



COS 484/584

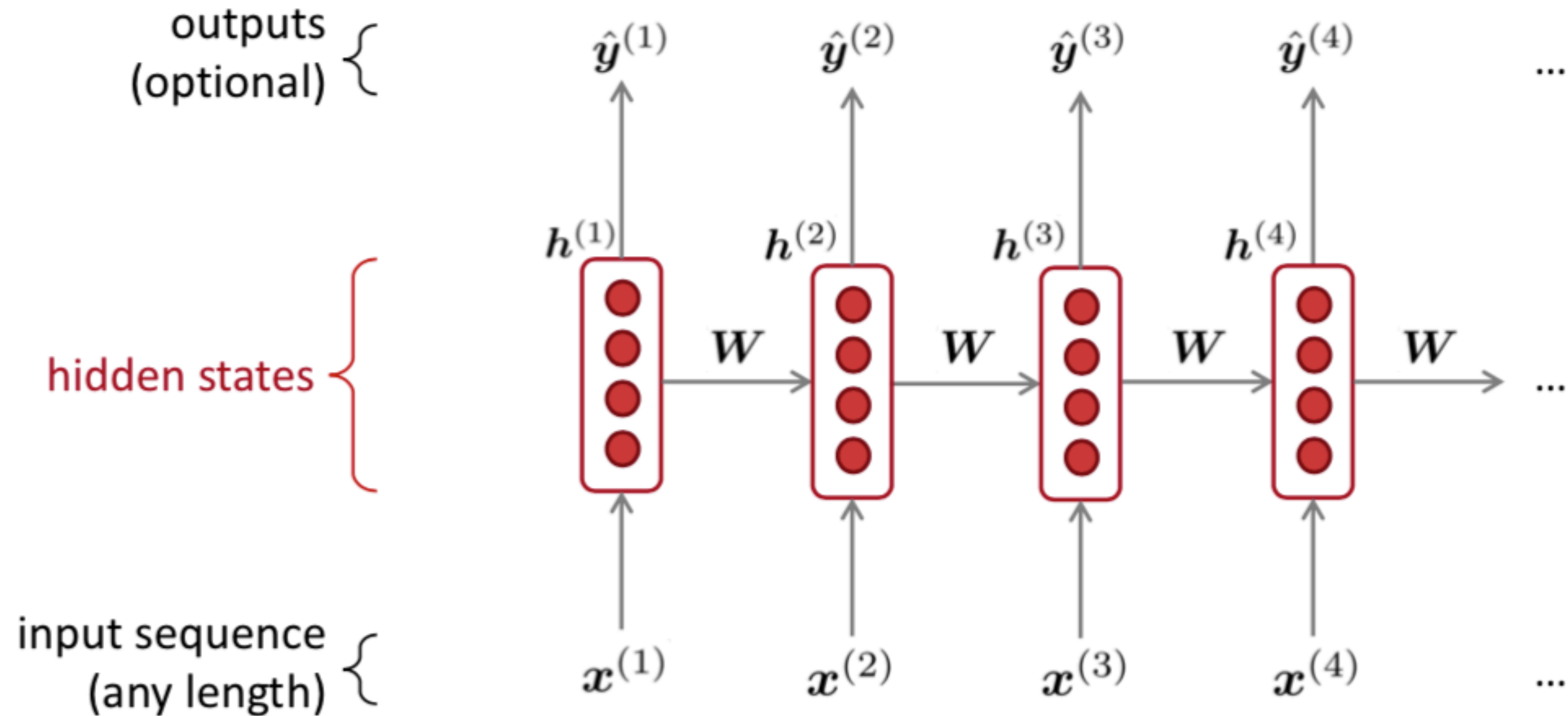
# LI 6: Neural Machine Translation - I

# 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

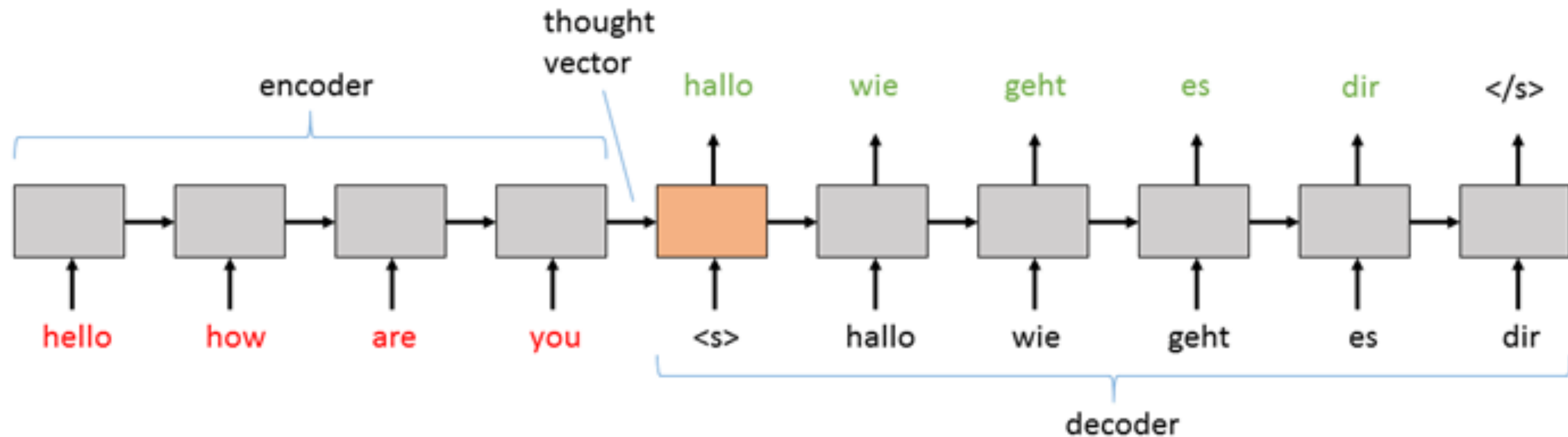
hidden states B) 128

C)  $\infty$

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

$N \rightarrow$  ...

