

	<i>Nous</i>	<i>allons</i>	<i>prendre</i>	<i>une</i>	<i>verre</i>
We'll					
have					
a					
drink					

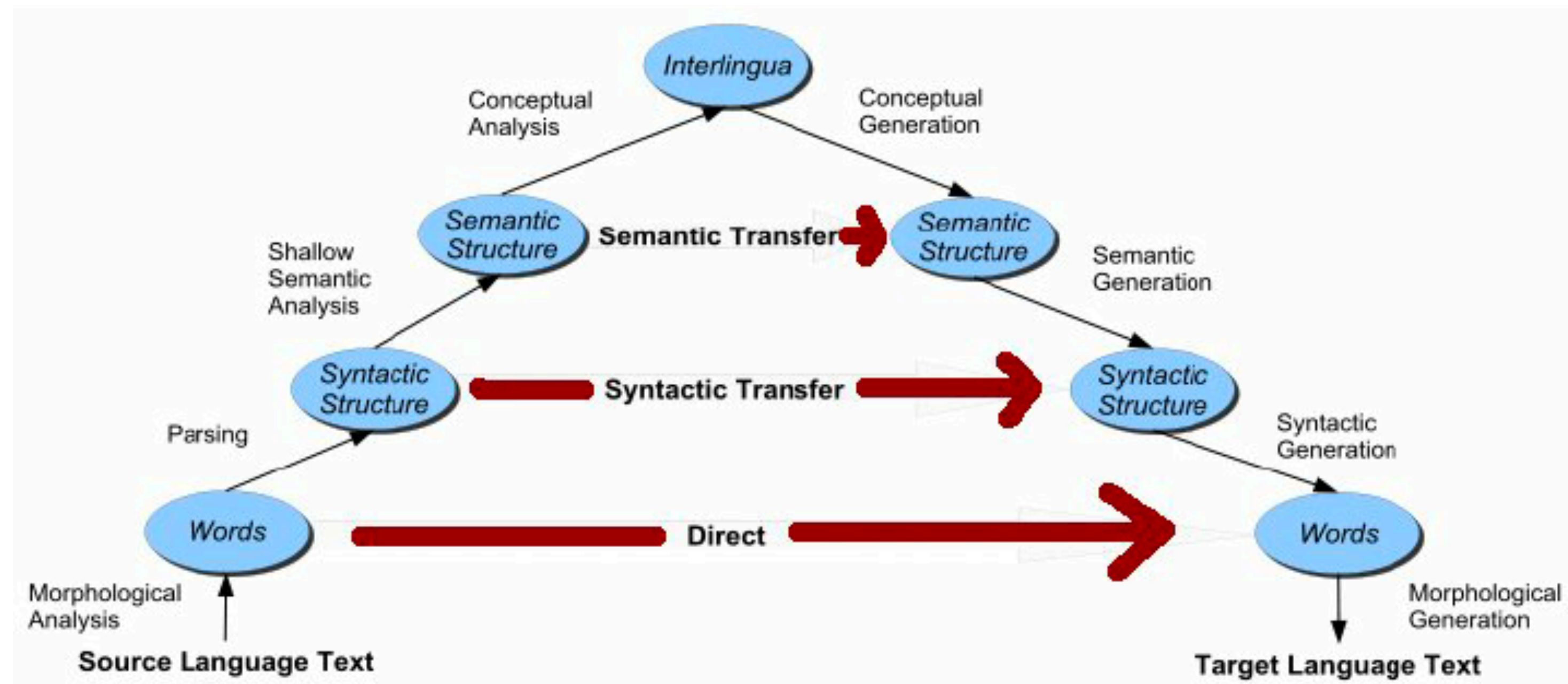
- Word-by-word translation is not sufficient in many cases
- Solution: build alignments and translation tables between multiword spans or "phrases"

Phrase-based MT

- Solution: build alignments and translation tables between multiword spans or “phrases”
- **Translations** condition on multi-word units and assign probabilities to multi-word units
- **Alignments** map from spans to spans

$$p(\mathbf{w}^{(s)} \mid \mathbf{w}^{(t)}, \mathcal{A}) = \prod_{((i,j),(k,\ell)) \in \mathcal{A}} p_{\mathbf{w}^{(s)}|\mathbf{w}^{(t)}}(\{w_{i+1}^{(s)}, w_{i+2}^{(s)}, \dots, w_j^{(s)}\} \mid \{w_{k+1}^{(t)}, w_{k+2}^{(t)}, \dots, w_\ell^{(t)}\})$$

Vauquois Pyramid



- Hierarchy of concepts and distances between them in different languages
- Lowest level: individual words/characters
- Higher levels: syntax, semantics
- Interlingua: Generic language-agnostic representation of meaning

Syntactic MT

- ▶ Rather than use phrases, use a *synchronous context-free grammar*: constructs “parallel” trees in two languages simultaneously

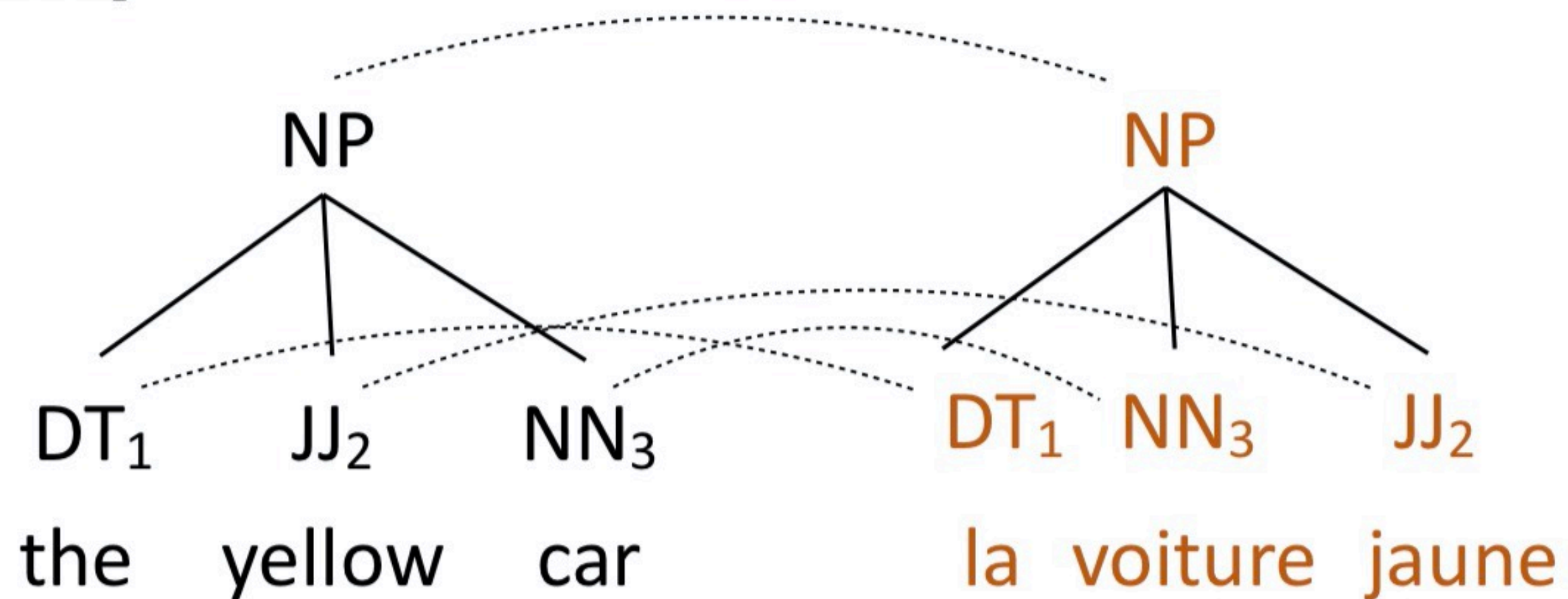
NP \rightarrow [DT₁ JJ₂ NN₃; DT₁ NN₃ JJ₂]

DT \rightarrow [the, la]

DT \rightarrow [the, le]

