

Deep Learning
Sequence to Sequence models:
Attention Models

Sequence-to-sequence modelling

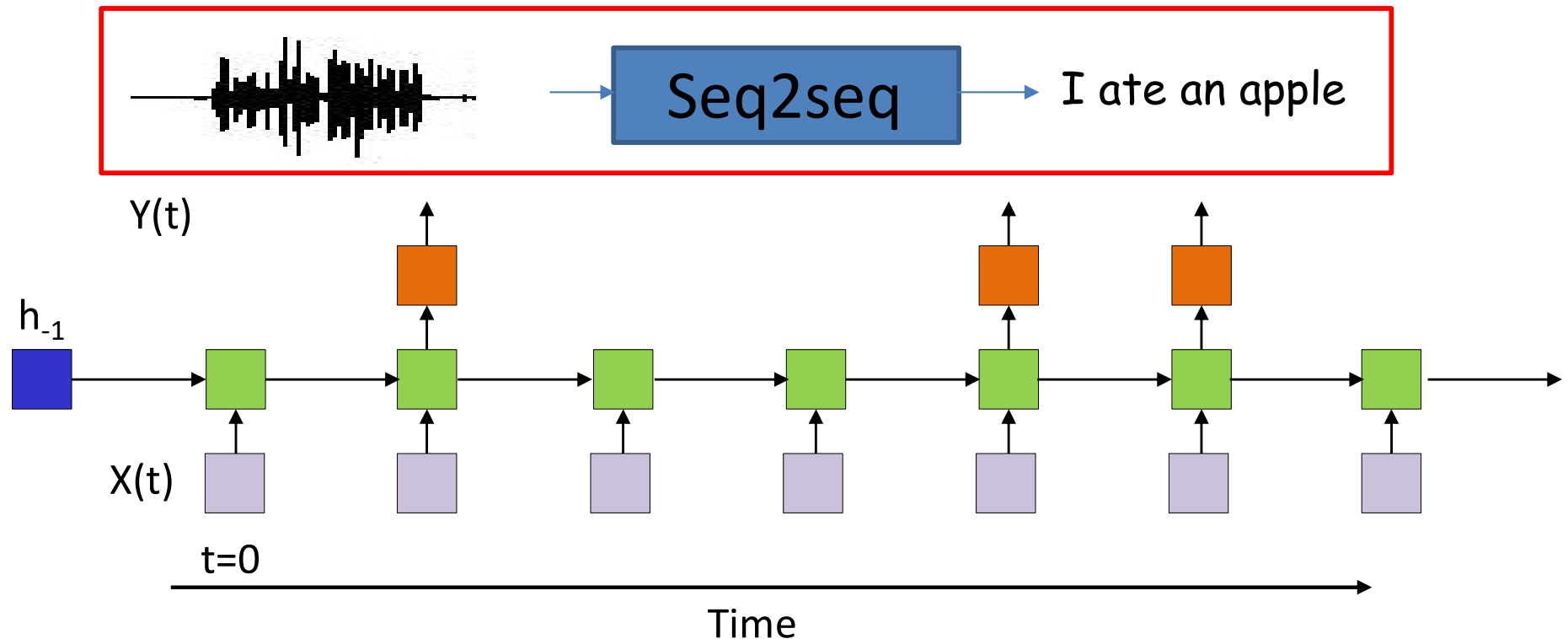
- Problem:
 - A sequence $X_1 \dots X_N$ goes in
 - A different sequence $Y_1 \dots Y_M$ comes out
- E.g.
 - Speech recognition: Speech goes in, a word sequence comes out
 - Alternately output may be phoneme or character sequence
 - Machine translation: Word sequence goes in, word sequence comes out
 - Dialog : User statement goes in, system response comes out
 - Question answering : Question comes in, answer goes out
- In general $N \neq M$
 - No synchrony between X and Y .

Sequence to sequence



- Sequence goes in, sequence comes out
- No notion of “time synchrony” between input and output
 - May even not even maintain order of symbols
 - E.g. “I ate an apple” → “Ich habe einen apfel gegessen”
 - Or even seem related to the input
 - E.g. “My screen is blank” → “Please check if your computer is plugged in.”


Recap: Have dealt with the “aligned” case: CTC



- The input and output sequences happen in the same order
 - Although they may be asynchronous
 - Order-correspondence, but no time synchrony
 - E.g. Speech recognition
 - The input speech corresponds to the phoneme sequence output

Today



- Sequence goes in, sequence comes out
- No notion of “time synchrony” between input and output
 - May even not even maintain order of symbols
 - E.g. “I ate an apple” → “Ich habe einen apfel gegessen”
 - Or even seem related to the input
 - E.g. “My screen is blank” → “Please check if your computer is plugged in.”

Recap: Predicting text

- Simple problem: Given a series of symbols (characters or words) $w_1 w_2 \dots w_n$, predict the next symbol (character or word) w_{n+1}

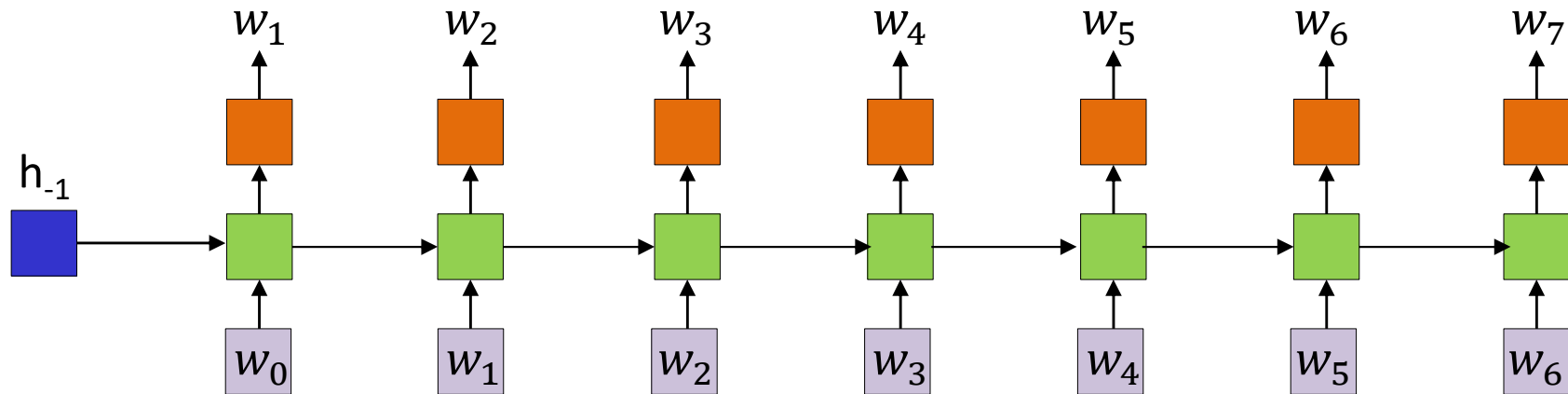
Language modelling using RNNs

Four score and seven years ???

A B R A H A M L I N C O L ??

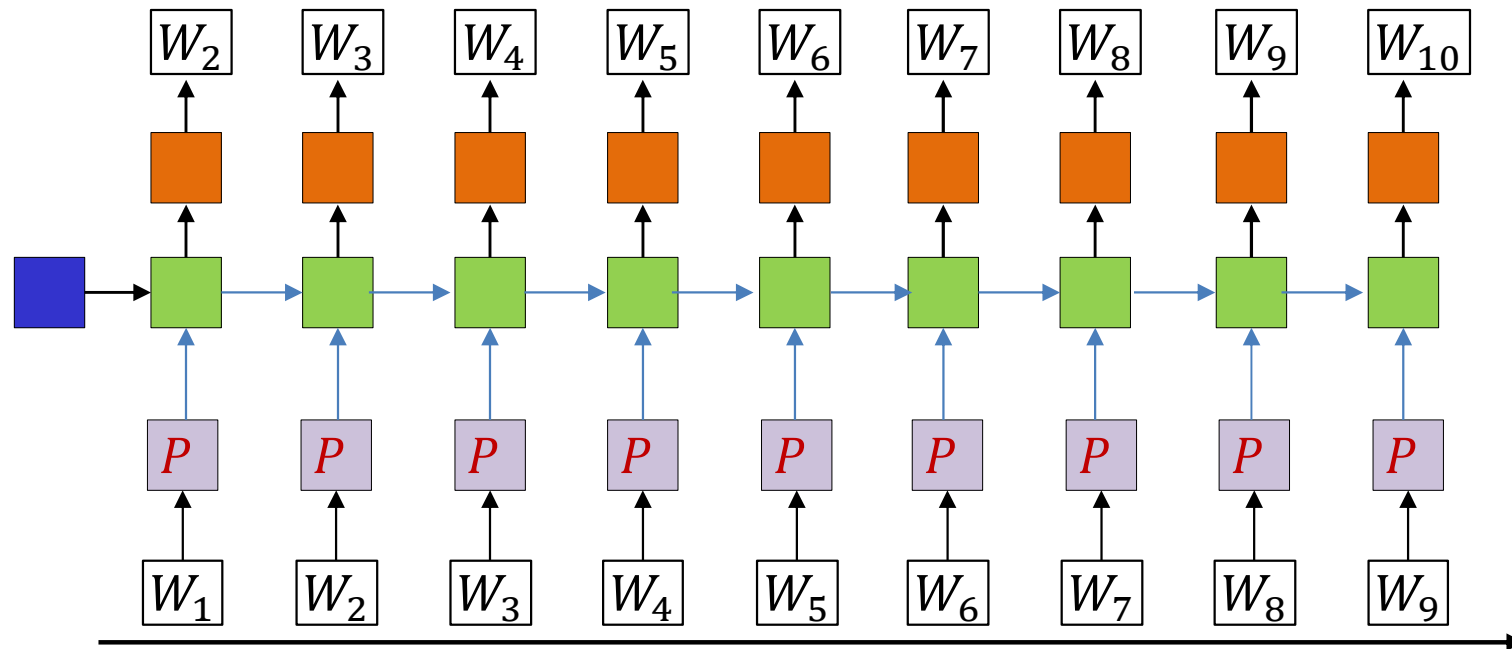
- Problem: Given a sequence of words (or characters) predict the next one
 - The problem of learning the sequential structure of language

Simple recurrence : Text Modelling



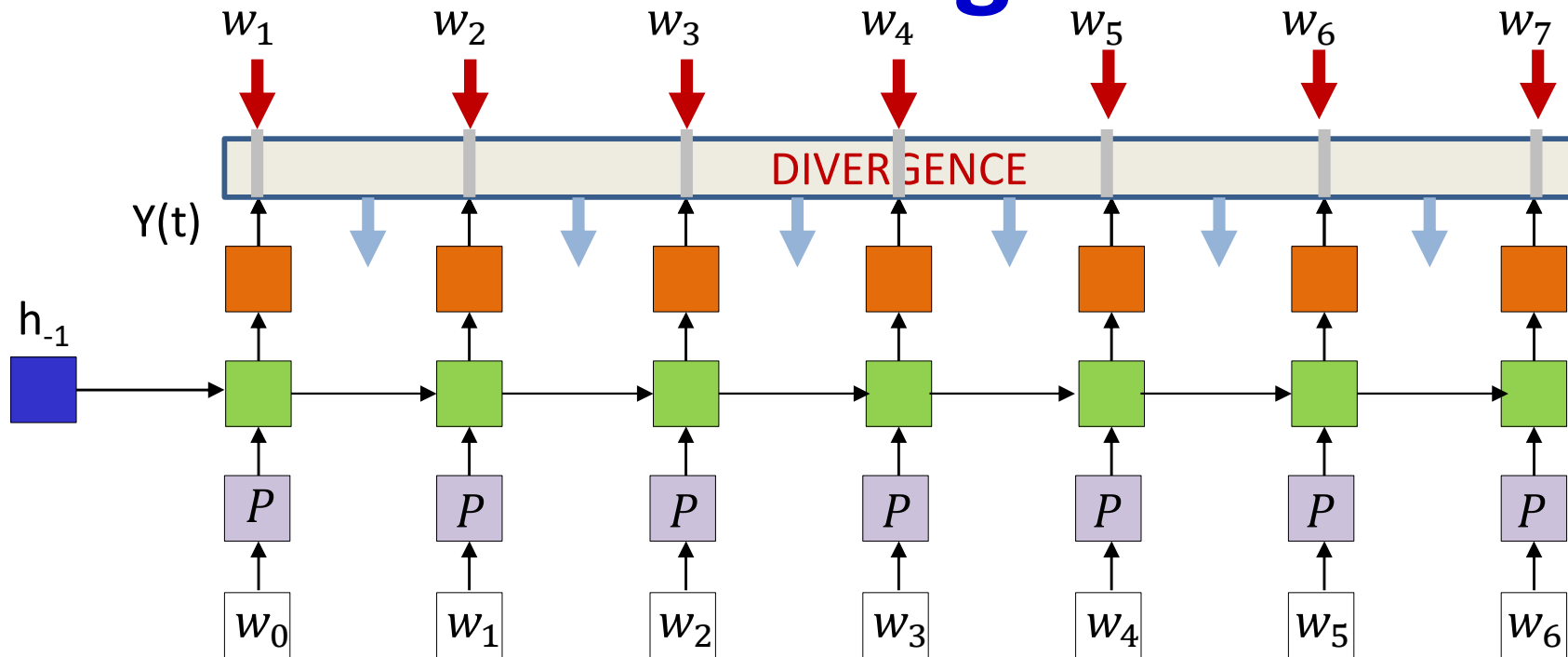
- Learn a model that can predict the next symbol given a sequence of symbols
 - Characters or words
- After observing inputs $w_0 \dots w_k$ it predicts w_{k+1}
 - In reality, outputs a probability distribution for w_{k+1}

Generating Language: The model



- Input: symbols as one-hot vectors
 - Dimensionality of the vector is the size of the “vocabulary”
 - Projected down to lower-dimensional “embeddings”
- The hidden units are (one or more layers of) LSTM units
- Output at each time: A probability distribution for the next word in the sequence
- All parameters are trained via backpropagation from a lot of text

Training



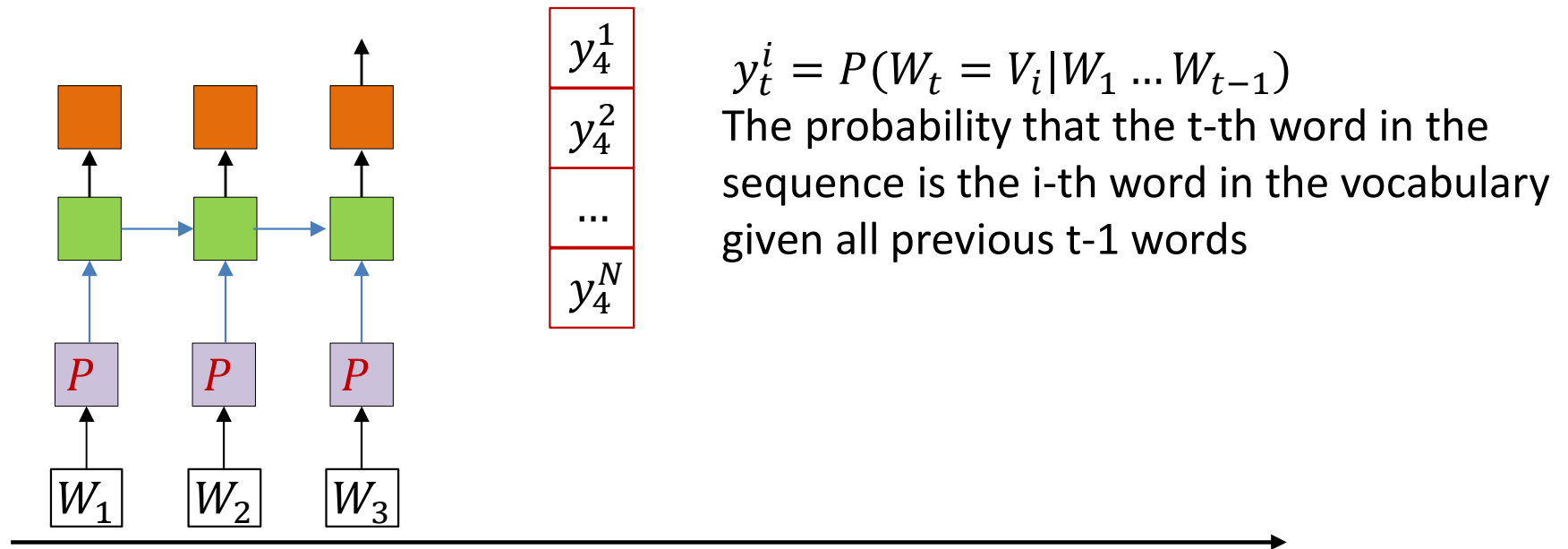
- Input: symbols as one-hot vectors
 - Dimensionality of the vector is the size of the “vocabulary”
 - Projected down to lower-dimensional “embeddings”
- Output: Probability distribution over symbols

$$Y(t, i) = P(V_i | w_0 \dots w_{t-1})$$
 - V_i is the i -th symbol in the vocabulary
- Divergence

The probability assigned to the correct next word

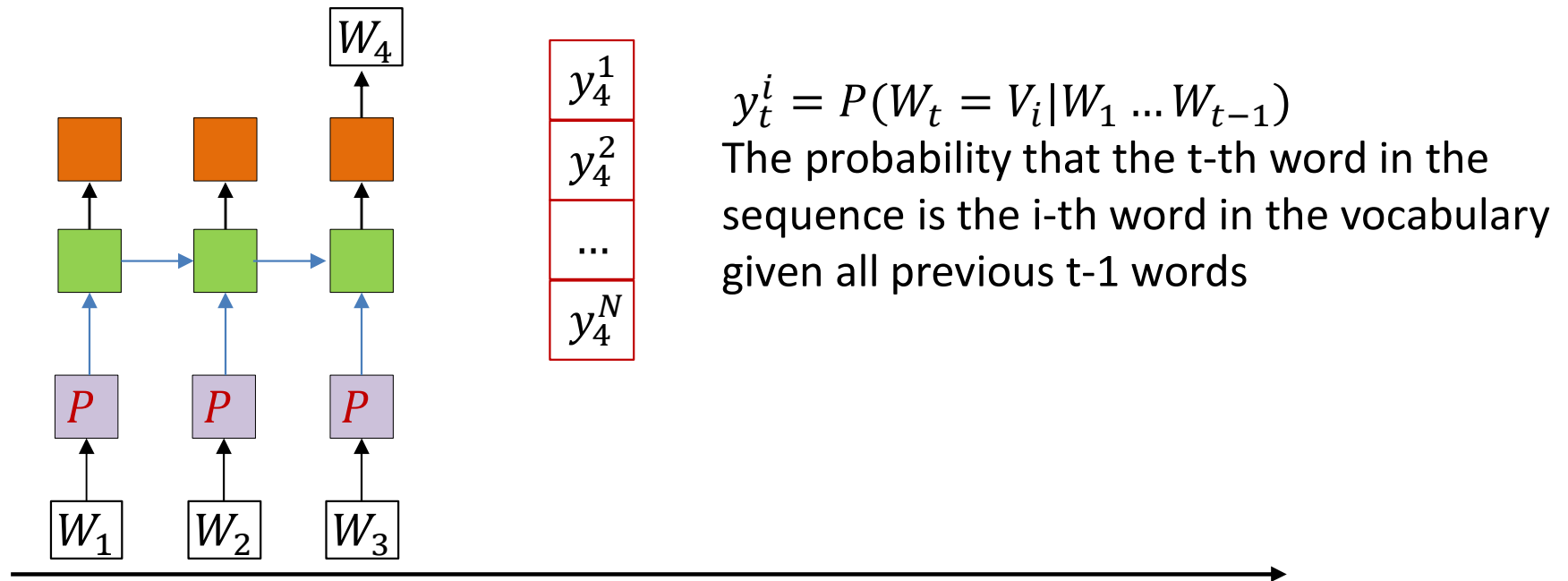
$$Div(w(1 \dots T), Y(0 \dots T - 1)) = \sum_t KL(w(t + 1), Y(t)) = - \sum_t \log Y(t, w_{t+1})$$

Generating Language: Synthesis



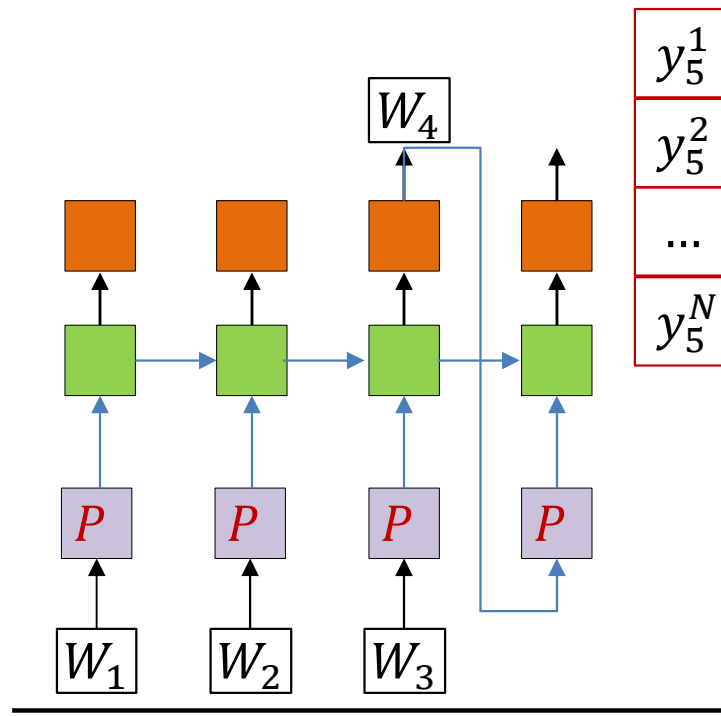
- On trained model : Provide the first few words
 - One-hot vectors
- After the last input word, the network generates a probability distribution over words
 - Outputs an N-valued probability distribution rather than a one-hot vector

Generating Language: Synthesis



- On trained model : Provide the first few words
 - One-hot vectors
- After the last input word, the network generates a probability distribution over words
 - Outputs an N-valued probability distribution rather than a one-hot vector
- Draw a word by sampling from the distribution
 - And set it as the next word in the series

Generating Language: Synthesis

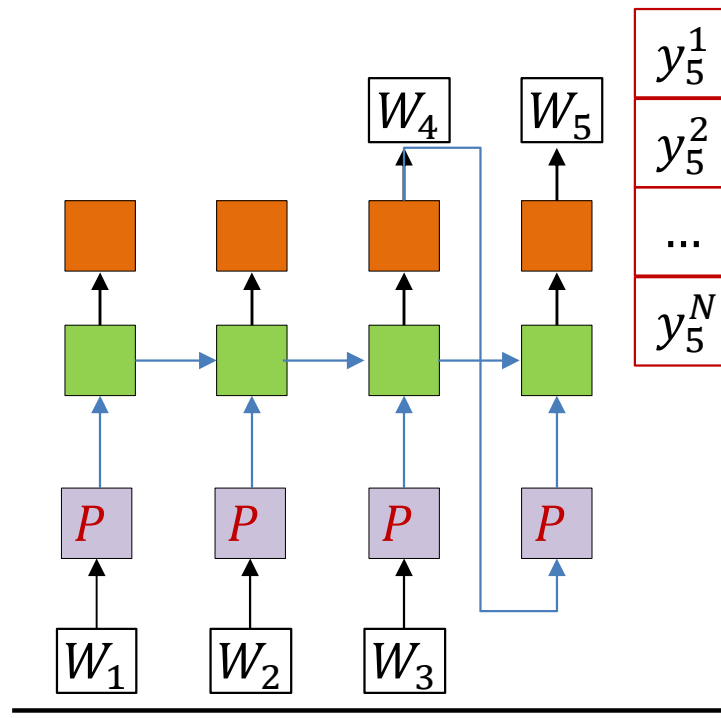


$$y_t^i = P(W_t = V_i | W_1 \dots W_{t-1})$$

The probability that the t -th word in the sequence is the i -th word in the vocabulary given all previous $t-1$ words

- Feed the drawn word as the next word in the series
 - And draw the next word by sampling from the output probability distribution

Generating Language: Synthesis

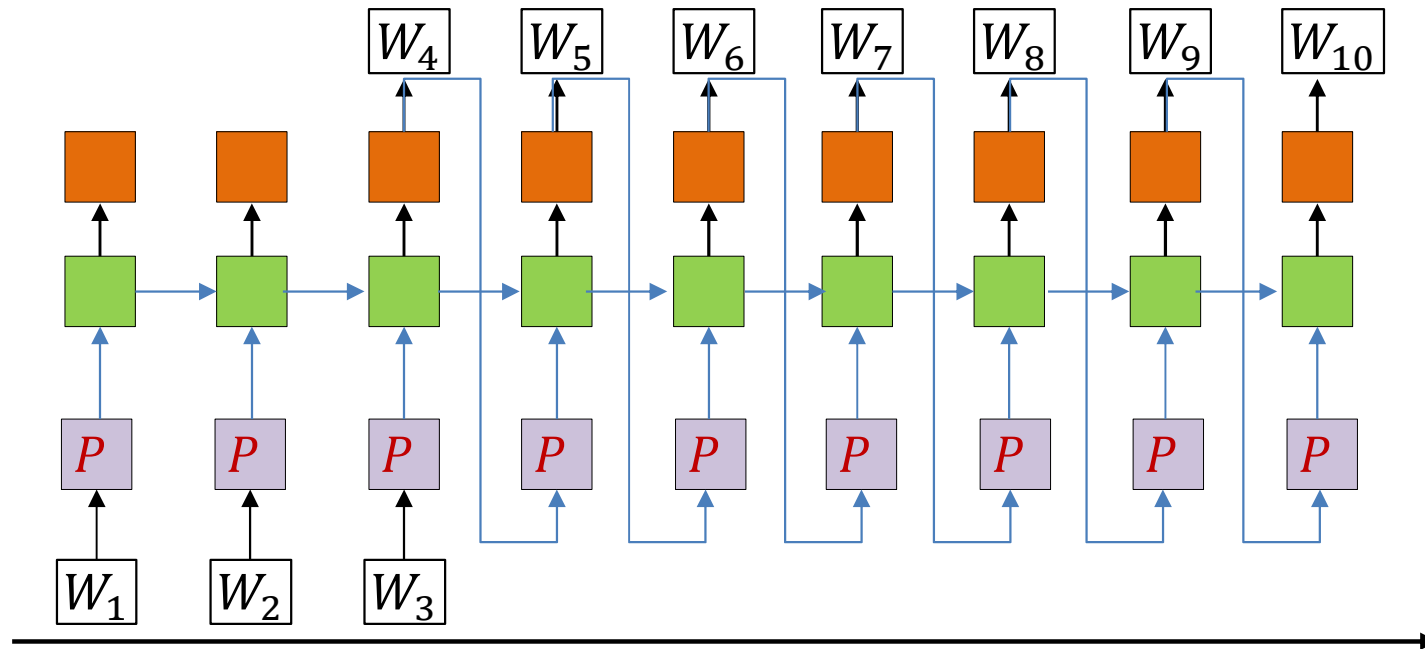


$$y_t^i = P(W_t = V_i | W_1 \dots W_{t-1})$$

The probability that the t -th word in the sequence is the i -th word in the vocabulary given all previous $t-1$ words

- Feed the drawn word as the next word in the series
 - And draw the next word by sampling from the output probability distribution

Generating Language: Synthesis



- Feed the drawn word as the next word in the series
 - And draw the next word by sampling from the output probability distribution
- **When do we stop?**

A note on beginnings and ends

- A sequence of words by itself does not indicate if it is a complete sentence or not

... four score and eight ...

- Unclear if this is the *start* of a sentence, the *end* of a sentence, or *both* (i.e. a complete sentence)
- To make it explicit, we will add two additional symbols (in addition to the words) to the base vocabulary
 - **<SOS>** : Indicates start of a sentence
 - **<EOS>** : Indicates end of a sentence

A note on beginnings and ends

- Some examples:

four score and eight

- This is clearly the middle of sentence

<sos> four score and eight

- This is a fragment from the *start* of a sentence

four score and eight <eos>

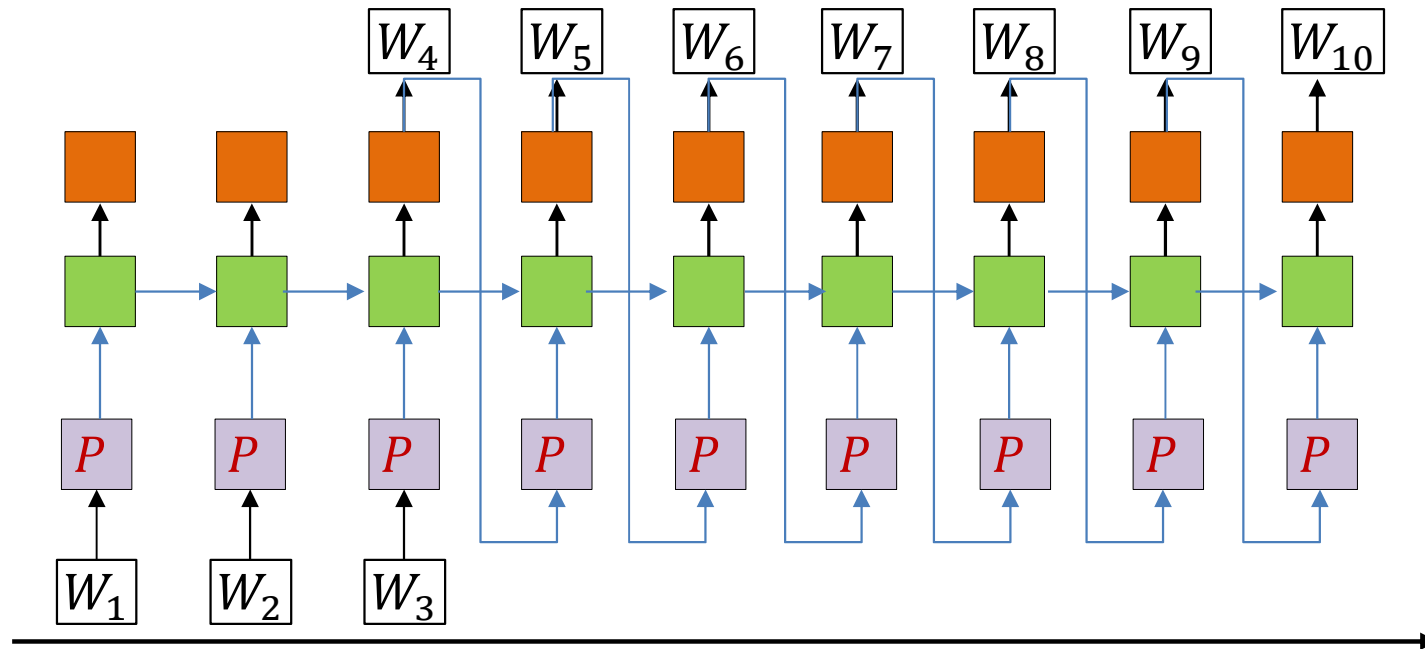
- This is the *end* of a sentence

<sos> four score and eight <eos>

- This is a full sentence

- In situations where the start of sequence is obvious, the <sos> may not be needed, but <eos> is required to terminate sequences
- Sometimes we will use a single symbol to represent both start and end of sentence, e.g just <eos> , or even a separate symbol, e.g. <s>

Generating Language: Synthesis



- Feed the drawn word as the next word in the series
 - And draw the next word by sampling from the output probability distribution
- Continue this process until we draw an `<eos>`
 - Or we decide to terminate generation based on some other criterion

Poll 1

- @



Poll 1

Which of the following is a complete sentence

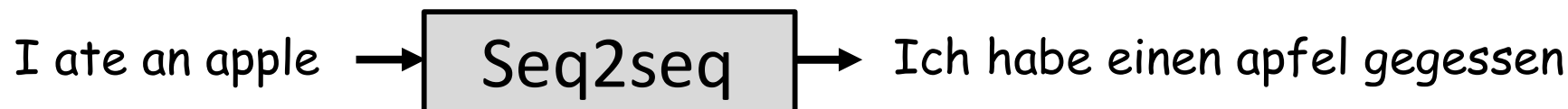
<eos> Hello World <eos>

<sos> Hello World

Hello World <eos>

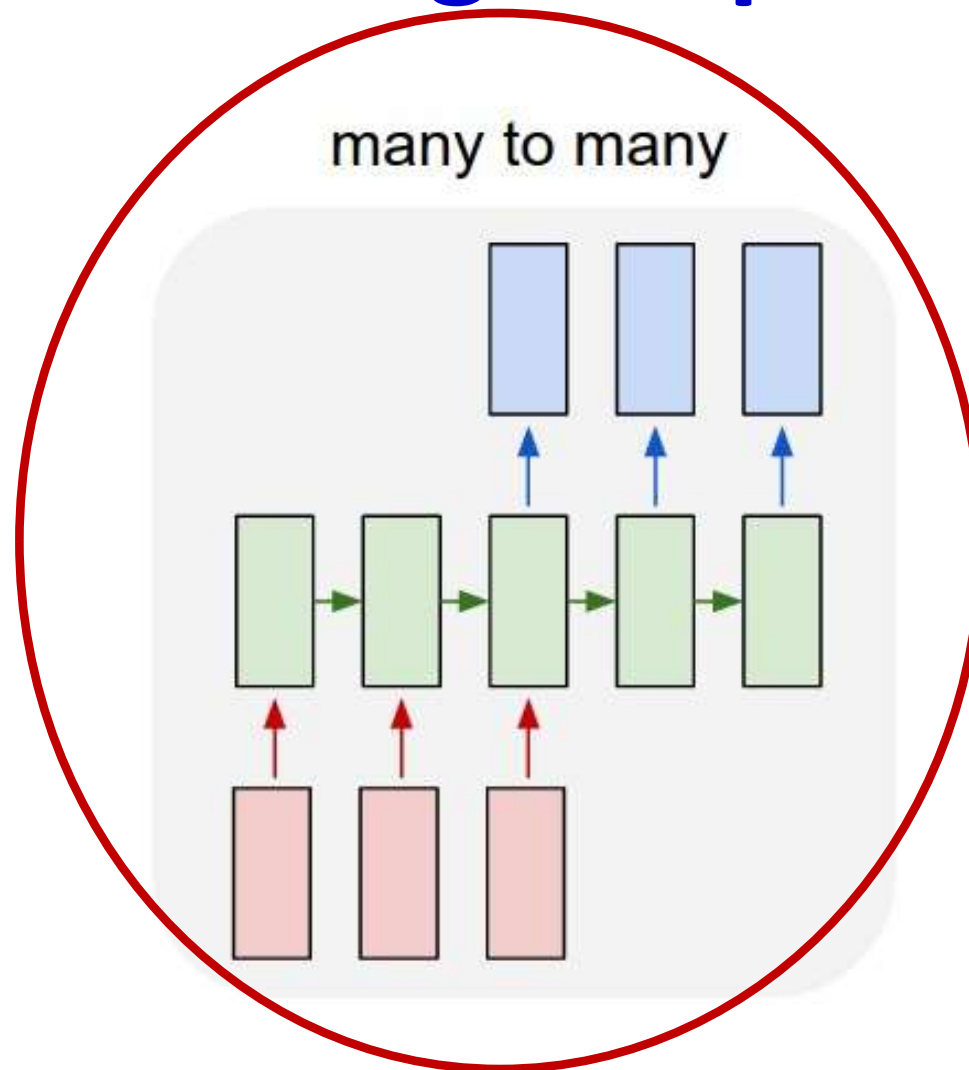
Hello World

Returning our problem



- Problem:
 - A sequence $X_1 \dots X_N$ goes in
 - A different sequence $Y_1 \dots Y_M$ comes out
- No expected synchrony between input and output

Modelling the problem

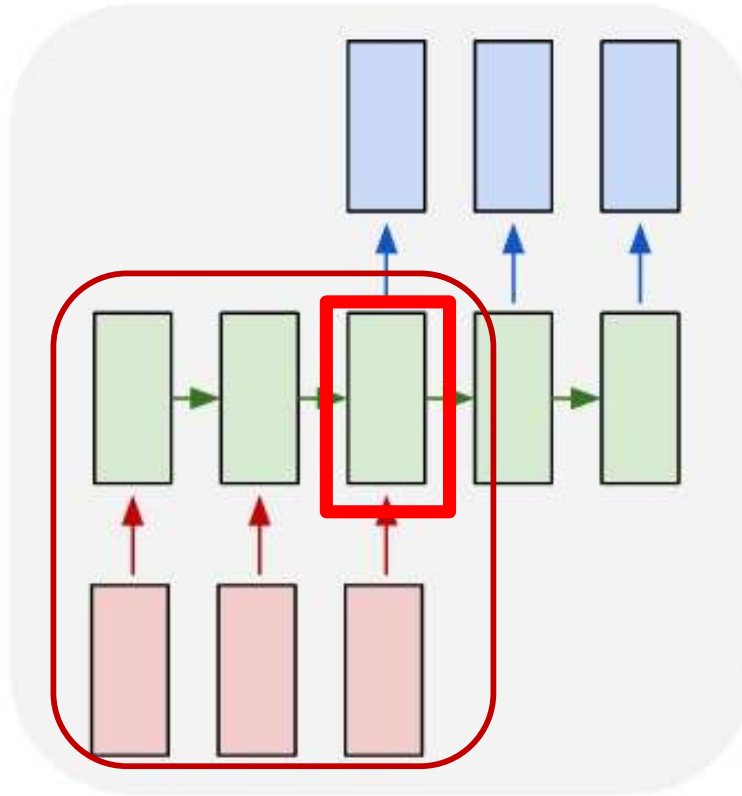


- *Delayed* sequence to sequence

Modelling the problem

many to many

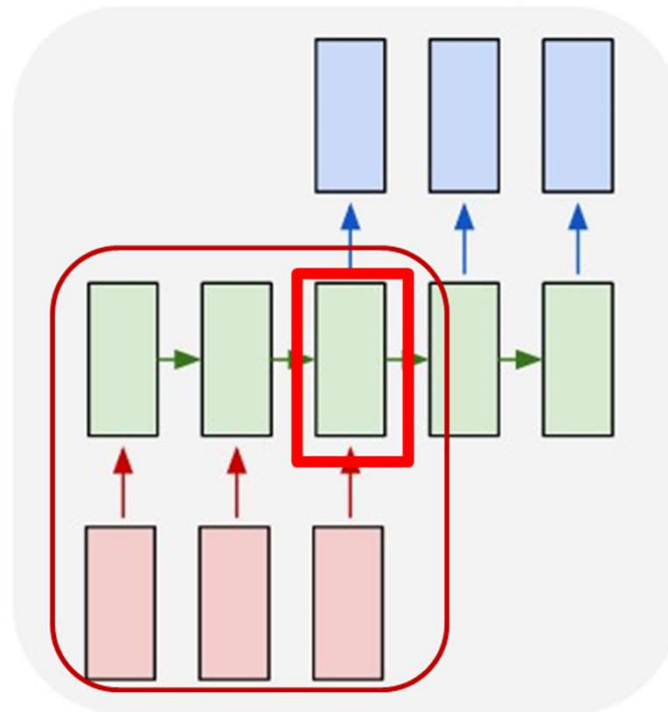
First process the input and generate a hidden representation for it



- *Delayed* sequence to sequence

Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1, l)$  is available for all layers
for  $t = 0:T-1$  # Including both ends of the index
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
 $H = h(T-1)$ 
```

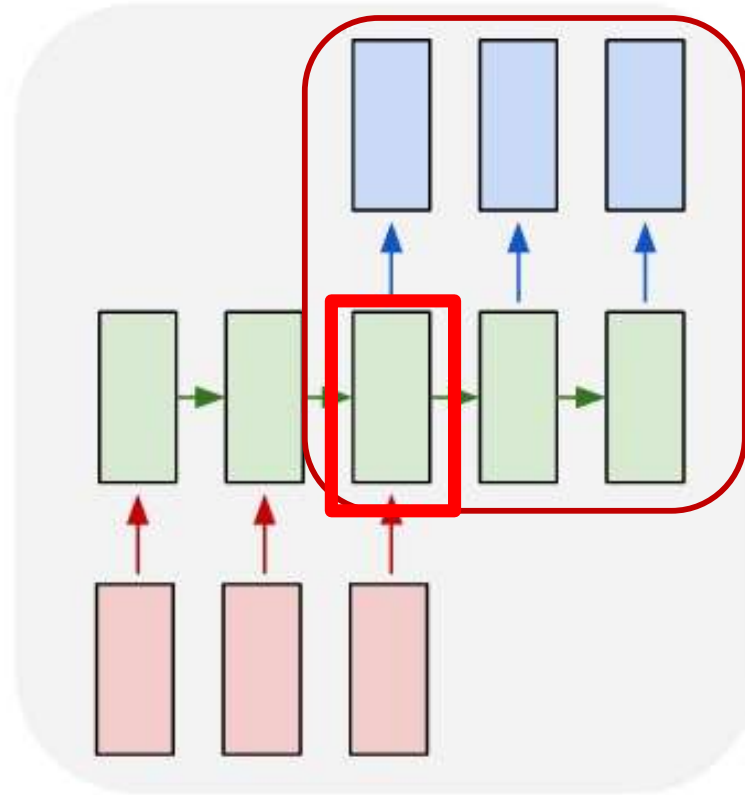


“RNN_input” may be a multi-layer RNN of any kind

Modelling the problem

many to many

First process the input and generate a hidden representation for it



Then use it to generate an output

- *Delayed* sequence to sequence

Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1,1)$  is available for all layers
for  $t = 0:T-1$  # Including both ends of the index
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
 $H = h(T-1)$ 
```

```
# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
 $t = 0$ 
 $h_{\text{out}}(0) = H$ 
do
     $t = t+1$ 
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1))$ 
     $y_{\text{out}}(t) = \text{draw\_word\_from}(y(t))$ 
until  $y_{\text{out}}(t) == \text{<eos>}$ 
```

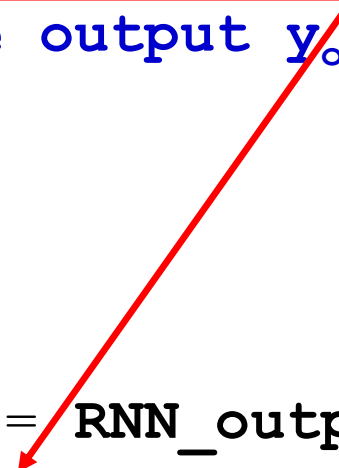
Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1,1)$  is available for all layers
for  $t = 0:T-1$  # Including both ends of the index
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
 $H = h(T-1)$ 
```

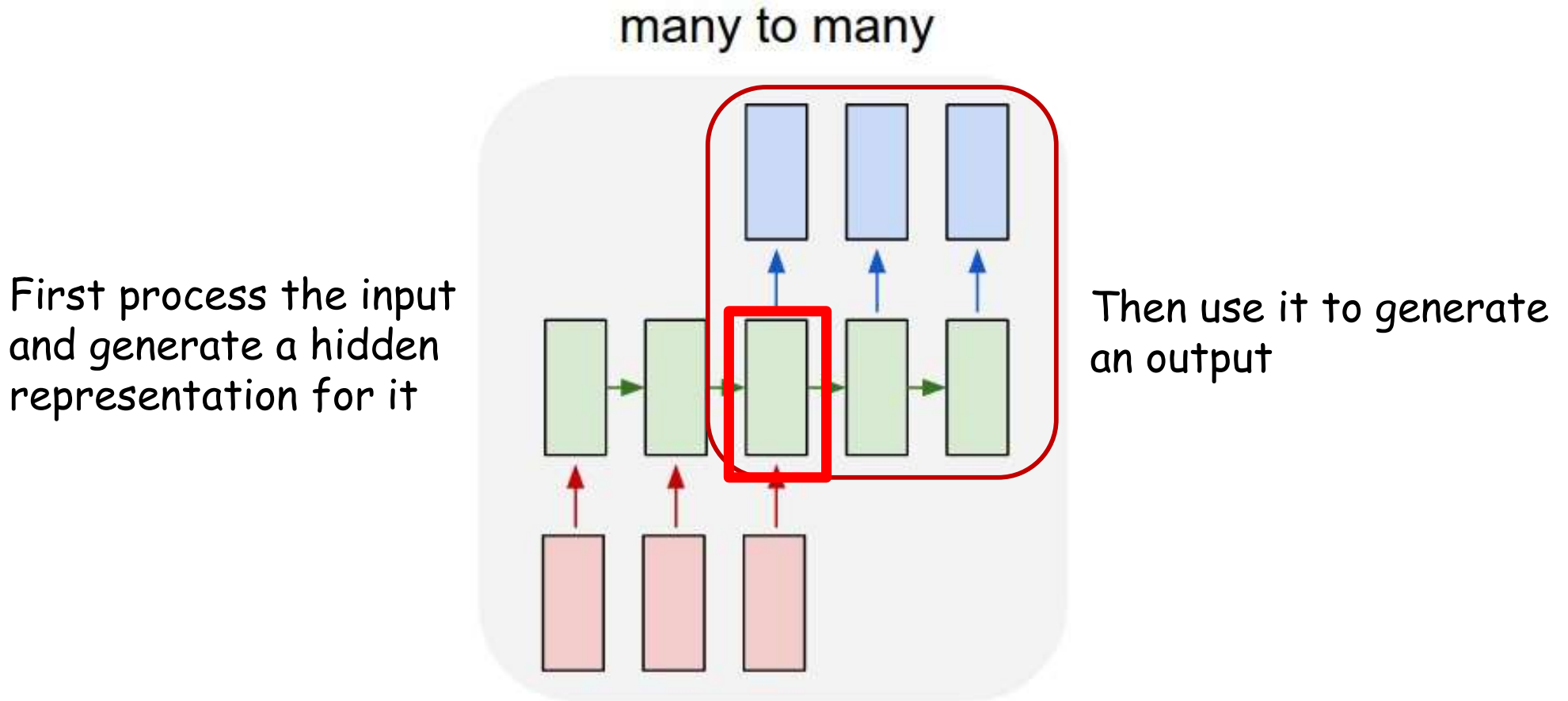
The output at each time is a probability distribution over symbols.

We draw a word from this distribution

```
# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
 $t = 0$ 
 $h_{\text{out}}(0) = H$ 
do
     $t = t+1$ 
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1))$ 
     $y_{\text{out}}(t) = \text{draw\_word\_from}(y(t))$ 
until  $y_{\text{out}}(t) == \text{<eos>}$ 
```



Modelling the problem



- *Problem:* Each word that is output depends only on current hidden state, and not on previous outputs

Pseudocode

Changing this output at time t does not affect the output at $t+1$

E.g. If we have drawn "It was a" vs "It was an", the probability that the next word is "dark" remains the same (dark must ideally not follow "an")

This is because the output at time t does not influence the computation at $t+1$

Now generate the output $y(1)$ $y(2)$

$t = 0$

$h_{\text{out}}(0) = H$

do

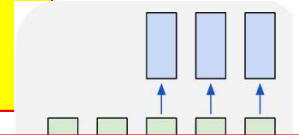
$t = t+1$

$[y(t), h_{\text{out}}(t)] = \text{RNN_output_step}(h_{\text{out}}(t-1))$

$y_{\text{out}}(t) = \text{draw_word_from}(y(t))$

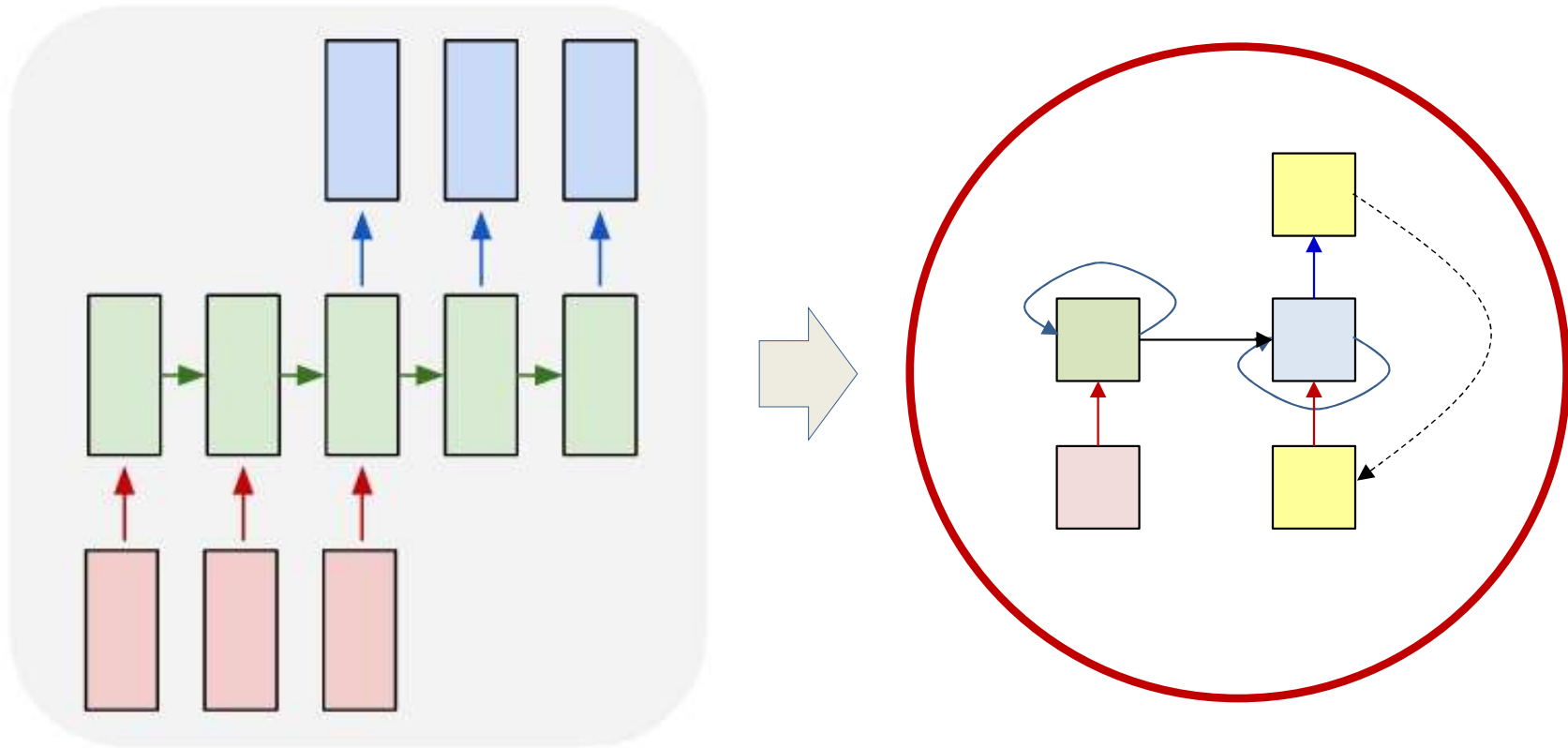
until $y_{\text{out}}(t) == \text{<eos>}$

The RNN recursion only considers the hidden state $h(t-1)$ from the previous time and not the actual output word $y_{\text{out}}(t-1)$



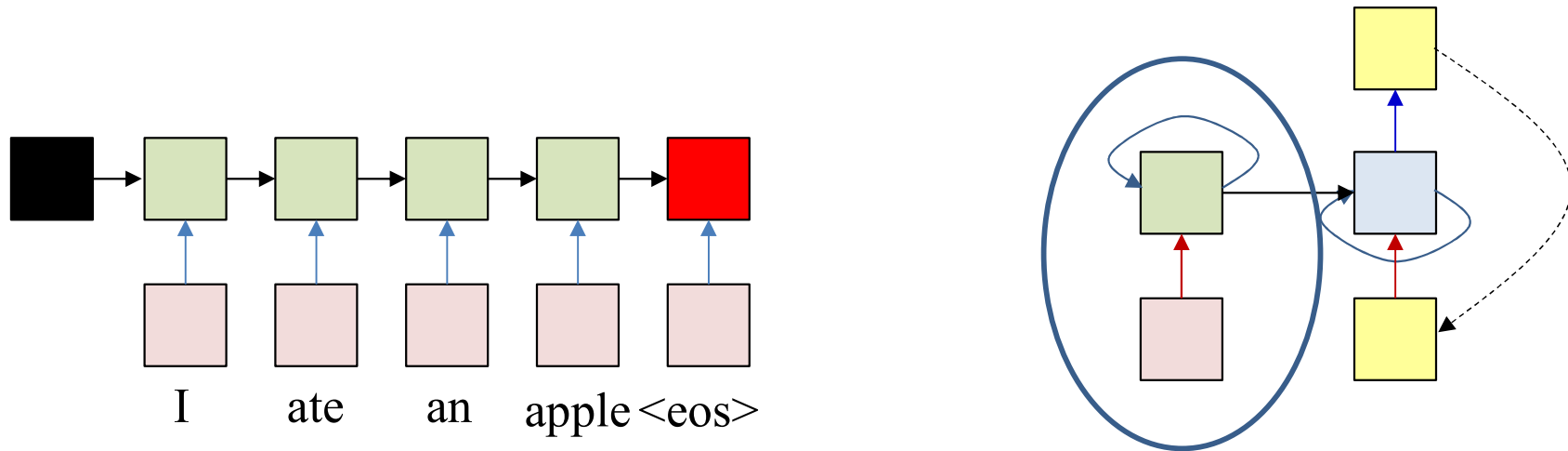
Modelling the problem

many to many



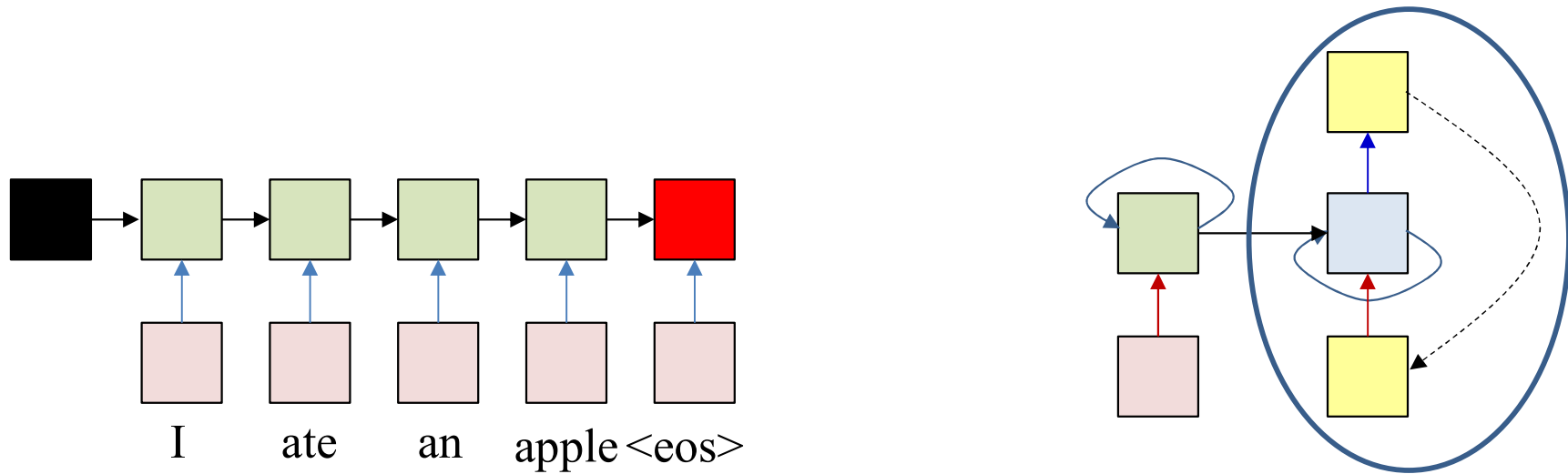
- *Delayed* sequence to sequence
 - Delayed *self-referencing* sequence-to-sequence

The “simple” translation model



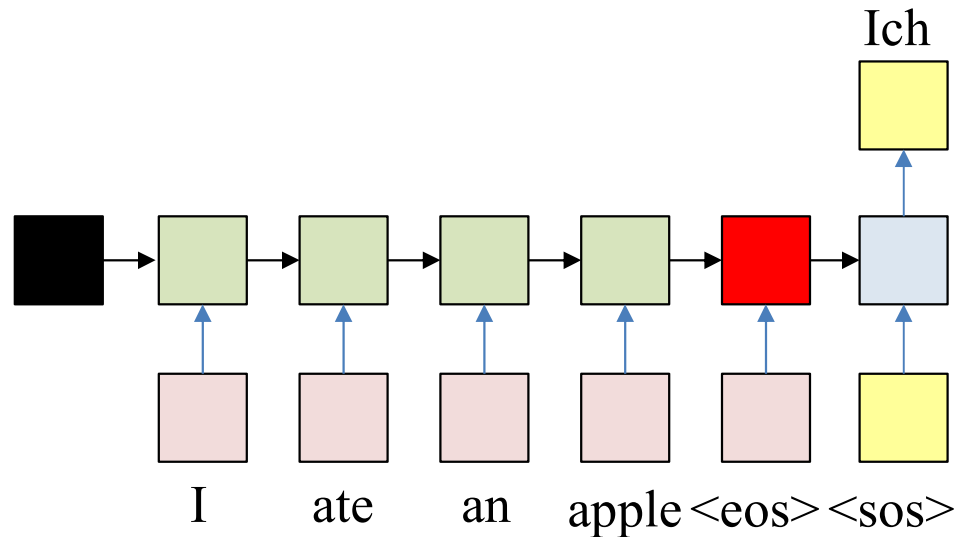
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence

The “simple” translation model



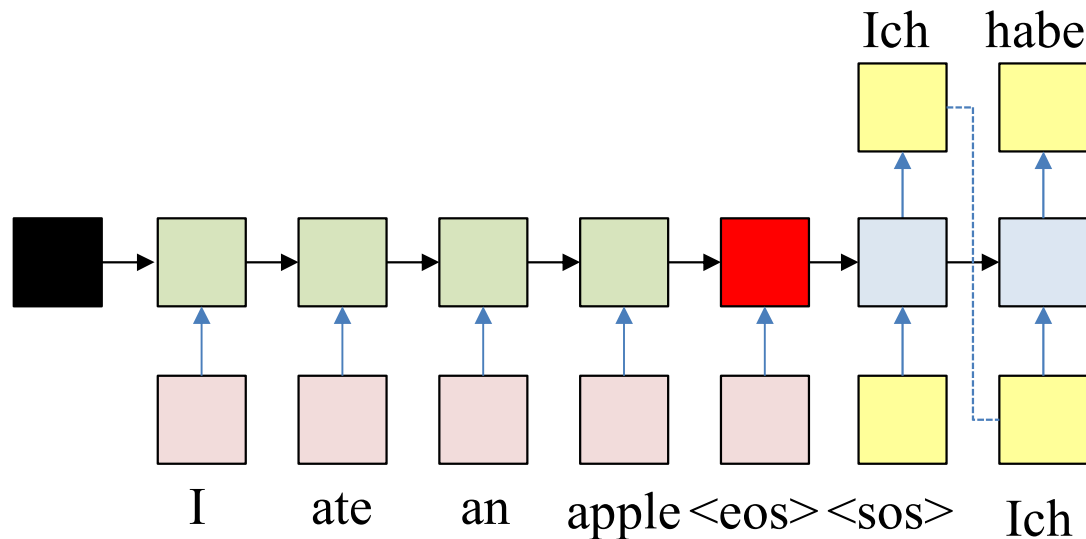
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence
- Subsequently a *second* RNN uses the hidden activation as initial state, and <sos> as initial symbol, to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

The “simple” translation model



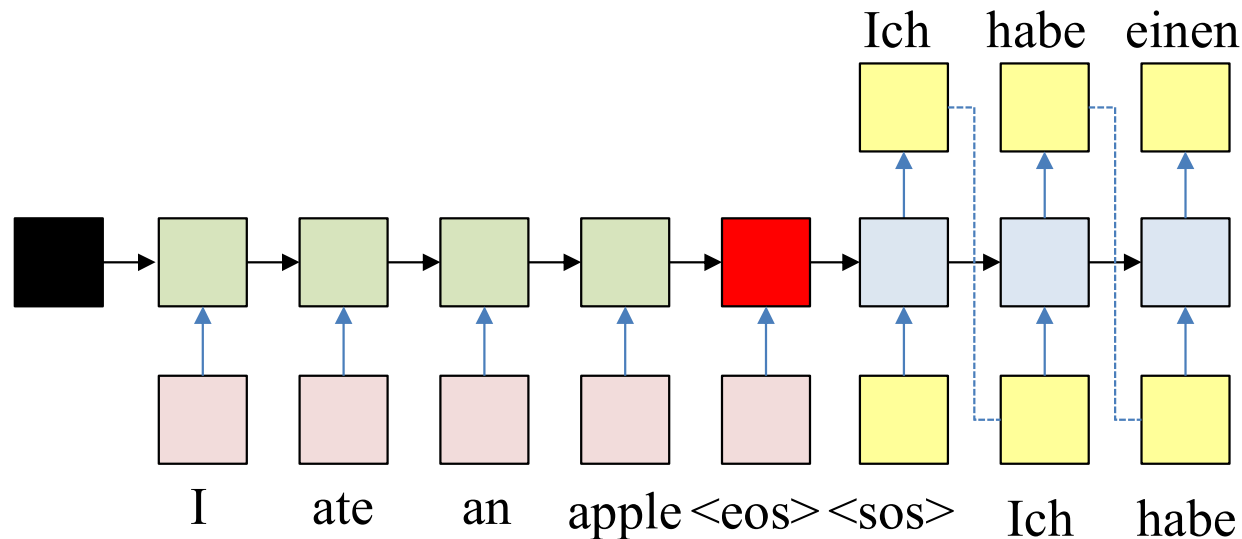
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence
- Subsequently a *second* RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

The “simple” translation model



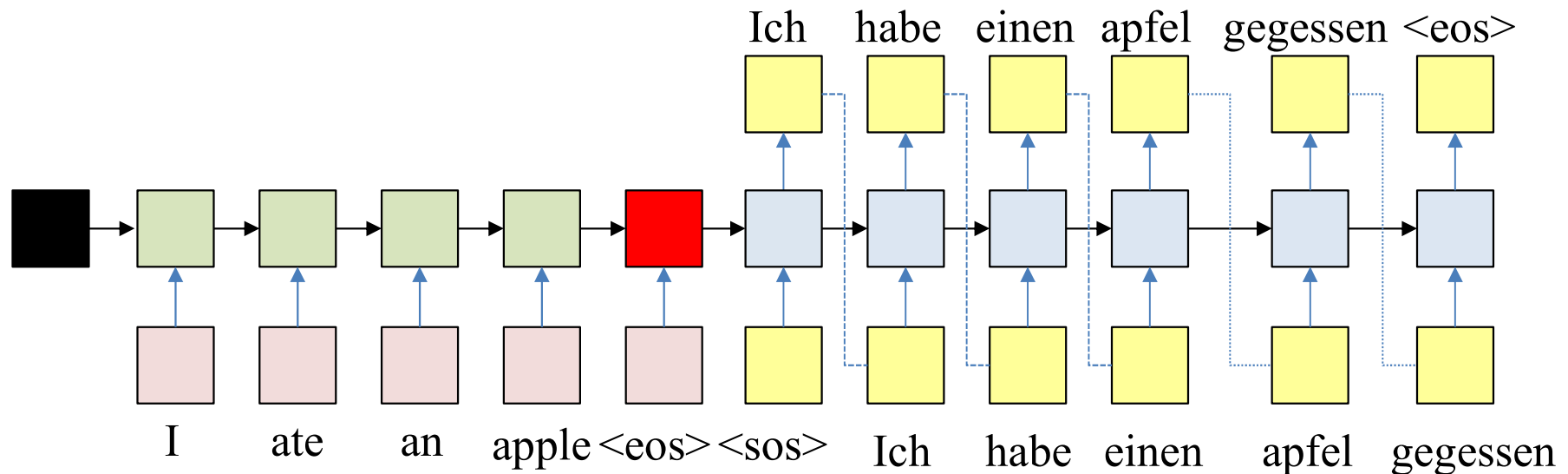
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence
- Subsequently a *second* RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

The “simple” translation model



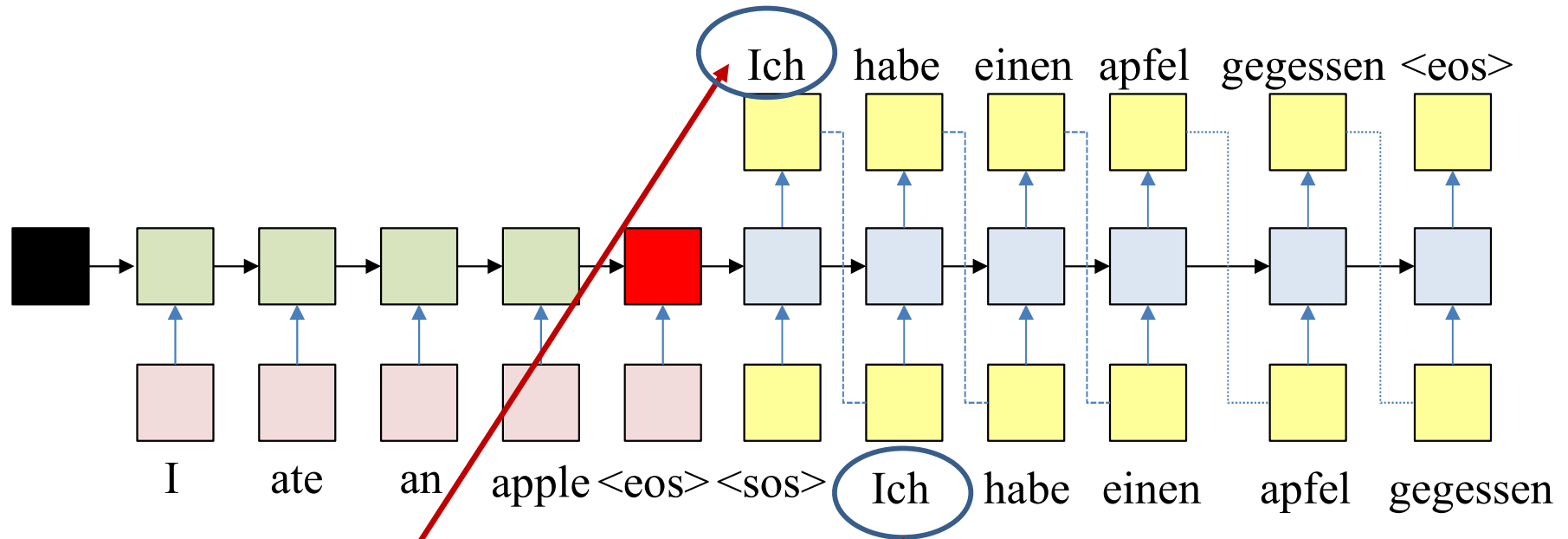
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence
- Subsequently a *second* RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

The “simple” translation model



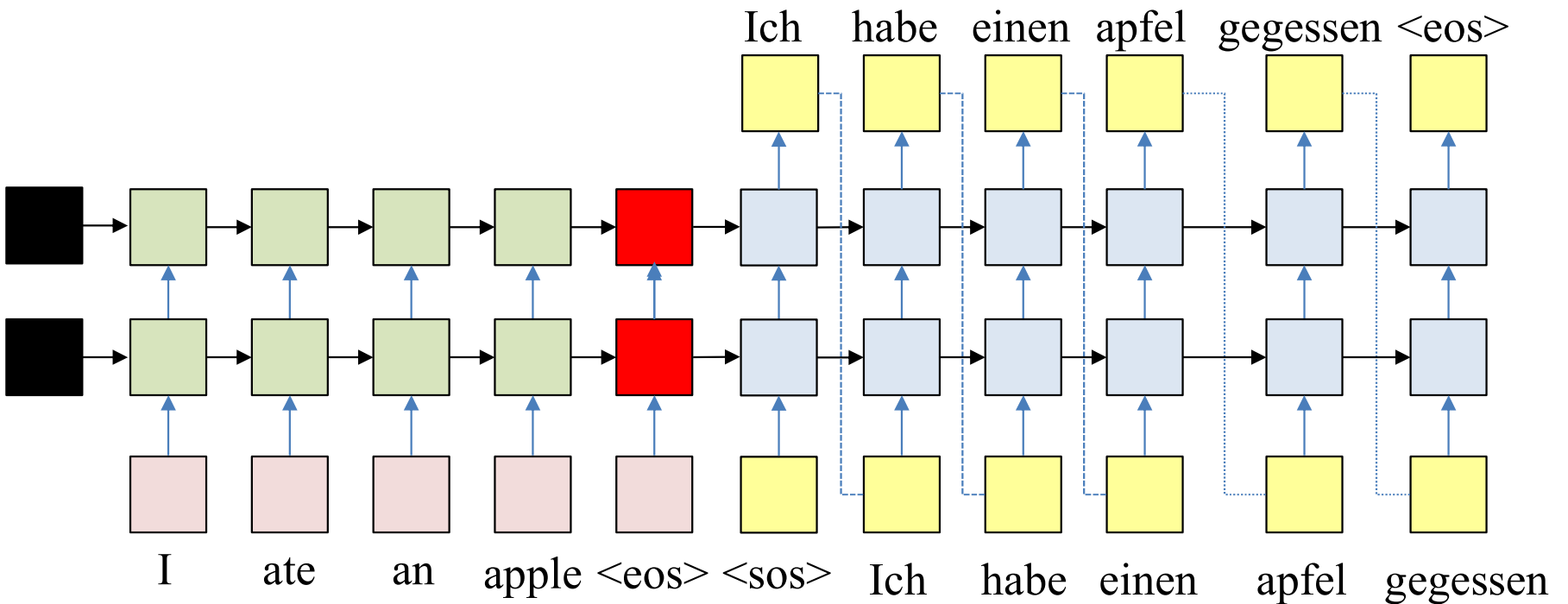
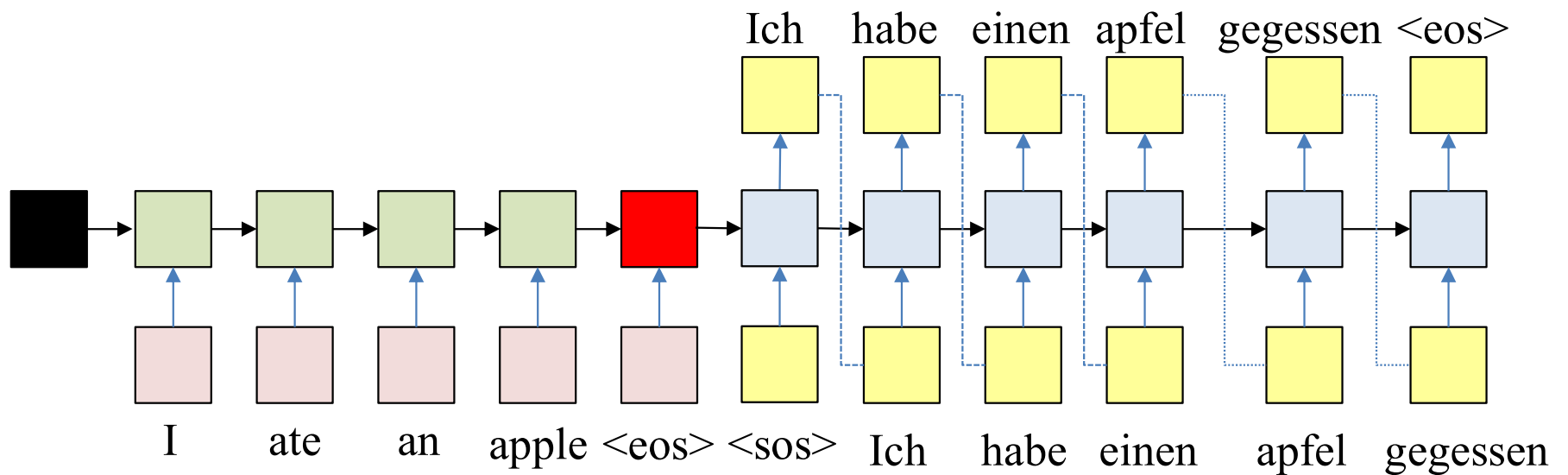
- The input sequence feeds into a recurrent structure
- The input sequence is terminated by an explicit <eos> symbol
 - The hidden activation at the <eos> “stores” all information about the sentence
- Subsequently a *second* RNN uses the hidden activation as initial state to produce a sequence of outputs
 - The output at each time becomes the input at the next time
 - Output production continues until an <eos> is produced

The “simple” translation model



Note that drawing a different word here

Would result in a different word being input here, and as a result the output here and subsequent outputs would all change



- We will illustrate with a single hidden layer, but the discussion generalizes to more layers

Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1,1)$  is available for all layers
t = 0
do
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
until  $x(t) == \text{"<eos>"}$ 
H =  $h(T-1)$ 

# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
t = 0
 $h_{\text{out}}(0) = H$ 

# Note: begins with a "start of sentence" symbol
#       <sos> and <eos> may be identical
 $y_{\text{out}}(0) = \text{<sos>}$ 
do
    t = t+1
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1), y_{\text{out}}(t-1))$ 
     $y_{\text{out}}(t) = \text{draw\_word\_from}(y(t))$ 
until  $y_{\text{out}}(t) == \text{<eos>}$ 
```

Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1,1)$  is available for all layers
t = 0
do
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
until  $x(t) == \text{"<eos>"}$ 
H =  $h(T-1)$ 
```

```
# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
```

```
t = 0
 $h_{\text{out}}(0) = H$ 
```

```
# Note: begins with a "start of sentence" symbol
```

```
# <sos> and <eos> may be identical
```

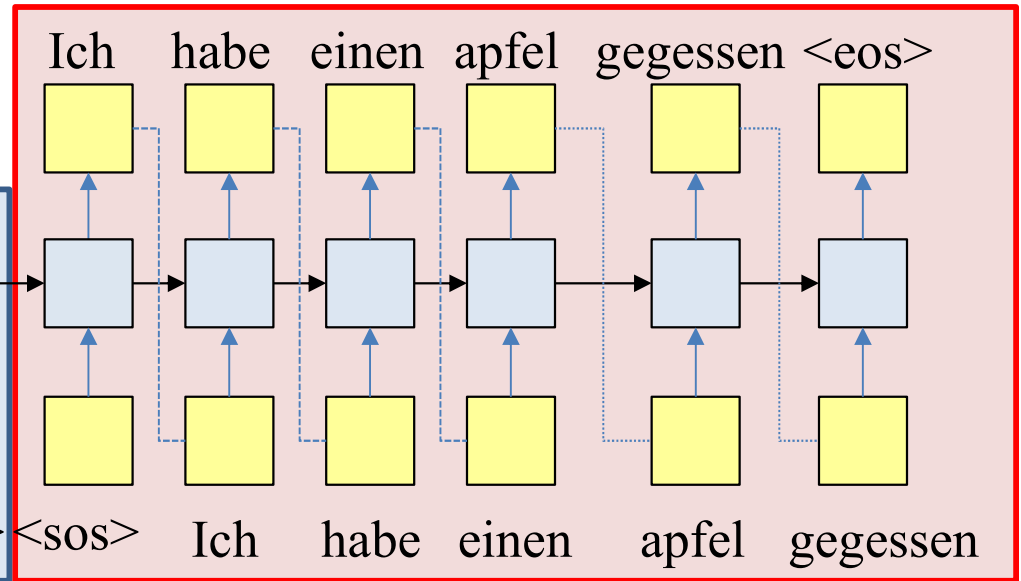
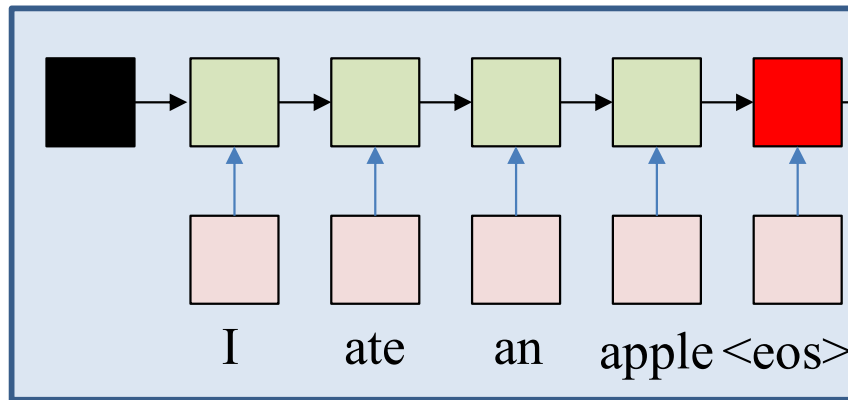
```
 $y_{\text{out}}(0) = \text{"<sos>"}$ 
```

```
do
    t = t+1
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1), y_{\text{out}}(t-1))$ 
     $y_{\text{out}}(t) = \text{draw\_word\_from}(y(t))$ 
until  $y_{\text{out}}(t) == \text{"<eos>"}$ 
```

Drawing a different word at t will change the next output since $y_{\text{out}}(t)$ is fed back as input

The “simple” translation model

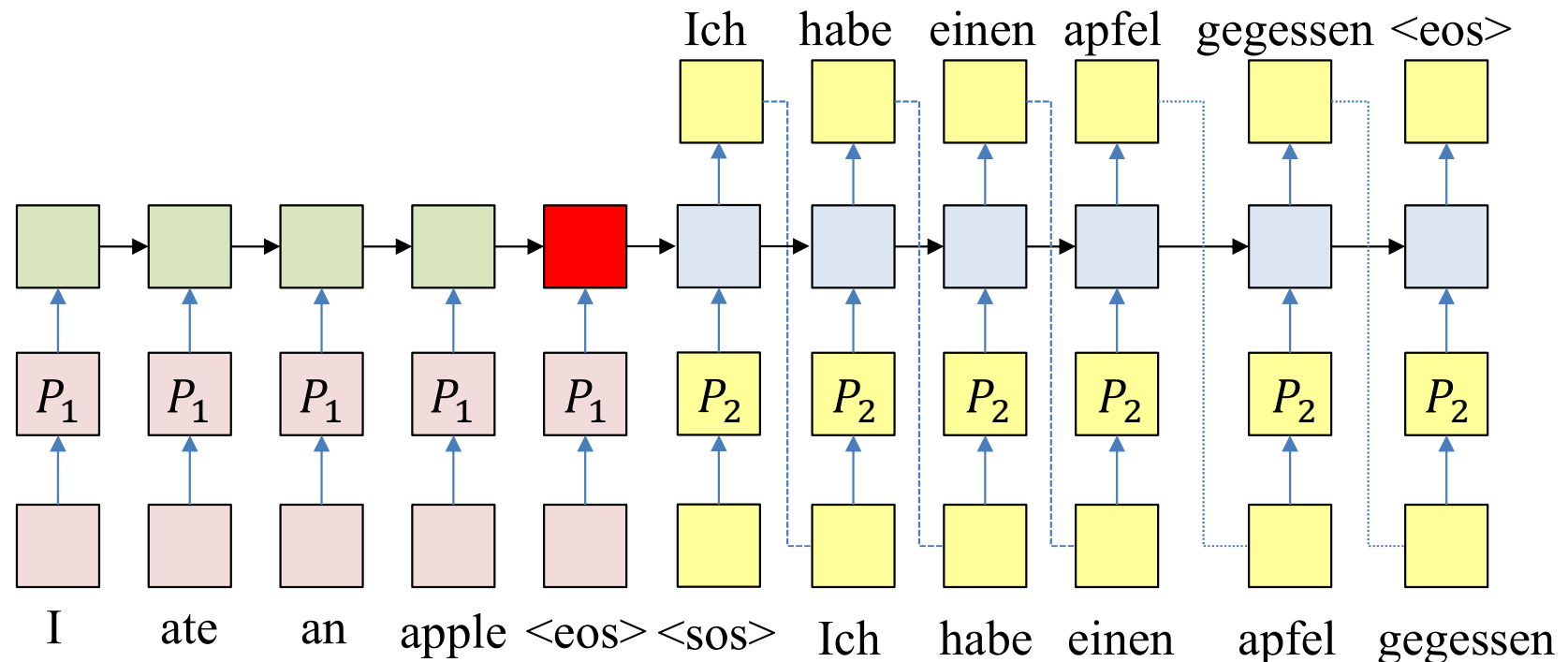
ENCODER



DECODER

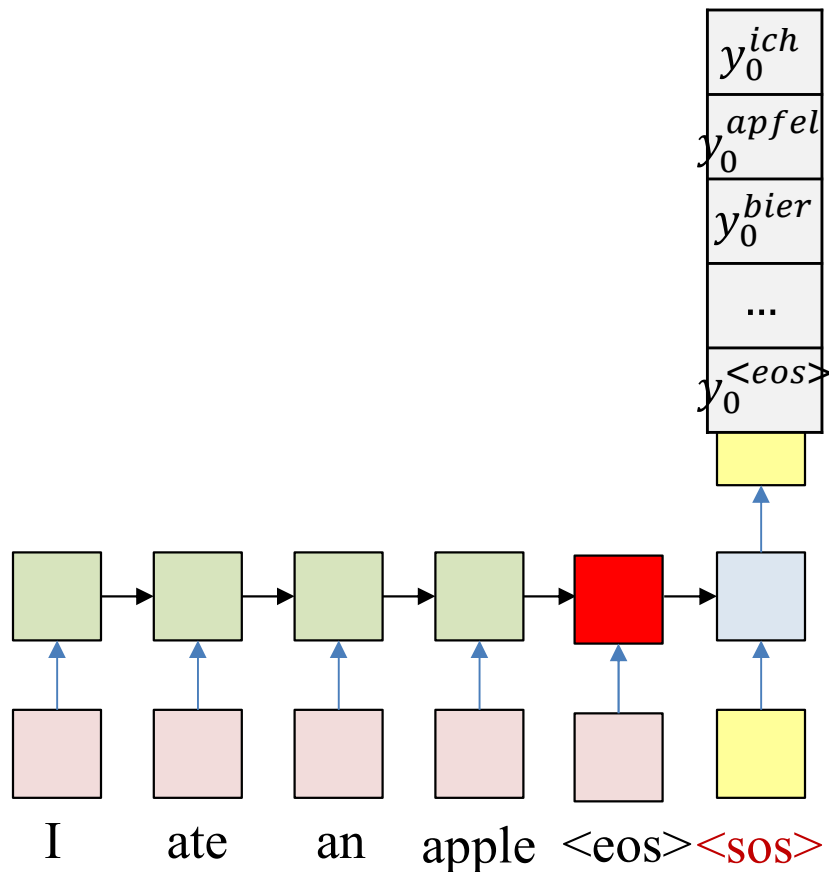
- The recurrent structure that extracts the hidden representation from the input sequence is the *encoder*
- The recurrent structure that utilizes this representation to produce the output sequence is the *decoder*

The “simple” translation model



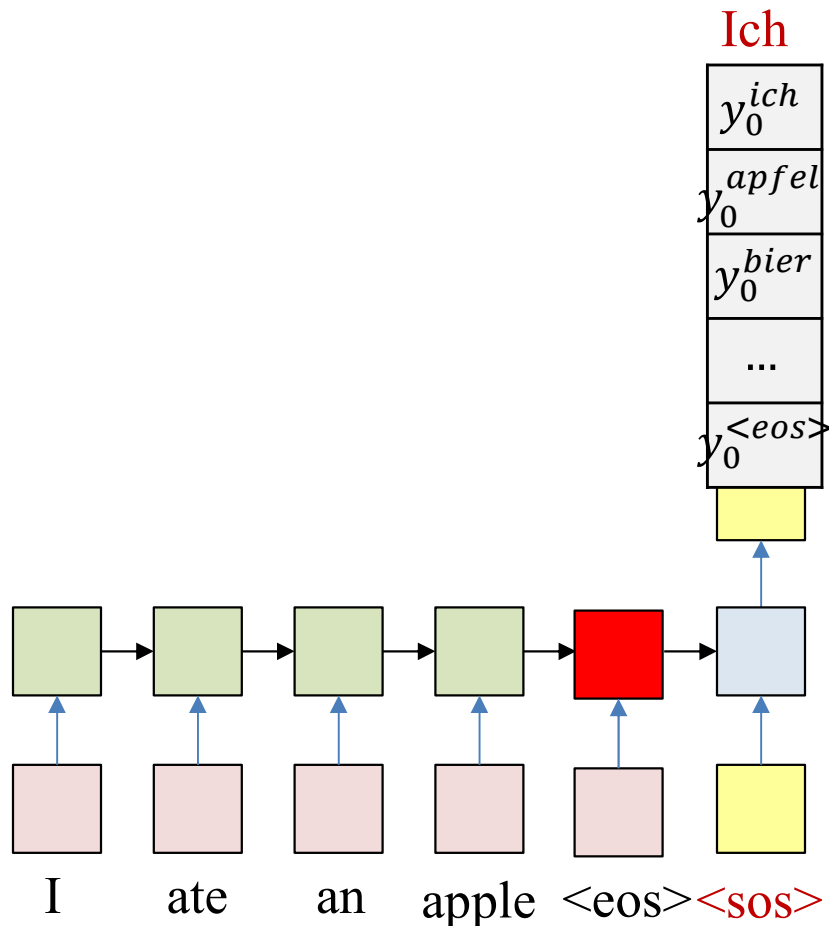
- A more detailed look: The one-hot word representations may be compressed via embeddings
 - Embeddings will be learned along with the rest of the net
 - In the following slides we will not represent the projection matrices

What the network actually produces



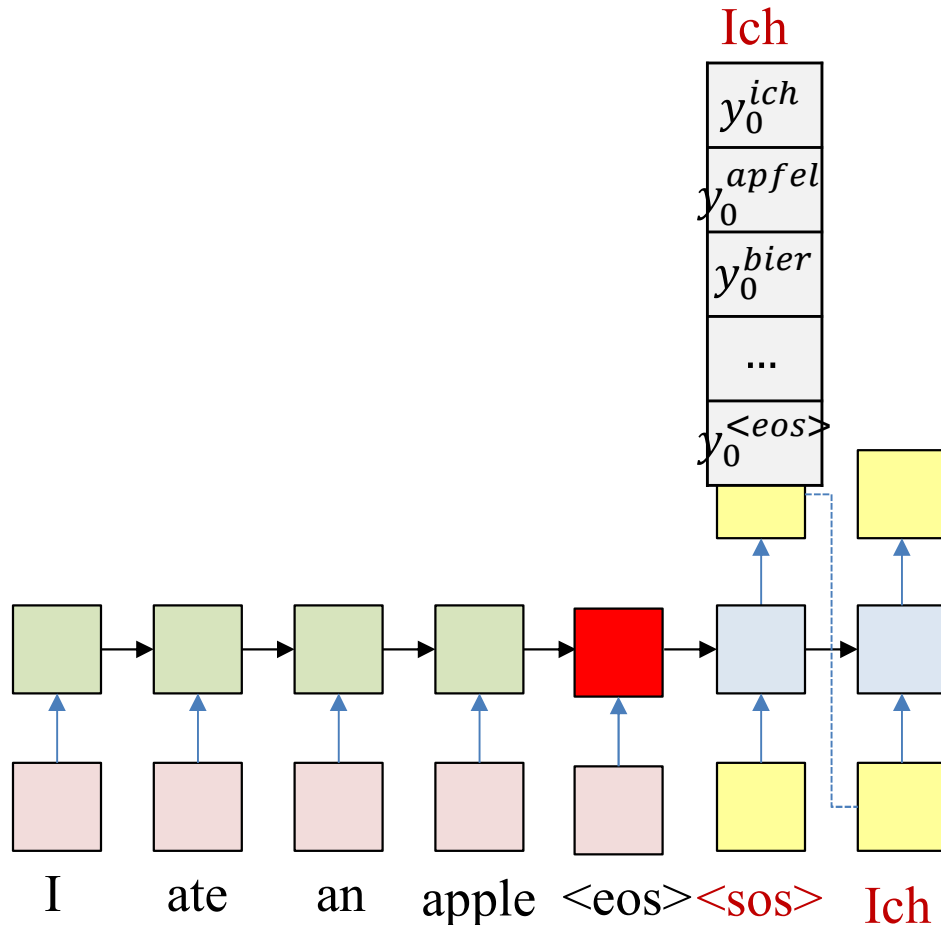
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



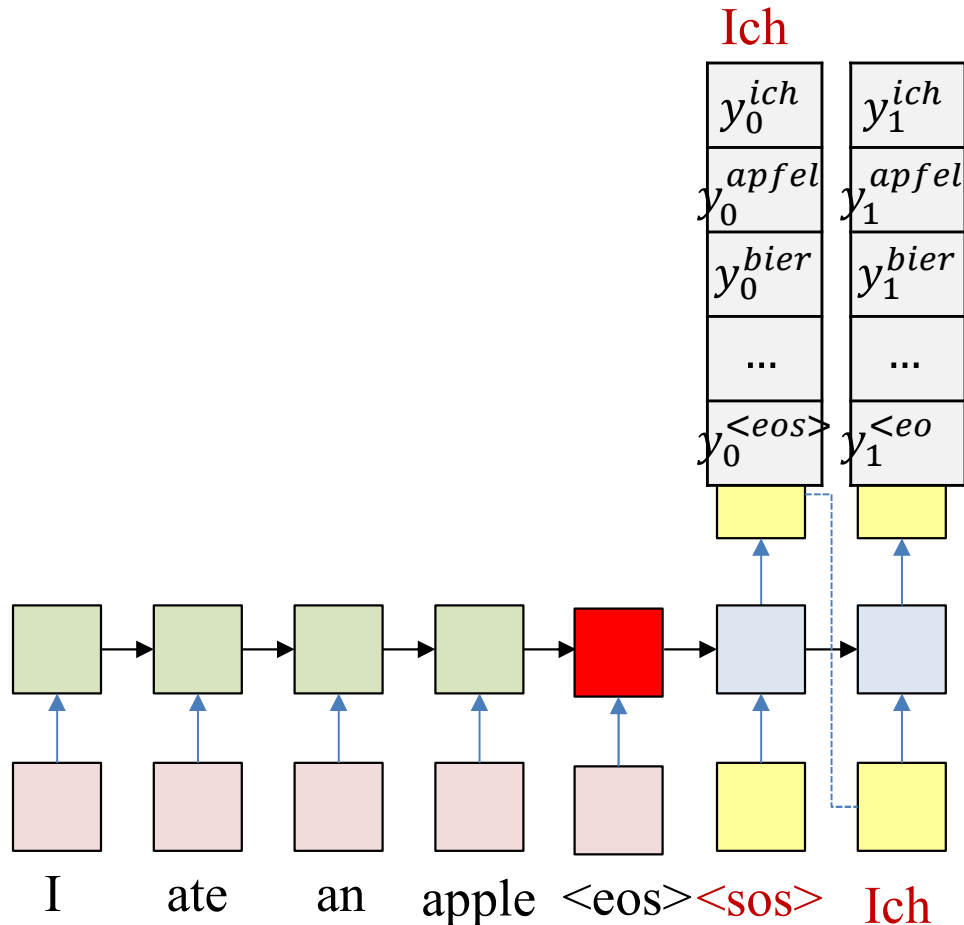
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



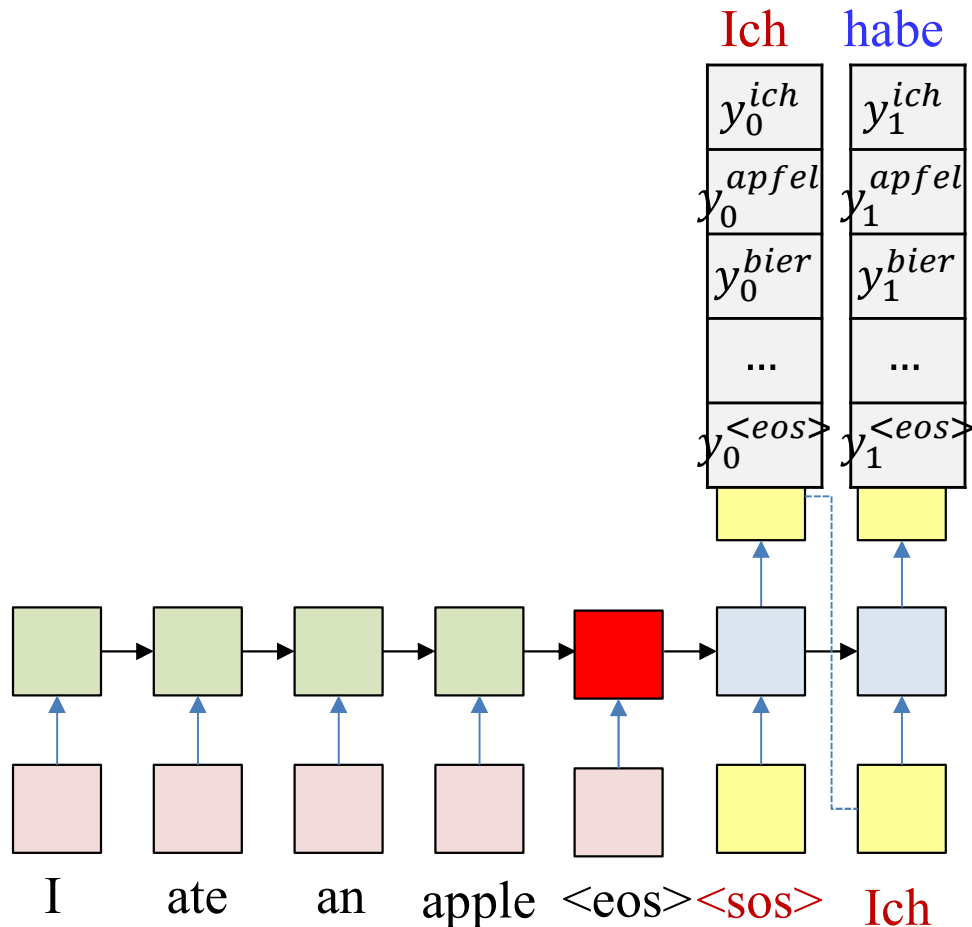
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



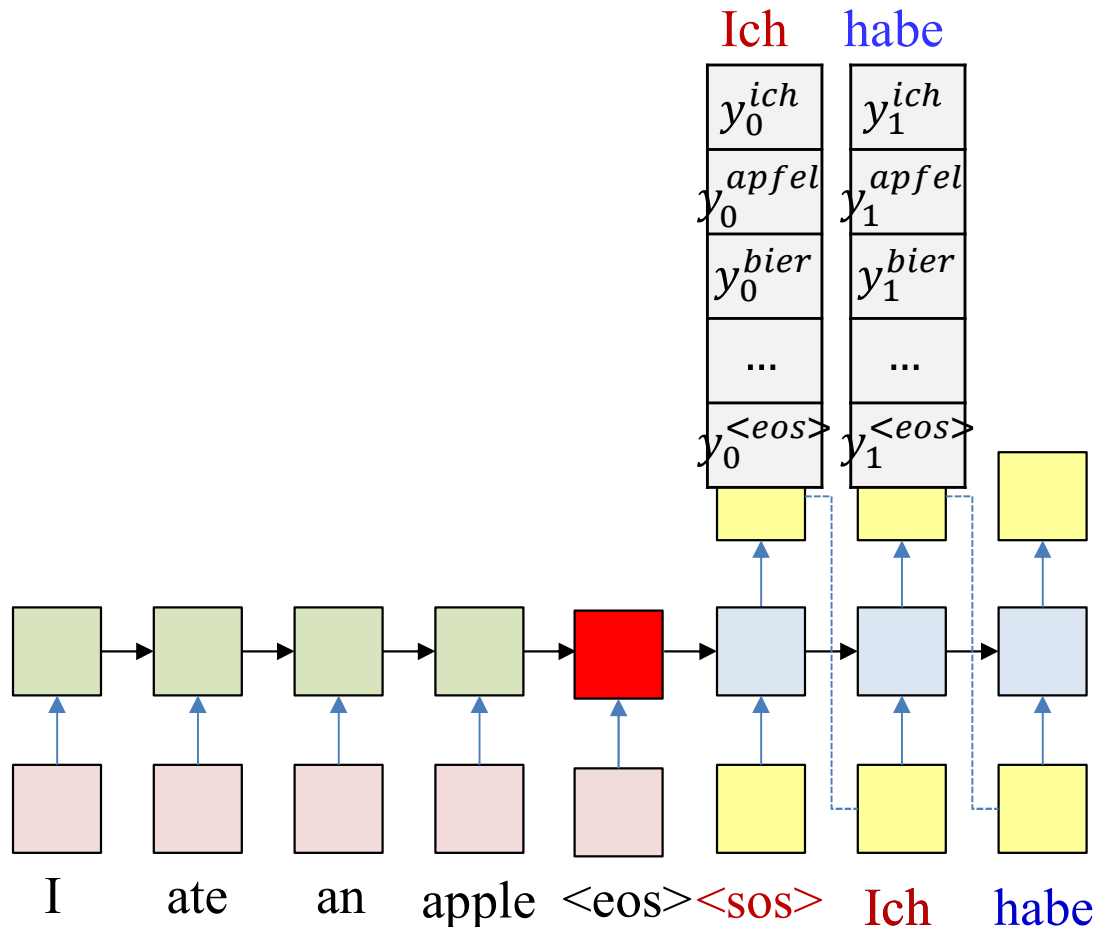
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



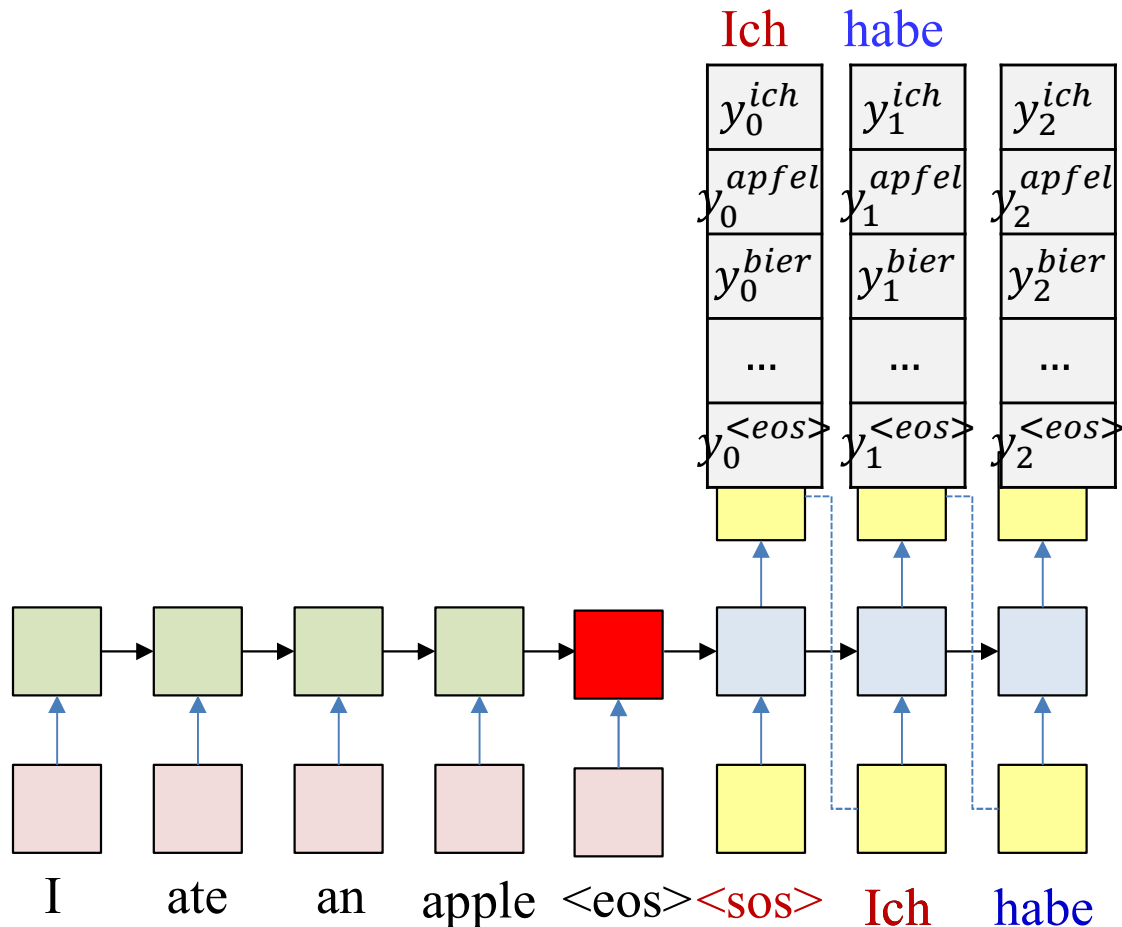
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



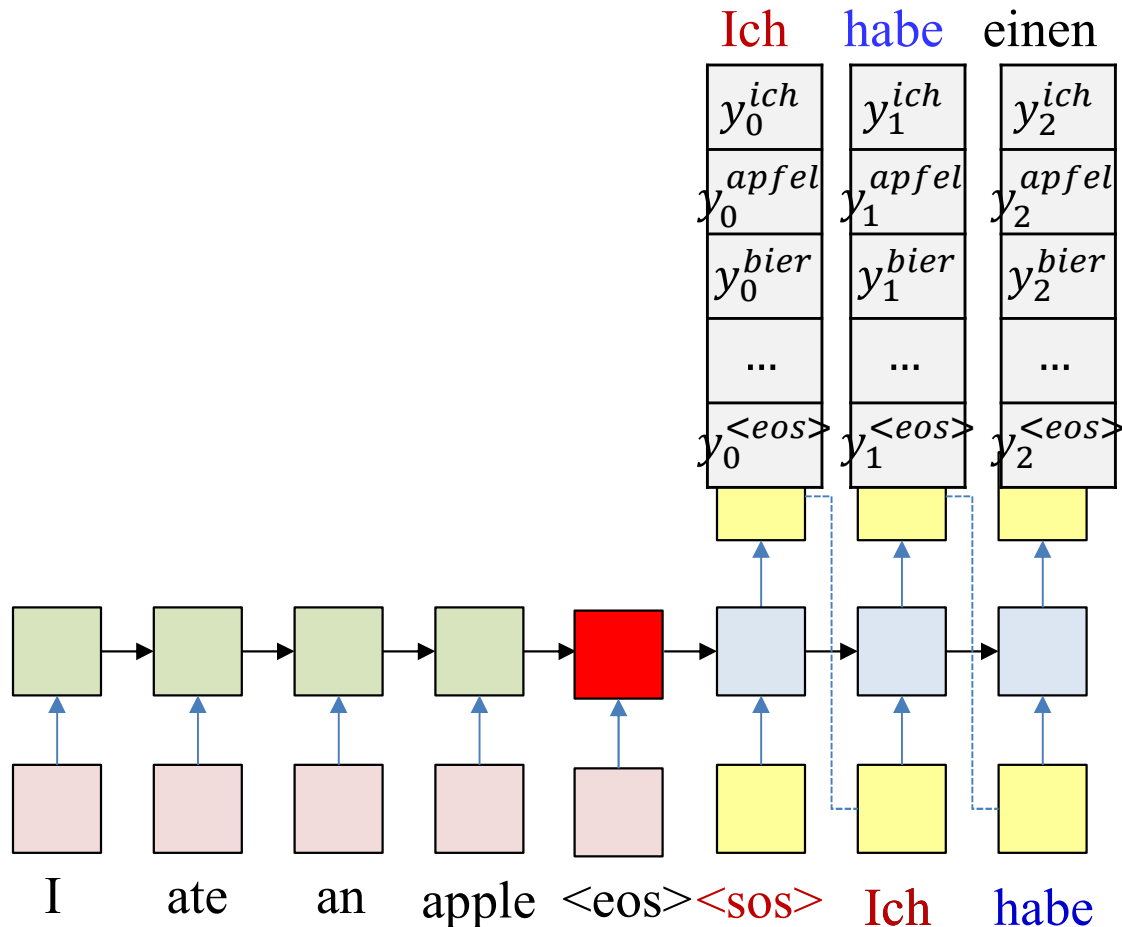
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



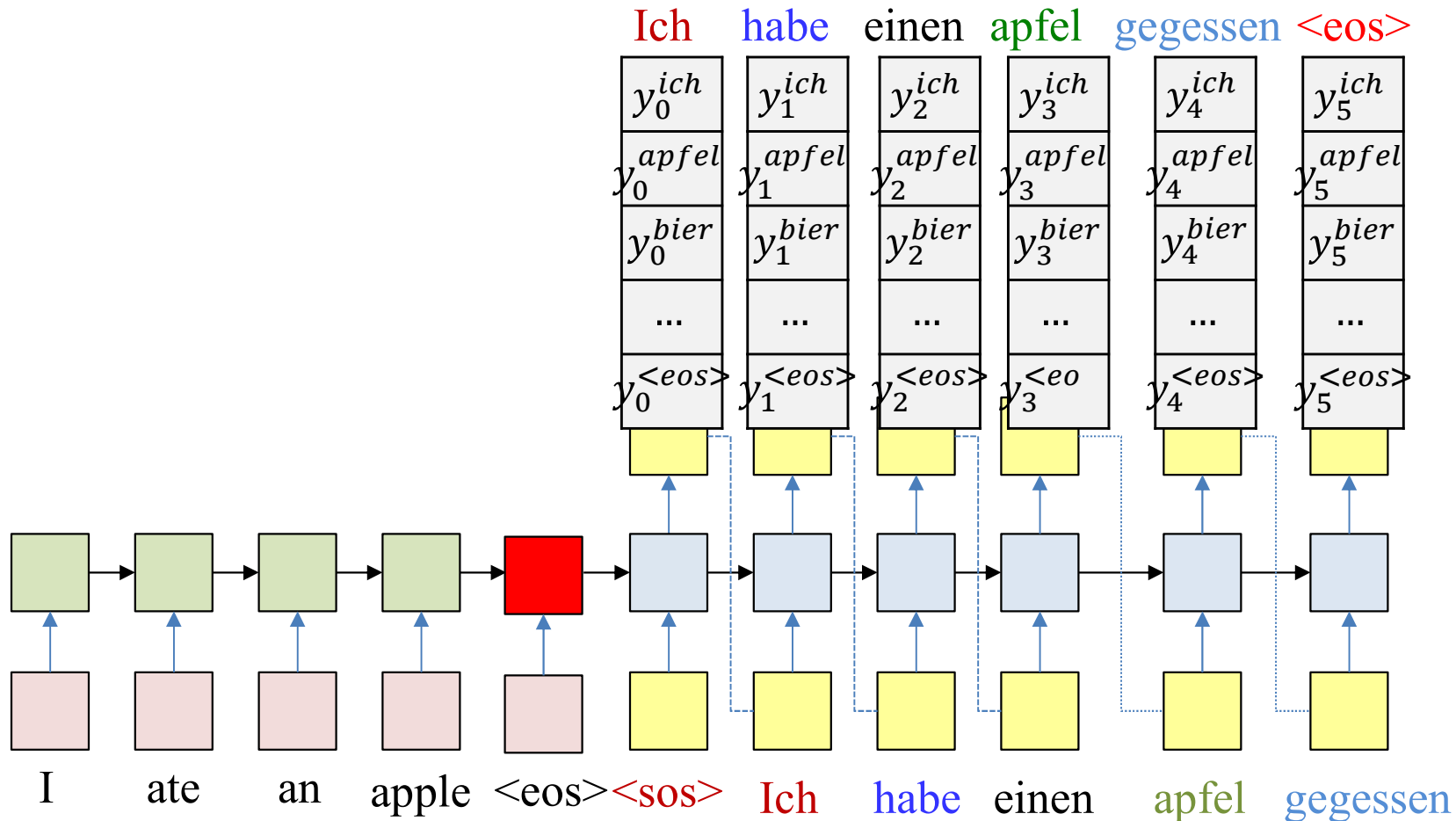
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



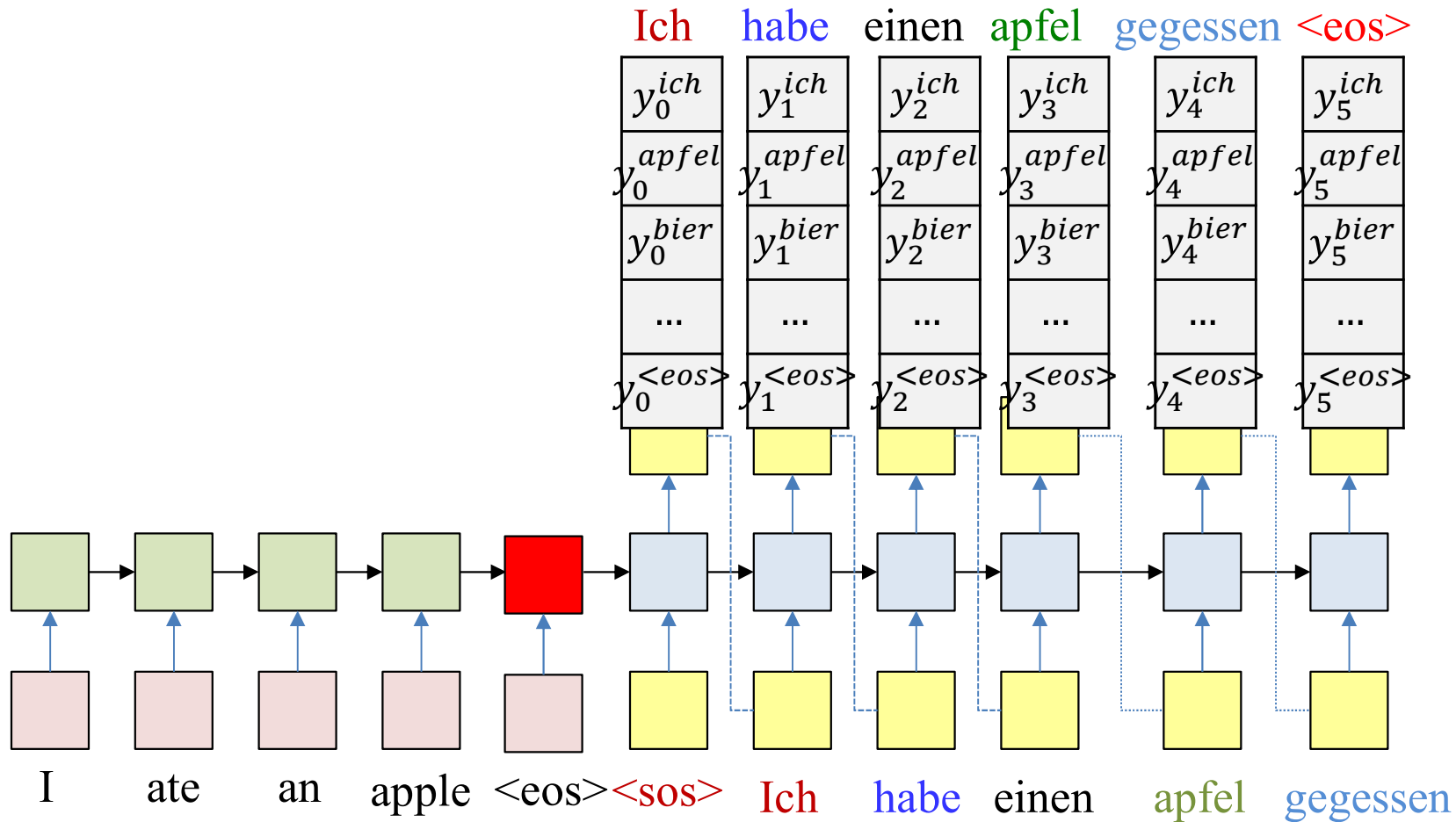
- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

What the network actually produces



- At each time k the network actually produces a probability distribution over the output vocabulary
 - $y_k^w = P(O_k = w | O_{k-1}, \dots, O_1, I_1, \dots, I_N)$
 - The probability given the entire input sequence I_1, \dots, I_N and the partial output sequence O_1, \dots, O_{k-1} until k
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time

Generating an output from the net



- At each time the network produces a probability distribution over words, given the entire input and entire output sequence so far
- At each time a word is *drawn* from the output distribution
- The drawn word is provided as input to the next time
- The process continues until an <eos> is generated

Pseudocode

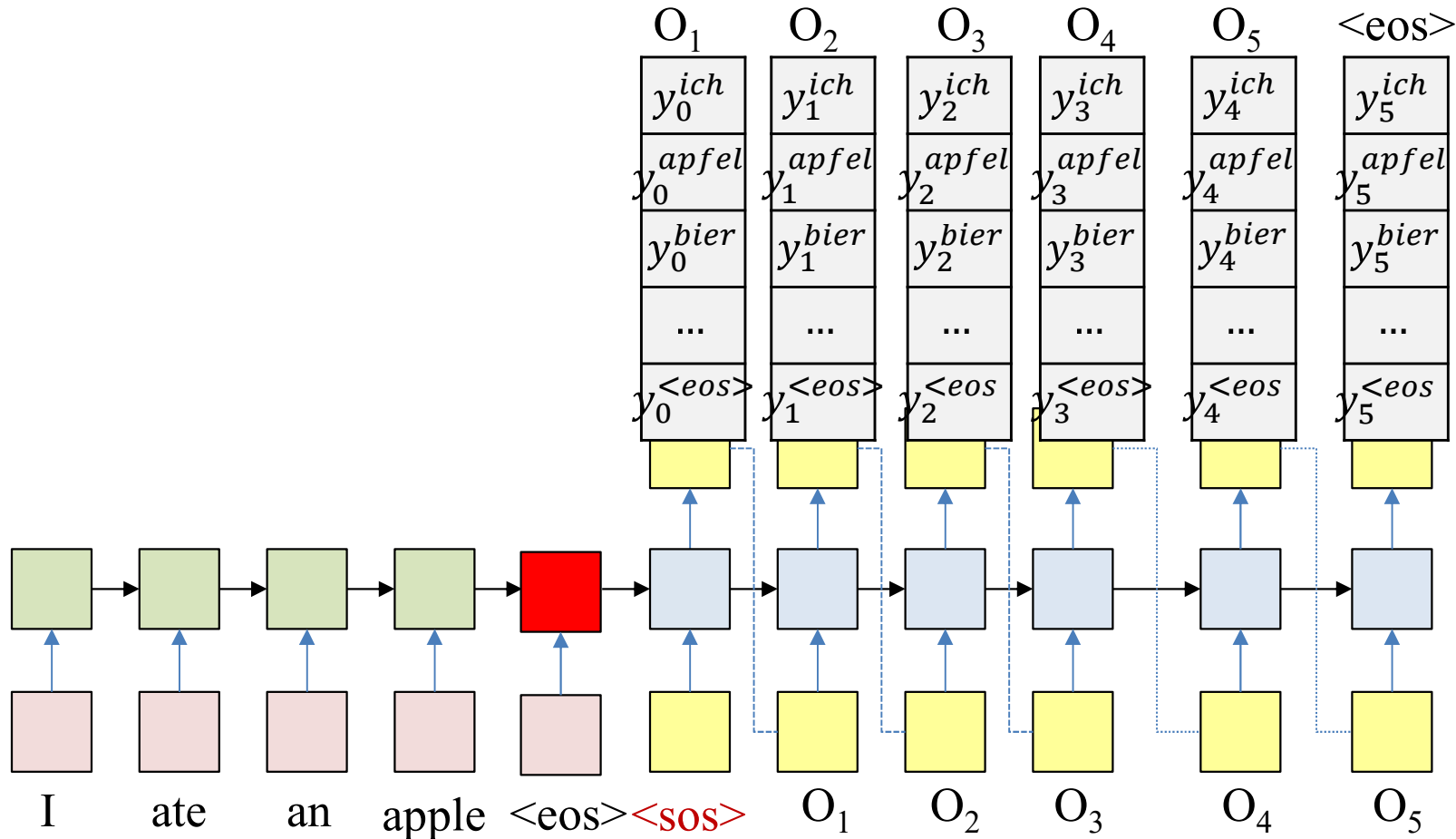
```
# First run the inputs through the network
# Assuming  $h(-1, l)$  is available for all layers
t = 0
do
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
until  $x(t) == \text{"<eos>"}$ 
H =  $h(T-1)$ 

# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
t = 0
 $h_{\text{out}}(0) = H$ 

# Note: begins with a "start of sentence" symbol
#      <sos> and <eos> may be identical
 $y_{\text{out}}(0) = \text{"<sos>"}$ 
do
    t = t+1
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1), y_{\text{out}}(t-1))$ 
     $y_{\text{out}}(t) = \text{draw\_word\_from}(y(t))$ 
until  $y_{\text{out}}(t) == \text{"<eos>"}$ 
```

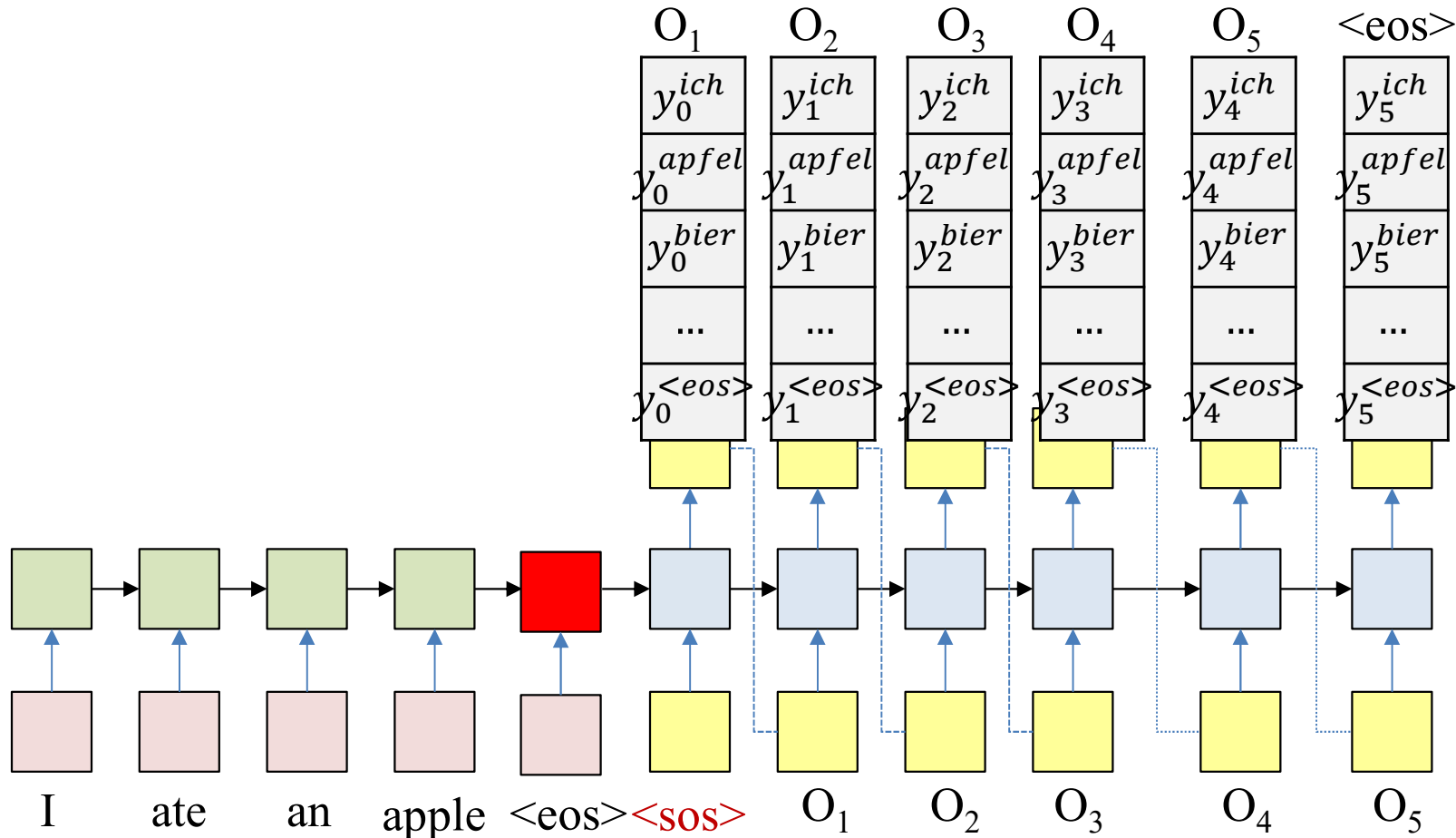
What is this magic operation?

The probability of the output



$$\begin{aligned}
 &P(O_1, \dots, O_L | I_1, \dots, I_N) \\
 &= P(O_1 | I_1, \dots, I_N) P(O_2 | O_1, I_1, \dots, I_N) P(O_3 | O_1, O_2, I_1, \dots, I_N) \dots P(O_L | O_1, \dots, O_{L-1}, I_1, \dots, I_N) \\
 &= y_1^{O_1} y_2^{O_2} \dots y_L^{O_L}
 \end{aligned}$$

The probability of the output

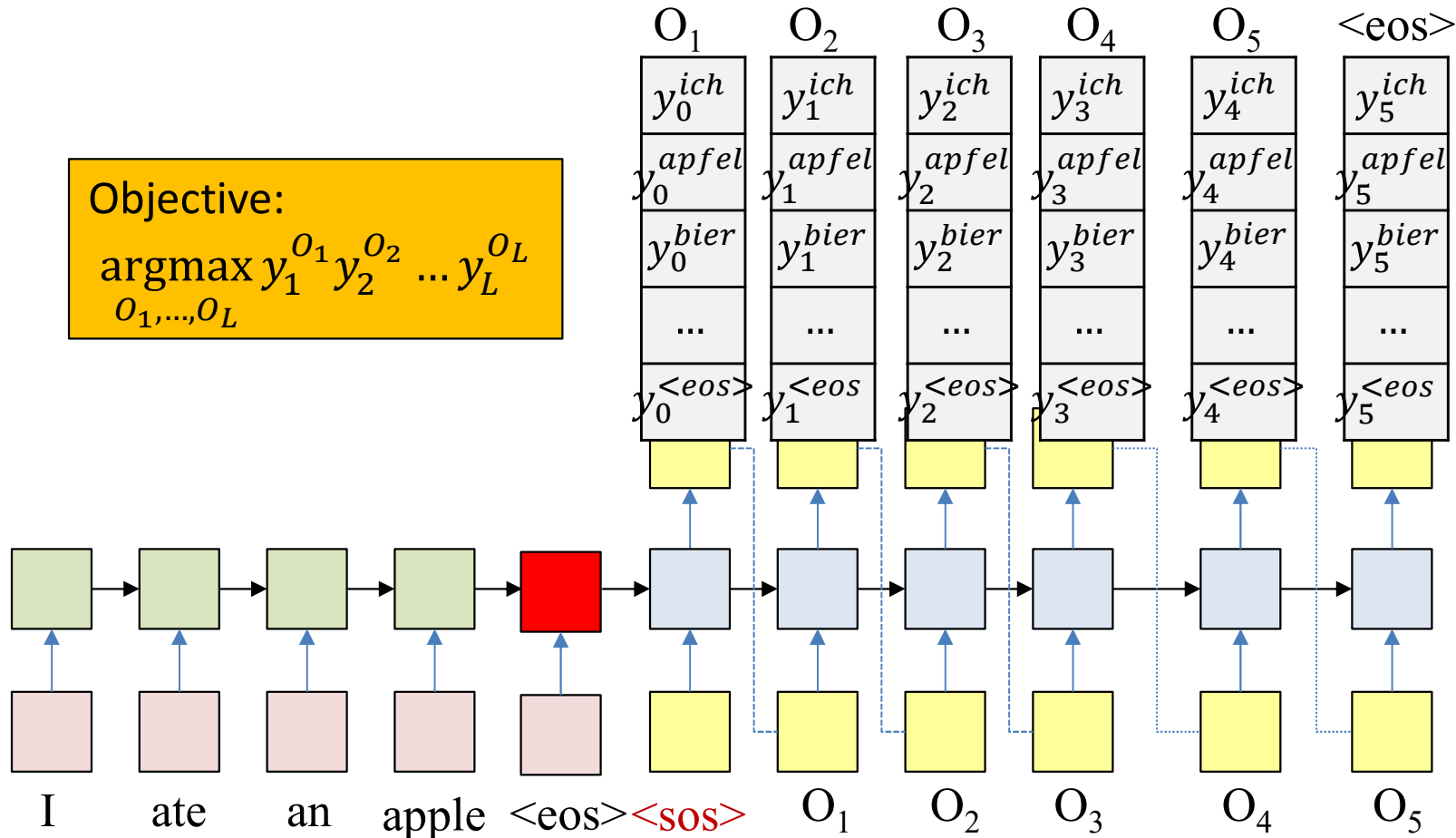


- The objective of drawing: Produce the most likely output (that ends in an <eos>)**

$$\operatorname{argmax}_{O_1, \dots, O_L} P(O_1, \dots, O_L | W_1^{\text{in}}, \dots, W_N^{\text{in}})$$

$$= \operatorname{argmax}_{O_1, \dots, O_L} y_1^{O_1} y_2^{O_2} \dots y_L^{O_L}$$

Greedy drawing



- So how do we draw words at each time to get the most likely word sequence?
- *Greedy* answer – select the most probable word at each time

Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1,1)$  is available for all layers
```

```
t = 0
```

```
do
```

```
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
```

```
until  $x(t) == \text{"<eos>"}$ 
```

```
H =  $h(T-1)$ 
```

```
# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
```

```
t = 0
```

```
 $h_{\text{out}}(0) = H$ 
```

```
# Note: begins with a "start of sentence" symbol
```

```
#      <sos> and <eos> may be identical
```

```
 $y_{\text{out}}(0) = \text{<sos>}$ 
```

```
do
```

```
    t = t+1
```