# Sequence to Sequence learning (Seq2seq)

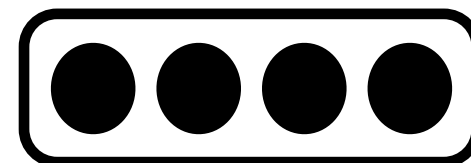


- **Encode** entire input sequence into a single vector (**using an RNN**)
- **Decode** one word at a time (**again, using an RNN!**)
- Beam search for better inference
- Learning is not trivial! (vanishing/exploding gradients)

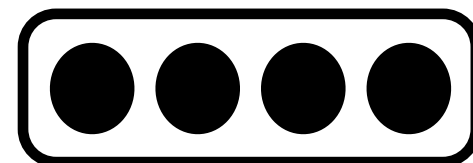
# Encoder

*Sentence: This cat is cute*

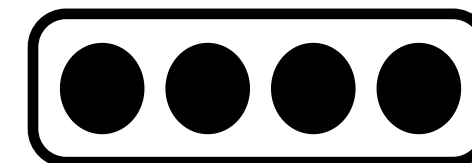
word  
embedding



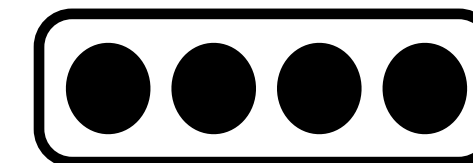
This



cat



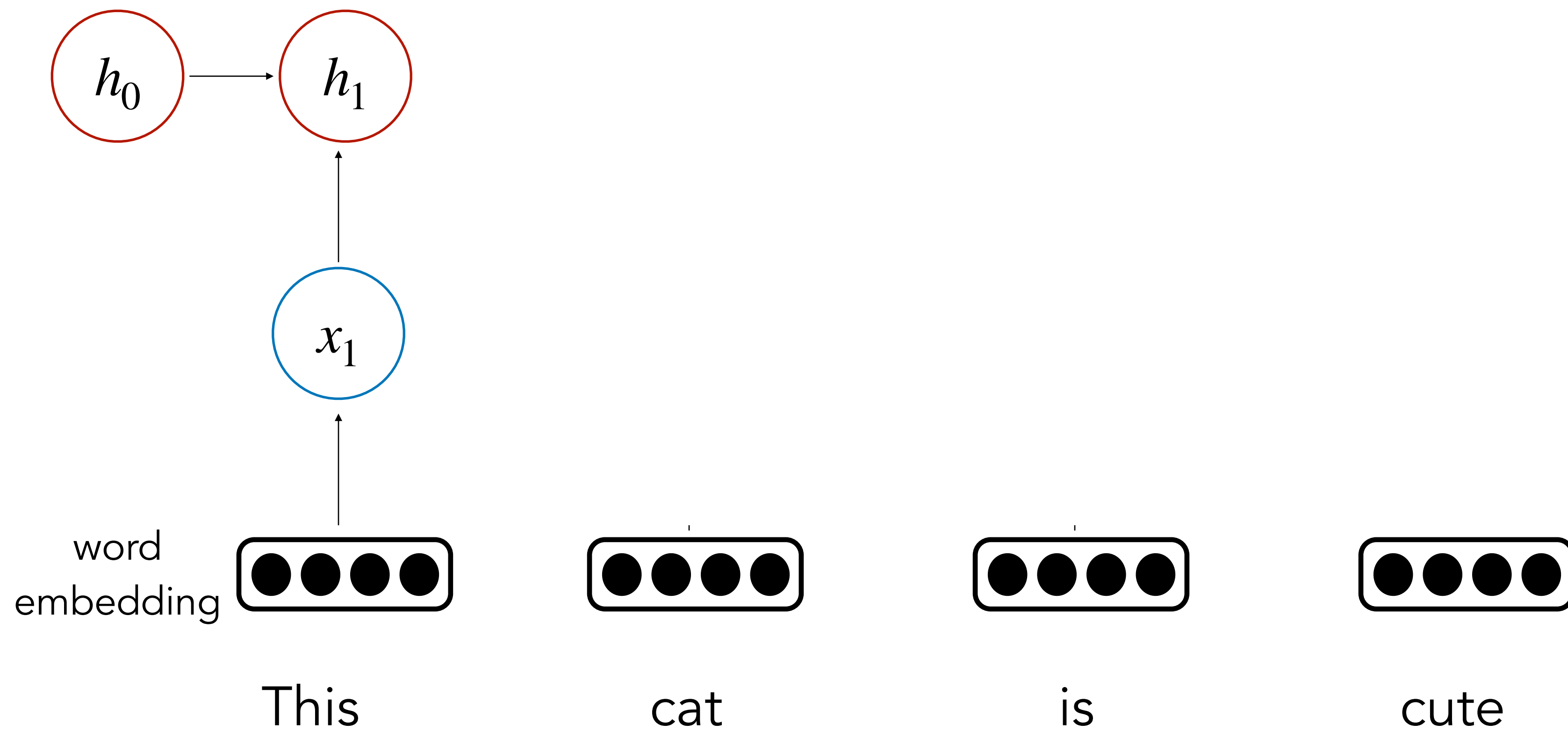
is



cute

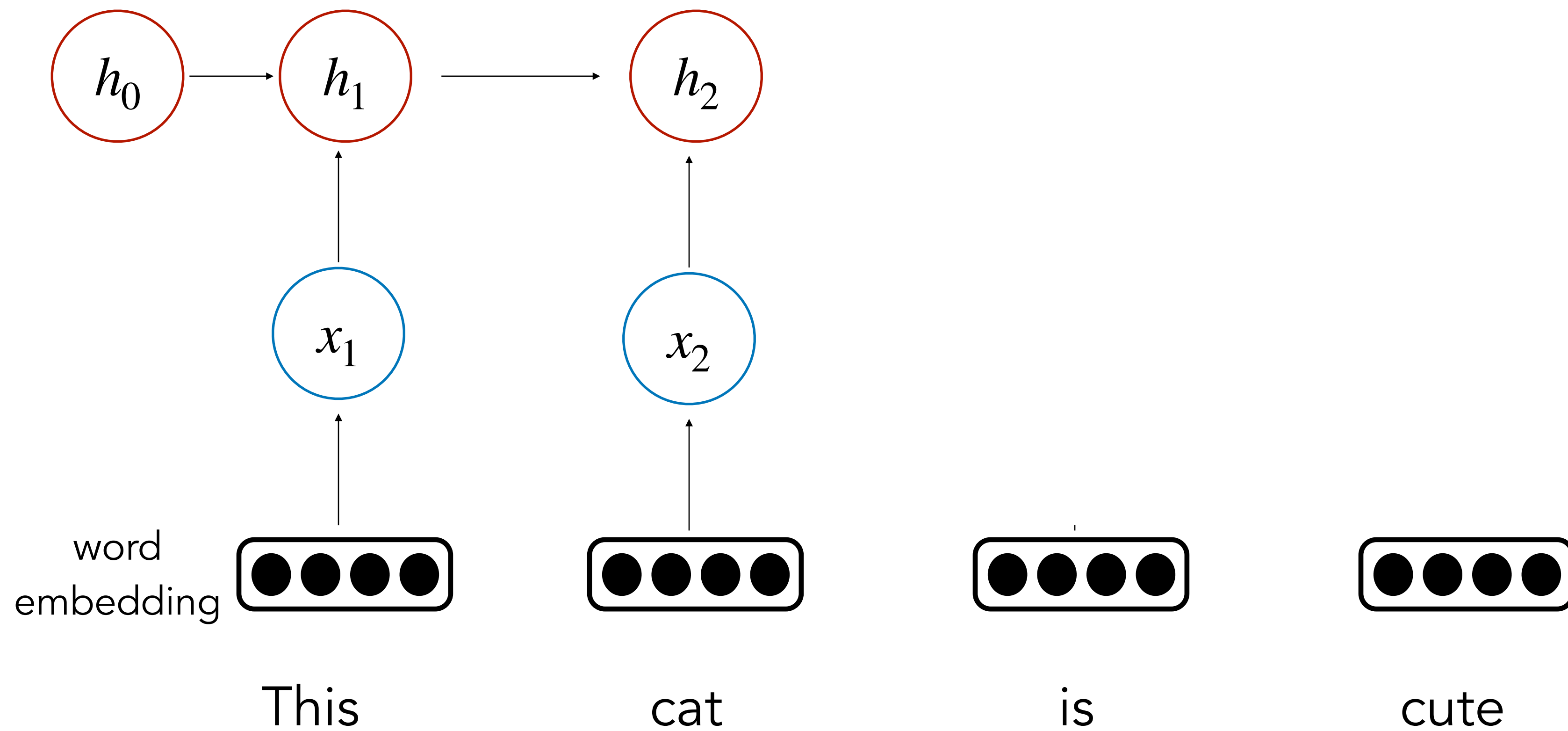
# Encoder

*Sentence: This cat is cute*



