

ATTENTIVE LANGUAGE MODELS

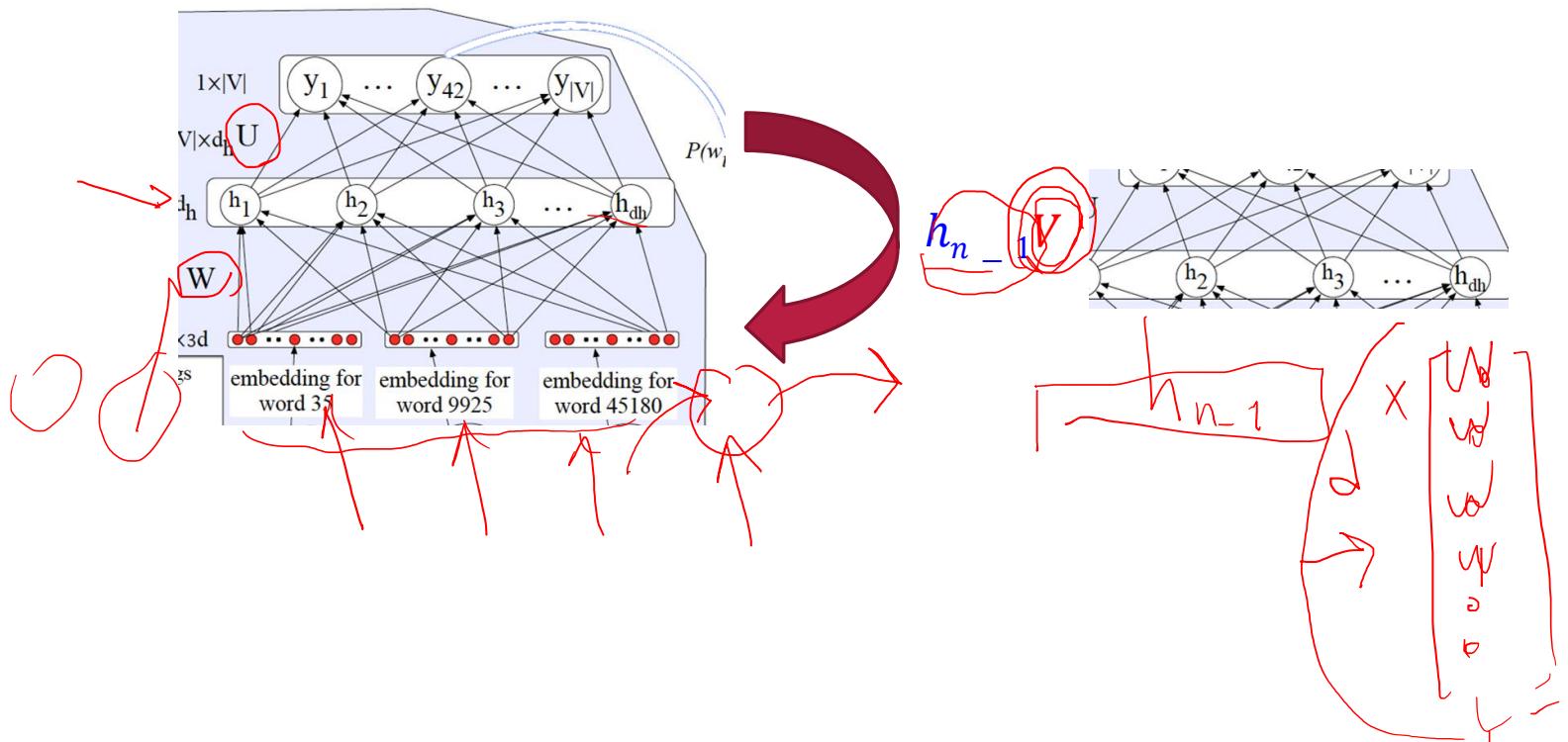
disclaimer:
the content is Hossein's understanding of methods

