



**THE OHIO STATE
UNIVERSITY**

CSE 5525: Foundations of Speech and Language Processing

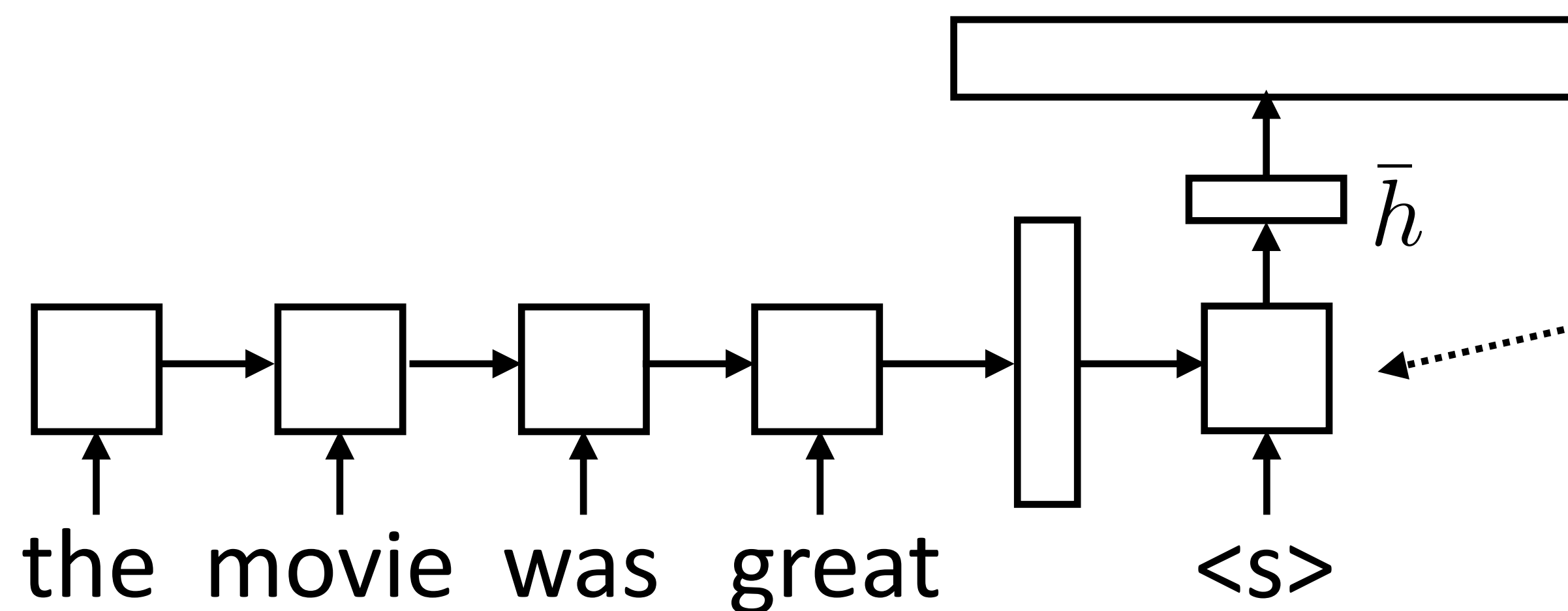
Sequence to sequence (seq2seq) II

Huan Sun (CSE@OSU)

Many thanks to Prof. Greg Durrett @ UT Austin for sharing his slides.

Recall: Seq2seq Model

- ▶ Generate next word conditioned on previous word as well as hidden state
- ▶ W size is $|\text{vocab}| \times |\text{hidden state}|$, softmax over entire vocabulary

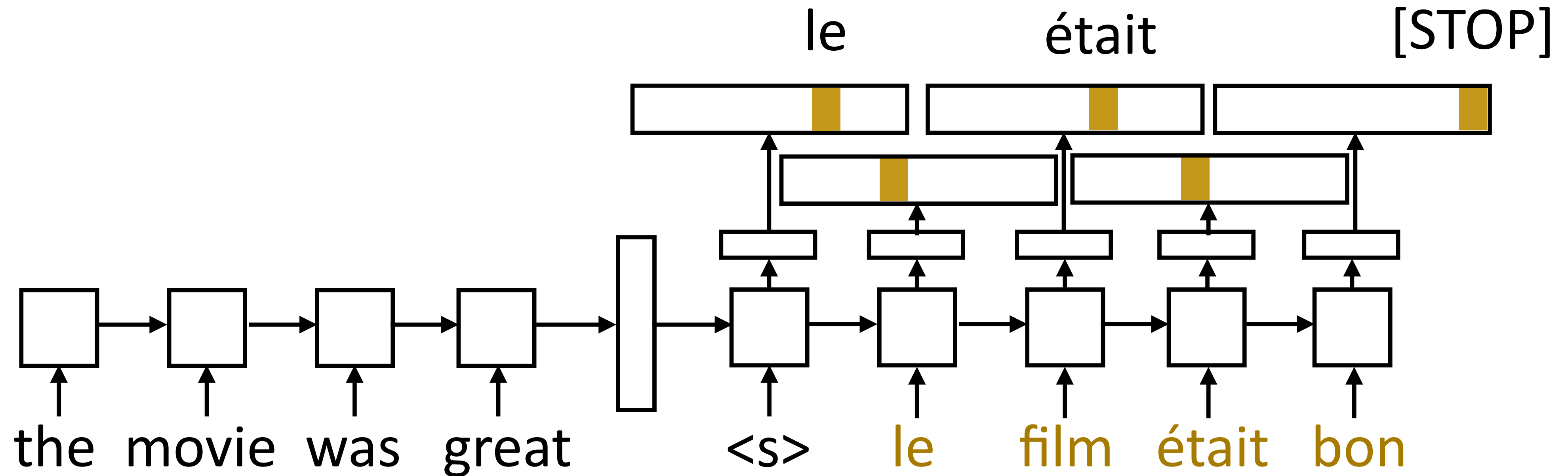


$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = \text{softmax}(W \bar{h})$$

$$P(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^n P(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$$

Decoder has separate parameters from encoder, so this can learn to be a language model (produce a plausible next word given current one)

Recall: Seq2seq Training



- ▶ Objective: maximize $\sum_{(\mathbf{x}, \mathbf{y})} \sum_{i=1}^n \log P(y_i^* | \mathbf{x}, y_1^*, \dots, y_{i-1}^*)$
- ▶ Teacher forcing: feed the correct word regardless of model's prediction (most typical way to train)

Recall: Semantic Parsing as Translation

“what states border Texas”



`lambda x (state (x) and border (x , e89)))`

- ▶ Write down a linearized form of the semantic parse, train seq2seq models to directly translate into this representation
- ▶ No need to have an explicit grammar, simplifies algorithms
- ▶ Might not produce well-formed logical forms, might require lots of data

This Lecture

- ▶ Attention for sequence-to-sequence models
- ▶ Copy mechanisms for copying words to the output
- ▶ Transformer architecture