# Encoder

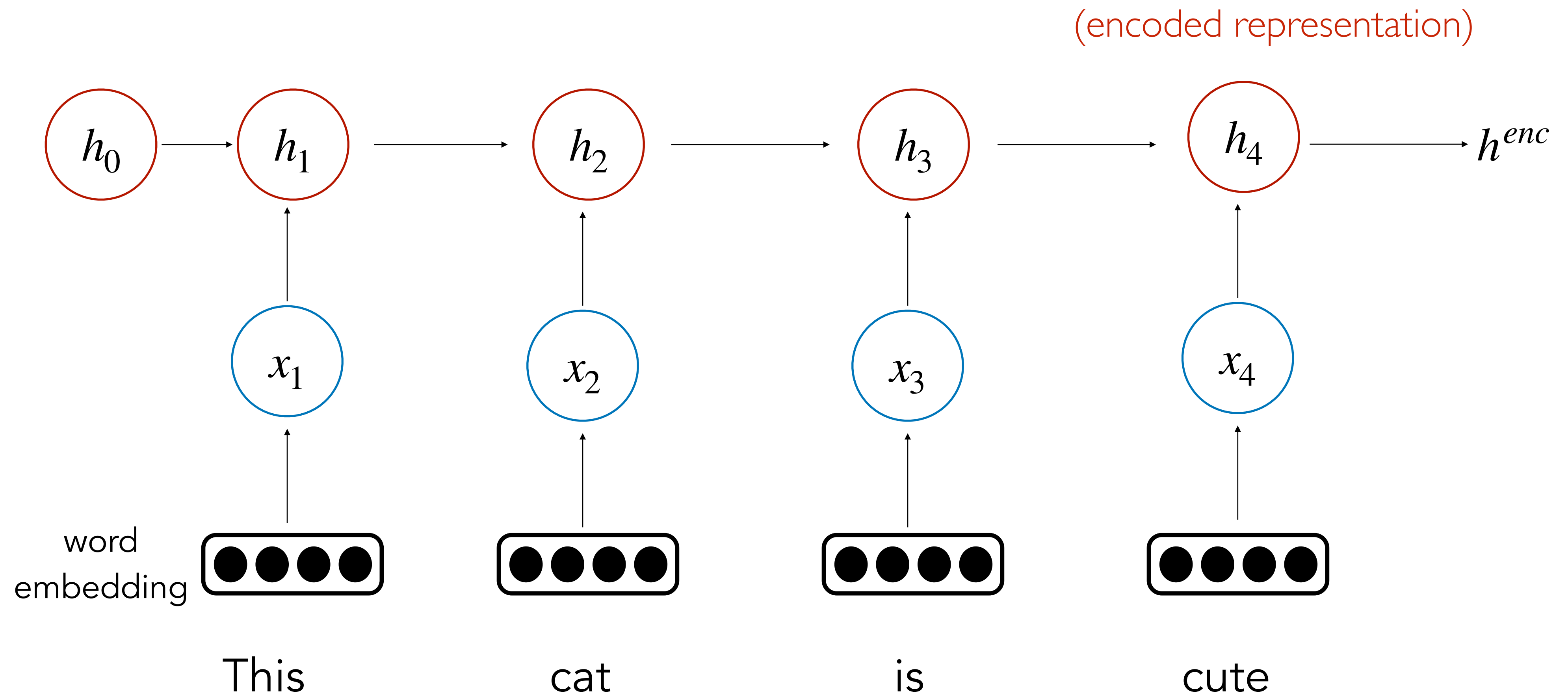
*Sentence: This cat is cute*





# Encoder

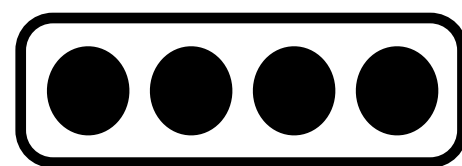
*Sentence: This cat is cute*



# Decoder

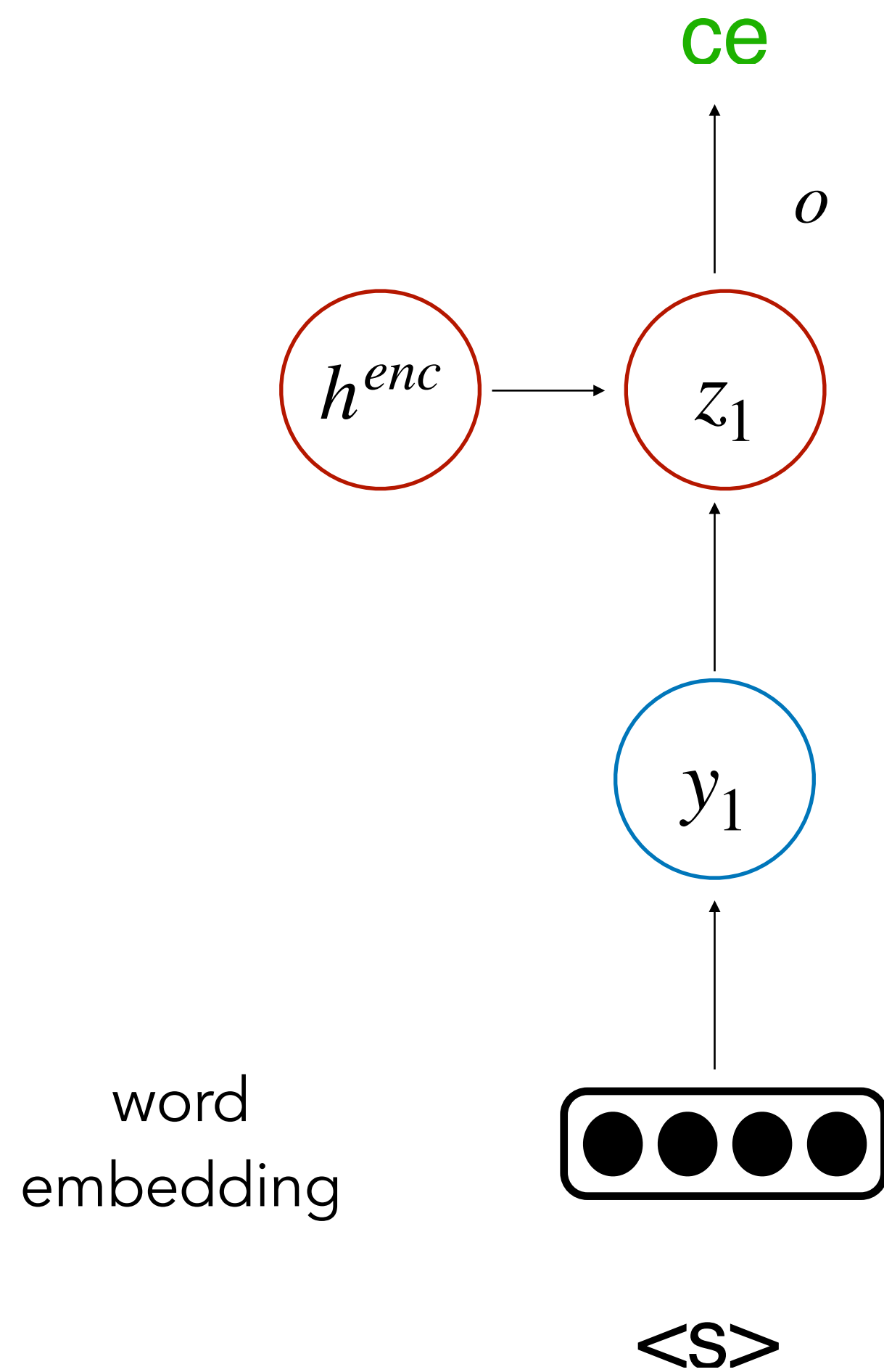
$h^{enc}$

word  
embedding

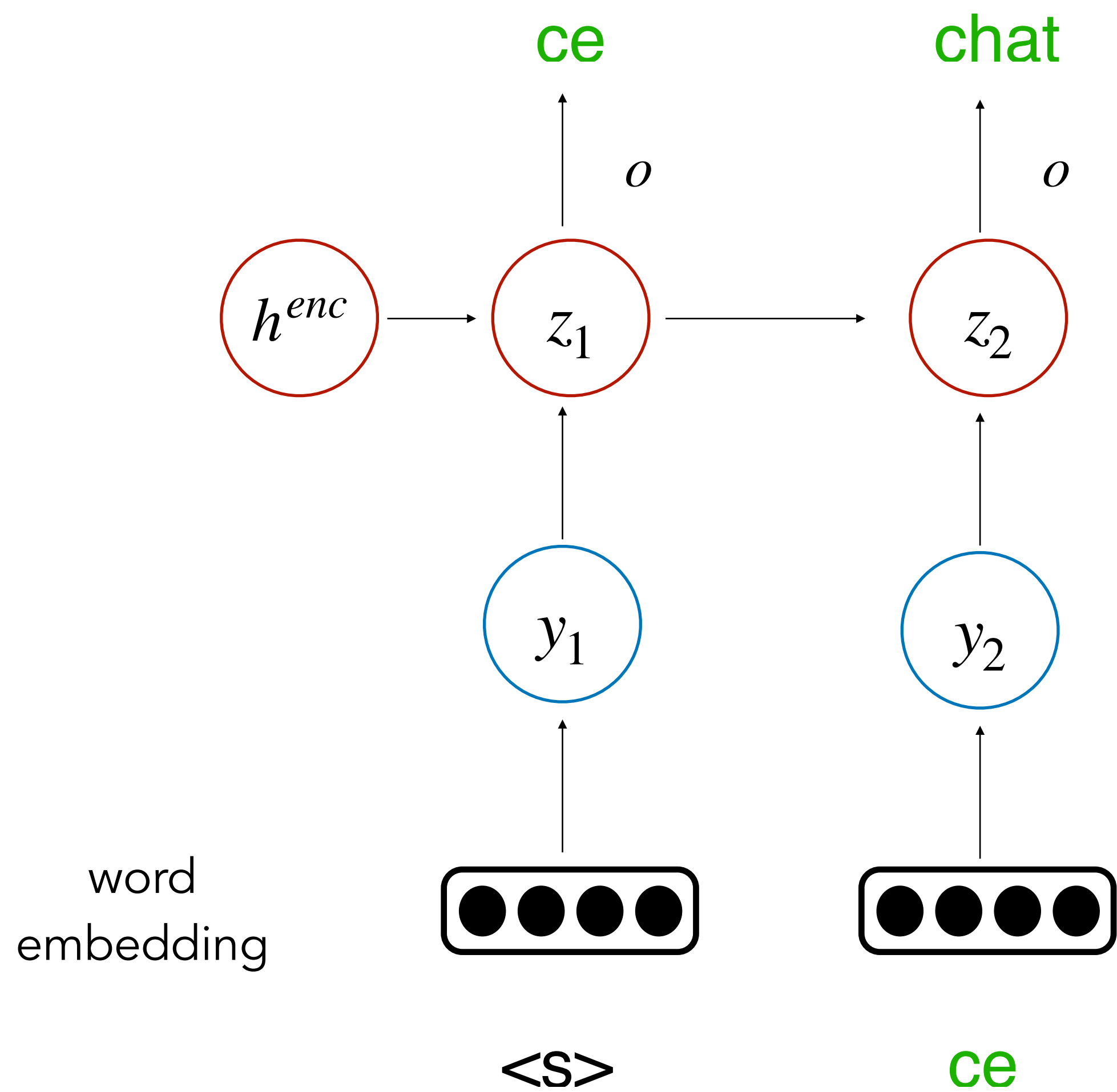


<S>

# Decoder

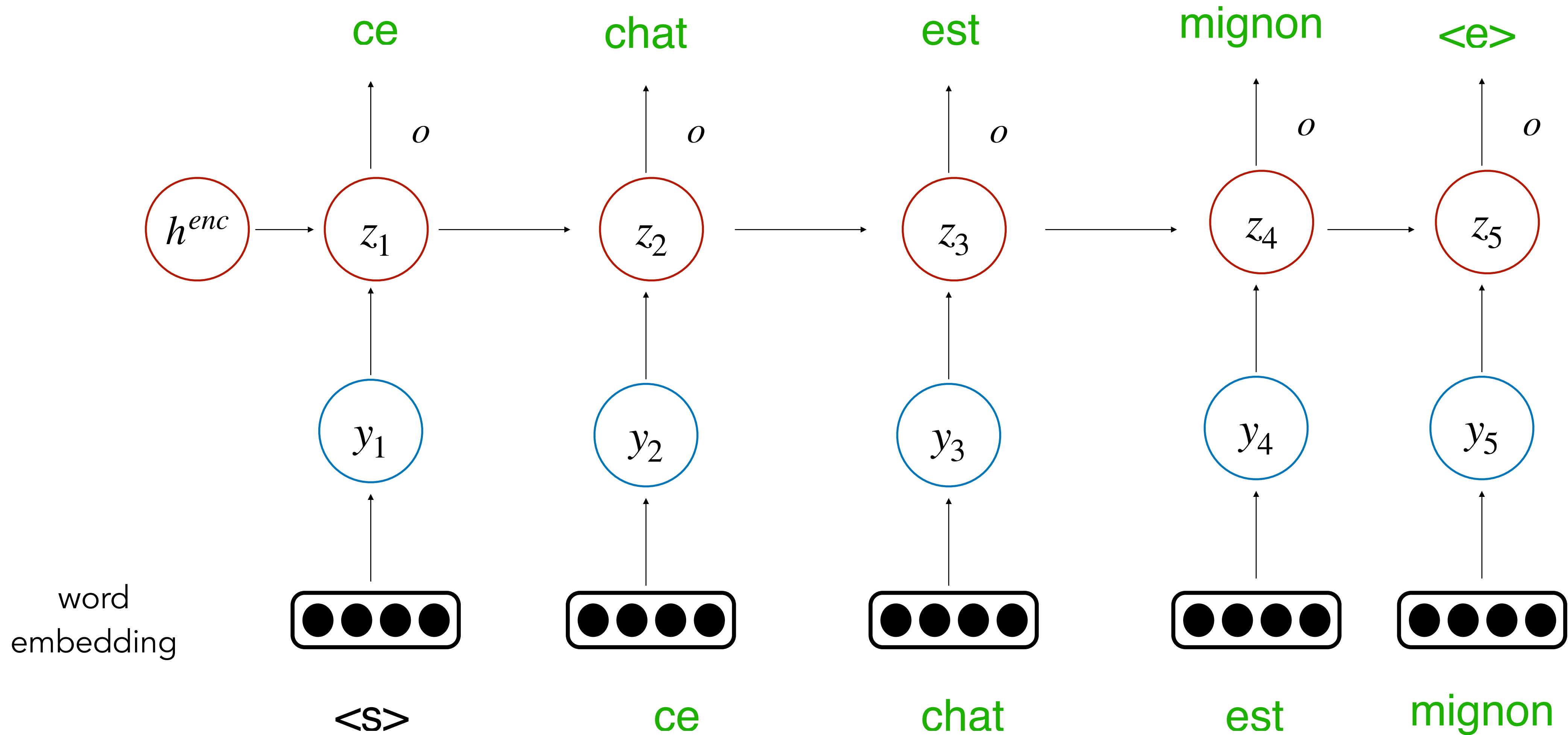