BERT: Bidirectional Encoder Representations from Transformers

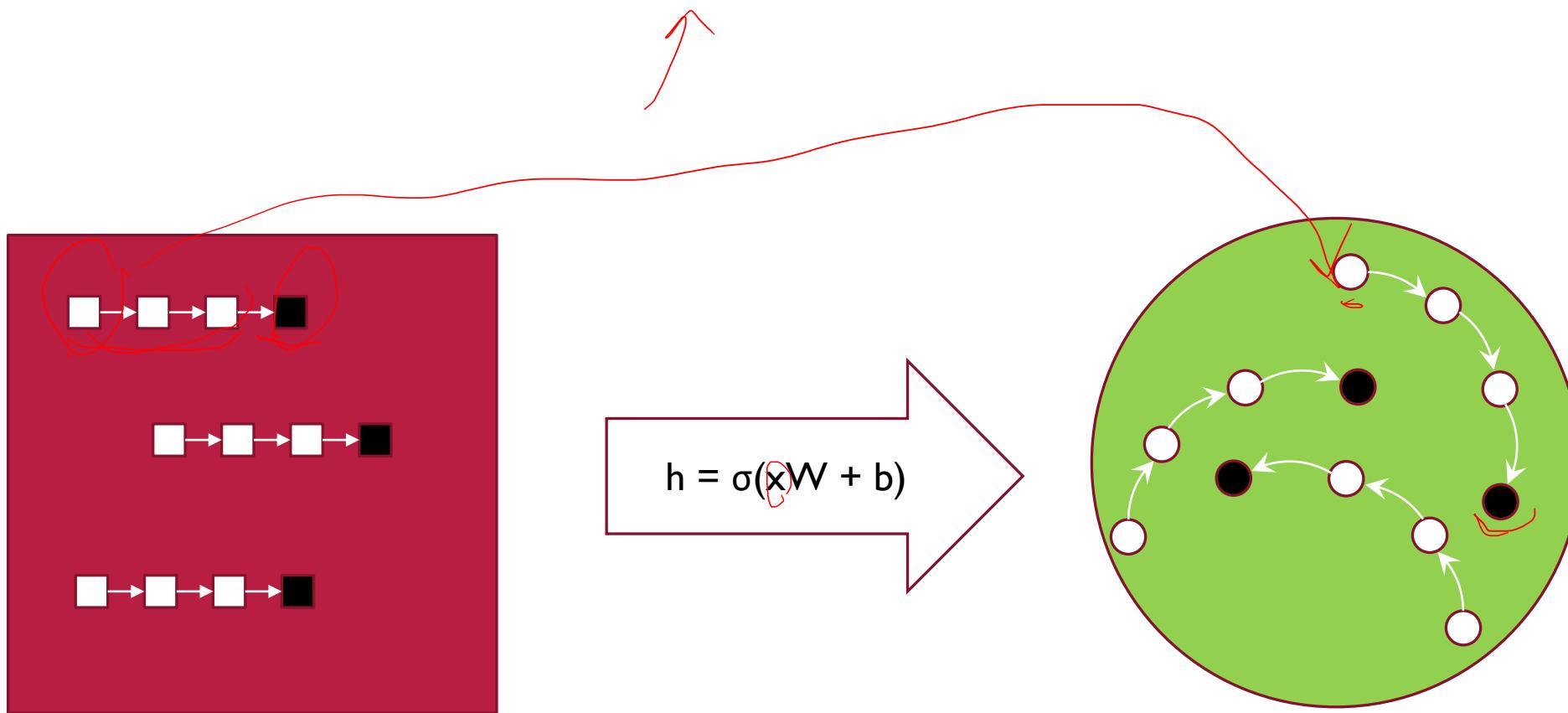
We finish the course without explaining BERT!

Recurrent Neural LM

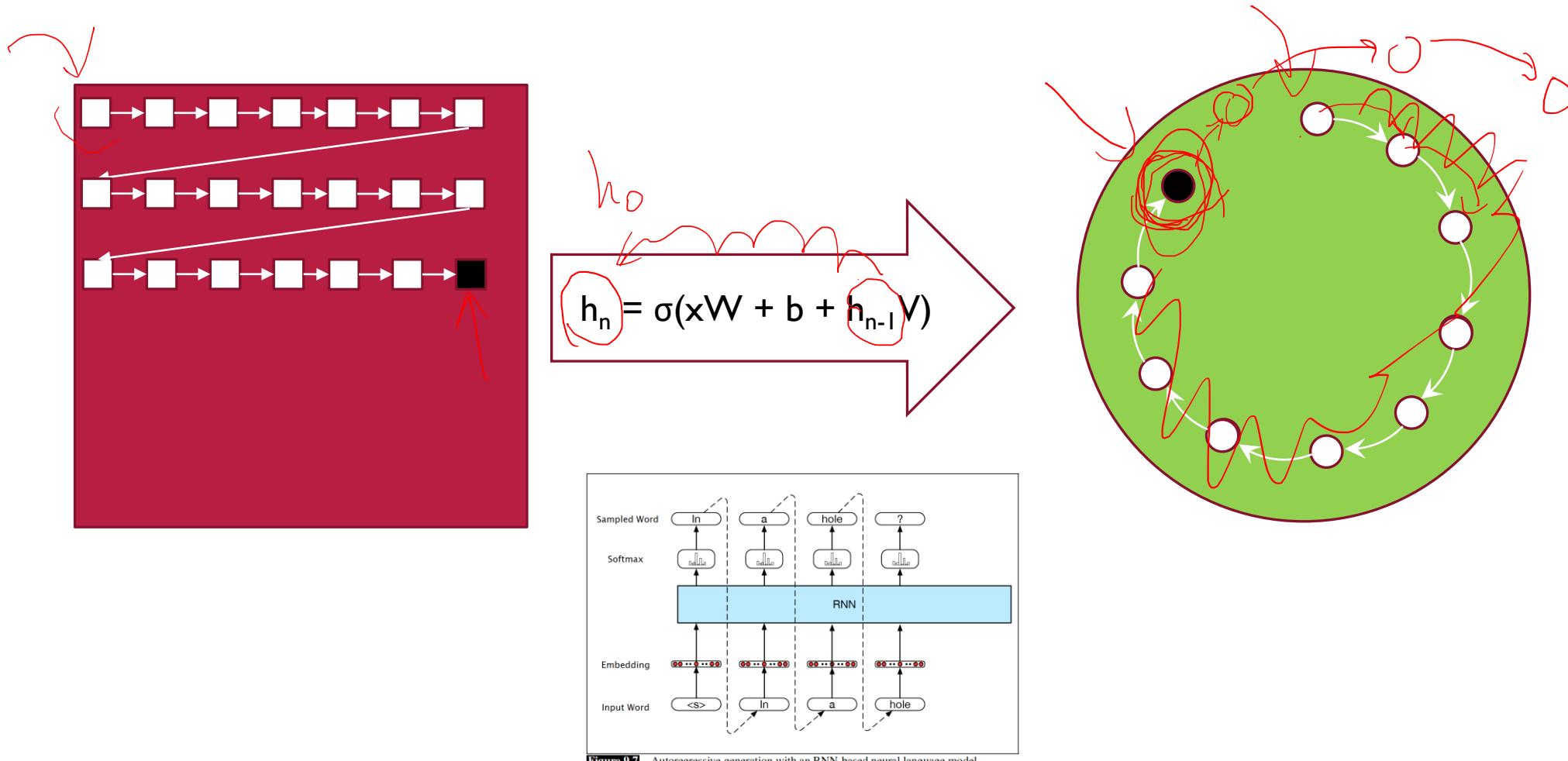
$$Y_n = \text{softmax}(h_n U)$$
$$h_n = \tanh(X_n W + b + h_{n-1} V)$$
$$X_n$$