Attention

Problems with Seq2seq Models

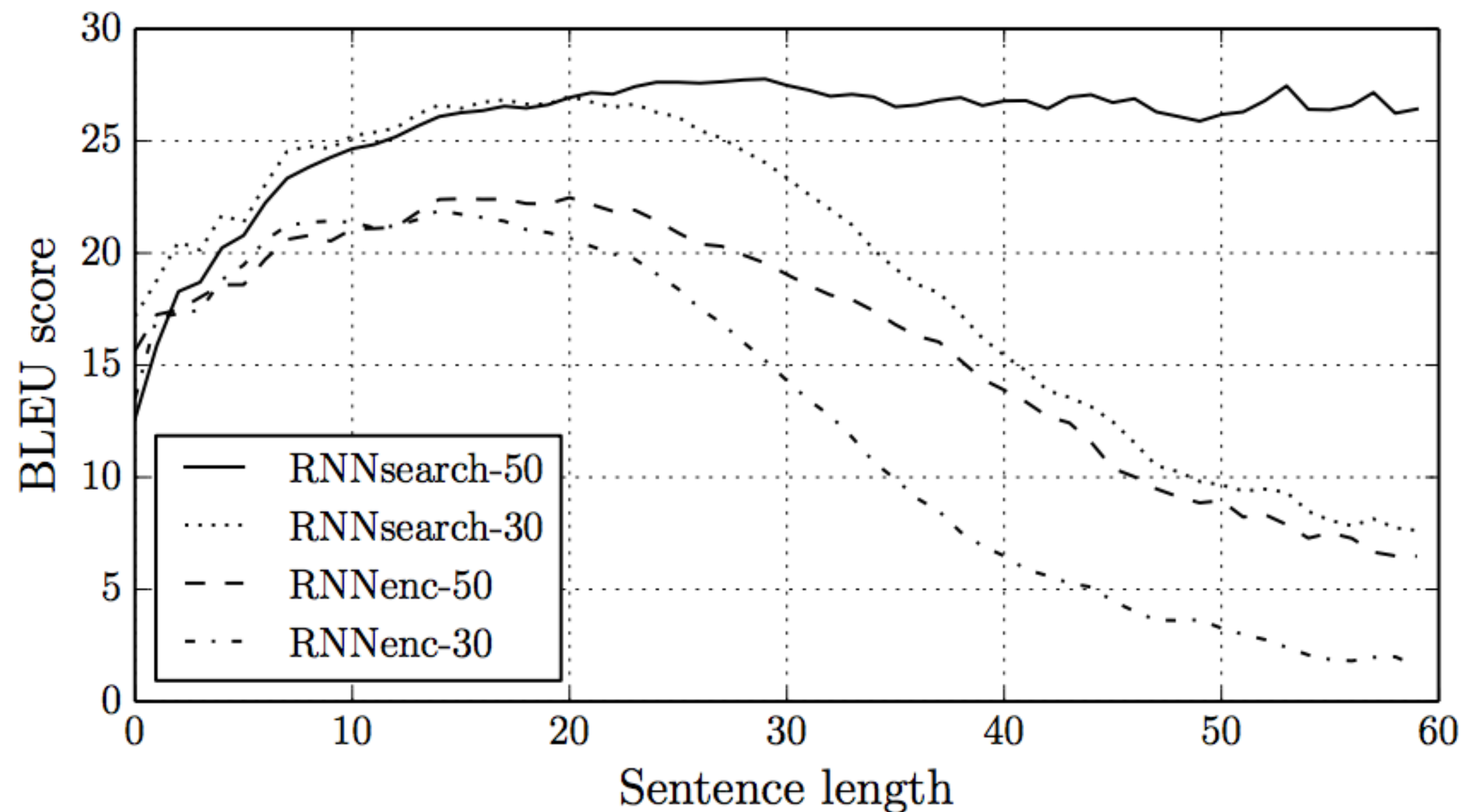
- ▶ Encoder-decoder models like to repeat themselves:

Un garçon joue dans la neige → A boy plays in the snow **boy plays boy plays**

- ▶ When/why does this happen?
 - ▶ Models trained poorly
 - ▶ LSTM state is not behaving as expected so it gets stuck in a “loop” of generating the same output tokens again and again
- ▶ Need some notion of input coverage or what input words we’ve translated

Problems with Seq2seq Models

- ▶ Bad at long sentences: 1) a fixed-size hidden representation doesn't scale; 2) LSTMs still have a hard time remembering for really long periods of time



RNNenc: the model we've discussed so far
RNNsearch: uses attention

Problems with Seq2seq Models

- ▶ Unknown words:

en: The ecotax portico in Pont-de-Buis, ... [truncated] ..., was taken down on Thursday morning

fr: Le portique écotaxe de Pont-de-Buis, ... [truncated] ..., a été démonté jeudi matin

nn: Le unk de unk à unk, ... [truncated] ..., a été pris le jeudi matin

- ▶ Encoding these rare words into a vector space is really hard
- ▶ In fact, we don't want to encode them, we want a way of directly looking back at the input and copying them (e.g., *Pont-de-Buis*)

Jean et al. (2015), Luong et al. (2015)

