



05.11.19 18:11

# Neural Machine Translation

## CMPT 413/825, Fall 2019



Jetic Gū  
School of Computing Science

# Overview

- Focus: Neural Machine Translation
- Architecture: Encoder-Decoder Neural Network
- Main Story:
  - Introduction
  - Encoder-Decoder Architecture: Sequence-to-Sequence
  - Attention Mechanisms
  - \*Copy Mechanism
  - \*BeamSearch
  - \*[Extra] Beyond Seq2Seq: Attention is all you need
  - \*[Extra] Beyond NMT

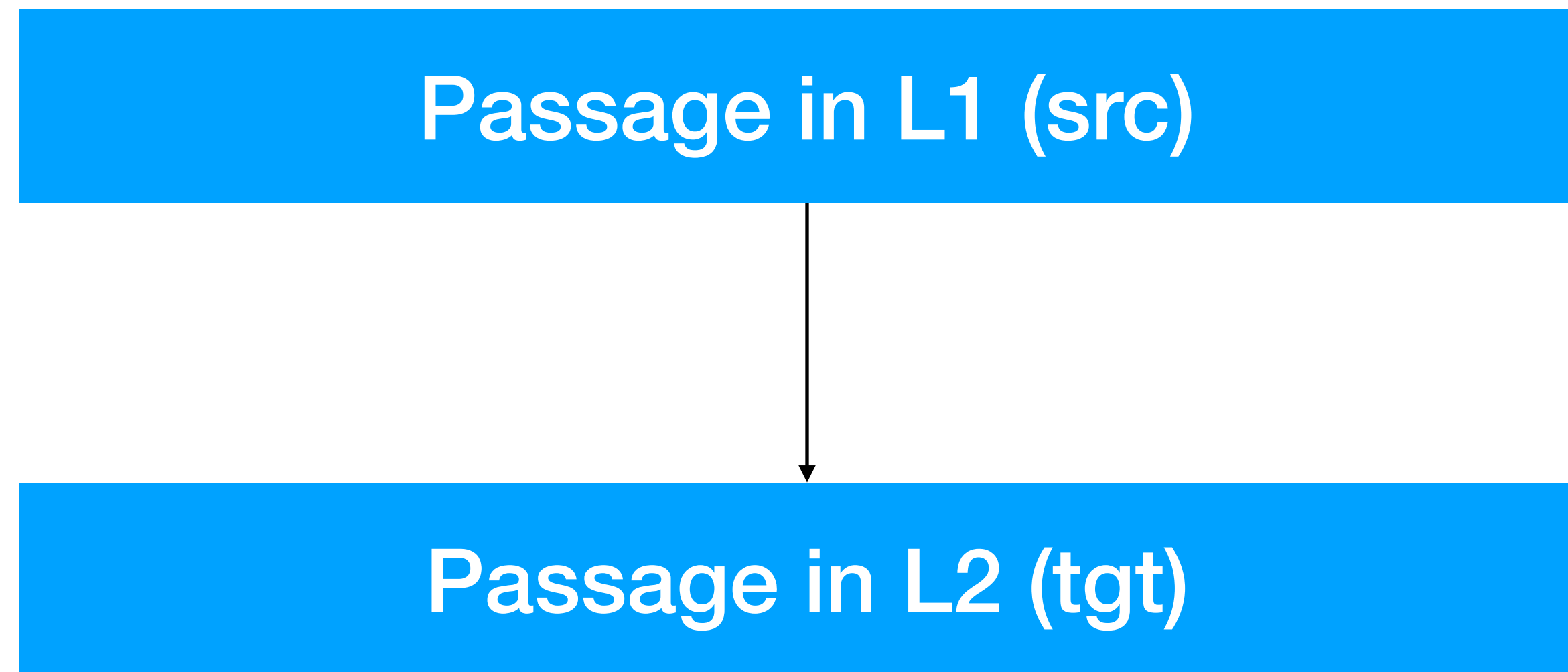
# Machine Translation

Was hat Prof. Sarkar gegessen?



What did Prof. Sarkar eat?

# Machine Translation



# Machine Translation

$$F = (f_1, f_2, \dots, f_{|F|})$$



$$E = (e_1, e_2, \dots, e_{|E|})$$

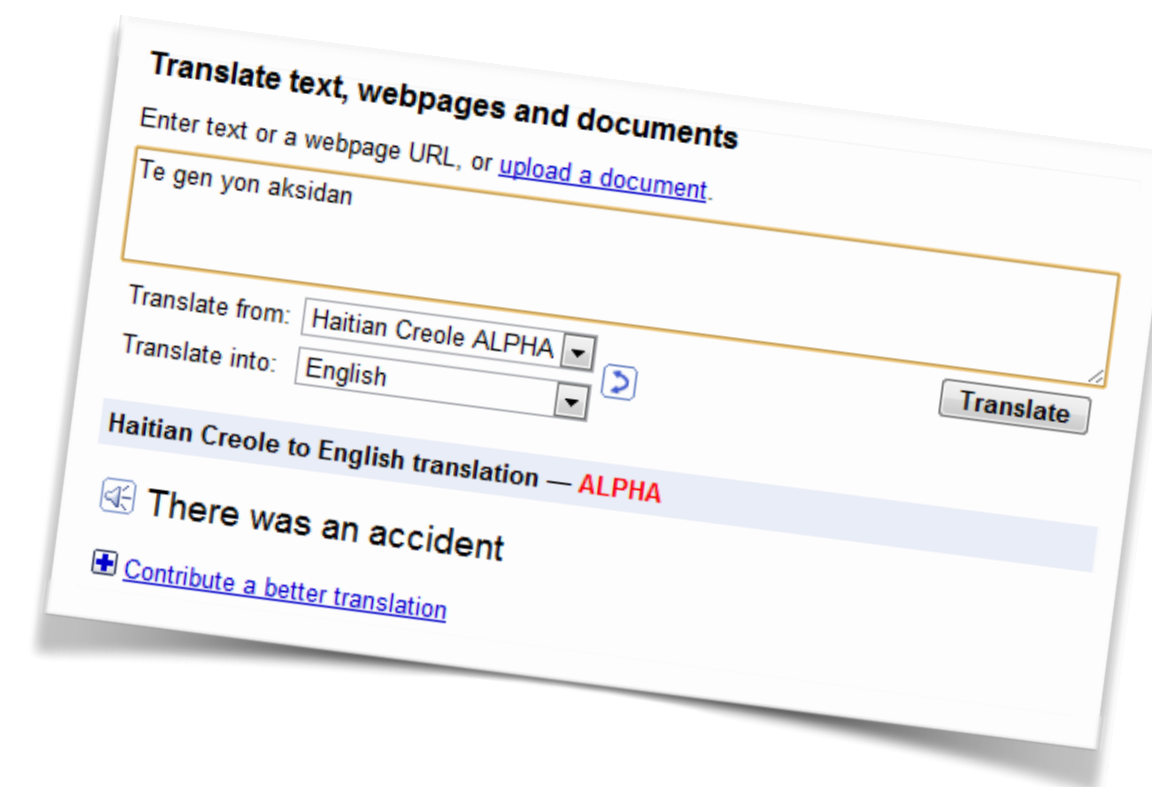
$$Pr(E | F)$$

# History

- Started in the 1950s: rule-based, tightly linked to formal linguistics theories
- 1980s: Statistical MT
- 2000s-2015: Statistical Phrase-Based MT
- 2015-Present: Neural Machine Translation

# History

- Started in the 1950s: rule-based, tightly linked to formal linguistics theories
- 1980s: Statistical MT
- 2000s-2015: **Statistical Phrase-Based MT**
- 2015-Present: Neural Machine Translation



# History

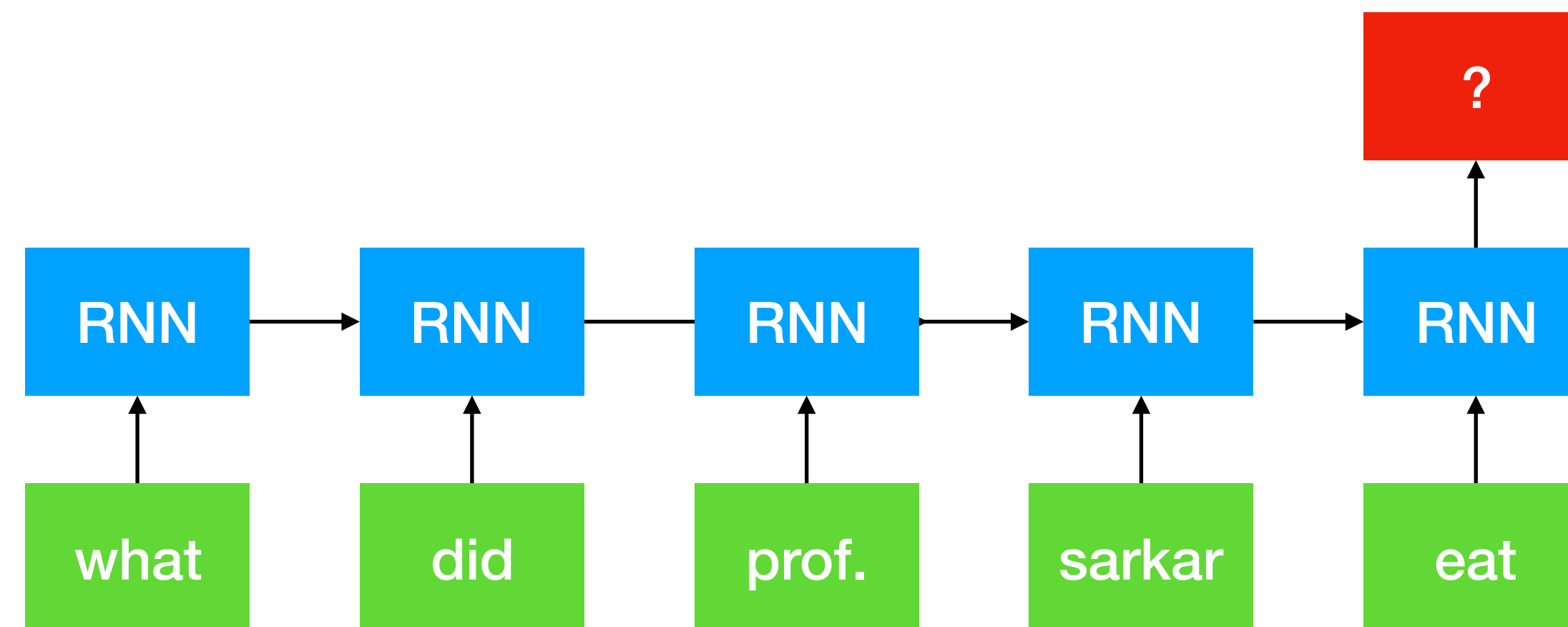
- Started in the 1950s: rule-based, tightly linked to formal linguistics theories
- 1980s: Statistical MT
- 2000s-2015: Statistical Phrase-Based MT
- 2015-Present: **Neural Machine Translation**
- ~2018-Present: **Neural Machine Translation + PBMT Hybrid**