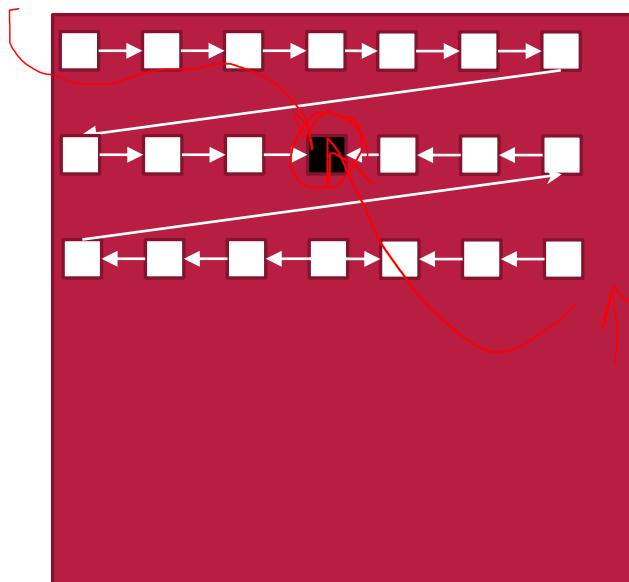
Neural LM



Recurrent Neural LM



Recurrent Neural LM: Bidirection



$$h_i = \sigma(xW + b + h_{i-1}V)$$
$$h_i = \sigma(xW' + b' + h_{i+1}V')$$

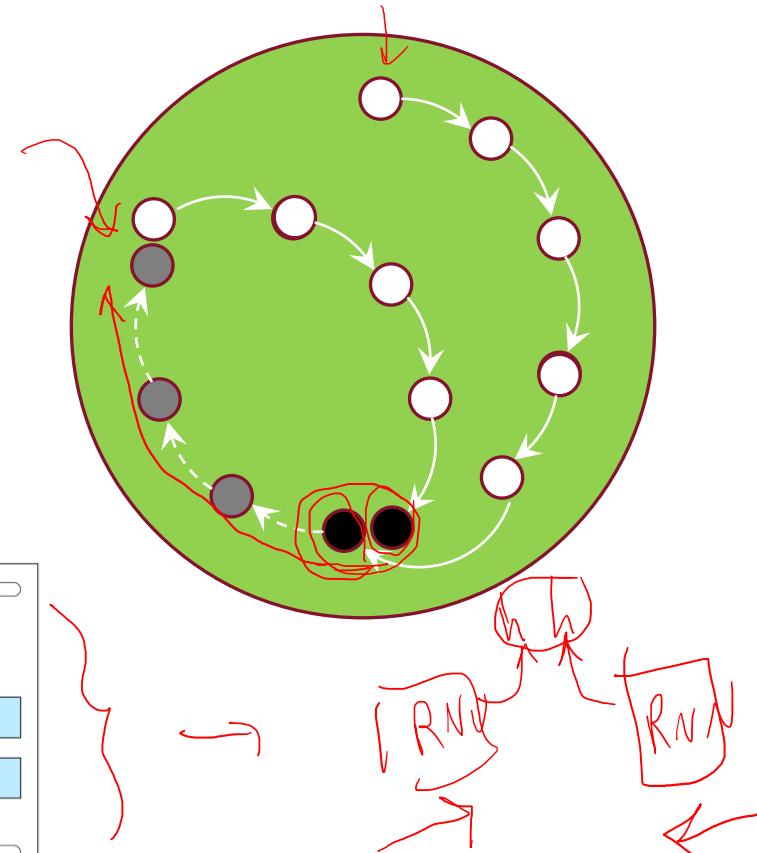
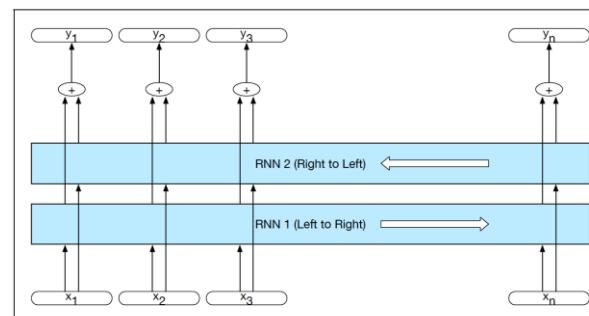
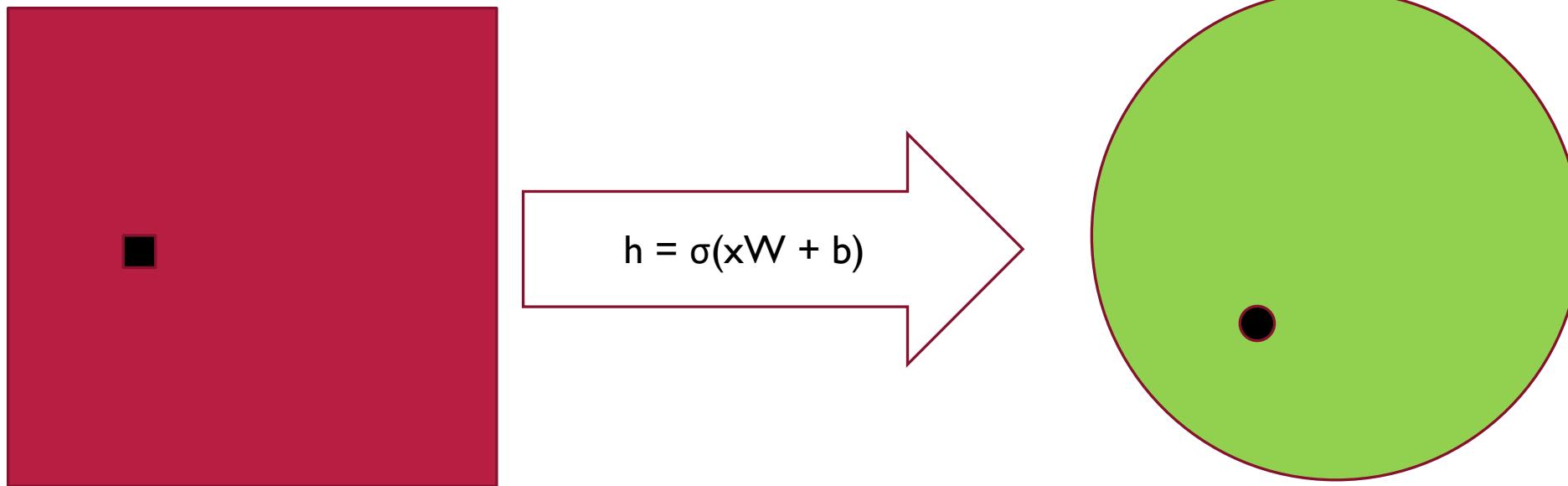