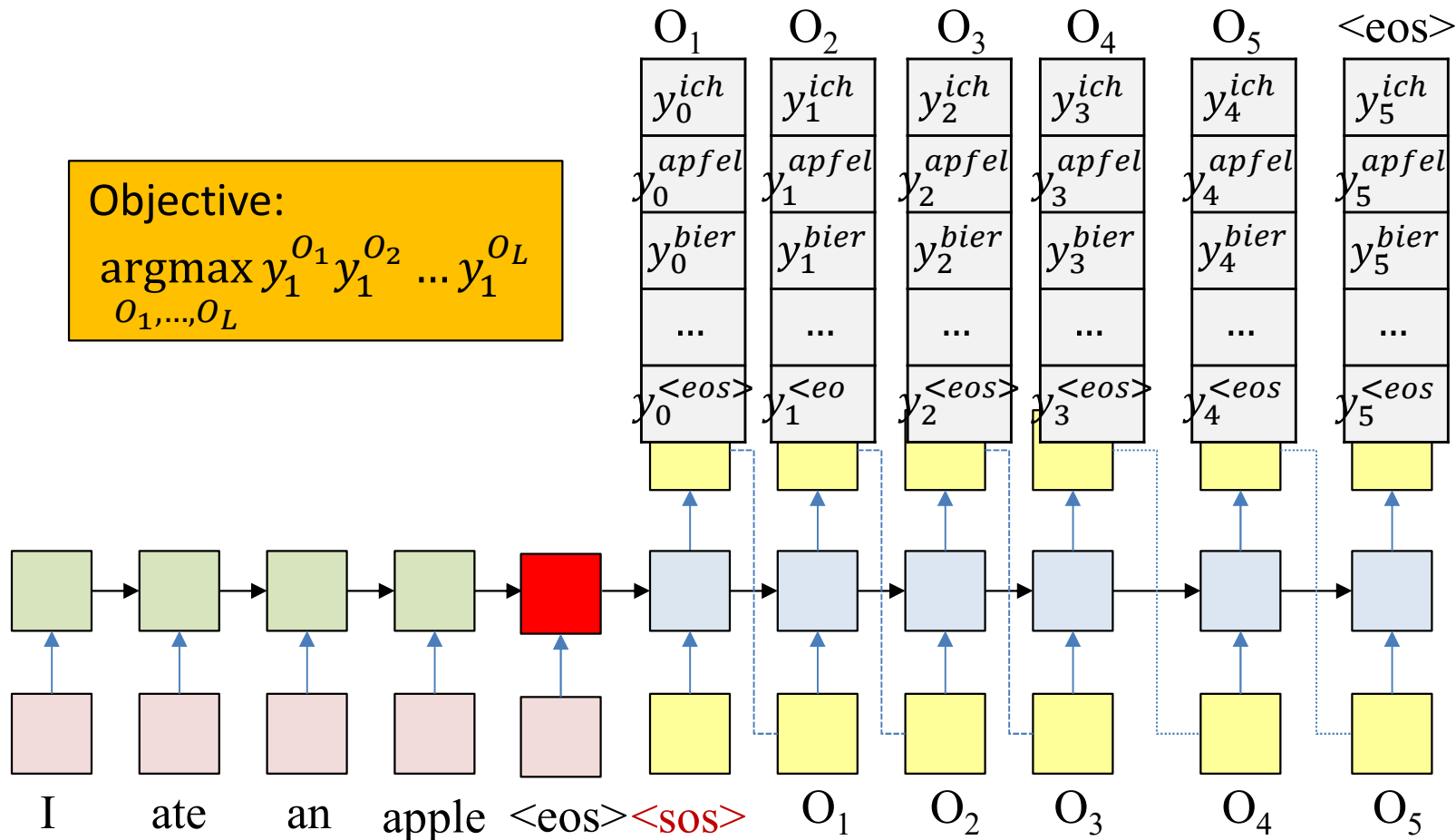
```
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1), y_{\text{out}}(t-1))$ 
```

```
     $y_{\text{out}}(t) = \text{argmax}_i(y(t, i))$ 
```

```
until  $y_{\text{out}}(t) == \text{<eos>}$ 
```

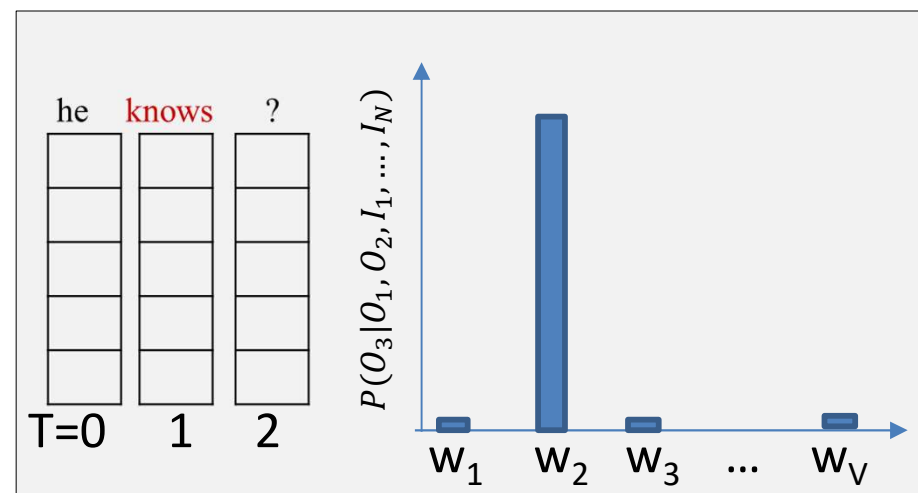
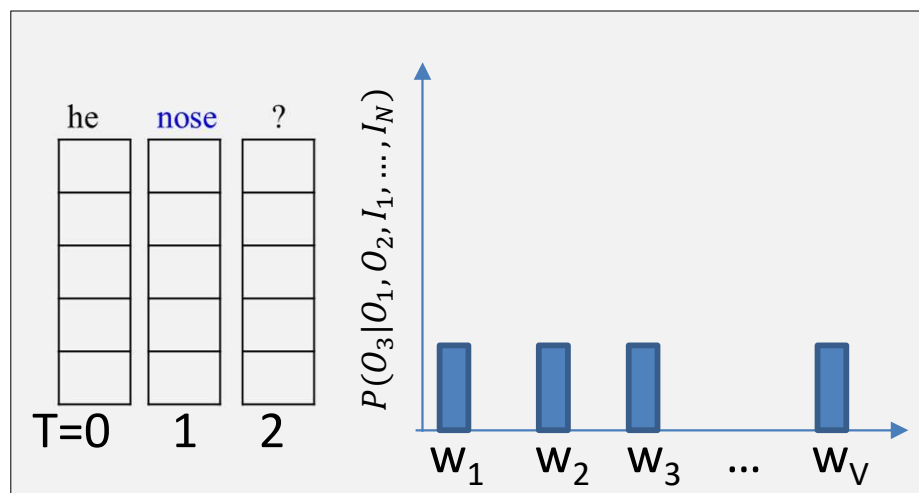
Select the most likely output at each time

Greedy drawing



- Cannot just pick the most likely symbol at each time
 - That may cause the distribution to be more “confused” at the next time
 - Choosing a different, less likely word could cause the distribution at the next time to be more peaky, resulting in a more likely output overall

Greedy is not good

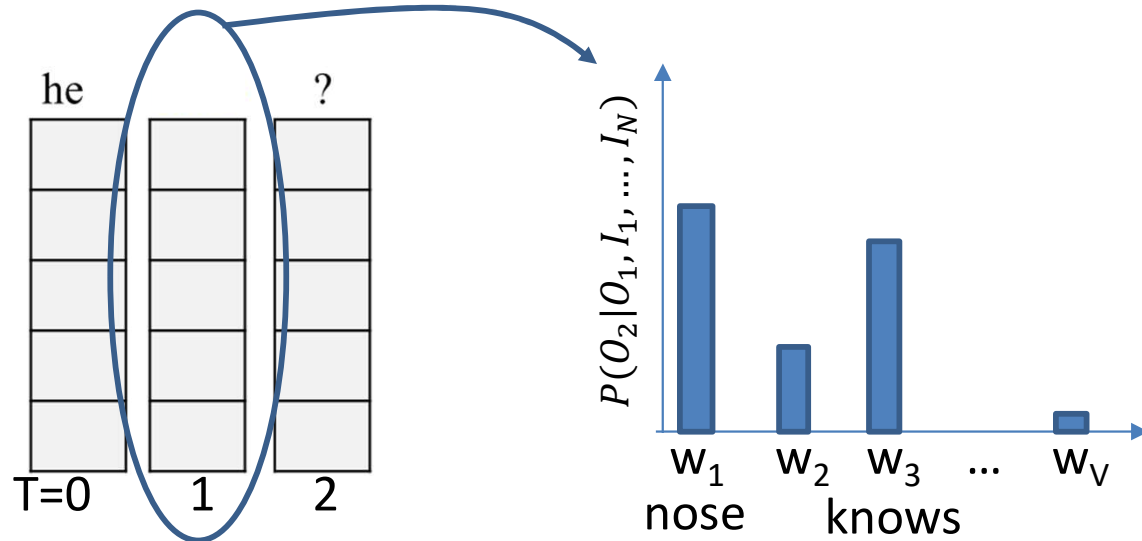


- Hypothetical example (from English speech recognition : Input is speech, output must be text)
- “Nose” has highest probability at $t=2$ and is selected
 - The model is very confused at $t=3$ and assigns low probabilities to many words at the next time
 - Selecting any of these will result in low probability for the entire 3-word sequence
- “Knows” has slightly lower probability than “nose”, but is still high and is selected
 - “he knows” is a reasonable beginning and the model assigns high probabilities to words such as “something”
 - Selecting one of these results in higher overall probability for the 3-word sequence

Greedy is not good

What should we have chosen at $t=2$??

Will selecting “nose” continue to have a bad effect into the distant future?

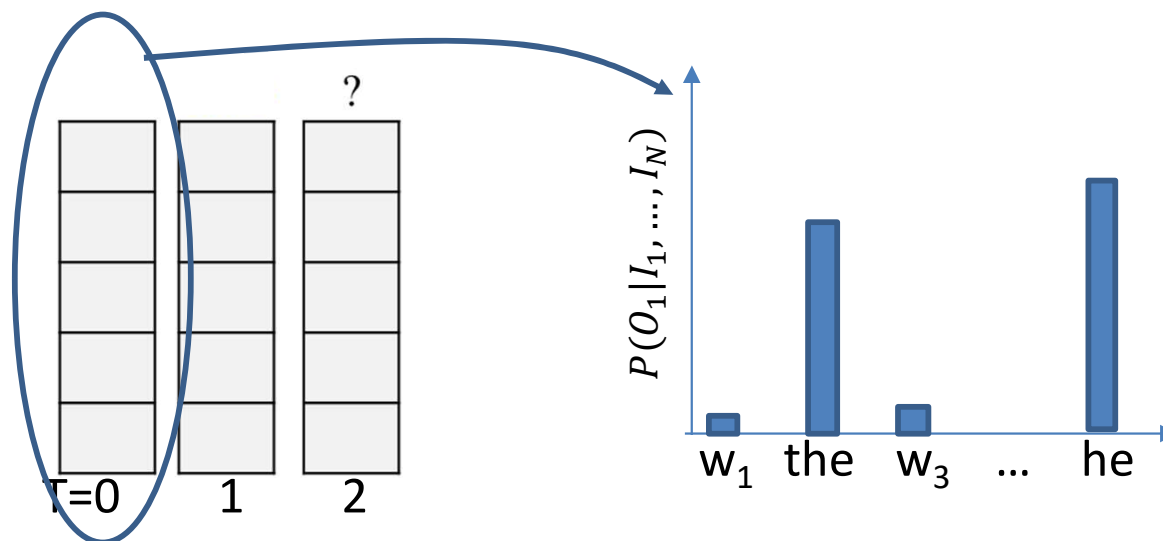


- Problem: Impossible to know a priori which word leads to the more promising future
 - Should we draw “nose” or “knows”?
 - Effect may not be obvious until several words down the line
 - Or the choice of the wrong word early may cumulatively lead to a poorer overall score over time

Greedy is not good

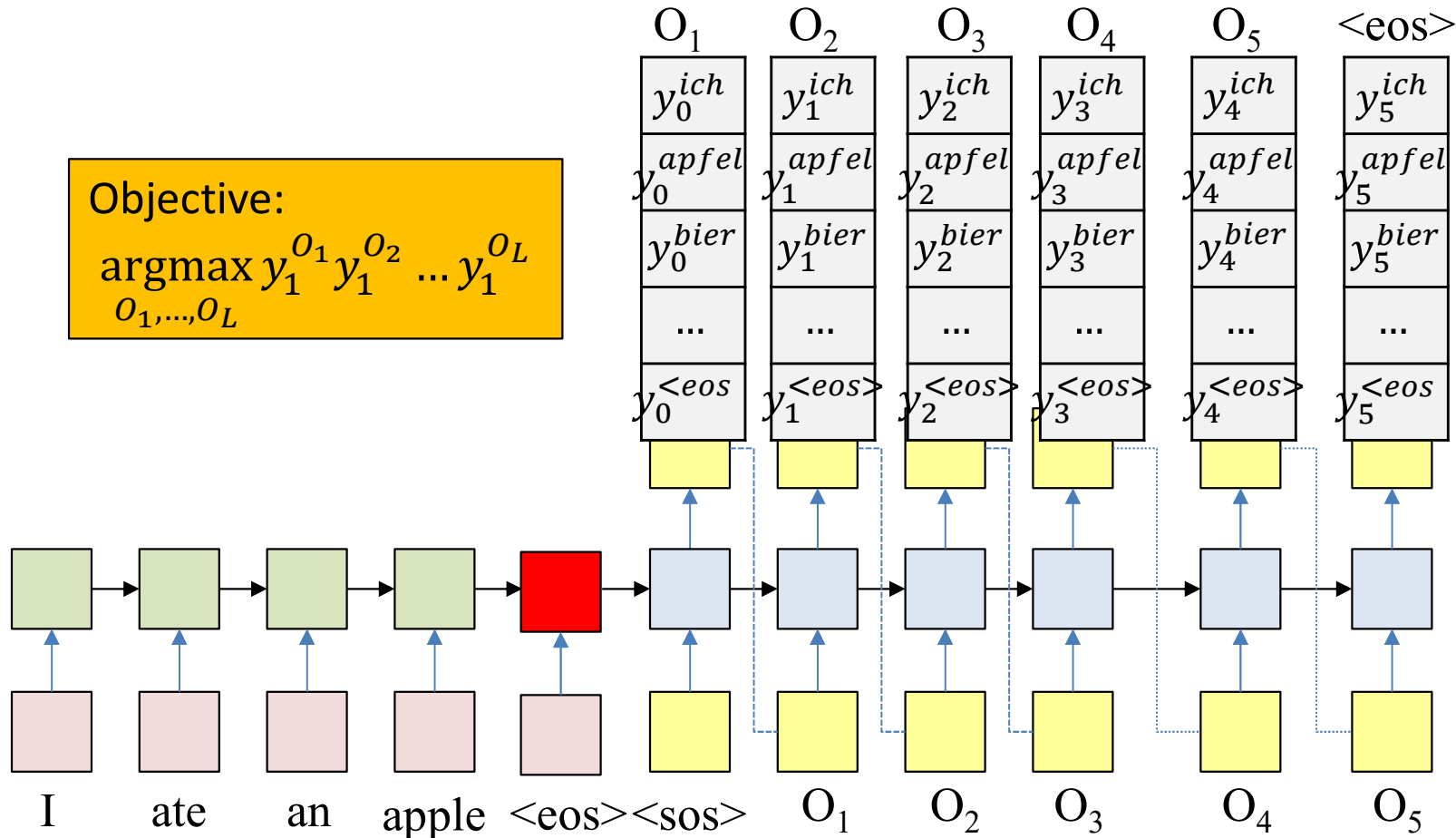
What should we
have chosen at $t=1$??

Choose “the” or “he”?



- Problem: Impossible to know a priori which word leads to the more promising future
 - Even earlier: Choosing the lower probability “the” instead of “he” at $T=0$ may have made a choice of “nose” more reasonable at $T=1$..
- In general, making a poor choice at any time commits us to a poor future
 - But we cannot know at that time the choice was poor

Drawing by random sampling



- Alternate option: Randomly draw a word at each time according to the output probability distribution

Pseudocode

```
# First run the inputs through the network
# Assuming  $h(-1,1)$  is available for all layers
t = 0
do
     $[h(t), \dots] = \text{RNN\_input\_step}(x(t), h(t-1), \dots)$ 
until  $x(t) == \text{"<eos>"}$ 
H =  $h(T-1)$ 
```

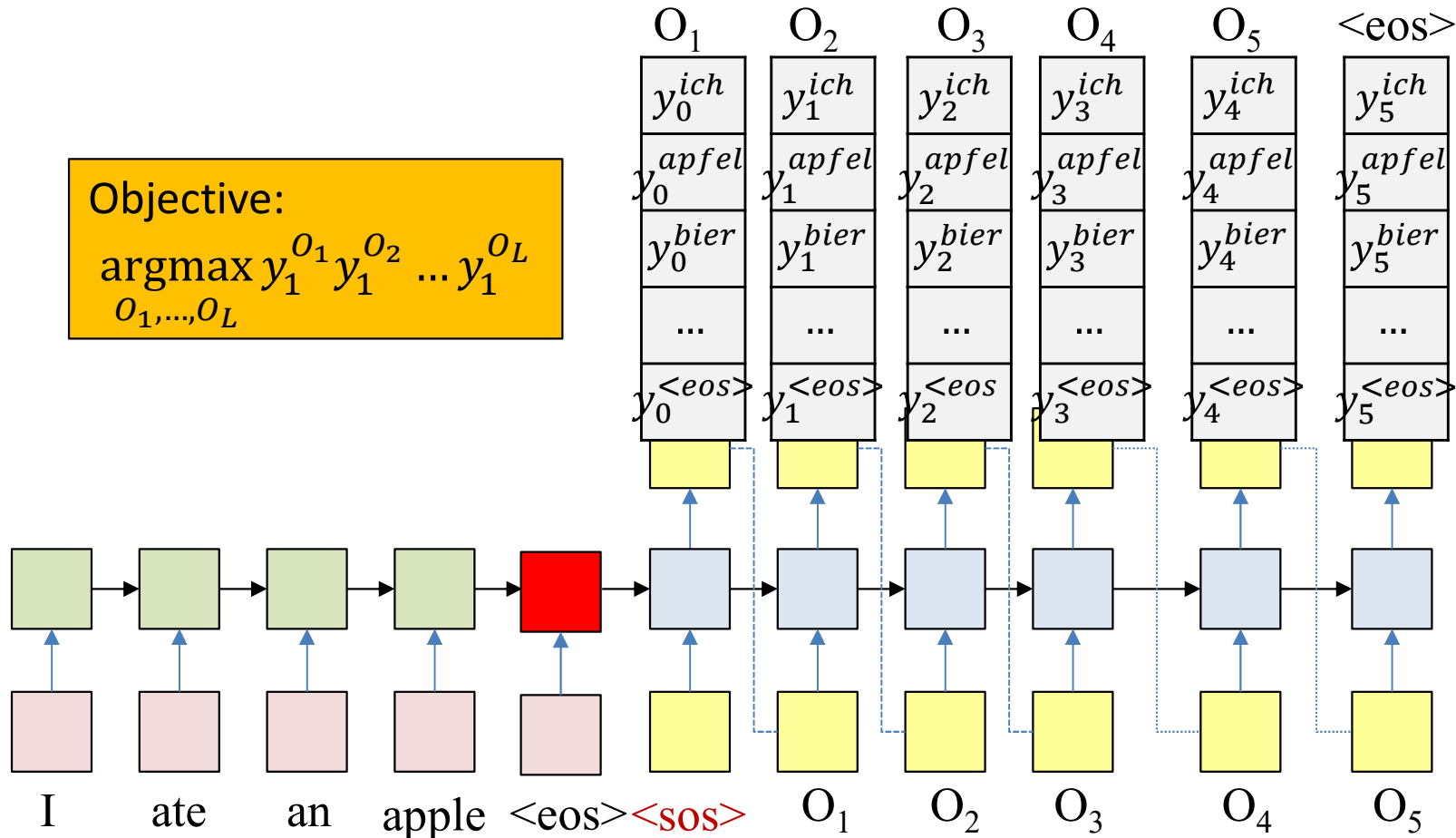
```
# Now generate the output  $y_{\text{out}}(1), y_{\text{out}}(2), \dots$ 
t = 0
 $h_{\text{out}}(0) = H$ 
```

```
# Note: begins with a "start of sentence" symbol
#       <sos> and <eos> may be identical
 $y_{\text{out}}(0) = \text{<sos>}$ 
do
```

```
    t = t+1
     $[y(t), h_{\text{out}}(t)] = \text{RNN\_output\_step}(h_{\text{out}}(t-1), y_{\text{out}}(t-1))$ 
     $y_{\text{out}}(t) = \text{sample}(y(t))$ 
until  $y_{\text{out}}(t) == \text{<eos>}$ 
```

Randomly sample from the output distribution.

Drawing by random sampling



- Alternate option: Randomly draw a word at each time according to the output probability distribution
 - Unfortunately, not guaranteed to give you the most likely output
 - May sometimes give you more likely outputs than greedy drawing though

Poll 2

- @ @ @



Poll 2

For greedy decoding, we choose the word that has been assigned the highest probability at each time (T/F)

- True
- False

In decoding through random sampling we randomly choose the next word according to the probability assigned to it by the decoder (T/F)

- True
- False

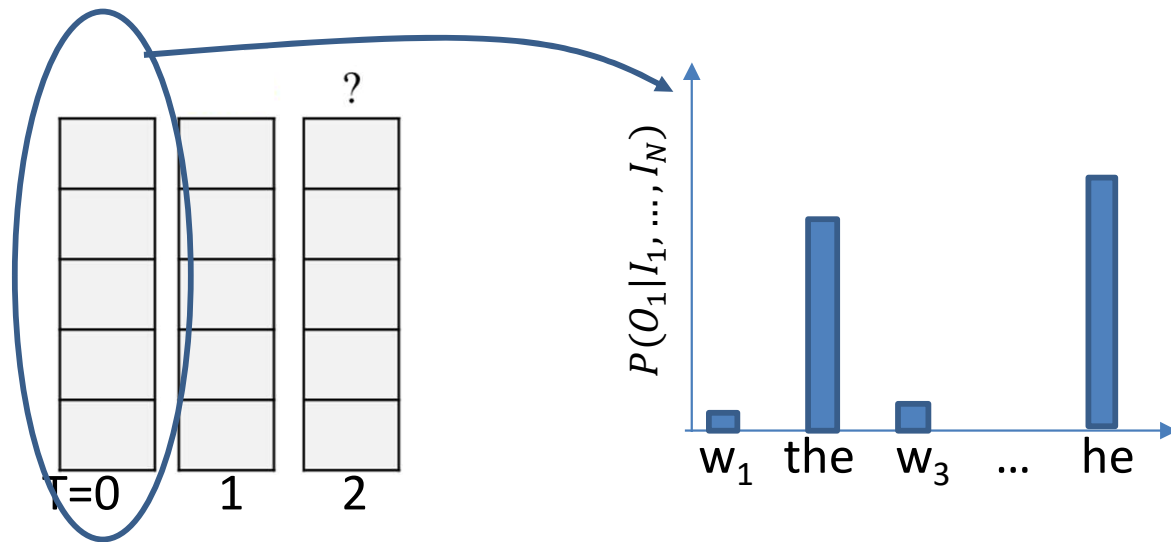
The procedure used for randomly sampling a word from a distribution has been presented in today's class (T/F)

- True
- False

Your choices can get you stuck

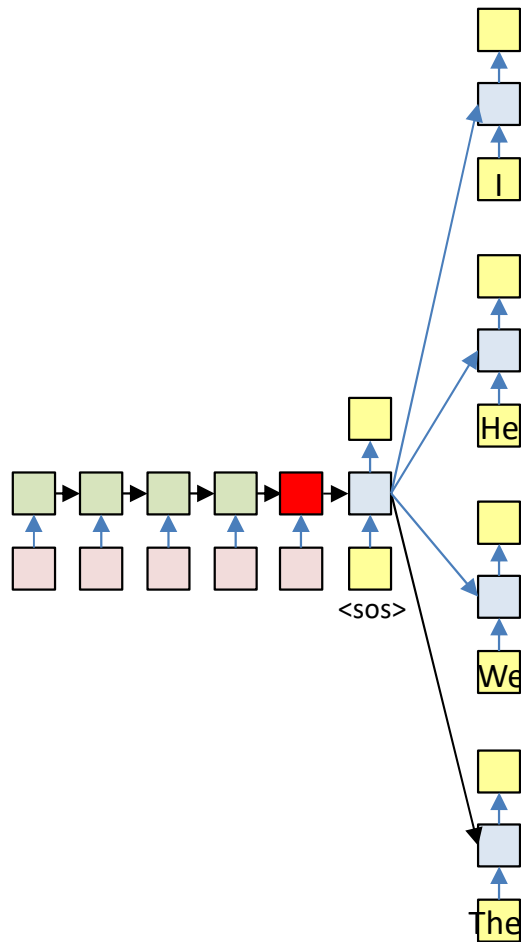
What should we
have chosen at $t=1$??

Choose “the” or “he”?



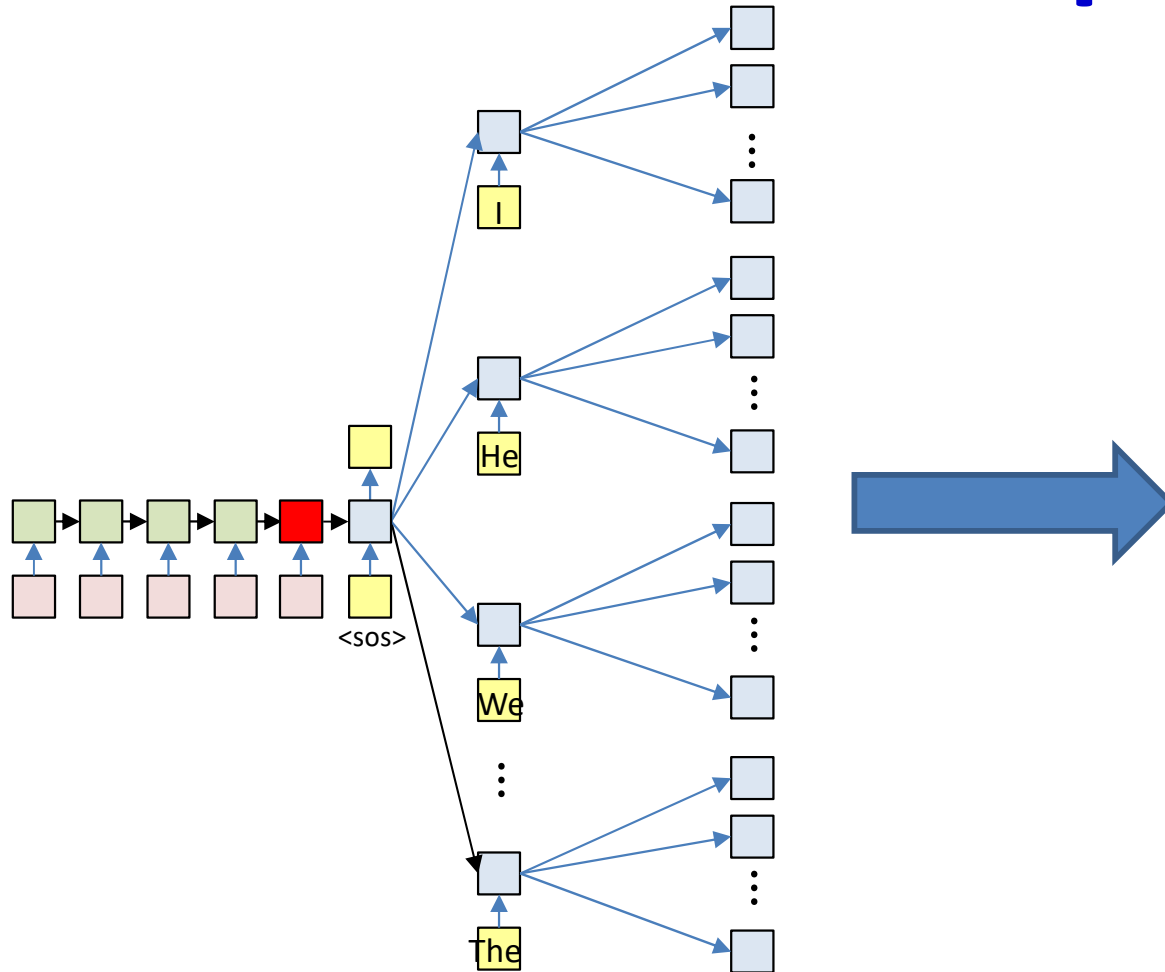
- Problem: making a poor choice at any time commits us to a poor future
 - But we cannot know at that time the choice was poor
- Solution: Don't choose..

Optimal Solution: Multiple choices



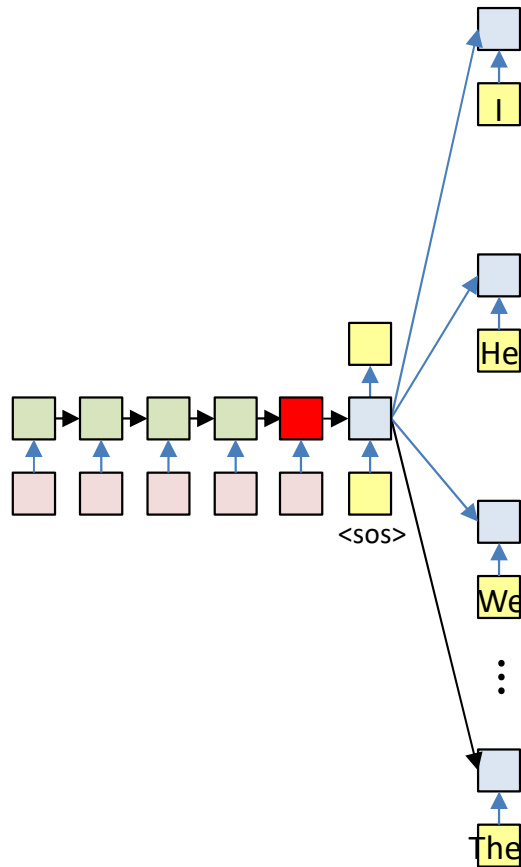
- Retain all choices and *fork* the network
 - With every possible word as input

Problem: Multiple choices



- **Problem:** This will blow up very quickly
 - For an output vocabulary of size V , after T output steps we'd have forked out V^T branches

Solution: Prune

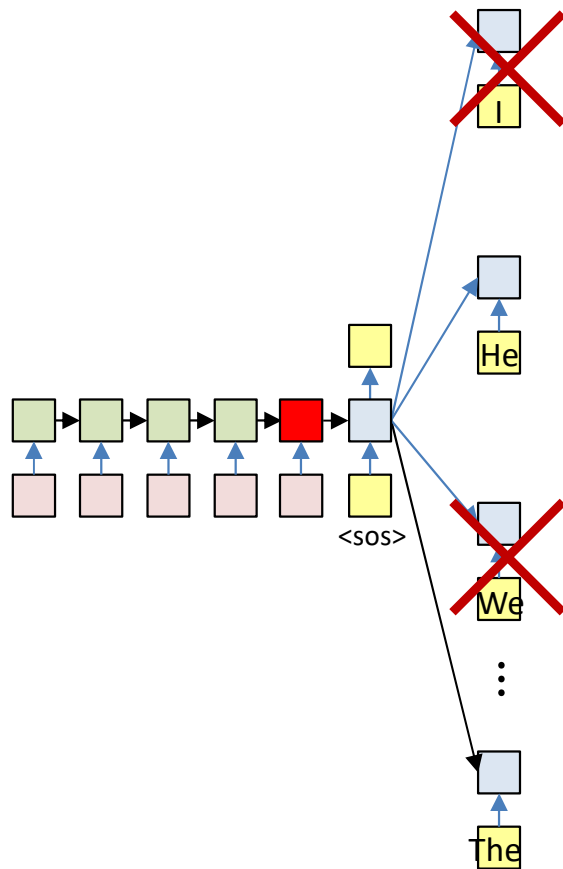


$$Top_K P(O_1|I_1, \dots, I_N)$$

- **Solution: Prune**

- At each time, retain only the top K scoring forks

Solution: Prune

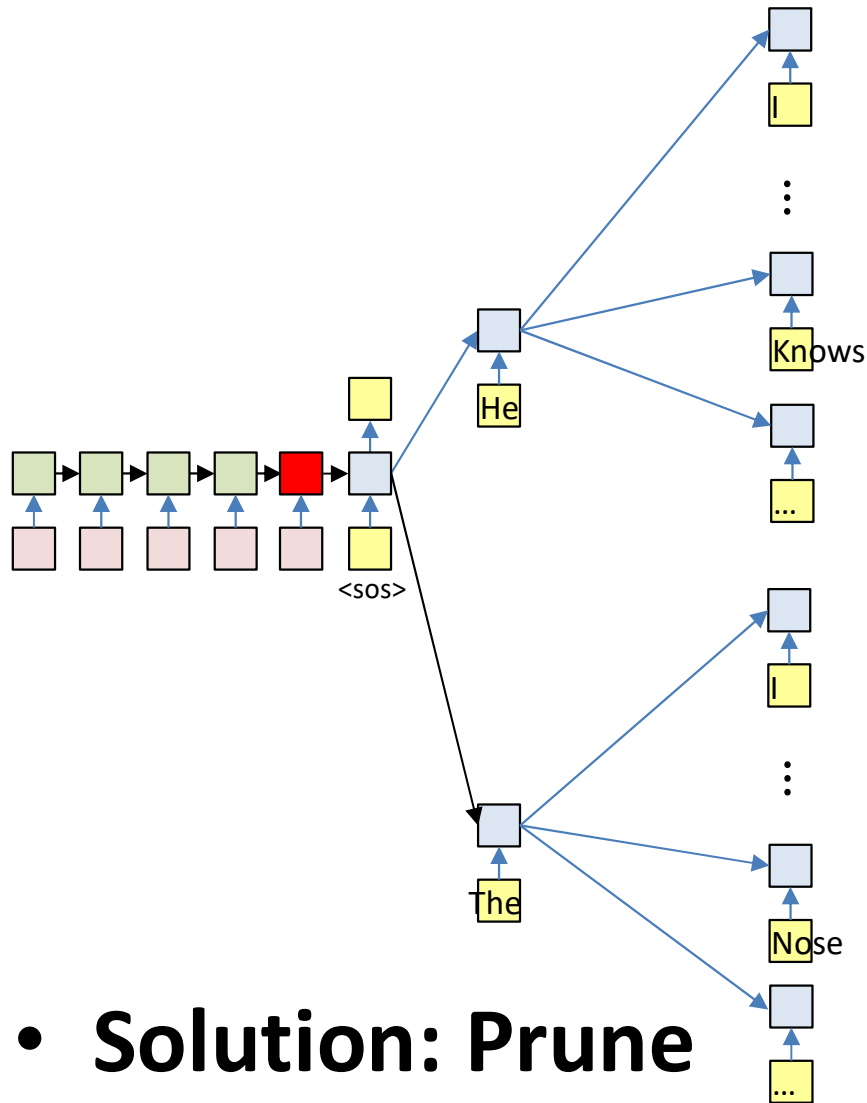


$$Top_K P(O_1 | I_1, \dots, I_N)$$

- **Solution: Prune**

- At each time, retain only the top K scoring forks

Solution: Prune



Note: based on product

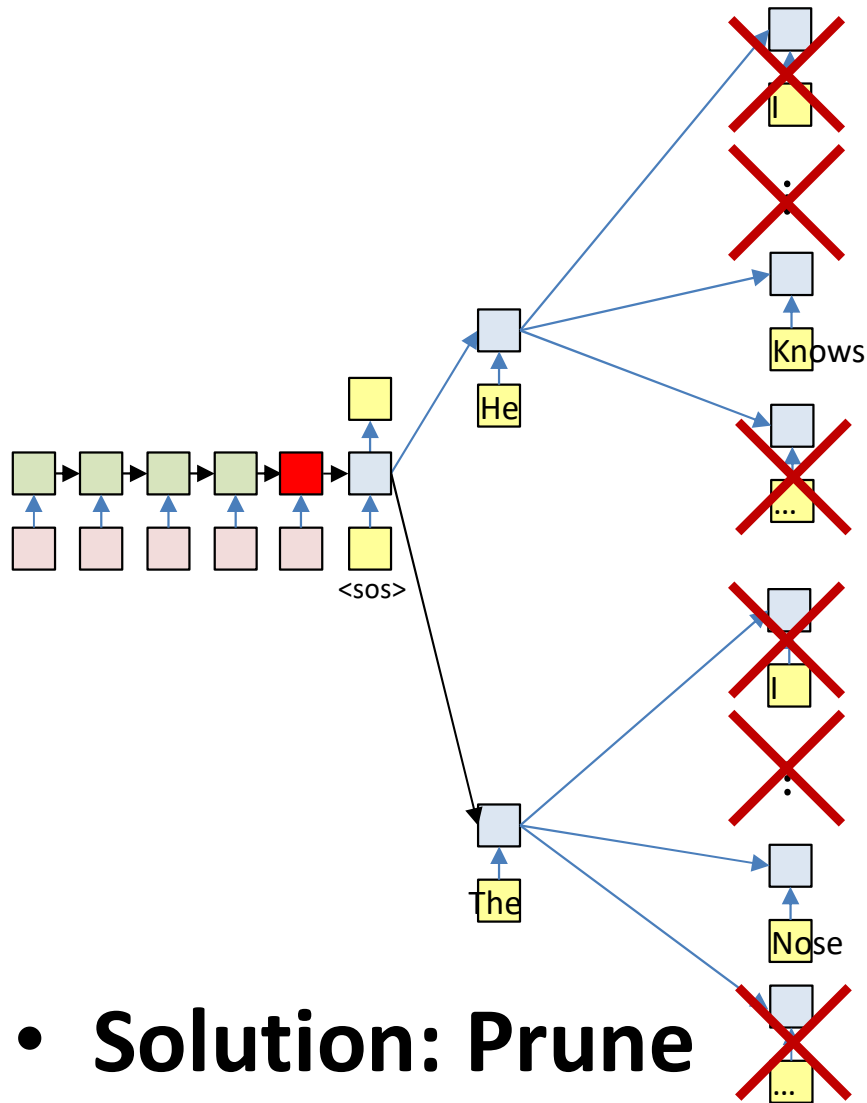
$$\text{Top}_K P(O_2 O_1 | I_1, \dots, I_N)$$

$$= \text{Top}_K P(O_2 | O_1, I_1, \dots, I_N) P(O_1 | I_1, \dots, I_N)$$

- **Solution: Prune**

- At each time, retain only the top K scoring forks

Solution: Prune



Note: based on product

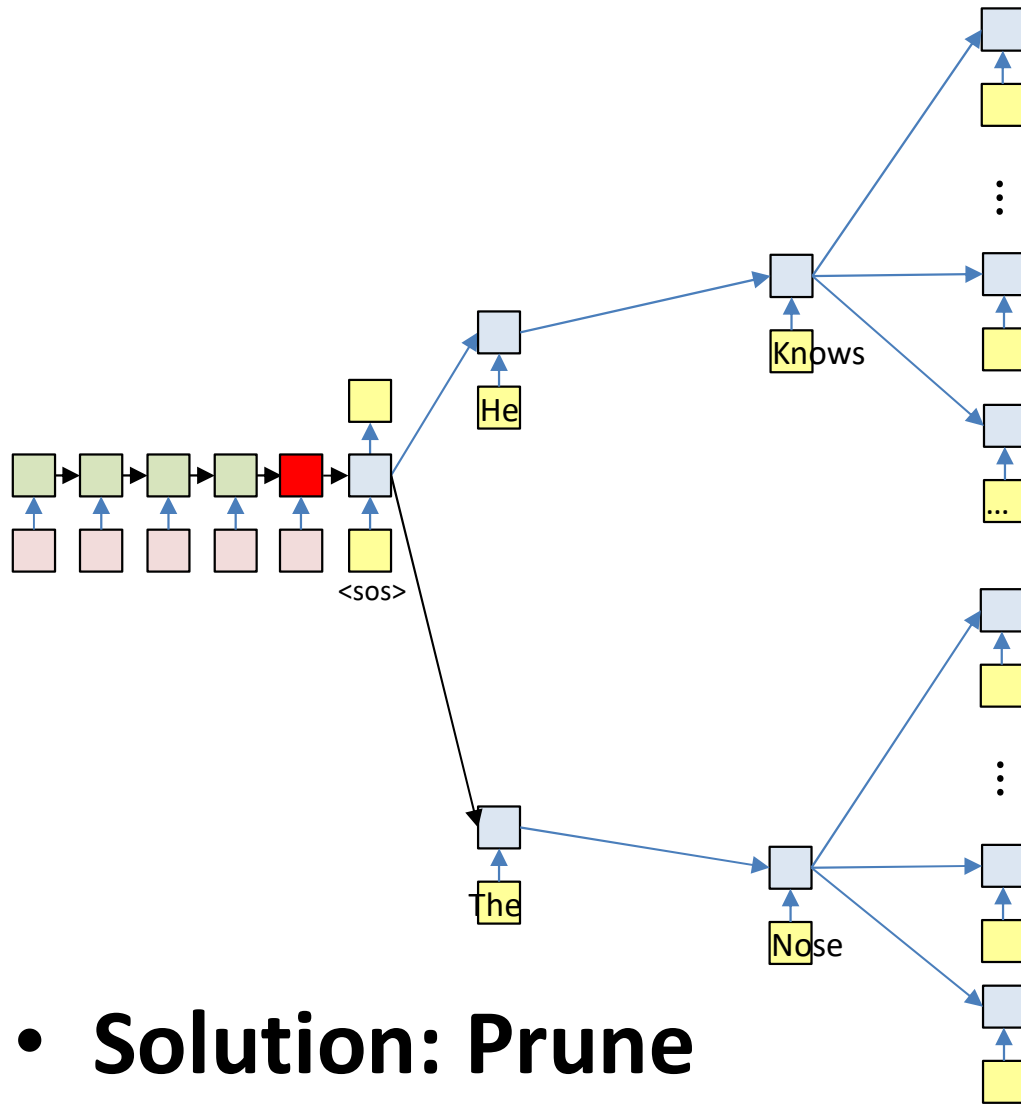
$$Top_K P(O_2 O_1 | I_1, \dots, I_N)$$

$$= Top_K P(O_2 | O_1, I_1, \dots, I_N) P(O_1 | I_1, \dots, I_N)$$

• Solution: Prune

- At each time, retain only the top K scoring forks

Solution: Prune

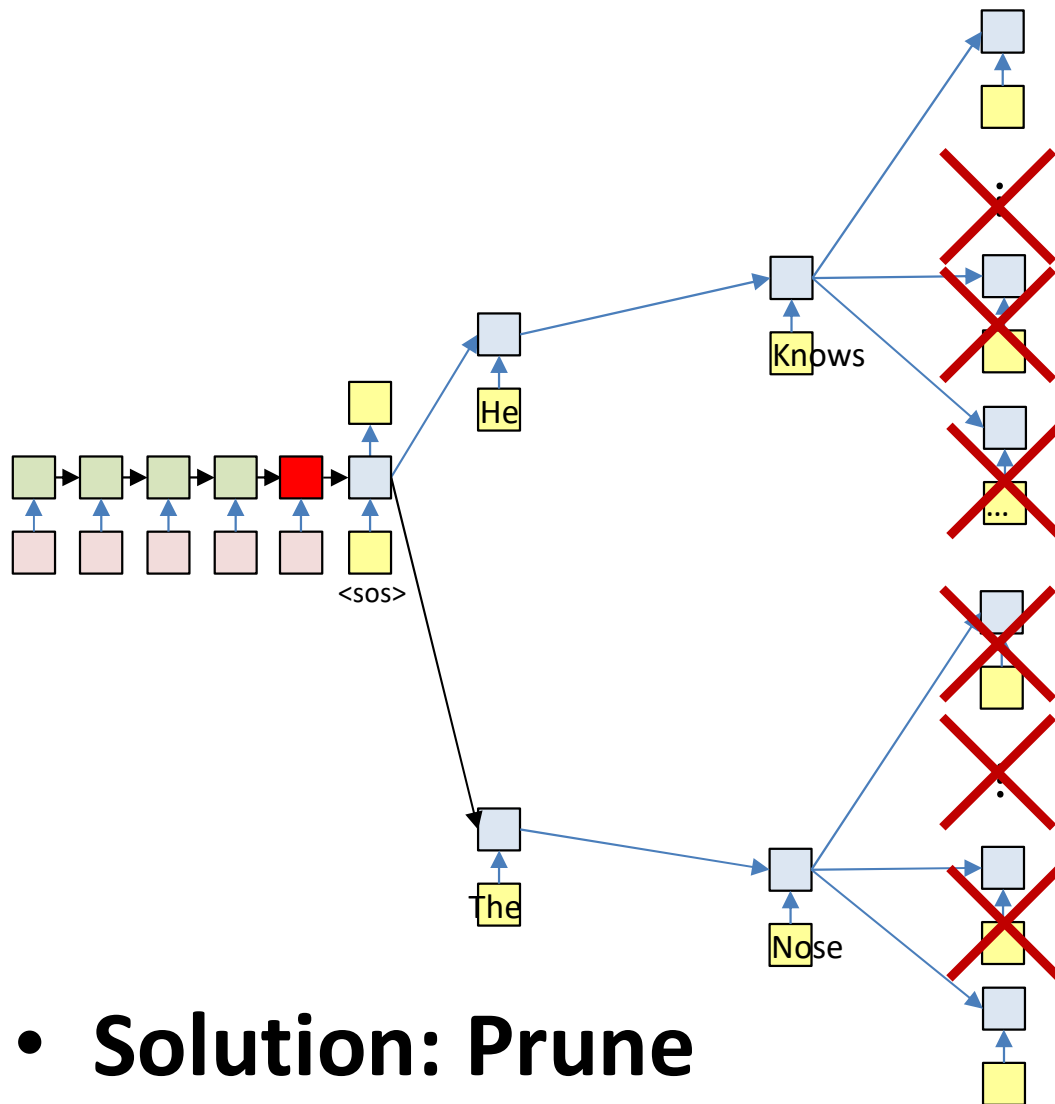


$$= \text{Top}_K P(O_3|O_1, O_2, I_1, \dots, I_N) \times \\ P(O_2|O_1, I_1, \dots, I_N) \times \\ P(O_1|I_1, \dots, I_N)$$

- **Solution: Prune**

- At each time, retain only the top K scoring forks

Solution: Prune

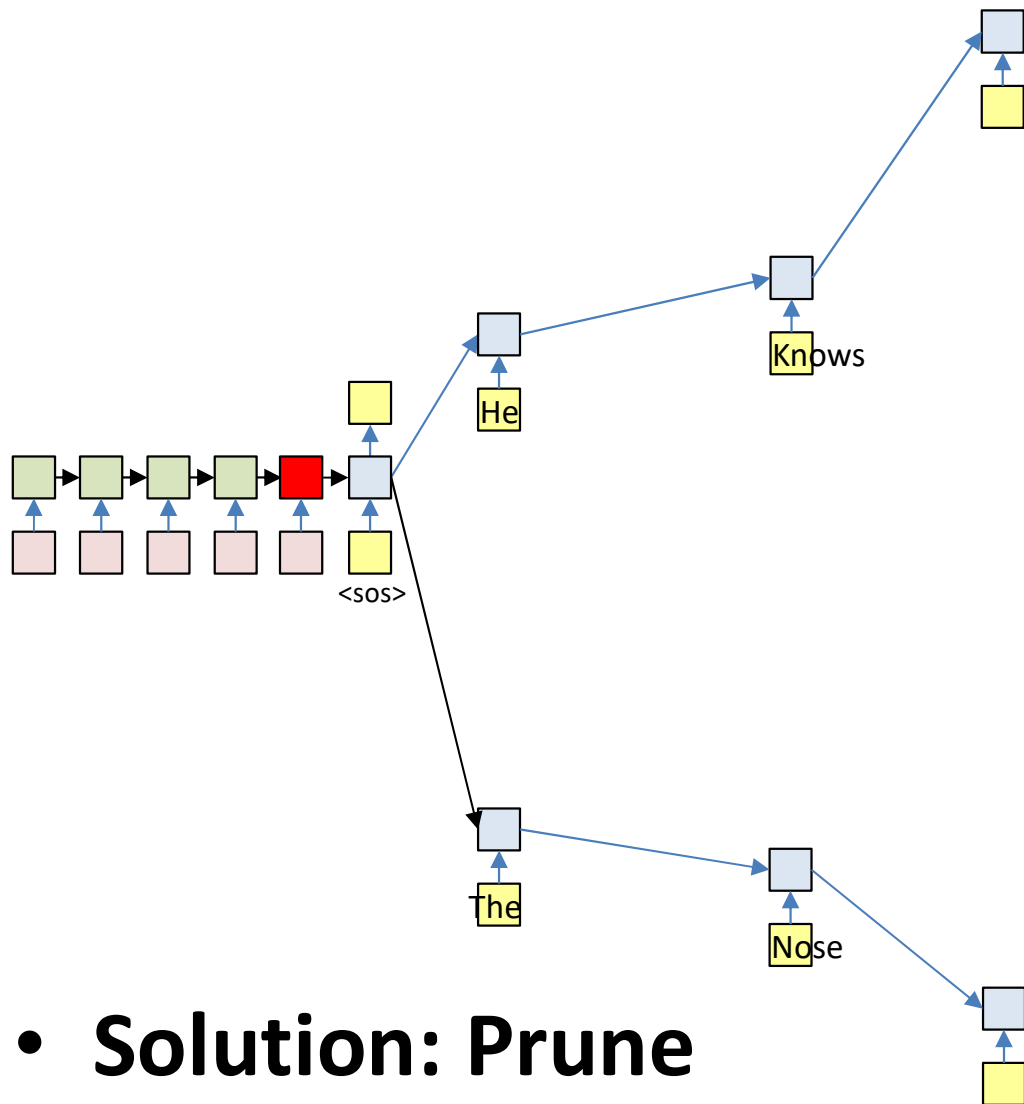


$$= \text{Top}_K P(O_2|O_1, O_2, I_1, \dots, I_N) \times \\ P(O_2|O_1, I_1, \dots, I_N) \times \\ P(O_1|I_1, \dots, I_N)$$

- **Solution: Prune**

- At each time, retain only the top K scoring forks

Solution: Prune

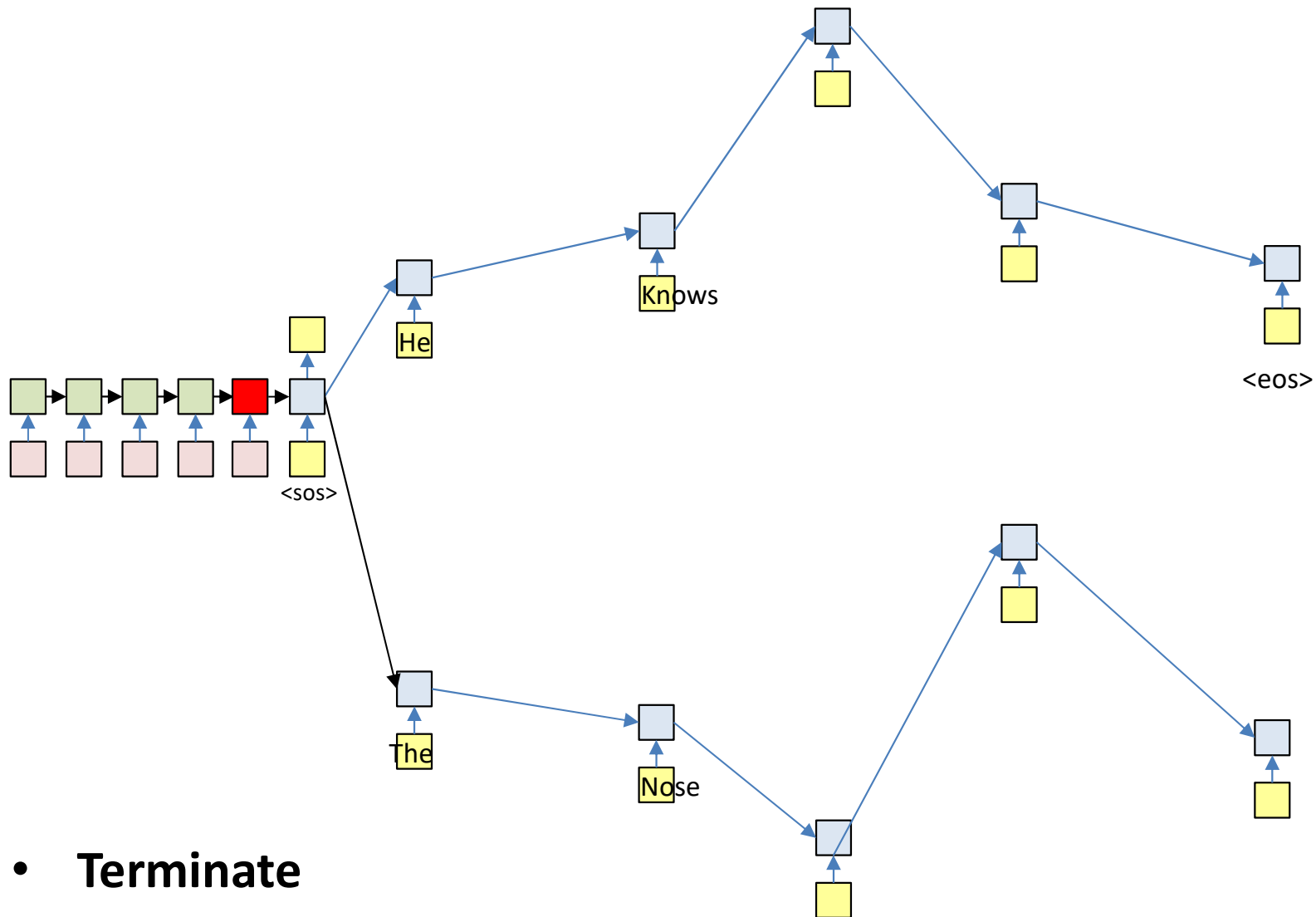


$$Top_K \prod_{t=1}^n P(O_t | O_1, \dots, O_{t-1}, I_1, \dots, I_N)$$

- **Solution: Prune**

- At each time, retain only the top K scoring forks

Terminate

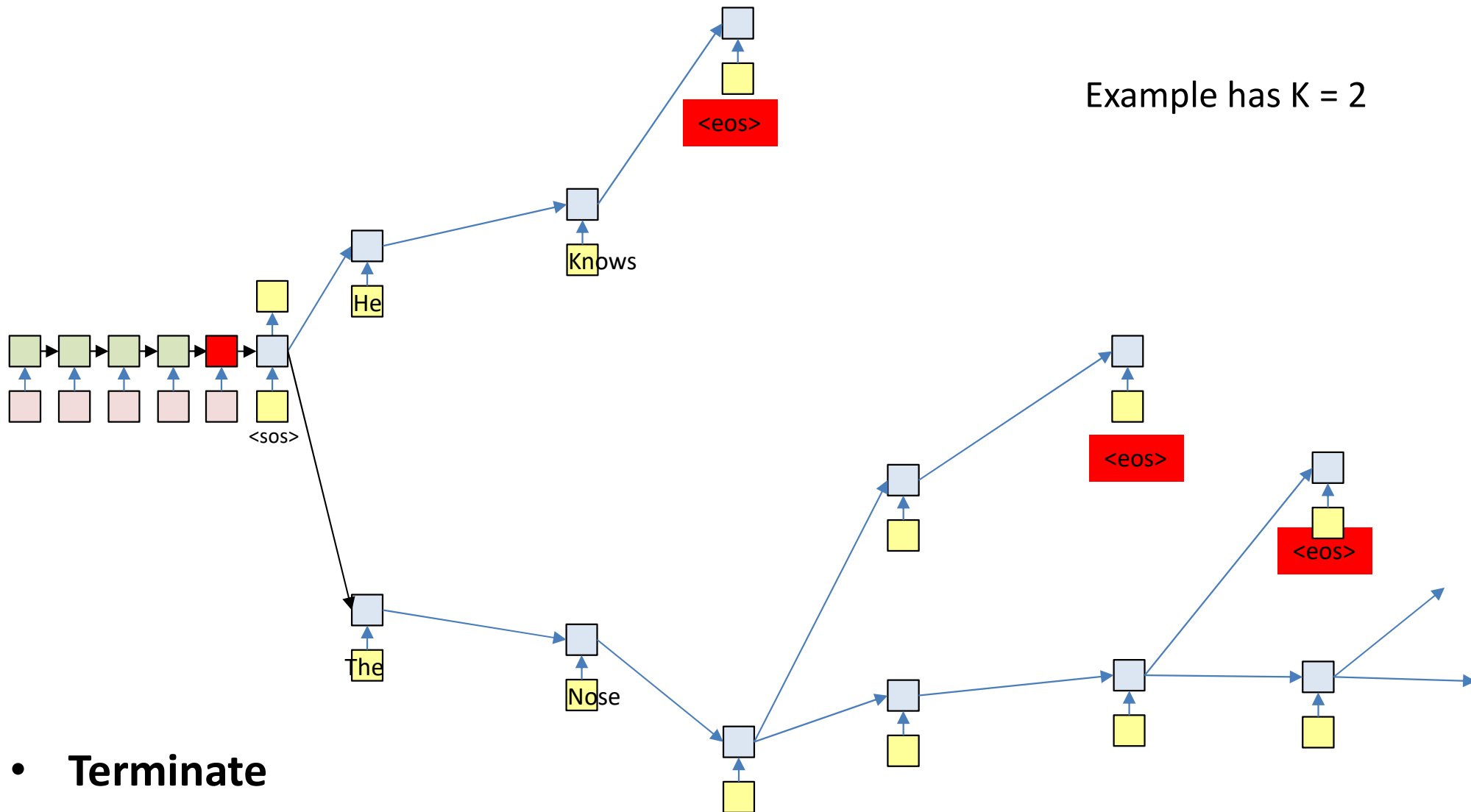


- **Terminate**

- When the current most likely path overall ends in <eos>

- Or continue producing more outputs (each of which terminates in <eos>) to get N-best outputs

Termination: <eos>



- **Terminate**

- Paths cannot continue once the output an <eos>

- So paths may be different lengths

- Select the most likely sequence ending in <eos> across *all* terminating sequences