# Decoder



# 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

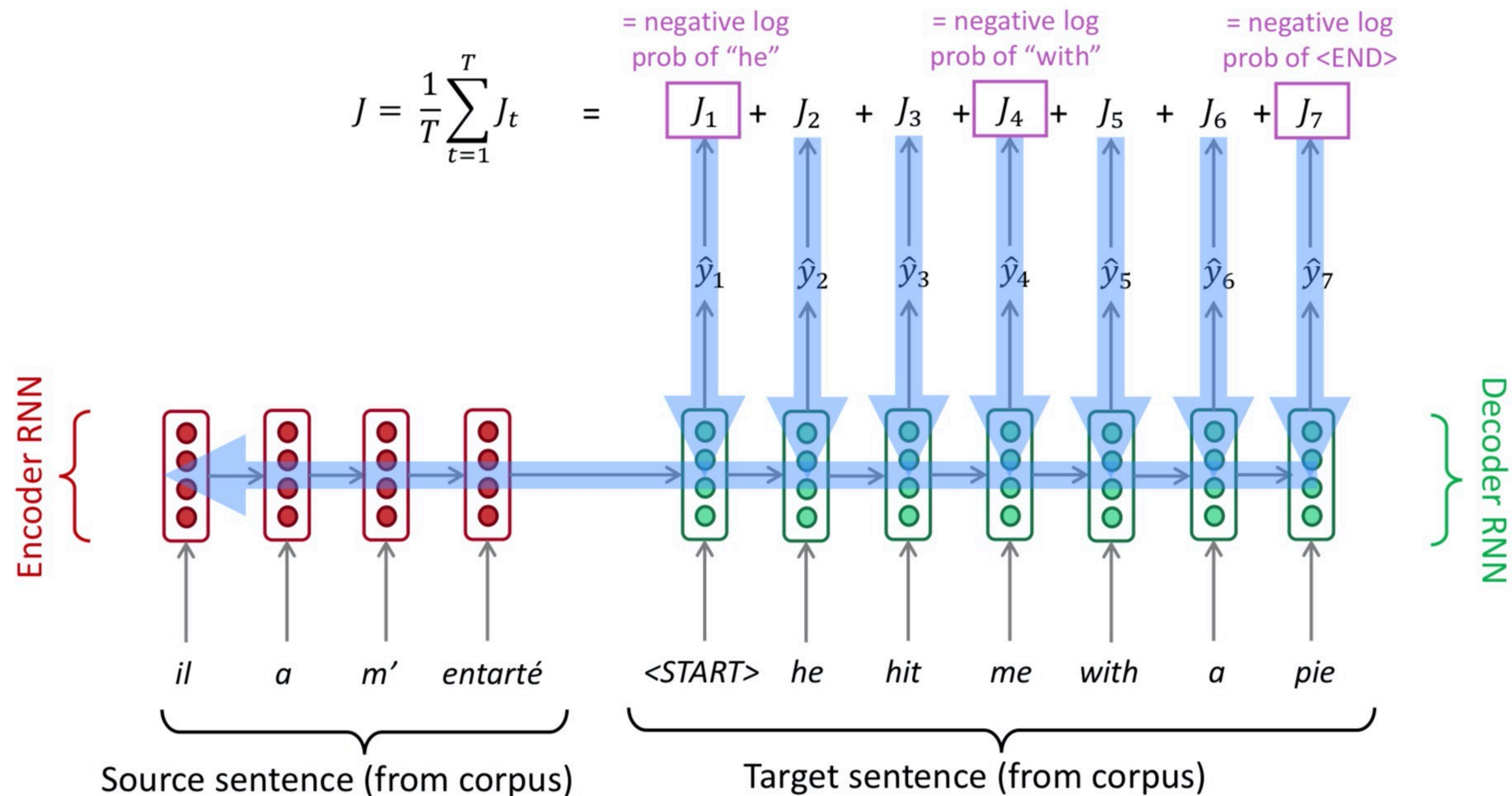
36M sentence pairs

*Russian:* Машинный перевод - это круто!



*English:* Machine translation is cool!

# Seq2seq training

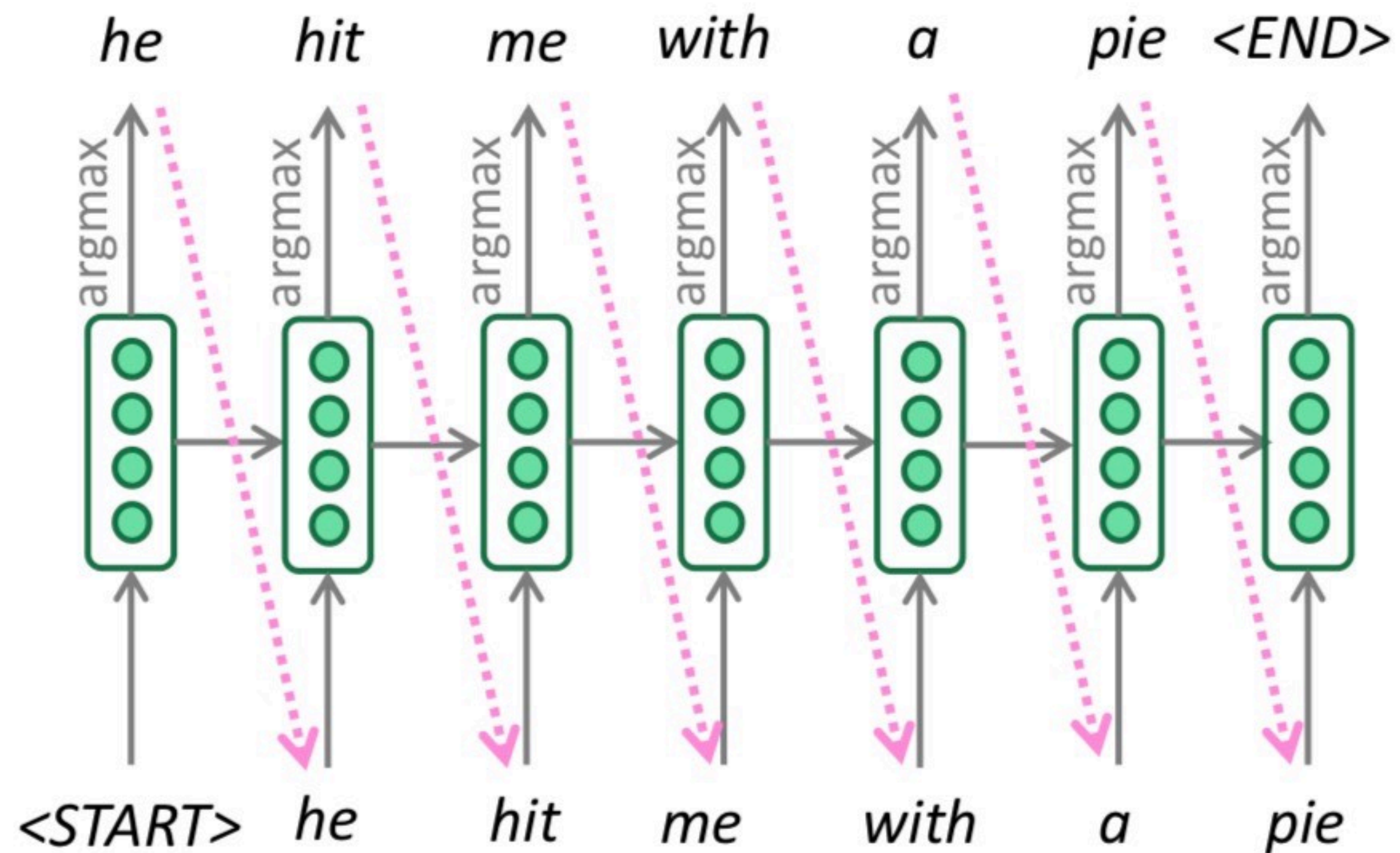


Seq2seq is optimized as a single system.  