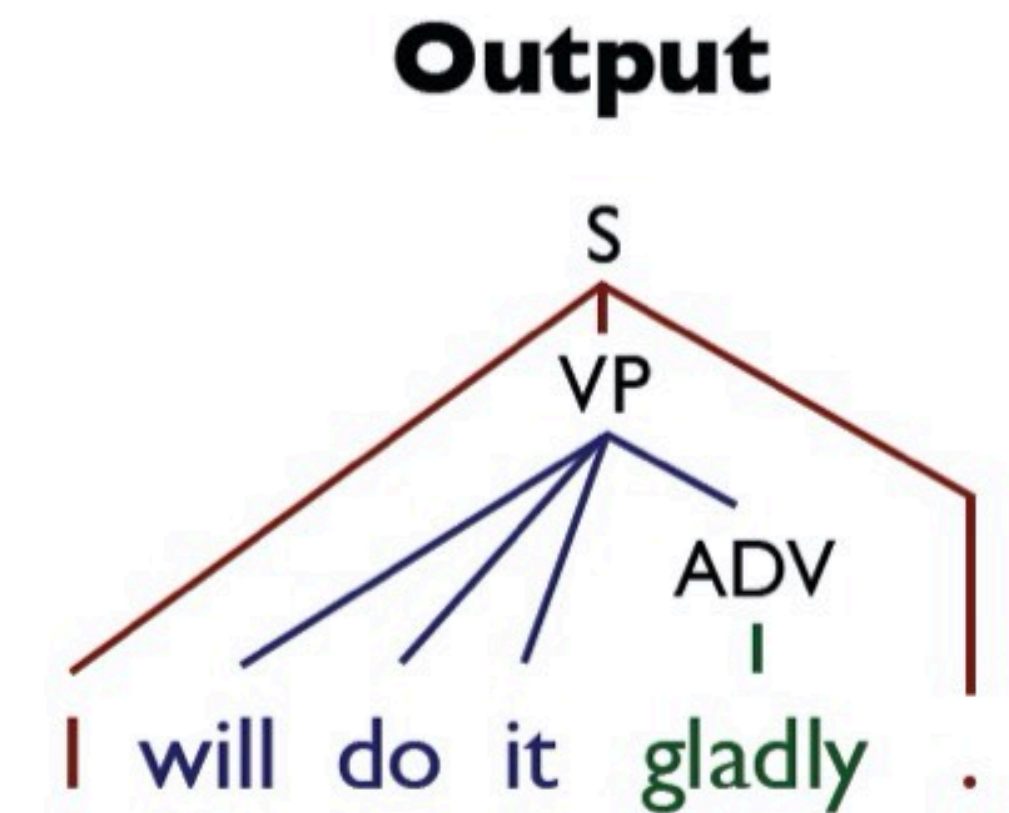
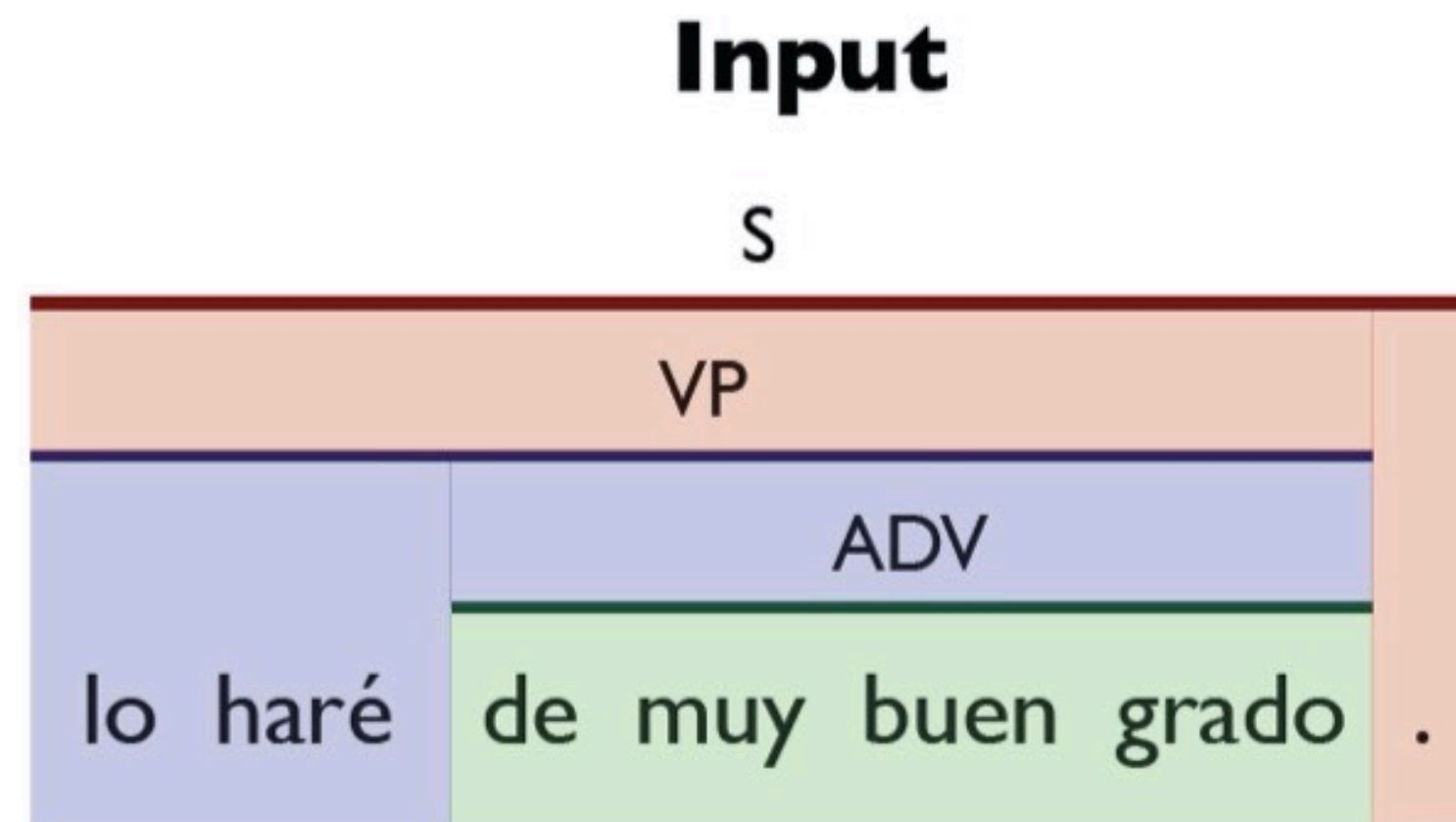
NN \rightarrow [car, voiture]

JJ \rightarrow [yellow, jaune]



- ▶ Assumes parallel syntax up to reordering
- ▶ Translation = parse the input with “half” the grammar, read off other half

Syntactic MT



Grammar

- ▶ Relax this by using lexicalized rules, like “syntactic phrases”
- ▶ Leads to HUGE grammars, parsing is slow

$S \rightarrow \langle VP . ; I VP . \rangle$ **OR** $S \rightarrow \langle VP . ; you VP . \rangle$
 $VP \rightarrow \langle lo haré ADV ; will do it ADV \rangle$
 $s \rightarrow \langle lo haré ADV . ; I will do it ADV . \rangle$
 $ADV \rightarrow \langle de muy buen grado ; gladly \rangle$