Pseudocode: Beam search

```
# Assuming encoder output H is available
path = <sos>
beam = {path}
pathscore = [path] = 1
state[path] = h[0] # Output of encoder
do # Step forward
    nextbeam = {}
    nextpathscore = []
    nextstate = {}
    for path in beam:
        cfin = path[end]
        hpath = state[path]
        [y,h] = RNN_output_step(hpath,cfin)
        for c in Symbolset
            newpath = path + c
            nextstate[newpath] = h
            nextpathscore[newpath] = pathscore[path]*y[c]
            nextbeam += newpath # Set addition
        end
    end
    beam, pathscore, state, bestpath = prune(nextstate,nextpathscore,nextbeam,bw)
until bestpath[end] = <eos>
```

Pseudocode: Prune

Note, there are smarter ways to implement this

```
function prune (state, score, beam, beamwidth)
    sortedscore = sort(score)
    threshold = sortedscore[beamwidth]

    prunedstate = {}
    prunedscore = []
    prunedbeam = {}

    bestscore = -inf
    bestpath = none
    for path in beam:
        if score[path] > threshold:
            prunedbeam += path # set addition
            prunedstate[path] = state[path]
            prunedscore[path] = score[path]
            if score[path] > bestscore
                bestscore = score[path]
                bestpath = path
            end
        end
    end
    return prunedbeam, prunedscore, prunedstate, bestpath
```


Poll 3

- @, @



Poll 3

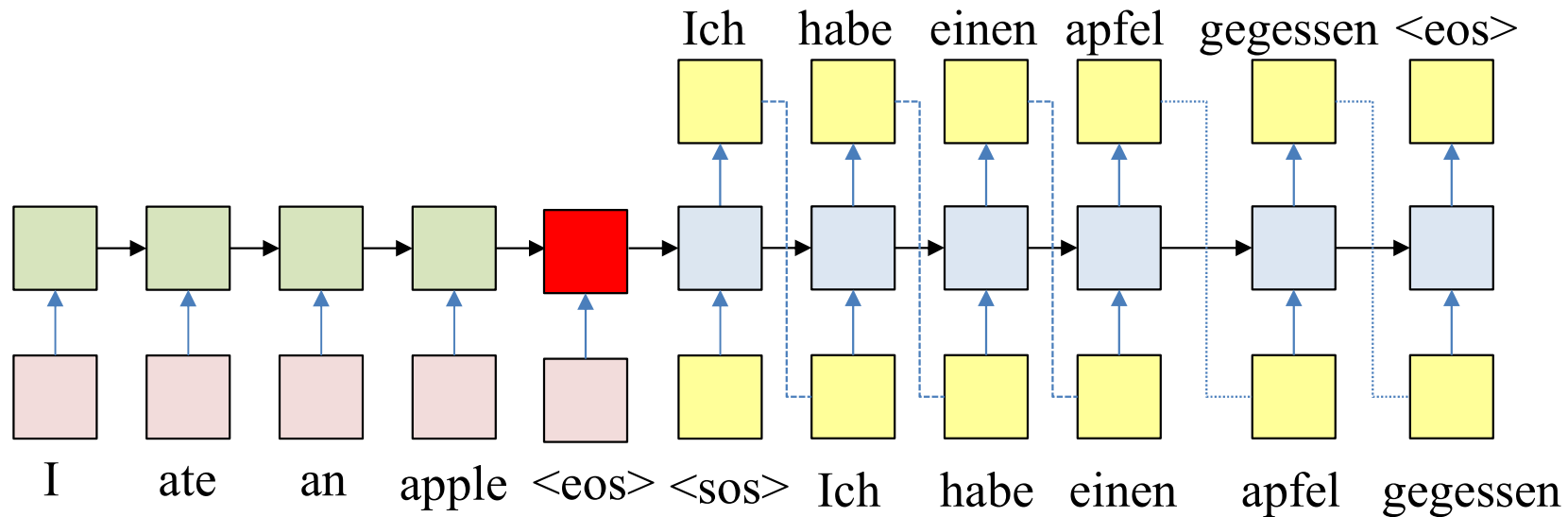
“Theoretically correct” decoding requires you to evaluate the *entire tree representing every possible word sequence* to select the best one (T/F)

- True
- False

Beam search is theoretically correct decoding (T/F)

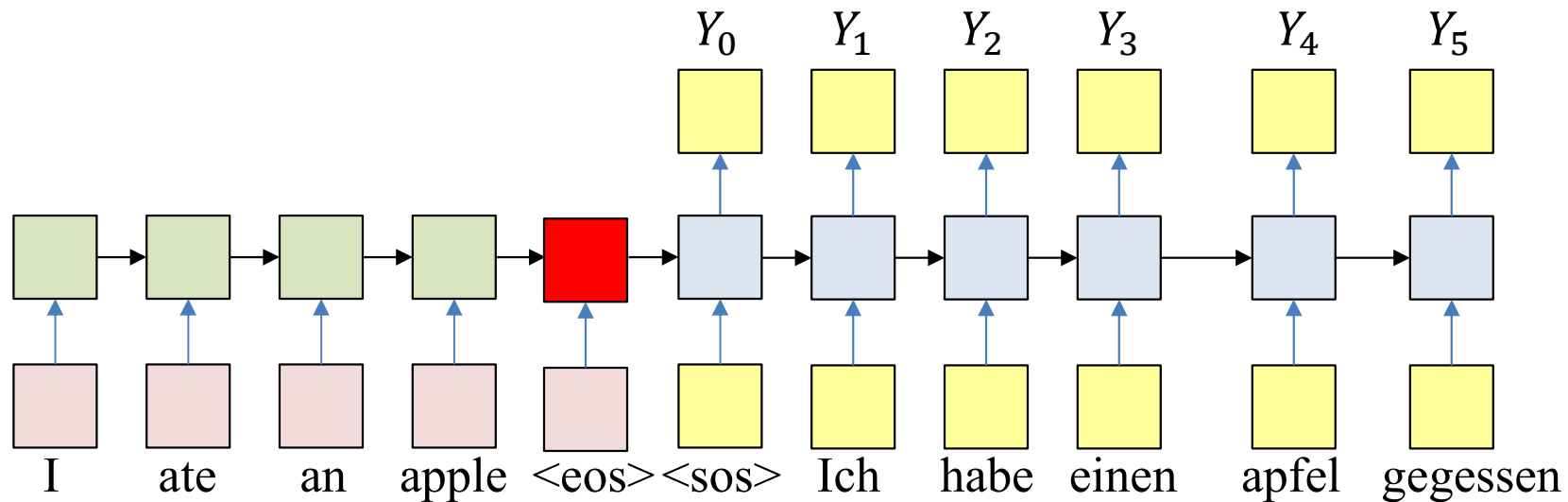
- True
- False

Training the system



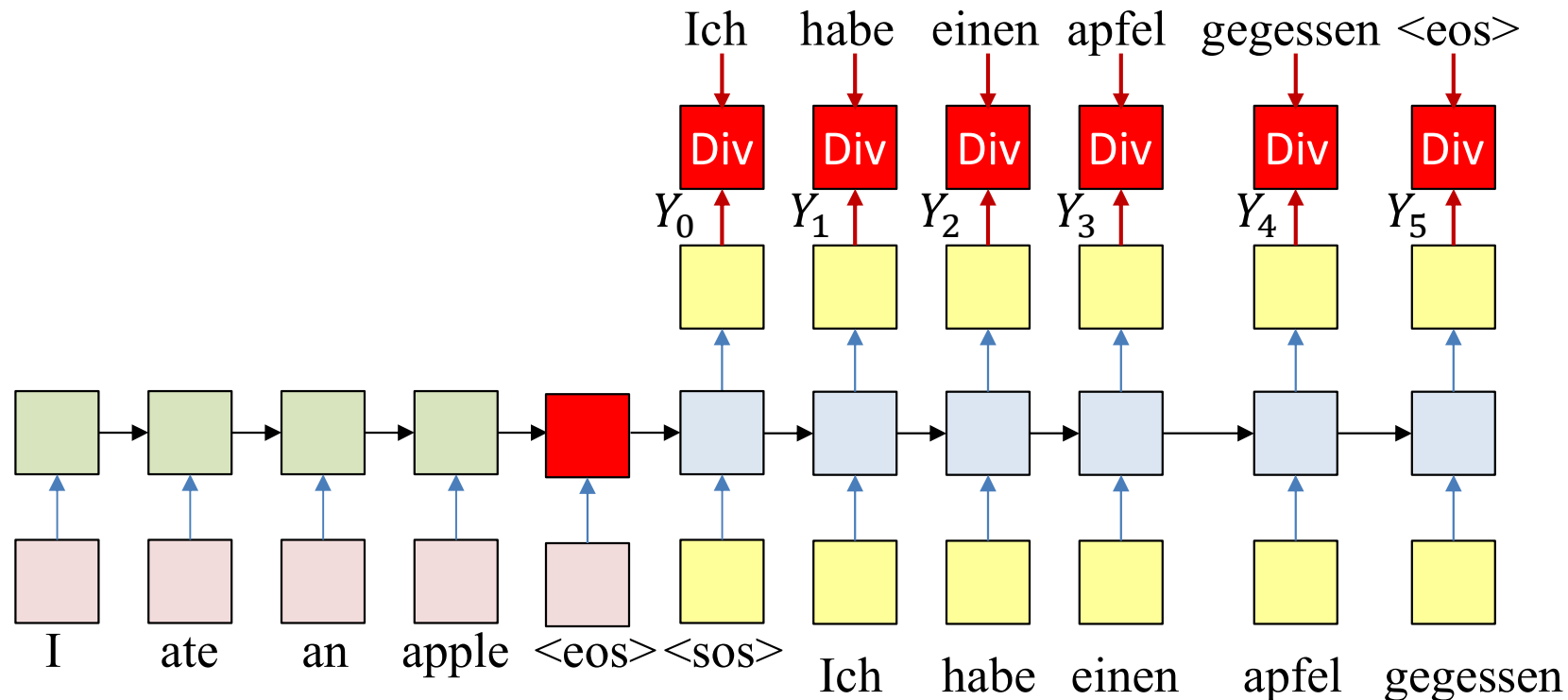
- Must learn to make predictions appropriately
 - Given "I ate an apple <eos>", produce "Ich habe einen apfel gegessen <eos>".

Training : Forward pass



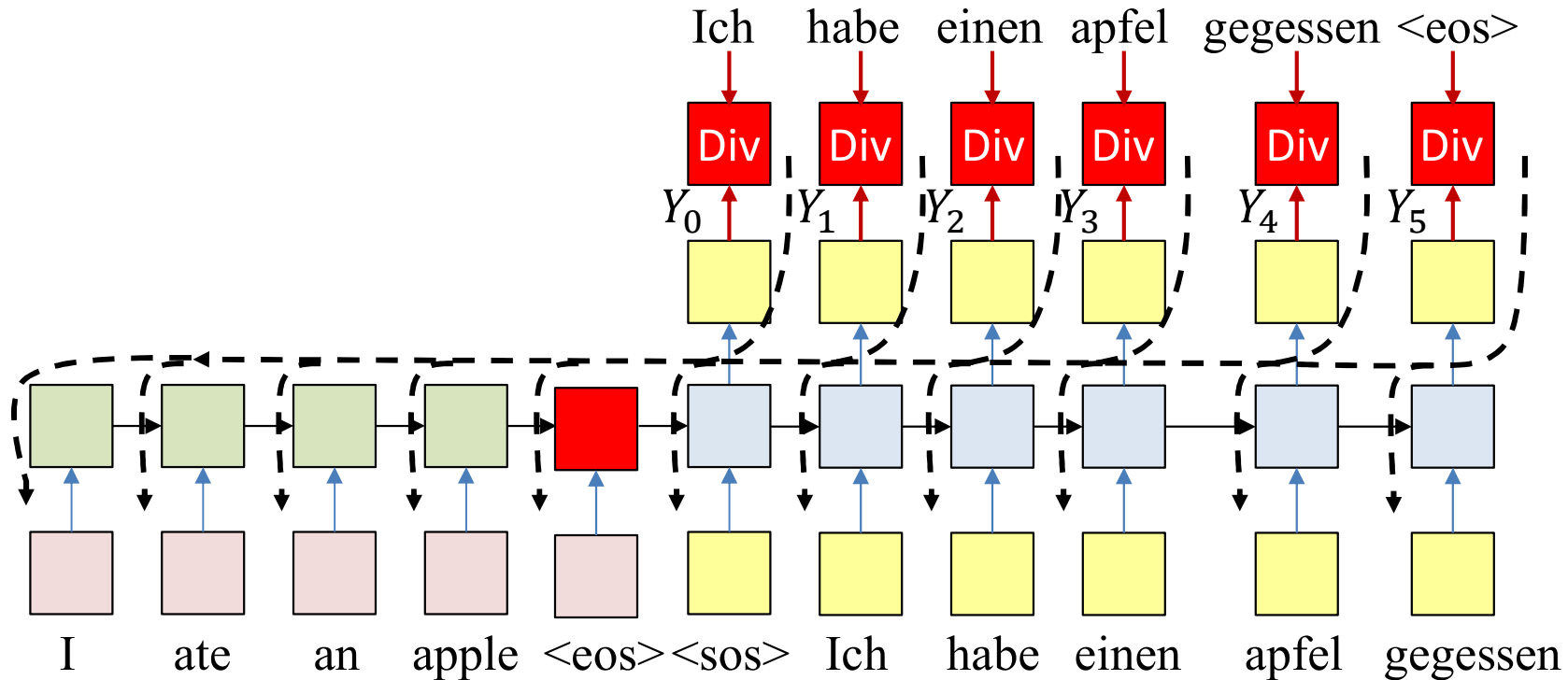
- Forward pass: Input the source and target sequences, sequentially
 - Output will be a probability distribution over target symbol set (vocabulary)

Training : Backward pass



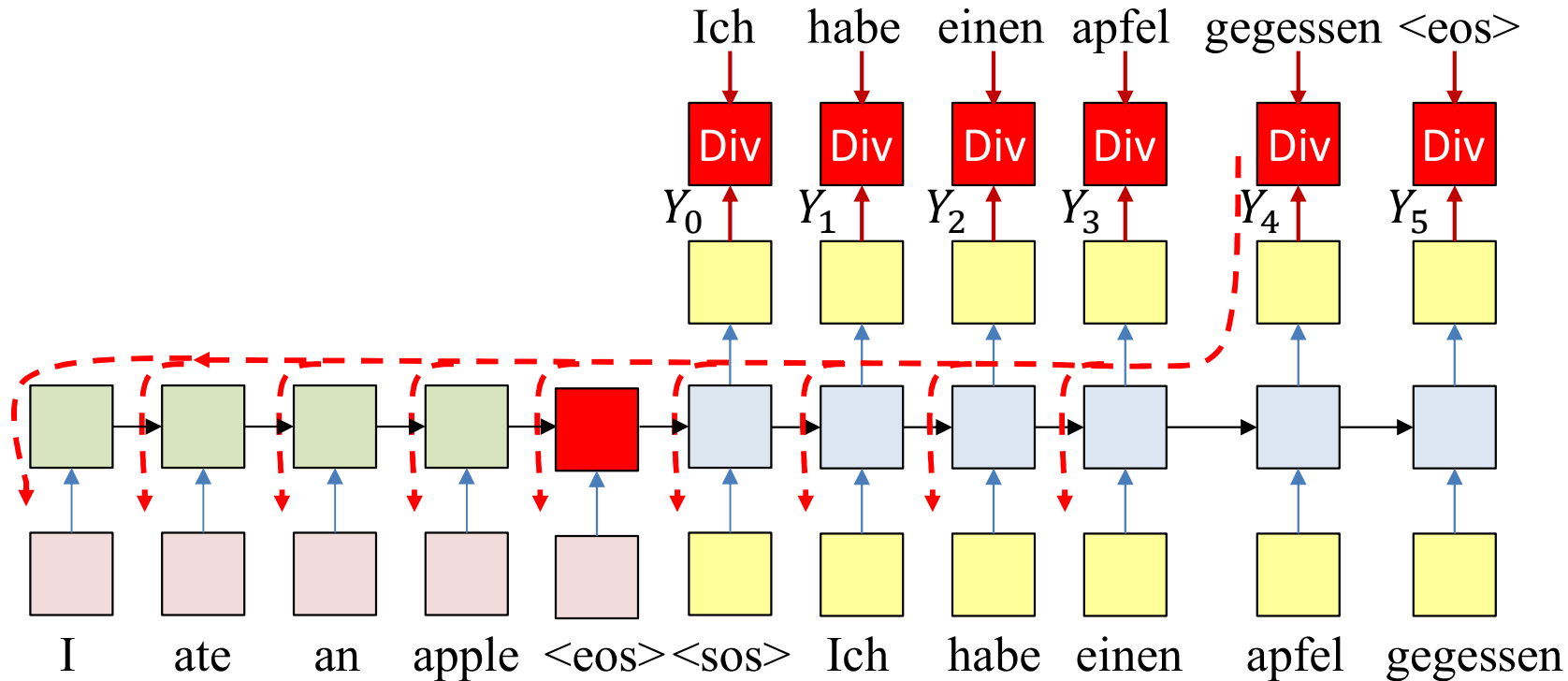
- Backward pass: Compute the divergence between the output distribution and target word sequence

Training : Backward pass



- Backward pass: Compute the divergence between the output distribution and target word sequence
- Backpropagate the derivatives of the divergence through the network to learn the net

Training : Backward pass

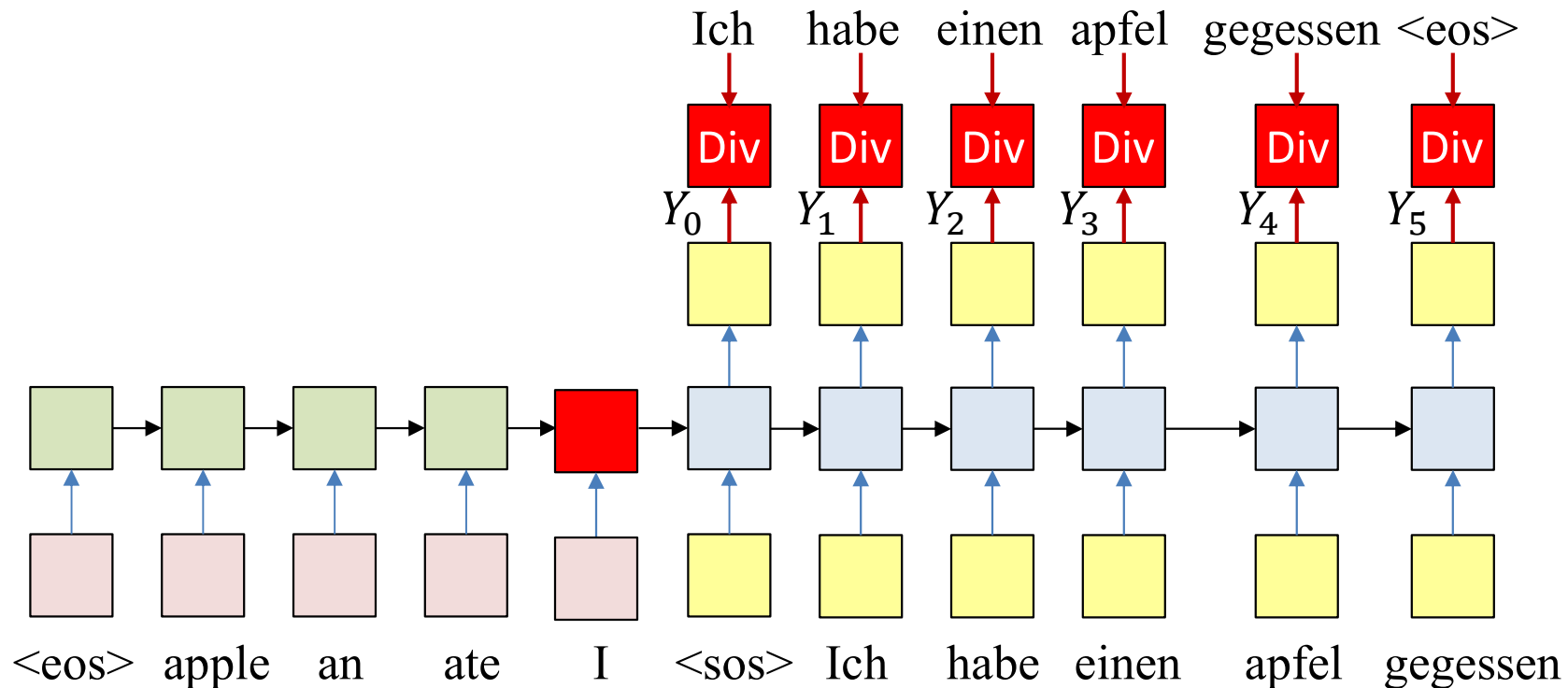


- In practice, if we apply SGD, we may randomly sample words from the output to actually use for the backprop and update
 - Typical usage: Randomly select one word from each input training instance (comprising an input-output pair)
 - For each iteration
 - Randomly select training instance: (input, output)
 - Forward pass
 - Randomly select a single output $y(t)$ and corresponding desired output $d(t)$ for backprop

Overall training

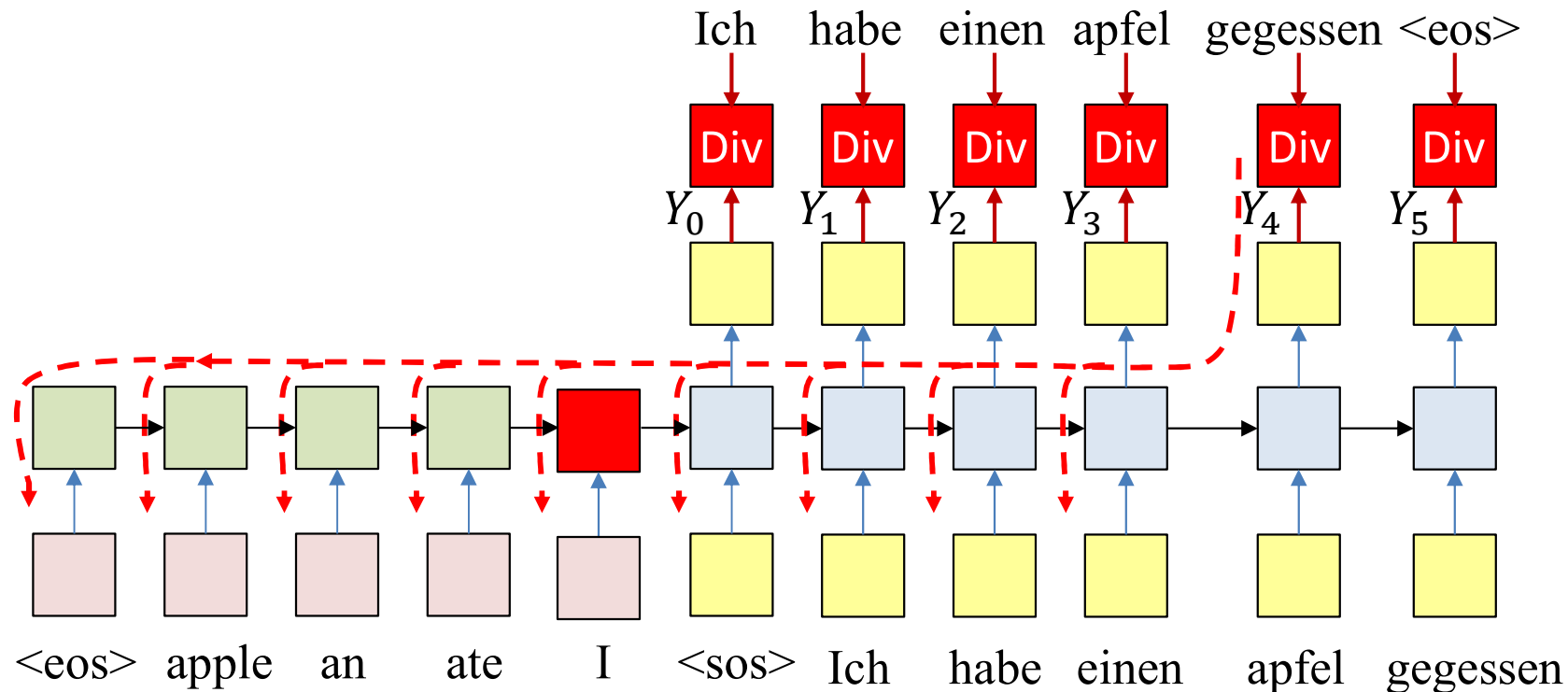
- Given several training instance (\mathbf{X}, \mathbf{D})
- For each training instance
 - Forward pass: Compute the output of the network for (\mathbf{X}, \mathbf{D})
 - Note, both \mathbf{X} *and* \mathbf{D} are used in the forward pass
 - Backward pass: Compute the divergence between selected words of the desired target \mathbf{D} and the actual output \mathbf{Y}
 - Propagate derivatives of divergence for updates
- Update parameters

Trick of the trade: Reversing the input



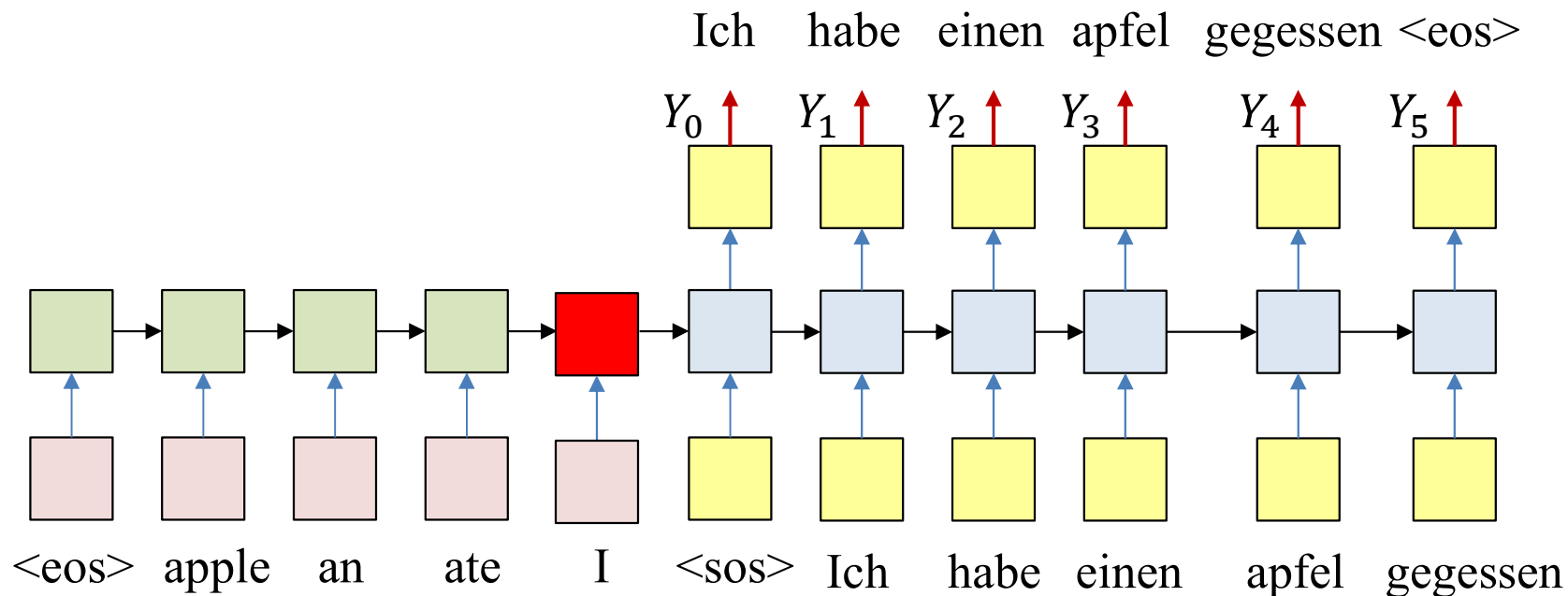
- Standard trick of the trade: The input sequence is fed *in reverse order*
 - Things work better this way

Trick of the trade: Reversing the input



- Standard trick of the trade: The input sequence is fed *in reverse order*
 - Things work better this way

Trick of the trade: Reversing the input



- Standard trick of the trade: The input sequence is fed *in reverse order*
 - Things work better this way
- *This happens both for training and during inference on test data*

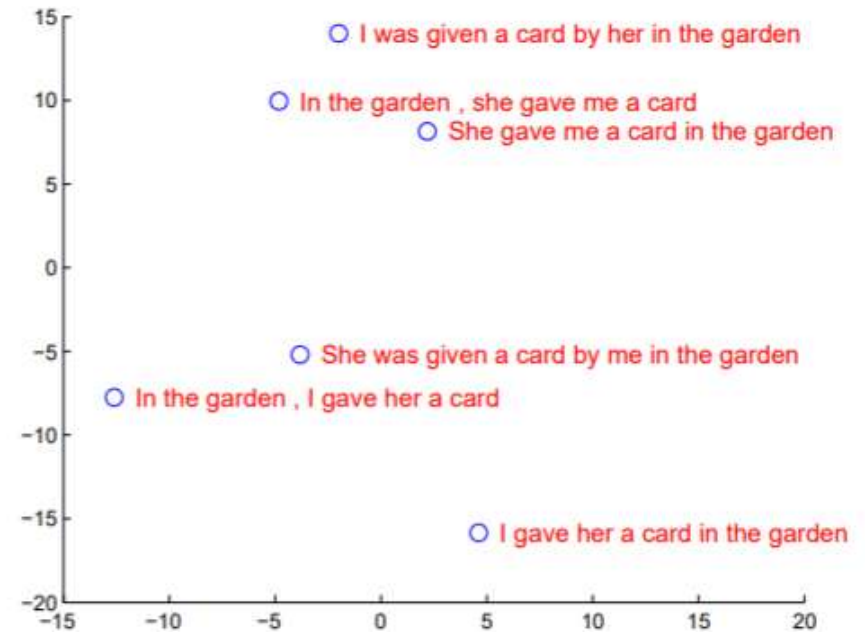
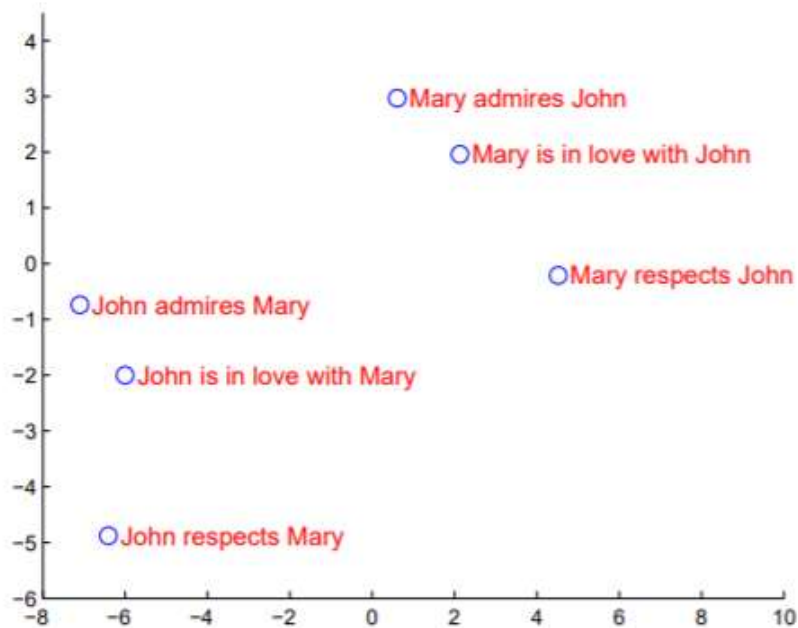
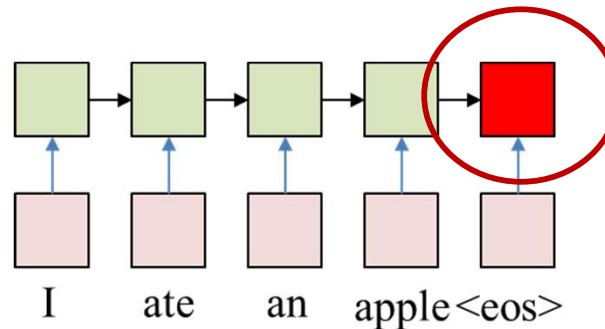
Overall training

- Given several training instance (\mathbf{X}, \mathbf{D})
- Forward pass: Compute the output of the network for (\mathbf{X}, \mathbf{D}) *with input in reverse order*
 - Note, both \mathbf{X} and \mathbf{D} are used in the forward pass
- Backward pass: Compute the divergence between the desired target \mathbf{D} and the actual output \mathbf{Y}
 - Propagate derivatives of divergence for updates

Applications

- Machine Translation
 - My name is Tom → Ich heiße Tom/Mein name ist Tom
- Automatic speech recognition
 - Speech recording → “My name is Tom”
- Dialog
 - “I have a problem” → “How may I help you”
- Image to text
 - Picture → Caption for picture

Machine Translation Example



- Hidden state clusters by meaning!
 - From “Sequence-to-sequence learning with neural networks”, Sutskever, Vinyals and Le

Machine Translation Example

| Type | Sentence |
|------------------|--|
| Our model | Ulrich UNK , membre du conseil d' administration du constructeur automobile Audi , affirme qu' il s' agit d' une pratique courante depuis des années pour que les téléphones portables puissent être collectés avant les réunions du conseil d' administration afin qu' ils ne soient pas utilisés comme appareils d' écoute à distance . |
| Truth | Ulrich Hackenberg , membre du conseil d' administration du constructeur automobile Audi , déclare que la collecte des téléphones portables avant les réunions du conseil , afin qu' ils ne puissent pas être utilisés comme appareils d' écoute à distance , est une pratique courante depuis des années . |
| Our model | “ Les téléphones cellulaires , qui sont vraiment une question , non seulement parce qu' ils pourraient potentiellement causer des interférences avec les appareils de navigation , mais nous savons , selon la FCC , qu' ils pourraient interférer avec les tours de téléphone cellulaire lorsqu' ils sont dans l' air ” , dit UNK . |
| Truth | “ Les téléphones portables sont véritablement un problème , non seulement parce qu' ils pourraient éventuellement créer des interférences avec les instruments de navigation , mais parce que nous savons , d' après la FCC , qu' ils pourraient perturber les antennes-relais de téléphonie mobile s' ils sont utilisés à bord ” , a déclaré Rosenker . |
| Our model | Avec la crémation , il y a un “ sentiment de violence contre le corps d' un être cher ” , qui sera “ réduit à une pile de cendres ” en très peu de temps au lieu d' un processus de décomposition “ qui accompagnera les étapes du deuil ” . |
| Truth | Il y a , avec la crémation , “ une violence faite au corps aimé ” , qui va être “ réduit à un tas de cendres ” en très peu de temps , et non après un processus de décomposition , qui “ accompagnerait les phases du deuil ” . |

- Examples of translation
 - From “Sequence-to-sequence learning with neural networks”, Sutskever, Vinyals and Le

Human Machine Conversation: Example

Machine: *what is the error that you are running please*

Human: *i am seeing an error related to vpn*

Machine: *what is the error message that you are getting when connecting to vpn using network connect ?*

Human: *connection refused or something like that*

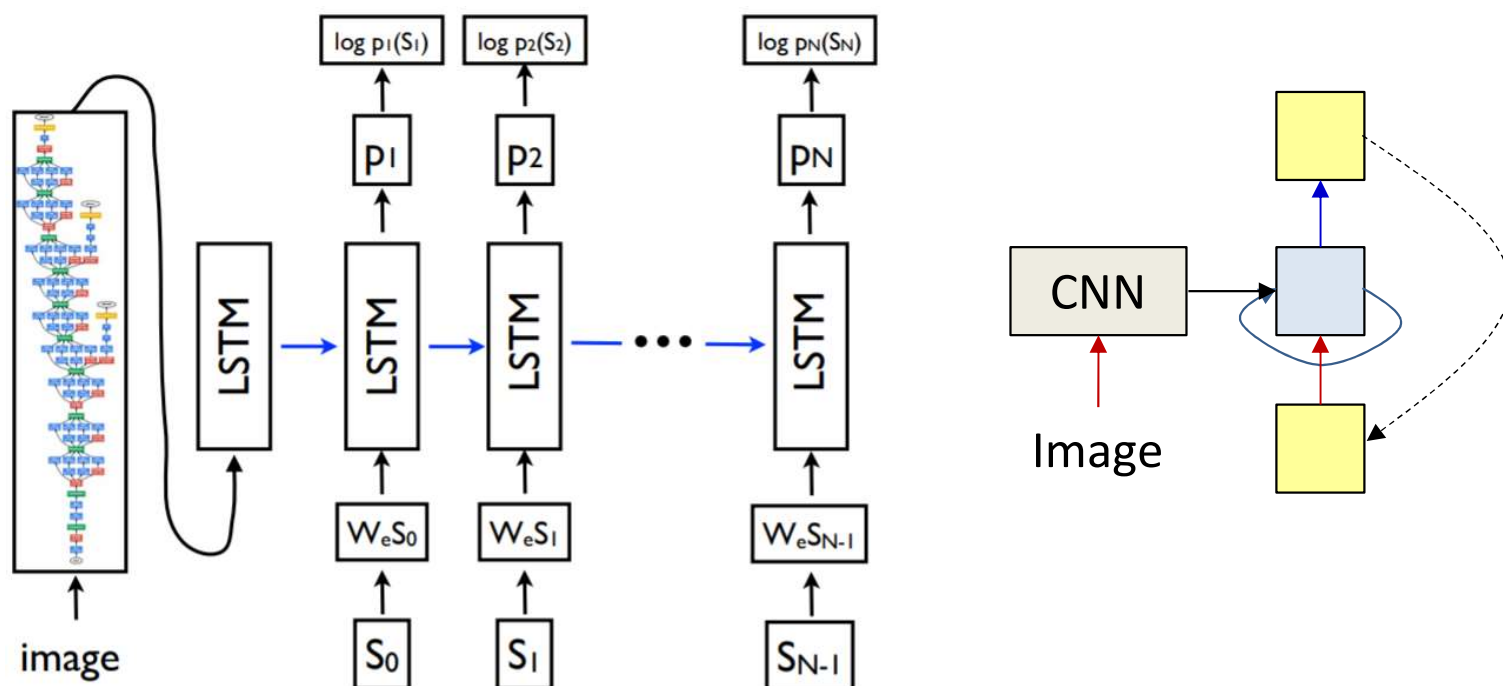
Machine: *may i know the version of network connect you connect ?*

Human: *i am not sure i know that*

Machine: *is the network connect prompting for any user certificate*

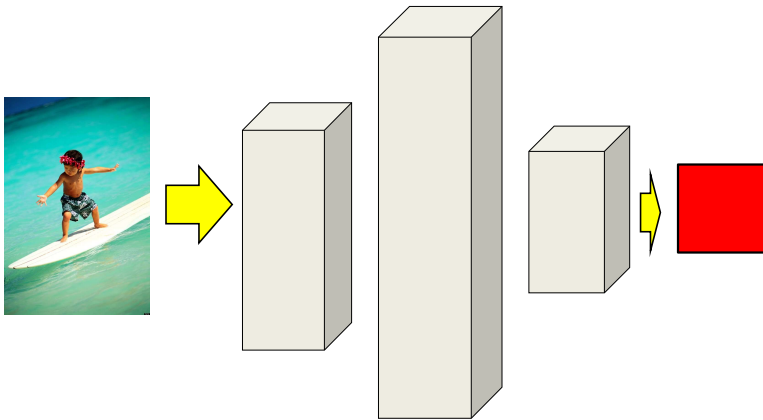
- From “A neural conversational model”, Oriin Vinyals and Quoc Le
- Trained on human-human conversations
- Task: Human text in, machine response out

Generating Image Captions



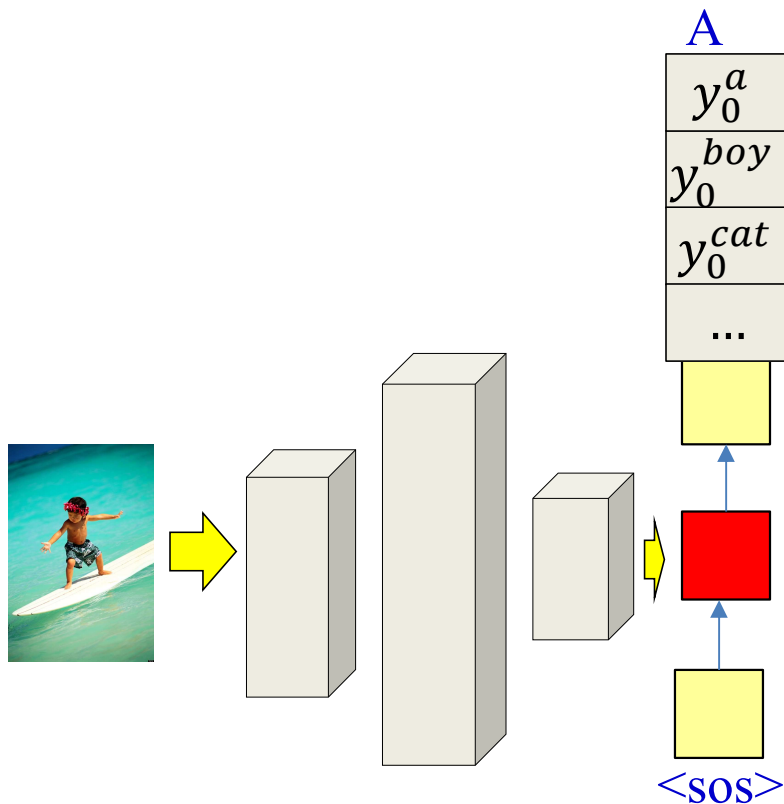
- Not really a seq-to-seq problem, more an image-to-sequence problem
- Initial state is produced by a state-of-art CNN-based image classification system
 - Subsequent model is just the decoder end of a seq-to-seq model
 - “Show and Tell: A Neural Image Caption Generator”, O. Vinyals, A. Toshev, S. Bengio, D. Erhan

Generating Image Captions



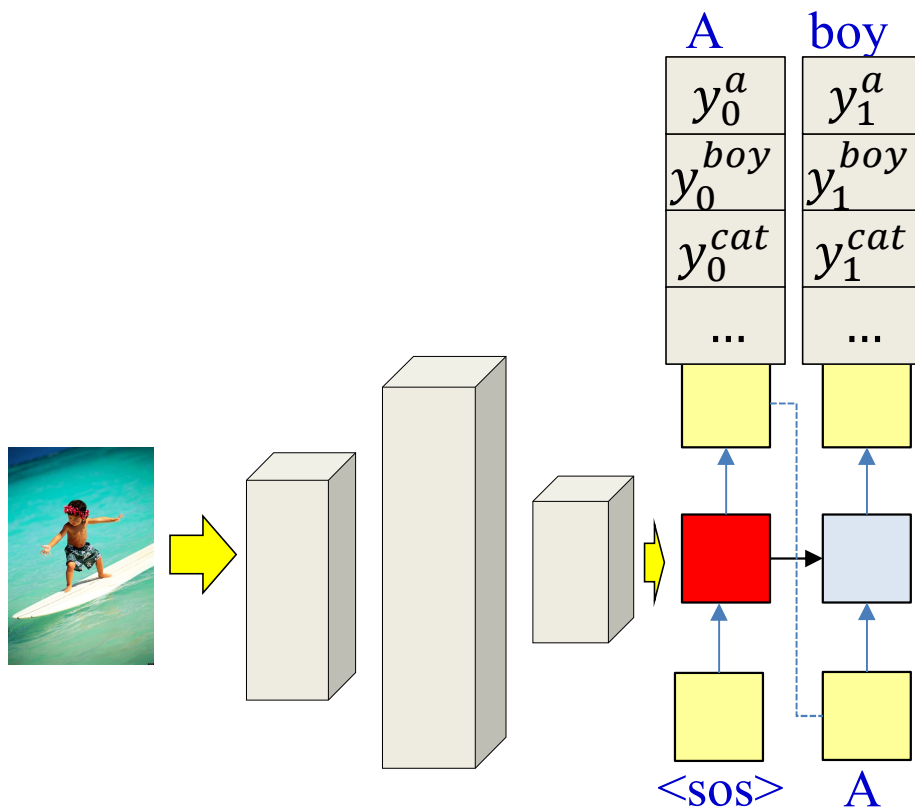
- Decoding: Given image
 - Process it with CNN to get output of classification layer

Generating Image Captions



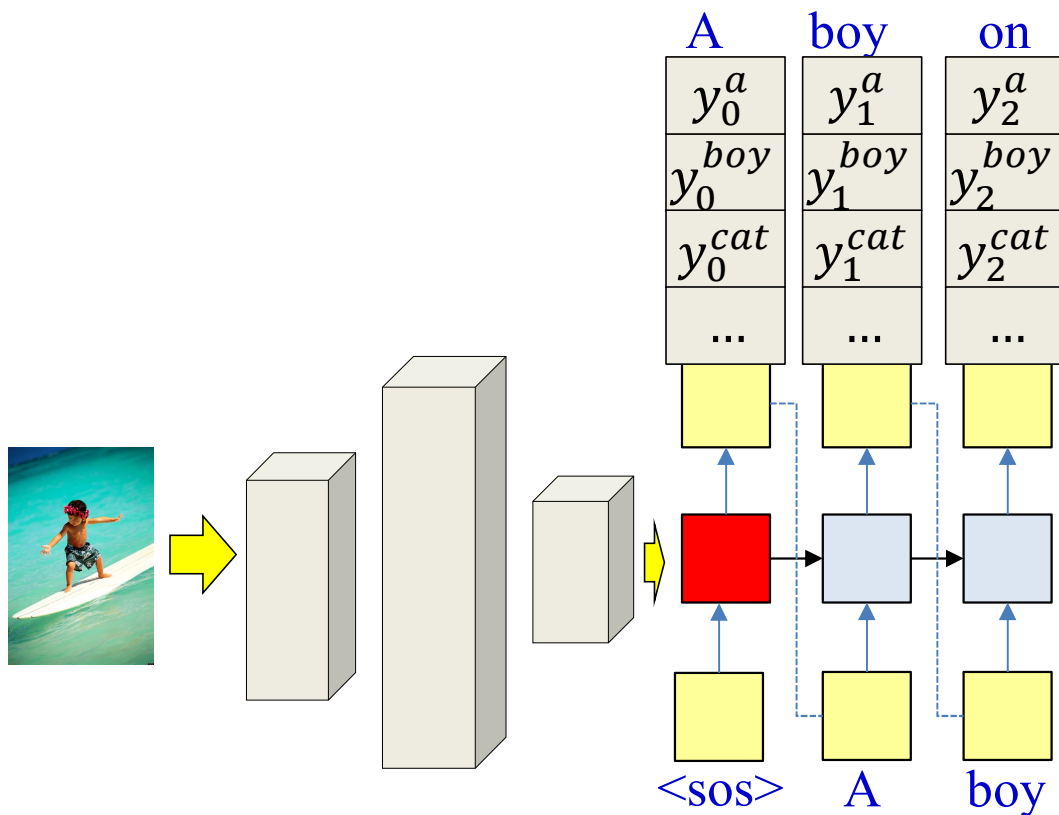
- Decoding: Given image
 - Process it with CNN to get output of classification layer
 - Sequentially generate words by drawing from the conditional output distribution $P(W_t | W_0 W_1 \dots W_{t-1}, Image)$
 - In practice, we can perform the beam search explained earlier

Generating Image Captions



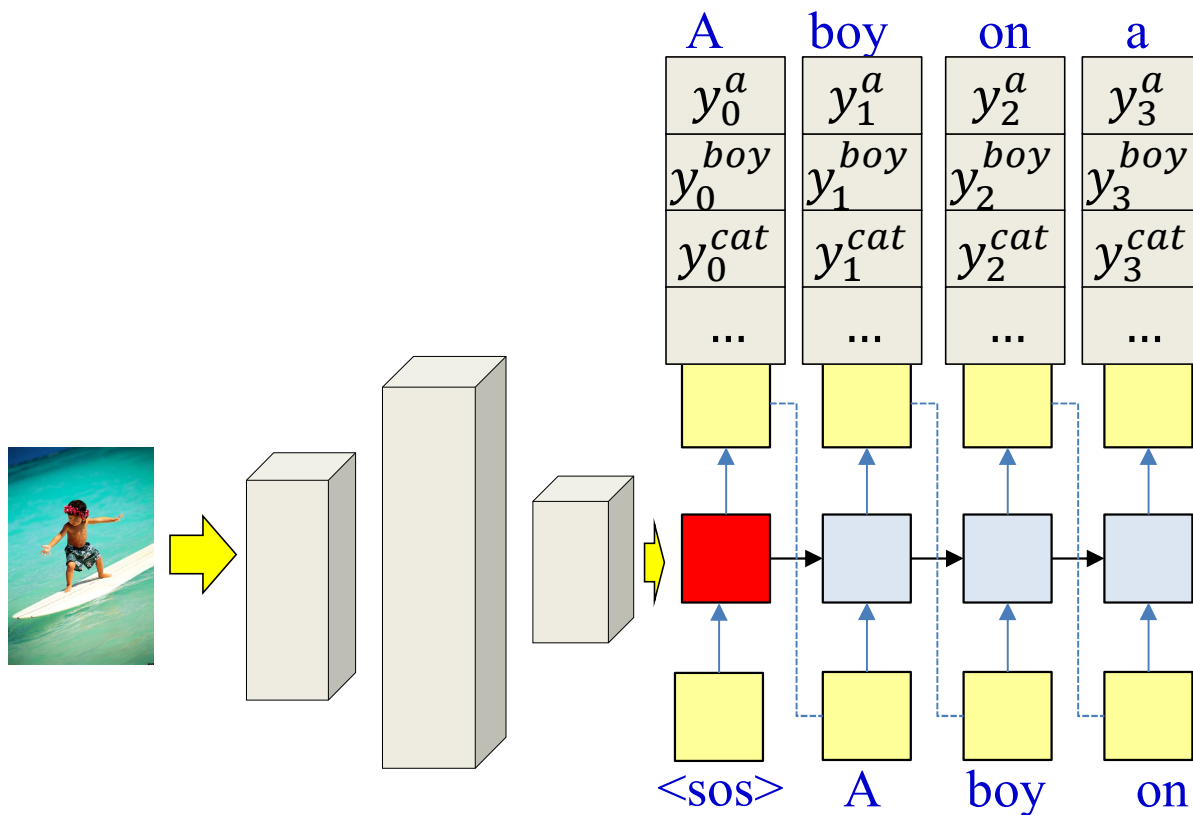
- Decoding: Given image
 - Process it with CNN to get output of classification layer
 - Sequentially generate words by drawing from the conditional output distribution $P(W_t | W_0 W_1 \dots W_{t-1}, Image)$
 - In practice, we can perform the beam search explained earlier

Generating Image Captions



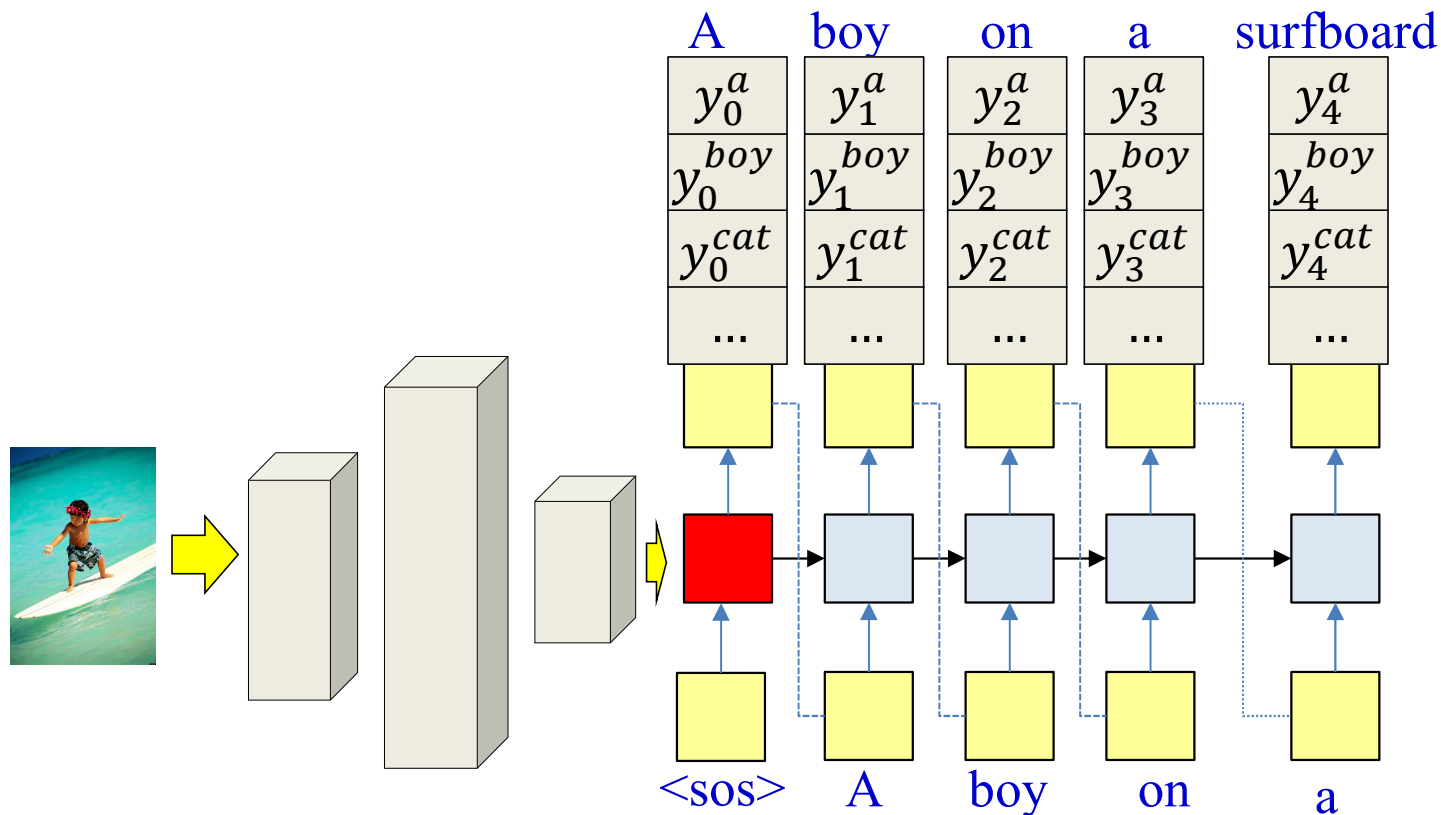
- Decoding: Given image
 - Process it with CNN to get output of classification layer
 - Sequentially generate words by drawing from the conditional output distribution $P(W_t | W_0 W_1 \dots W_{t-1}, Image)$
 - In practice, we can perform the beam search explained earlier

Generating Image Captions



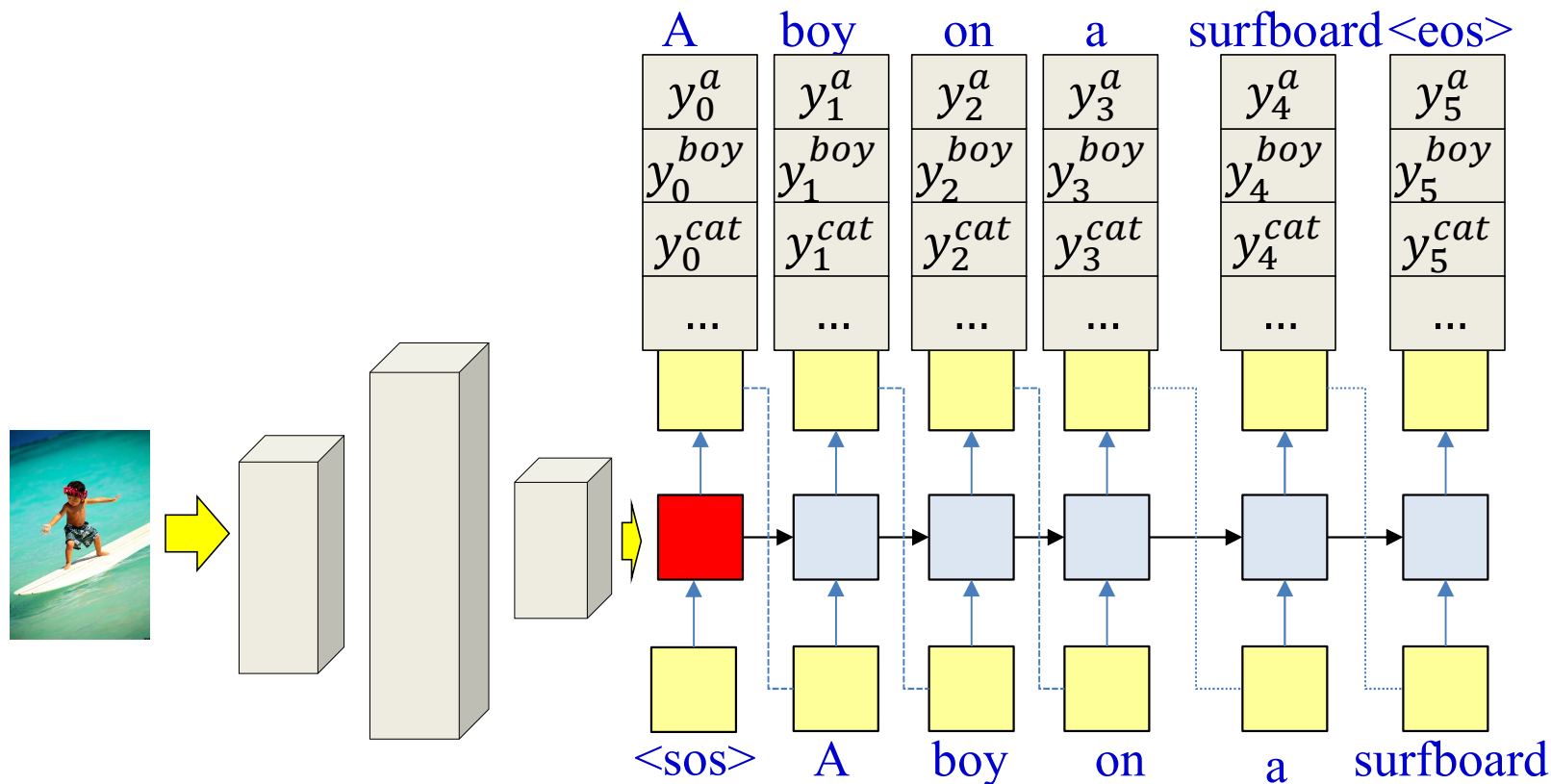
- Decoding: Given image
 - Process it with CNN to get output of classification layer
 - Sequentially generate words by drawing from the conditional output distribution $P(W_t | W_0 W_1 \dots W_{t-1}, Image)$
 - In practice, we can perform the beam search explained earlier

Generating Image Captions



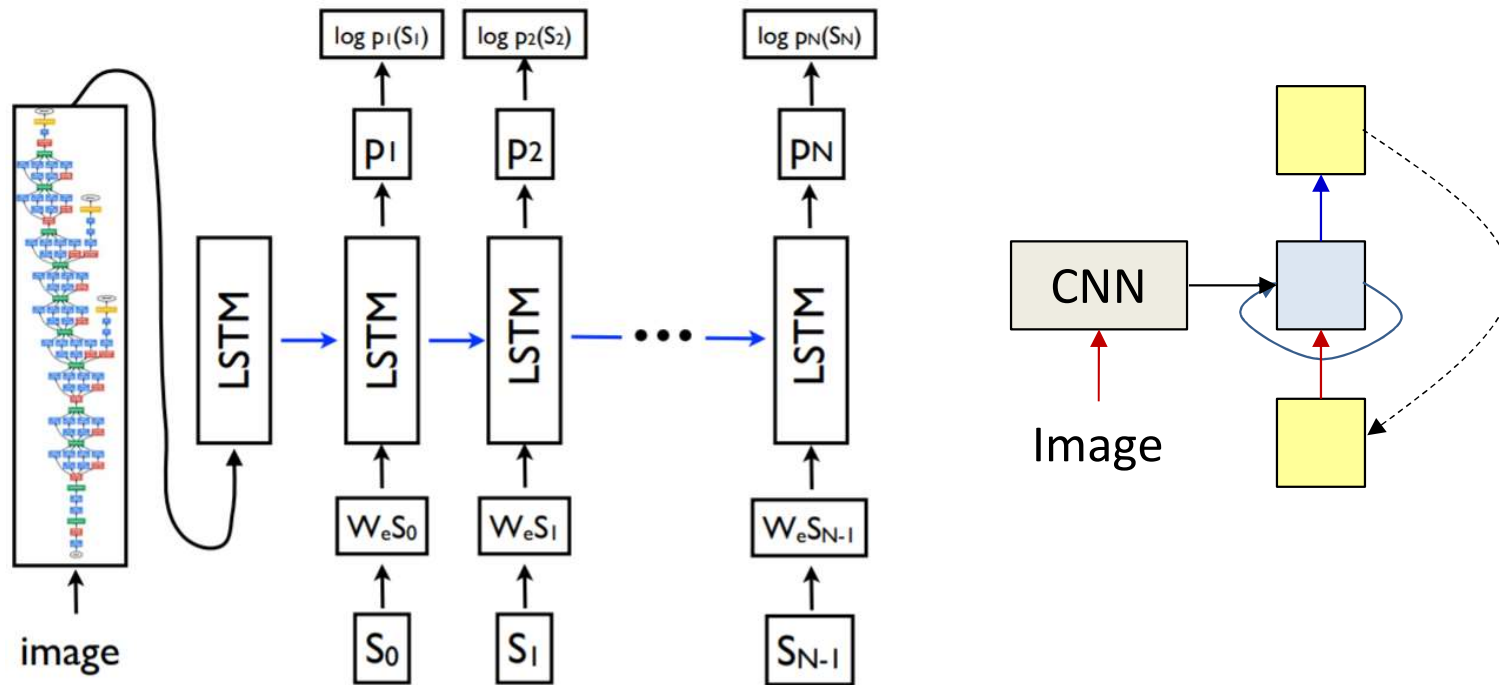
- Decoding: Given image
 - Process it with CNN to get output of classification layer
 - Sequentially generate words by drawing from the conditional output distribution $P(W_t | W_0 W_1 \dots W_{t-1}, Image)$
 - In practice, we can perform the beam search explained earlier

Generating Image Captions

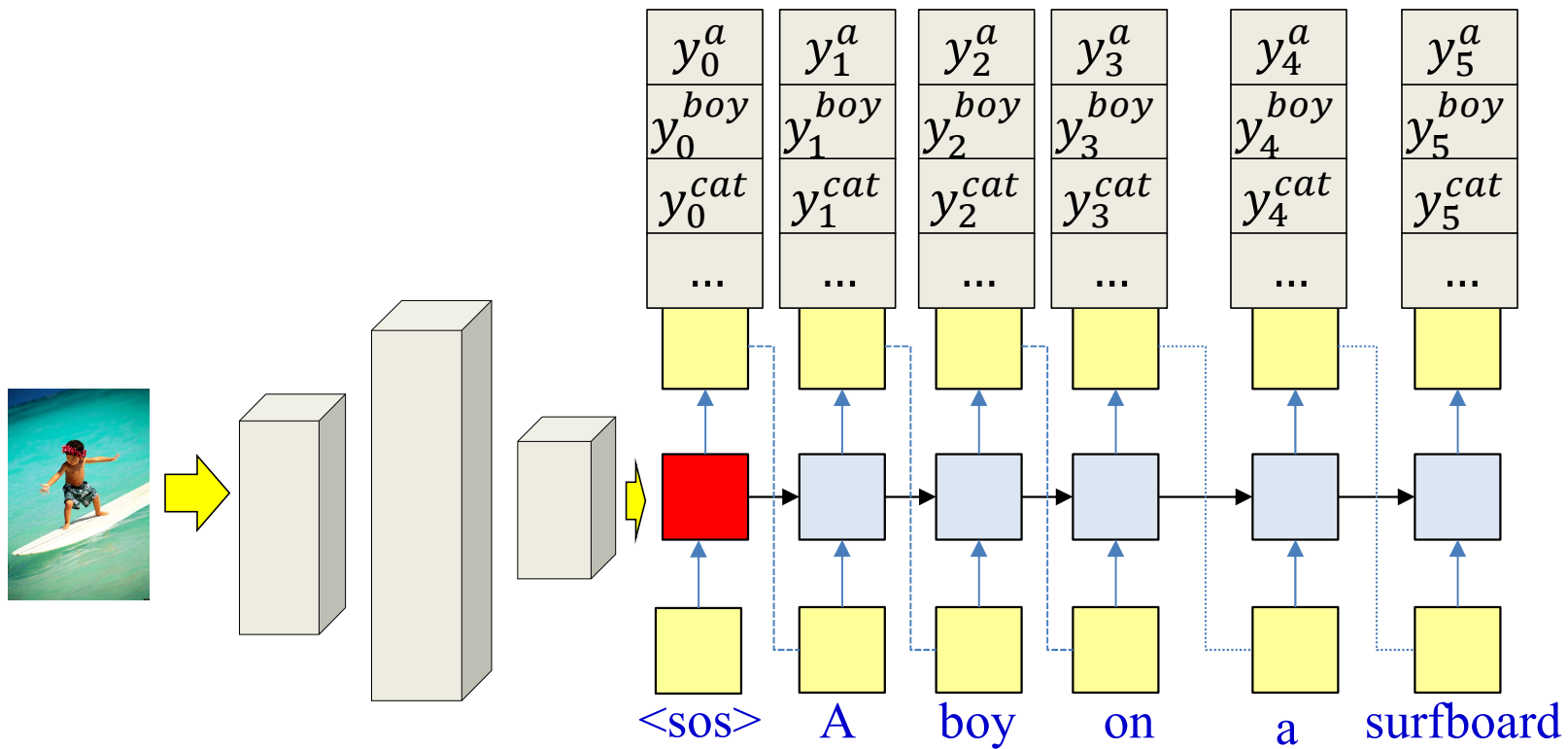


- Decoding: Given image
 - Process it with CNN to get output of classification layer
 - Sequentially generate words by drawing from the conditional output distribution $P(W_t | W_0 W_1 \dots W_{t-1}, Image)$
 - In practice, we can perform the beam search explained earlier

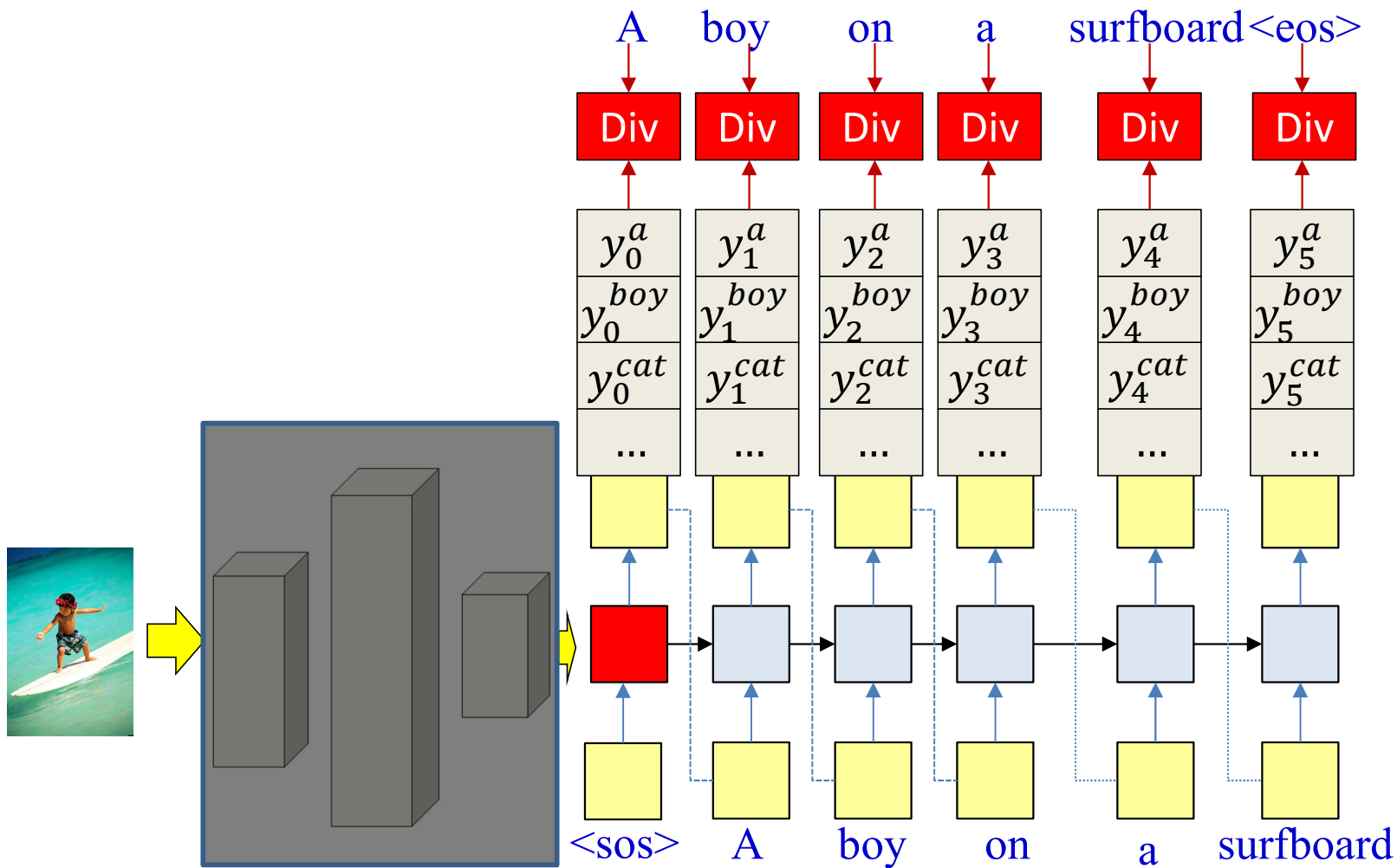
Training



- **Training:** Given several (Image, Caption) pairs
 - The image network is pretrained on a large corpus, e.g. image net



- **Training:** Given several (Image, Caption) pairs
 - The image network is pretrained on a large corpus, e.g. image net
- **Forward pass:** Produce output distributions given the image and caption



- **Training:** Given several (Image, Caption) pairs
 - The image network is pretrained on a large corpus, e.g. image net
- **Forward pass:** Produce output distributions given the image and caption
- **Backward pass:** Compute the divergence w.r.t. training caption, and backpropagate derivatives
 - All components of the network, including final classification layer of the image classification net are updated
 - The CNN portions of the image classifier are not modified (transfer learning)

Examples from Vinyals et al.