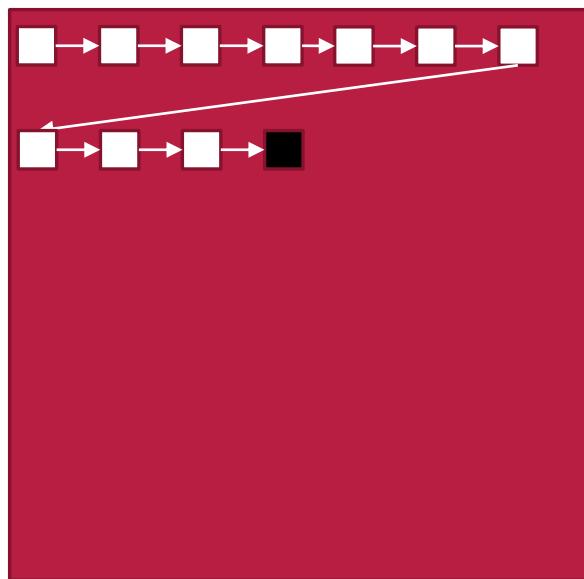
Figure 9.11 A bidirectional RNN. Separate models are trained in the forward and backward directions with the output of each model at each time point concatenated to represent the state of affairs at that point in time. The box wrapped around the forward and backward network emphasizes the modular nature of this architecture.