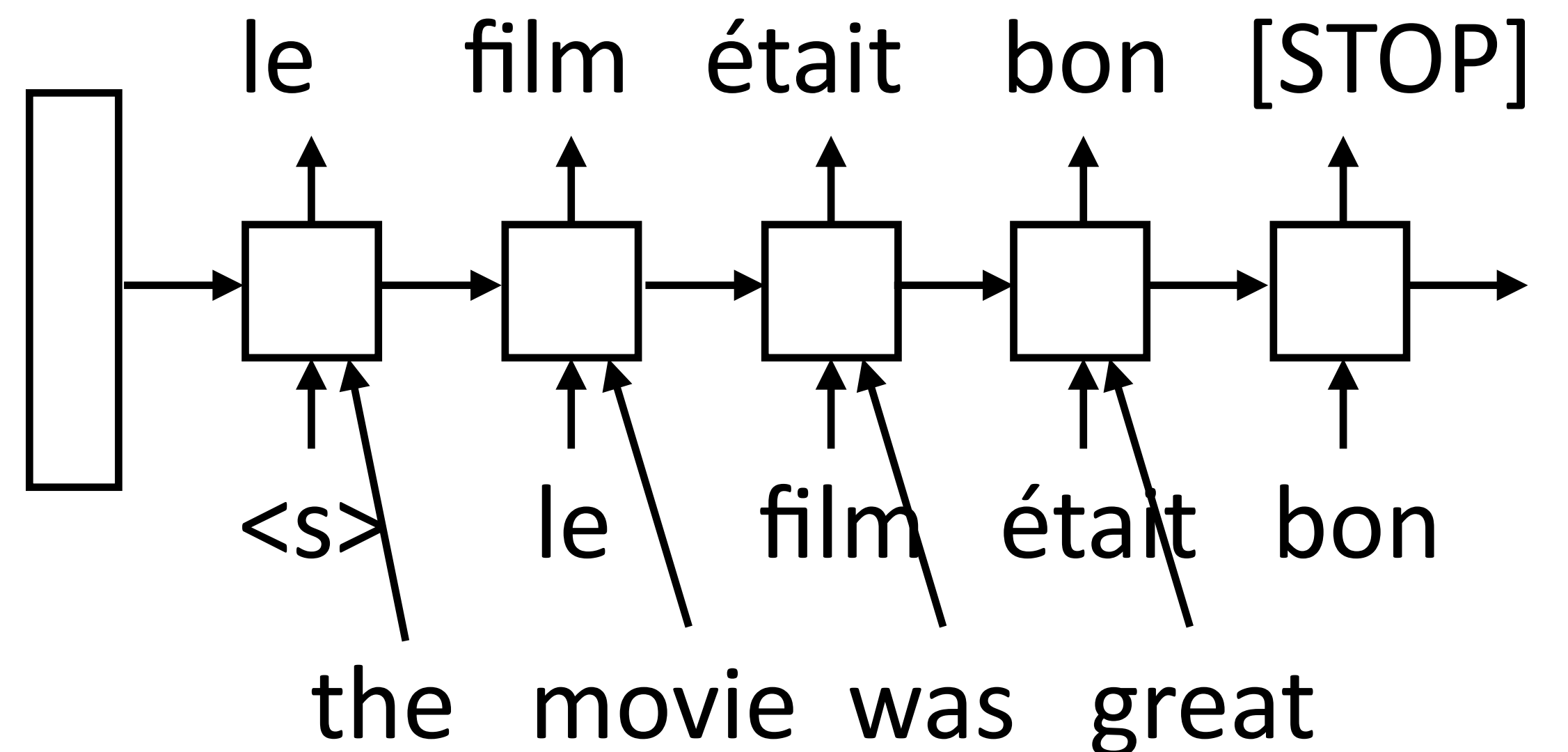
Aligned Inputs

- ▶ Suppose we knew the source and target would be word-by-word translated

the movie was great
/ / / /
le film était bon

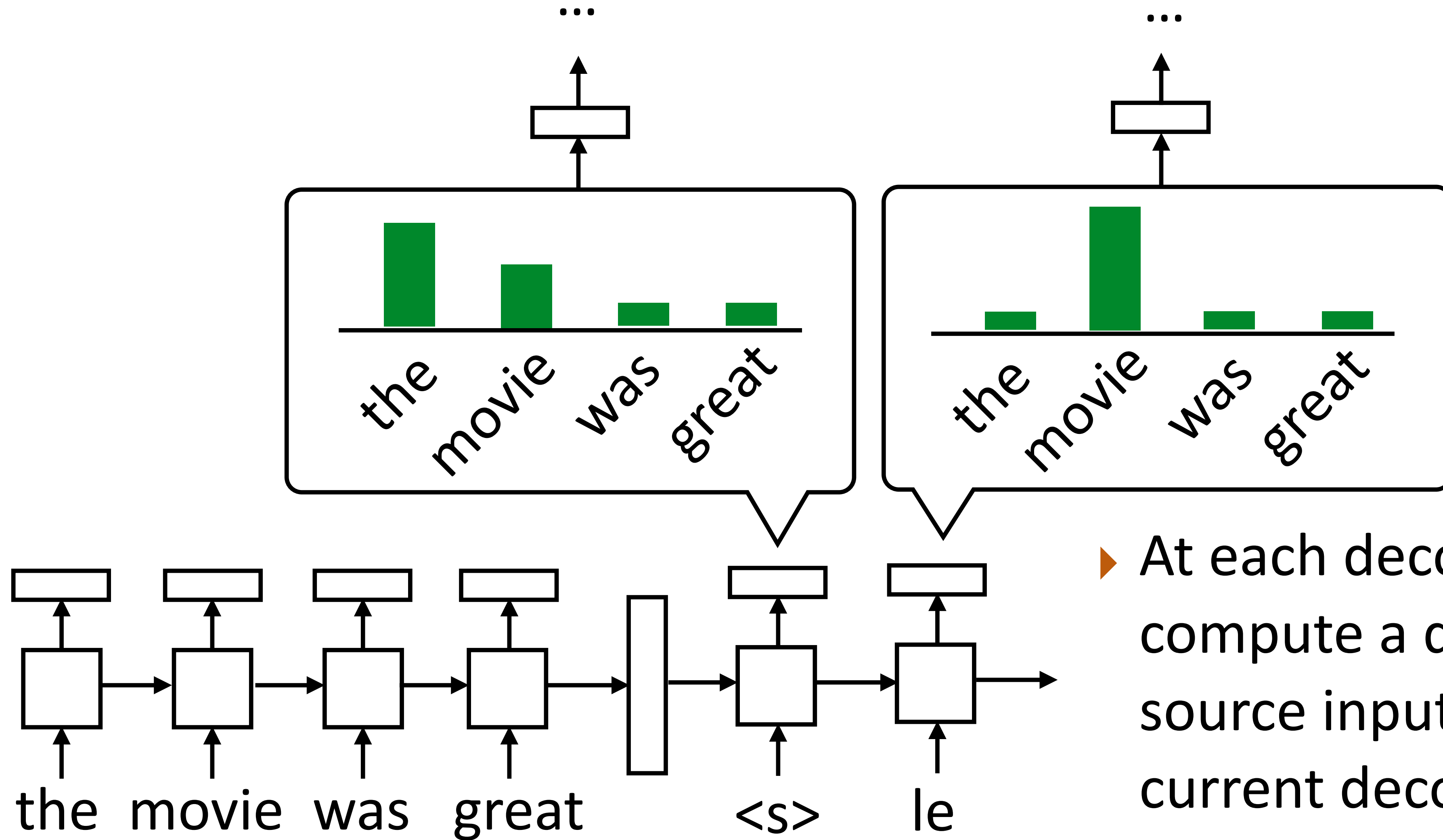
- ▶ Can look at the corresponding input word when translating — this could scale!

- ▶ Much less burden on the hidden state



- ▶ How can we achieve this without hardcoding it?

Attention

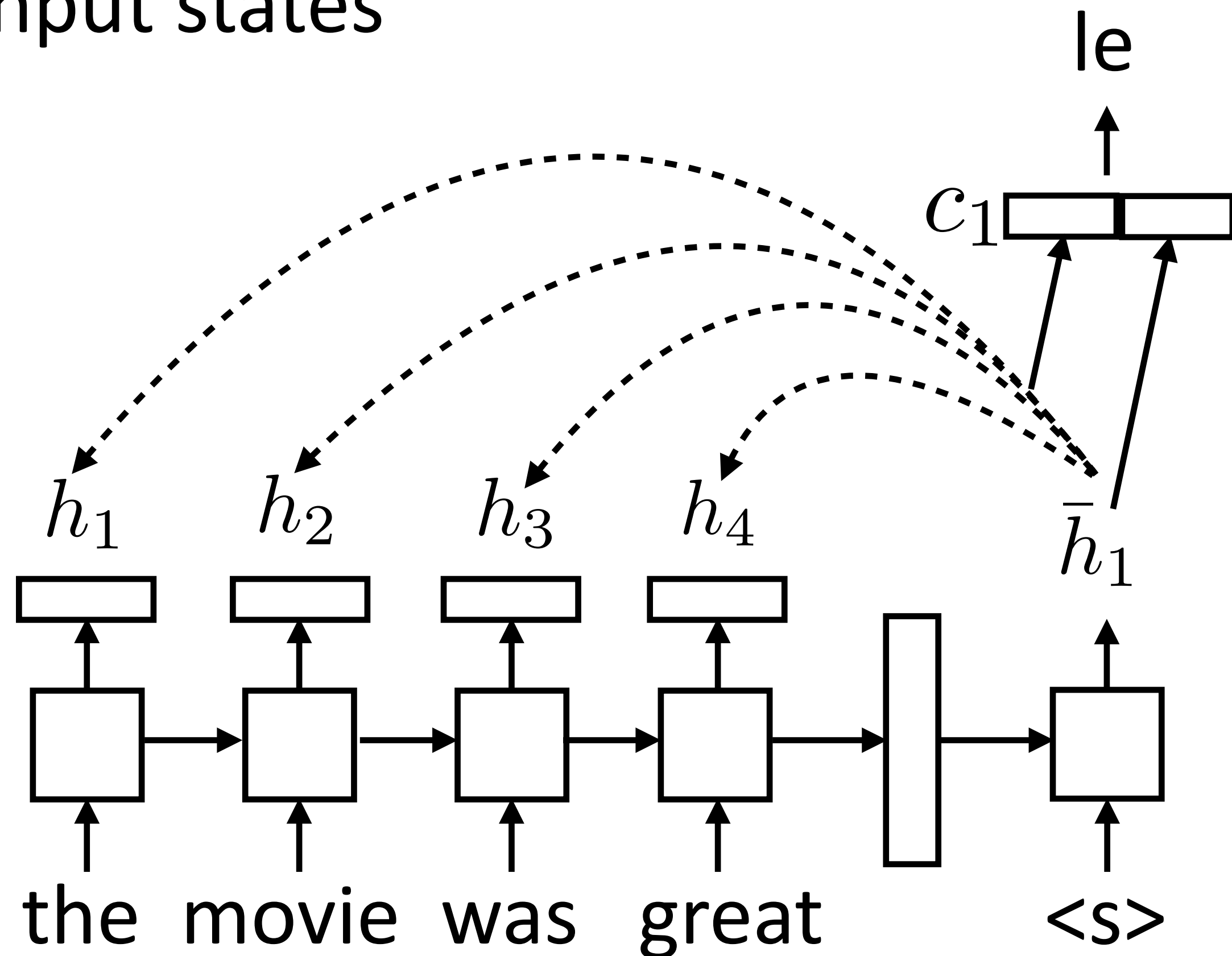


- ▶ At each decoder state, compute a distribution over source inputs based on current decoder state, use that in output layer

Attention

- ▶ For each decoder state, compute weighted sum of input states

- ▶ No attn: $P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = \text{softmax}(W \bar{h}_i)$