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"a young boy is holding a baseball bat."



"a cat is sitting on a couch with a remote control."



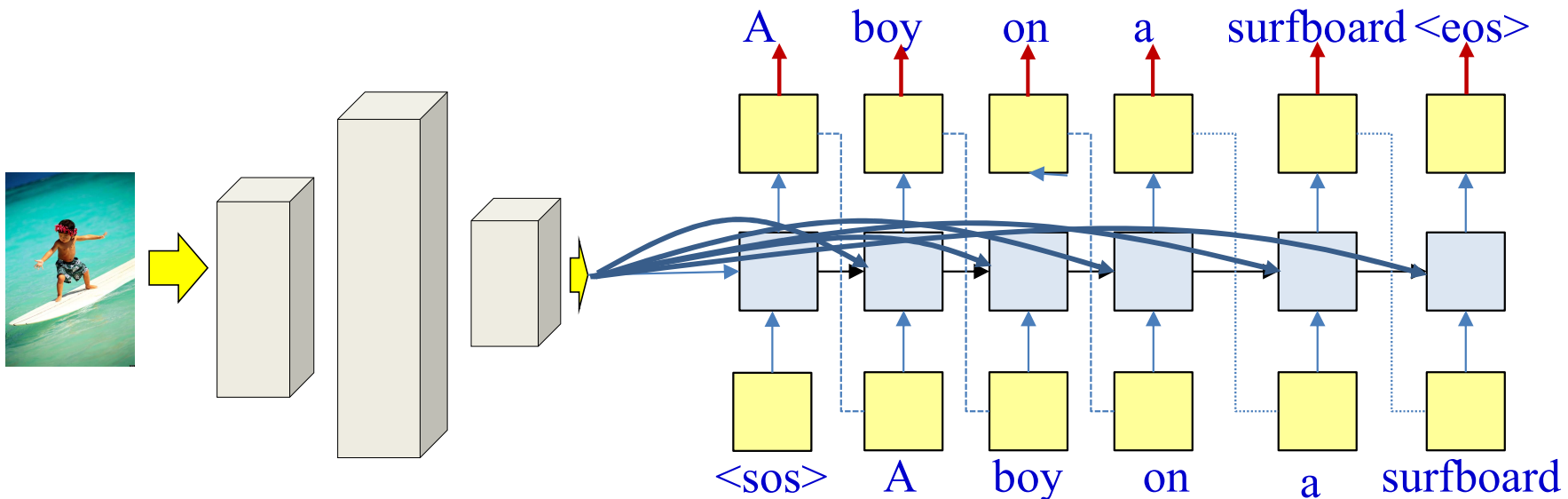
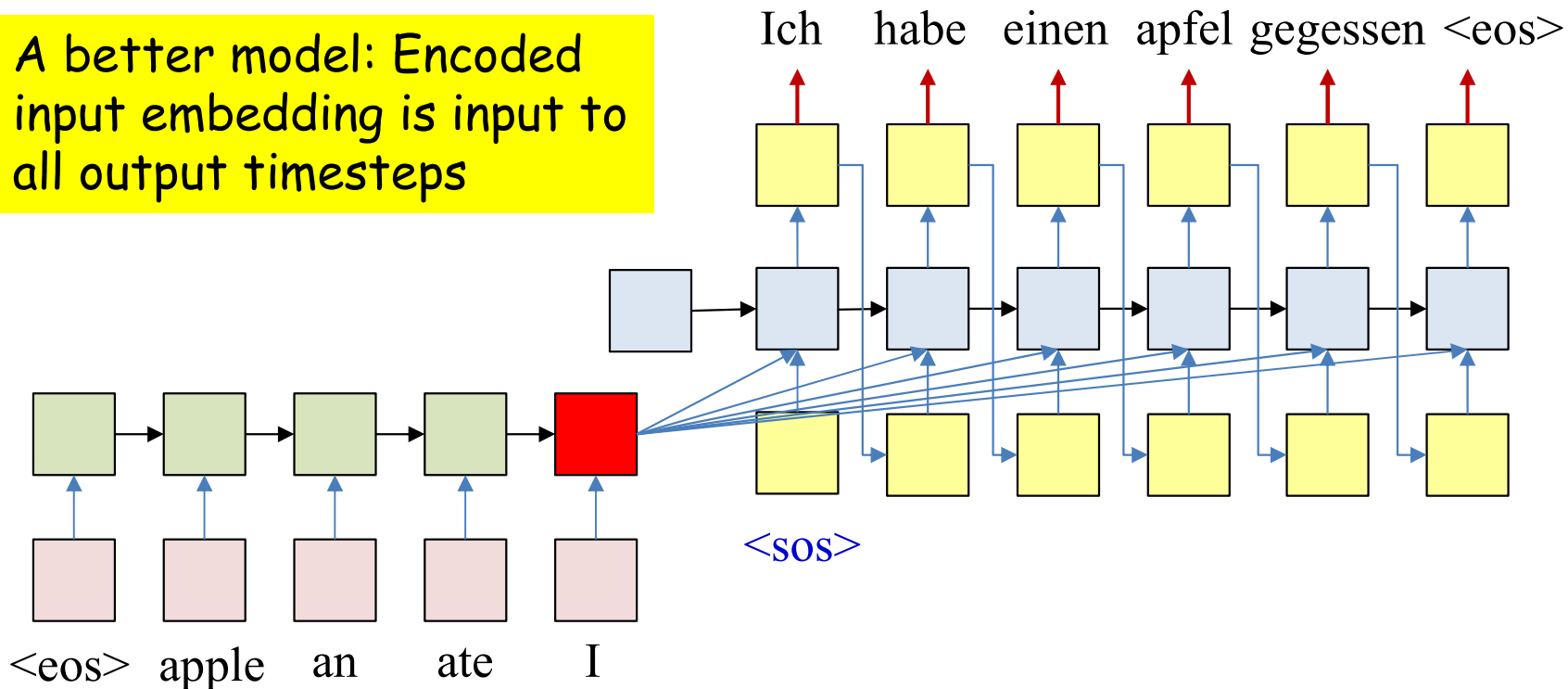
"a woman holding a teddy bear in front of a mirror."



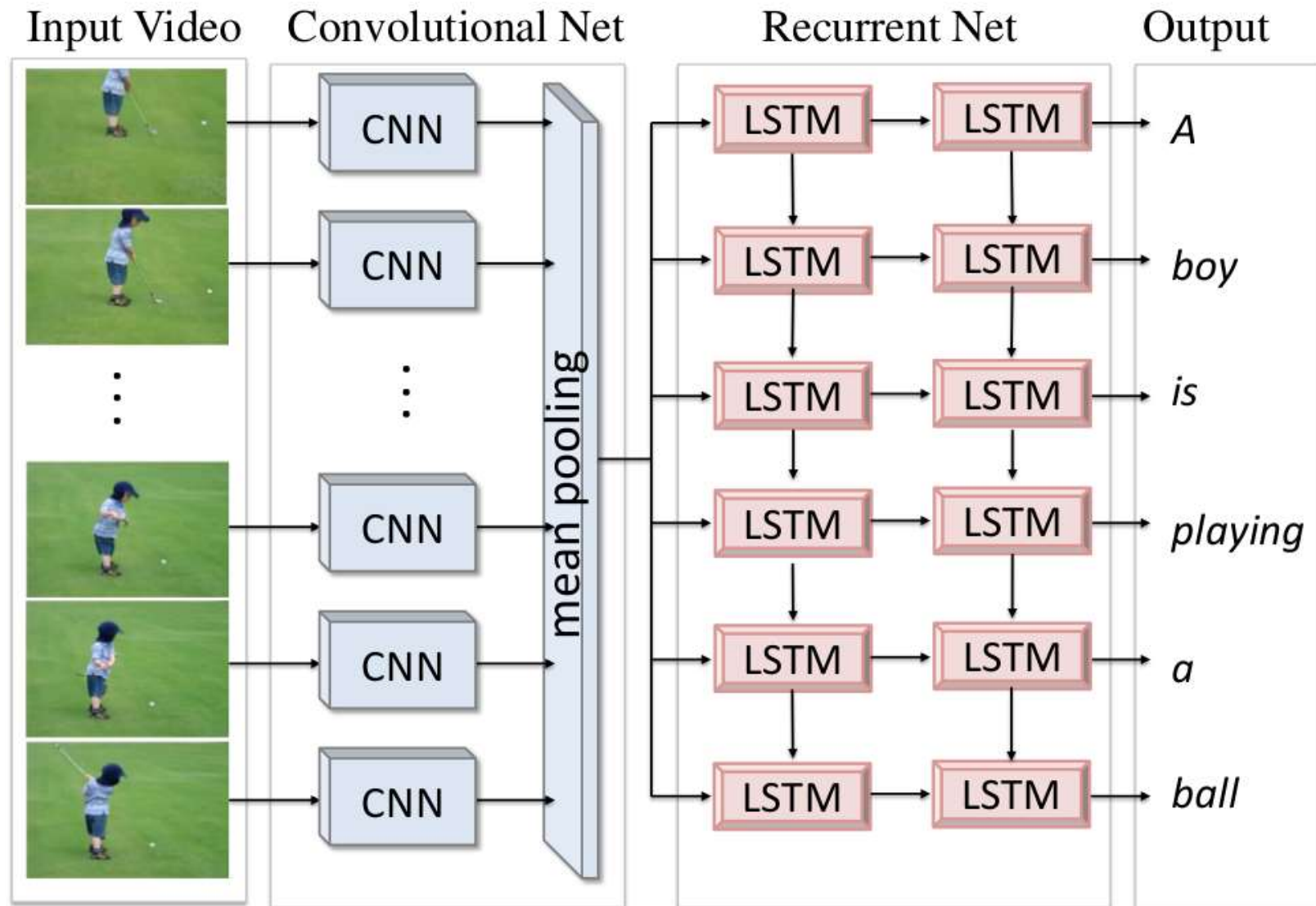
"a horse is standing in the middle of a road." ¹⁰⁸

Variants

A better model: Encoded input embedding is input to all output timesteps



Translating Videos to Natural Language Using Deep Recurrent Neural Networks



Translating Videos to Natural Language Using Deep Recurrent Neural Networks

Subhashini Venugopalan, Huijun Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, Kate Saenko¹¹⁰
North American Chapter of the Association for Computational Linguistics, Denver, Colorado, June 2015.

Pseudocode

```
# Assuming encoded input  $H$  (from text, image, video)
# is available
# Now generate the output  $y_{out}(1), y_{out}(2), \dots$ 
 $t = 0$ 
 $h_{out}(0) = H$     # Encoder embedding


# Note: begins with a "start of sentence" symbol
#      <sos> and <eos> may be identical
 $y_{out}(0) = \text{<sos>}$ 
do
     $t = t + 1$ 
     $[y(t), h_{out}(t)] = \text{RNN\_output\_step}(h_{out}(t-1), y_{out}(t-1), H)$ 
     $y_{out}(t) = \text{generate}(y(t))$  # Beam search, random, or greedy
until  $y_{out}(t) == \text{<eos>}$ 
```

Pseudocode

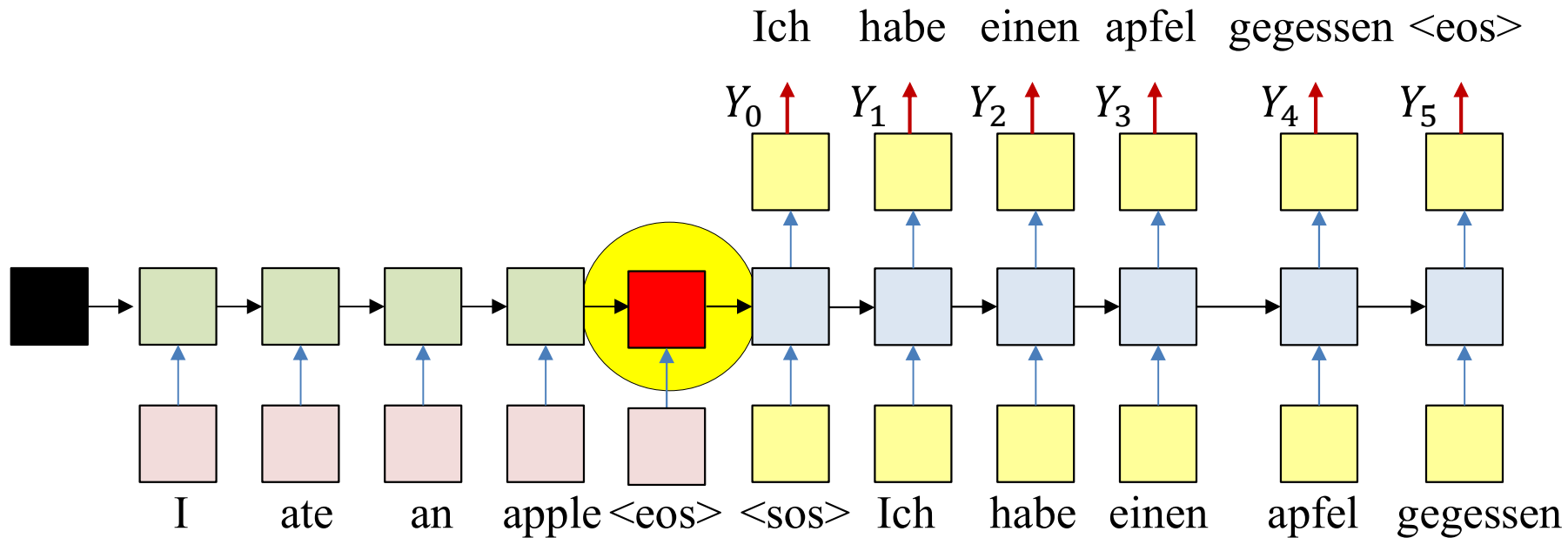
```
# Assuming encoded input  $H$  (from text, image, video)
# is available
# Now generate the output  $y_{out}(1), y_{out}(2), \dots$ 
 $t = 0$ 
 $h_{out}(0) = H$     # Encoder embedding

# Note: begins with a "start of sentence" symbol
#      <sos> and <eos> may be identical
 $y_{out}(0) = \text{<sos>}$ 
do
     $t = t + 1$ 
     $[y(t), h_{out}(t)] = \text{RNN\_output\_step}(h_{out}(t-1), y_{out}(t-1), H)$ 
     $y_{out}(t) = \text{generate}(y(t))$  # Beam search, random, or greedy
until  $y_{out}(t) == \text{<eos>}$ 
```

Also consider
encoder embedding

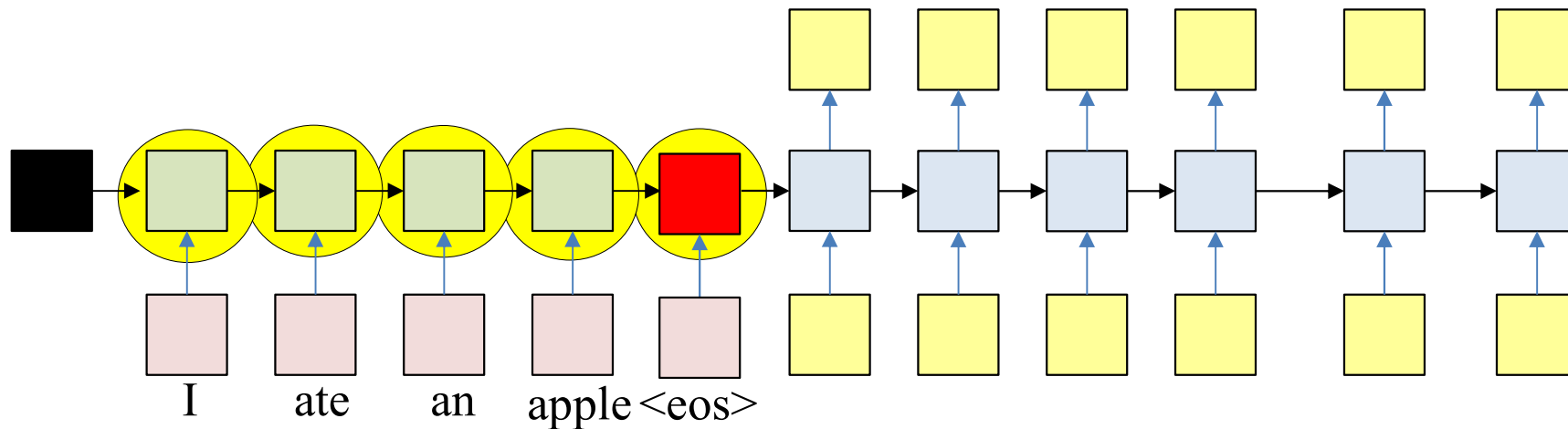


A problem with this framework



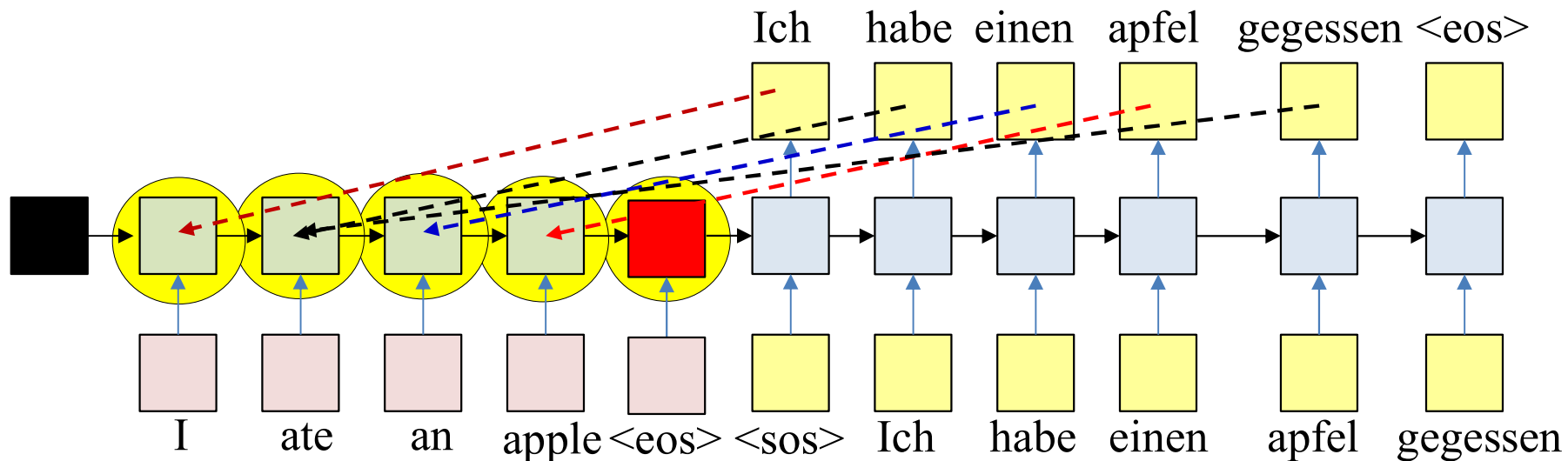
- All the information about the input sequence is embedded into a *single* vector
 - The “hidden” node layer at the end of the input sequence
 - This one node is “overloaded” with information
 - Particularly if the input is long

A problem with this framework



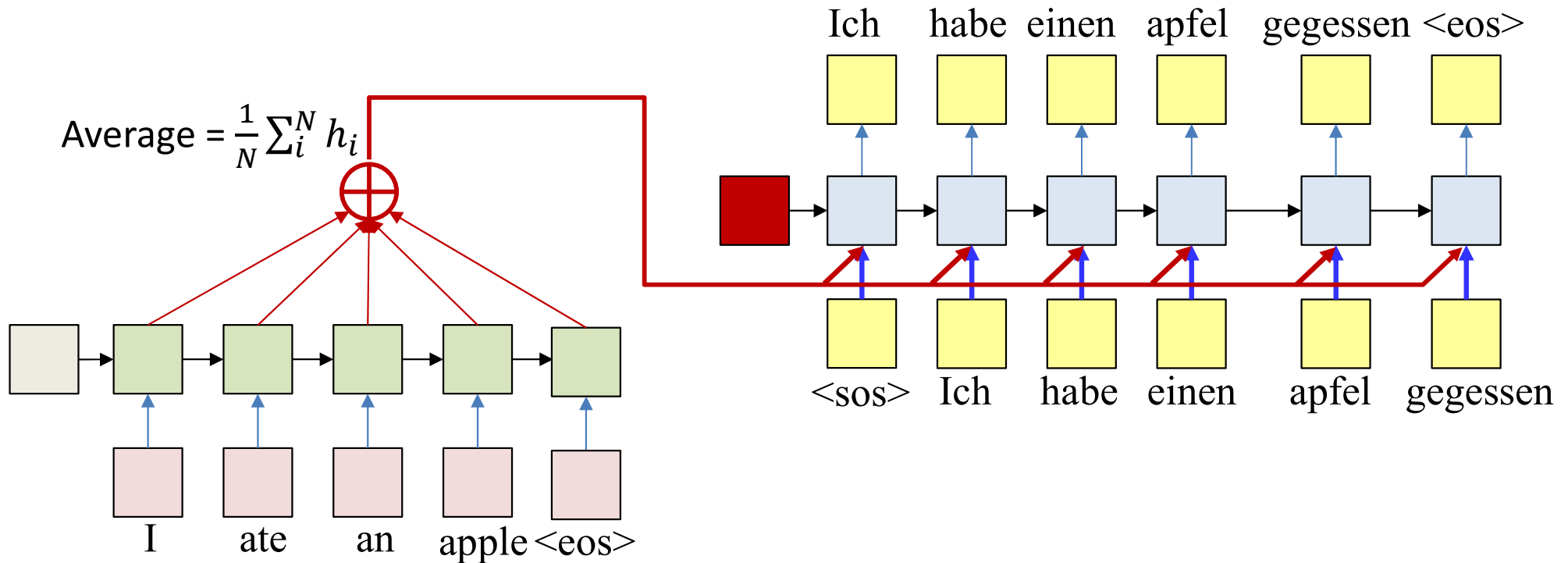
- In reality: *All* hidden values carry information
 - Some of which may be diluted by the time we get to the final state of the encoder

A problem with this framework



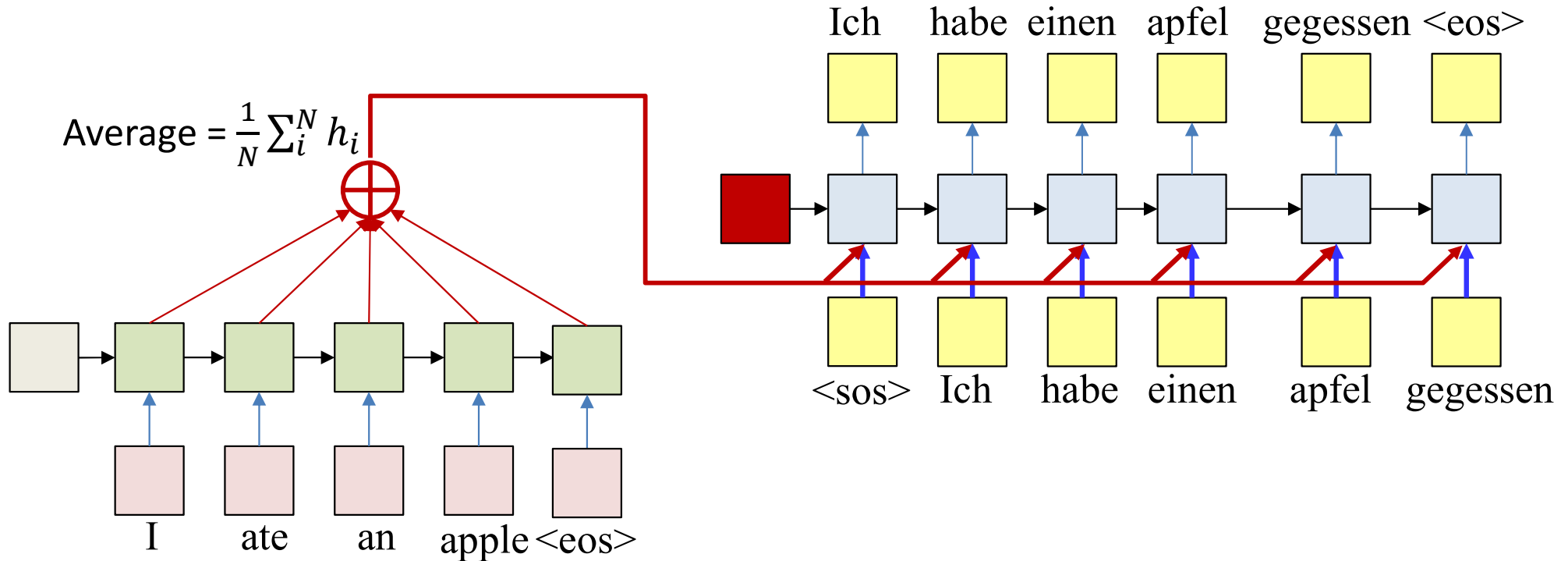
- In reality: *All* hidden values carry information
 - Some of which may be diluted by the time we get to the final state of the encoder
- *Every* output is related to the input directly
 - Not sufficient to have the encoder hidden state to *only* the initial state of the decoder
 - Misses the direct relation of the outputs to the inputs

Using all input hidden states



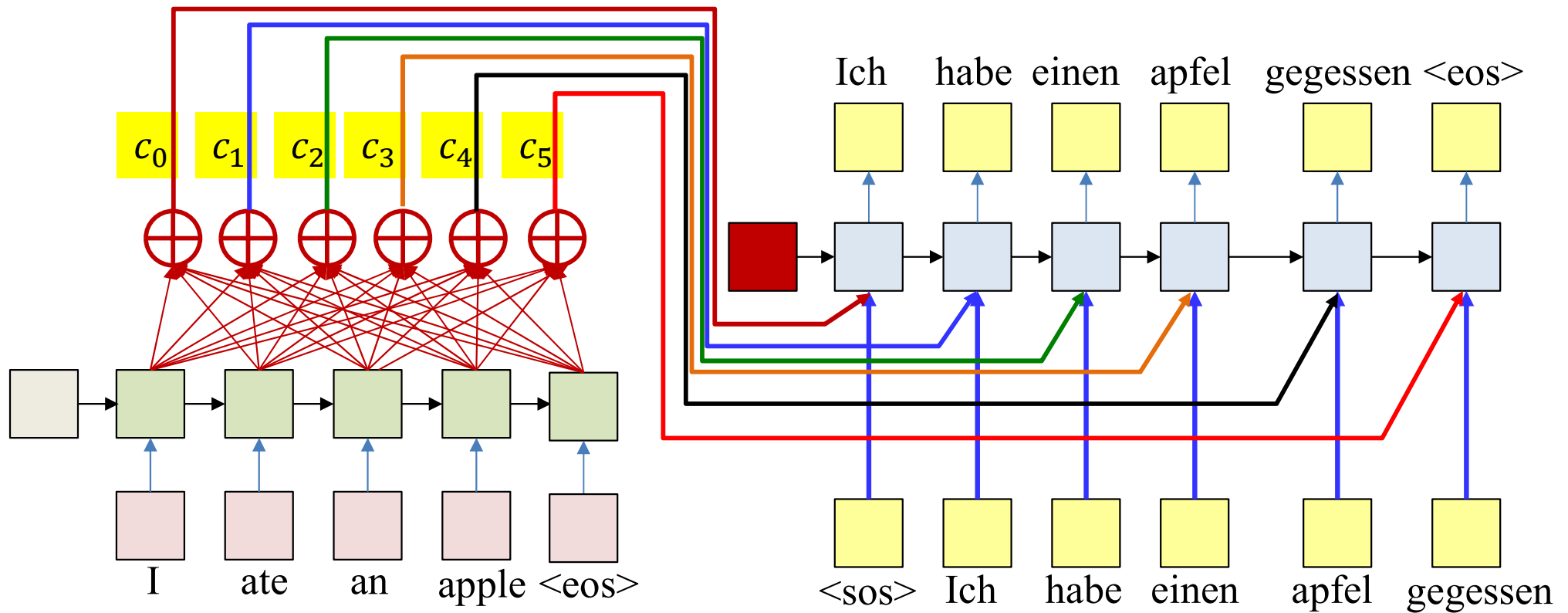
- Simple solution: Compute the average of all encoder hidden states
- Input this average to every stage of the decoder
- The initial decoder hidden state is now separate from the encoder
 - And may be a learnable parameter

Using all input hidden states



- **Problem:** The average applies the same weight to every input
- It supplies the same average to every output word
- In practice, different outputs may be related to different inputs
 - E.g. "Ich" is most related to "I", and "habe" and "gegessen" are both most related to "ate"

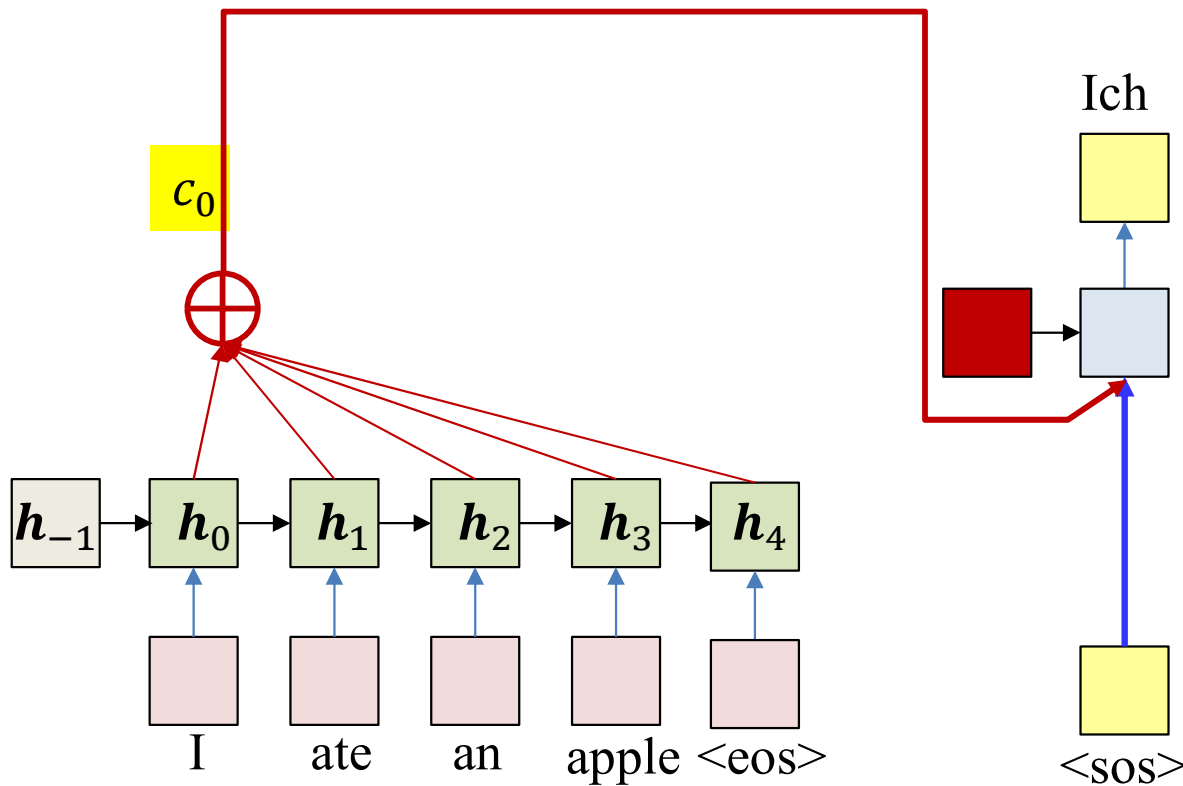
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_t = \frac{1}{N} \sum_i^N w_i(t) h_i$$

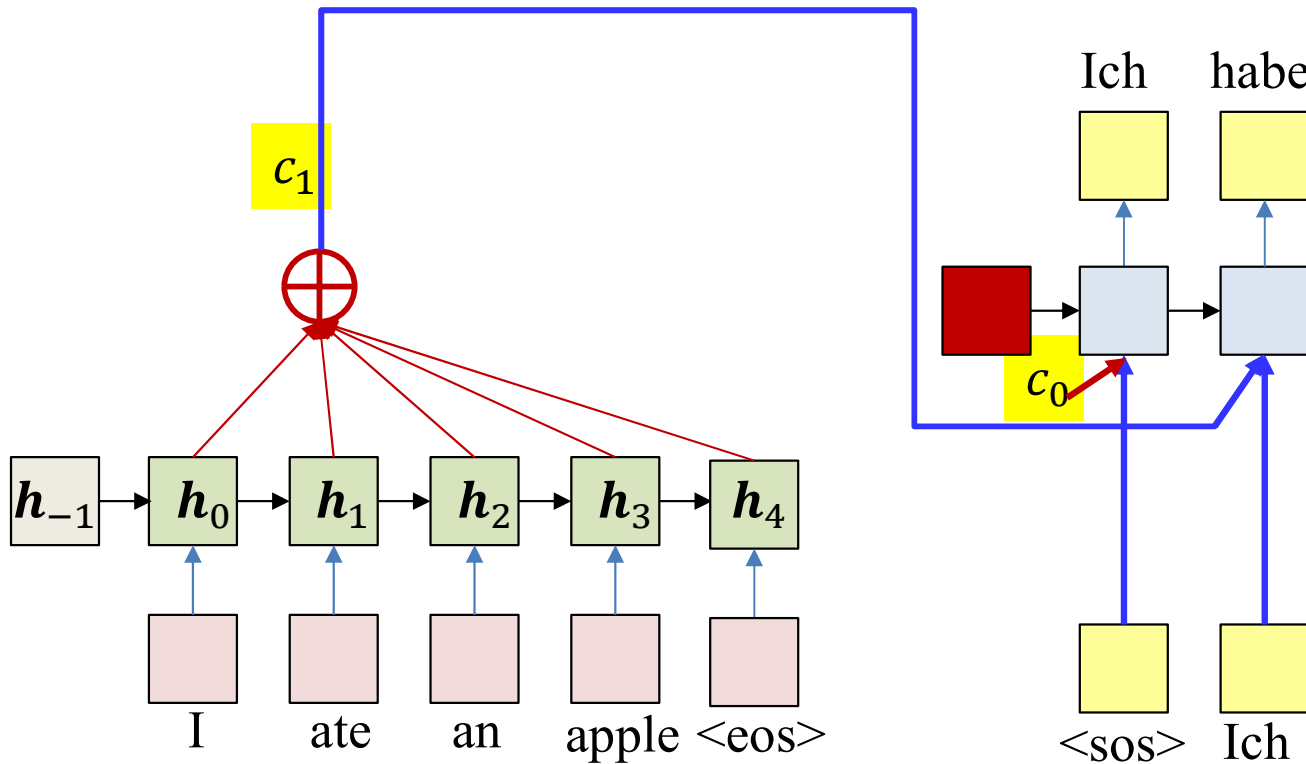
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_0 = \frac{1}{N} \sum_i^N w_i(0) h_i$$

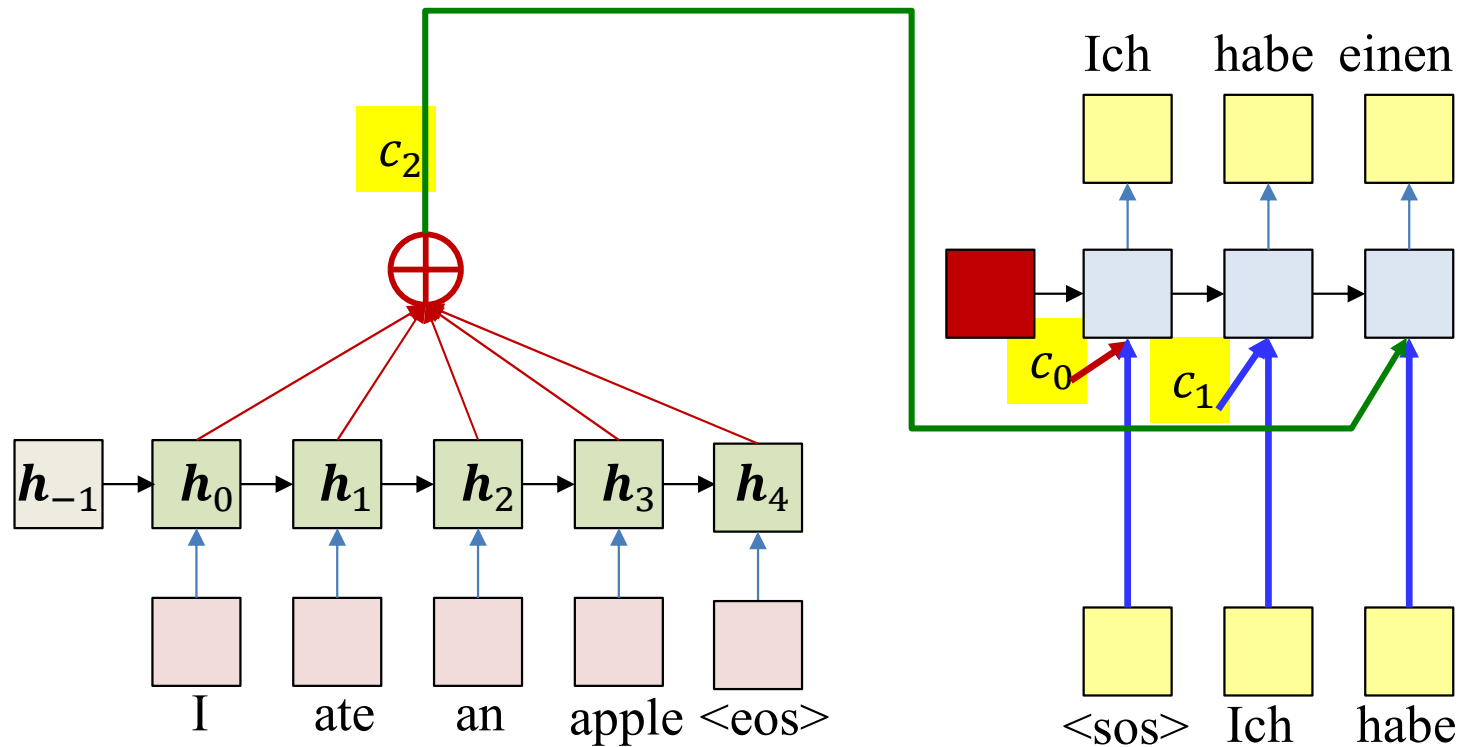
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_1 = \frac{1}{N} \sum_i^N w_i(1) h_i$$

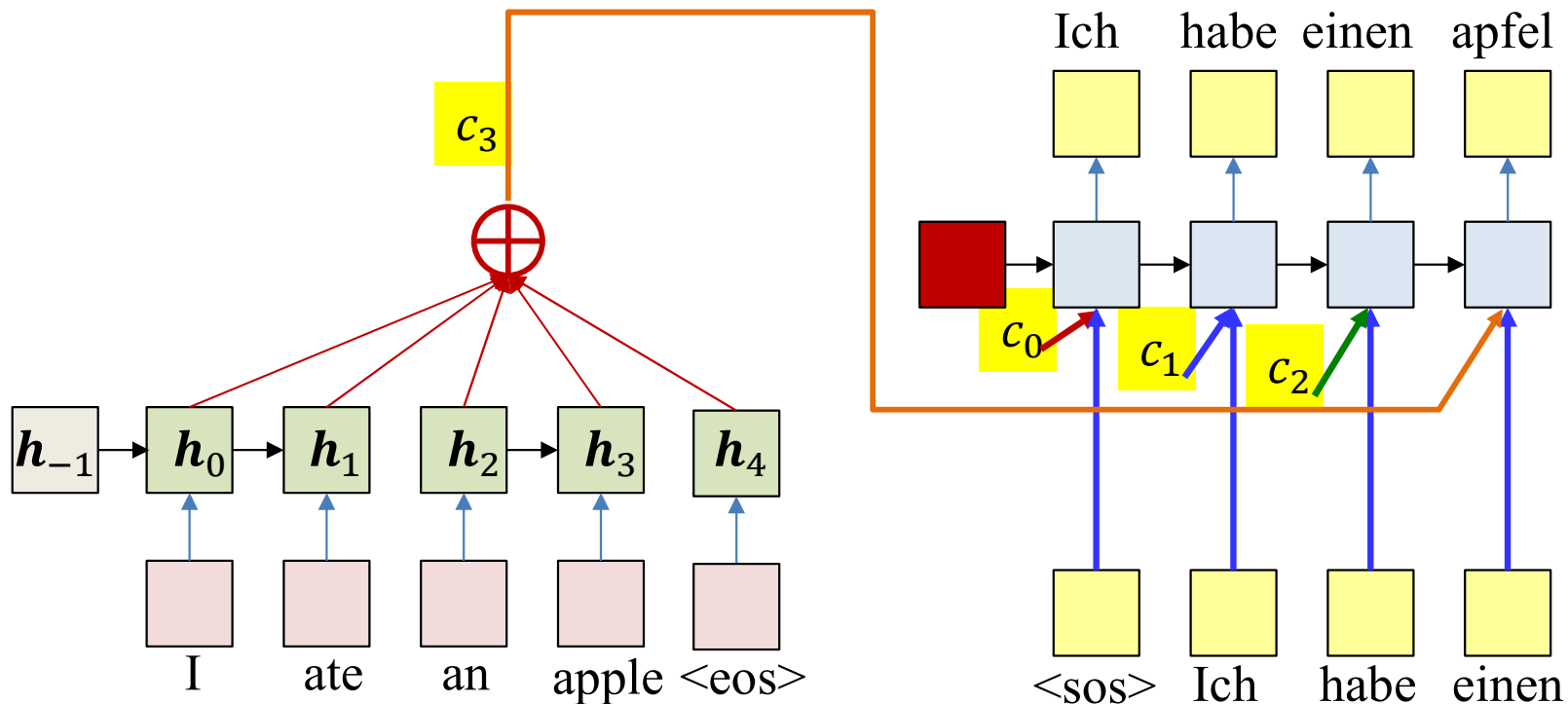
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_2 = \frac{1}{N} \sum_i^N w_i(2) h_i$$

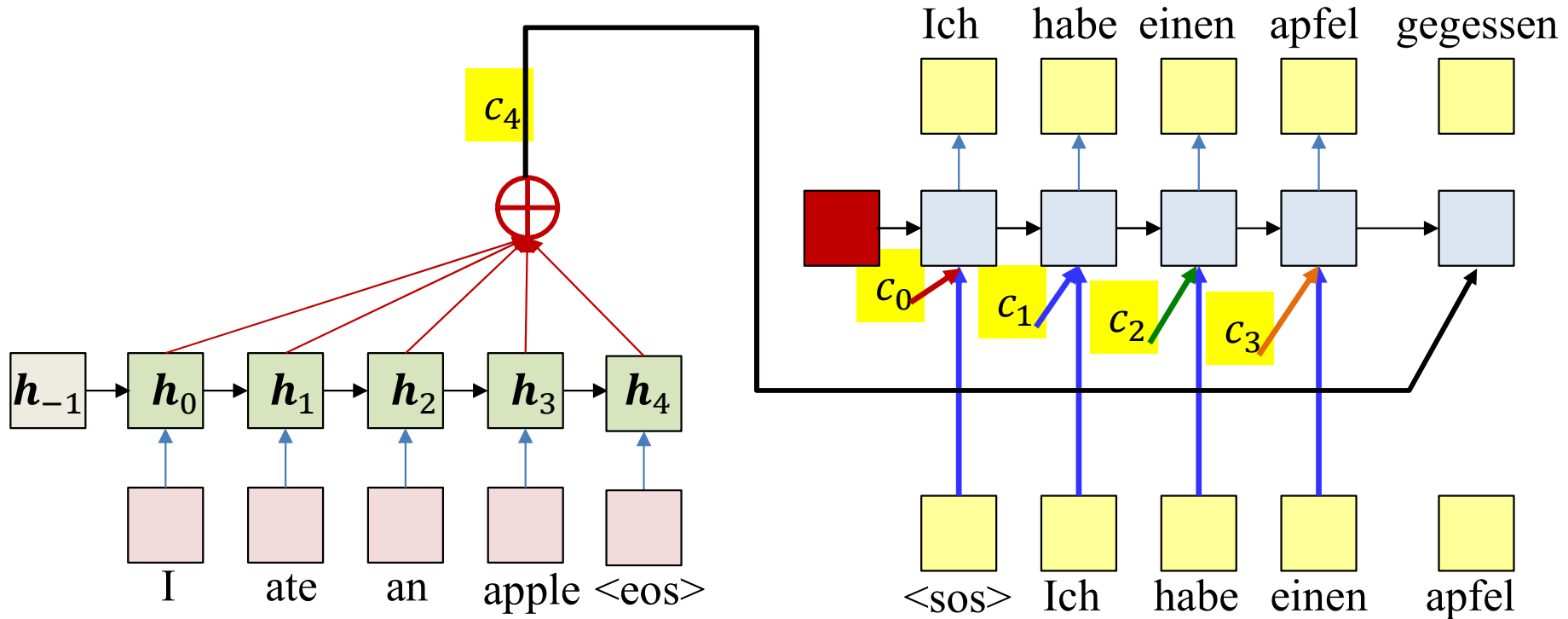
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_3 = \frac{1}{N} \sum_i^N w_i(3) h_i$$

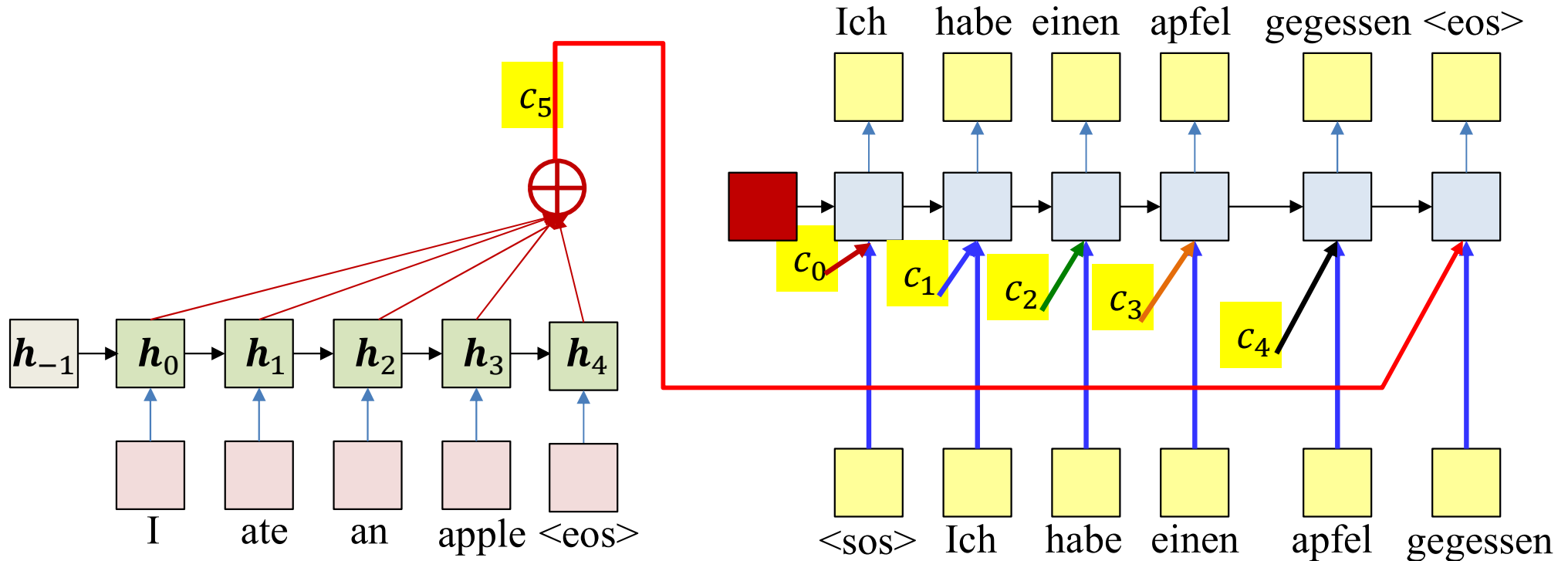
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_k = \frac{1}{N} \sum_i^N w_i(k) h_i$$

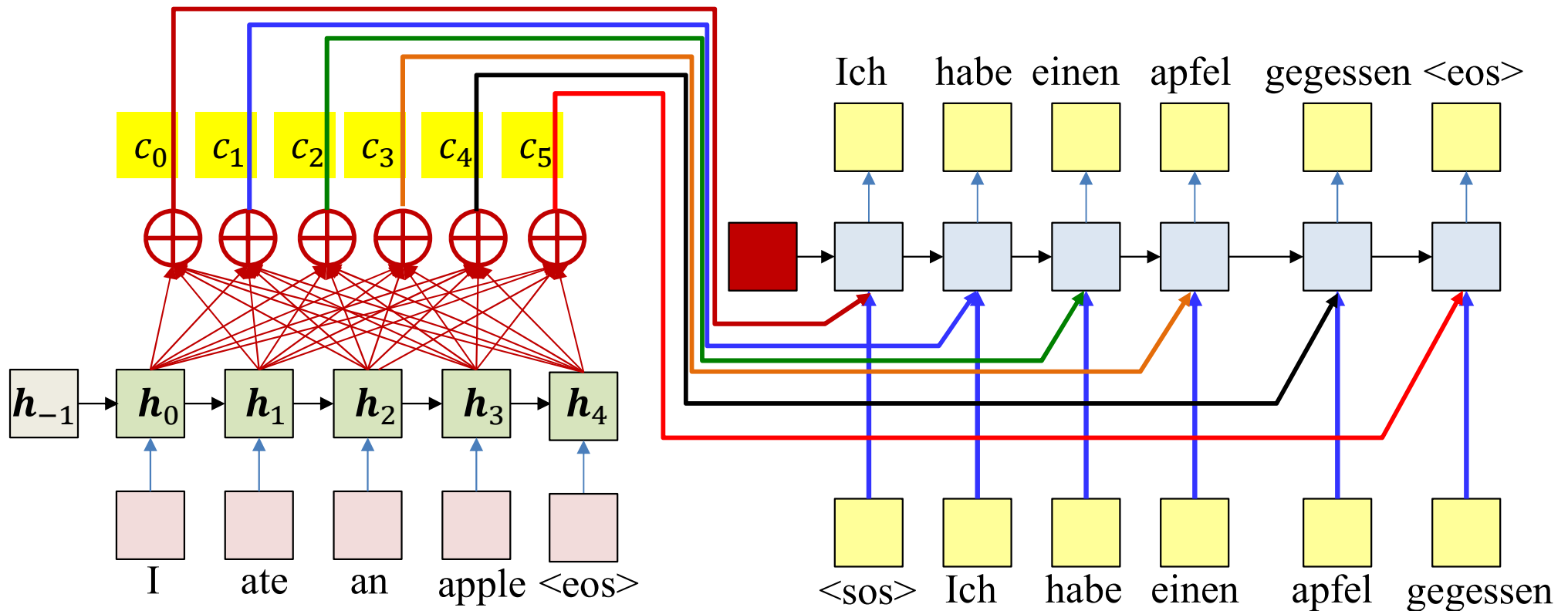
Using all input hidden states



- **Solution:** Use a *different* weighted average for each output word
 - The weighted average provided for the kth output word is:

$$c_5 = \frac{1}{N} \sum_i^N w_i(5) h_i$$

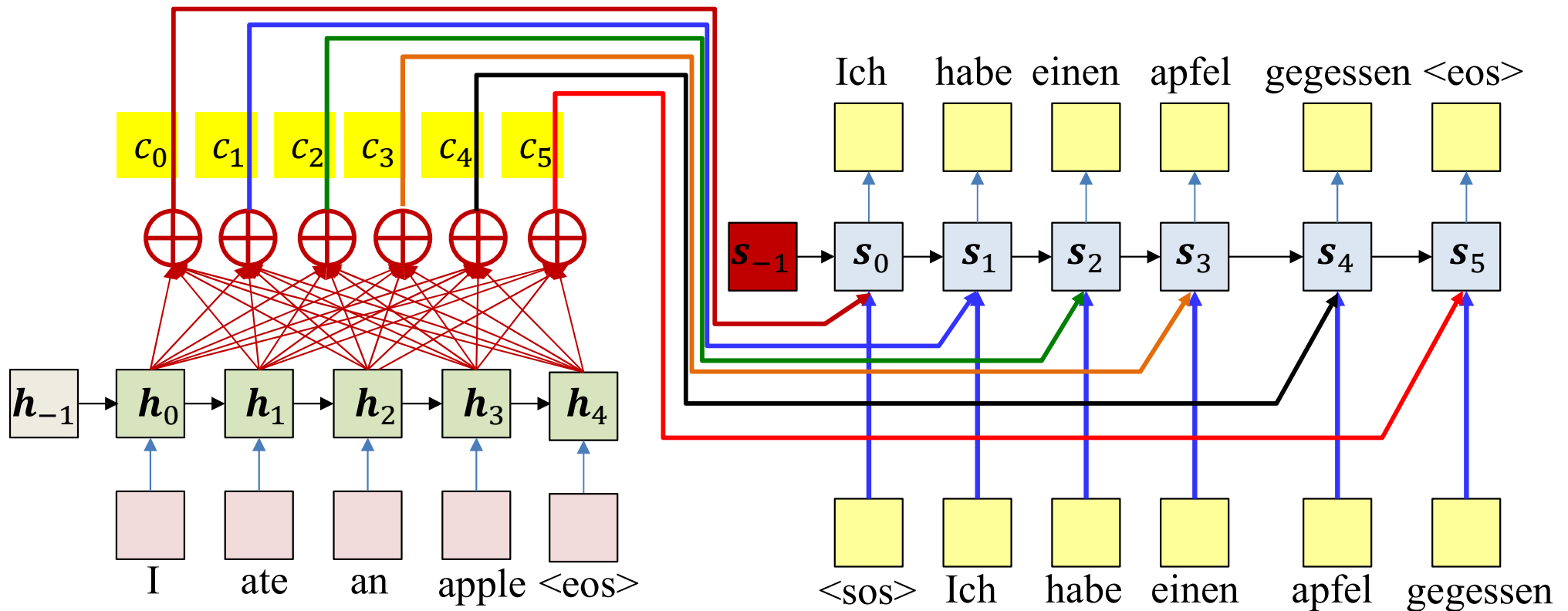
Using all input hidden states



$$c_t = \frac{1}{N} \sum_i^N w_i(t) h_i$$

- This solution will work if the weights w_{ki} can somehow be made to “focus” on the right input word
 - E.g., when predicting the word “apfel”, $w_3(4)$, the weight for “apple” must be high while the rest must be low
- How do we generate such weights??

