CMU 11-785   Introduction to Deep Learning, Fall 2020

Lecture 18

# Neural Machine Translation & Attention

TAVE Research DL001

Changdae Oh

2021.05.02

# Topics

❖ Seq2Seq for Machine Translation
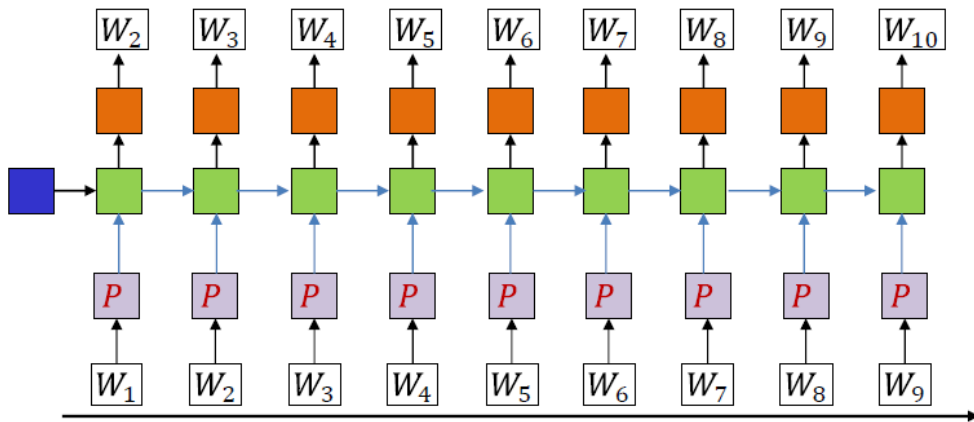
❖ Decoding & Training process

❖ Attention model

# Topics

❖ Seq2Seq for Machine Translation

❖ Decoding & Training process
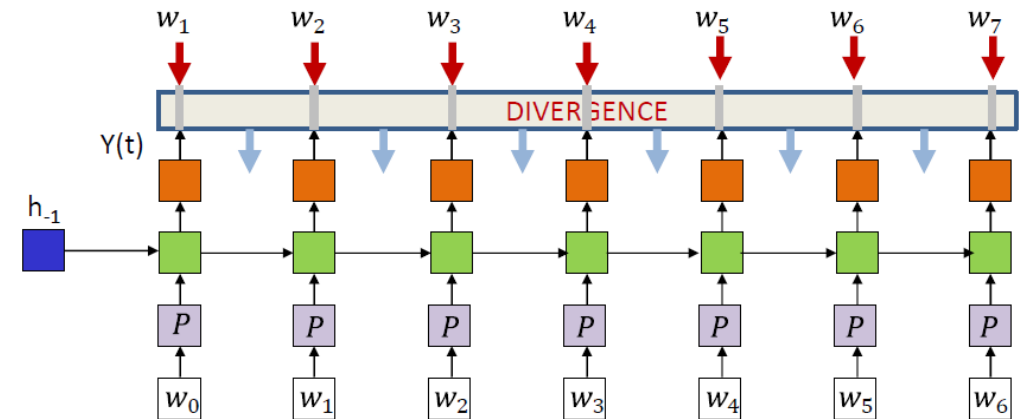
❖ Attention model

# Seq2Seq for Machine Translation

## Recap

### Language Modeling



### Training LM



- Learn a model that can predict the next symbol given a sequence of symbols
- After observing inputs $w_0$ , ... ,$w_k$ (one-hot vectors) it predicts $w_{k+1}$ (probability distribution)
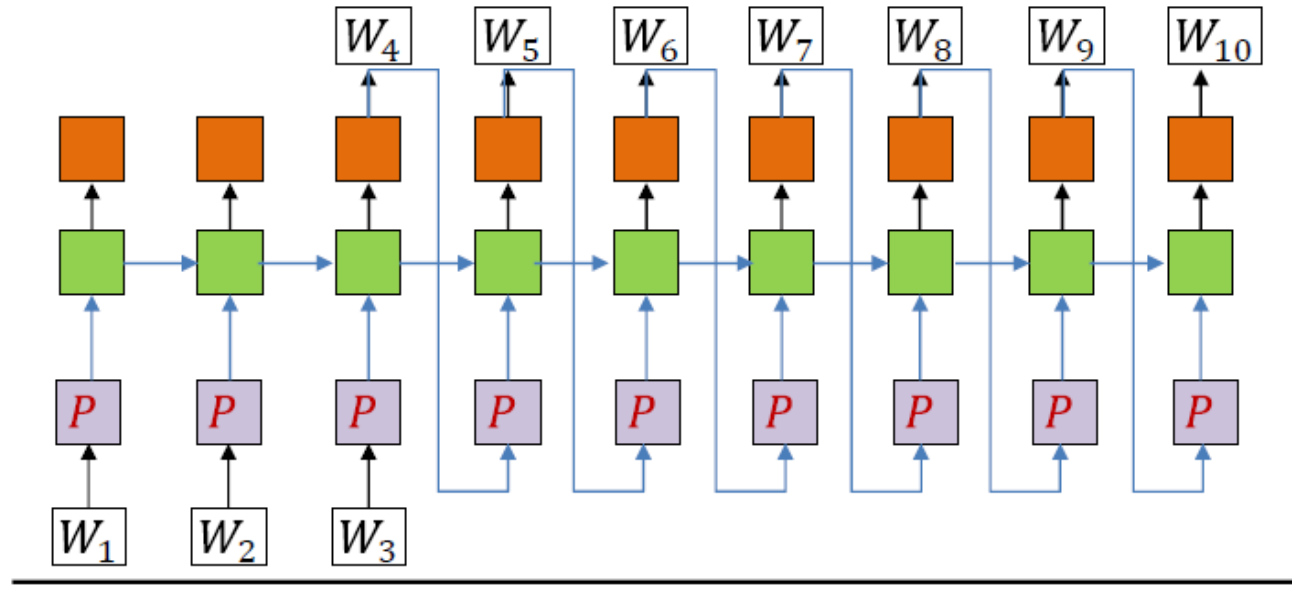
$$Y(t, i) = P(V_i | w_0 \ldots w_{t-1})$$