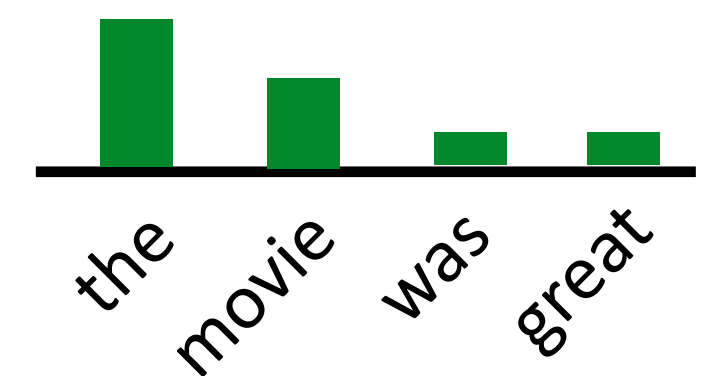
$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = \text{softmax}(W [c_i; \bar{h}_i])$$

$$c_i = \sum_j \alpha_{ij} h_j$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

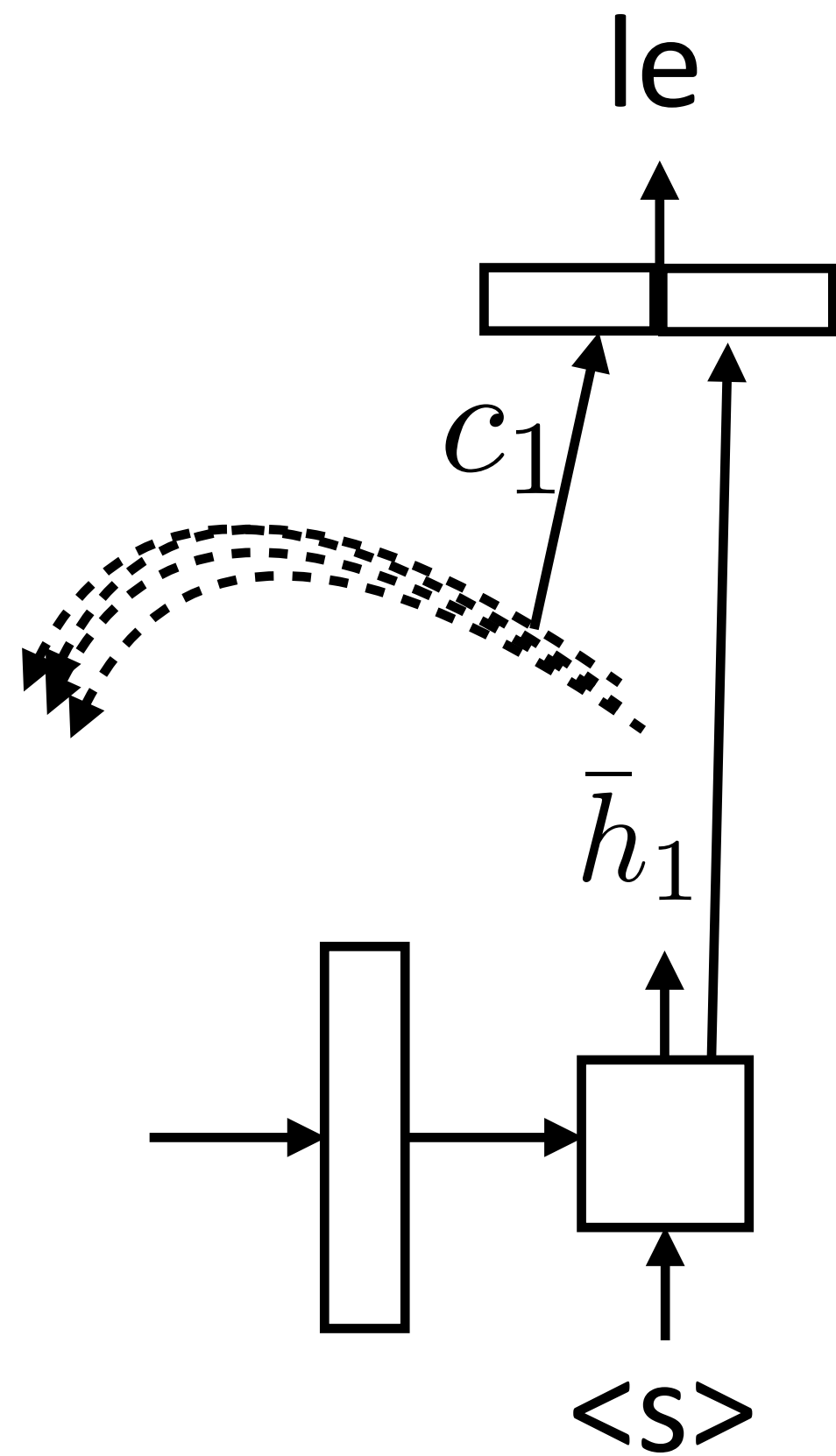
$$e_{ij} = f(\bar{h}_i, h_j)$$

- ▶ Weighted sum of input hidden states (vector)



- ▶ Some function f (TBD)

Attention



$$c_i = \sum_j \alpha_{ij} h_j$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

$$e_{ij} = f(\bar{h}_i, h_j)$$

potential choices:

$$f(\bar{h}_i, h_j) = \bar{h}_i^\top \cdot h_j$$

► Luong+ (2015): dot product

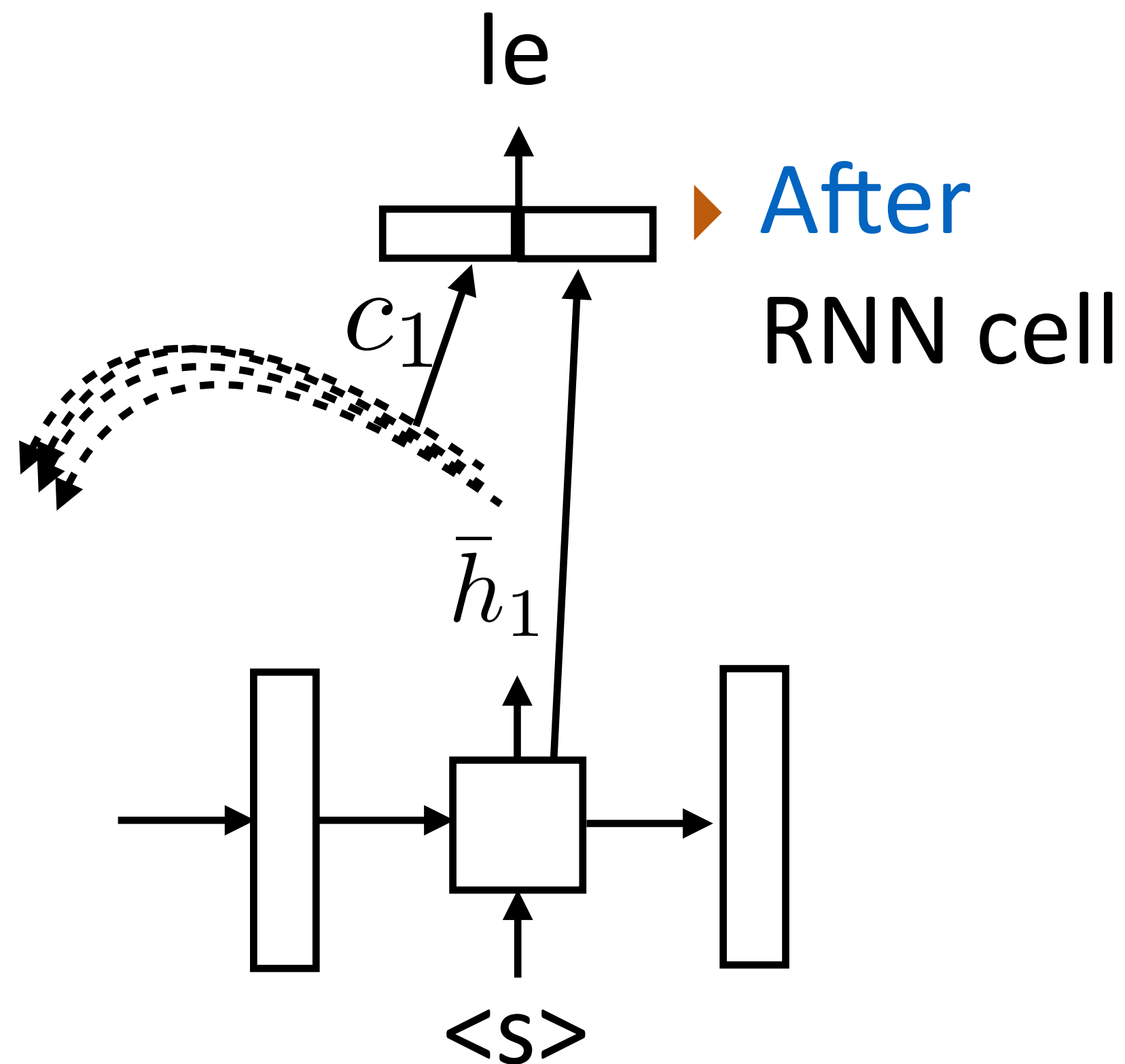
$$f(\bar{h}_i, h_j) = \bar{h}_i^\top W h_j$$