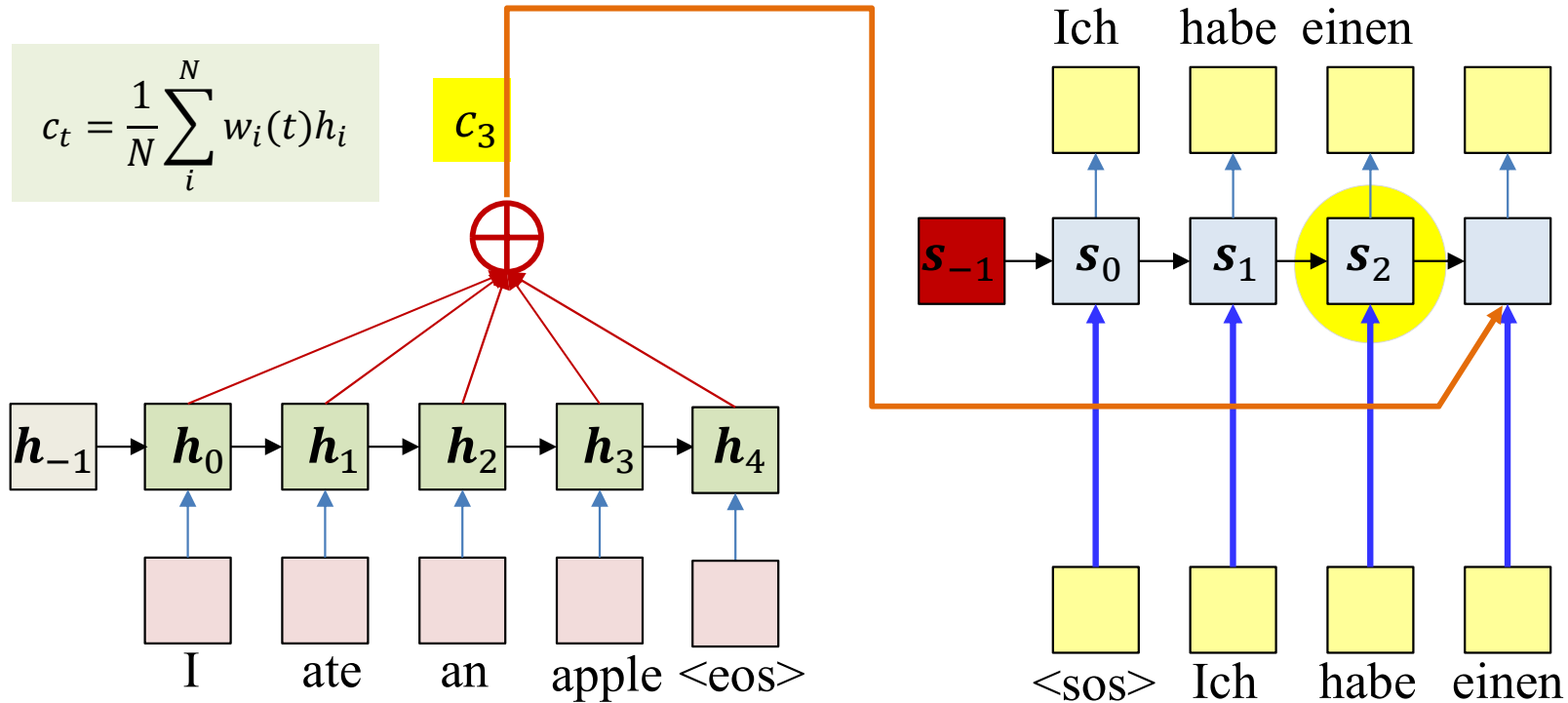
Attention Models



$$c_t = \frac{1}{N} \sum_i^N w_i(t) h_i$$

- **Attention weights:** The weights $w_i(t)$ are dynamically computed as functions of decoder state
 - Expectation: if the model is well-trained, this will automatically “highlight” the correct input
- But how are these computed?

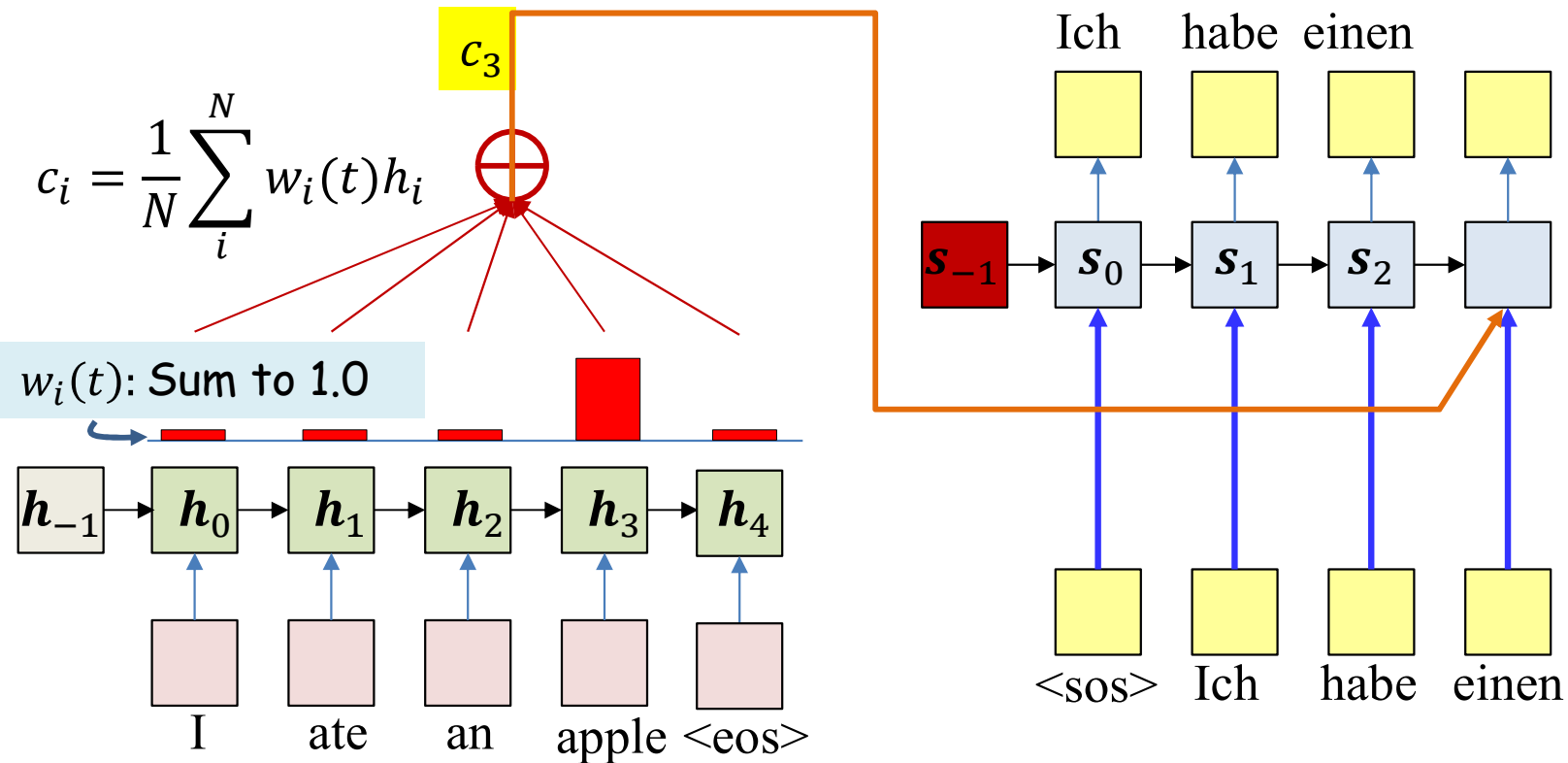
Attention weights at time t



- The “attention” weights $w_i(t)$ at time t must be computed from available information at time t
- The primary information is s_{t-1} (the state at time $t - 1$)
 - Also, the input word at time t , but generally not used for simplicity

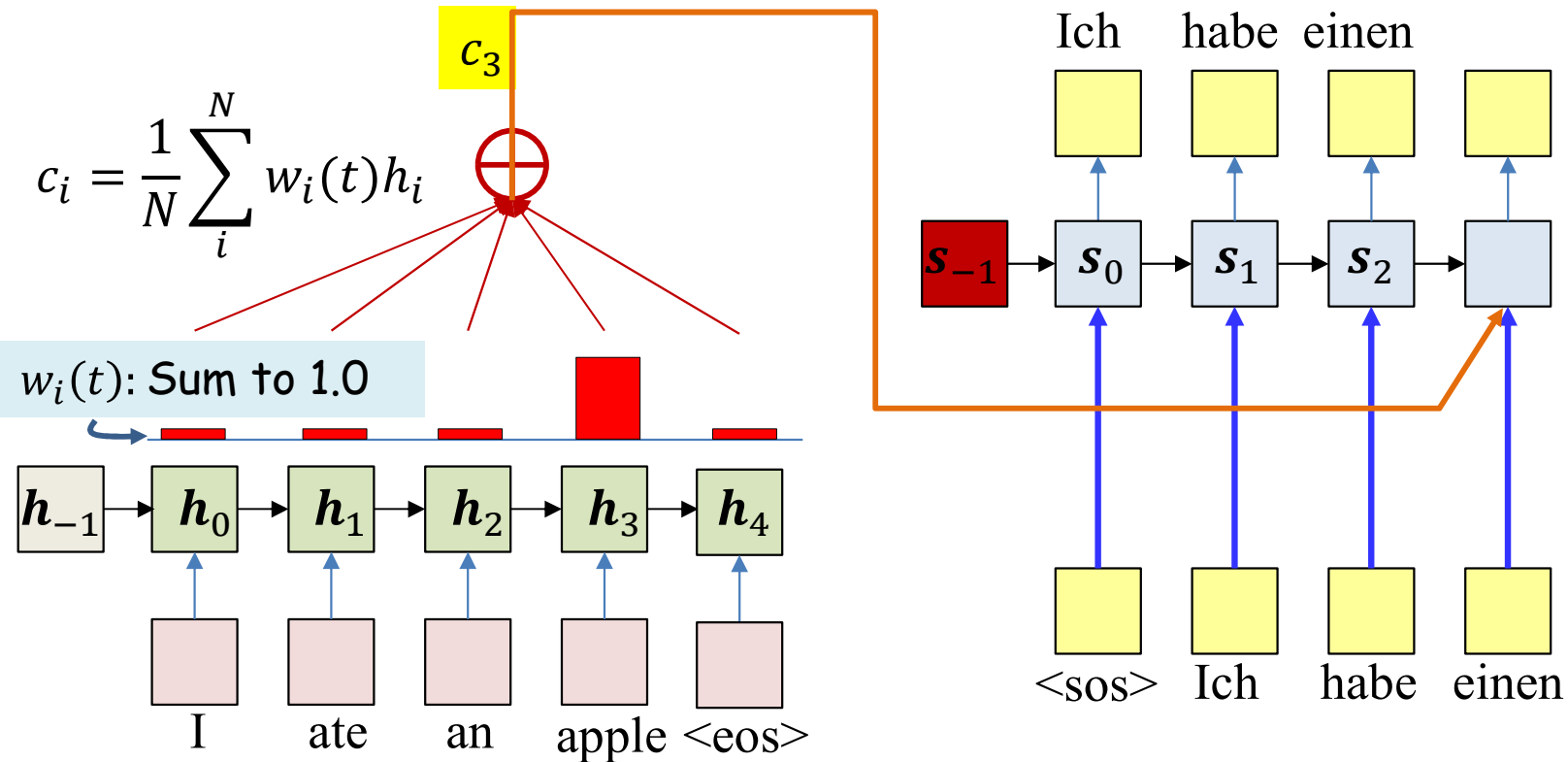
$$w_i(t) = a(\mathbf{h}_i, \mathbf{s}_{t-1})$$

Requirement on attention weights



- The weights $w_i(t)$ must be positive and sum to 1.0
 - I.e. be a distribution
 - Ideally, they must be high for the most relevant inputs for the i th output and low elsewhere

Requirement on attention weights



- The weights $w_i(t)$ must be positive and sum to 1.0
 - I.e. be a distribution
 - Ideally, they must be high for the most relevant inputs for the i th output and low elsewhere
- Solution: A two step weight computation
 - First compute *raw* weights (which could be +ve or -ve)
 - Then softmax them to convert them to a distribution

$$e_i(t) = g(\mathbf{h}_i, \mathbf{s}_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

Poll 4

- @, @, @



Poll 4

The attention framework computes a different “context” vector at each output step (T/F)

- **True**
- False

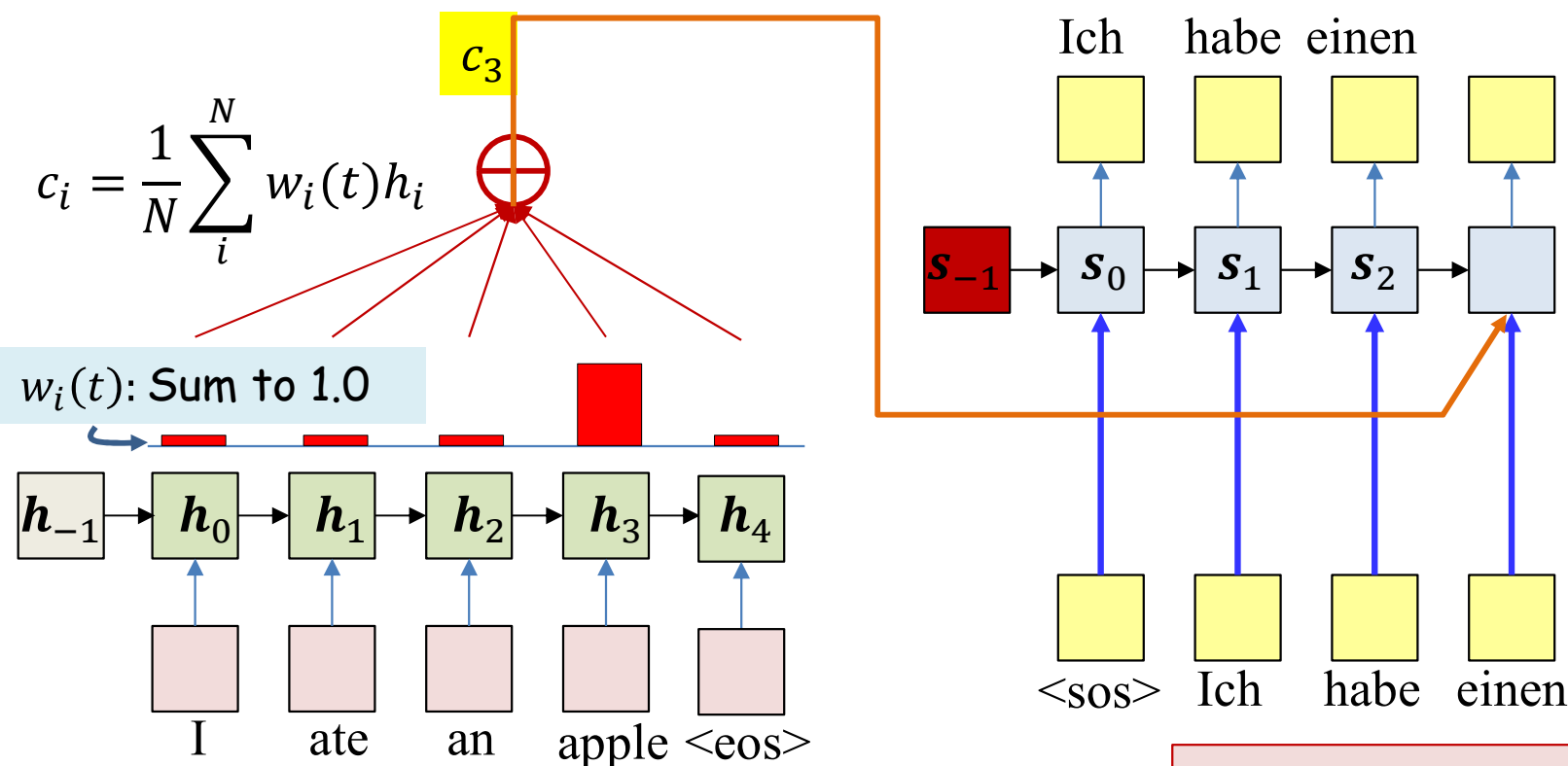
The context vector is chosen as the hidden (encoder) representation of the input word that is assigned the highest attention weight (T/F)

- True
- **False**

The attention weight to any input word is a function of the hidden encoder representation of the word and the most recent decoder state (T/F)

- **True**
- False

Requirement on attention weights



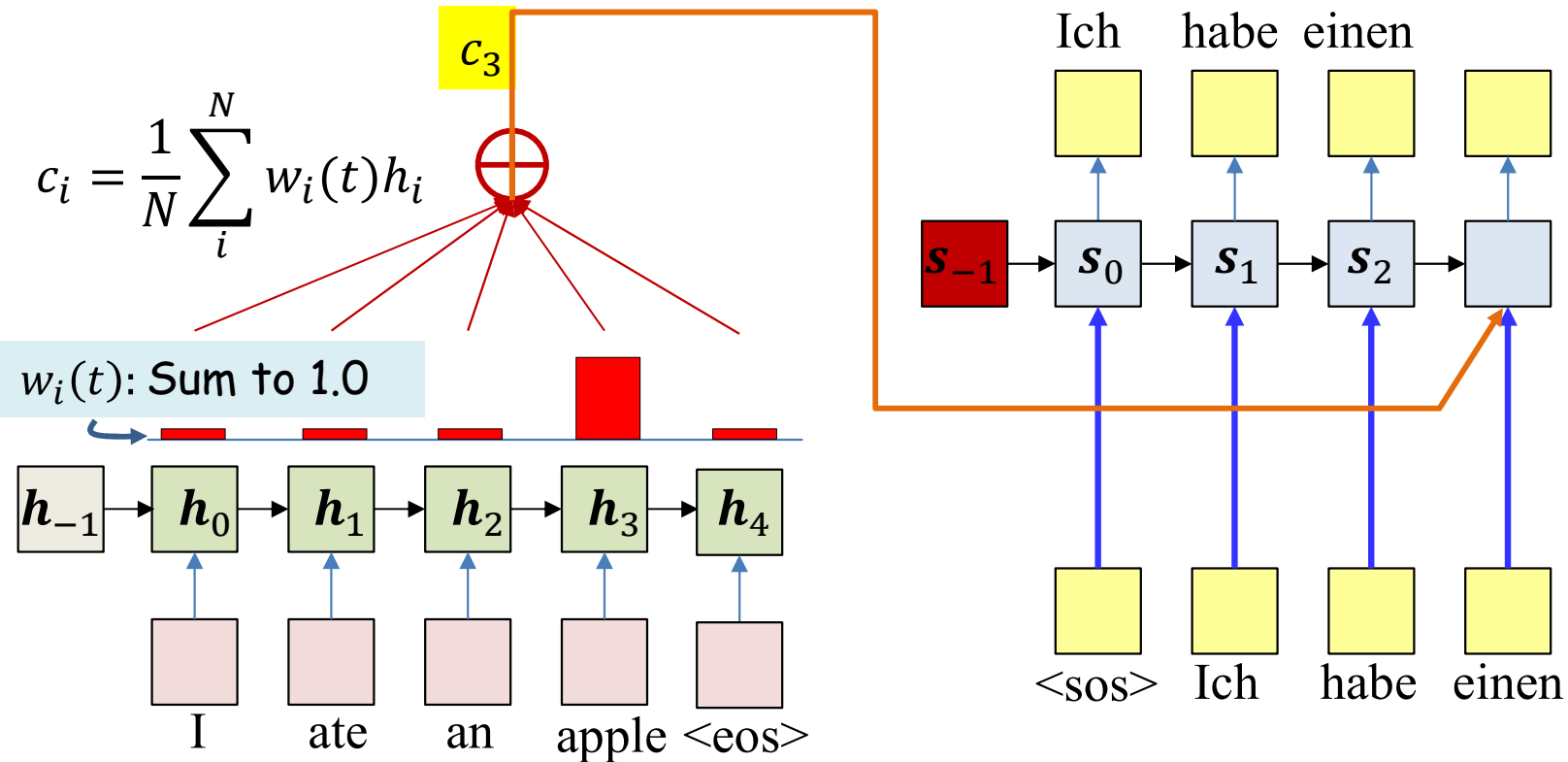
- The weights $w_i(t)$ must be positive and sum to 1.0
 - I.e. be a distribution
 - Ideally, they must be high for the most relevant inputs for the i th output and low elsewhere
- Solution: A two step weight computation
 - First compute *raw* weights (which could be +ve or -ve)
 - Then softmax them to convert them to a distribution

What is this function?

$$e_i(t) = g(\mathbf{h}_i, \mathbf{s}_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

Attention weights



- Typical options for $g()$ (**variables in red must be learned**)

$$g(h_i, s_{t-1}) = h_i^T s_{t-1}$$

$$g(h_i, s_{t-1}) = h_i^T \mathbf{W}_g s_{t-1}$$

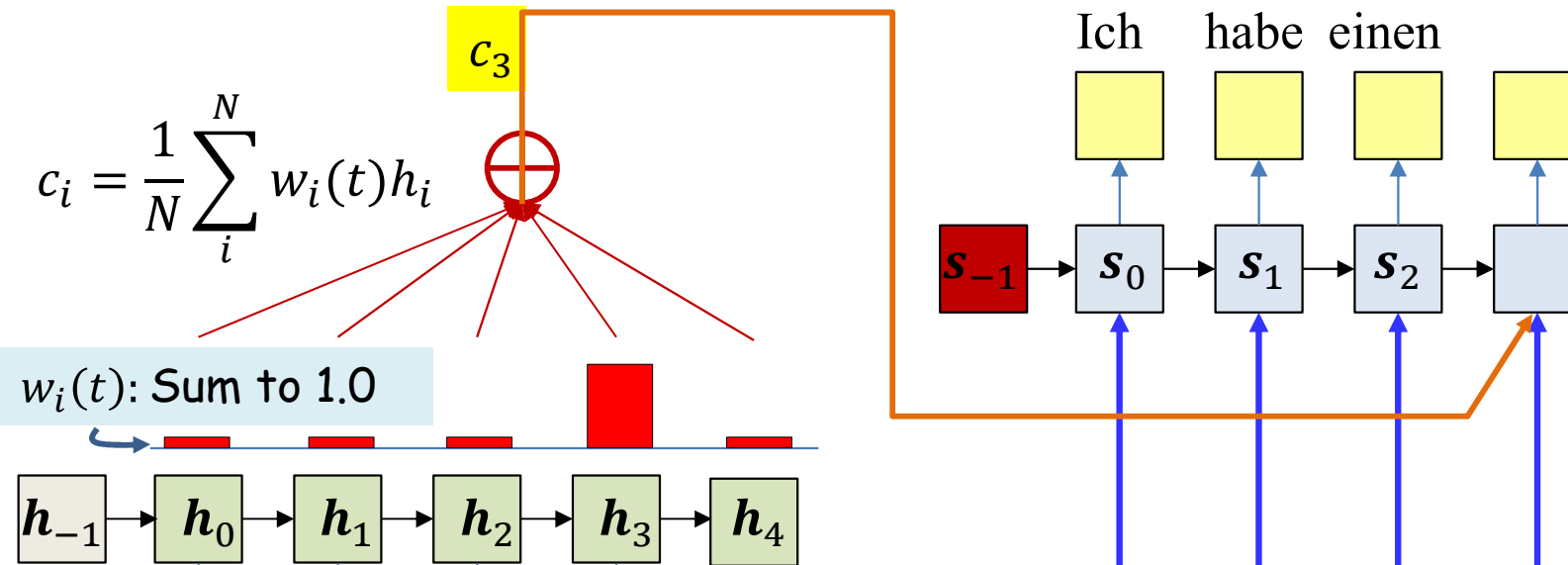
$$g(h_i, s_{t-1}) = \mathbf{v}_g^T \tanh \left(\mathbf{W}_g \begin{bmatrix} h_i \\ s_{t-1} \end{bmatrix} \right)$$

$$g(h_i, s_{t-1}) = \text{MLP}([h_i, s_{t-1}])$$

$$e_i(t) = g(h_i, s_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

Attention weights



Let's consider a typical conversion process assuming this model as an example

- Typical options for $g()$ (**variables in red must be learned**)

$$g(\mathbf{h}_i, \mathbf{s}_{t-1}) = \mathbf{h}_i^T \mathbf{s}_{t-1}$$

$$g(\mathbf{h}_i, \mathbf{s}_{t-1}) = \mathbf{h}_i^T \mathbf{W}_g \mathbf{s}_{t-1}$$

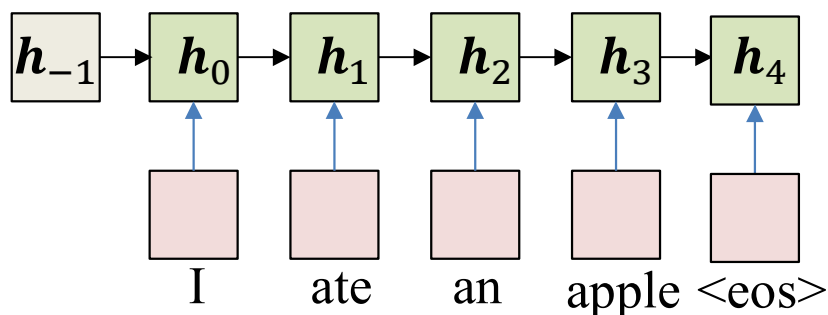
$$g(\mathbf{h}_i, \mathbf{s}_{t-1}) = \mathbf{v}_g^T \tanh \left(\mathbf{W}_g \begin{bmatrix} \mathbf{h}_i \\ \mathbf{s}_{t-1} \end{bmatrix} \right)$$

$$g(\mathbf{h}_i, \mathbf{s}_{t-1}) = \text{MLP}([\mathbf{h}_i, \mathbf{s}_{t-1}])$$

$$e_i(t) = g(\mathbf{h}_i, \mathbf{s}_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

Converting an input: Inference



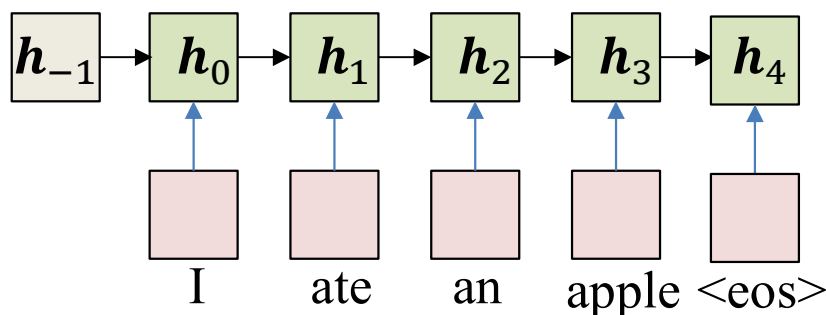
- Pass the input through the encoder to produce hidden representations h_i

Converting an input: Inference

This may be

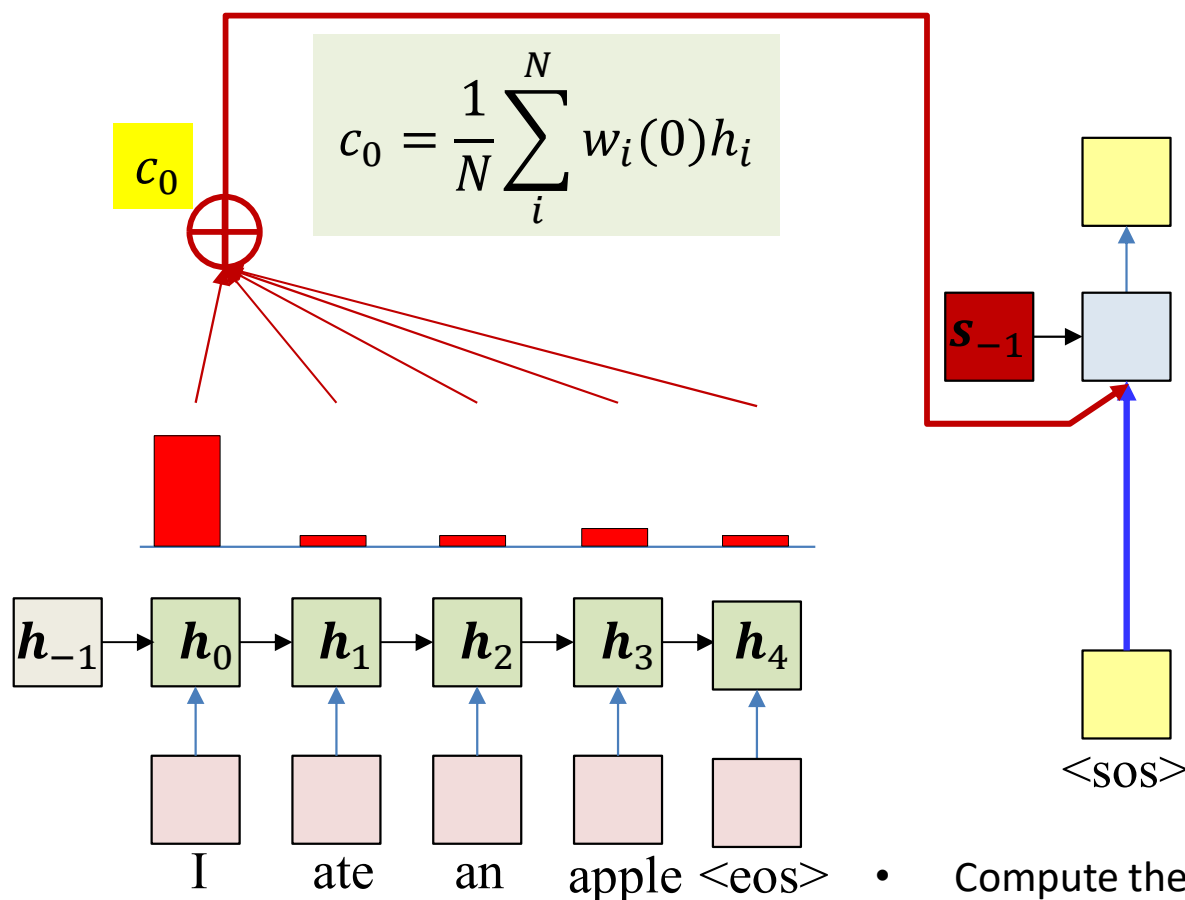
- a learned parameter, or
- Or just set to some fixed value, e.g. a vector of 1s or 0s, or
- Or the average of all the encoder embeddings: $\text{mean}(h_0, \dots, h_4)$
- Or $W_{init} \text{mean}(h_0, \dots, h_4)$ where W_{init} is a learned parameter

s_{-1}



- Pass the input through the encoder to produce hidden representations h_i

Converting an input: Inference



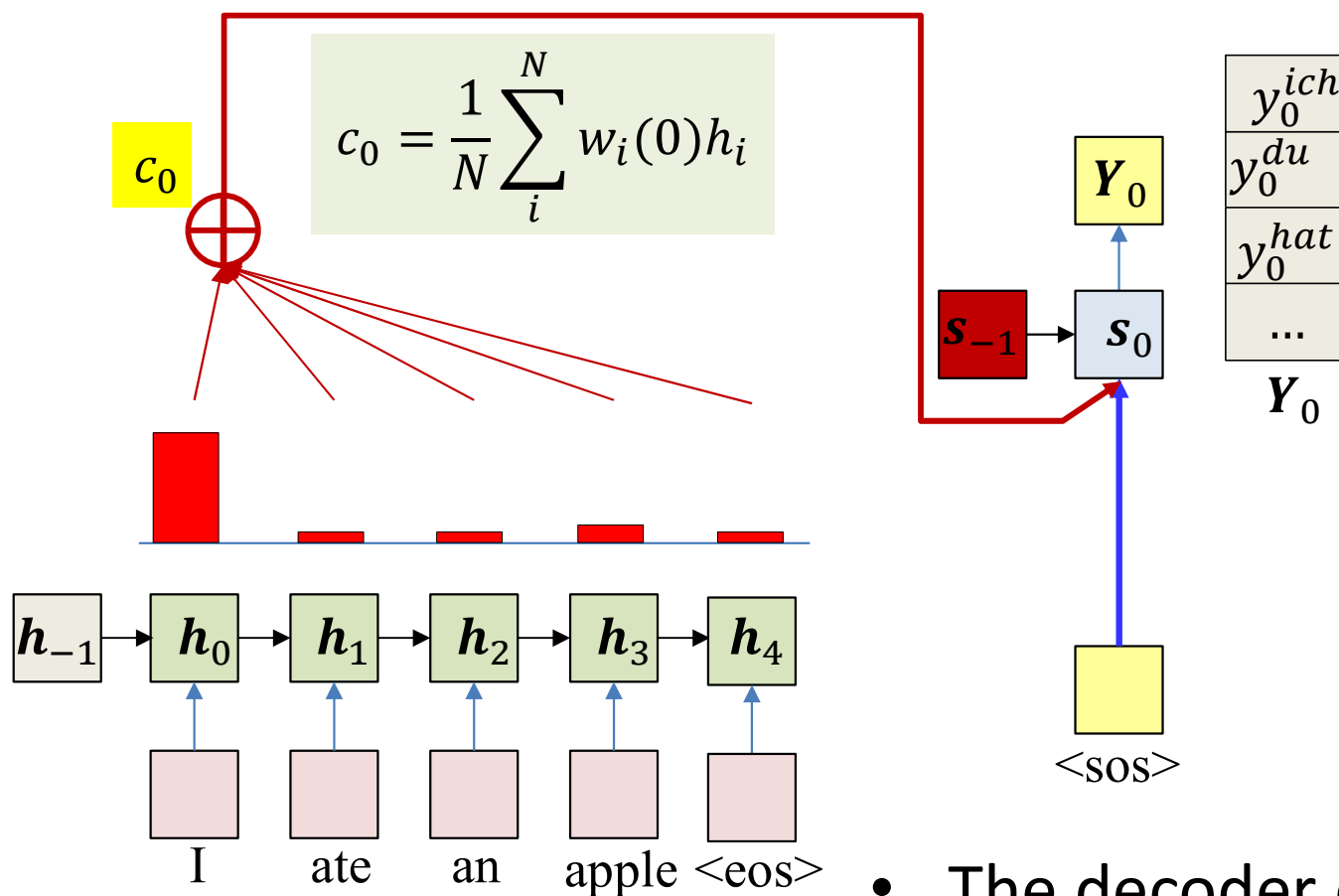
$$g(\mathbf{h}_i, \mathbf{s}_{-1}) = \mathbf{h}_i^T \mathbf{W}_g \mathbf{s}_{-1}$$

$$e_i(0) = g(\mathbf{h}_i, \mathbf{s}_{-1})$$

$$w_i(0) = \frac{\exp(e_i(0))}{\sum_j \exp(e_j(0))}$$

- Compute the attention weights $w_i(0)$ for the first output using s_{-1}
 - Will be a distribution over the input words
- Compute “context” c_0
 - Weighted sum of input word hidden states
- Input c_0 and <sos> to the decoder at time 0
 - <sos> because we are starting a new sequence
 - In practice we will enter the *embedding* of <sos>

Converting an input: Inference



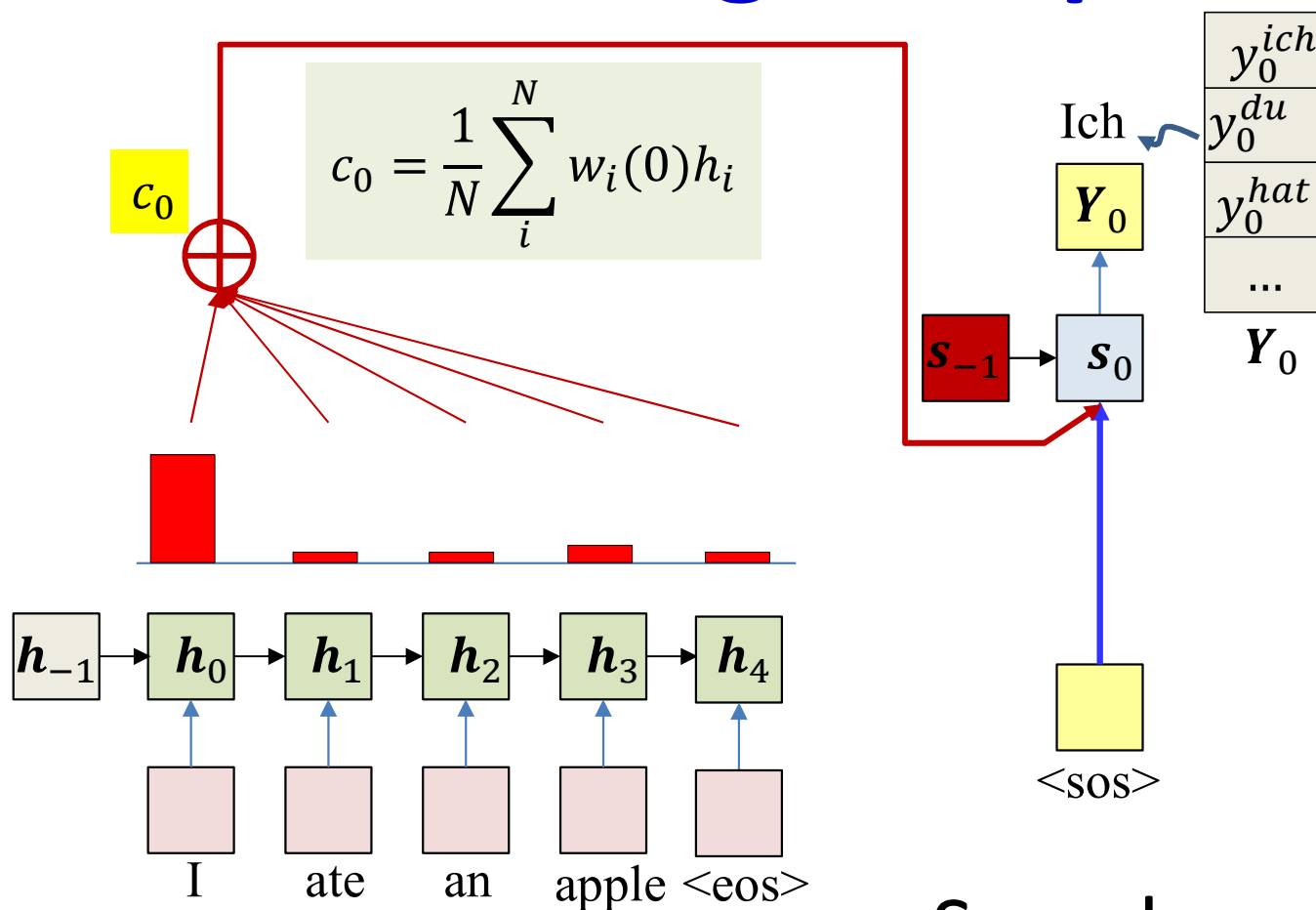
$$g(\mathbf{h}_i, \mathbf{s}_{-1}) = \mathbf{h}_i^T \mathbf{W}_g \mathbf{s}_{-1}$$

$$e_i(0) = g(\mathbf{h}_i, \mathbf{s}_{-1})$$

$$w_i(0) = \frac{\exp(e_i(0))}{\sum_j \exp(e_j(0))}$$

- The decoder computes
 - s_0
 - A probability distribution over the output vocabulary
 - Output of softmax output layer

Converting an input: Inference



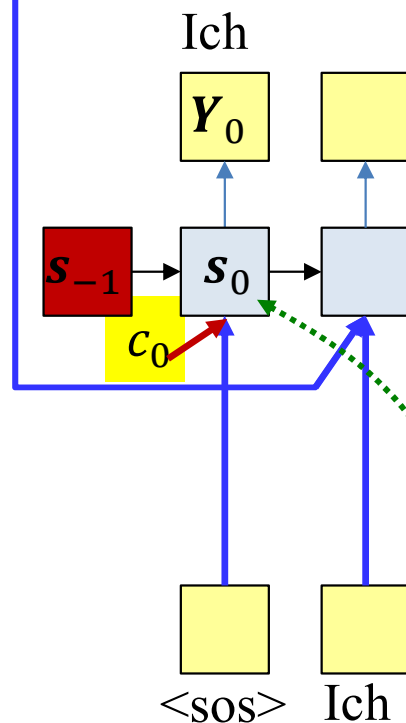
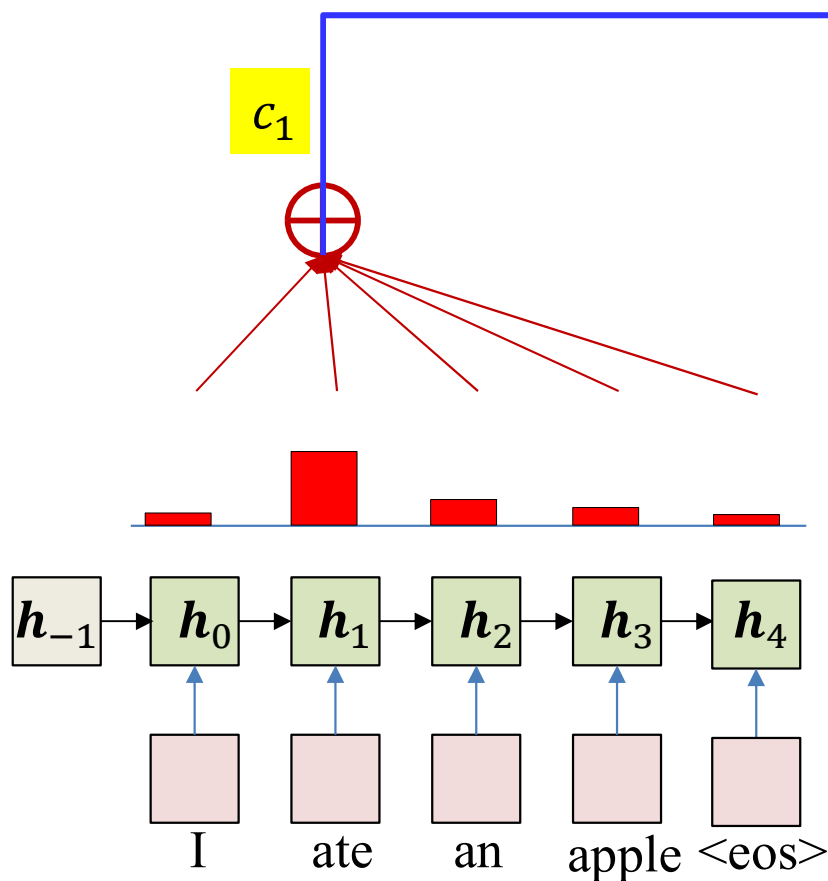
$$g(h_i, s_{-1}) = h_i^T \mathbf{W}_g s_{-1}$$

$$e_i(0) = g(h_i, s_{-1})$$

$$w_i(0) = \frac{\exp(e_i(0))}{\sum_j \exp(e_j(0))}$$

- Sample a word from the output distribution

Converting an input: Inference



- Compute the attention weights $w_i(1)$ over all inputs for the *second* output using s_0
 - Compute raw weights, followed by softmax
- Compute “context” c_1
 - Weighted sum of input hidden representations
- Input c_1 and first output word to the decoder
 - In practice we enter the *embedding* of the word

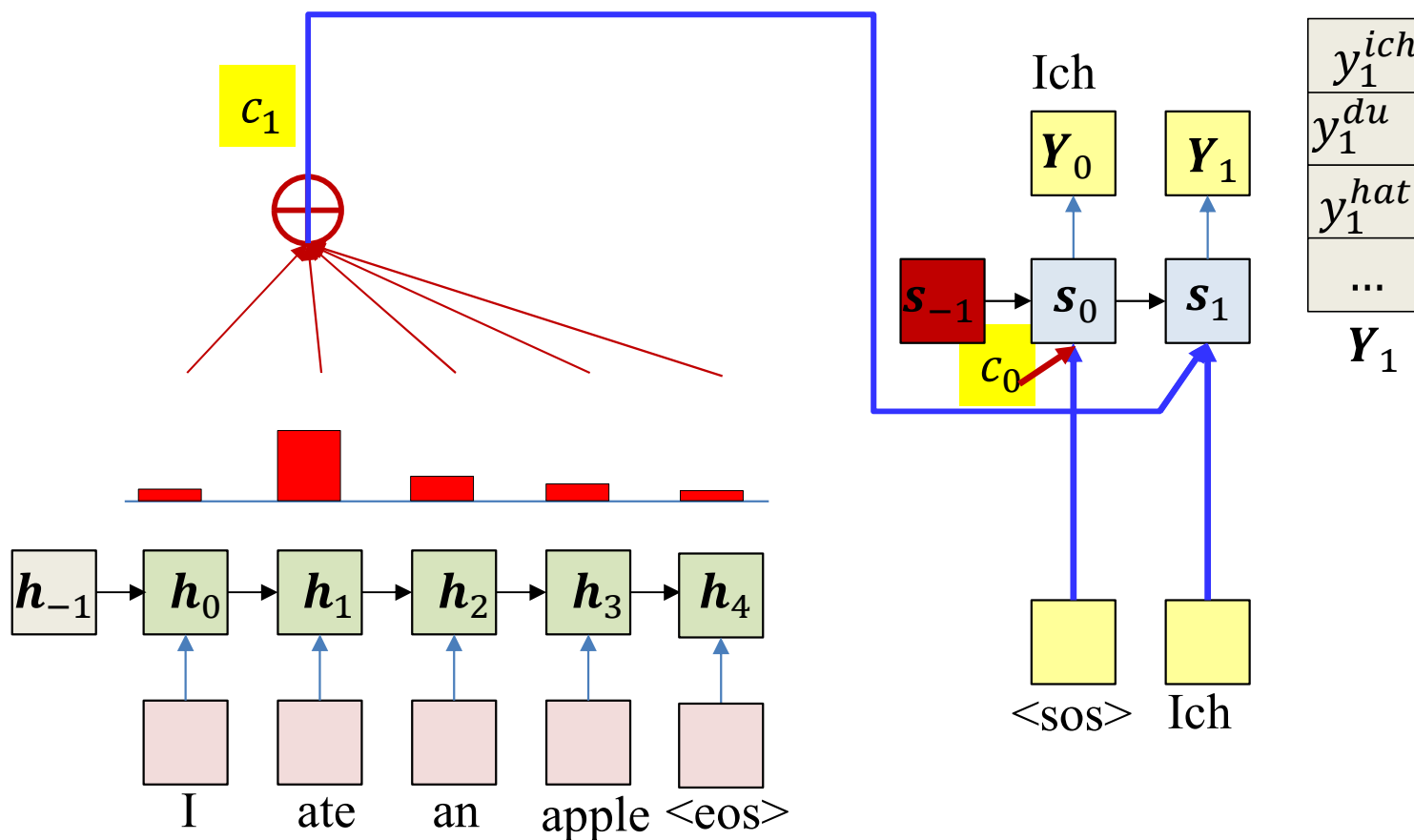
$$g(\mathbf{h}_i, \mathbf{s}_0) = \mathbf{h}_i^T \mathbf{W}_g \mathbf{s}_0$$

$$e_i(1) = g(\mathbf{h}_i, \mathbf{s}_0)$$

$$w_i(1) = \frac{\exp(e_i(1))}{\sum_j \exp(e_j(1))}$$

$$c_1 = \frac{1}{N} \sum_i^N w_i(1) \mathbf{h}_i$$

Converting an input: Inference



- The decoder computes
 - s_1
 - A probability distribution over the output vocabulary

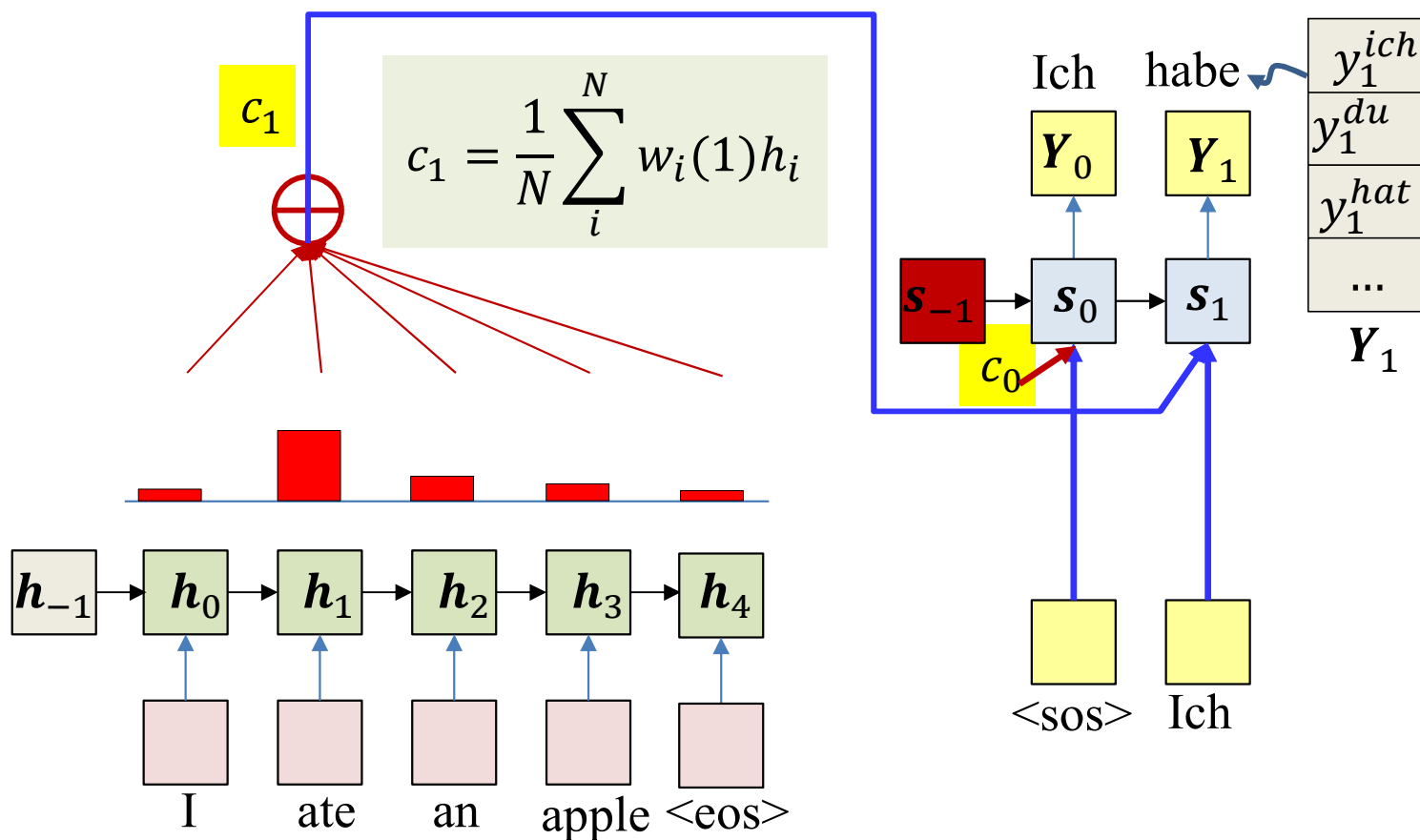
$$g(\mathbf{h}_i, \mathbf{s}_0) = \mathbf{h}_i^T \mathbf{W}_g \mathbf{s}_0$$

$$e_i(1) = g(\mathbf{h}_i, \mathbf{s}_0)$$

$$w_i(1) = \frac{\exp(e_i(1))}{\sum_j \exp(e_j(1))}$$

$$c_1 = \frac{1}{N} \sum_i^N w_i(1) \mathbf{h}_i$$

Converting an input: Inference



- Sample the second word from the output distribution

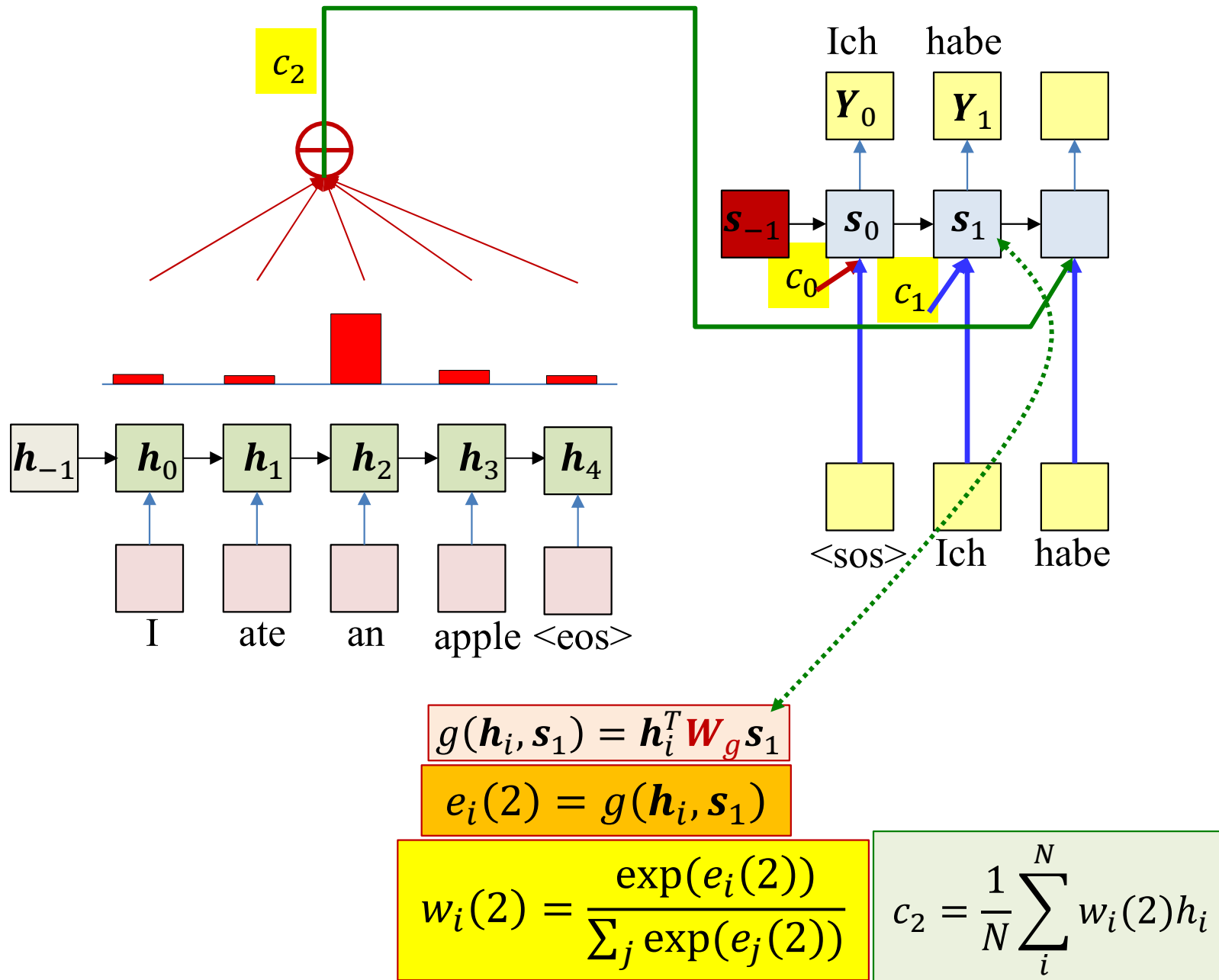
$$g(h_i, s_0) = h_i^T \mathbf{W}_g s_0$$

$$e_i(1) = g(h_i, s_0)$$

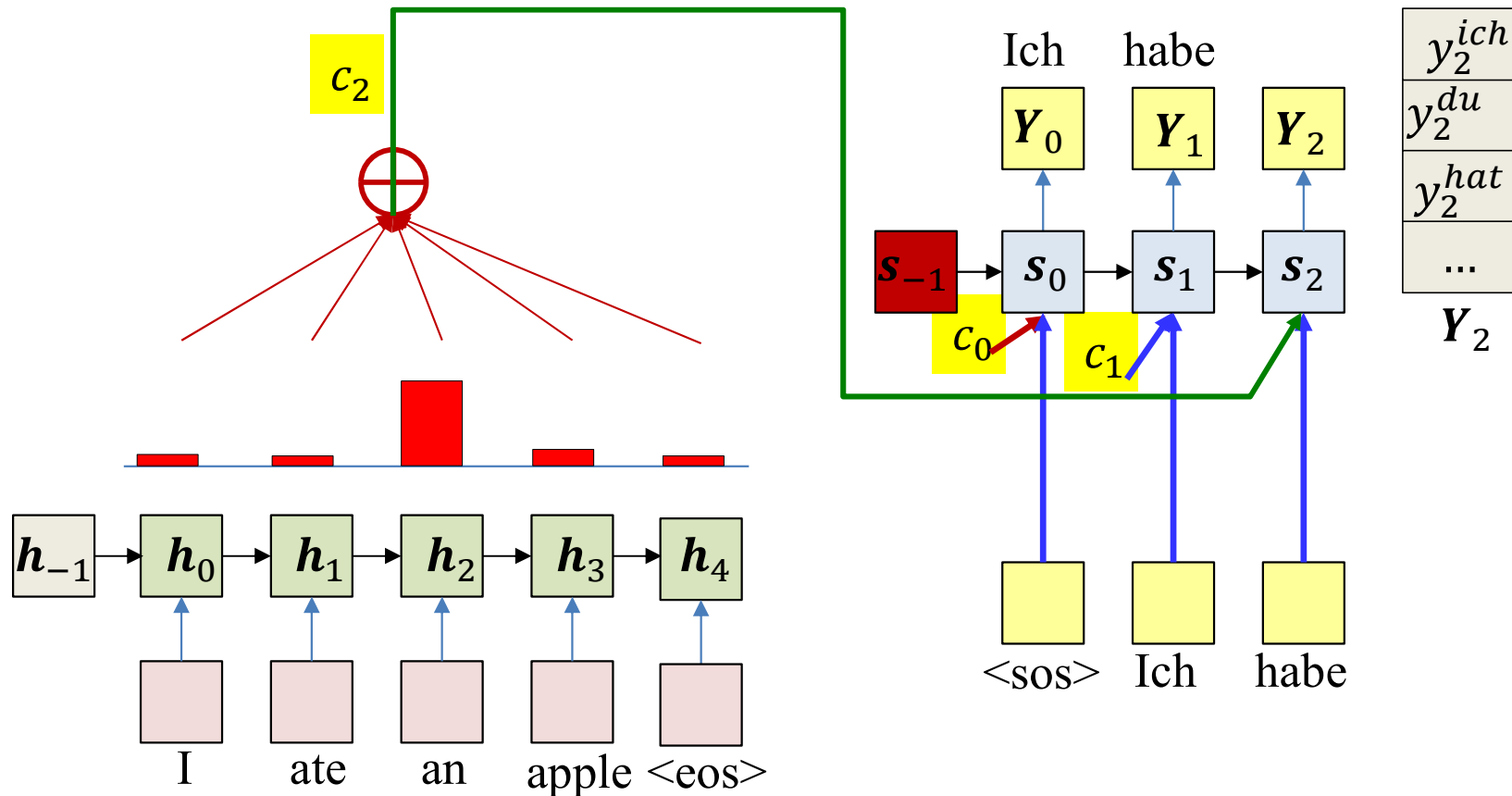
$$w_i(1) = \frac{\exp(e_i(1))}{\sum_j \exp(e_j(1))}$$