# Recap: Generative Neural LM

$$Pr(e' = e_t | e_{<t})$$

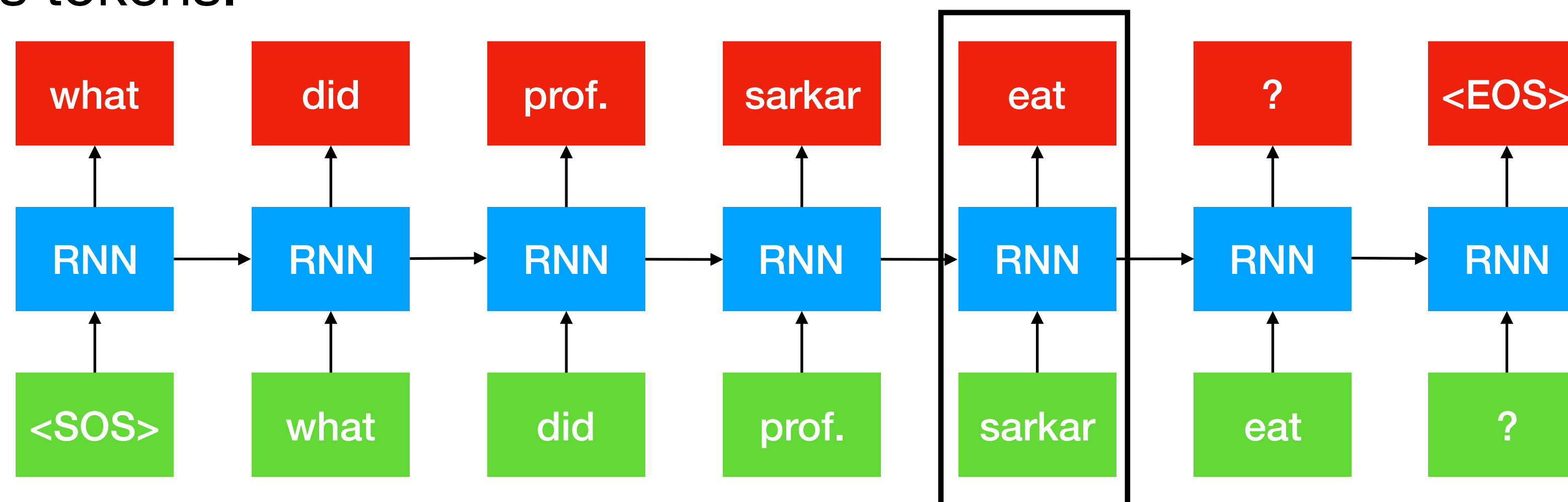
- Next token's probability distribution over the vocabulary, conditioned on previous tokens.



# Recap: Generative Neural LM

$$Pr(E) = \prod_t Pr(e' = e_t | e_{<t})$$

- Next token's probability distribution over the vocabulary, conditioned on previous tokens.



$$Pr(e' = e_5 | y_{1:4}) = Pr(e' = \text{'eat'} | \text{'what did prof. sarkar'})$$