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

Probability assigned to the correct next word

4
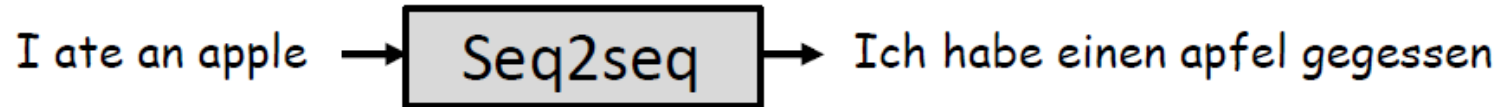
# Seq2Seq for Machine Translation

## Generating Language



- Feed the drawn word as the next word in the series
- Continue until the model draws an <eos>

# Seq2Seq for Machine Translation

## object

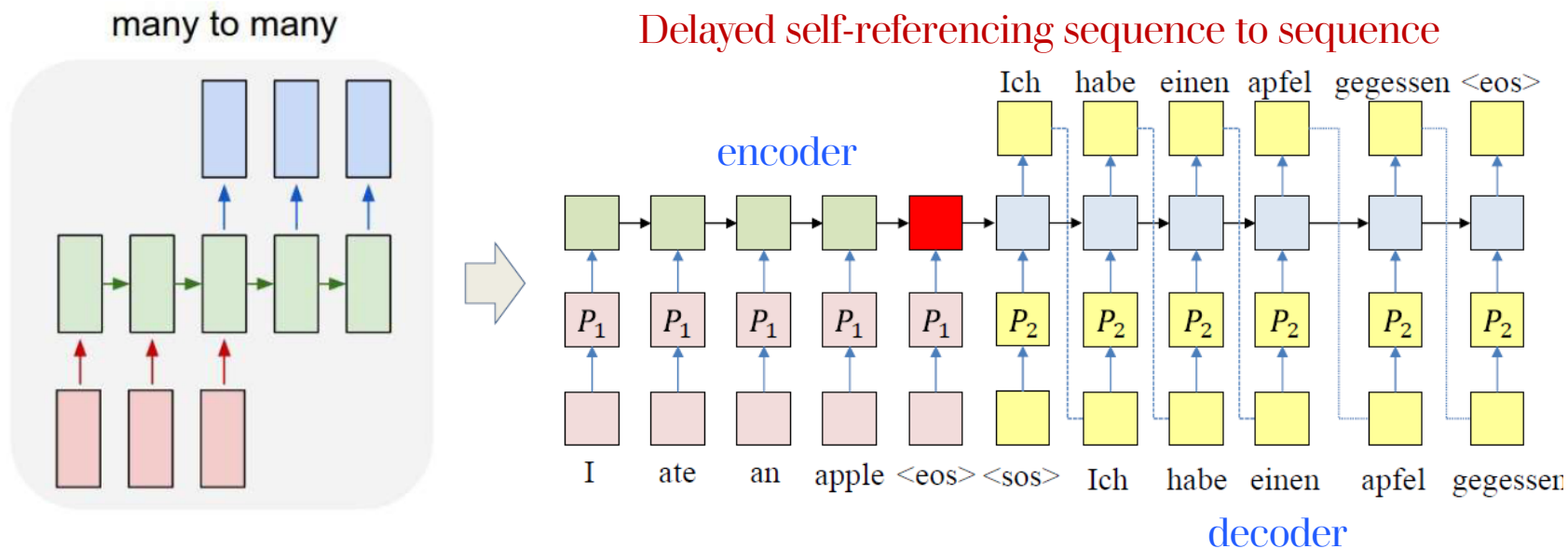I ate an apple → Seq2seq → Ich habe einen apfel gegessen

## problem

- No expected synchrony between input and output
- So, we can't solve the problem well by using only one RNN