$$c_1 = \frac{1}{N} \sum_i^N w_i(1) h_i$$

Converting an input: Inference



Converting an input: Inference



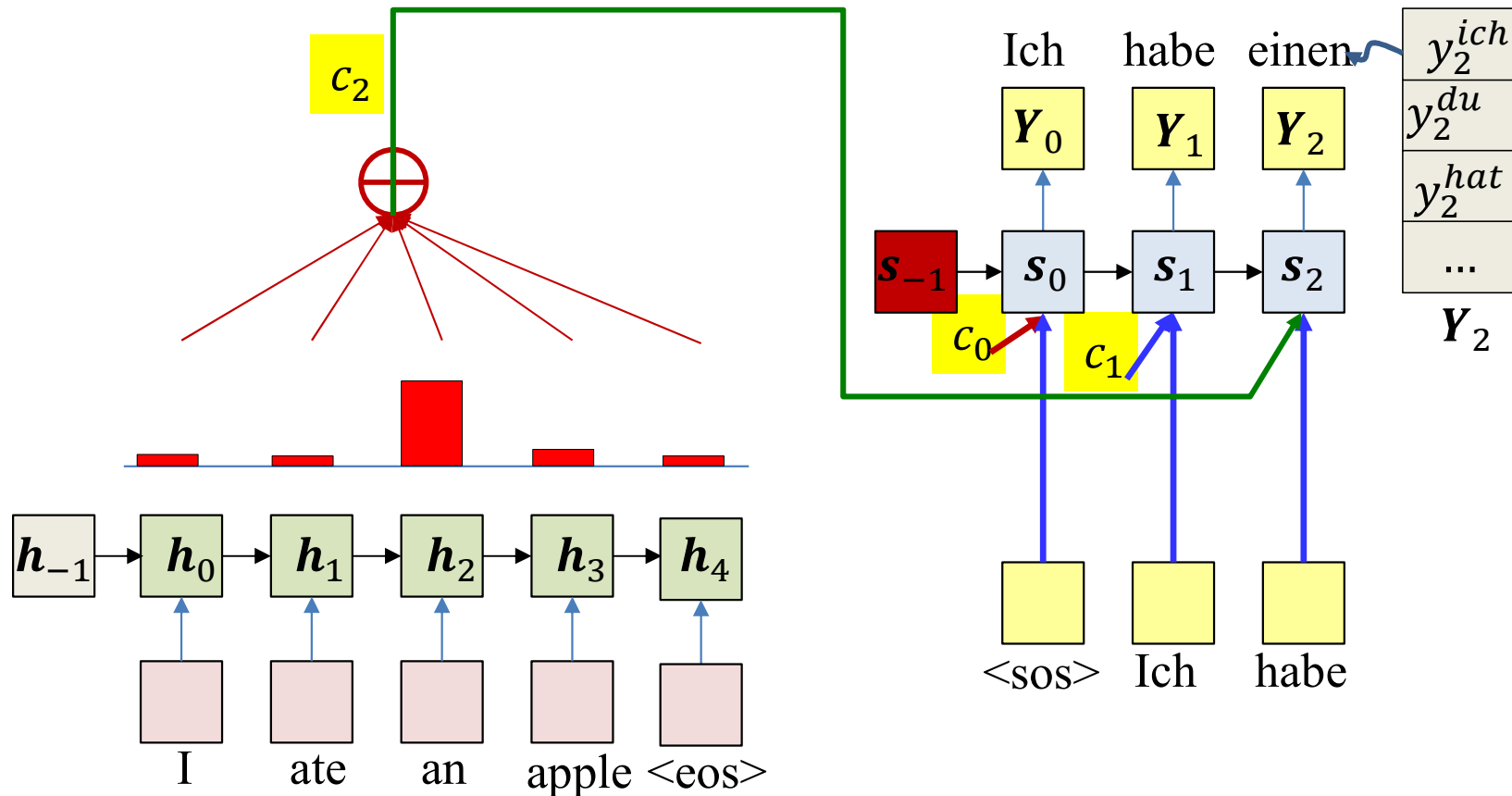
$$g(h_i, s_1) = h_i^T W_g s_1$$

$$e_i(2) = g(h_i, s_1)$$

$$w_i(2) = \frac{\exp(e_i(2))}{\sum_j \exp(e_j(2))}$$

$$c_2 = \frac{1}{N} \sum_i^N w_i(2) h_i$$

Converting an input: Inference



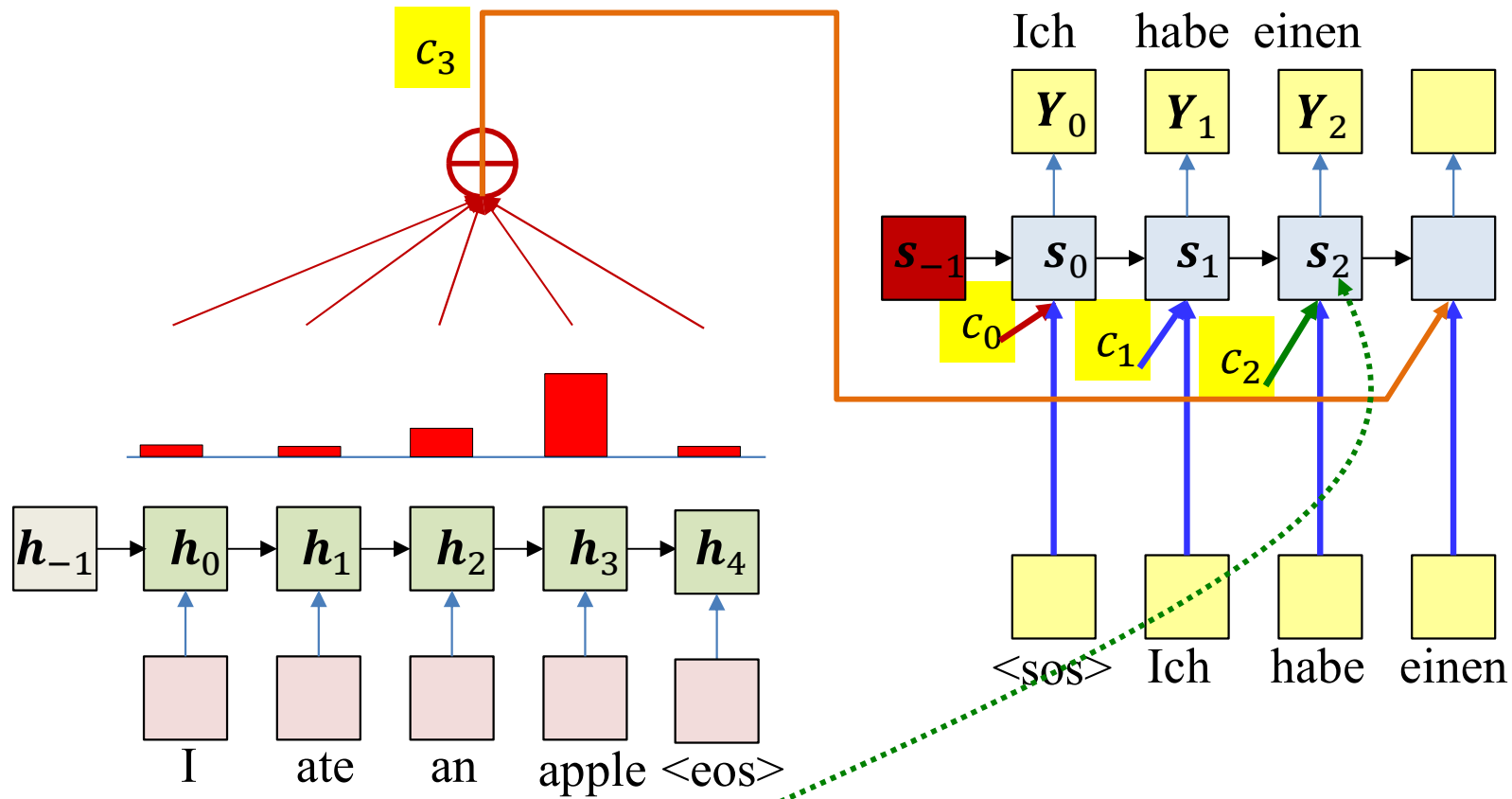
$$g(h_i, s_1) = h_i^T W_g s_1$$

$$e_i(2) = g(h_i, s_1)$$

$$w_i(2) = \frac{\exp(e_i(2))}{\sum_j \exp(e_j(2))}$$

$$c_2 = \frac{1}{N} \sum_i^N w_i(2) h_i$$

Converting an input: Inference

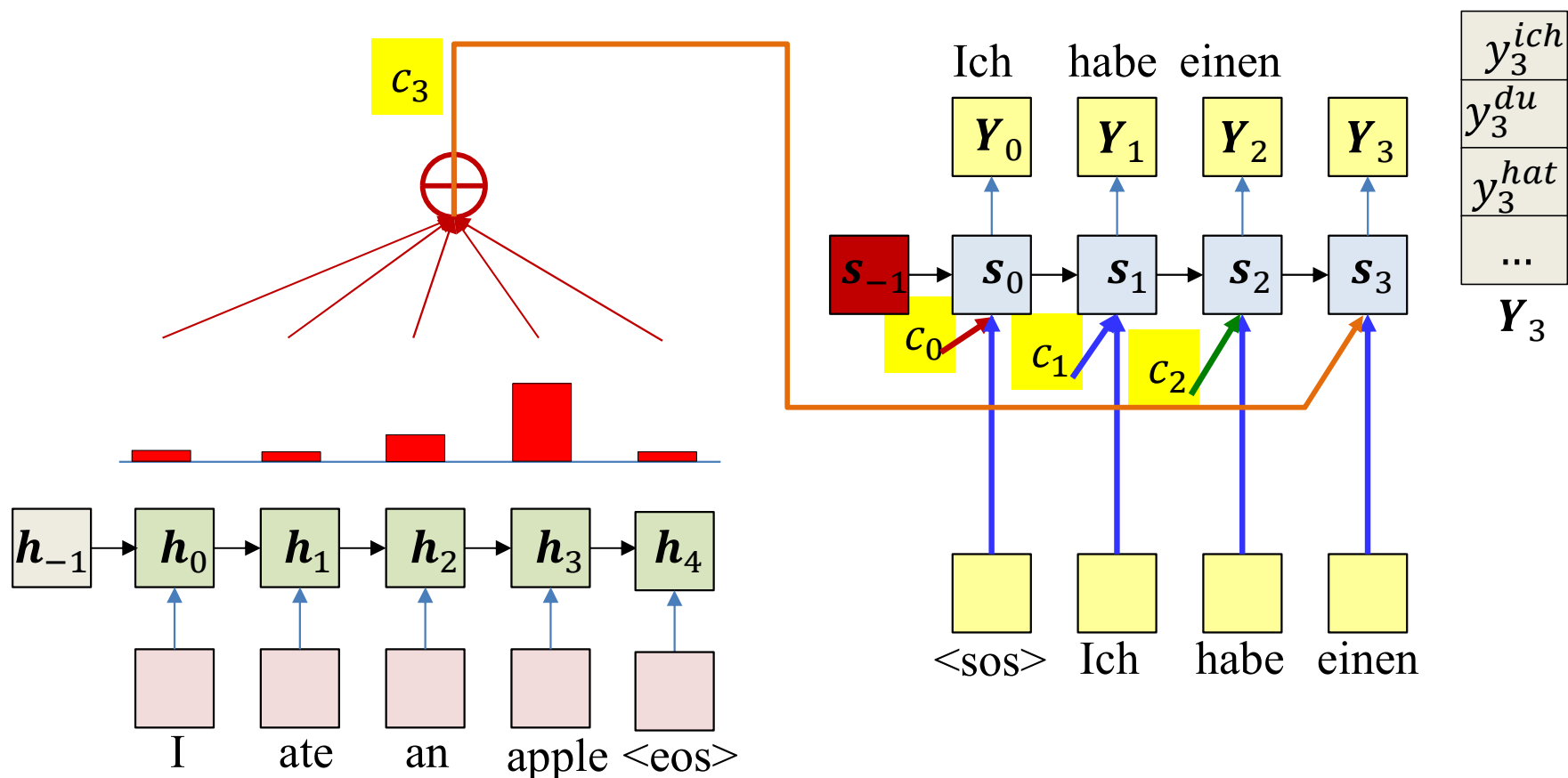


$$e_i(3) = g(h_i, s_2)$$

$$w_i(3) = \frac{\exp(e_i(3))}{\sum_j \exp(e_j(3))}$$

$$c_3 = \frac{1}{N} \sum_i^N w_i(3) h_i$$

Converting an input: Inference

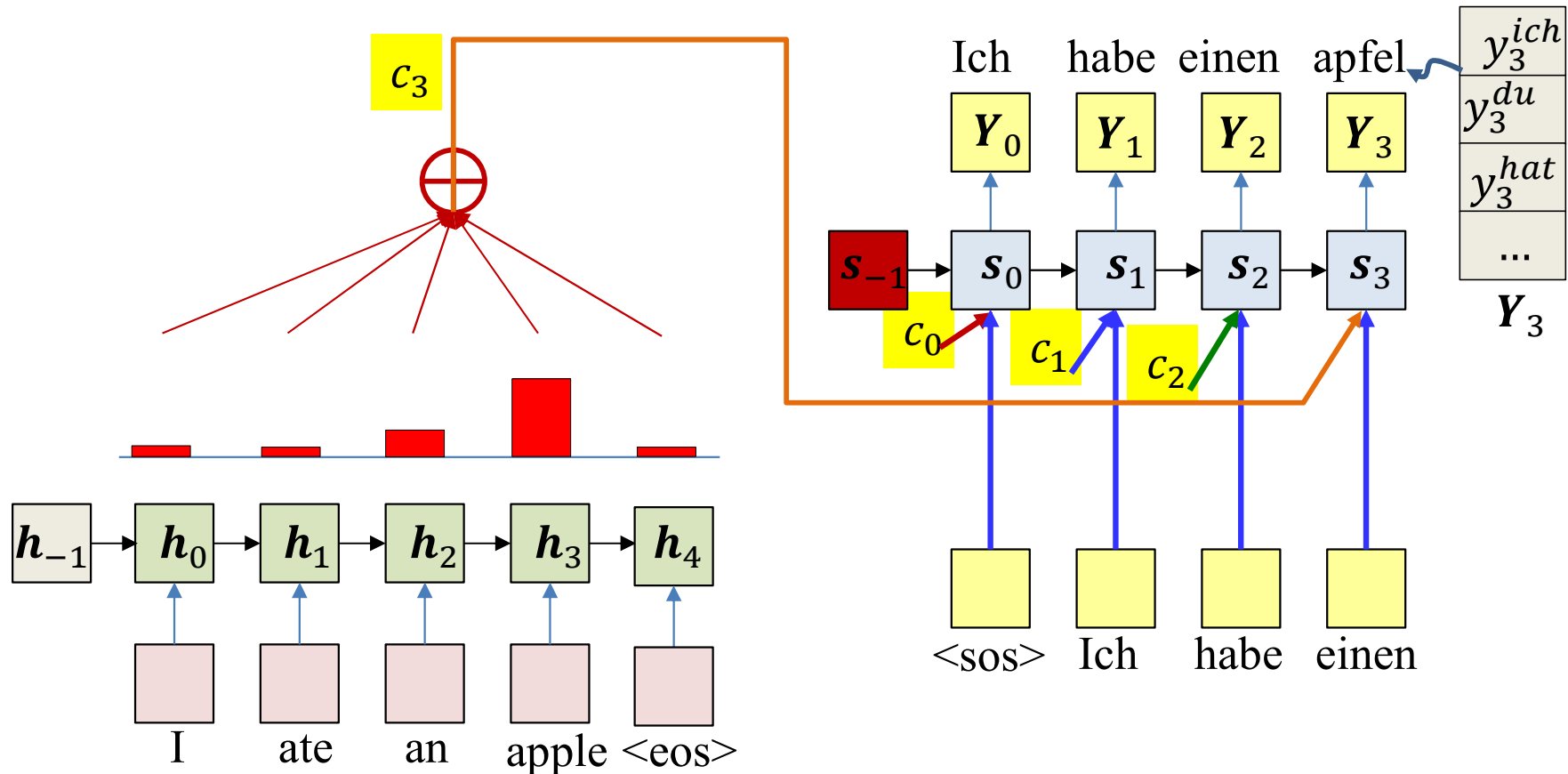


$$e_i(3) = g(h_i, s_2)$$

$$w_i(3) = \frac{\exp(e_i(3))}{\sum_j \exp(e_j(3))}$$

$$c_3 = \frac{1}{N} \sum_i^N w_i(3) h_i$$

Converting an input: Inference

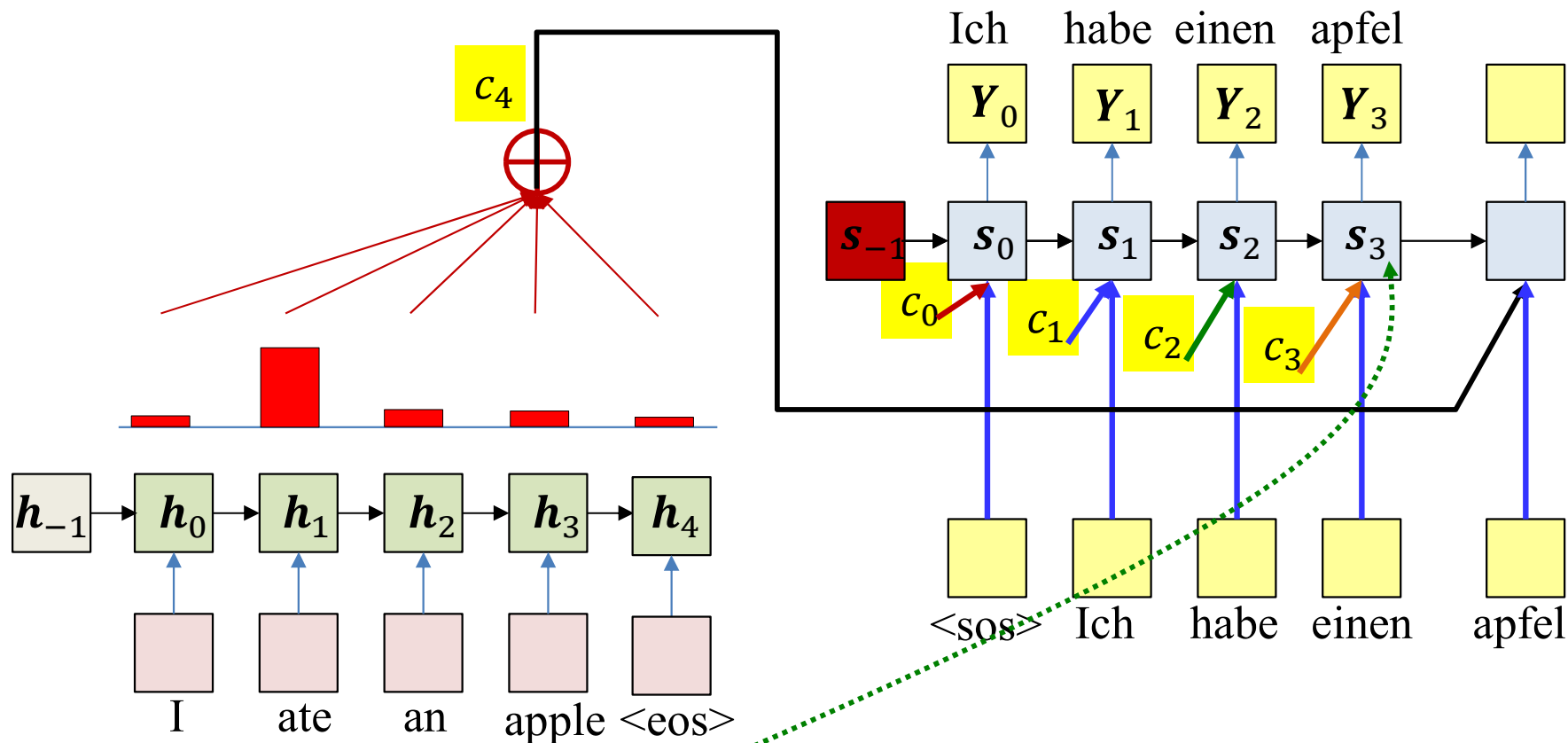


$$e_i(3) = g(h_i, s_2)$$

$$w_i(3) = \frac{\exp(e_i(3))}{\sum_j \exp(e_j(3))}$$

$$c_3 = \frac{1}{N} \sum_i^N w_i(3) h_i$$

Converting an input: Inference

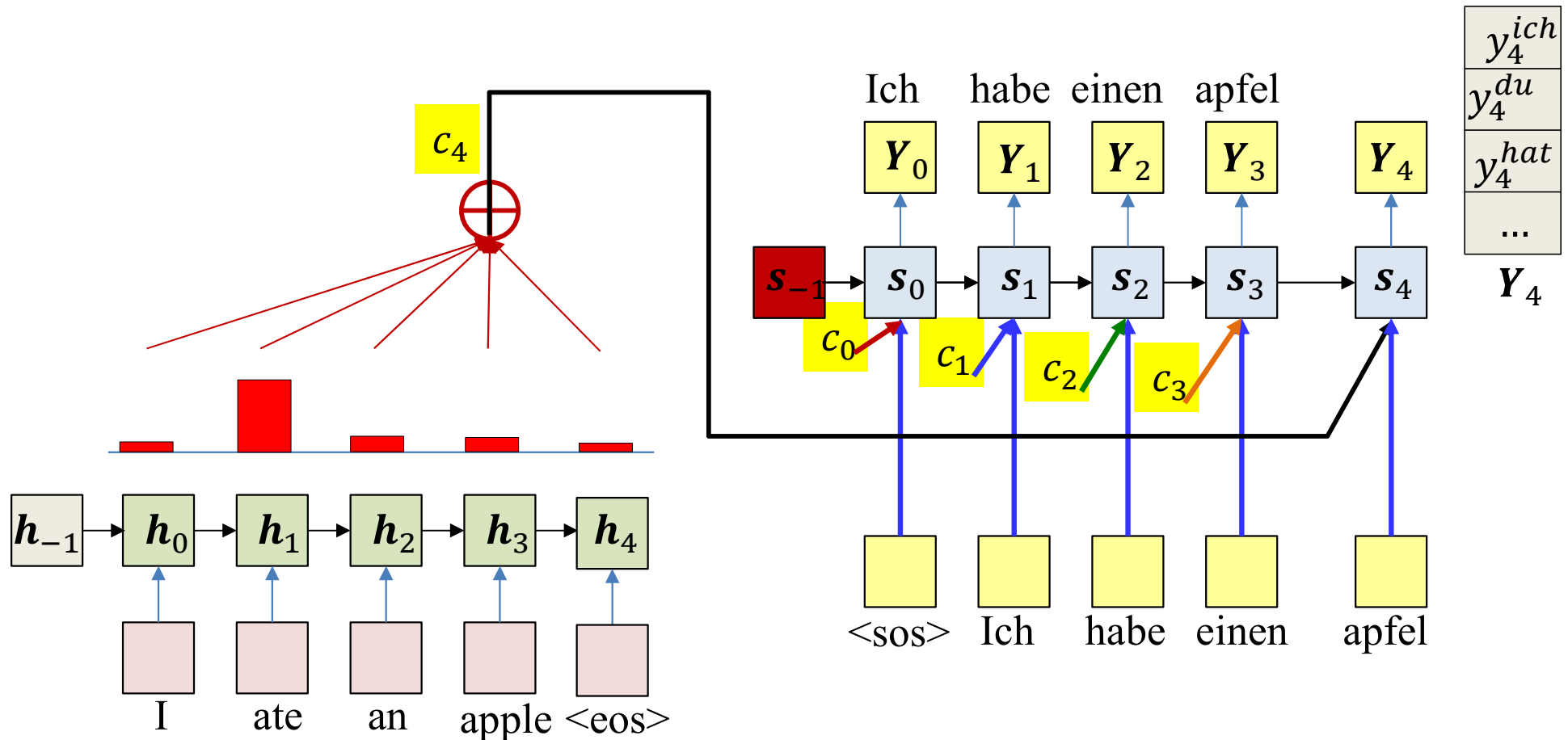


$$e_i(4) = g(h_i, s_3)$$

$$w_i(4) = \frac{\exp(e_i(4))}{\sum_j \exp(e_j(4))}$$

$$c_4 = \frac{1}{N} \sum_i^N w_i(4) h_i$$

Converting an input: Inference

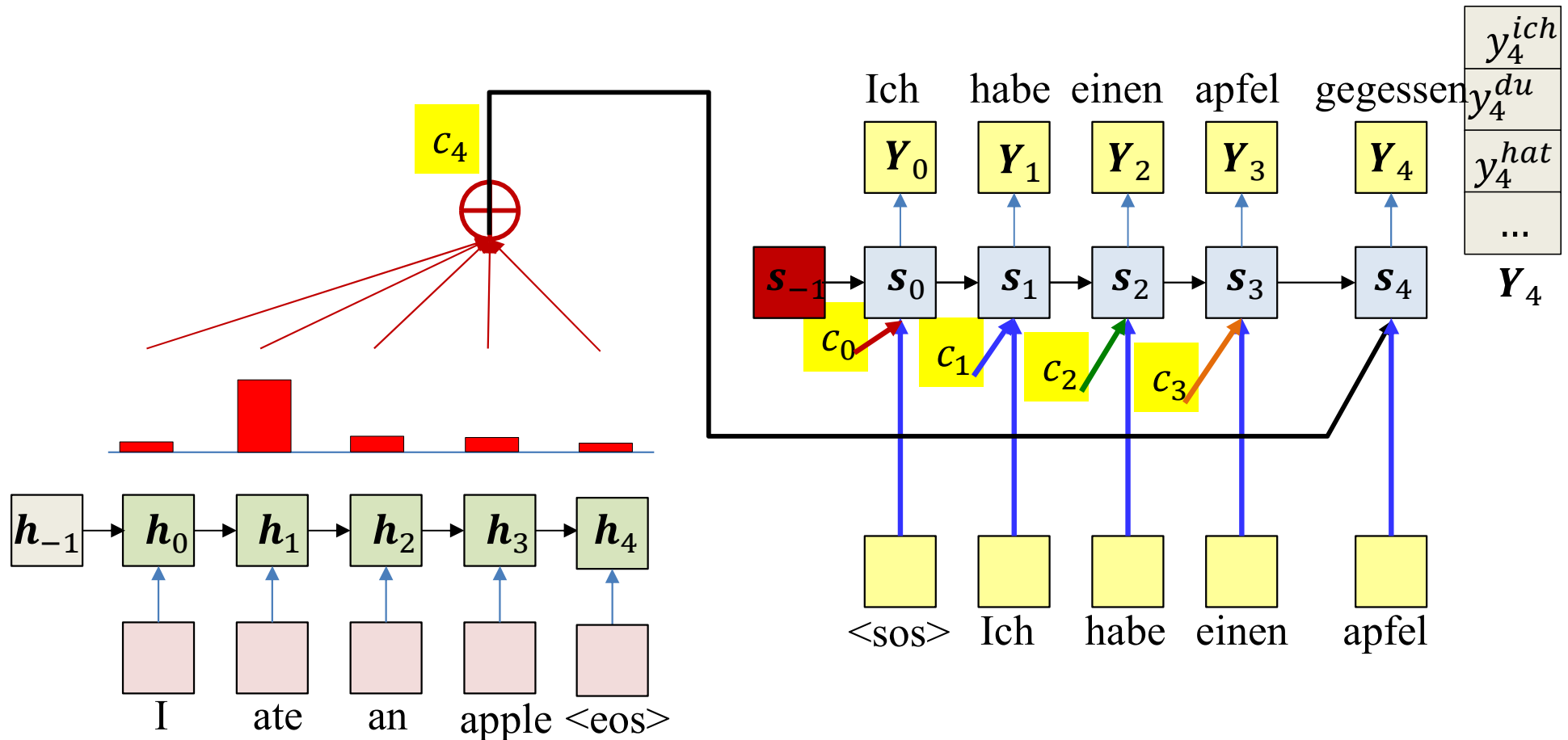


$$e_i(4) = g(h_i, s_3)$$

$$w_i(4) = \frac{\exp(e_i(4))}{\sum_j \exp(e_j(4))}$$

$$c_4 = \frac{1}{N} \sum_i^N w_i(4) h_i$$

Converting an input: Inference

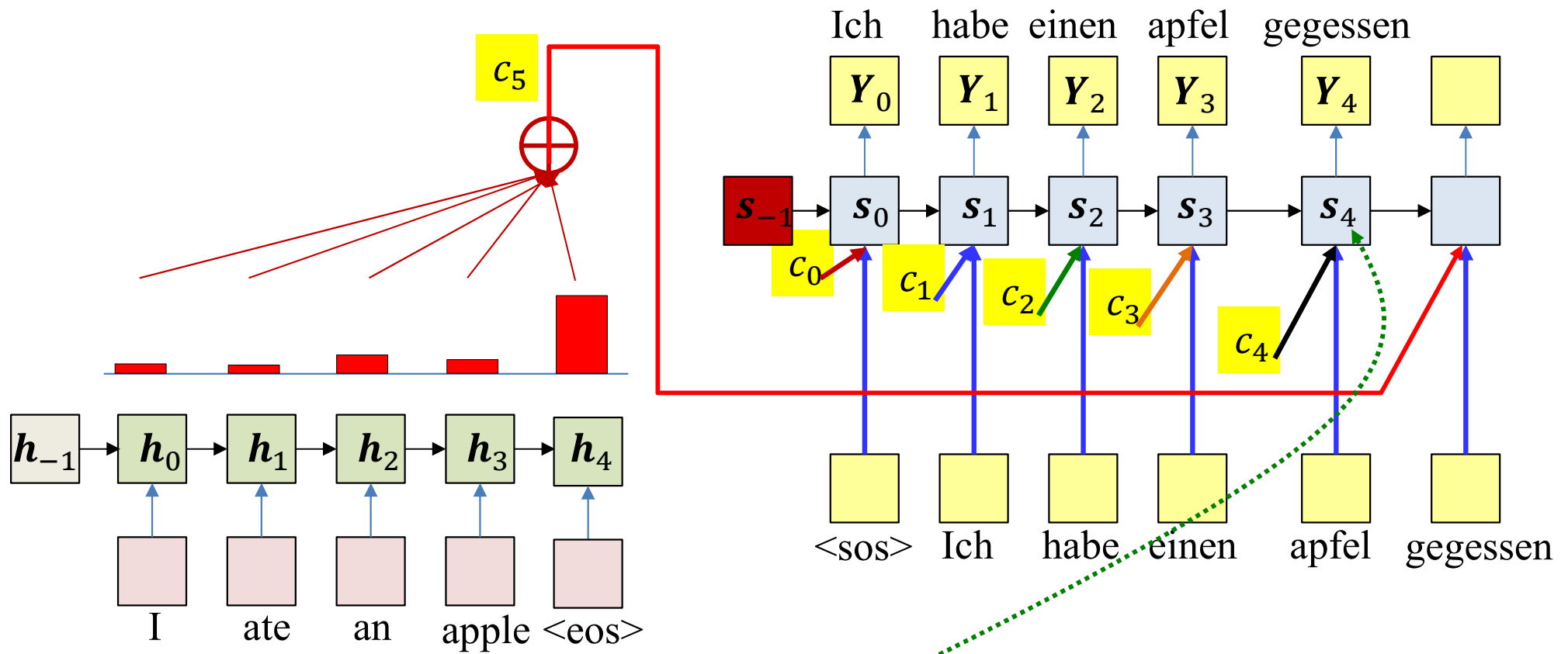


$$e_i(4) = g(h_i, s_3)$$

$$w_i(4) = \frac{\exp(e_i(4))}{\sum_j \exp(e_j(4))}$$

$$c_4 = \frac{1}{N} \sum_i^N w_i(4) h_i$$

Converting an input: Inference

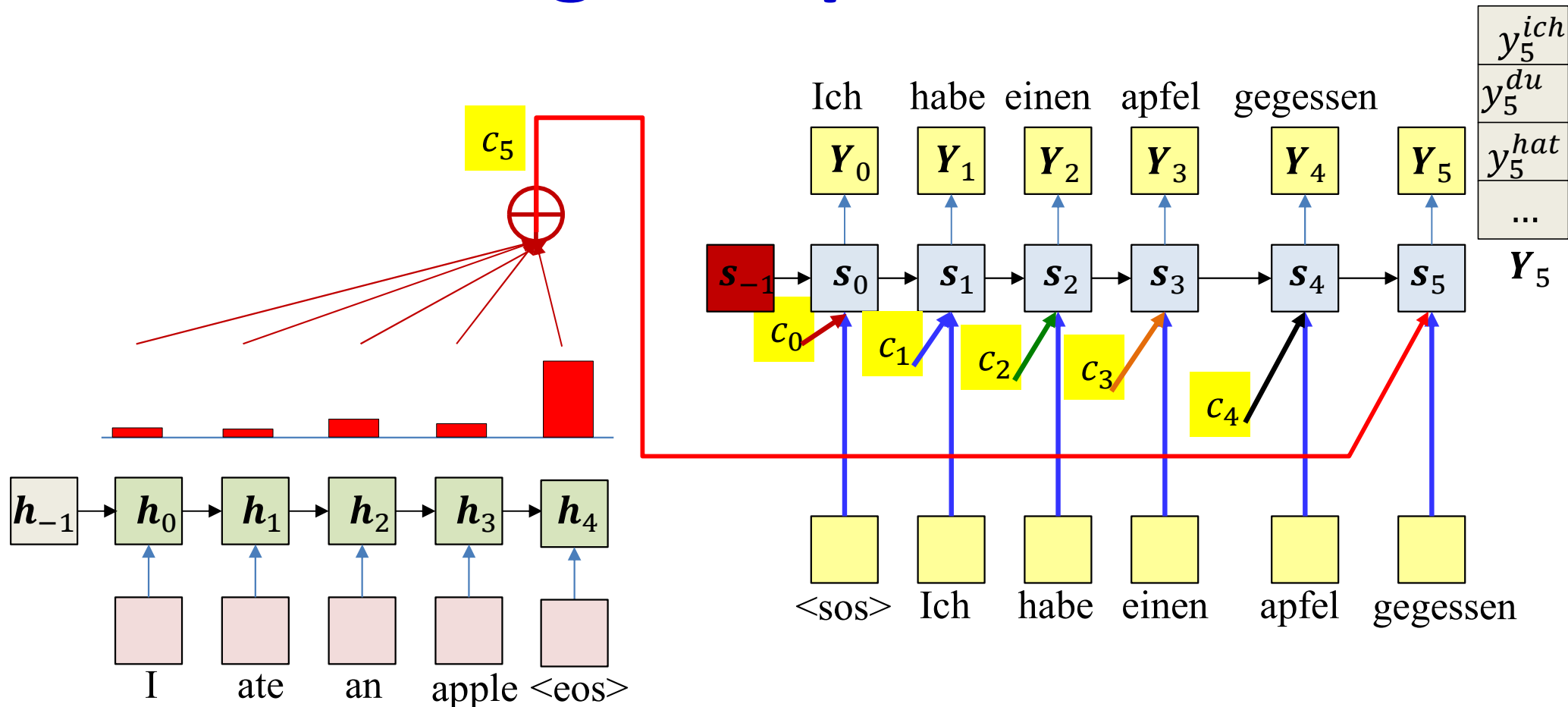


$$e_i(5) = g(h_i, s_4)$$

$$w_i(5) = \frac{\exp(e_i(5))}{\sum_j \exp(e_j(5))}$$

$$c_5 = \frac{1}{N} \sum_i^N w_i(5) h_i$$

Converting an input: Inference

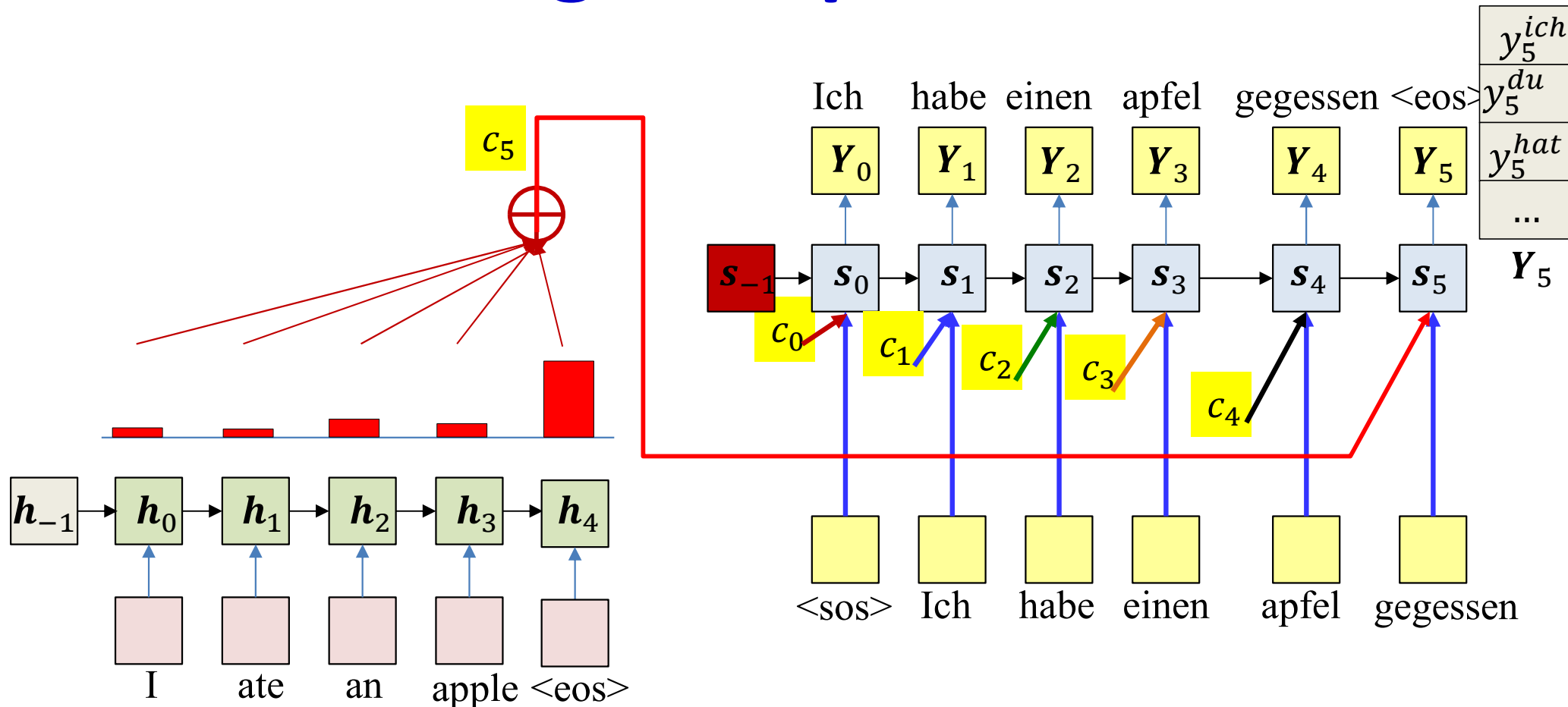


$$e_i(5) = g(h_i, s_4)$$

$$w_i(5) = \frac{\exp(e_i(5))}{\sum_j \exp(e_j(5))}$$

$$c_5 = \frac{1}{N} \sum_i^N w_i(5) h_i$$

Converting an input: Inference



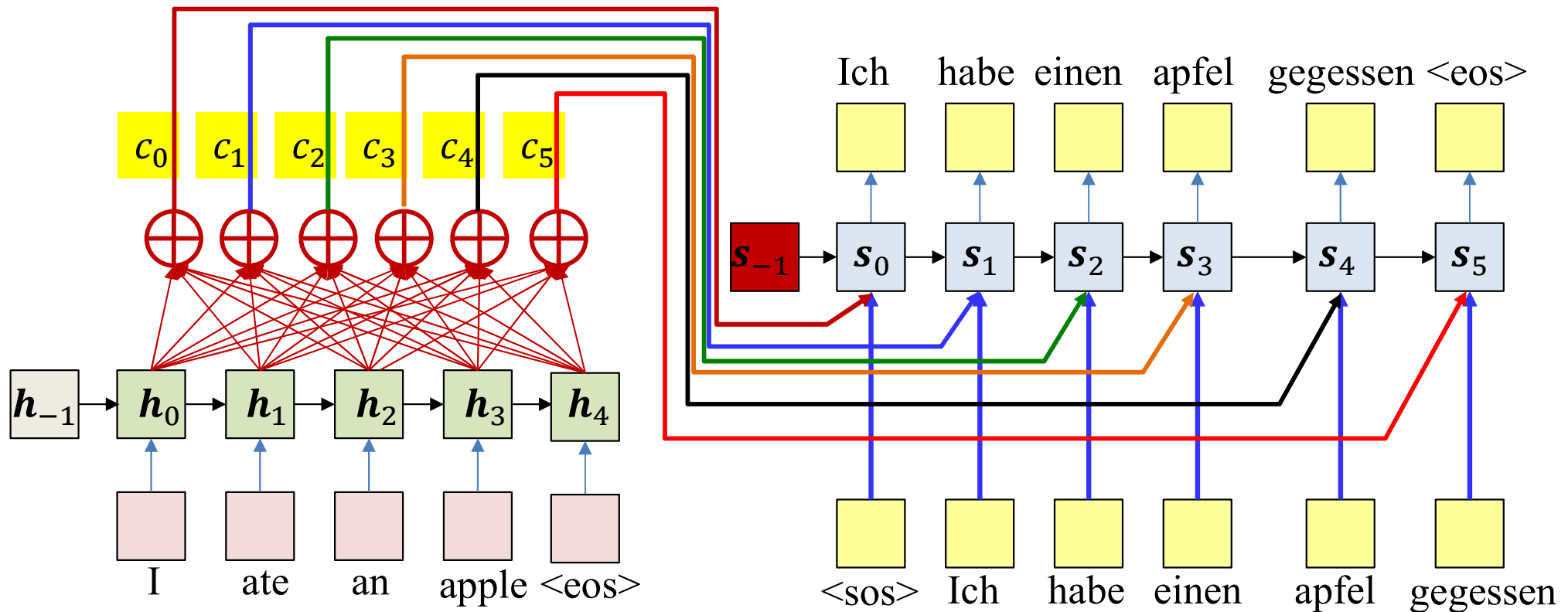
Continue this process until
<eos> is drawn

$$e_i(5) = g(h_i, s_4)$$

$$w_i(5) = \frac{\exp(e_i(5))}{\sum_j \exp(e_j(5))}$$

$$c_5 = \frac{1}{N} \sum_i^N w_i(5) h_i$$

Attention-based decoding



$$e_i(t) = g(\mathbf{h}_i, \mathbf{s}_{t-1})$$

$$w_i(t) = \frac{\exp(e_i(t))}{\sum_j \exp(e_j(t))}$$

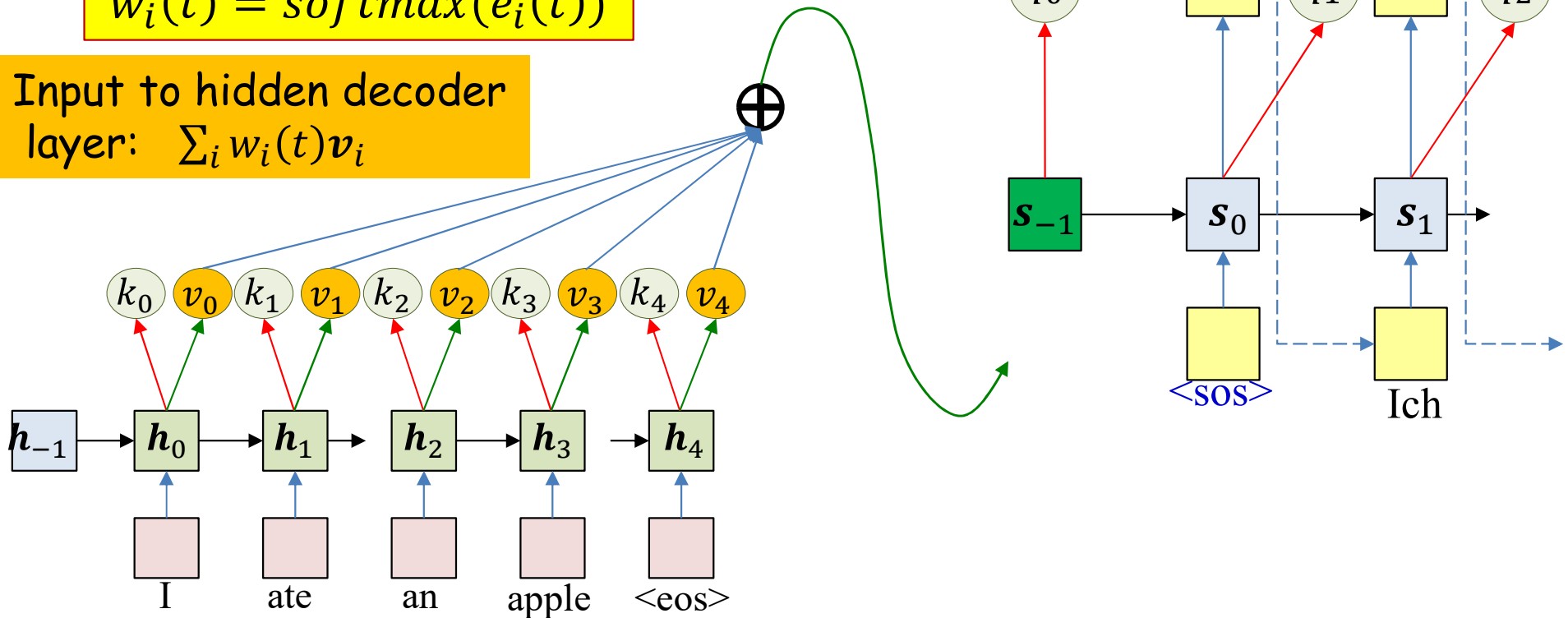
$$\mathbf{c}_t = \frac{1}{N} \sum_i^N w_i(t) \mathbf{h}_i$$

Modification: Query key value

$$e_i(t) = g(k_i, q_t)$$

$$w_i(t) = \text{softmax}(e_i(t))$$

Input to hidden decoder layer:
 $\sum_i w_i(t) v_i$



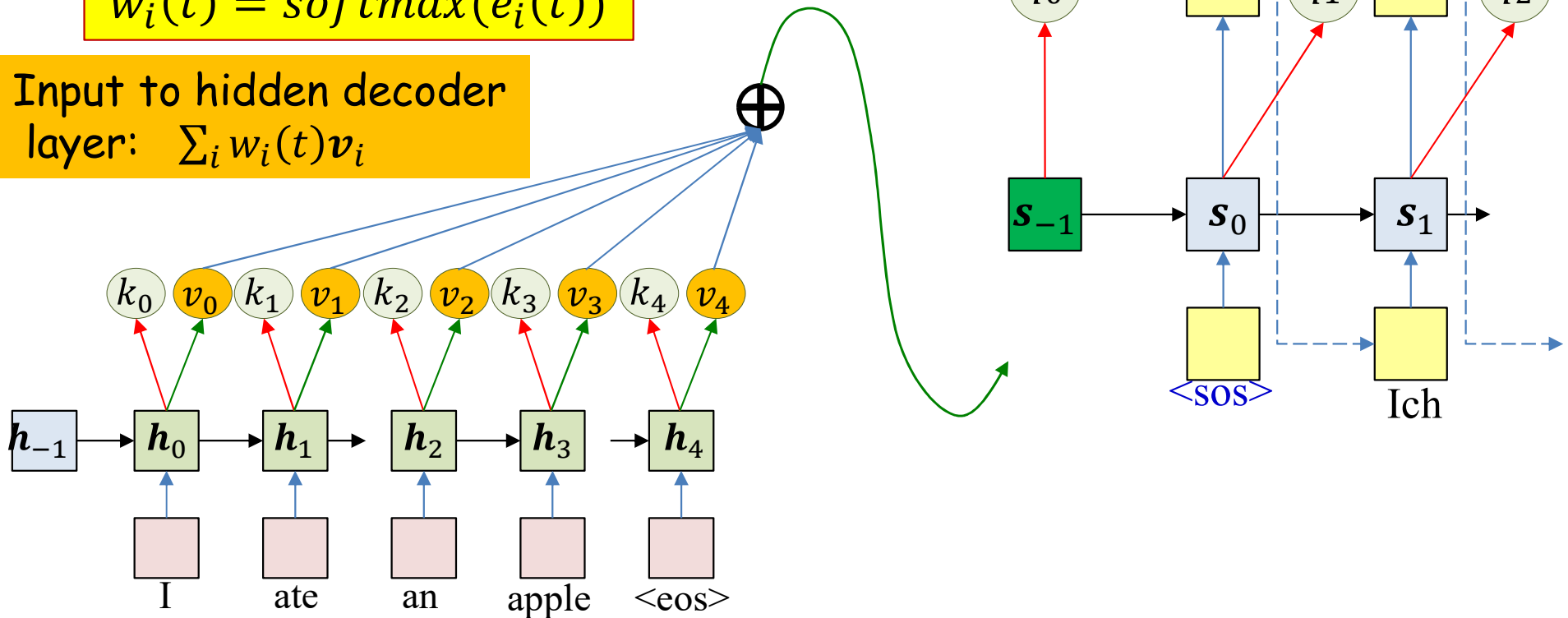
- Encoder outputs an explicit "key" and "value" at each input time
 - Key is used to evaluate the importance of the input at that time, for a given output
- Decoder outputs an explicit "query" at each output time
 - Query is used to evaluate which inputs to pay attention to
- The weight is a function of key and query
- The actual context is a weighted sum of value

Modification: Query key value

$$e_i(t) = g(k_i, q_t)$$

$$w_i(t) = \text{softmax}(e_i(t))$$

Input to hidden decoder layer:
 $\sum_i w_i(t) v_i$



Special case: $k_i = v_i = h_i$

$$q_t = s_{t-1}$$

We will continue using this assumption in the following slides but in practice the query-key-value format is used

Pseudocode

```
# Assuming encoded input
#      (K,V) = [kenc[0]... kenc[T]], [venc[0]... venc[T]]
# is available

t = -1
hout[-1] = 0    # Initial Decoder hidden state
q[0] = 0        # Initial query

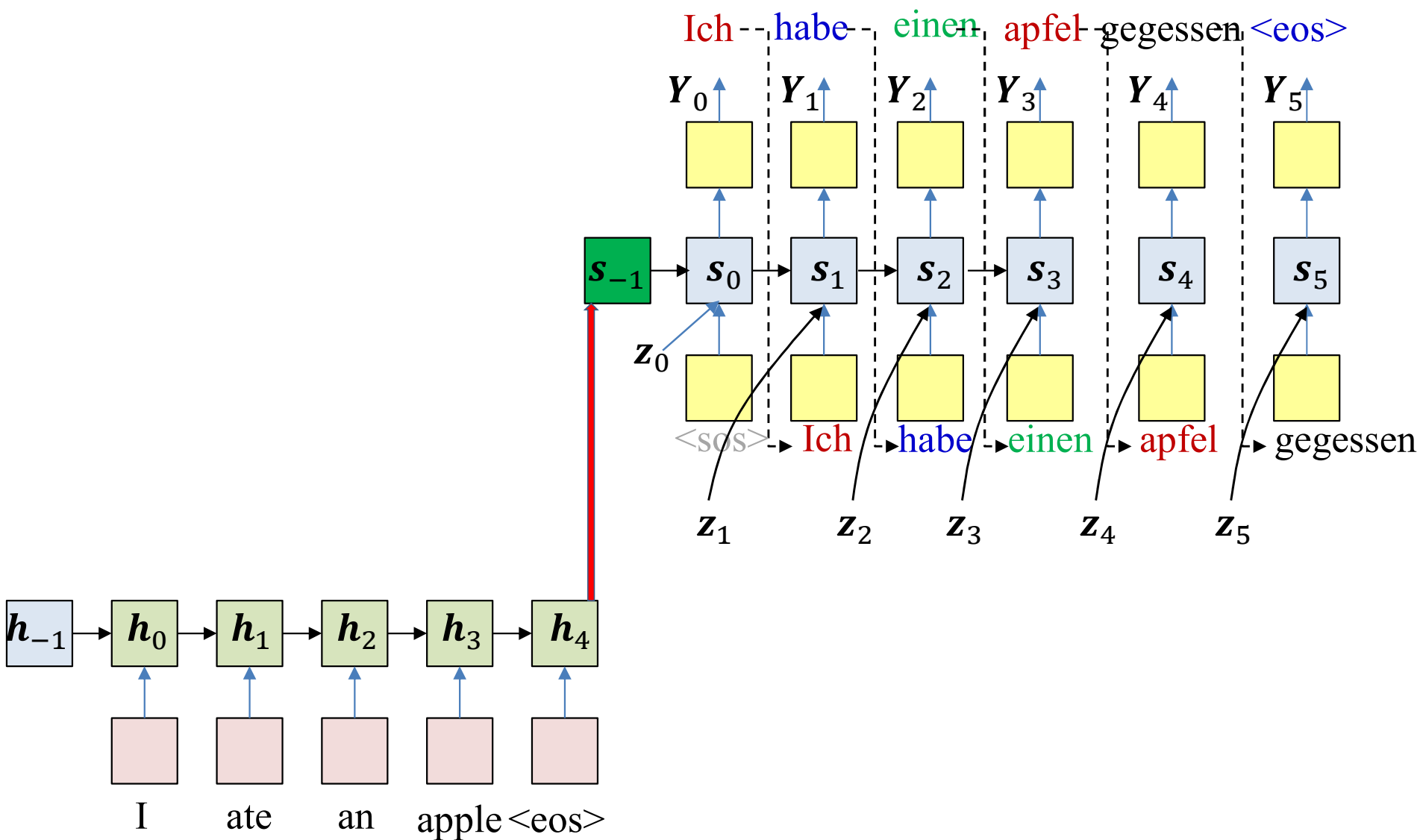
# Note: begins with a "start of sentence" symbol
#      <sos> and <eos> may be identical
Yout[0] = <sos>
do
    t = t+1
    C = compute_context_with_attention(q[t], K, V)
    y[t], hout[t], q[t+1] = RNN_decode_step(hout[t-1], yout[t-1], C)
    yout[t] = generate(y[t]) # Random, or greedy
until yout[t] == <eos>
```

Pseudocode : Computing context with attention

```
# Takes in previous state, encoder states, outputs attention-weighted context
function compute_context_with_attention(q, K, V)
    # First compute attention
    e = []
    for t = 1:T # Length of input
        e[t] = raw_attention(q, K[t])
    end
    maxe = max(e) # subtract max(e) from everything to prevent underflow
    a[1..T] = exp(e[1..T] - maxe) # Component-wise exponentiation
    suma = sum(a) # Add all elements of a
    a[1..T] = a[1..T]/suma

    C = 0
    for t = 1..T
        C += a[t] * V[t]
    end

    return C
```

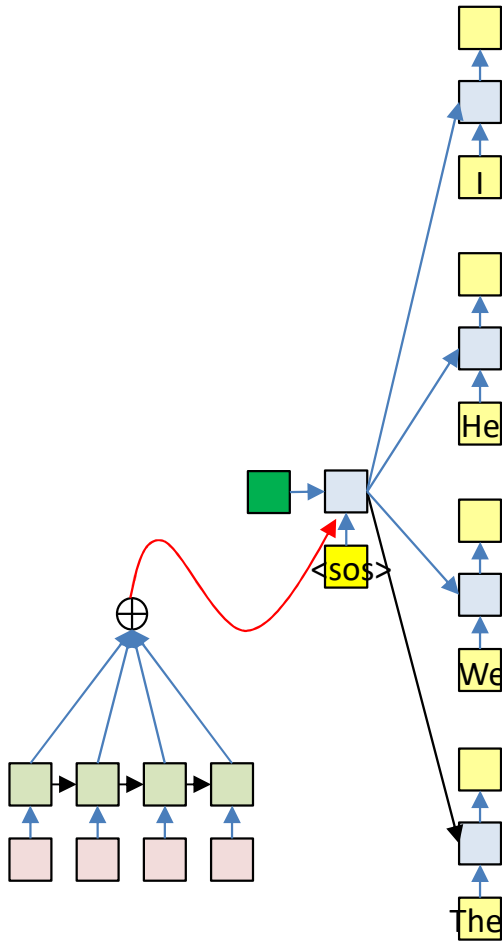


- As before, the objective of drawing: Produce the most likely output (that ends in an <eos>)

$$\operatorname{argmax}_{O_1, \dots, O_L} y_1^{O_1} y_1^{O_2} \dots y_1^{O_L}$$

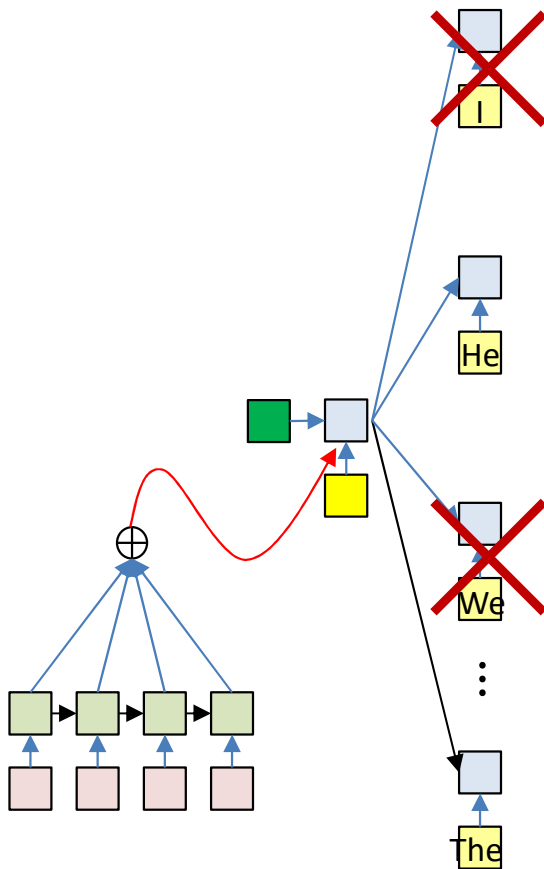
- Simply selecting the most likely symbol at each time may result in suboptimal output

Solution: Multiple choices



- Retain all choices and *fork* the network
 - With every possible word as input

To prevent blowup: Prune

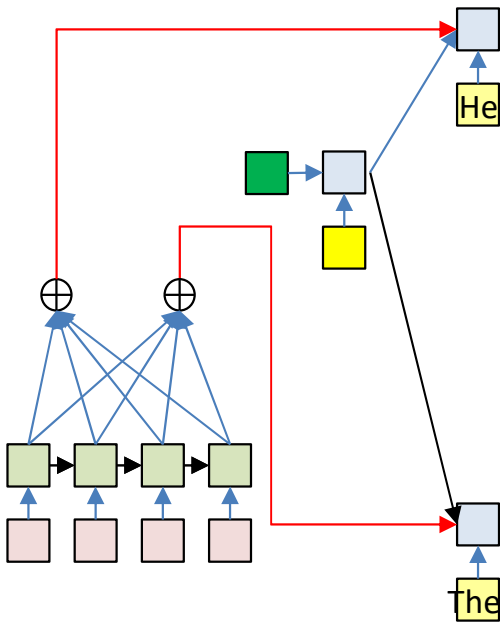


$$Top_K P(O_1 | I_1, \dots, I_N)$$

- **Prune**

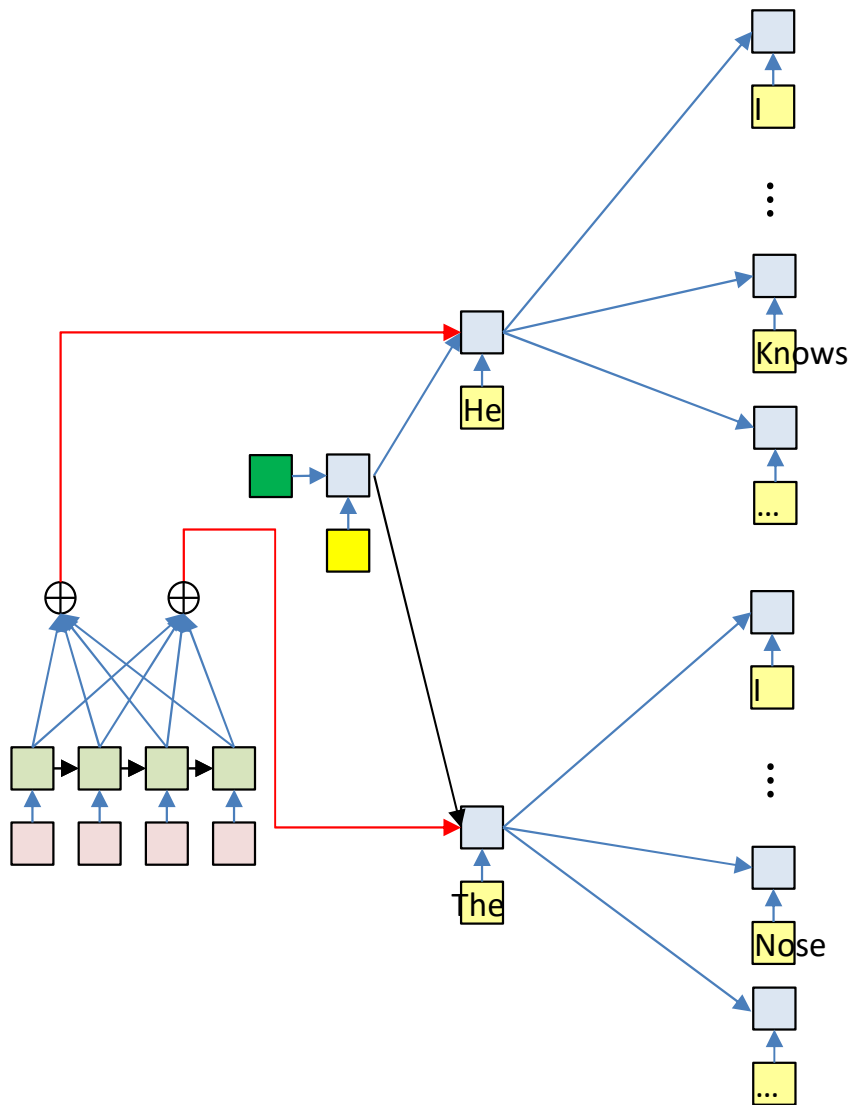
- At each time, retain only the top K scoring forks