► Luong+ (2015): bilinear

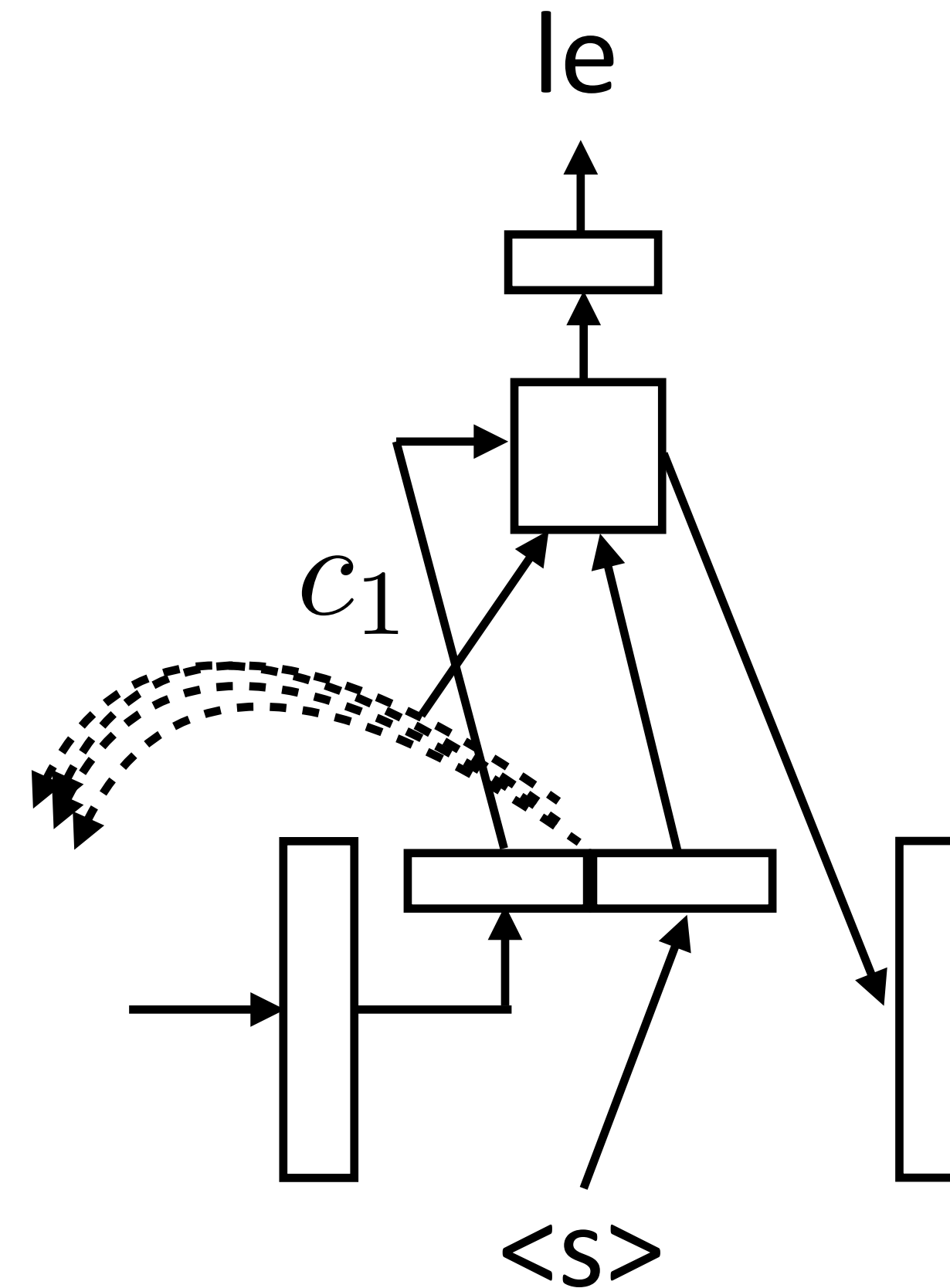
► Note that this all uses outputs of hidden layers

Alternatives

- ▶ When do we compute attention? Can compute before or after RNN cell



▶ **After**
RNN cell



- ▶ **Before** RNN cell; this one is a little more convoluted and less standard

Luong et al. (2015):

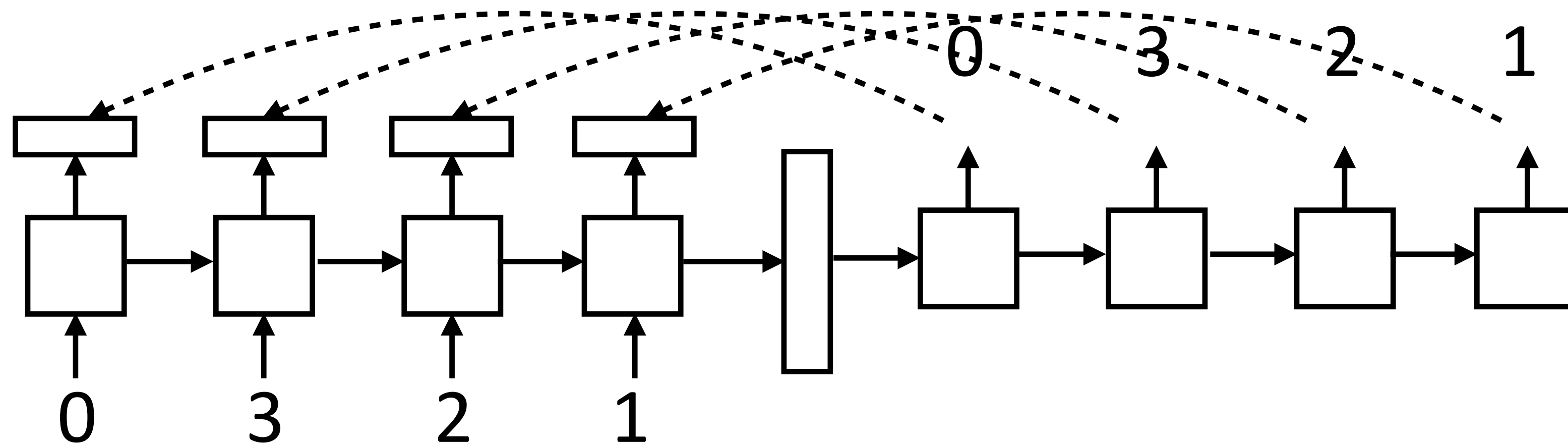
<https://arxiv.org/pdf/1508.04025.pdf>

Bahdanau et al. (2015)

<https://arxiv.org/pdf/1409.0473.pdf>

What can attention do?