Seq2Seq

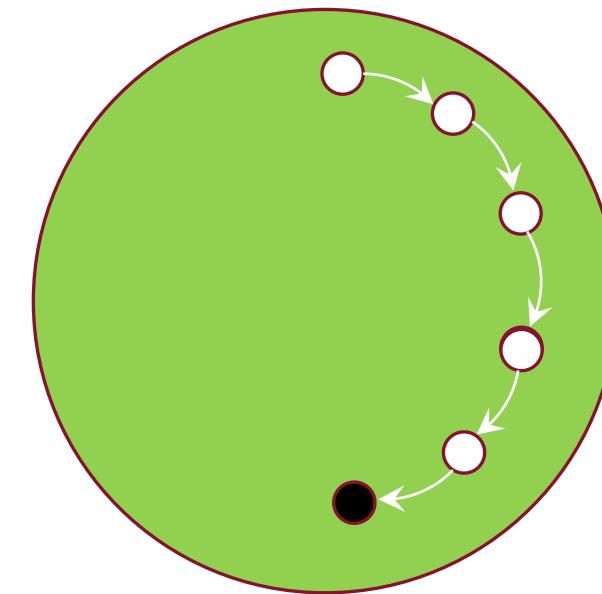
Neural LM



Recurrent Neural LM



$$h_i = \sigma(xW + b + h_{i-1}V)$$



Recurrent Neural LM: Encoder-Decoder

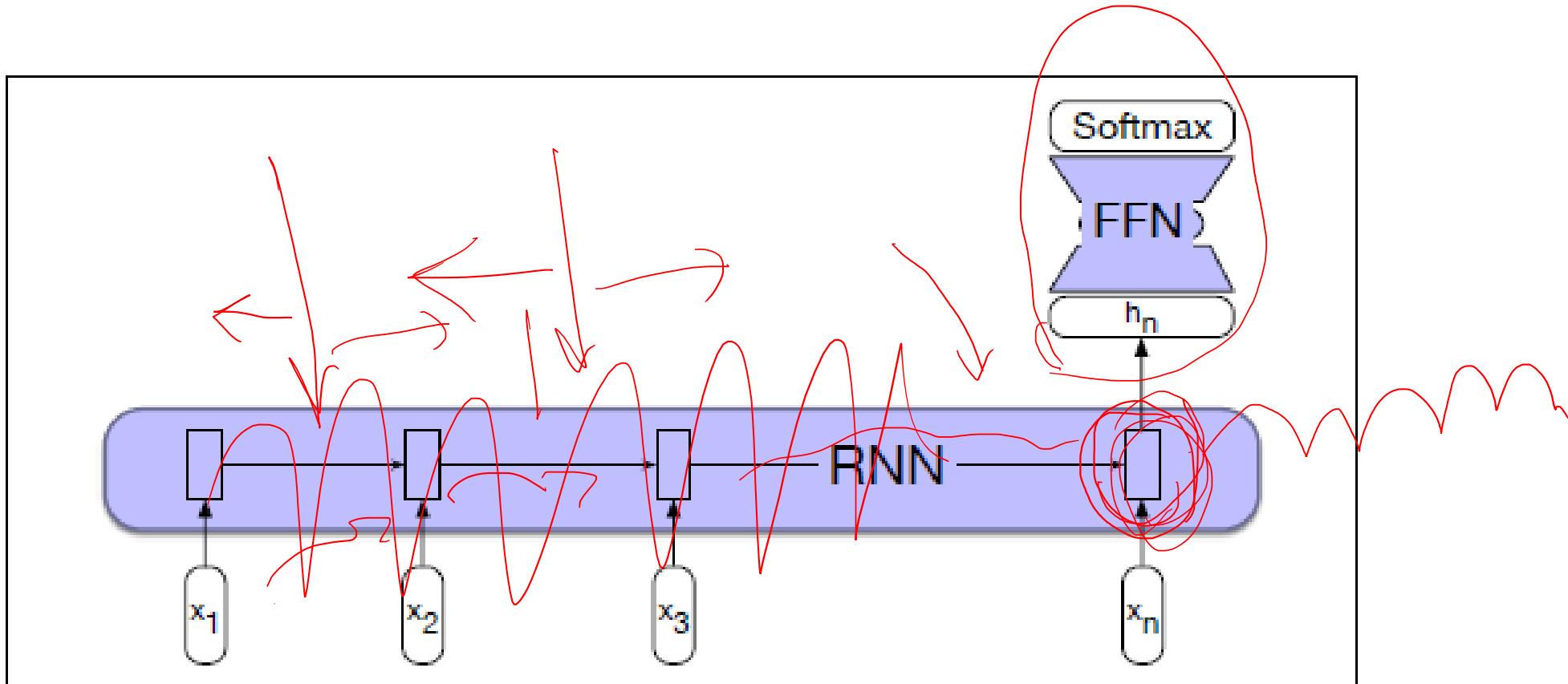
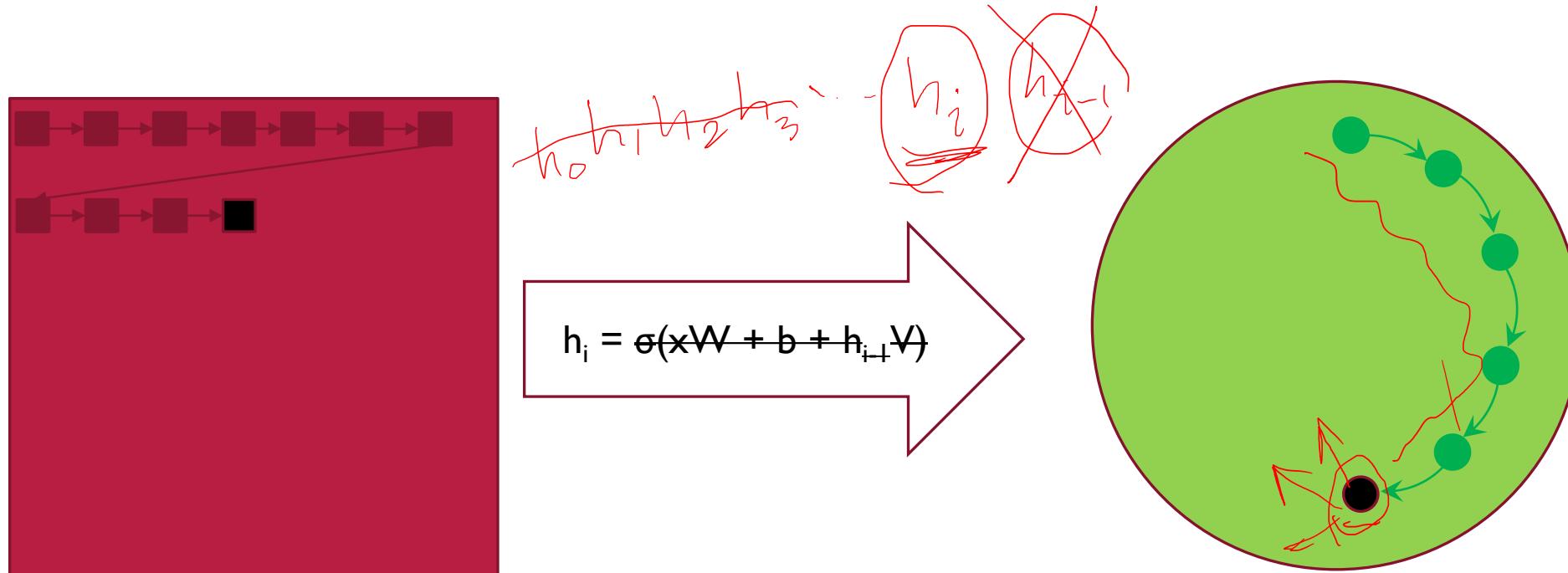


Figure 9.8 Sequence classification using a simple RNN combined with a feedforward network. The final hidden state from the RNN is used as the input to a feedforward network that performs the classification.