Decoding



- At each time, retain only the top K scoring forks

Decoding

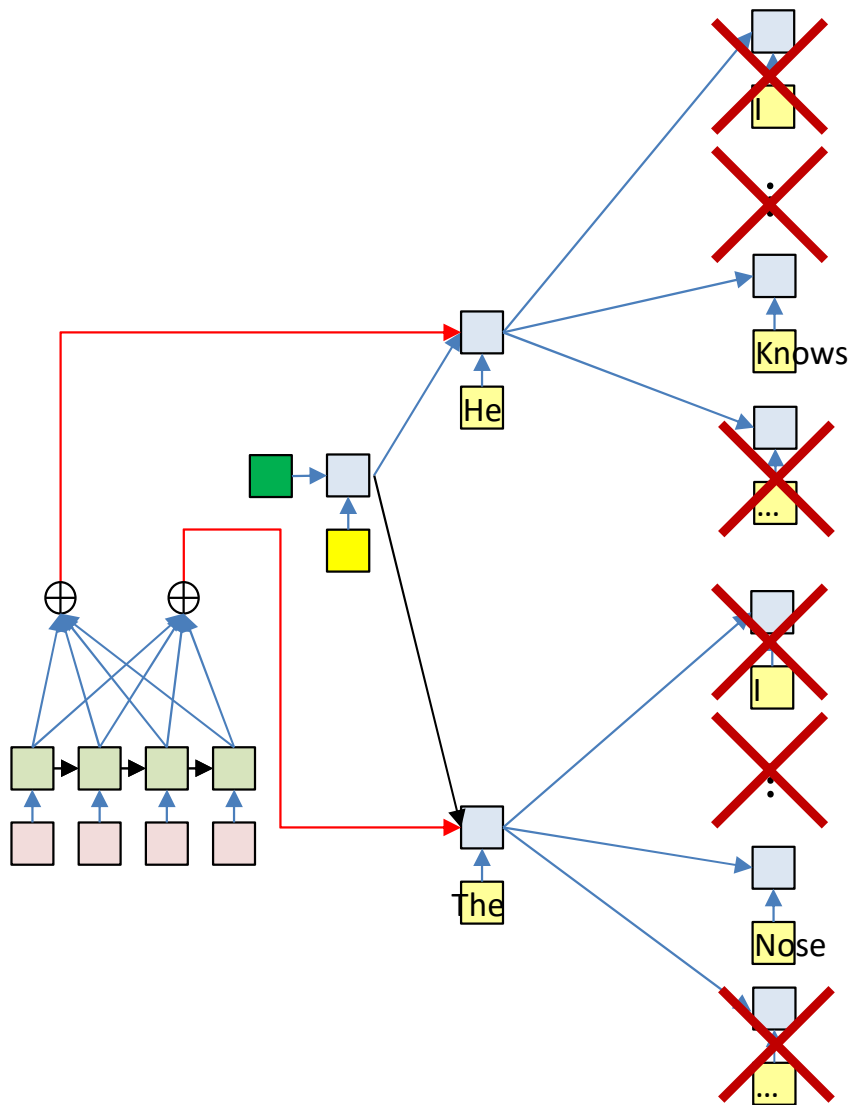


Note: based on product

$$\begin{aligned} & \text{Top}_K P(O_2 O_1 | I_1, \dots, I_N) \\ &= \text{Top}_K P(O_2 | O_1, I_1, \dots, I_N) P(O_1 | I_1, \dots, I_N) \end{aligned}$$

- At each time, retain only the top K scoring forks

Decoding



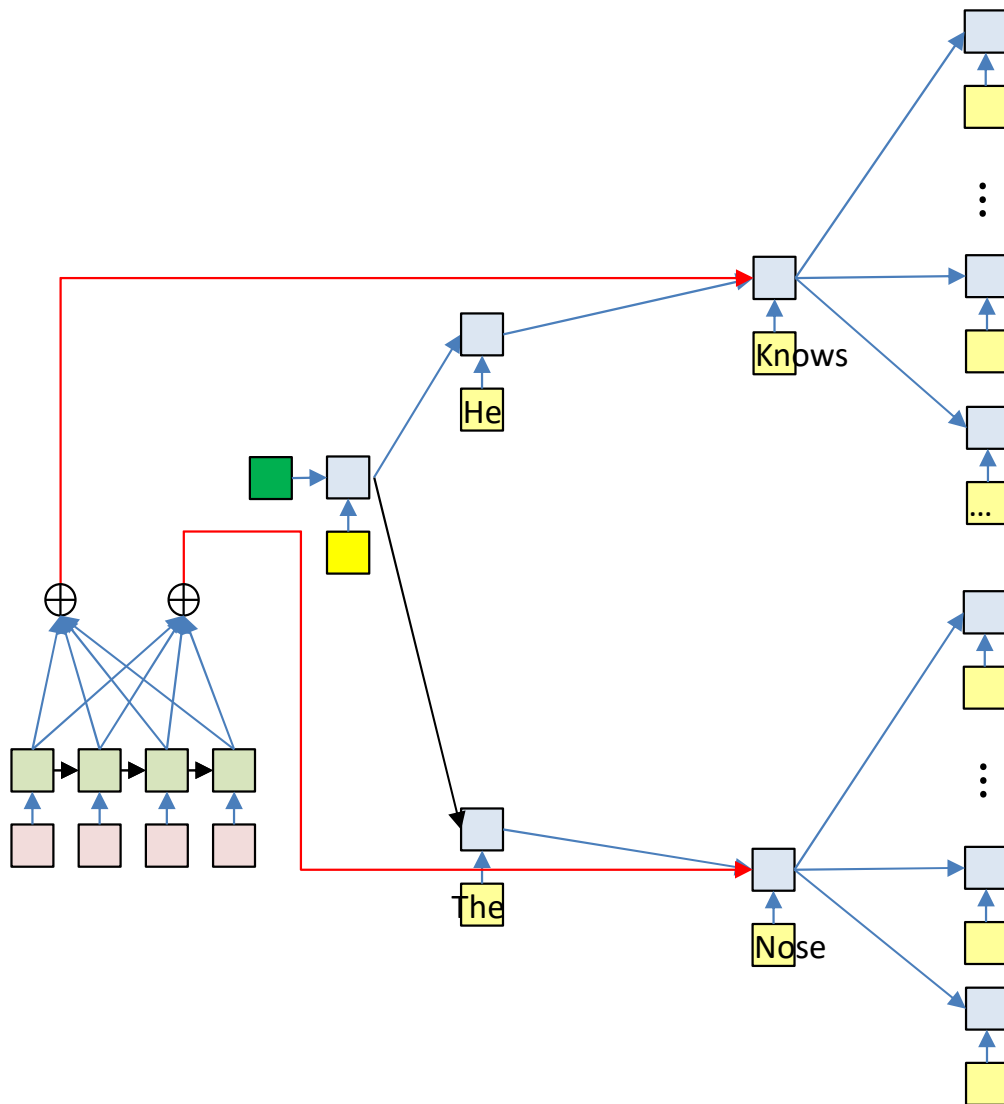
Note: based on product

$$Top_K P(O_2 O_1 | I_1, \dots, I_N)$$

$$= Top_K P(O_2 | O_1, I_1, \dots, I_N) P(O_1 | I_1, \dots, I_N)$$

- At each time, retain only the top K scoring forks

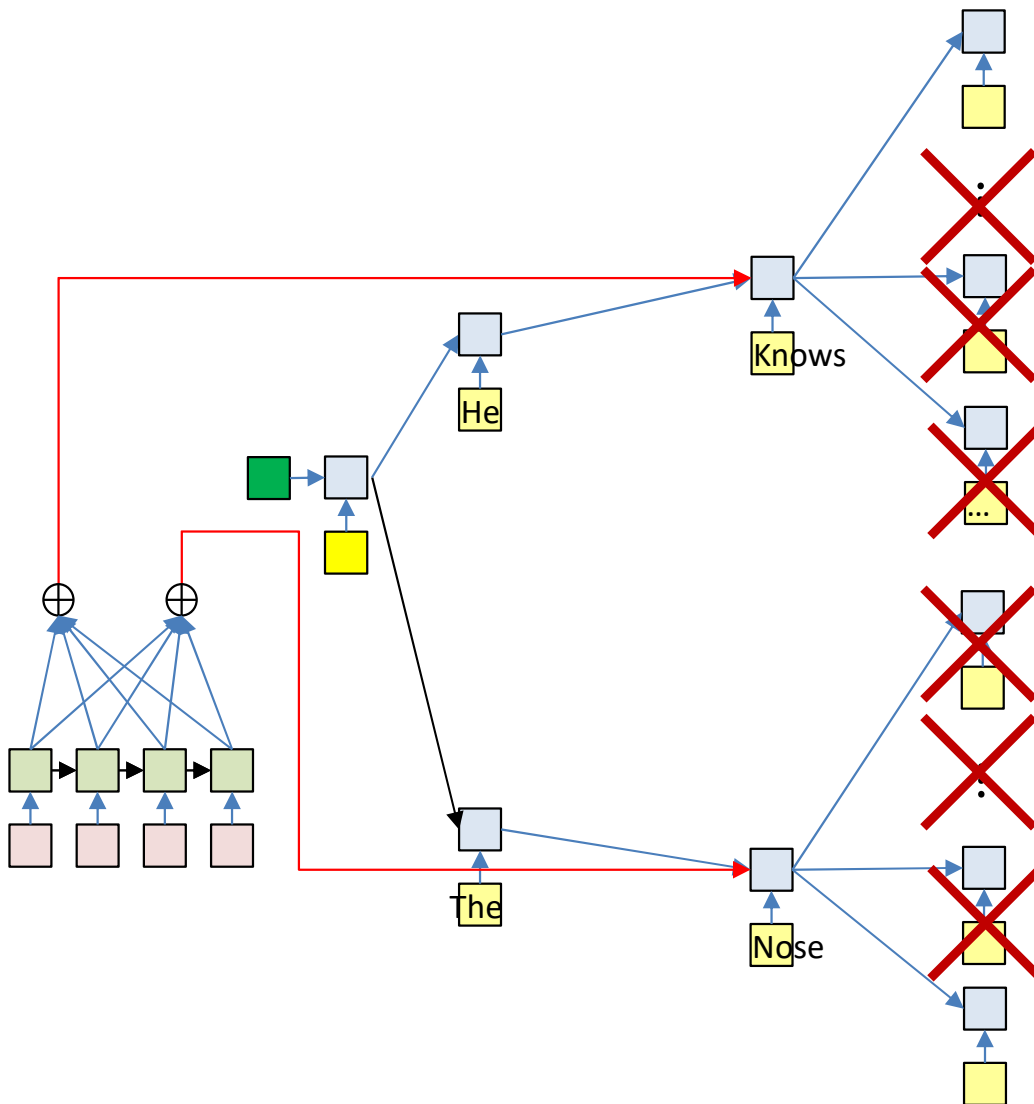
Decoding



$$= \text{Top}_K P(O_2|O_1, I_1, \dots, I_N) \times \\ P(O_2|O_1, I_1, \dots, I_N) \times \\ P(O_1|I_1, \dots, I_N)$$

- At each time, retain only the top K scoring forks

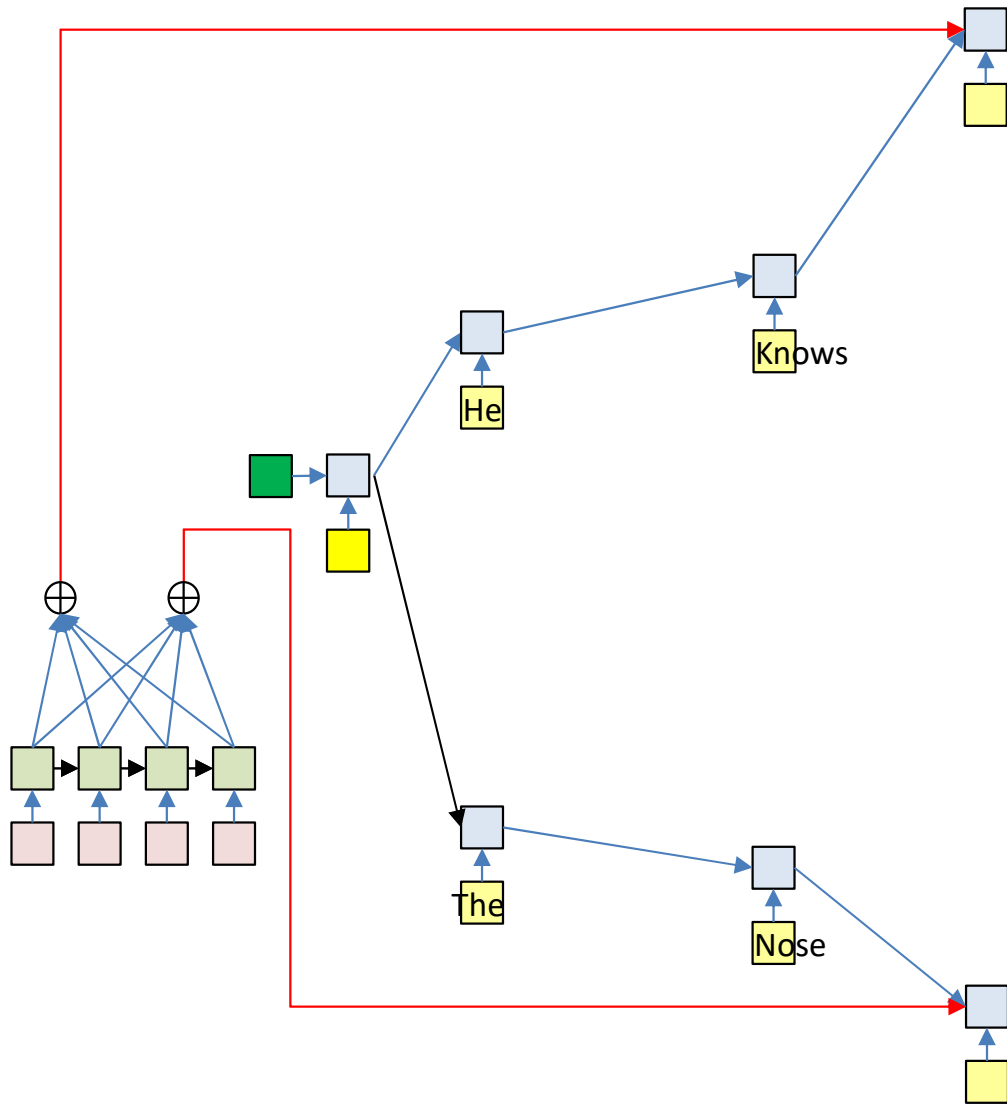
Decoding



$$= \text{Top}_K P(O_2|O_1, I_1, \dots, I_N) \times P(O_2|O_1, I_1, \dots, I_N) \times P(O_1|I_1, \dots, I_N)$$

- At each time, retain only the top K scoring forks

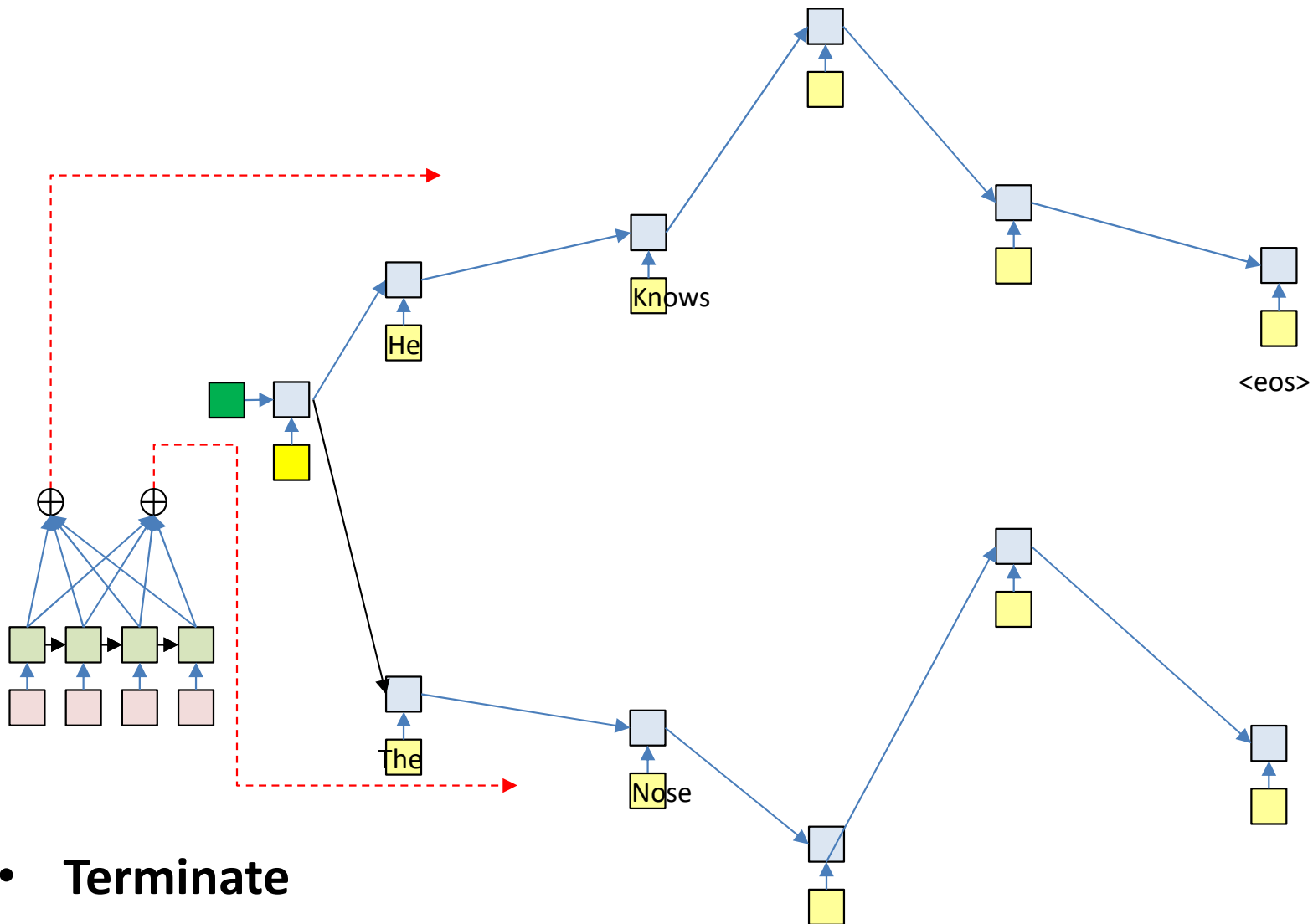
Decoding



$$Top_K \prod_{t=1}^n P(O_t | O_1, \dots, O_{t-1}, I_1, \dots, I_N)$$

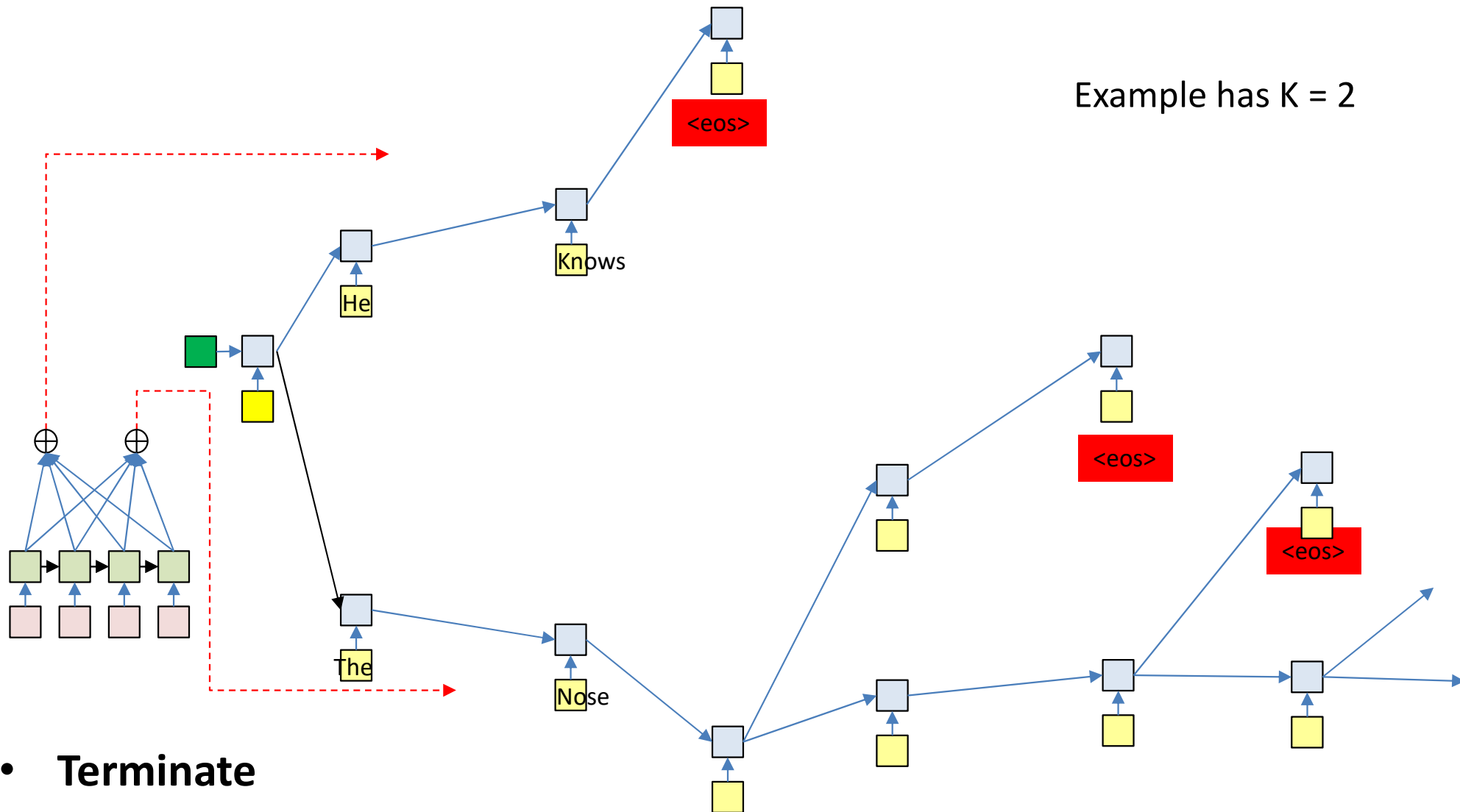
- At each time, retain only the top K scoring forks

Terminate



- **Terminate**
 - When the current most likely path overall ends in <eos>
 - Or continue producing more outputs (each of which terminates in <eos>) to get N-best outputs

Termination: <eos>



- **Terminate**

- Paths cannot continue once the output an <eos>
- So paths may be different lengths
 - Select the most likely sequence ending in <eos> across *all* terminating sequences

Pseudocode: Beam search

```
# Assuming encoder output  $H = h_{in}[1] \dots h_{in}[T]$  is available
path = <sos>
beam = {path}
pathscore = [path] = 1
state[path] = h[0] # initial state (computed using your favorite method)
do # Step forward
    nextbeam = {}
    nextpathscore = []
    nextstate = {}
    for path in beam:
        cfin = path[end]
        hpath = state[path]
        C = compute_context_with_attention(hpath, H)
        y, h = RNN_decode_step(hpath, cfin, C)
        for c in Symbolset
            newpath = path + c
            nextstate[newpath] = h
            nextpathscore[newpath] = pathscore[path]*y[c]
            nextbeam += newpath # Set addition
        end
    end
    beam, pathscore, state, bestpath = prune(nextstate, nextpathscore, nextbeam)
until bestpath[end] = <eos>
```

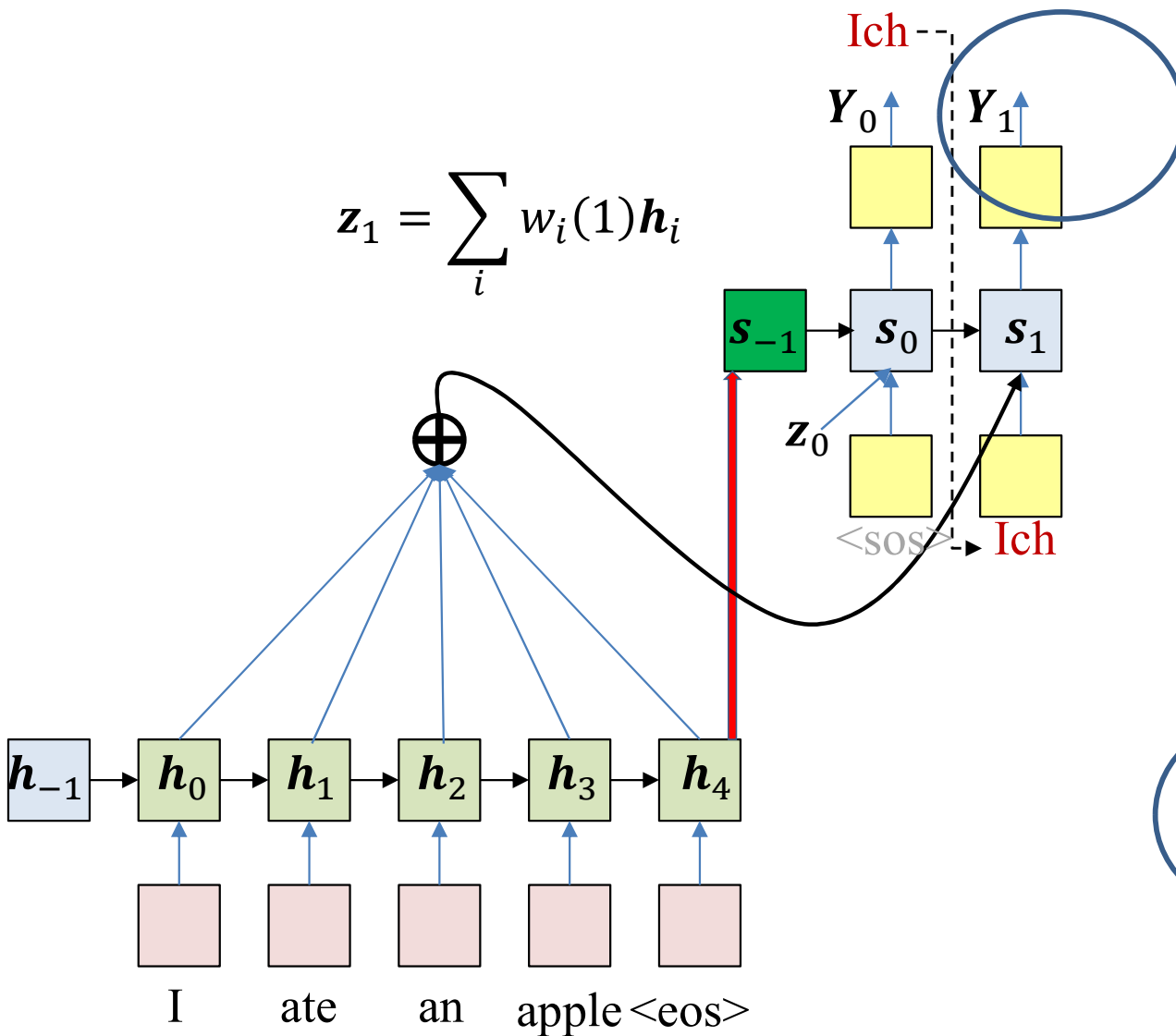
Pseudocode: Beam search

```
# Assuming encoder output  $H = h_{in}[1] \dots h_{in}[T]$  is available
path = <sos>
beam = {path}
pathscore = [path] = 1
state[path] = h[0] # computed using your favorite method
context[path] = compute_context_with_attention(h[0], H)
do # Step forward
    nextbeam = {}
    nextpathscore = []
    nextstate = {}
    nextcontext = {}
    for path in beam:
        cfin = path[end]
        hpath = state[path]
        C = context[path]
        y, h = RNN_decode_step(hpath, cfin, C)
        nextC = compute_context_with_attention(h, H)
        for c in Symbolset
            newpath = path + c
            nextstate[newpath] = h
            nextcontext[newpath] = nextC
            nextpathscore[newpath] = pathscore[path]*y[c]
            nextbeam += newpath # Set addition
        end
    end
    beam, pathscore, state, context, bestpath =
        prune (nextstate, nextpathscore, nextbeam, nextcontext)
until bestpath[end] = <eos>
```

Slightly more efficient.

Does not perform redundant context computation

What does the attention learn?



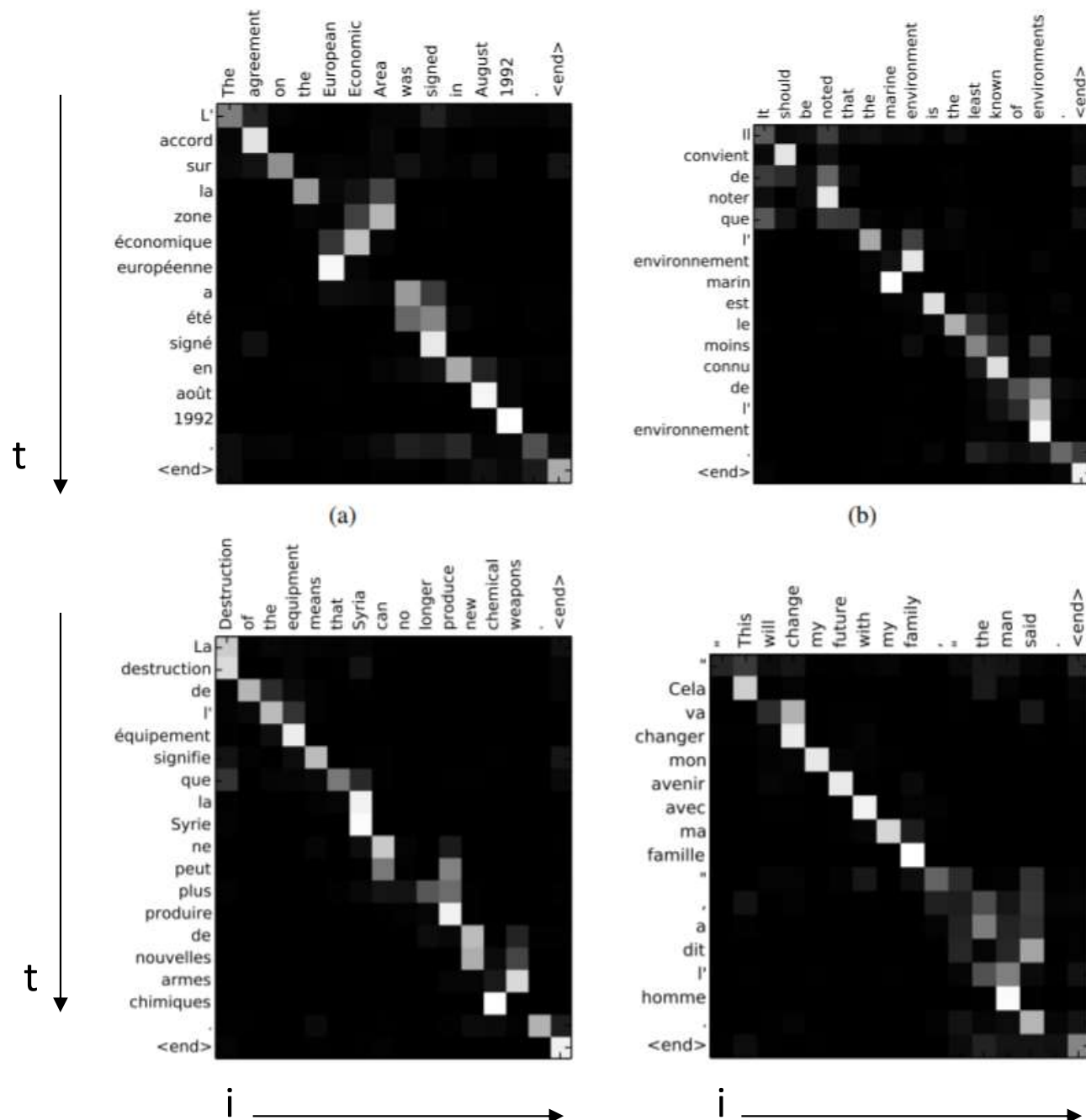
$$g(h_i, s_0) = h_i^T W_g s_0$$

$$e_i(1) = g(h_i, s_0)$$

$$w_i(1) = \frac{\exp(e_i(1))}{\sum_j \exp(e_j(1))}$$

- The key component of this model is the attention weight
 - It captures the relative importance of each position in the input to the current output

“Alignments” example: Bahdanau et al.



Plot of $w_i(t)$

Color shows value (white is larger)

Note how most important input words for any output word get automatically highlighted

The general trend is somewhat linear because word order is roughly similar in both languages

Translation Examples

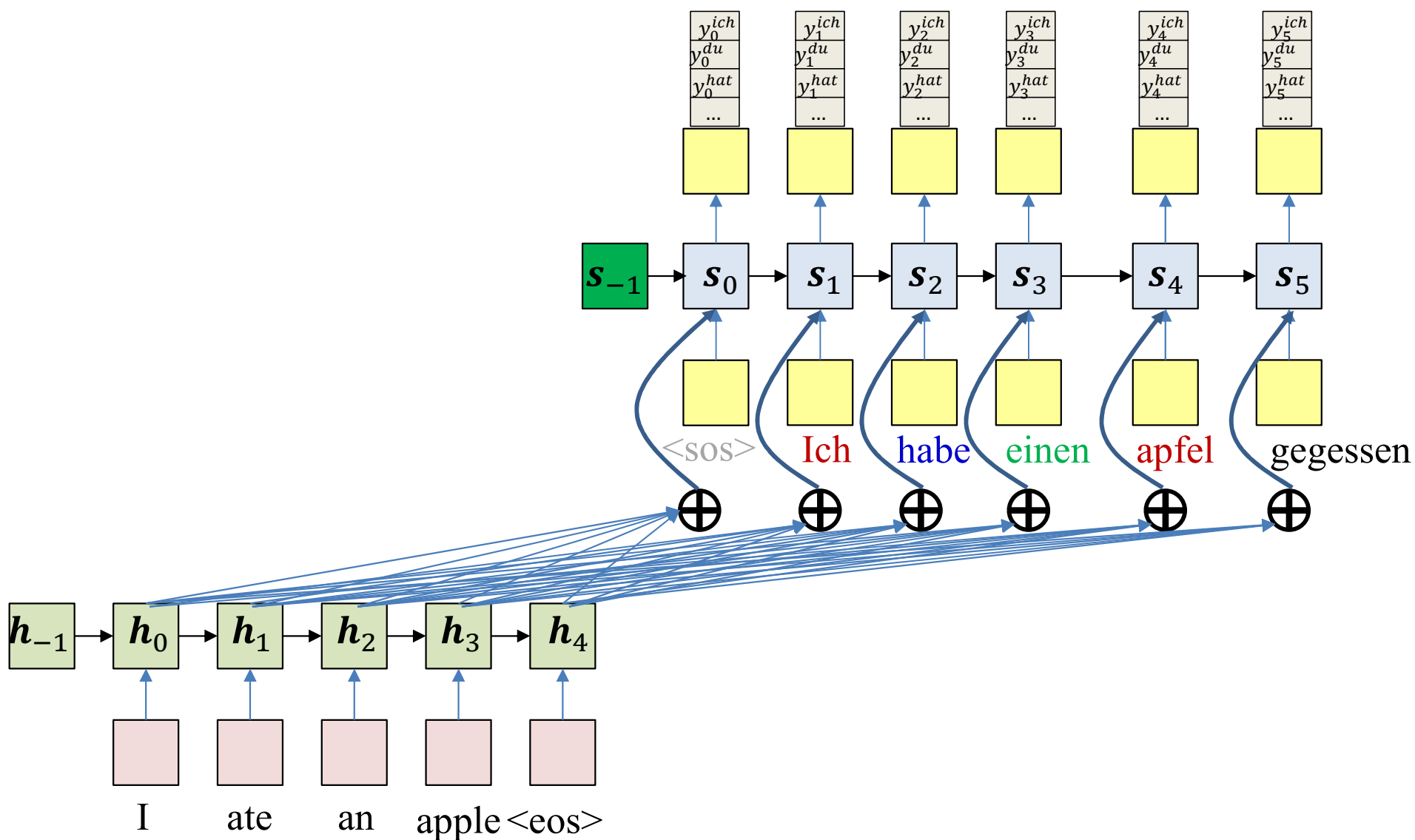
| | |
|------------------|--|
| Source | An admitting privilege is the right of a doctor to admit a patient to a hospital or a medical centre to carry out a diagnosis or a procedure, based on his status as a health care worker at a hospital. |
| Reference | Le privilège d'admission est le droit d'un médecin, en vertu de son statut de membre soignant d'un hôpital, d'admettre un patient dans un hôpital ou un centre médical afin d'y délivrer un diagnostic ou un traitement. |
| RNNenc-50 | Un privilège d'admission est le droit d'un médecin de reconnaître un patient à l'hôpital ou un centre médical d'un diagnostic ou de prendre un diagnostic en fonction de son état de santé. |
| RNNsearch-50 | Un privilège d'admission est le droit d'un médecin d'admettre un patient à un hôpital ou un centre médical pour effectuer un diagnostic ou une procédure, selon son statut de travailleur des soins de santé à l'hôpital. |
| Google Translate | Un privilège admettre est le droit d'un médecin d'admettre un patient dans un hôpital ou un centre médical pour effectuer un diagnostic ou une procédure, fondée sur sa situation en tant que travailleur de soins de santé dans un hôpital. |

| | |
|------------------|--|
| Source | This kind of experience is part of Disney's efforts to "extend the lifetime of its series and build new relationships with audiences via digital platforms that are becoming ever more important," he added. |
| Reference | Ce type d'expérience entre dans le cadre des efforts de Disney pour "étendre la durée de vie de ses séries et construire de nouvelles relations avec son public grâce à des plateformes numériques qui sont de plus en plus importantes", a-t-il ajouté. |
| RNNenc-50 | Ce type d'expérience fait partie des initiatives du Disney pour "prolonger la durée de vie de ses nouvelles et de développer des liens avec les lecteurs numériques qui deviennent plus complexes. |
| RNNsearch-50 | Ce genre d'expérience fait partie des efforts de Disney pour "prolonger la durée de vie de ses séries et créer de nouvelles relations avec des publics via des plateformes numériques de plus en plus importantes", a-t-il ajouté. |
| Google Translate | Ce genre d'expérience fait partie des efforts de Disney à "étendre la durée de vie de sa série et construire de nouvelles relations avec le public par le biais des plates-formes numériques qui deviennent de plus en plus important", at-il ajouté. |

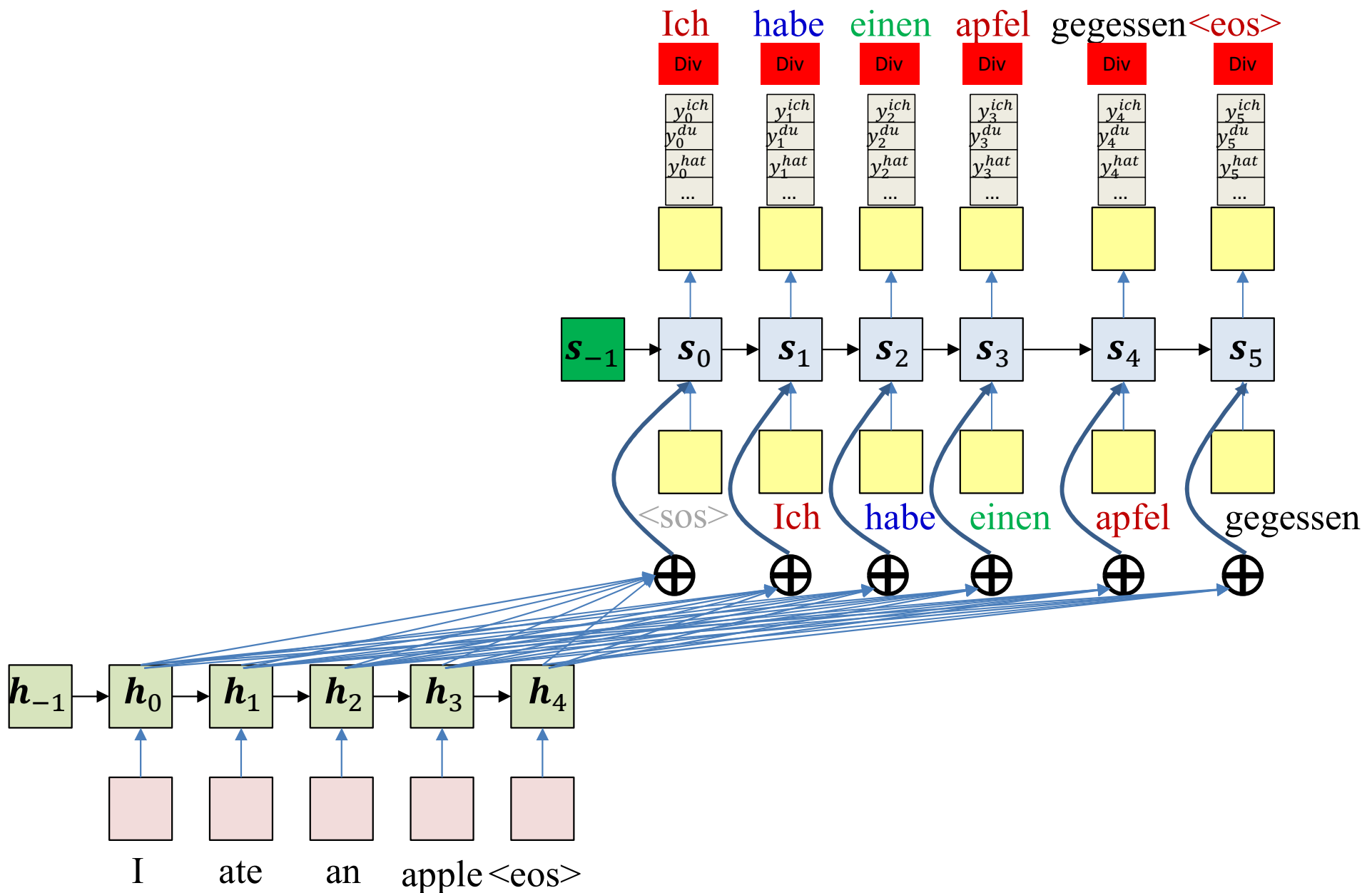
- Bahdanau et al. 2016

Training the network

- We have seen how a trained network can be used to compute outputs
 - Convert one sequence to another
- Lets consider training..

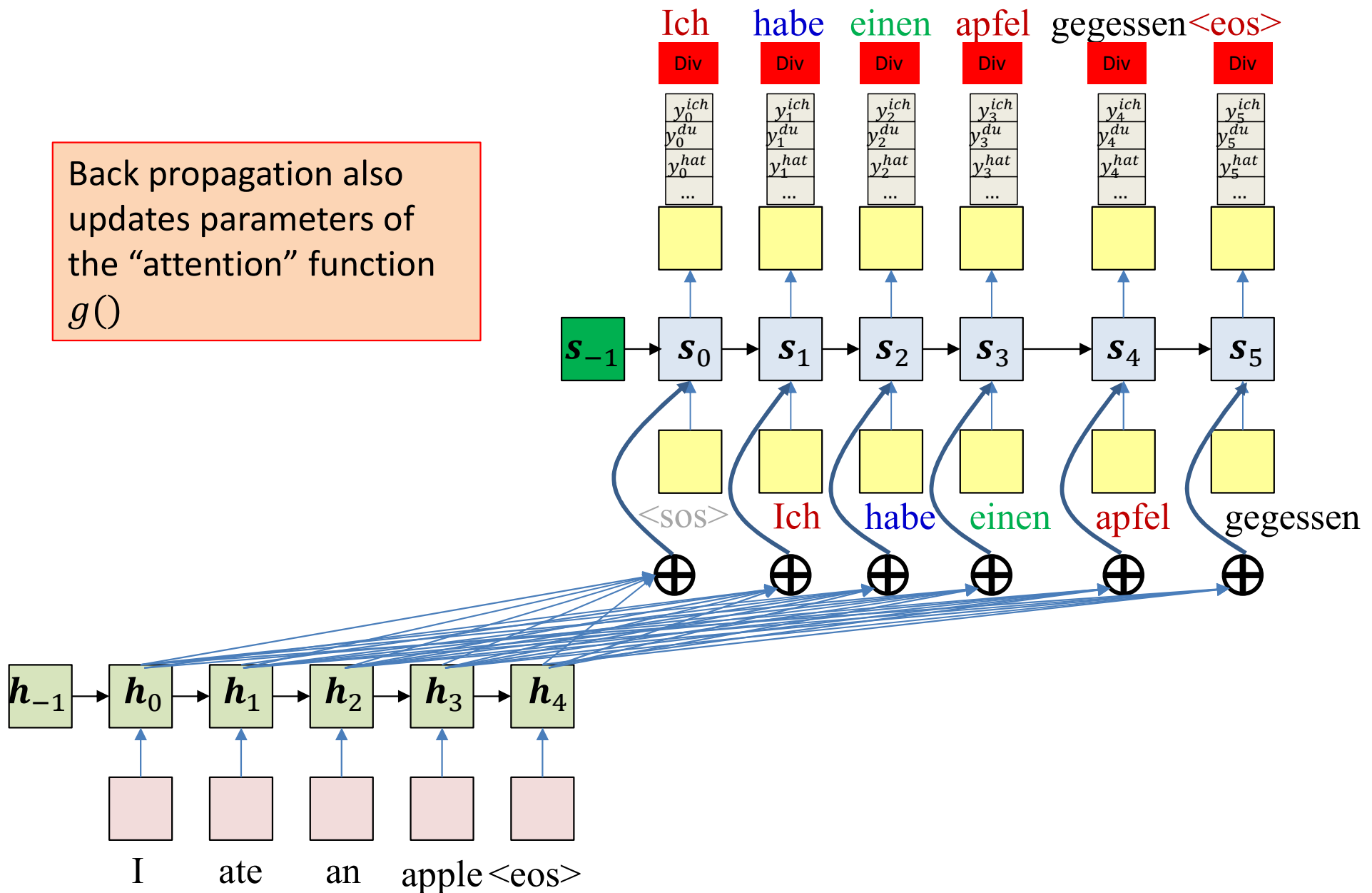


- Given training input (source sequence, target sequence) pairs
- **Forward pass:** Pass the actual Pass the input sequence through the encoder
 - At each time the output is a probability distribution over words



- **Backward pass:** Compute a divergence between target output and output distributions
 - Backpropagate derivatives through the network

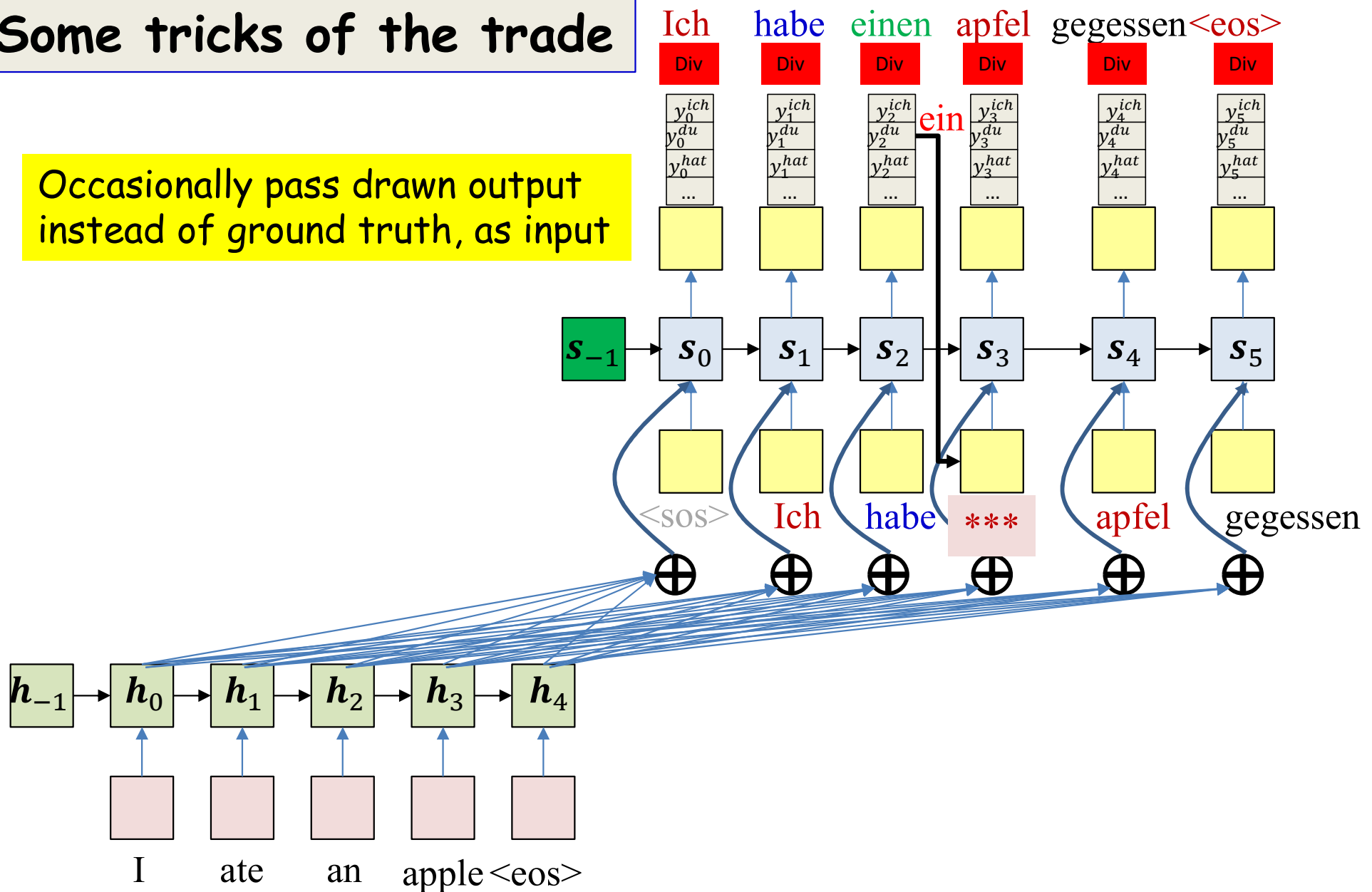
Back propagation also updates parameters of the “attention” function $g()$



- **Backward pass:** Compute a divergence between target output and output distributions
 - Backpropagate derivatives through the network

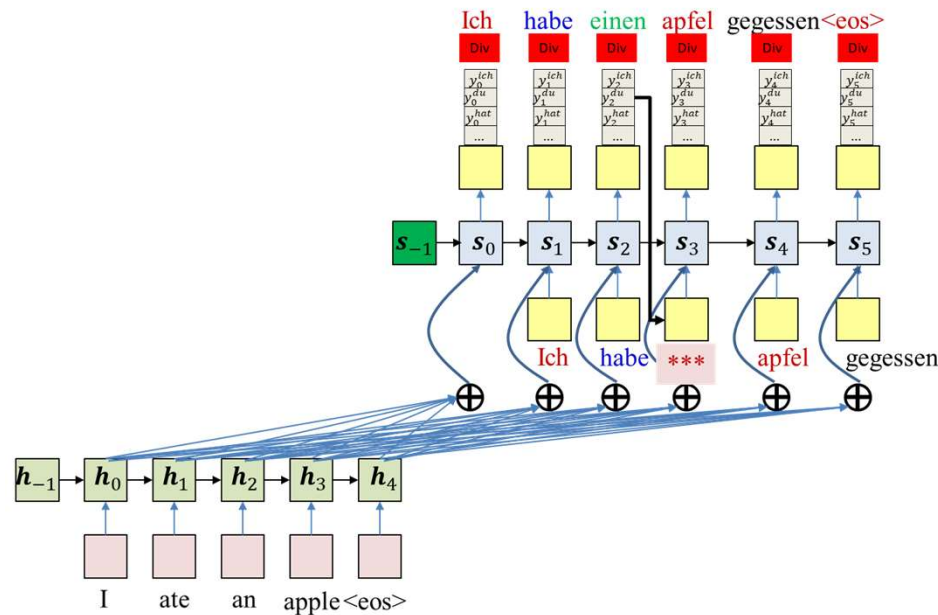
Some tricks of the trade

Occasionally pass drawn output instead of ground truth, as input



- **Backward pass:** Compute a divergence between target output and output distributions
 - Backpropagate derivatives through the network

Gumbel Noise trick

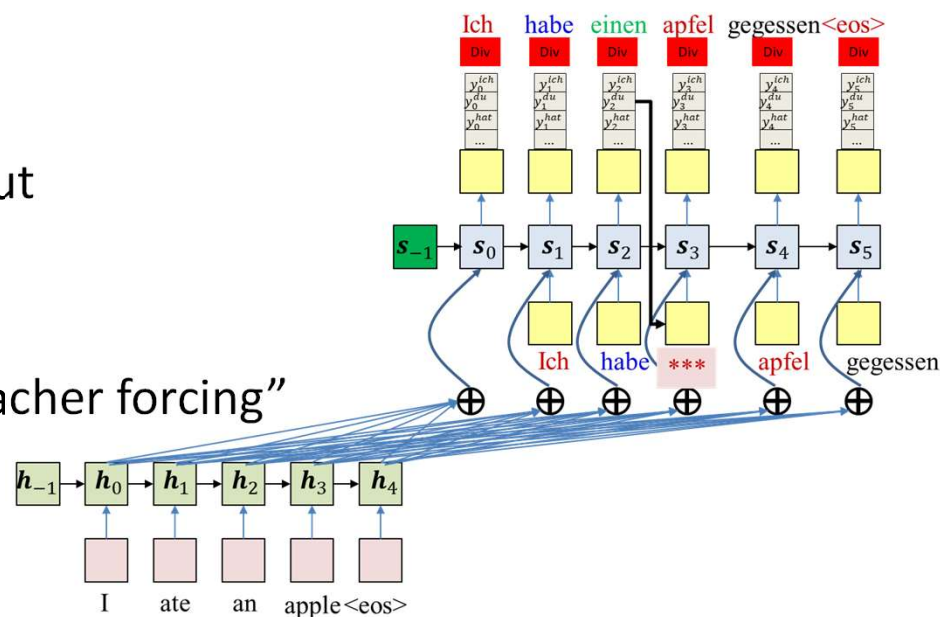


- Sampling is not differentiable
- The “Gumbel noise” trick:
 - “Reparametrization” :

$$RandomSample(Y) = \underset{i}{\operatorname{argmax}}(G_i + \log(Y))$$
 - G_i is drawn from the standard Gumbel distribution $Gumbel(0,1)$
- The “argmax” can be replaced by a “softmax”, making the process differentiable w.r.t. network outputs
 - $decoderoutput(t) = softmax(G_i + \log(Y(t)))$
- $\nabla_{Y(t)} decoderoutput(t)$ is employed in the chain rule to pass derivatives from $t+1$ back to $Y(t)$

Tricks of the trade...

- Teacher forcing:
 - Ideally we only use the decoder output during inference
 - This will not be stable
 - Passing in ground truth instead is “teacher forcing”



- Sampling the output:
 - Sample the system output and
 - as input during training for only some of the time
- The “Gumbel noise” trick:
 - Sampling is not differentiable, and gradients cannot be passed through it
 - The “Gumbel noise” approach recasts sampling as computing the argmax of a Gumbel distribution, with the network output as parameters
 - The “argmax” can be replaced by a “softmax”, making the process differentiable w.r.t. network outputs

Various extensions

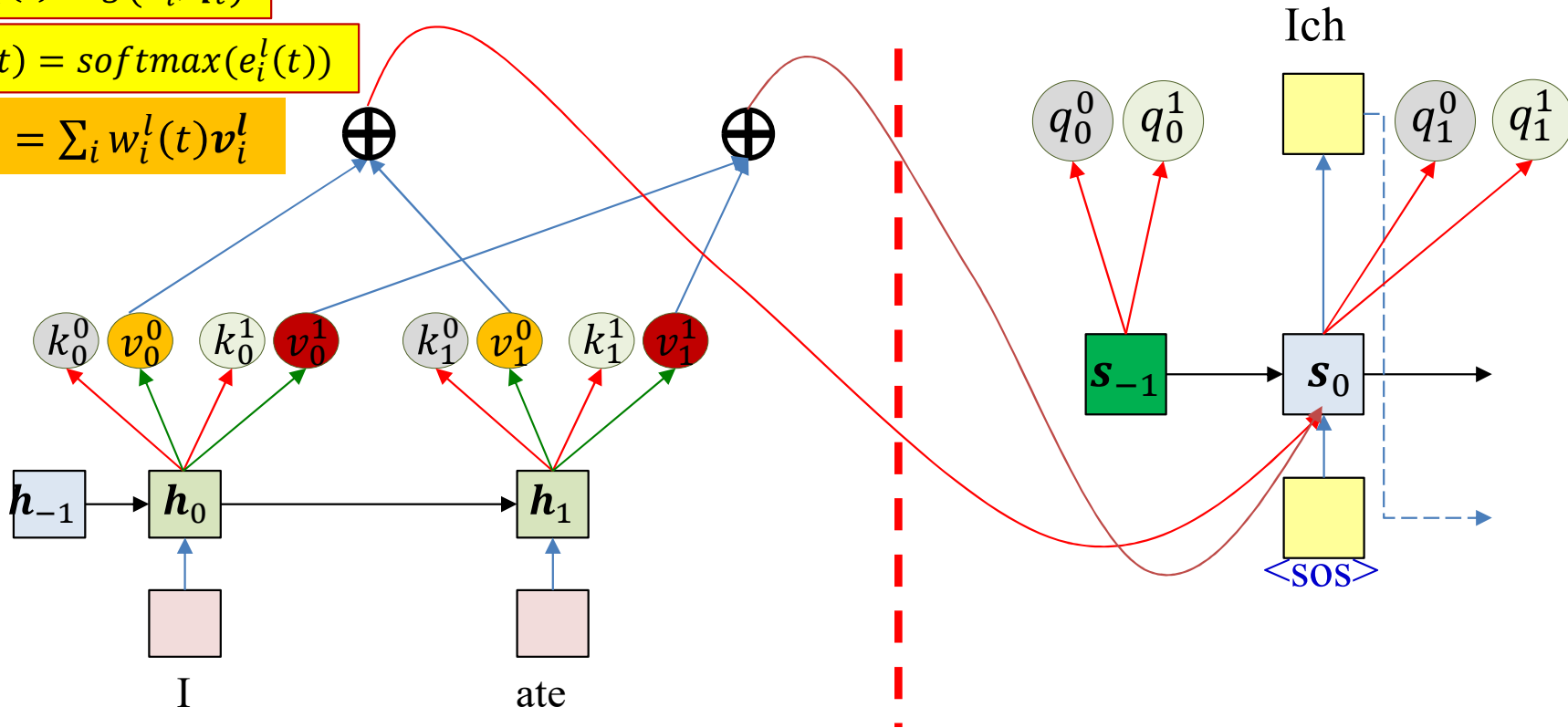
- Bidirectional processing of input sequence
 - Bidirectional networks in encoder
 - E.g. “Neural Machine Translation by Jointly Learning to Align and Translate”, Bahdanau et al. 2016
- Attention: Local attention vs global attention
 - E.g. “Effective Approaches to Attention-based Neural Machine Translation”, Luong et al., 2015
 - Other variants

Extensions: Multihead attention

$$e_i^l(t) = g(k_i^l, q_t^l)$$

$$w_i^l(t) = \text{softmax}(e_i^l(t))$$

$$C_t^l = \sum_i w_i^l(t) v_i^l$$



- Have multiple query/key/value sets.
 - Each attention “head” uses one of these sets
 - The combined contexts from all heads are passed to the decoder
- Each “attender” focuses on a different aspect of the input that’s important for the decode

Some impressive results..

- Attention-based models are currently responsible for the state of the art in many sequence-conversion systems
 - Machine translation
 - Input: Text in source language
 - Output: Text in target language
 - Speech recognition
 - Input: Speech audio feature vector sequence
 - Output: Transcribed word or character sequence

Attention models in image captioning



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

- “Show attend and tell: Neural image caption generation with visual attention”, Xu et al., 2016
- Encoder network is a convolutional neural network
 - Filter outputs at each location are the equivalent of \mathbf{h}_i in the regular sequence-to-sequence model

In closing

- Have looked at various forms of sequence-to-sequence models
- Generalizations of recurrent neural network formalisms
- For more details, please refer to papers
 - Post on piazza if you have questions
- Will appear in HW4: **Speech recognition with attention models**