## model

many to many

Delayed self-referencing sequence to sequence

encoder

Ich   habe   einen   apfel   gegessen <eos>

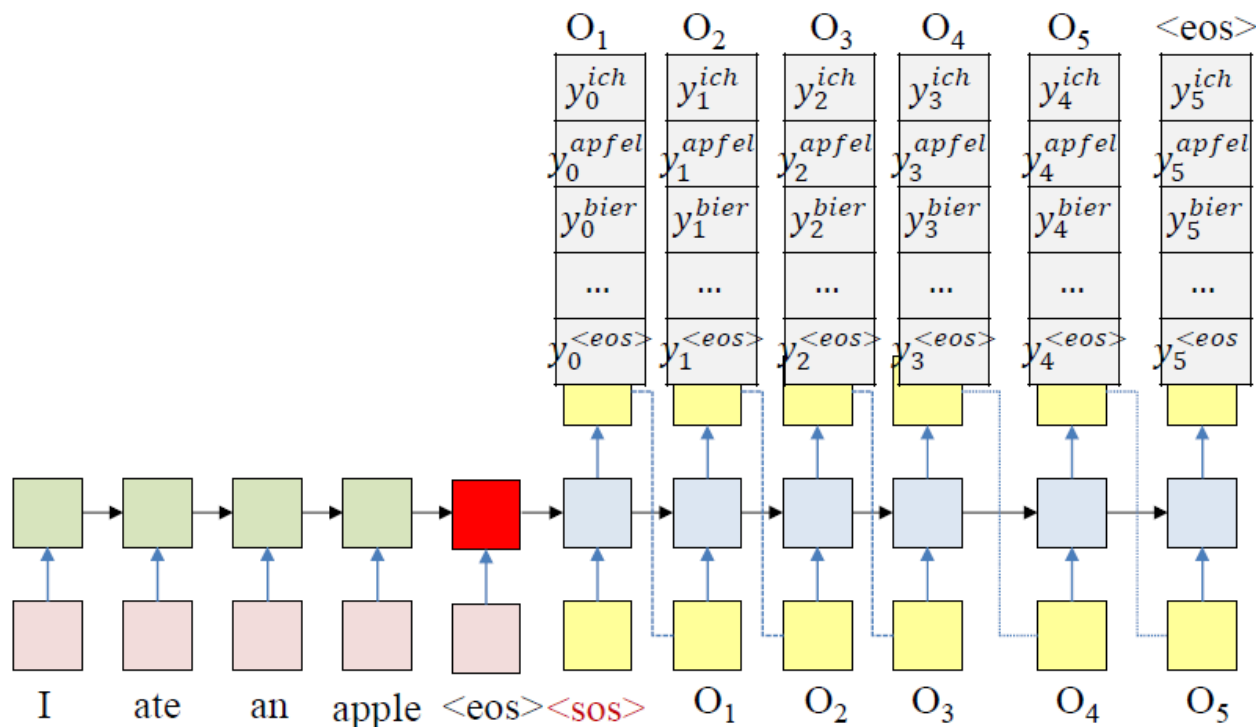I     ate     an    apple <eos> <sos>   Ich   habe   einen   apfel   gegesser

decoder

# Topics

❖ Seq2Seq for Machine Translation

❖ **Decoding & Training process**

❖ Attention model

# Decoding & Training process

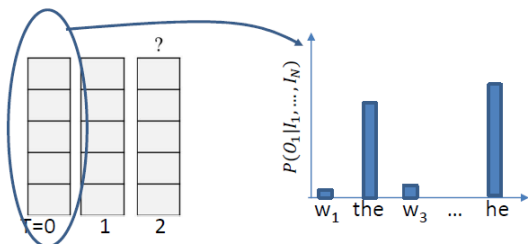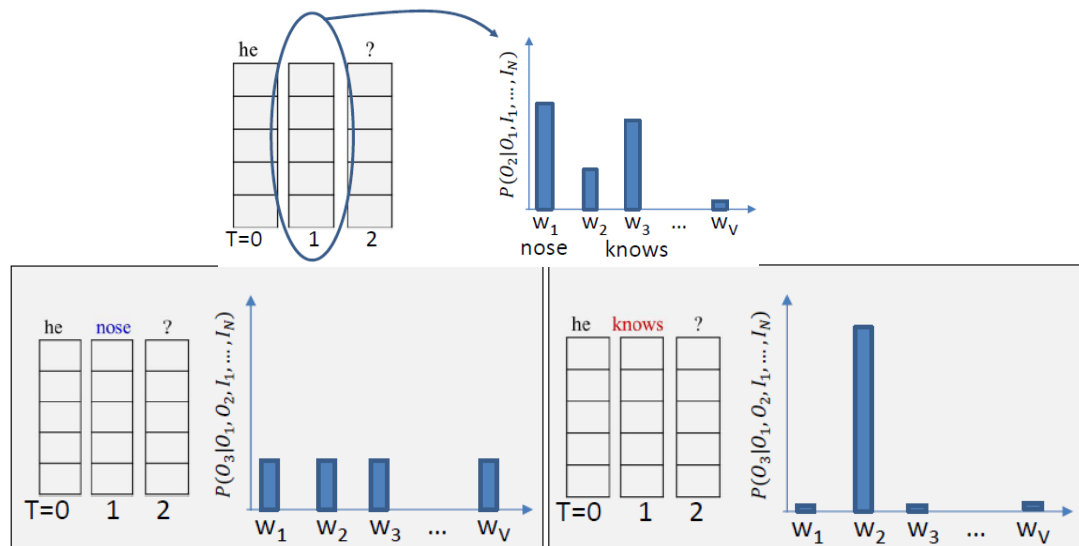Generating outputs



- Goal : produce the most likely output

$$\underset{O_1,\dots,O_L}{\operatorname{argmax}} P\big(O_1, \dots, O_L \mid W_1^{in}, \dots, W_N^{in}\big)$$