# Conditional Generative Neural LM

$$F = (f_1, f_2, \dots, f_{|F|})$$



$$E = (e_1, e_2, \dots, e_{|E|})$$

$$\underline{Pr(E | F) = \prod_t Pr(e'_t = e_t | F, e_{<t})}$$

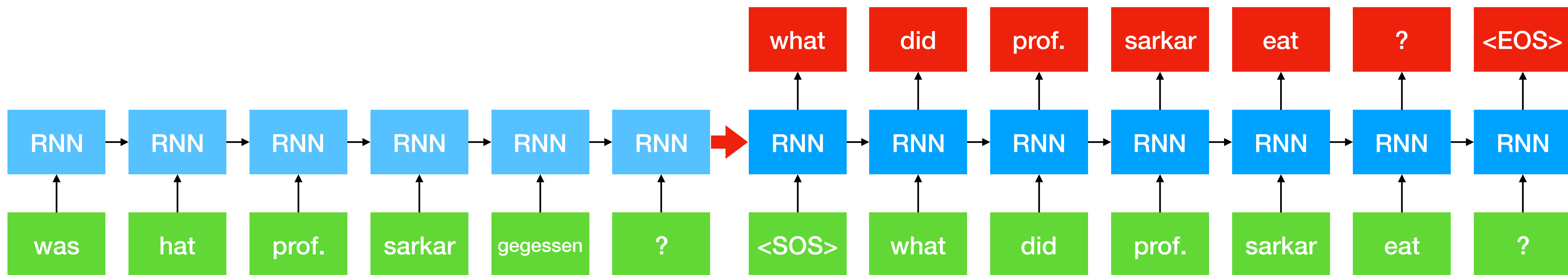
CLM

# Conditional Generative Neural LM

$$\frac{Pr(E) = \prod_t Pr(e' = e_t | e_{<t})}{\text{LM}}$$

$$\frac{Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})}{\text{CLM}}$$

# Conditional Generative Neural LM



$$\underline{Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})}$$

CLM

Concept

# Conditional Generative Neural LM

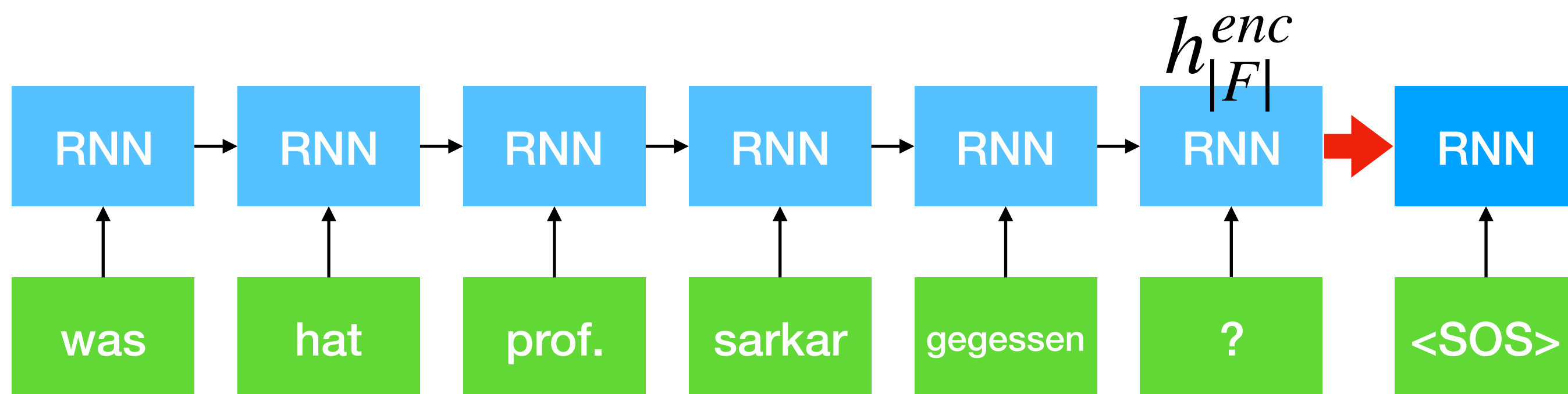
$$h_0^{dec} = h_{|F|}^{enc}$$

$$W \in \mathbb{R} \text{ [?]}$$

$$h_1^{dec} = RNN^{dec}(\text{'<SOS>'}, h_0^{dec})$$

$$b \in \mathbb{R} \text{ [?]}$$

$$e'_1 = \operatorname{argmax}_e (\operatorname{softmax}(Wh_1^{dec} + b))$$



$$\underline{Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})}$$

CLM

# Conditional Generative Neural LM

$$h_0^{dec} = h_{|F|}^{enc}$$

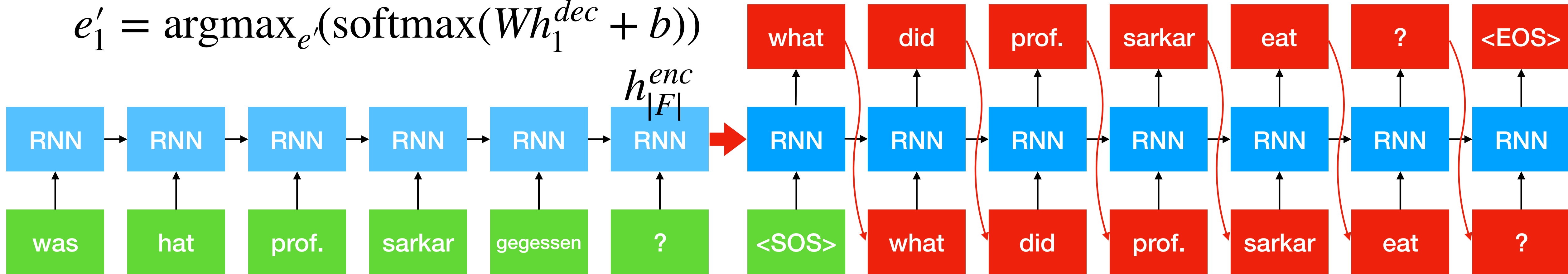
$$W \in \mathbb{R}^{|V_E| \times d}$$

$$b \in \mathbb{R}^{|V_E|}$$

$$h_1^{dec} = RNN^{dec}(\text{'what'}, h_0^{dec})$$

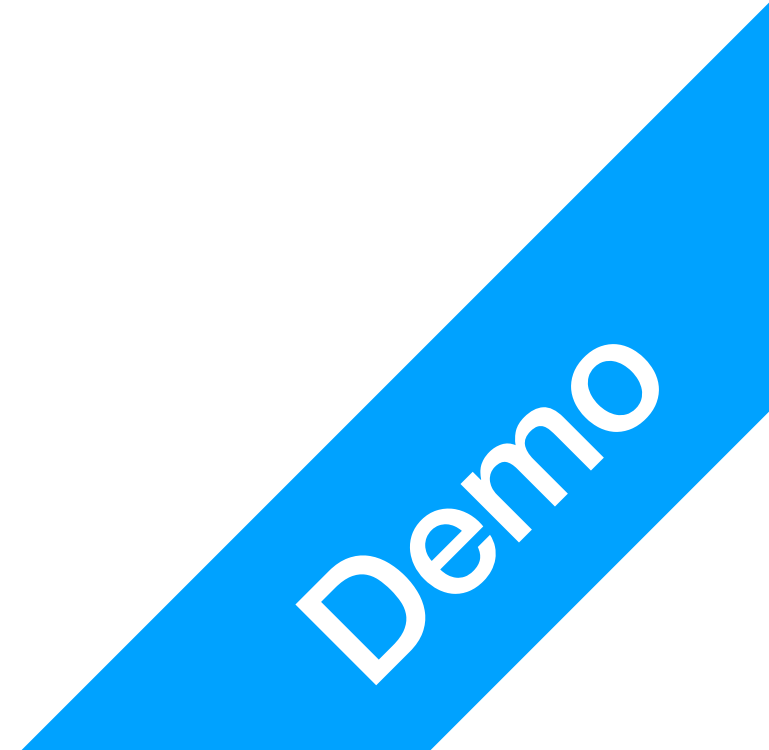
$$h_1^{dec} = RNN^{dec}(\text{'<SOS>'}, h_0^{dec})$$

$$e'_1 = \operatorname{argmax}_e (\operatorname{softmax}(Wh_1^{dec} + b))$$

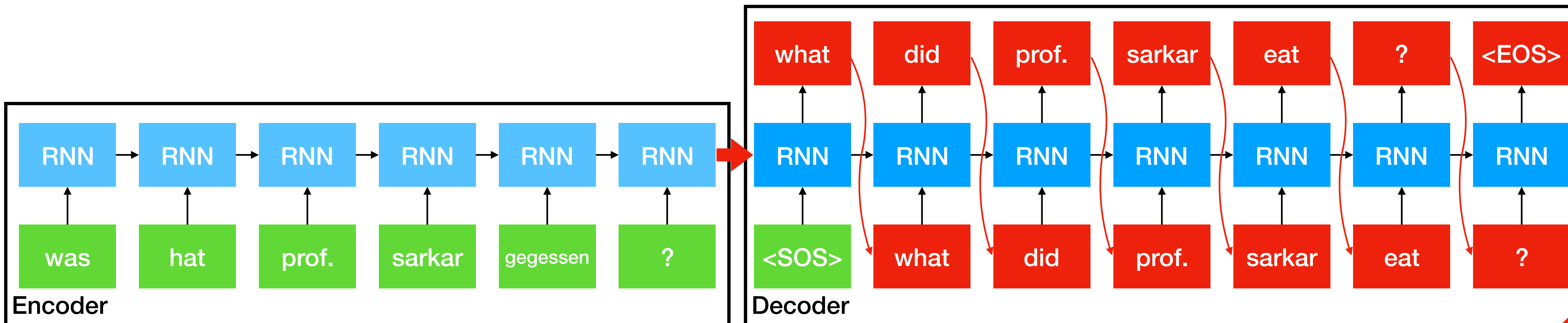


$$\underline{Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})}$$

CLM



# Sequence-to-Sequence



$$\Pr(E | F) = \prod_t \Pr(e' = e_t | F, e_{<t})$$

CLM

Concept

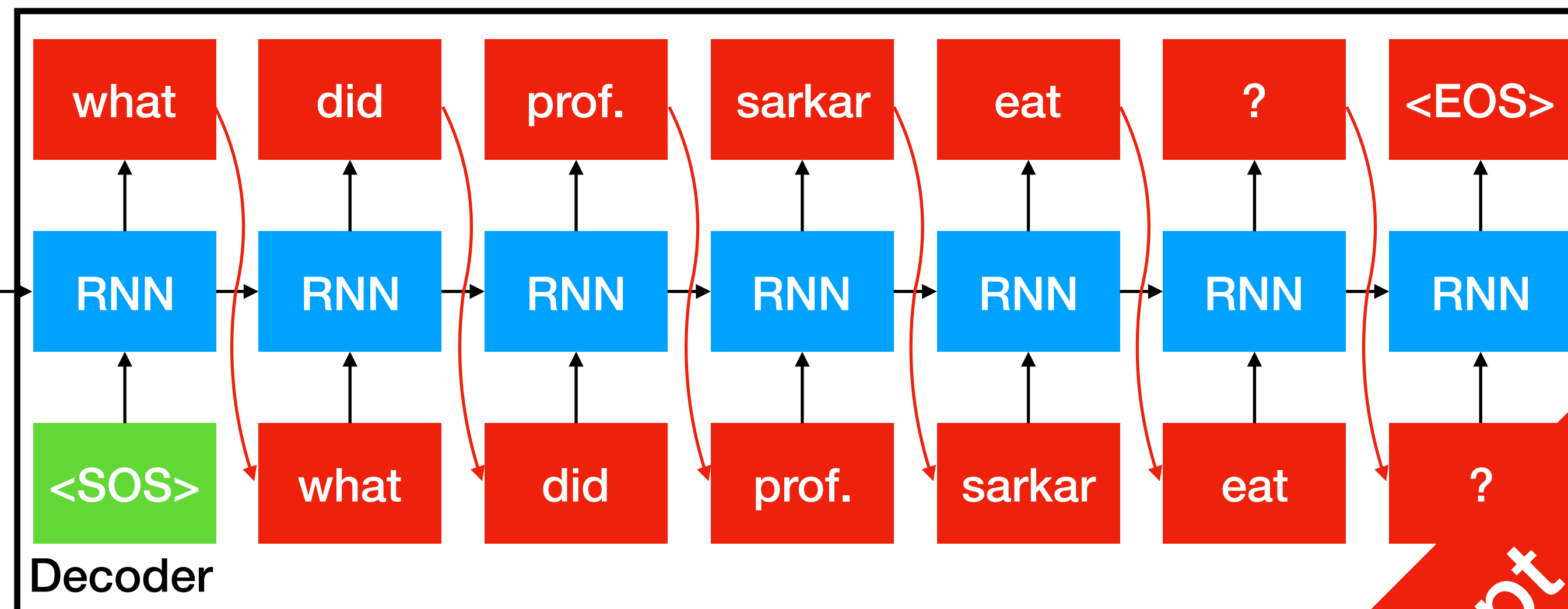
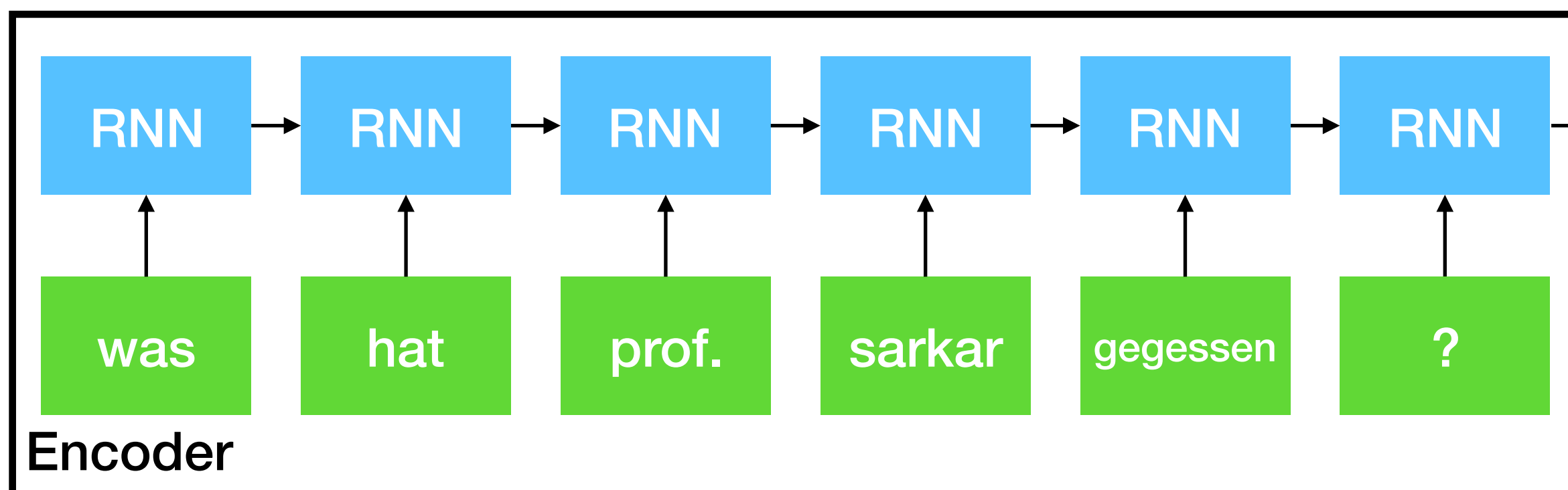