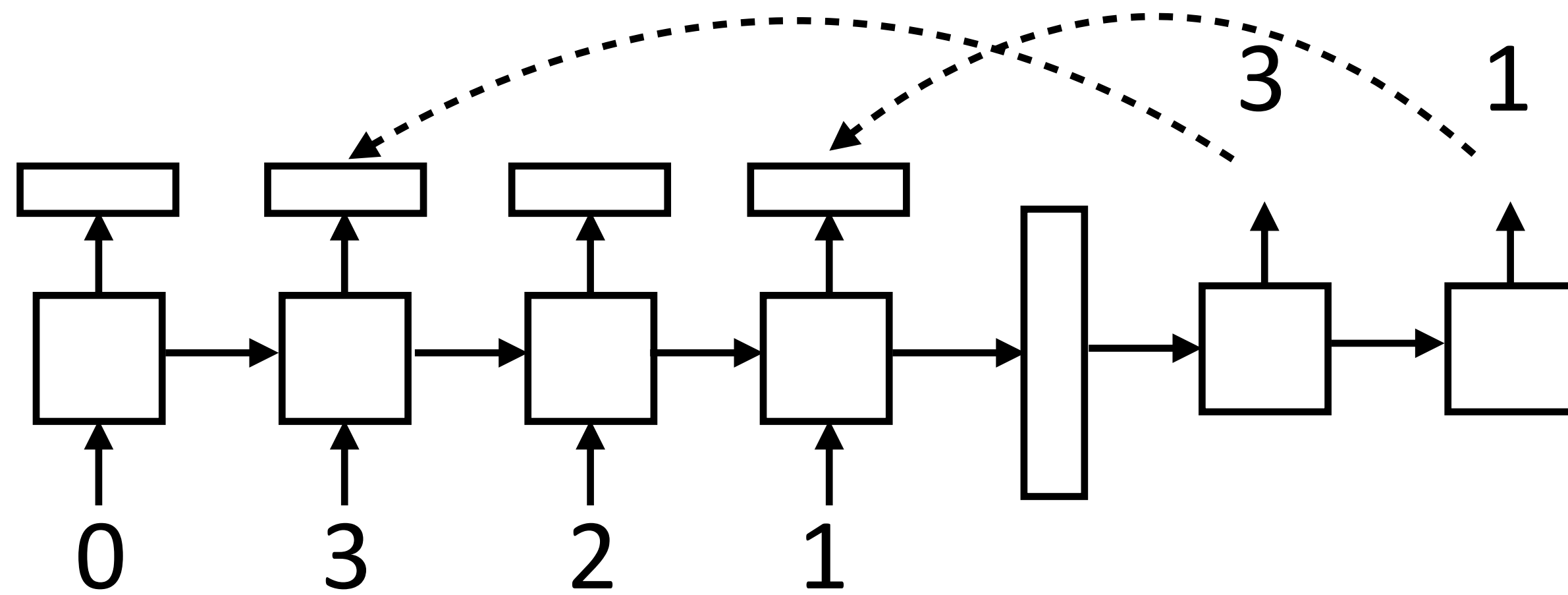
- ▶ Learning to copy — how might this work?



- ▶ LSTM can learn to count with the right weight matrix
- ▶ This is a kind of position-based addressing

What can attention do?

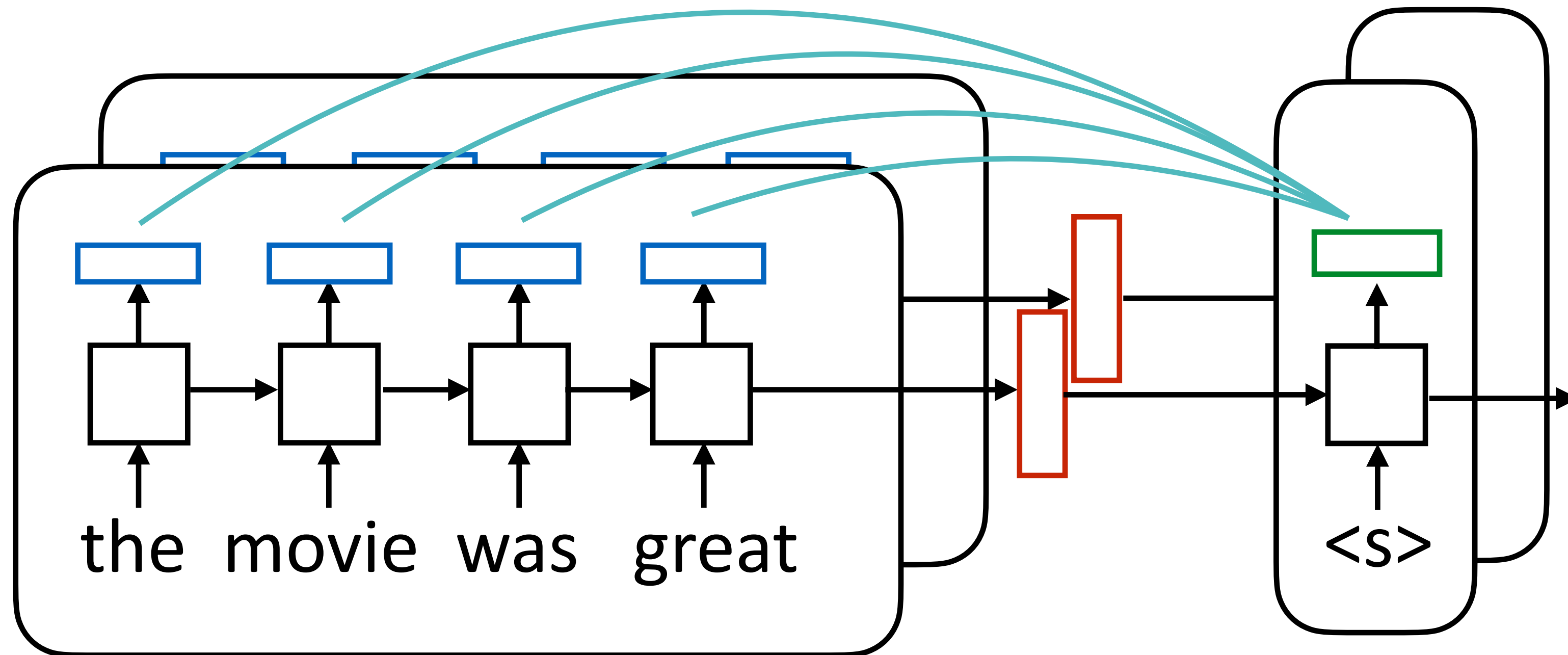
- ▶ Learning to subsample tokens



- ▶ Need to count (for ordering) and also determine which tokens are in/out
- ▶ Content-based addressing

Batching Attention

token outputs: batch size x sentence length x hidden size



hidden state: batch size x hidden size

$$e_{ij} = f(\bar{h}_i, h_j)$$
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{j'} \exp(e_{ij'})}$$

sentence outputs:
batch size x hidden size

attention scores = batch size x sentence length

$c =$ batch size x hidden size

$$c_i = \sum_j \alpha_{ij} h_j$$

- ▶ Make sure tensors are the right size!