Backpropagation operates "*end-to-end*".

(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

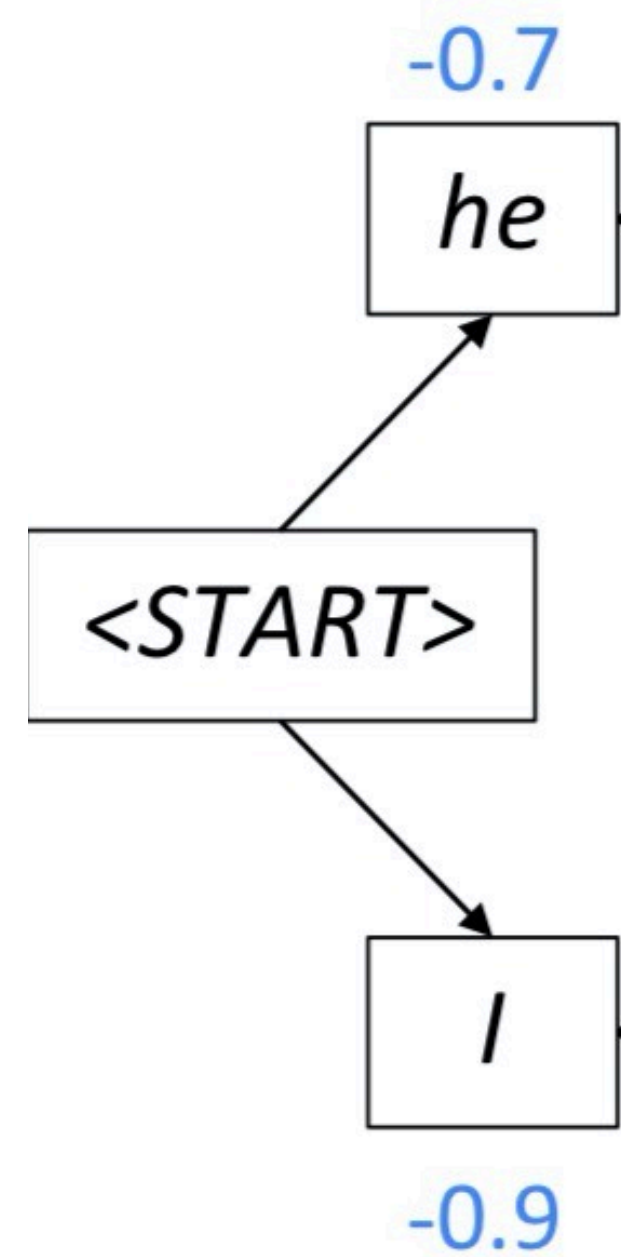
- ▶ **Key idea:** At every step, keep track of the **k most probable** partial translations (hypotheses)
- ▶ Score of each hypothesis = log probability of sequence so far