# Sequence-to-Sequence

- S2S: can the two RNNs have the same vocab?
- S2S: can the two RNNs be the same?
- S2S: are there other ways to connect the two RNNs?

$$h_0^{dec} = h_{|F|}^{enc} \in \mathbb{R}^d$$

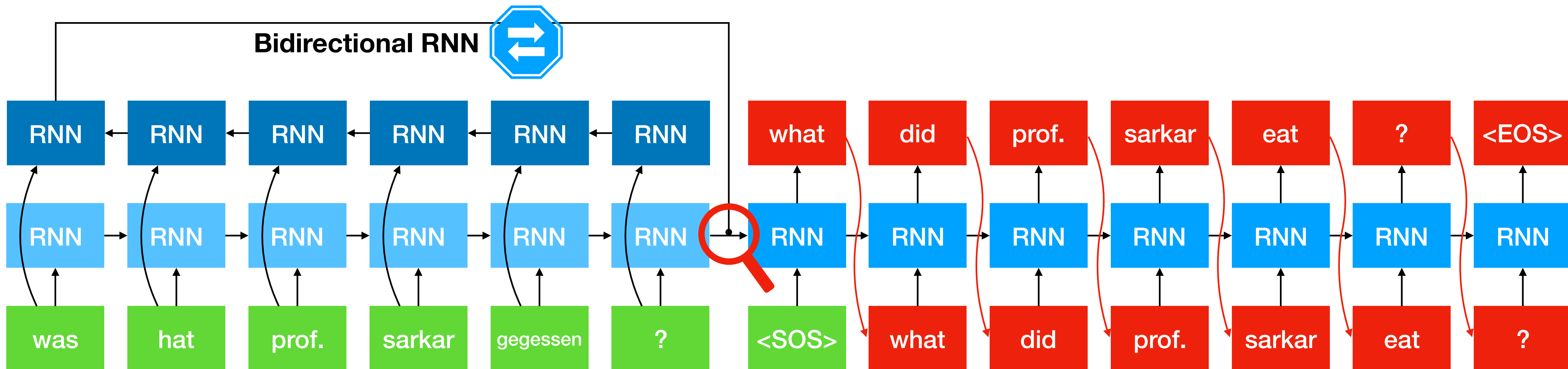


$$Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})$$

CLM

Concept

# Sequence-to-Sequence



e.g.  $h_0^{dec} = [\vec{h}_{|F|}^{enc}; \overleftarrow{h}_{|F|}^{enc}] \in \mathbb{R}^{2d}$

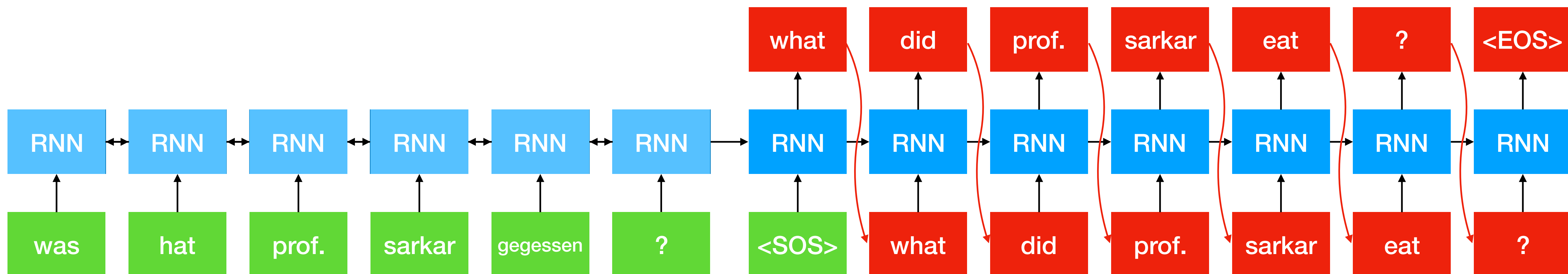
e.g.  $h_0^{dec} = MLP([\vec{h}_{|F|}^{enc}; \overleftarrow{h}_{|F|}^{enc}]) \in \mathbb{R}^d$

$$Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})$$

CLM

Concept

# Sequence-to-Sequence



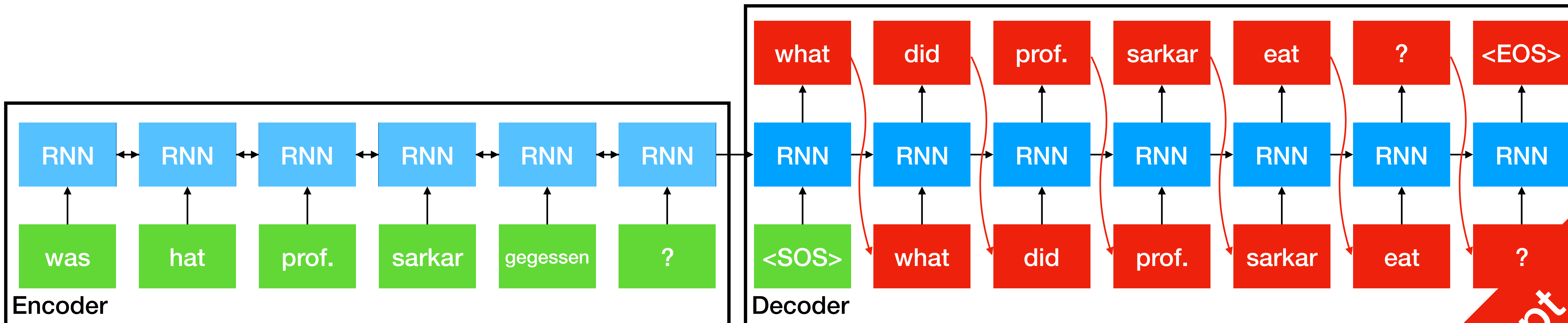
$$Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})$$

**CLM**



# Sequence-to-Sequence

- Difference between a Generative LM and a Conditional Generative LM?
- Difference between Encoder-Decoder and Conditional Generative LM?
- Difference between Encoder-Decoder and Sequence-to-Sequence?



$$Pr(E | F) = \prod_t Pr(e' = e_t | F, e_{<t})$$

**CLM**

Concept

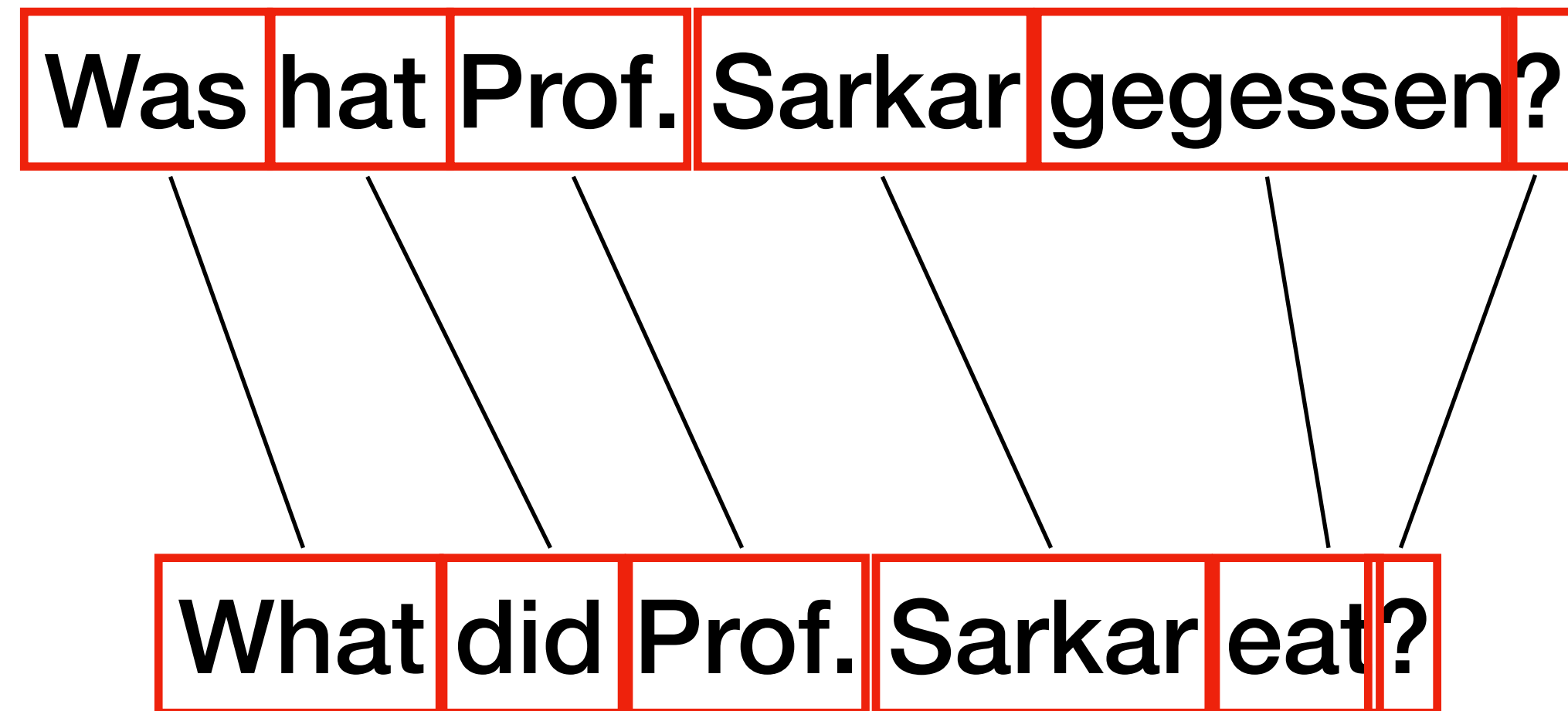
# Summary

- From Generative Neural Language Model (LM) to Conditional Generative Neural Language Model (CLM)
- Encoder-Decoder Architecture
- Vanilla sequence-to-sequence Model
- Common approaches to seq2seq: BiRNN (e.g. BiLSTM, BiGRU) encoder, RNN decoder (LSTM, GRU)