Recurrent Neural LM: Encoder-Decoder

A representation for the whole history



All history is **encoded** in h_i . We don't need ~~history~~!
 h_i is enough to **decode** the whole future.

Encoder-RNN

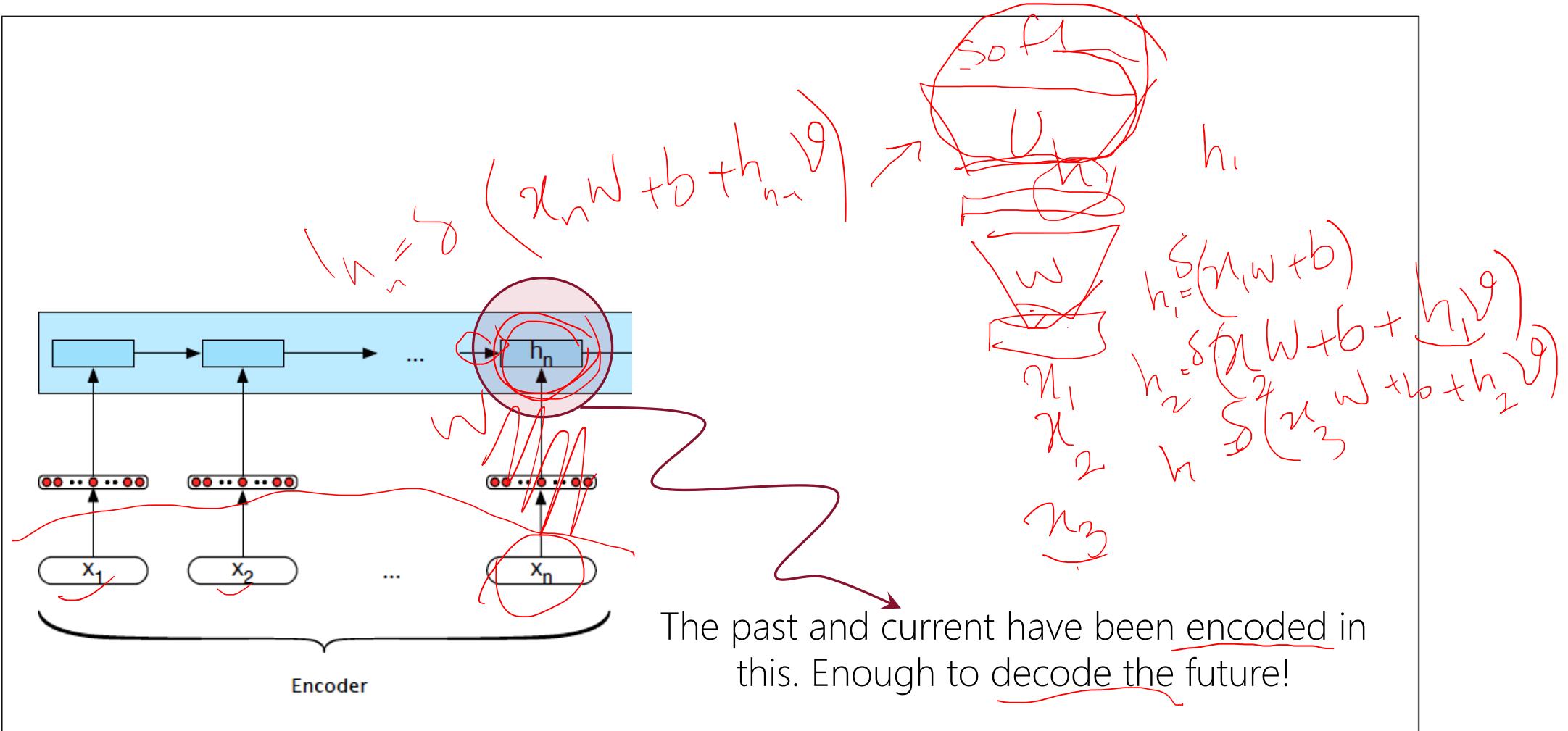
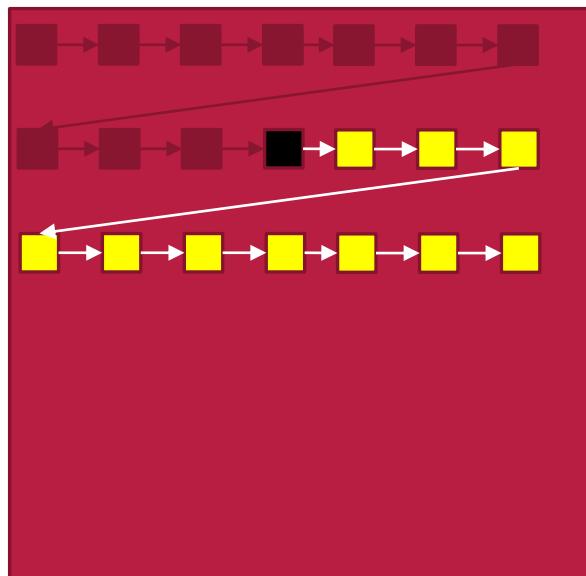
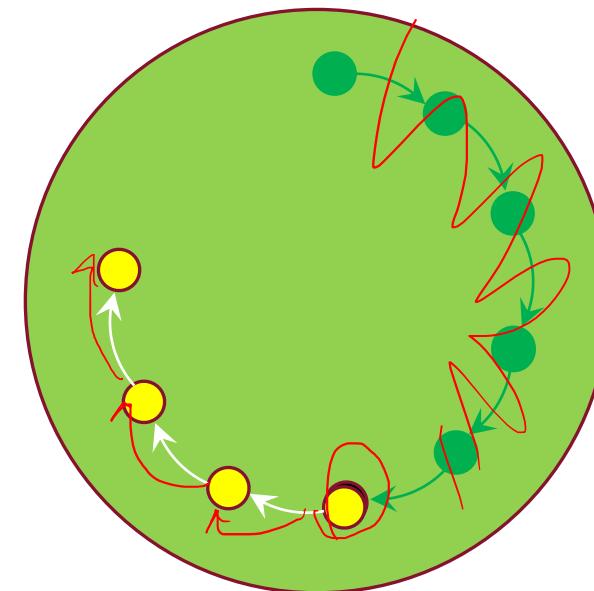