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

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

# Beam decoding

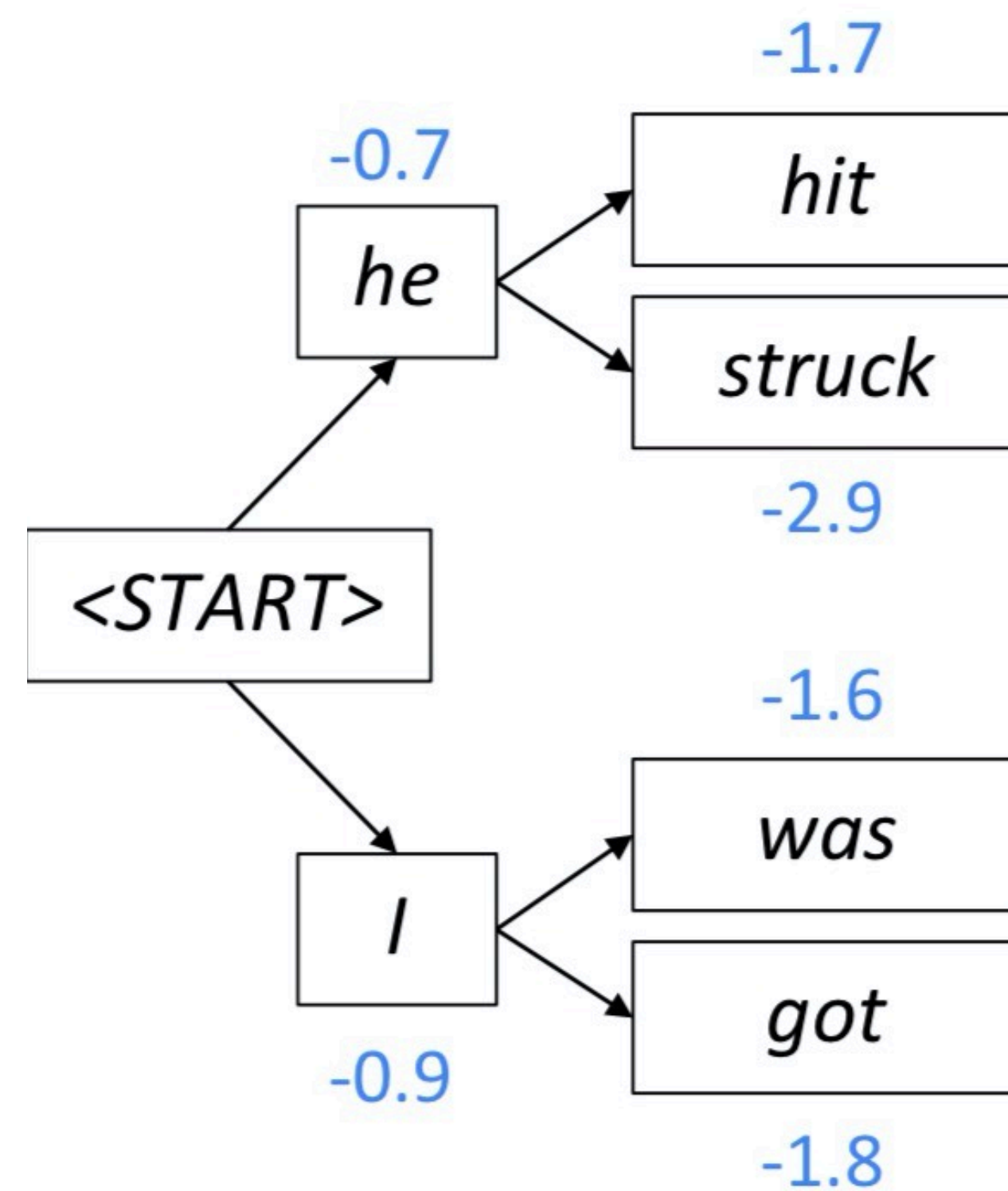
Beam size =  $k = 2$ . Blue numbers =  $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



(slide credit: Abigail See)