# Is Vanilla Seq2Seq Good Enough?

- Relies on  $h_{|F|}^{enc} \in \mathbb{R}^d$  to represent entire source sentence
  - word-level information, semantics, style, sentiment, etc.
- Problematic when facing longer source sentences

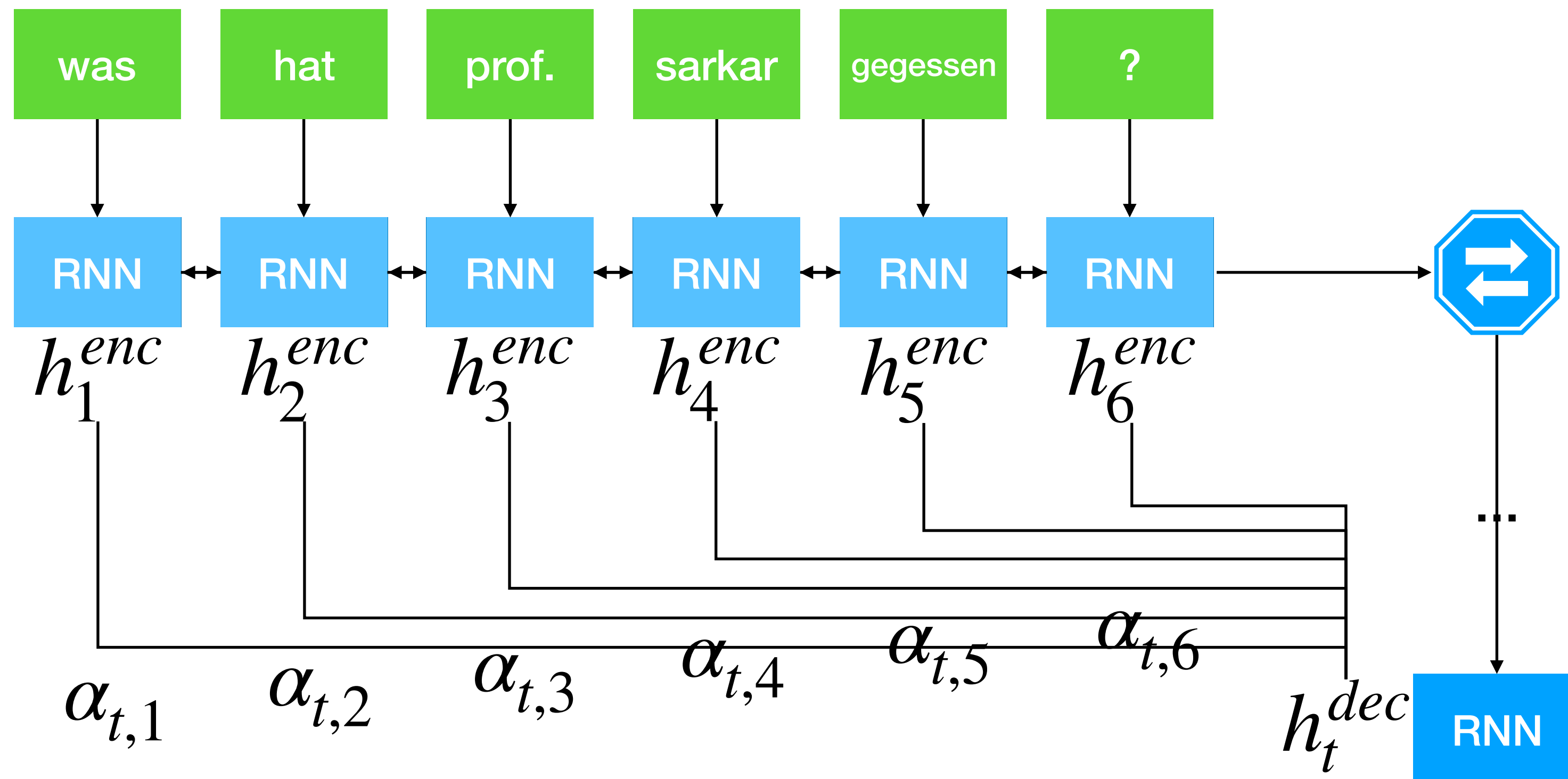
# Is Vanilla Seq2Seq Good Enough?



- Phrase-Based MT: explicit alignment
- Attention: provide decoder-relevant contextual information at each step  $t$



# Attention



$$* f(x, y) = V \tanh(W'x + W''y)$$

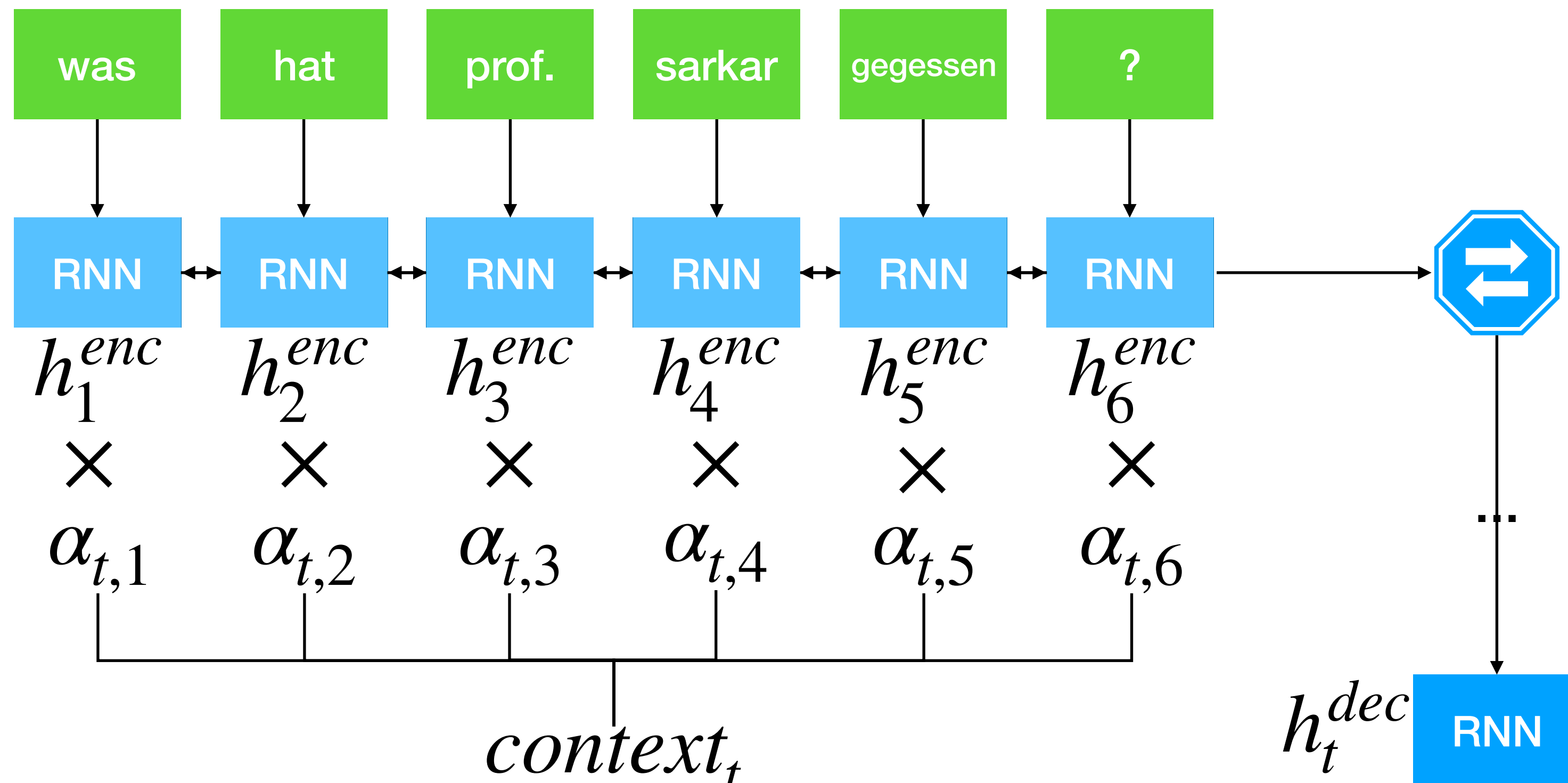
$$score_{t,i} = f(h_i^{enc}, h_t^{dec})$$