Luong et al. (2015)

Results

- ▶ Machine translation: BLEU score of 14.0 on English-German -> 16.8 with attention, 19.0 with smarter attention (we'll come back to this later)
- ▶ Summarization/headline generation: bigram recall from 11% -> 15%
- ▶ Semantic parsing: ~30-50% accuracy -> 70+% accuracy on Geoquery

Luong et al. (2015)

Chopra et al. (2016)

Jia and Liang (2016)

Copying Input/Pointers

Unknown Words

en: The ecotax portico in Pont-de-Buis, ... [truncated] ..., was taken down on Thursday morning

fr: Le portique écotaxe de Pont-de-Buis, ... [truncated] ..., a été démonté jeudi matin

nn: Le unk de unk à unk, ... [truncated] ..., a été pris le jeudi matin

- ▶ Want to be able to copy named entities like Pont-de-Buis

$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = \text{softmax}(W[c_i; \bar{h}_i])$$

from attention from RNN
hidden state

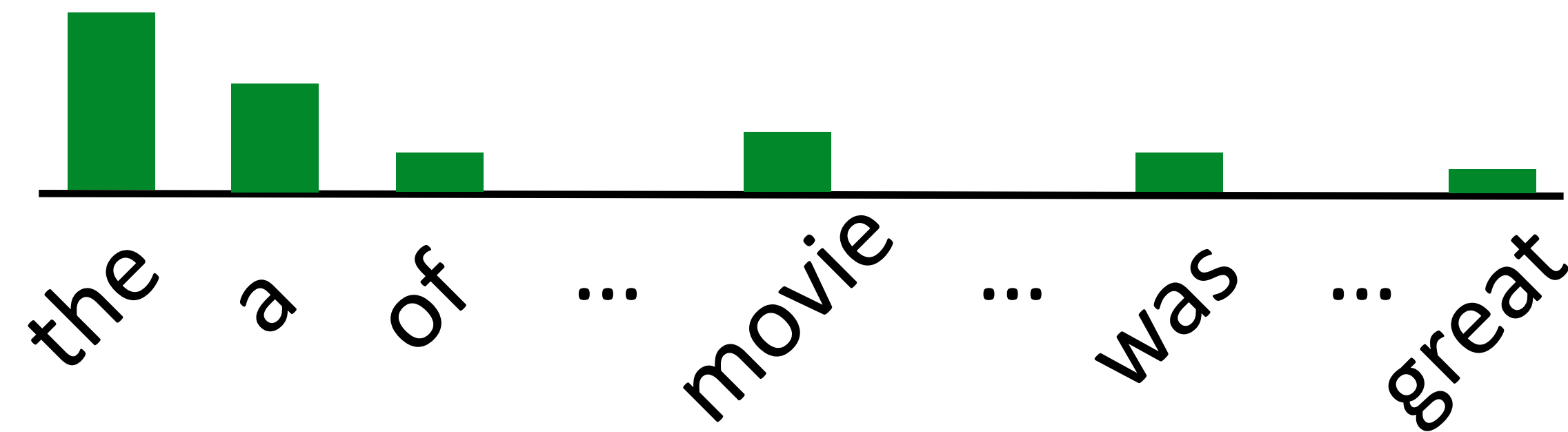
- ▶ Problems: target word has to be in the vocabulary, attention + RNN need to generate good embedding to pick it

Jean et al. (2015), Luong et al. (2015)

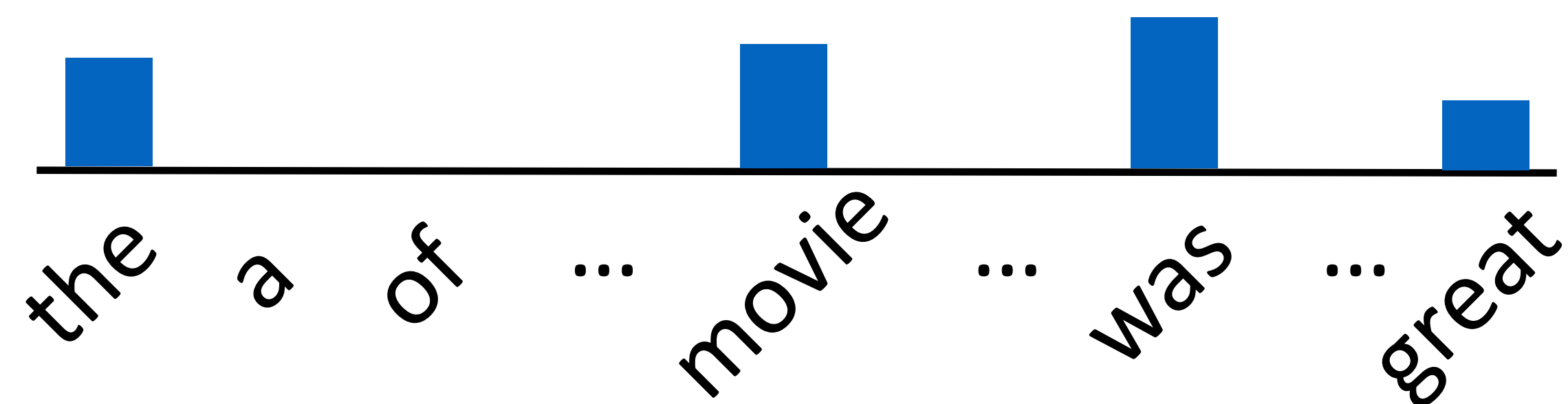
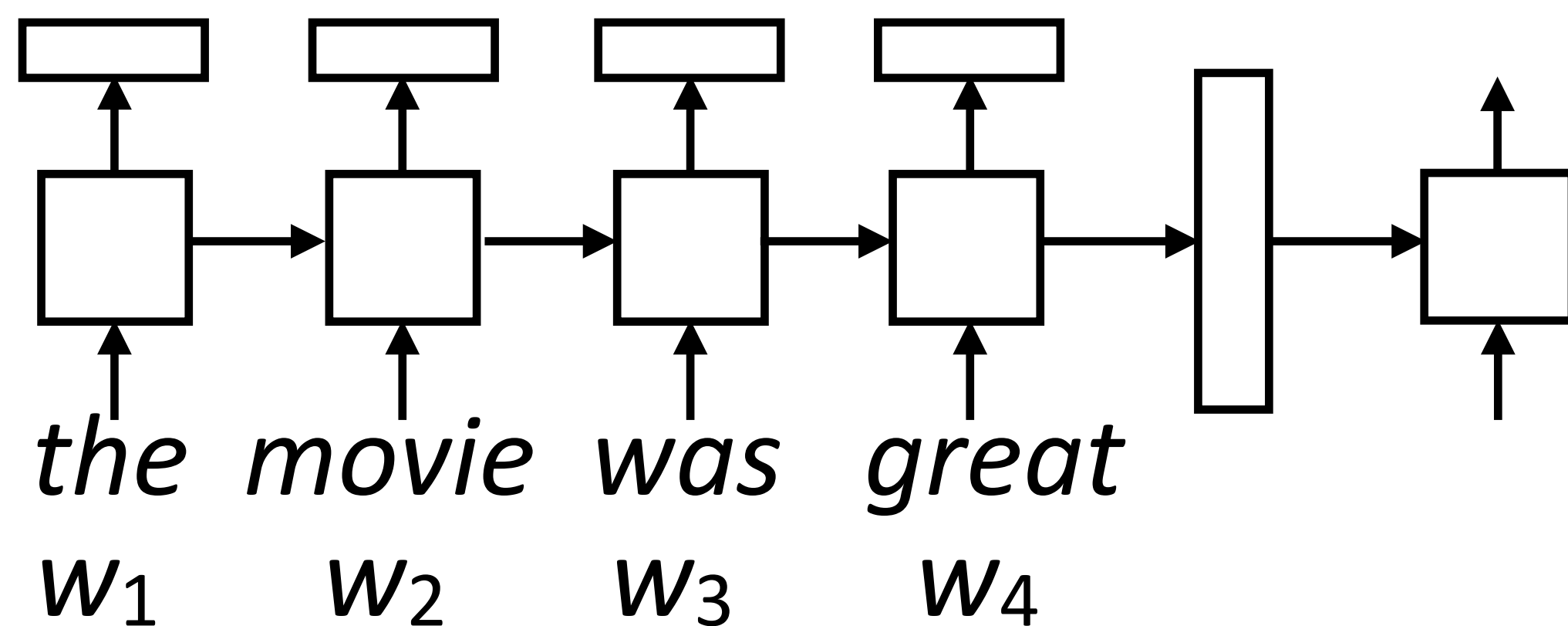
Pointer Networks

$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = \text{softmax}(W[c_i; \bar{h}_i])$$

- ▶ Standard decoder (P_{vocab}): softmax over vocabulary, all words get >0 prob
- ▶ Pointer network: predict from *source words* instead of *target vocab*



$$P_{\text{pointer}}(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) \propto \begin{cases} h_j^\top V \bar{h}_i & \text{if } y_i = w_j \\ 0 & \text{otherwise} \end{cases}$$