Figure 10.3 Basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN serves as the context for the decoder in its role as h_0 in the decoder RNN.

Recurrent Neural LM: Encoder-Decoder



$$\begin{aligned} h_{i+1} &= \sigma(xW + b + h_iV) \\ h_{i+2} &= \sigma(xW + b + h_{i+1}V) \\ &\dots \end{aligned}$$



Decoder-RNN (Test: Autoregressive)

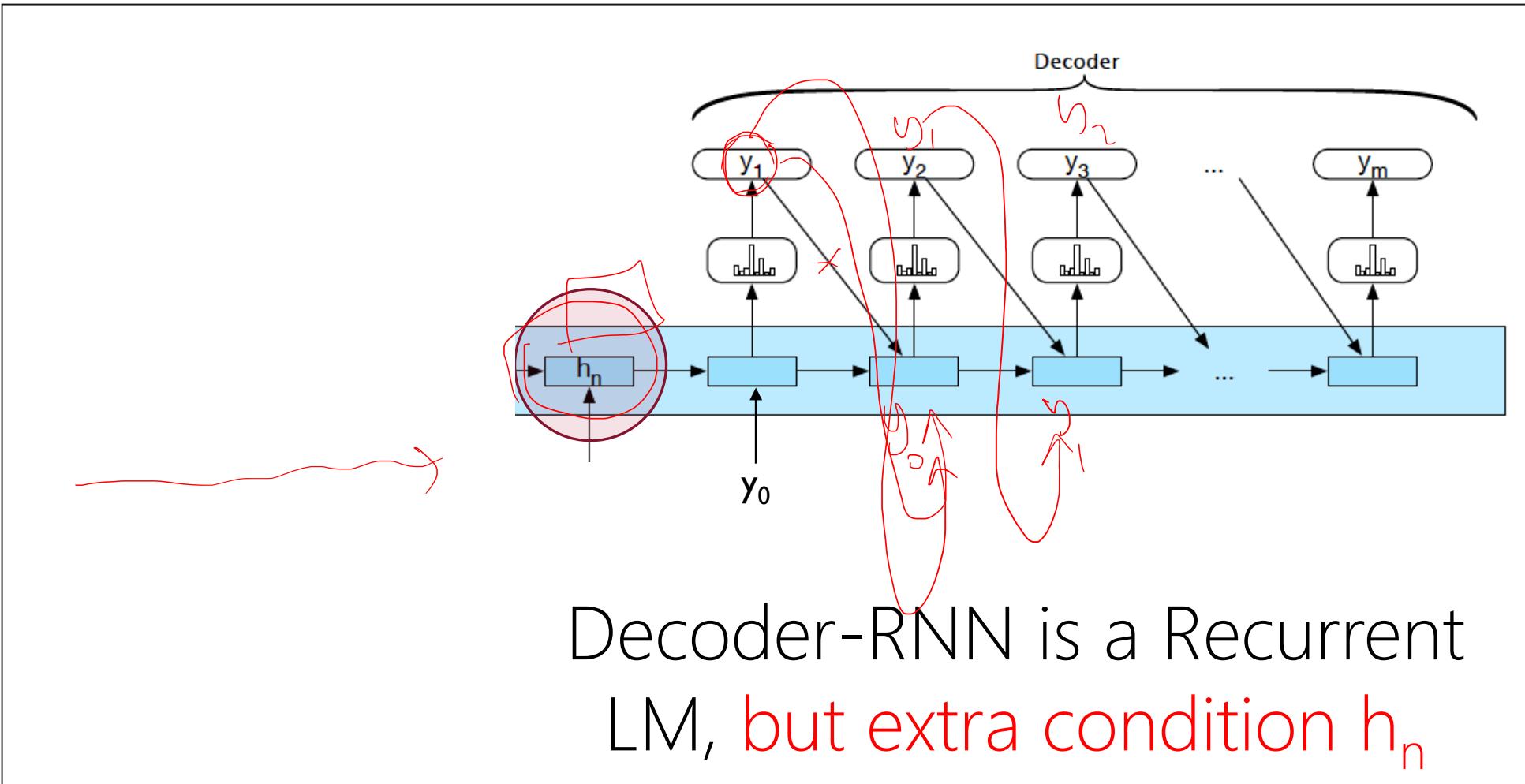


Figure 10.3 Basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN serves as the context for the decoder in its role as h_0 in the decoder RNN.