$$\alpha_t = \text{softmax}(score_{t,i})$$





# Attention



$$*f(x, y) = V \tanh(W'x + W''y)$$

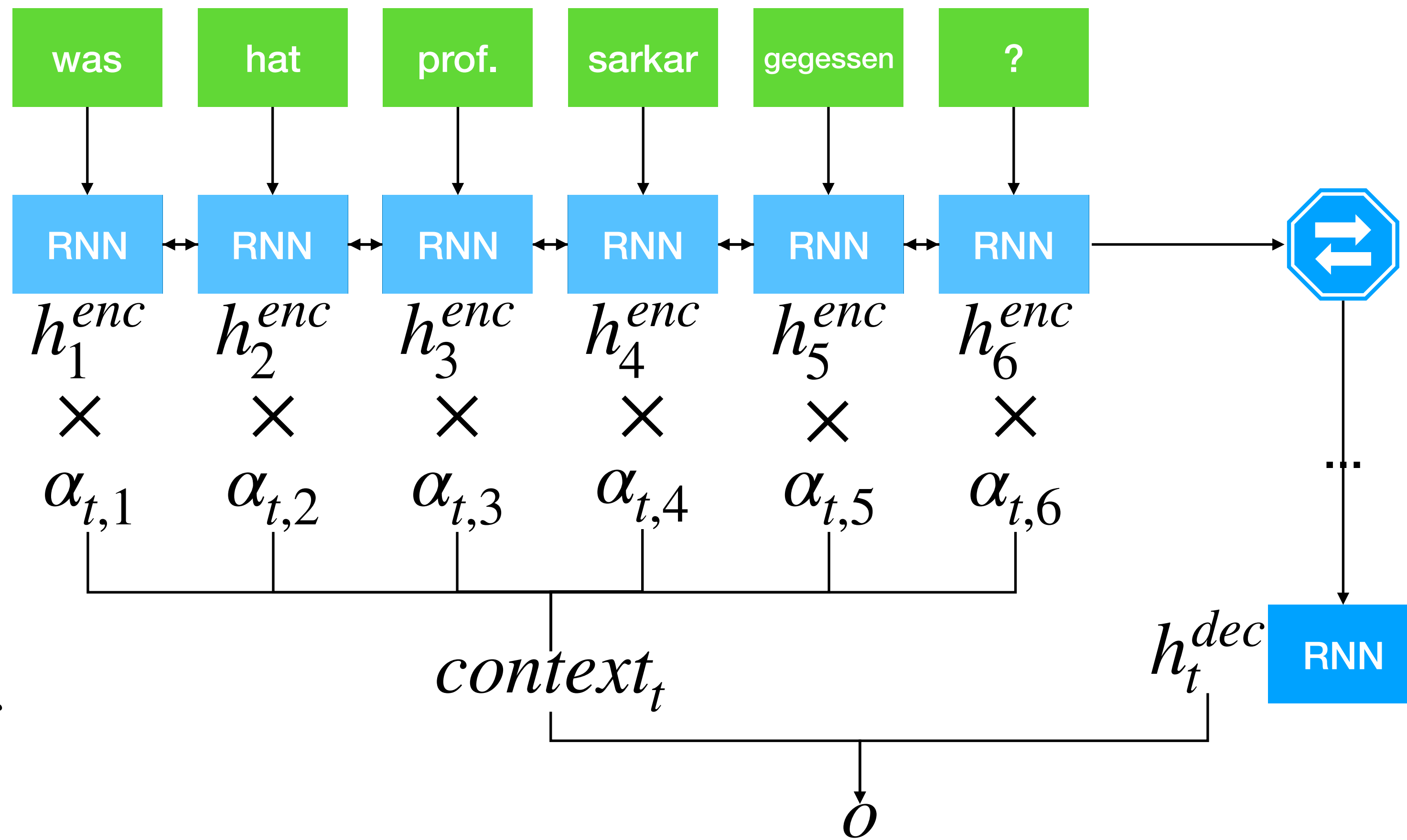
$$score_{t,i} = f(h_i^{enc}, h_t^{dec})$$

$$\alpha_t = \text{softmax}(score_t)$$

$$context_t = \sum_i \alpha_{t,i} h_i^{enc}$$



# Attention



$$*f(x, y) = V \tanh(W'x + W''y)$$

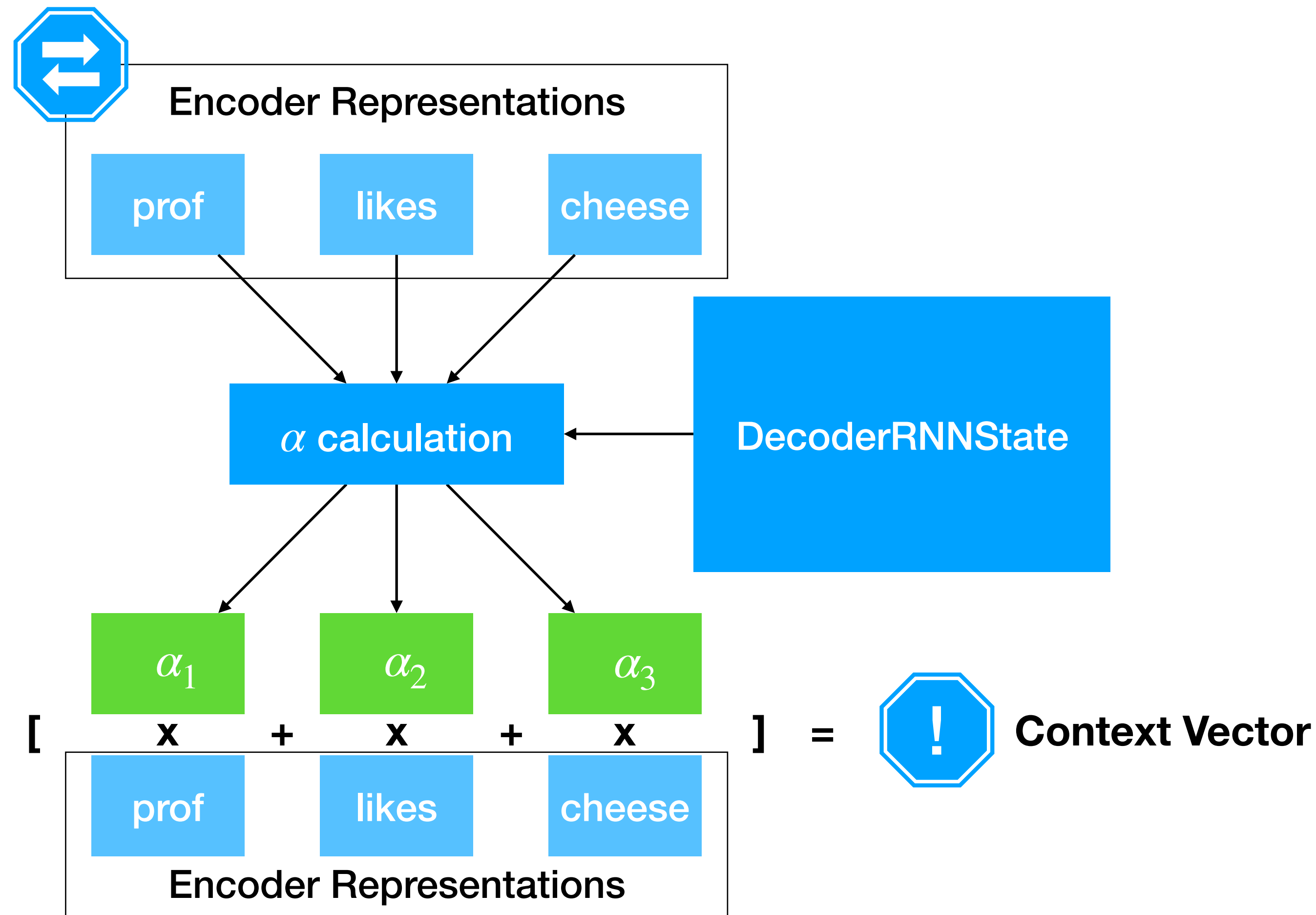
$$score_{t,i} = f(h_i^{enc}, h_t^{dec})$$

$$\alpha_t = \text{softmax}(score_{t,i})$$

$$context_t = \sum_i \alpha_{t,i} h_i^{enc}$$

$$o = \text{softmax}(W[context_t; h_t^{dec}] + b)$$

# Attention



$$* f(x, y) = V \tanh(W'x + W''y)$$

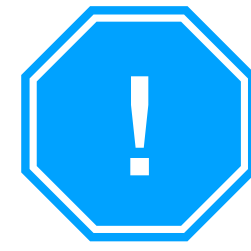
$$score_{t,i} = f(h_i^{enc}, h_t^{dec})$$

$$\alpha_t = \text{softmax}(score_{t,i})$$

$$context_t = \sum_i \alpha_{t,i} h_i^{enc}$$

$$o = \text{softmax}(W[context_t; h_t^{dec}] + b)$$

Concept



# Attention

- Why does Attention work?
- What is in the context vector?
- What is in  $\alpha$ ?

$$\begin{aligned}score_{t,i} &= f(h_i^{enc}, h_t^{dec}) \\ \alpha_t &= \text{softmax}(score_{t,i}) \\ context_t &= \sum_i \alpha_{t,i} h_i^{enc}\end{aligned}$$