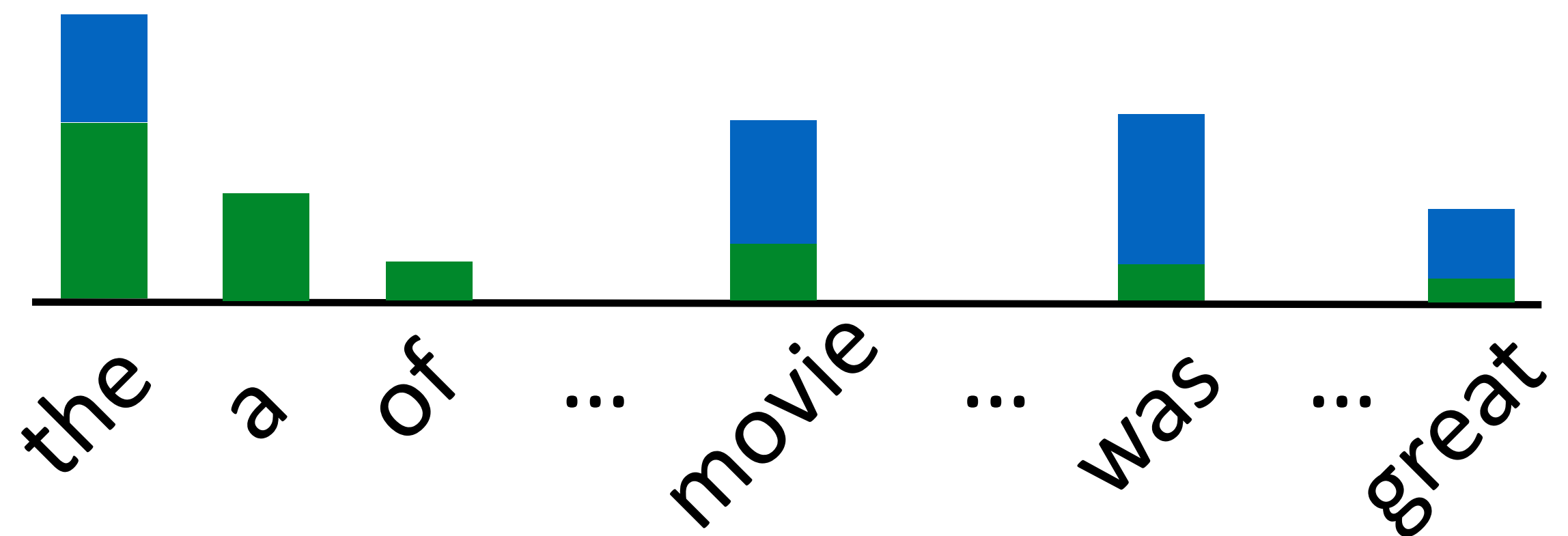


Pointer Generator Mixture Models

- ▶ Define the decoder model as a mixture model of the P_{vocab} and P_{pointer} models (previous slide)

$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = P(\text{copy})P_{\text{pointer}} + (1 - P(\text{copy}))P_{\text{vocab}}$$

- ▶ Predict $P(\text{copy})$ based on decoder state, input, etc.
- ▶ Marginalize over copy variable during training and inference
- ▶ Model will be able to both generate and copy, flexibly adapt between the two



Copying

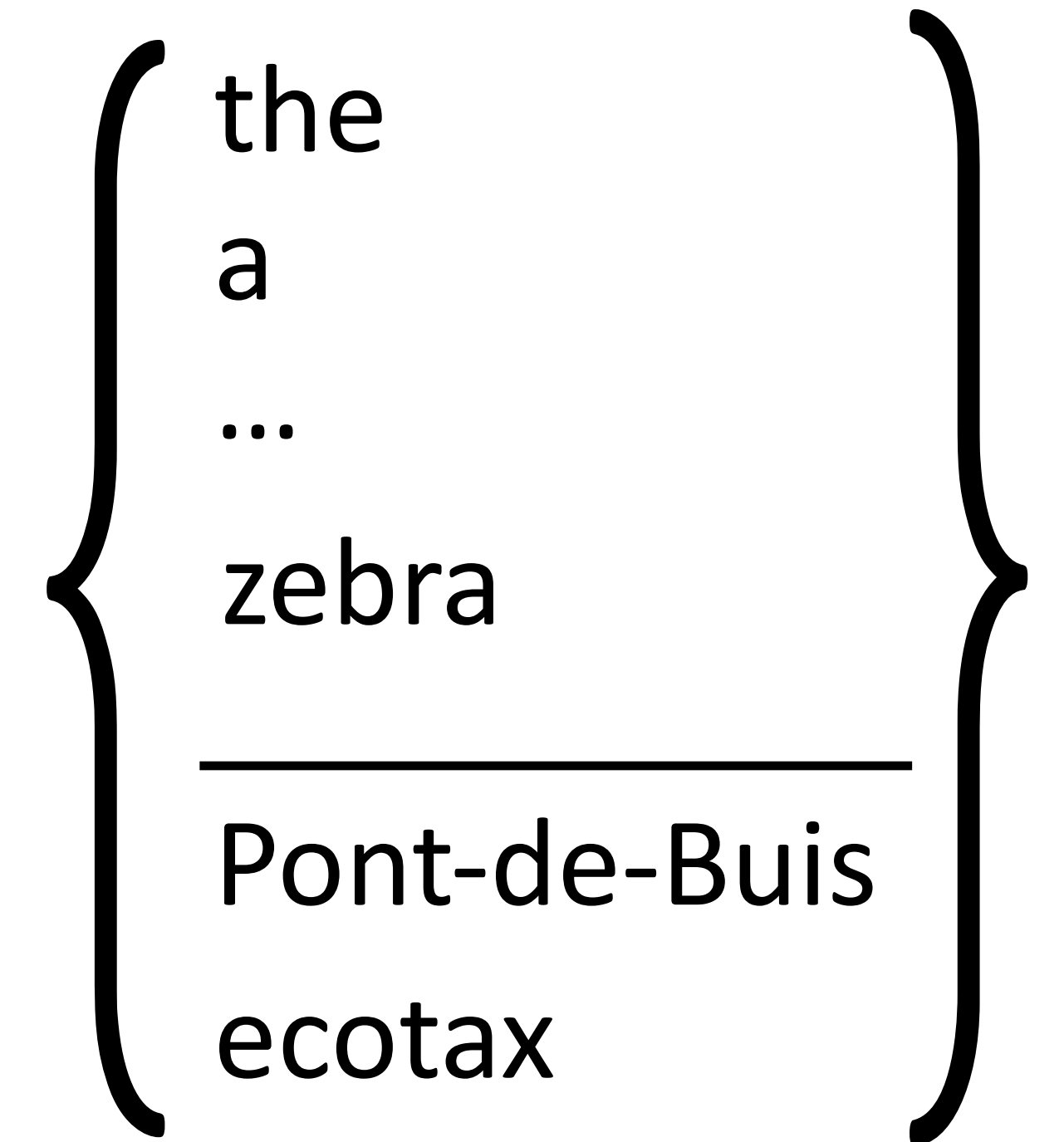
en: The ecotax portico in Pont-de-Buis , ... [truncated] ..

fr: Le portique écotaxe de Pont-de-Buis , ... [truncated] .

nn: Le unk de unk à unk , ... [truncated] ..., a été pris

- ▶ Some words we may want to copy may not be in the fixed output vocab (*Pont-de-Buis*)

- ▶ Solution: expand the vocabulary dynamically. New words can only be predicted by copying (always 0 probability under P_{vocab})



Results

	GEO	ATIS
No Copying	74.6	69.9
With Copying	85.0	76.3

- ▶ For semantic parsing, copying tokens from the input (*texas*) can be very useful
- ▶ Copying typically helps a bit, but attention captures most of the benefit. However, vocabulary expansion is critical for some tasks (machine translation)

Takeaways

- ▶ Attention is very helpful for seq2seq models
- ▶ Explicitly copying input can be beneficial as well
- ▶ Transformers are strong models we'll come back to later
- ▶ <https://arxiv.org/abs/1706.03762> Vaswani et al. (2017)