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

- Greedy drawing
- Random sampling
- Beam search

# Decoding & Training process

## Problems of greedy drawing



- Impossible to know a priori which word leads to the more promising future

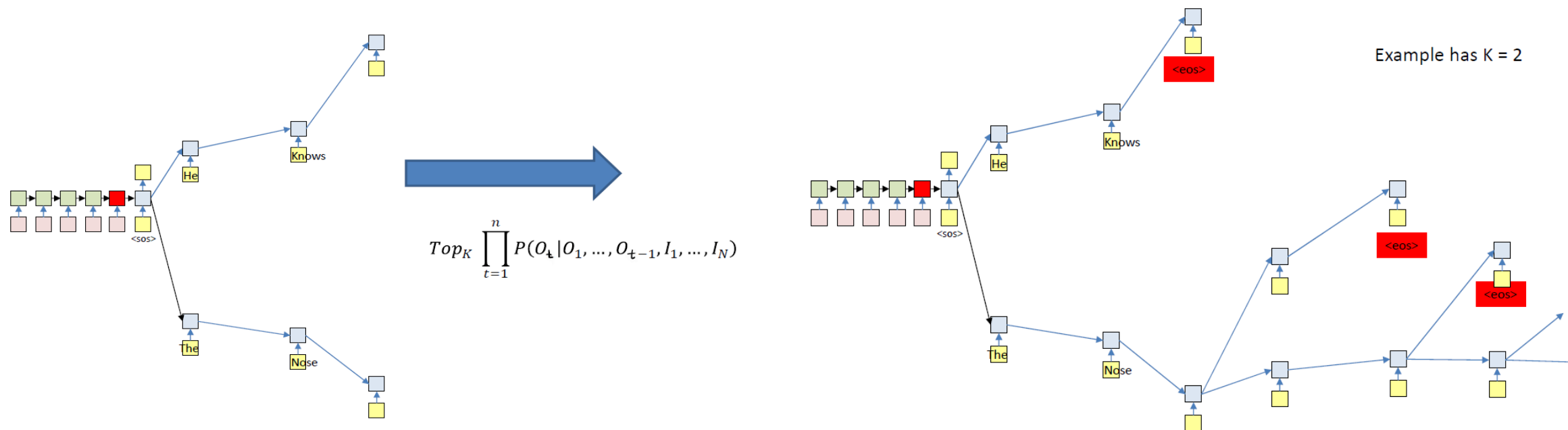## Drawing by random sampling

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

- Sometimes give more likely output than greedy

- But, not guaranteed to give the most likely output

- Making a poor choice at any time commits us to poor future

- But we cannot know at that time the choice was poor

Solution : don't choose
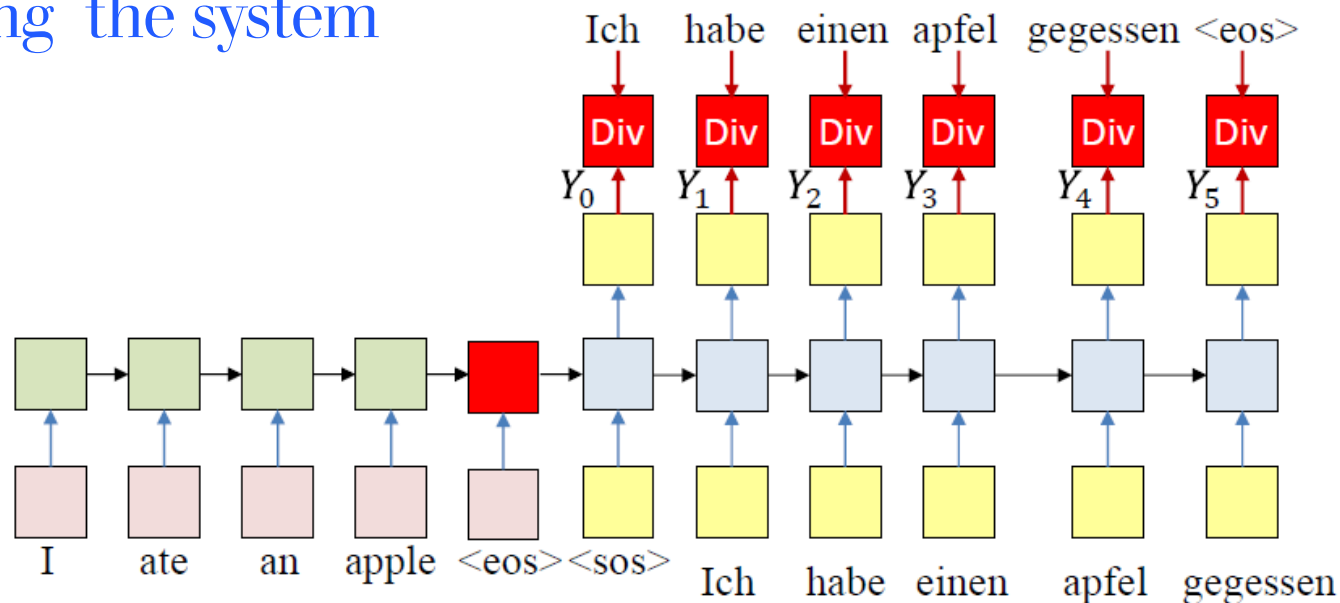
# Decoding & Training process

Multiple choices & pruning : **beam search**



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

Example has K = 2

- At each time, retain only the top K scoring forks
- Terminate when the current most likely path overall ends in <eos>
  - select the most likely sequence ending in <eos> across all terminating seqences

# Decoding & Training process

✓ Can reversing the input seq

✓ Can randomly choose some output words for backprop gradients

## forward

- Input source seq to encoder & target seq (ground truth) to decoder (Use teacher forcing) => easier training, calculate divergence

- Compute the divergence between the output distribution and target word sequence
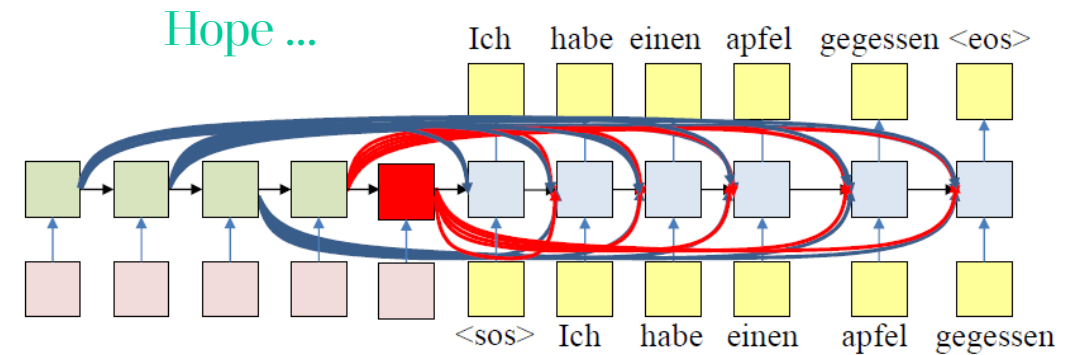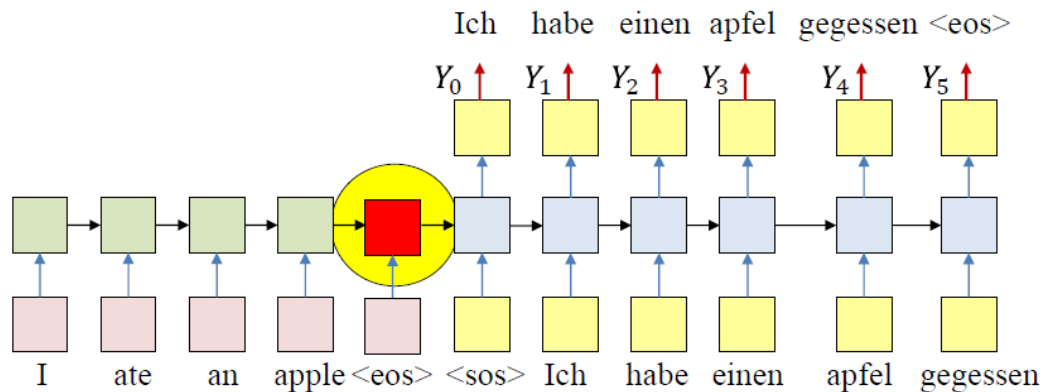
## backward