Concept

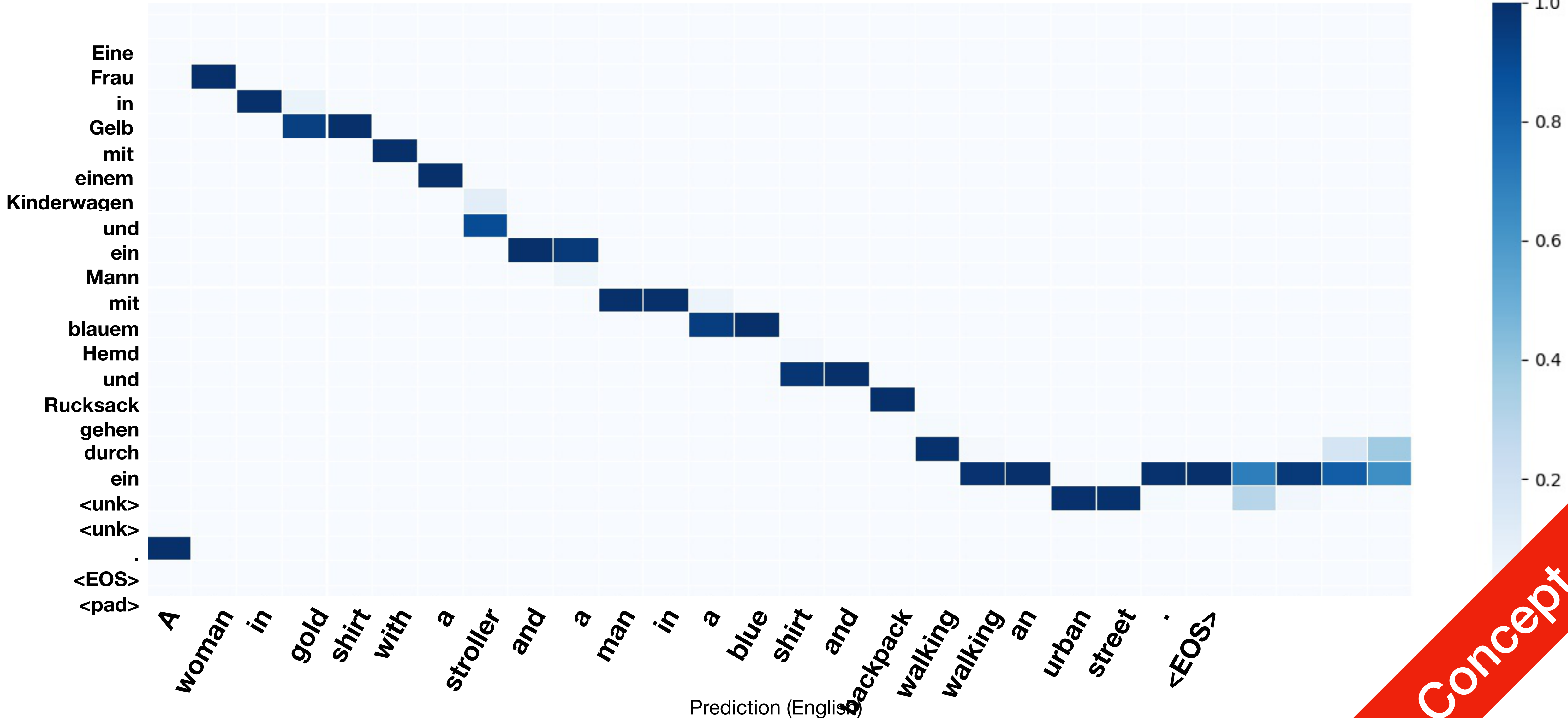


P2  
Attention

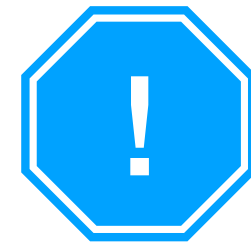


# Attention

Source (German)



Concept



# Attention

- Why does Attention work?
- What is in the context vector?
- What is in  $\alpha$ ?
- ~~Alignment<sup>1</sup>!~~

$$\begin{aligned}score_{t,i} &= f(h_i^{enc}, h_t^{dec}) \\ \alpha_t &= \text{softmax}(score_{t,i}) \\ context_t &= \sum_i \alpha_{t,i} h_i^{enc}\end{aligned}$$

Concept

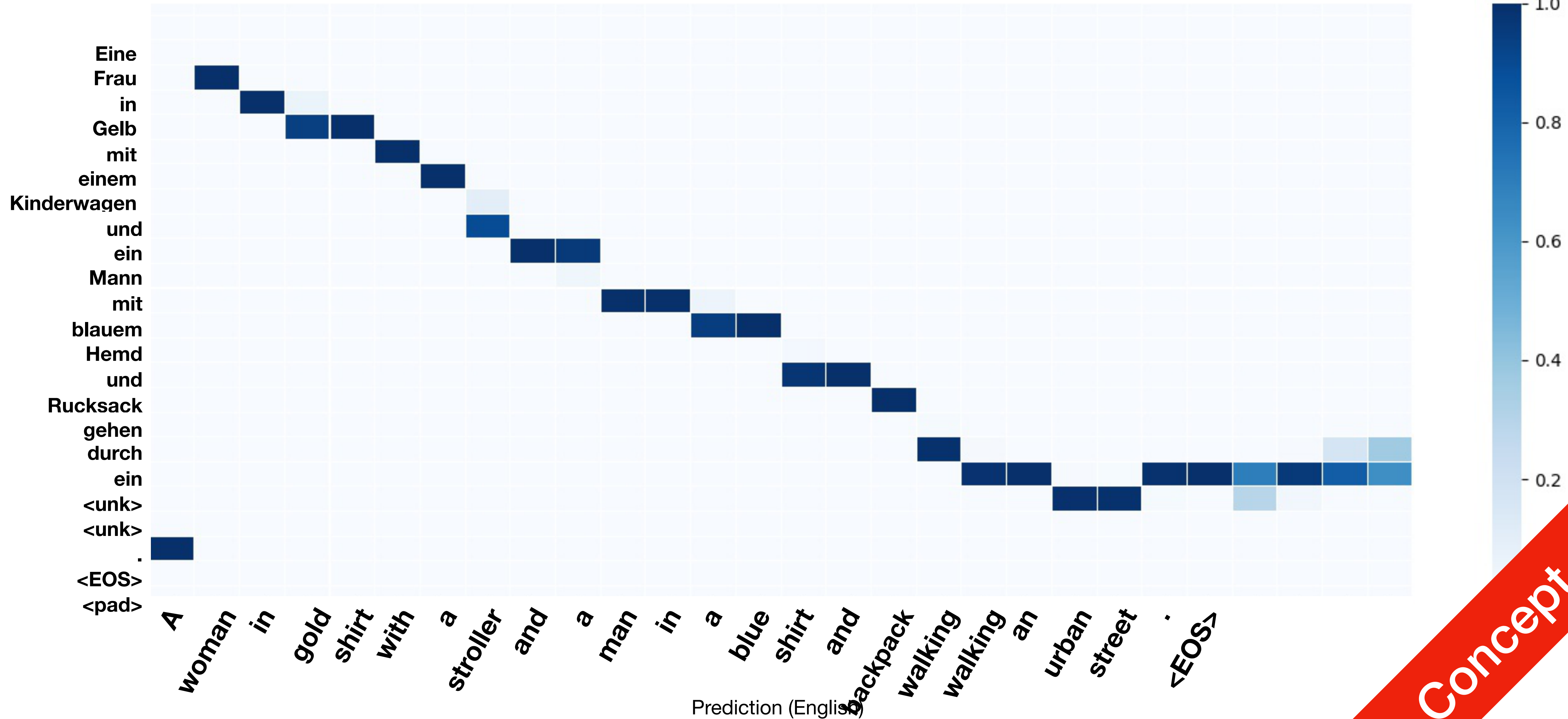


P2  
Attention



# Attention

Source (German)