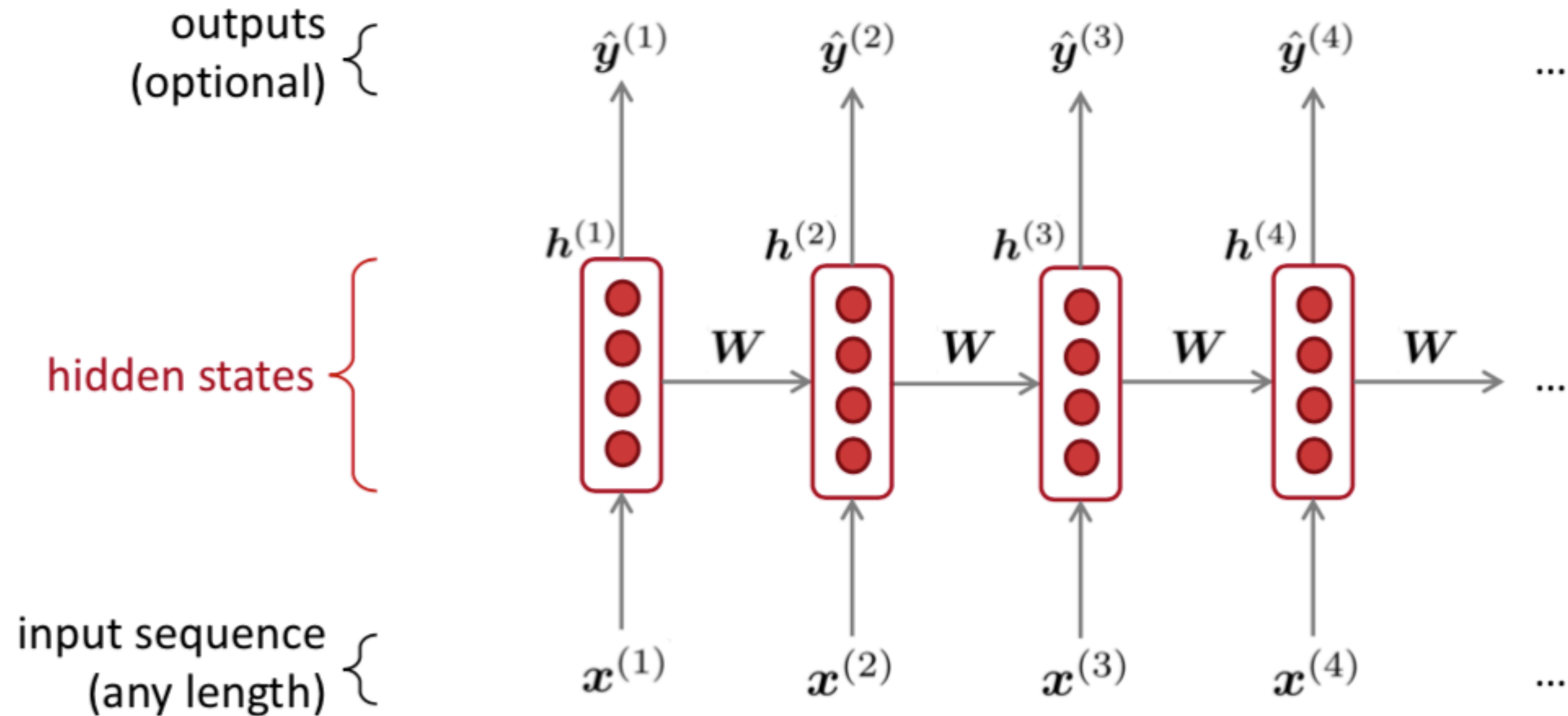
Neural Machine Translation

Neural Machine Translation

- ▶ A **single neural network** is used to translate from source to target language
- ▶ Architecture: Encoder-Decoder
 - ▶ Two main components:
 - ▶ **Encoder**: Convert source sentence (input) into a vector/matrix
 - ▶ **Decoder**: Convert encoding into a sentence in target language (output)

Recall: RNNs

$$\mathbf{h}_t = g(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b}) \in \mathbb{R}^d$$



Recall: RNNs



$$\mathbf{h}_t = g(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{x}_t + \mathbf{b}) \in \mathbb{R}^d$$

outputs
(optional) { $\hat{y}^{(1)}$ $\hat{y}^{(2)}$ $\hat{y}^{(3)}$ $\hat{y}^{(4)}$...

What is the maximum sequence length an RNN could theoretically take as input?

A) 10

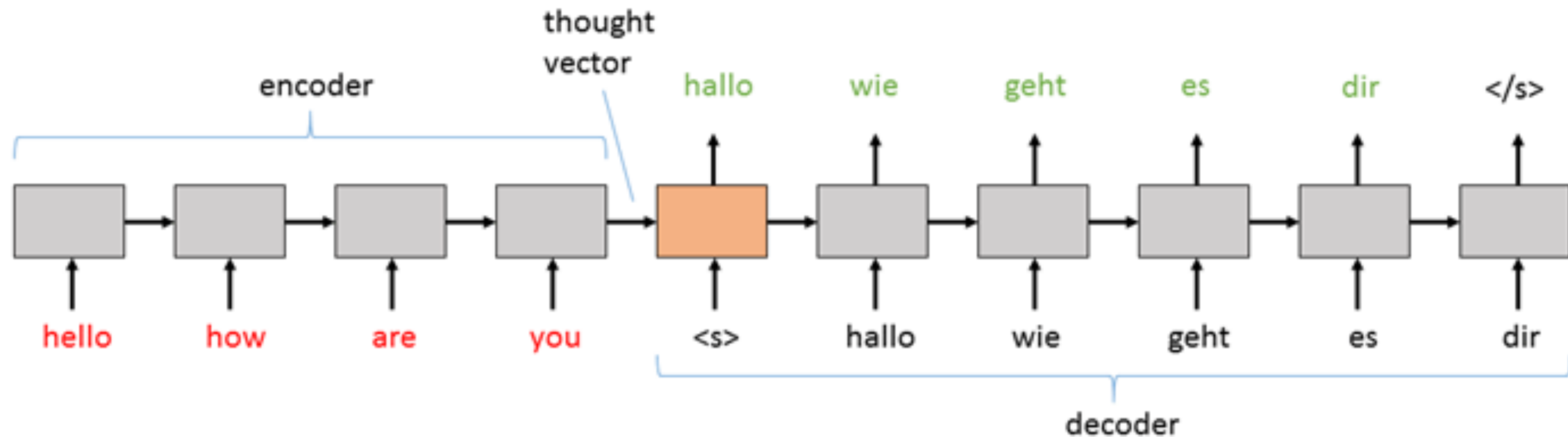
hidden states B) 128

C) ∞

input sequence
(any length) { $x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$...

$N \rightarrow$...

Sequence to Sequence learning (Seq2seq)

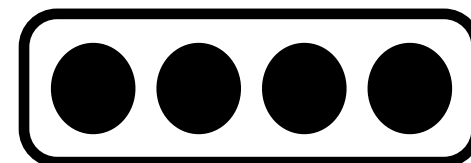


- **Encode** entire input sequence into a single vector (**using an RNN**)
- **Decode** one word at a time (**again, using an RNN!**)
- Beam search for better inference
- Learning is not trivial! (vanishing/exploding gradients)