- Backprop DIV through whole end2end network

11

# Topics

❖ Seq2Seq for Machine Translation

❖ Decoding & Training process

❖ **Attention model**

# Attention model

Problem with naïve enc-dec net
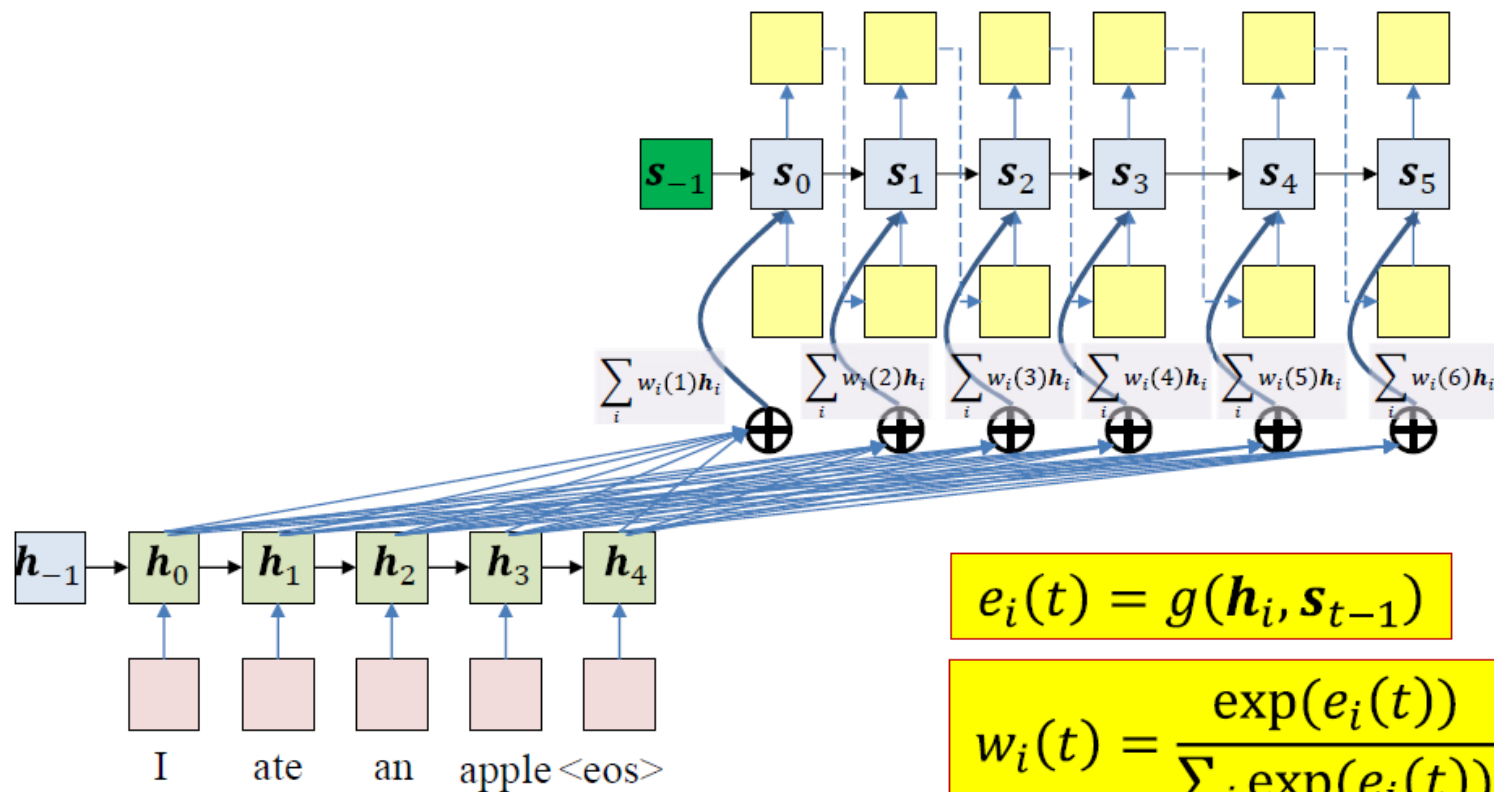


- Information bottleneck
  (all the information about the input seq is embedded into a single vector)

- In reality : all hidden values carry information
  -> some of which may be diluted downstream

Feasible solution :