Concept







# Attention

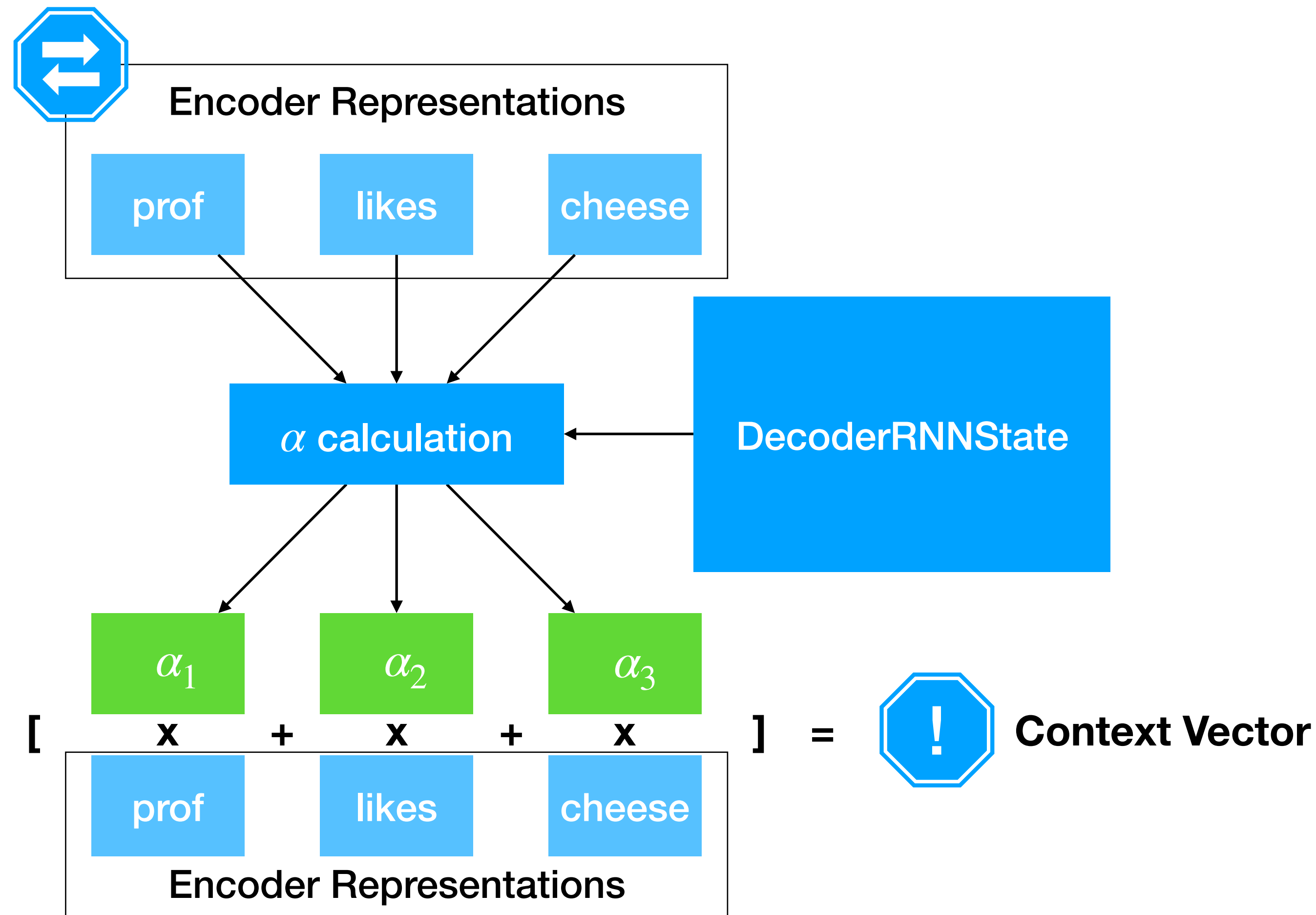
- Why does Attention work?
- What is in the context vector?
- What is in  $\alpha$ ?
  - ~~It's alignment~~<sup>1</sup> !
  - learns to refer to useful information in src
  - similar to human attention: we pay attention to whatever is needed

$$score_{t,i} = f(h_i^{enc}, h_t^{dec})$$

$$\alpha_t = \text{softmax}(score_{t,i})$$

$$context_t = \sum_i \alpha_{t,i} h_i^{enc}$$

# Attention



$$* f(x, y) = V \tanh(W'x + W''y)$$

$$score_{t,i} = f(h_i^{enc}, h_t^{dec})$$

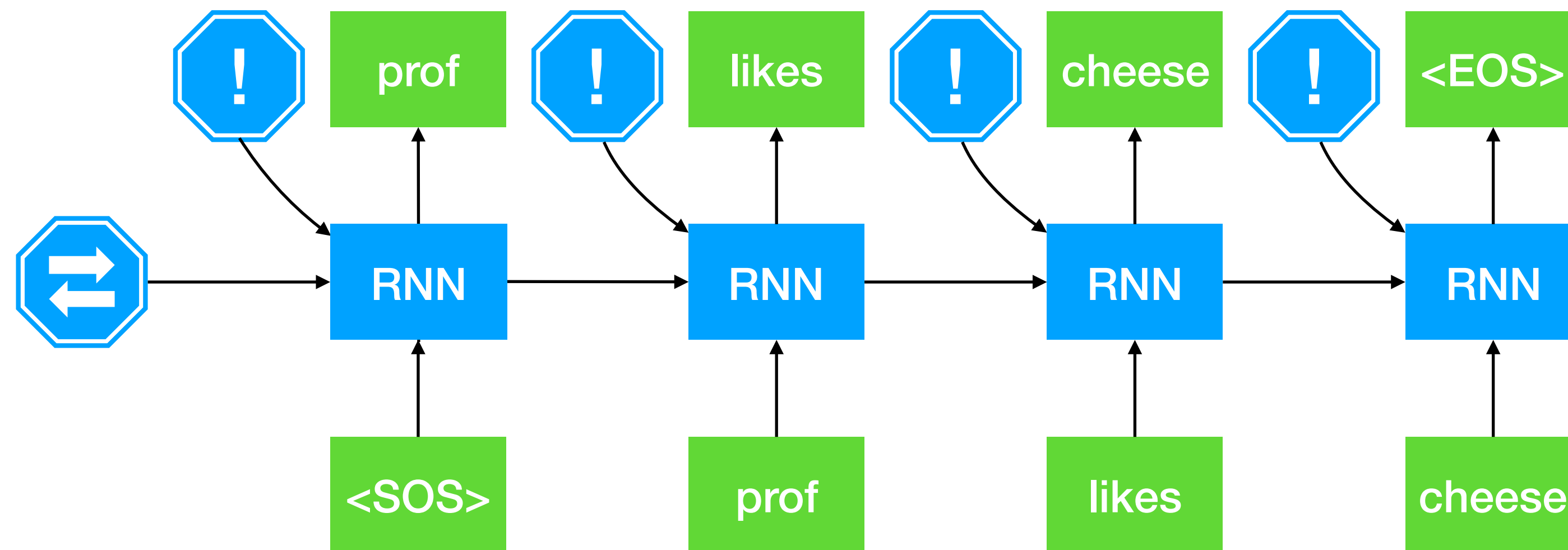
$$\alpha_t = \text{softmax}(score_{t,i})$$

$$context_t = \sum_i \alpha_{t,i} h_i^{enc}$$

$$o = \text{softmax}(W[context_t; h_t^{dec}] + b)$$

Concept

# Decoding with Attention

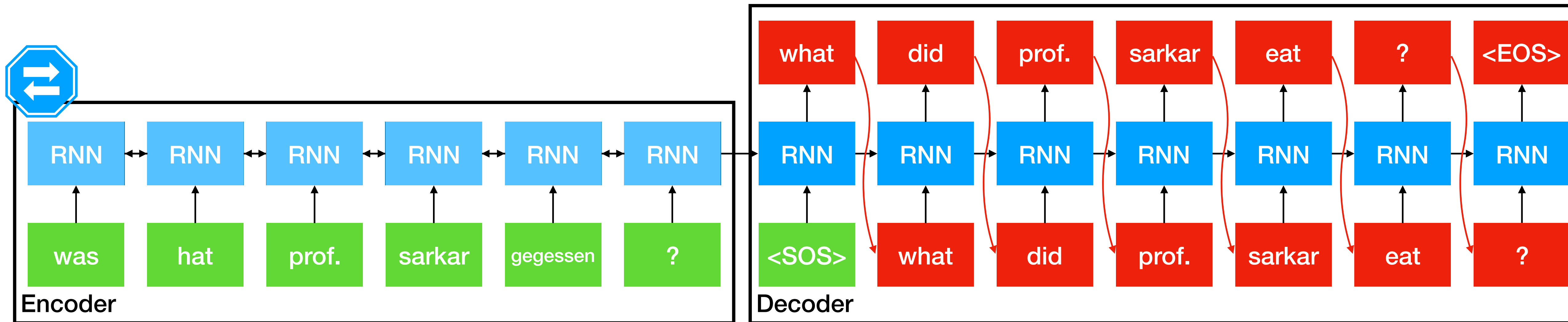


# A Tad More on Attention $\alpha$

$$\text{context}_t = \sum_i \alpha_{t,i} h_i^{\text{enc}} \quad \alpha_t = \text{softmax}(\text{score}_{t,i}) \quad \text{score}_{t,i} = f(h_i^{\text{enc}}, h_t^{\text{dec}})$$

- $f(x, y) = V \tanh(W'x + W''y)$
- $f(x, y) = x^T y$
- $f(x, y) = x^T W y$
- ... ..

# Sequence-to-Sequence

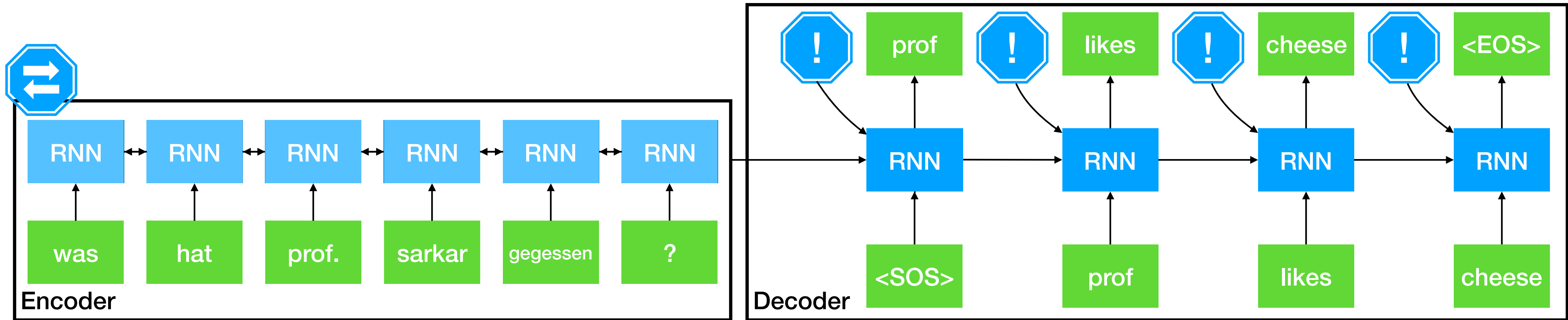


$$\Pr(E | F) = \prod_t \Pr(e'_t = e_t | F, e_{<t})$$

CLM

Review

# Attention



Review

"HW4 is out. Good luck."

**Next Tuesday:**

- Copy Mechanism
- BeamSearch
- [Extra] Beyond Seq2Seq: Attention is all you need
- [Extra] Beyond NMT