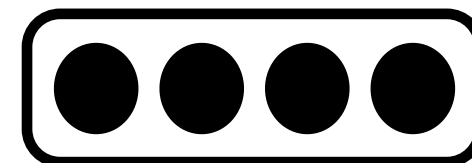
Encoder

Sentence: This cat is cute

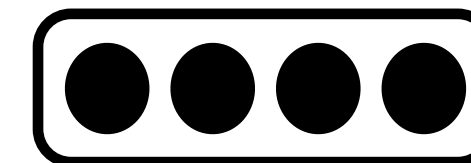
word
embedding



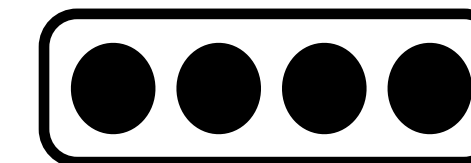
This



cat



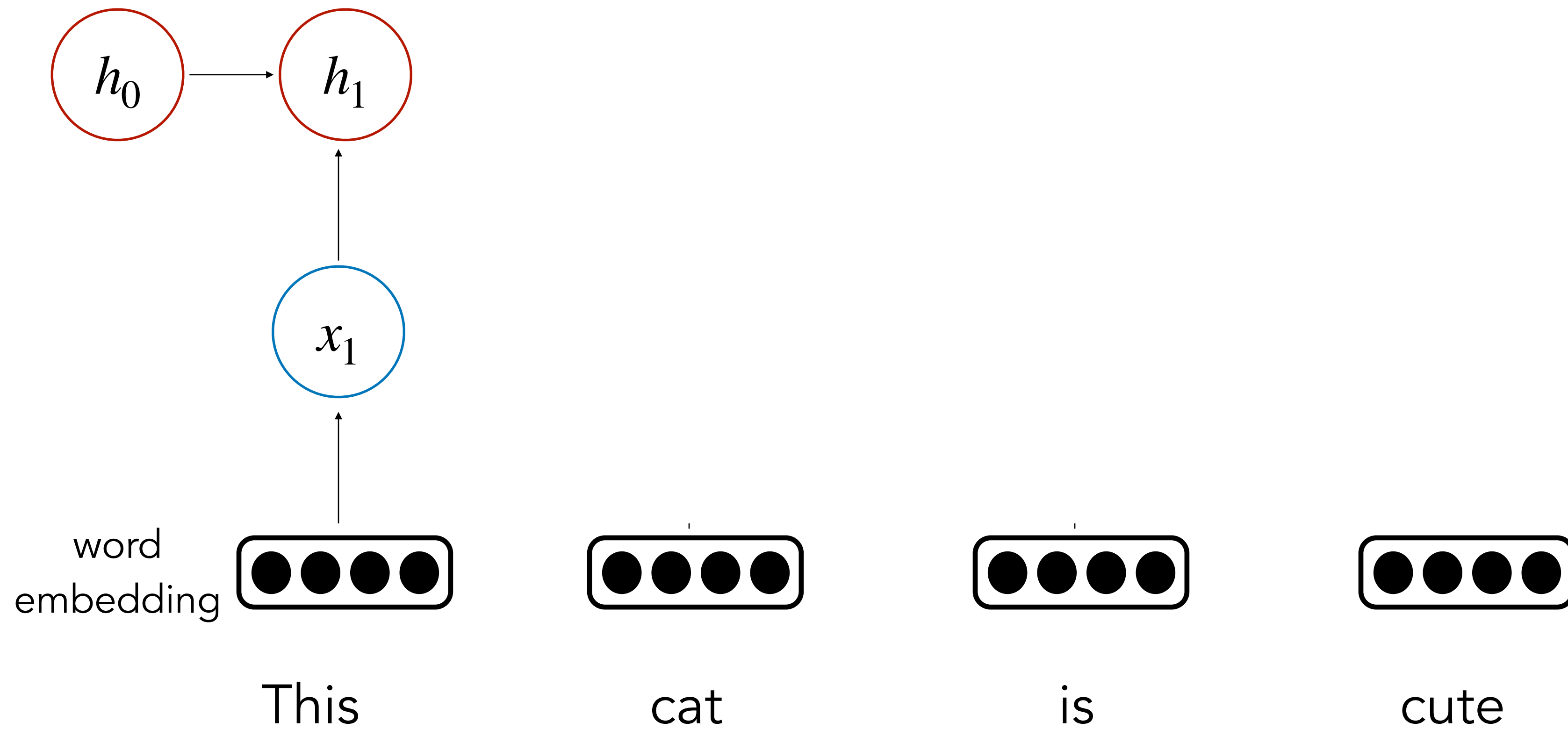
is



cute

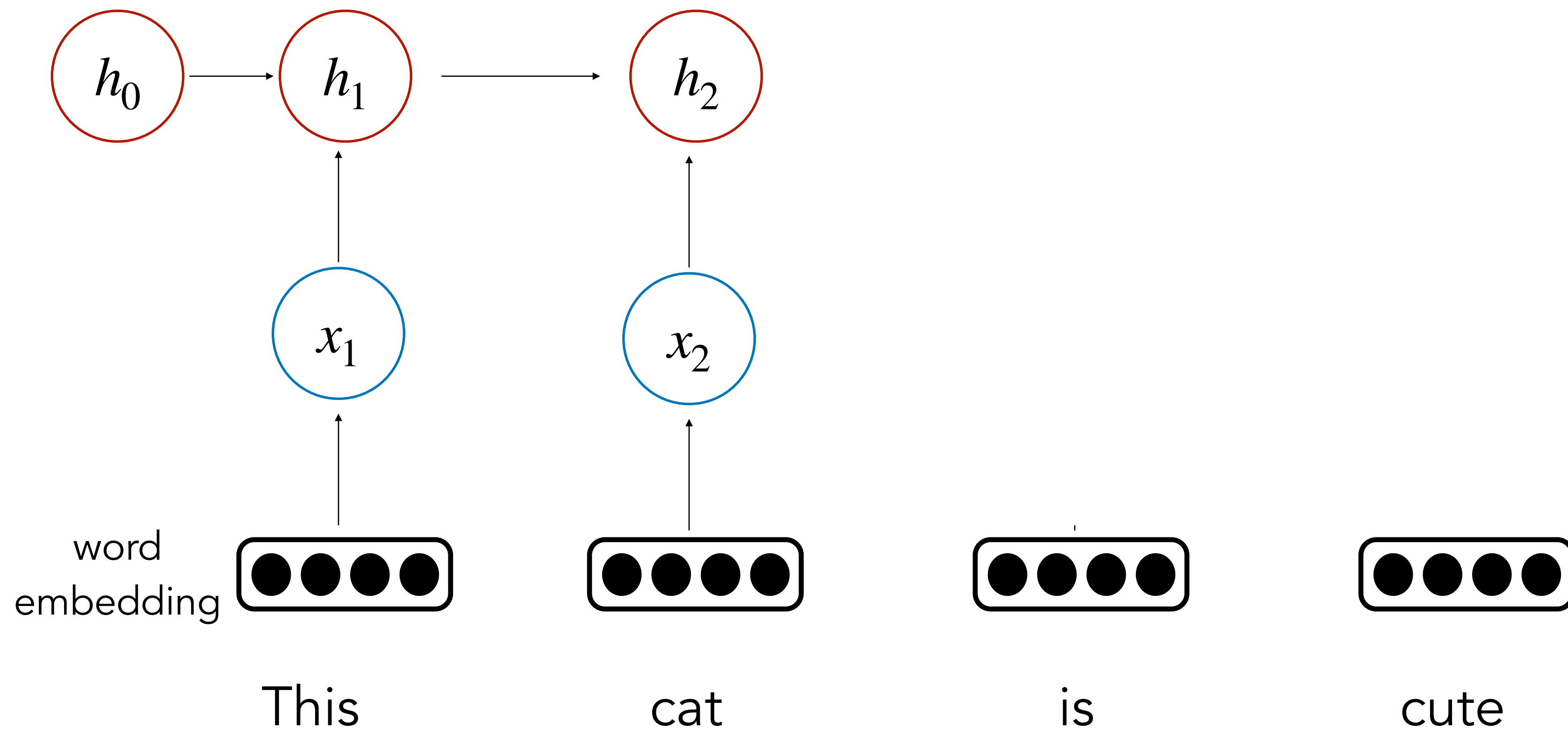
Encoder

Sentence: This cat is cute



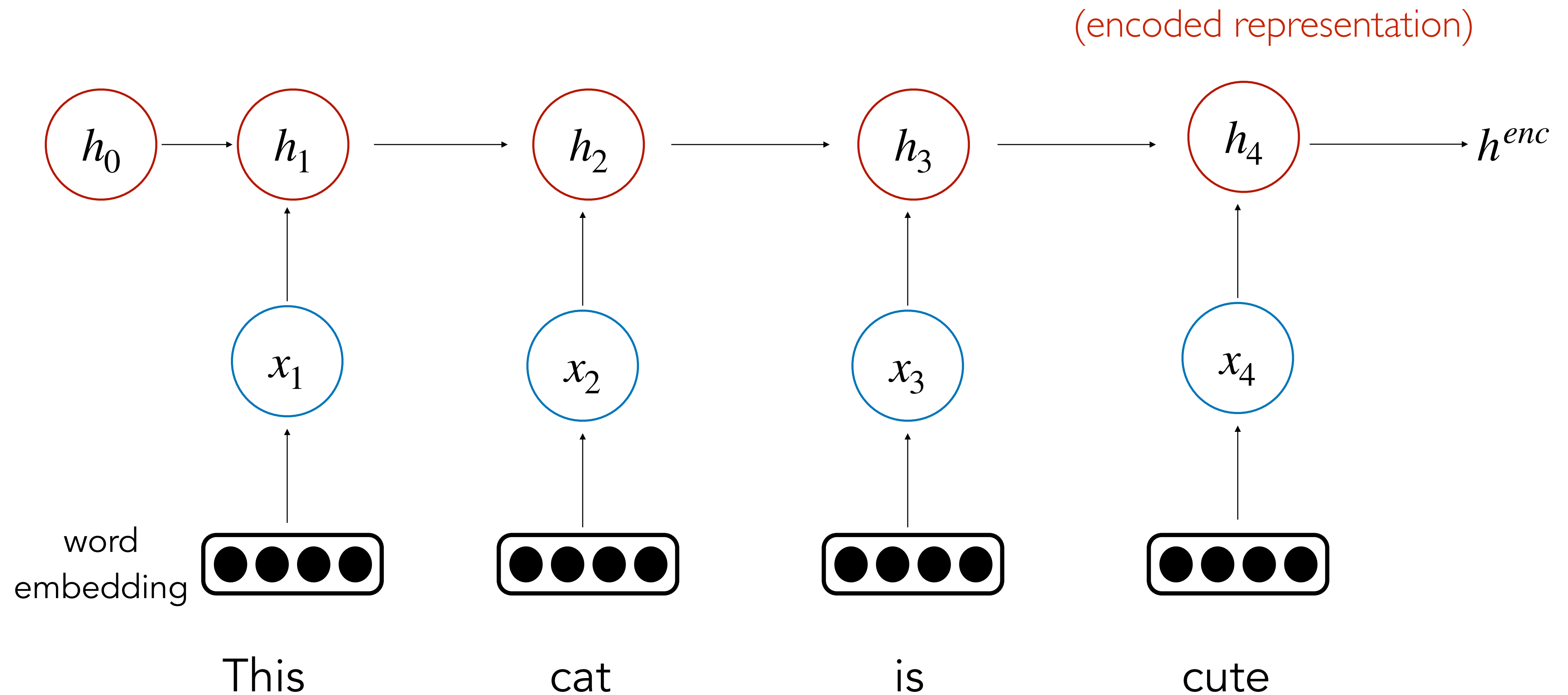
Encoder