**Attention mechanism**

# Attention model



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

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

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

$$g(h_i, s_{t-1}) = h_i^T W_g s_{t-1}$$
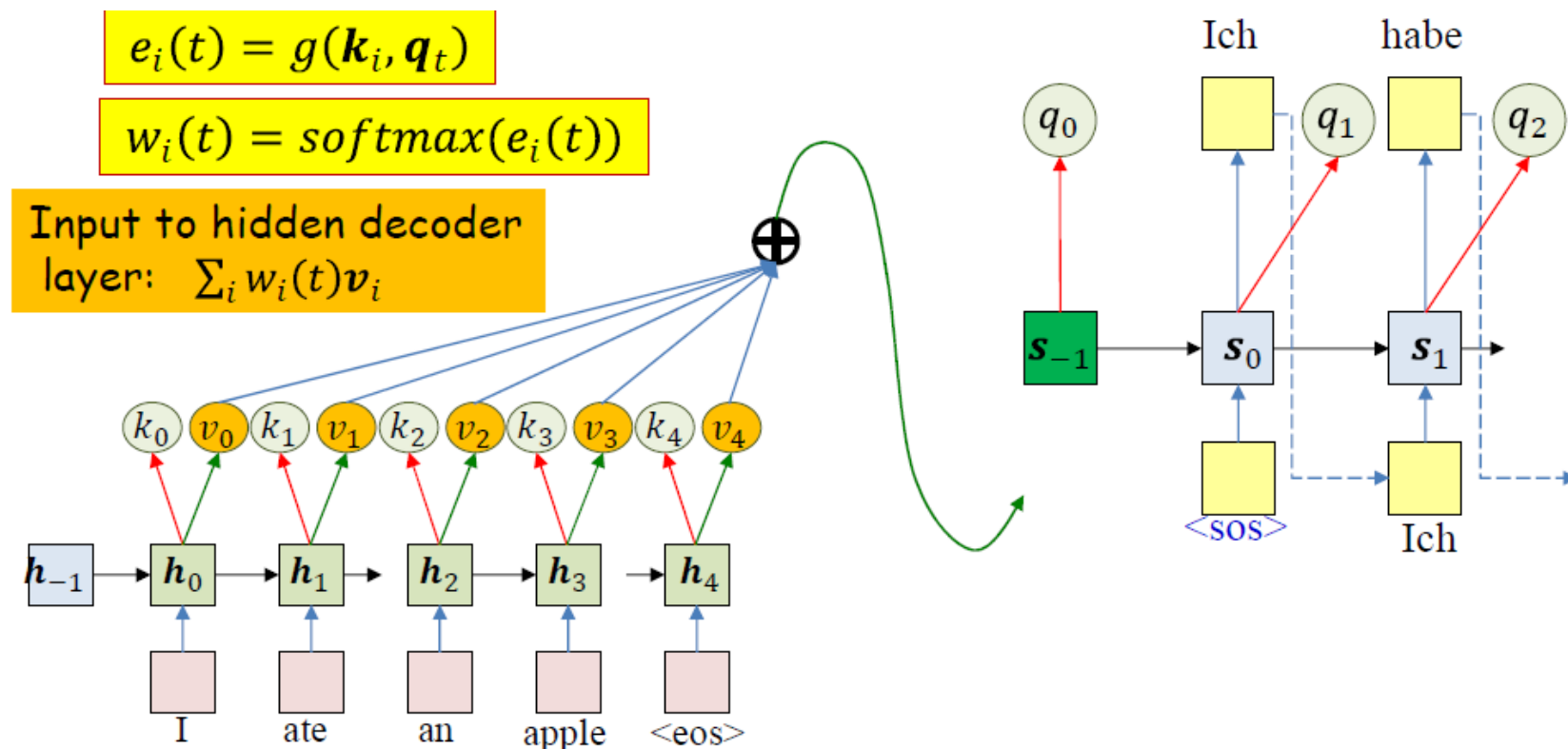
$$g(h_i, s_{t-1}) = v_g^T tanh\left(W_g \begin{bmatrix} h_i \\ s_{t-1} \end{bmatrix}\right)$$

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

- e : attention scores, w : attention weights
- Weights vary by output time
  (time-varying weight that specifies relationship of output time to input time)
- The weights are a distribution over the input