# Beam decoding

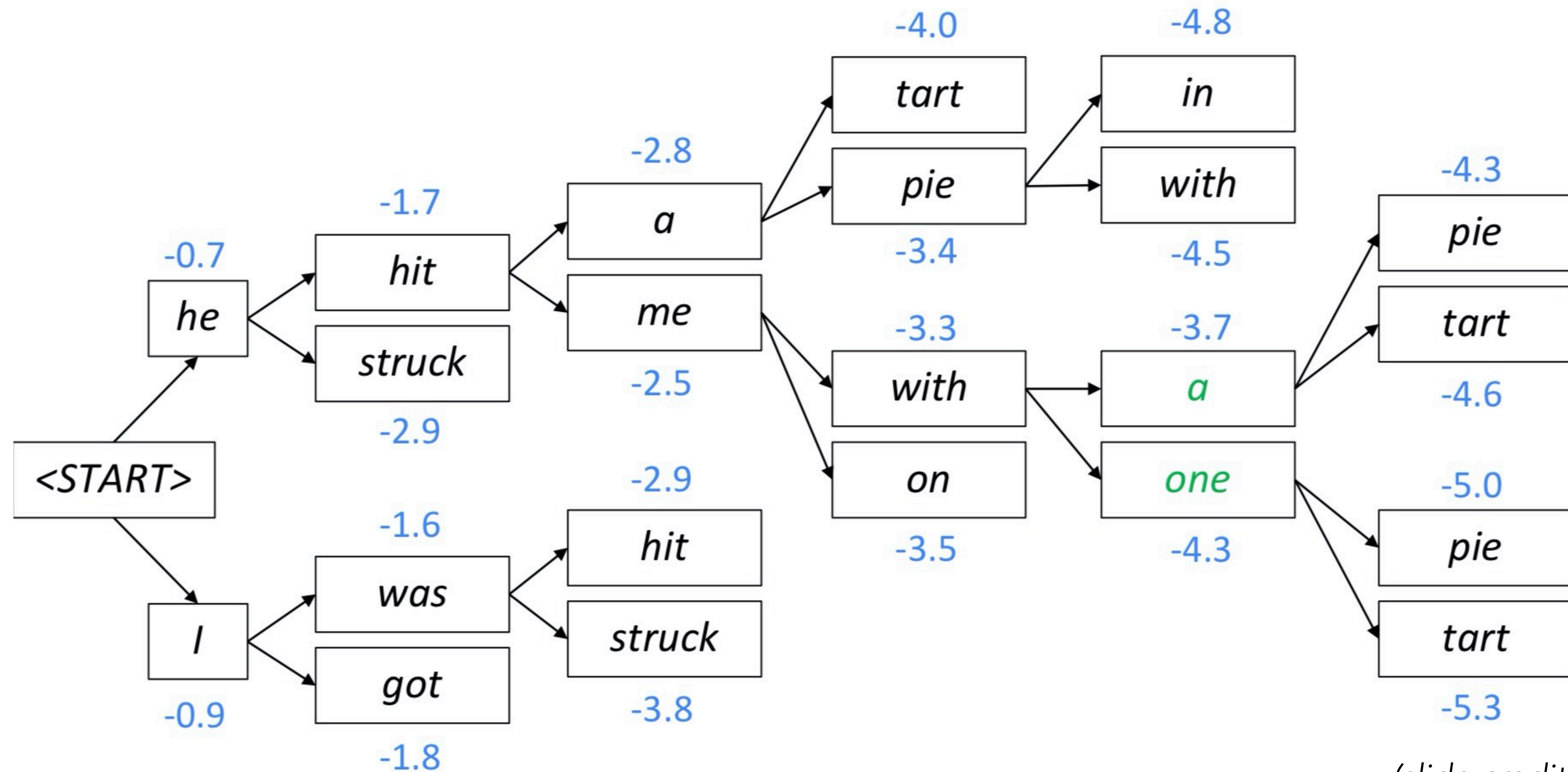
Beam size =  $k = 2$ . Blue numbers =  $\text{score}(y_1, \dots, y_t) = \sum_{i=1}^t \log P_{\text{LM}}(y_i | y_1, \dots, y_{i-1}, x)$



(slide credit: Abigail See)

# Beam decoding

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

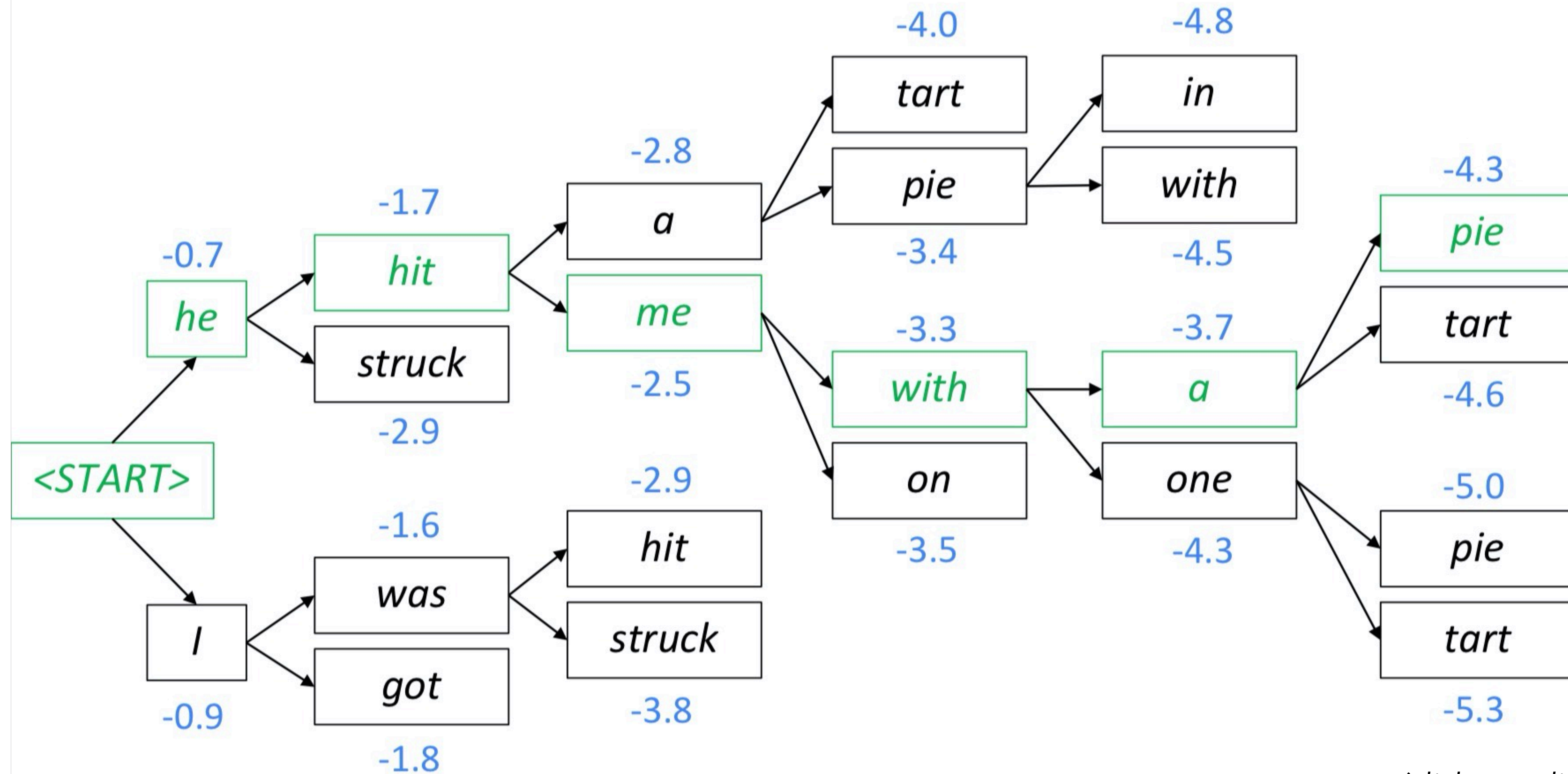


(slide credit: Abigail See)



# 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