Sentence: This cat is cute



Encoder

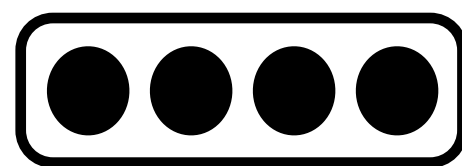
Sentence: This cat is cute



Decoder

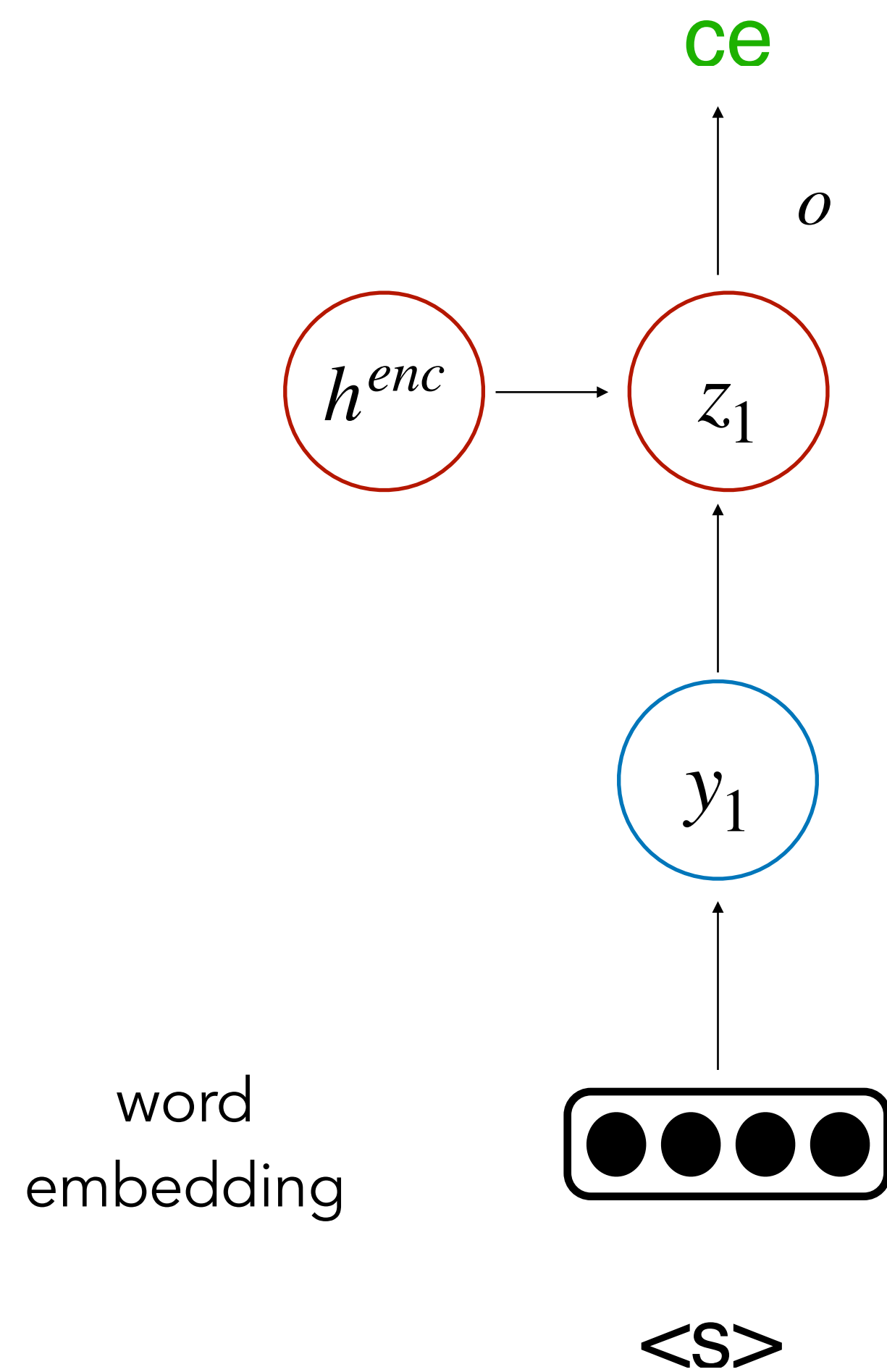
h^{enc}

word
embedding

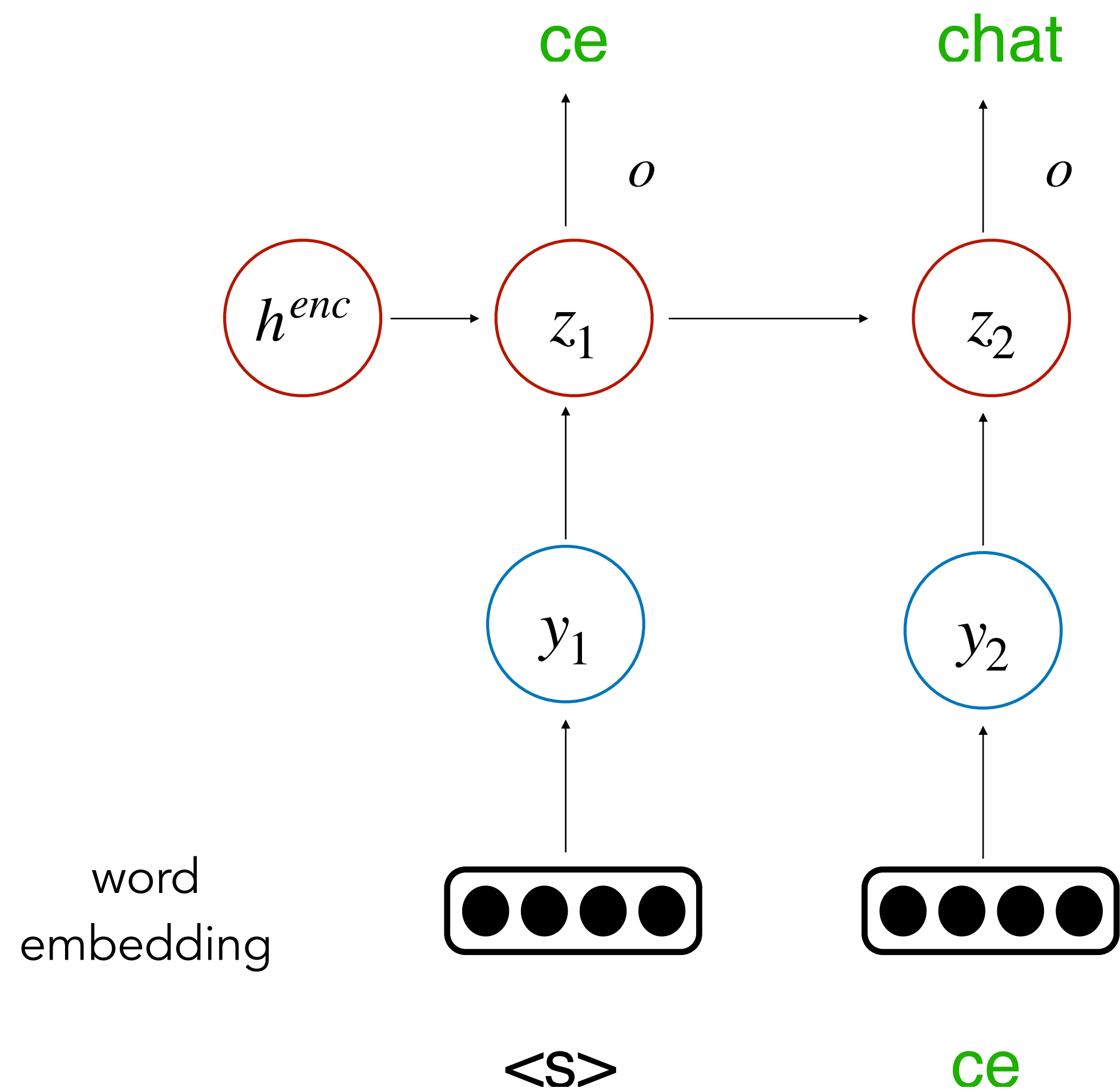


<S>

Decoder

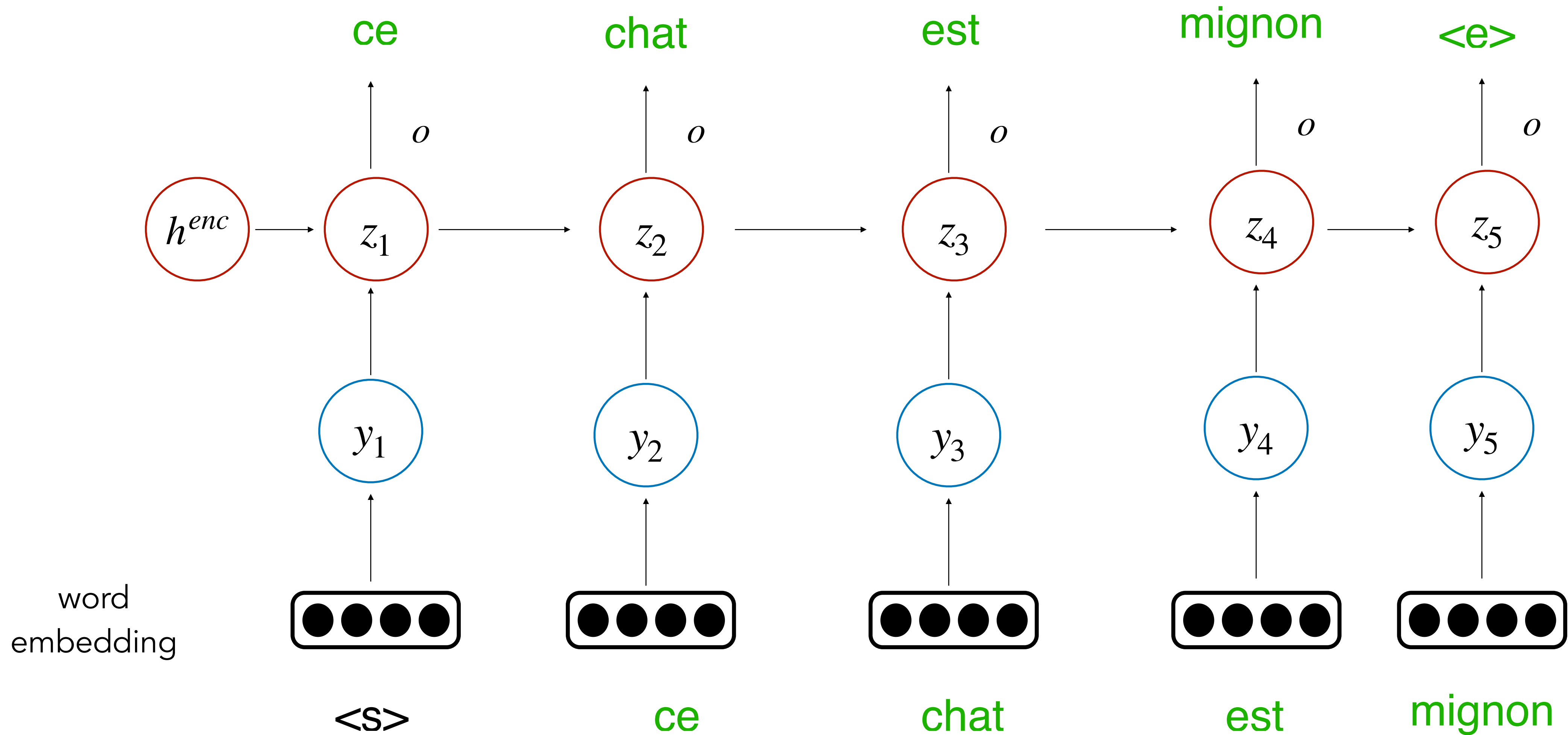


Decoder



Decoder

- A conditioned language model



Seq2seq training