14

# Attention model

<Query - Key - Value> : Generalize the Attention

$$e_i(t) = g(\boldsymbol{k}_i, \boldsymbol{q}_t)$$

$$w_i(t) = softmax(e_i(t))$$

Input to hidden decoder
layer: $\sum_i w_i(t)\boldsymbol{v}_i$

Ich          habe

$q_0$   $q_1$   $q_2$

$\boldsymbol{s}_{-1}$   $\boldsymbol{s}_0$   $\boldsymbol{s}_1$

<sos>

Ich

$k_0$ $v_0$ $k_1$ $v_1$ $k_2$ $v_2$ $k_3$ $v_3$ $k_4$ $v_4$

$\boldsymbol{h}_{-1}$ $\boldsymbol{h}_0$ $\boldsymbol{h}_1$ $\boldsymbol{h}_2$ $\boldsymbol{h}_3$ $\boldsymbol{h}_4$
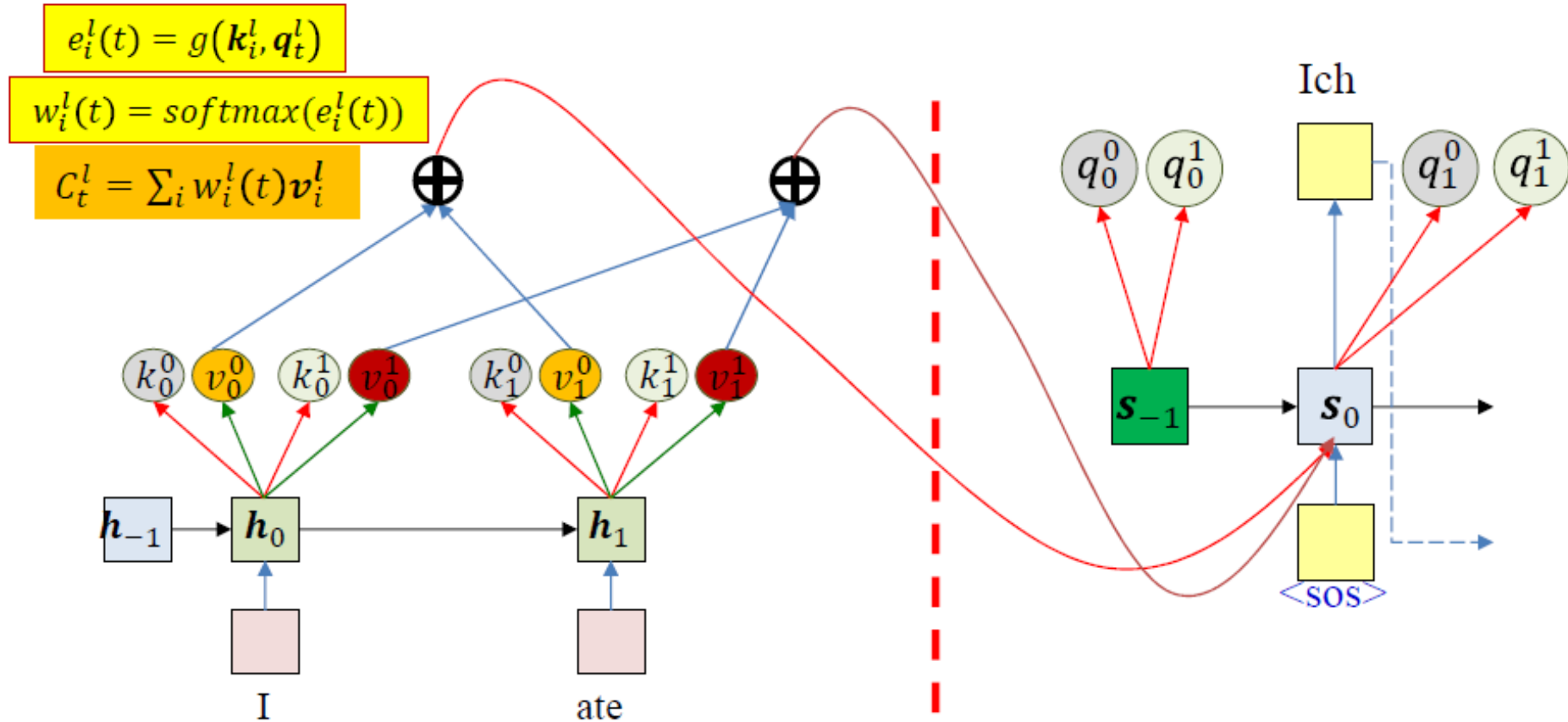
I     ate     an     apple     <eos>

- The weight is a function of key and query
- The actual context is a weighted sum of value

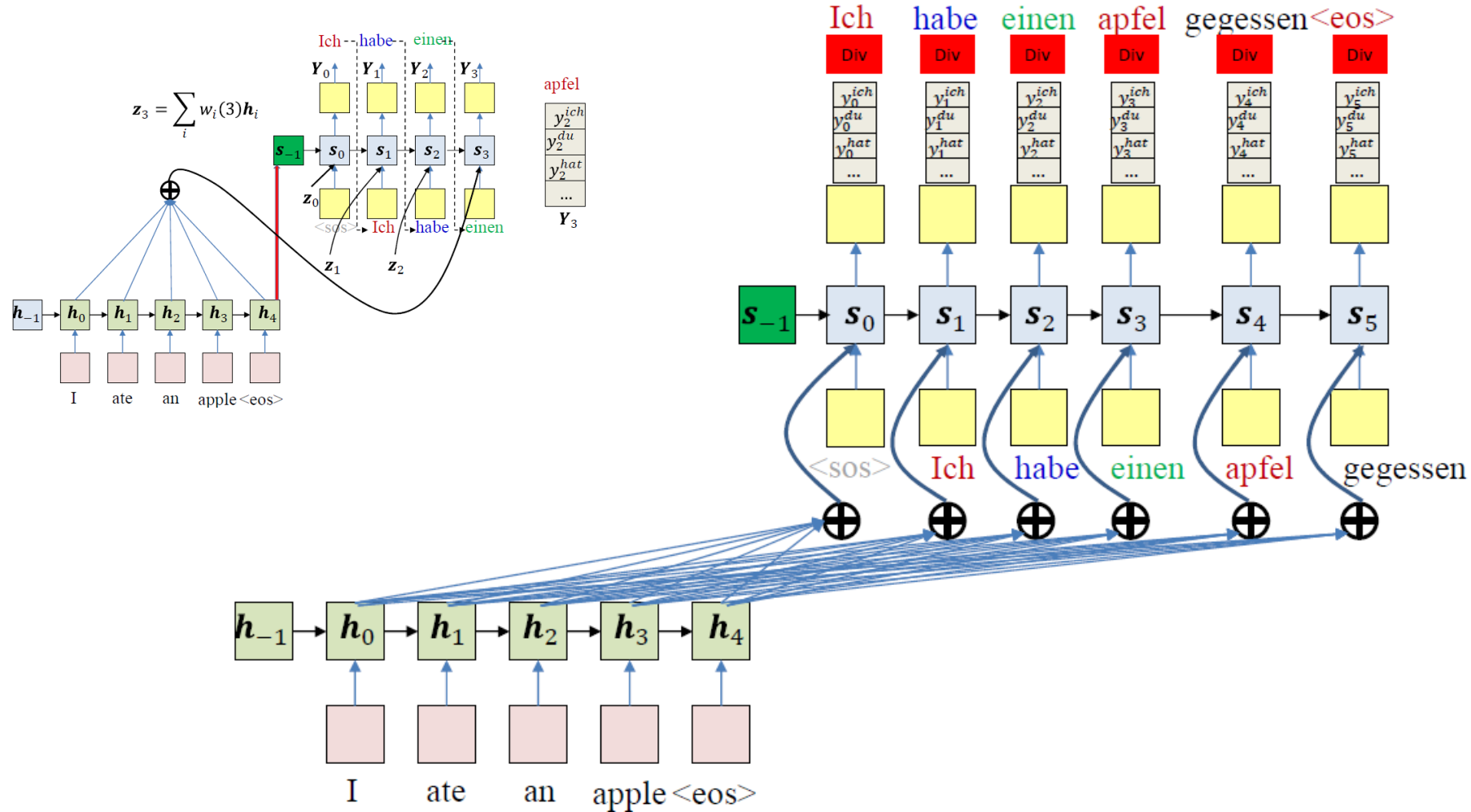Special case:  $k_i = v_i = h_i$
$q_t = s_{t-1}$