- ▶ Similar to training a language model!

- ▶ Minimize cross-entropy loss:

$$\sum_{t=1}^T -\log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

- ▶ Back-propagate gradients through *both decoder and encoder*
- ▶ Need a really big corpus

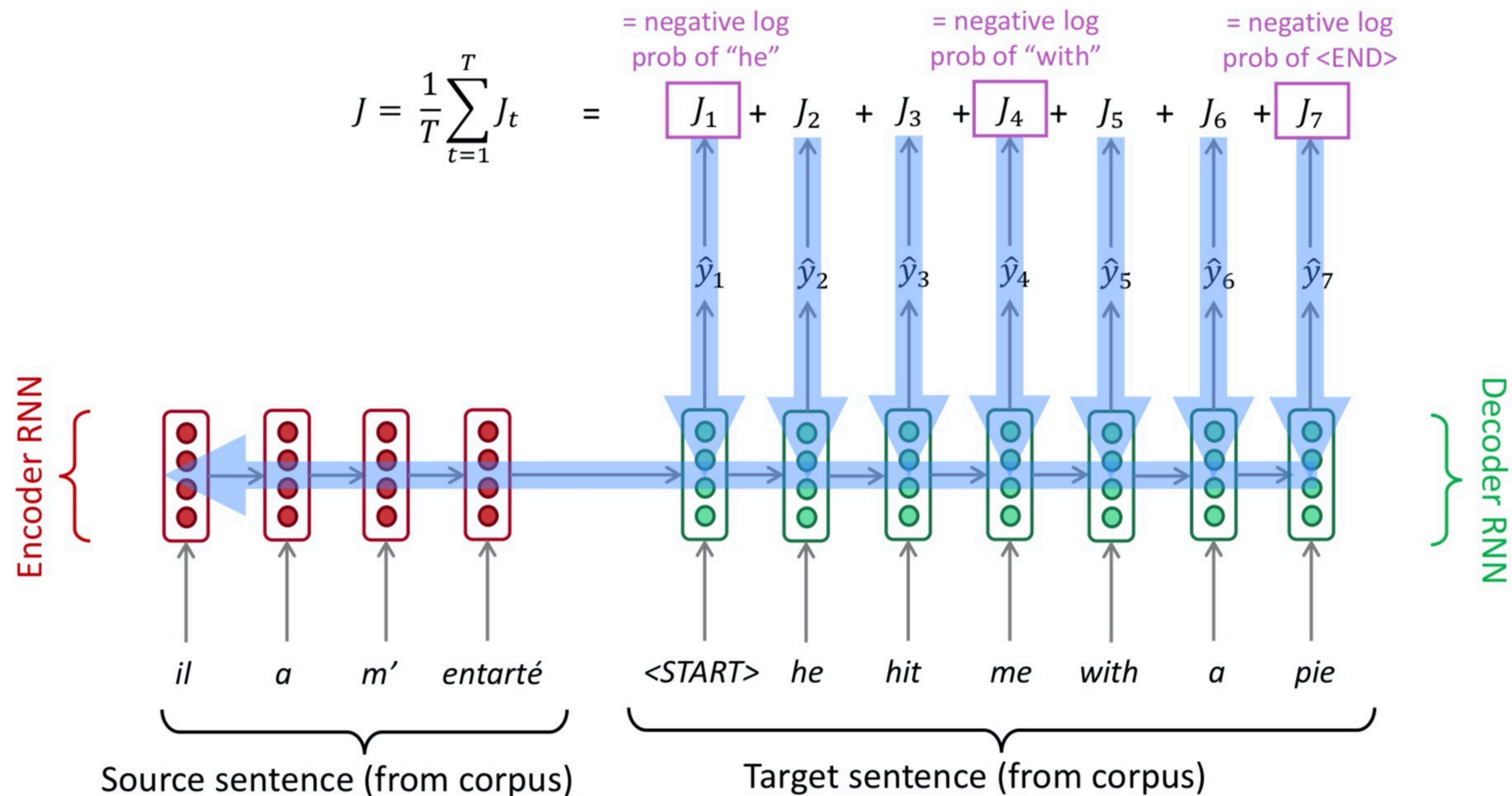
36M sentence pairs

Russian: Машинный перевод - это круто!



English: Machine translation is cool!

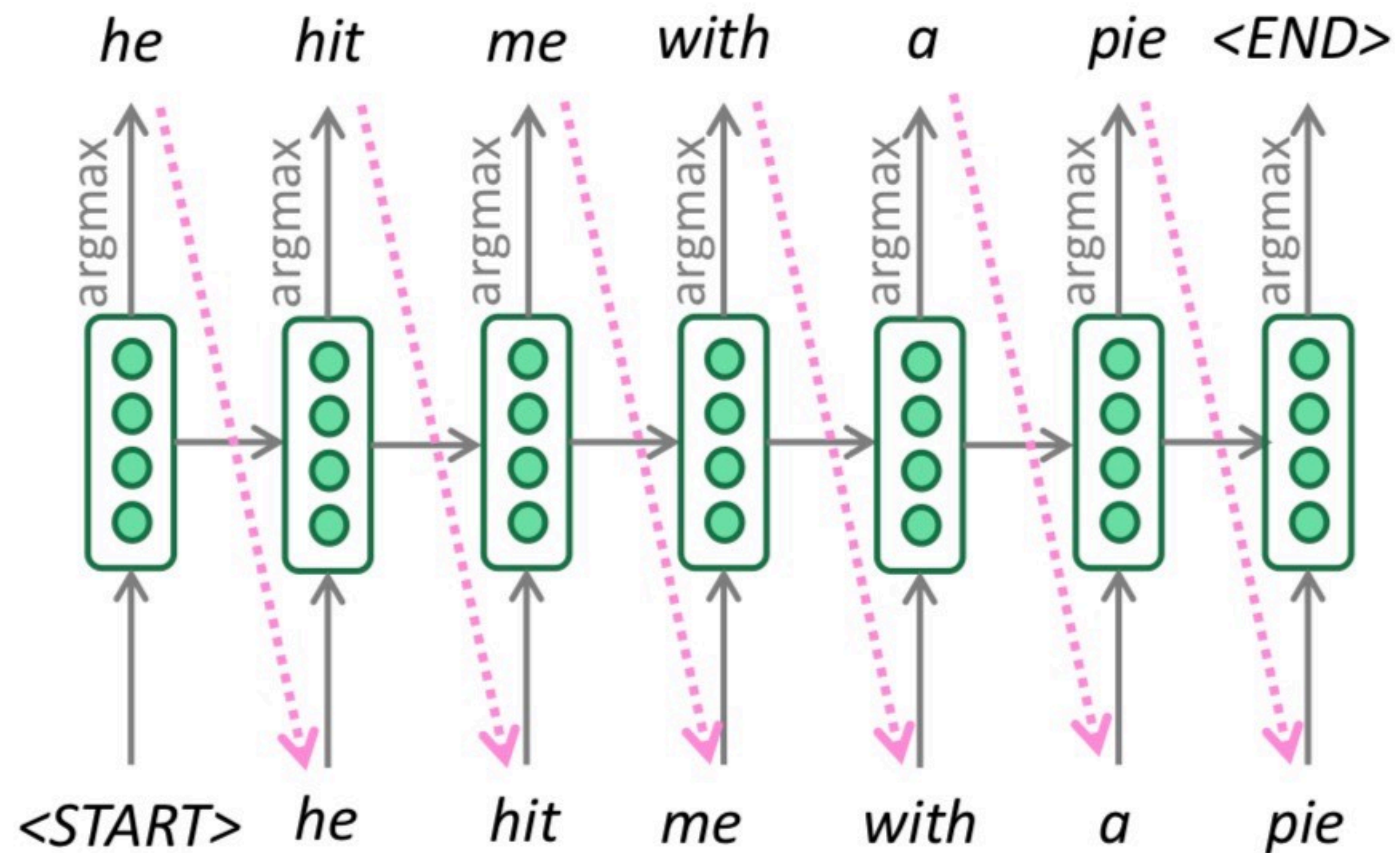
Seq2seq training



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

(slide credit: Abigail See)

Greedy decoding



- ▶ Compute argmax at every step of decoder to generate word
- ▶ What's wrong?

Exhaustive search?



- ▶ Find $\arg \max_{y_1, \dots, y_T} P(y_1, \dots, y_T | x_1, \dots, x_n)$
- ▶ Requires computing all possible sequences

V - Vocabulary
T - length of sequence

What is the complexity of doing this search?

A) $O(VT)$

B) $O(V^T)$

C) $O(T^V)$

A middle ground: Beam search

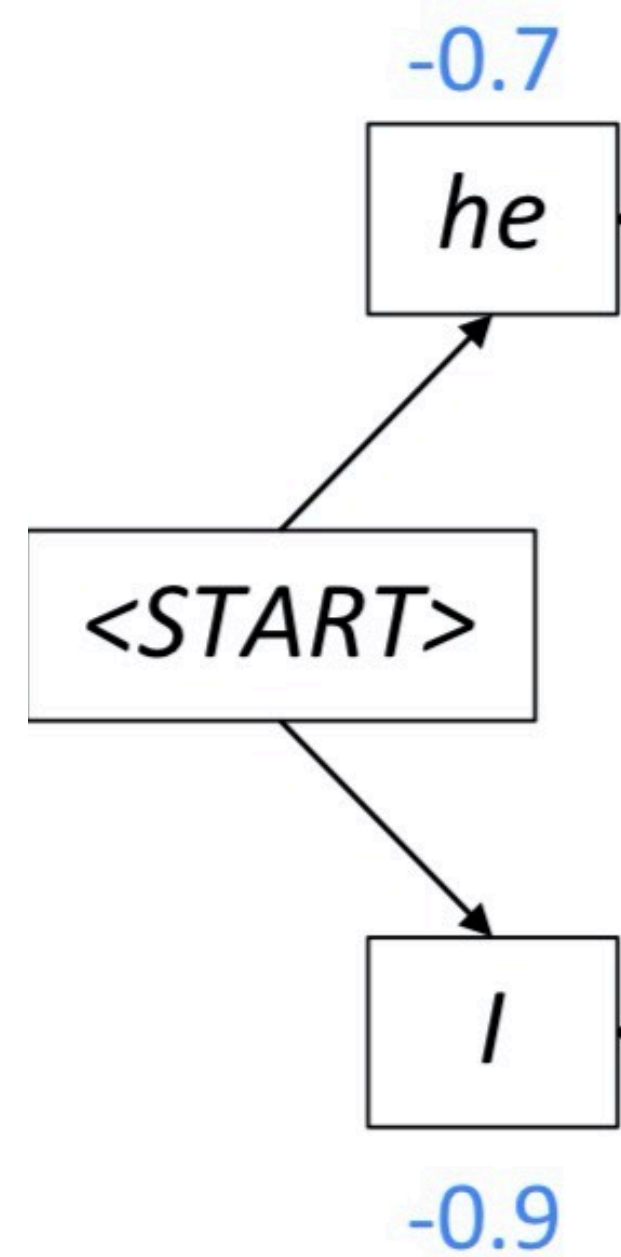
- ▶ **Key idea:** At every step, keep track of the **k most probable** partial translations (hypotheses)
- ▶ Score of each hypothesis = log probability of sequence so far

$$\sum_{t=1}^j \log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

- ▶ Not guaranteed to be optimal
- ▶ More efficient than exhaustive search

Beam decoding

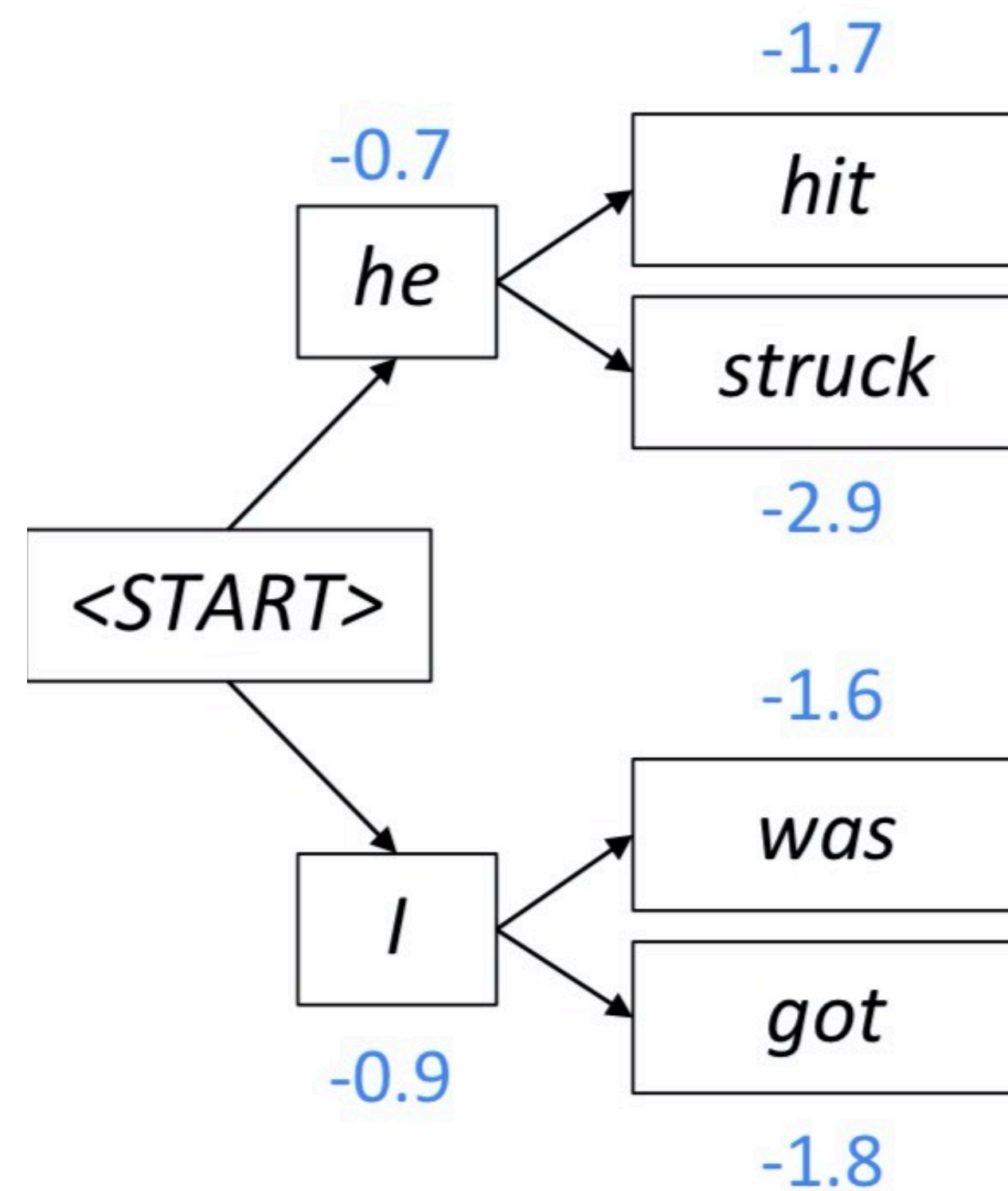
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