# Attention model

Multi-head attention



$$e_i^l(t) = g(\boldsymbol{k}_i^l, \boldsymbol{q}_t^l)$$

$$w_i^l(t) = softmax(e_i^l(t))$$

$$c_t^l = \sum_i w_i^l(t) \boldsymbol{v}_i^l$$

- Can have multiple Q/K/V sets (each attention head uses one of these sets)
- Each attender focuses on a different aspect of the input
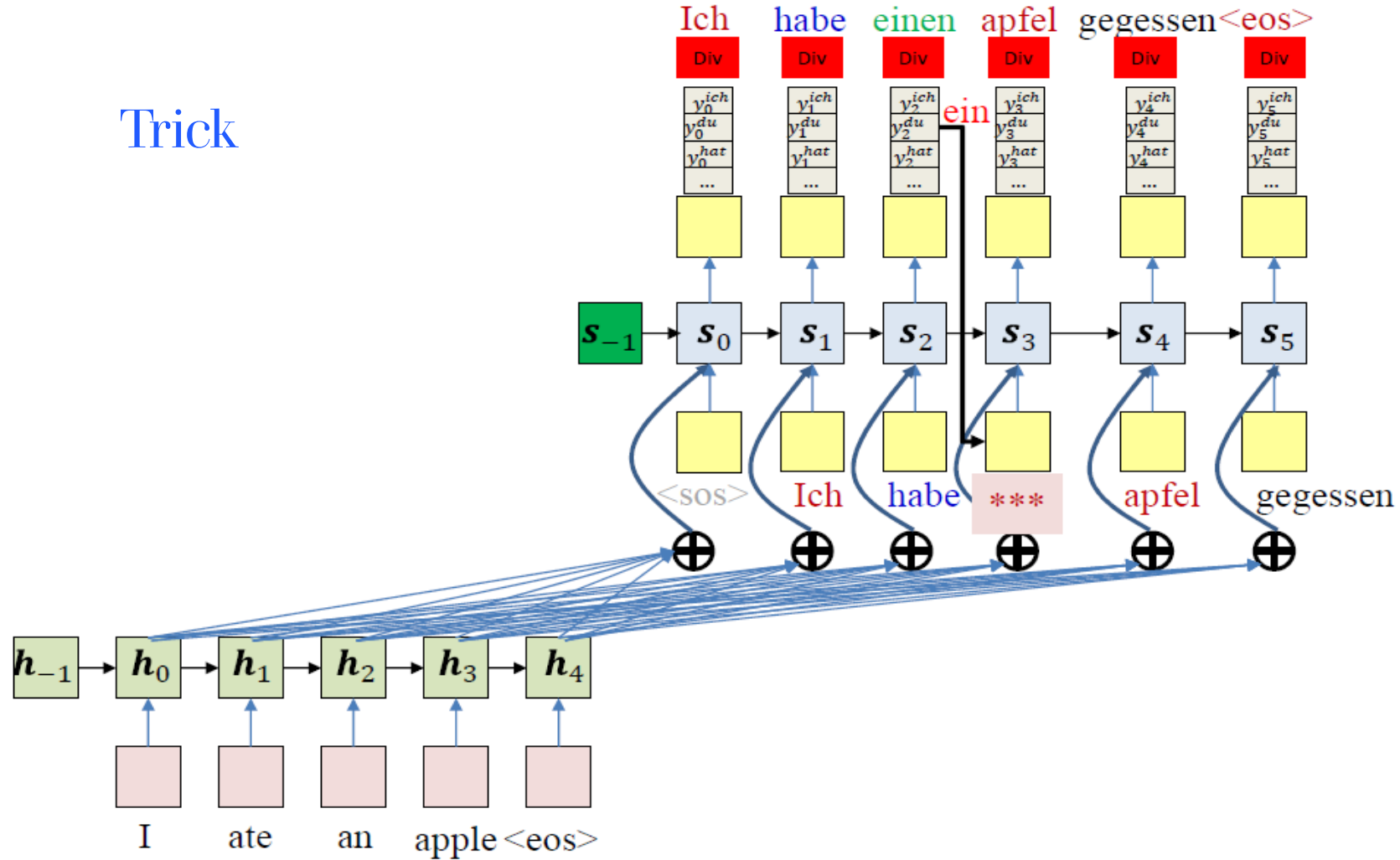
# Attention model

## Inference & train



- If attention function is parametric ,
  back propagation also updates parameters of the attention function

# Attention model

- Pass drawn output instead of ground truth, as input