(slide credit: Abigail See)

Beam decoding

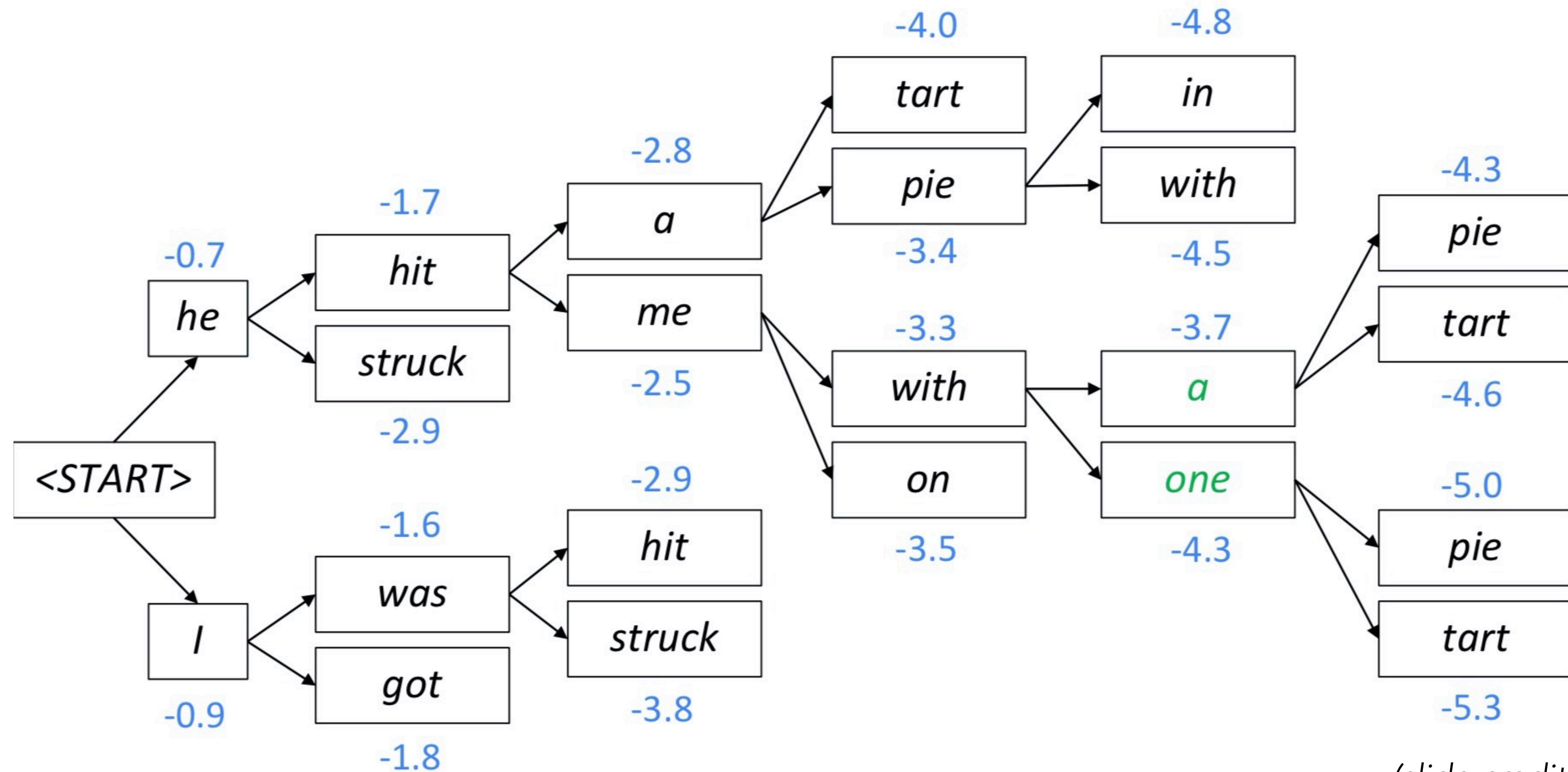
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



(slide credit: Abigail See)

Beam decoding

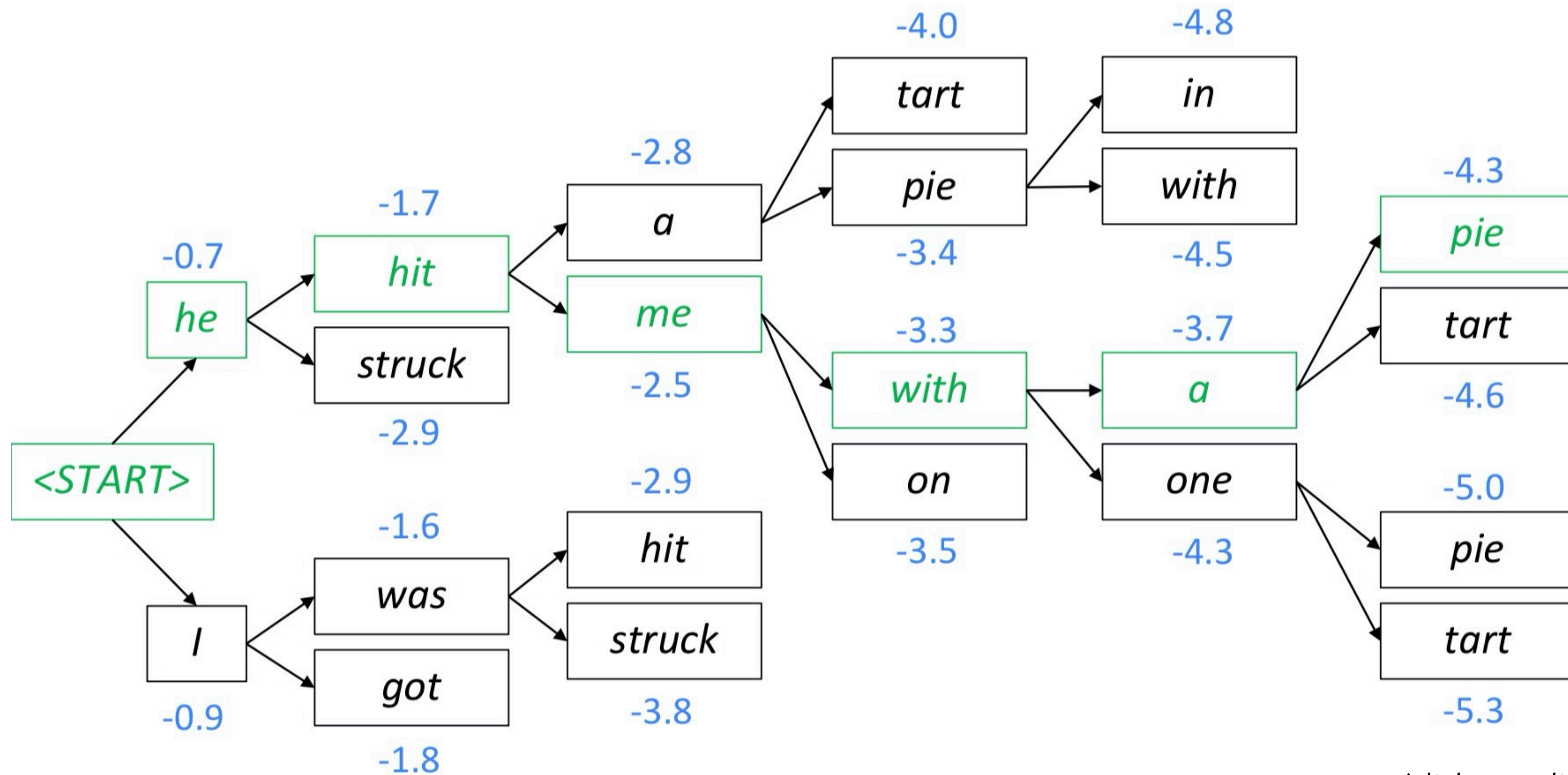
Beam size = $k = 2$. Blue numbers = $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



(slide credit: Abigail See)

Backtrack

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



(slide credit: Abigail See)

Beam decoding

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

$$\frac{1}{T} \sum_{t=1}^T \log P(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_n)$$

- ▶ Otherwise shorter hypotheses have higher scores

NMT vs SMT

Pros



Cons