$$\frac{P(w_n | w_1, \dots, w_{n-1})}{P(w_i | w_{i-2}, w_{i-1})} = P(w_i | w_{i-1}, w_{i-2})$$

Conditional Language Modeling

The probability of the next word, given the previous words and conditioning context x:

$$p(w_i | x, w_1, w_2, \dots, w_{i-1})$$

Conditional Language Modeling

Condition x	Language Model
An author	A document <u>written by that author</u>
A topic label	An article about <u>that topic</u> RNN
{SPAM, NOT_SPAM}	An email <u>...</u>
A sentence in French	Its English translation RNN
An image	A text description of the image RNN
A document	Its summary RNN
Meteorological measurements	A weather report
A question	Its answer
Acoustic signal	Transcription of speech
X-Ray	Medical report
Code	Documentation help
Code in C	Code in Java

(Handwritten annotations in red)

- Condition x: An author → A document written by that author
- Condition x: A topic label → An article about that topic RNN
- Condition x: {SPAM, NOT_SPAM} → An email ...
- Condition x: A sentence in French → Its English translation RNN
- Condition x: An image → A text description of the image RNN
- Condition x: A document → Its summary RNN
- Condition x: Meteorological measurements → A weather report
- Condition x: A question → Its answer
- Condition x: Acoustic signal → Transcription of speech
- Condition x: X-Ray → Medical report
- Condition x: Code → Documentation help
- Condition x: Code in C → Code in Java

Decoder-RNN (Train: Teacher Forcing)

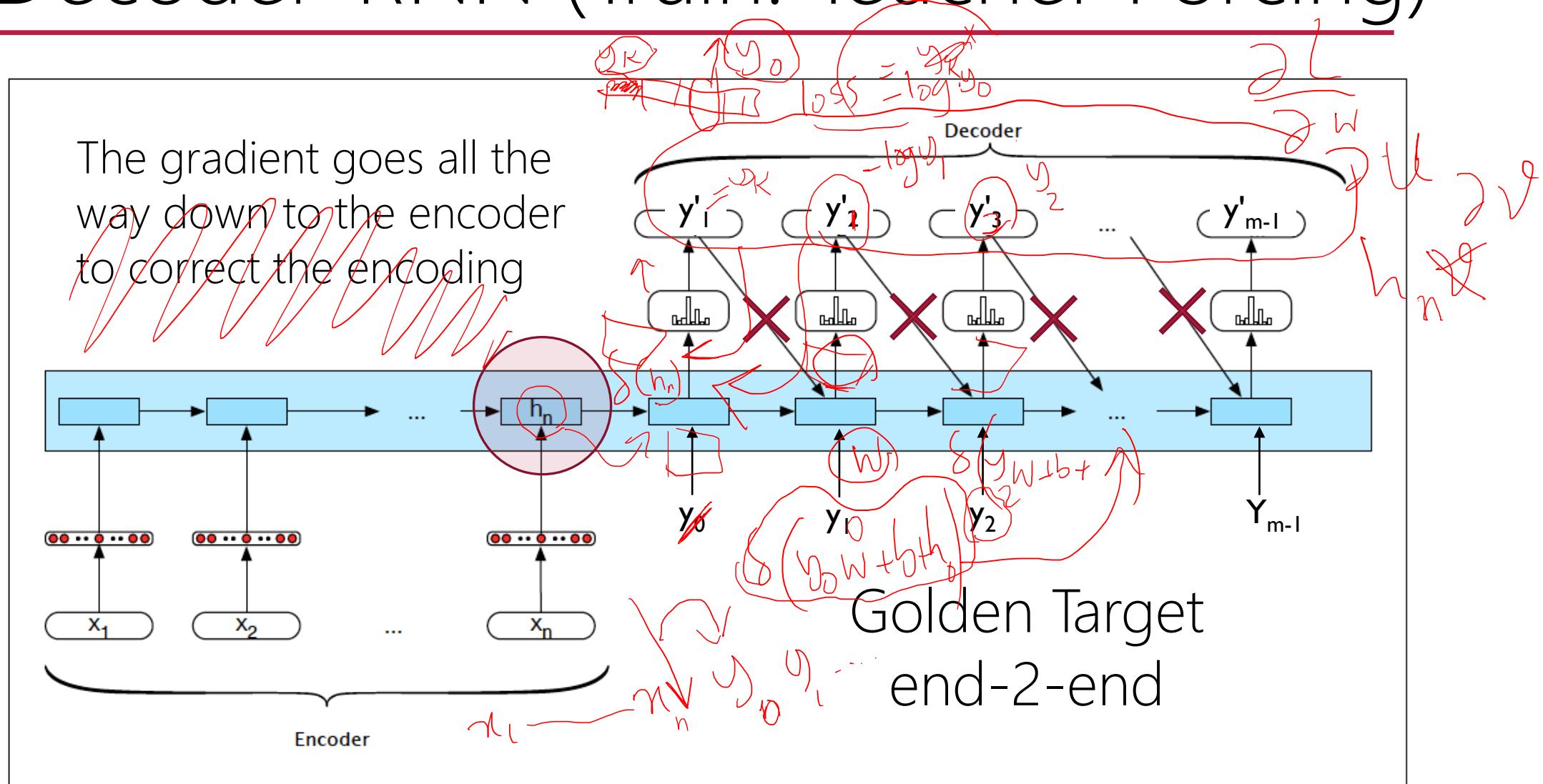


Figure 10.3 Basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN serves as the context for the decoder in its role as h_0 in the decoder RNN.

Decoder-RNN (Train: Teacher Forcing)

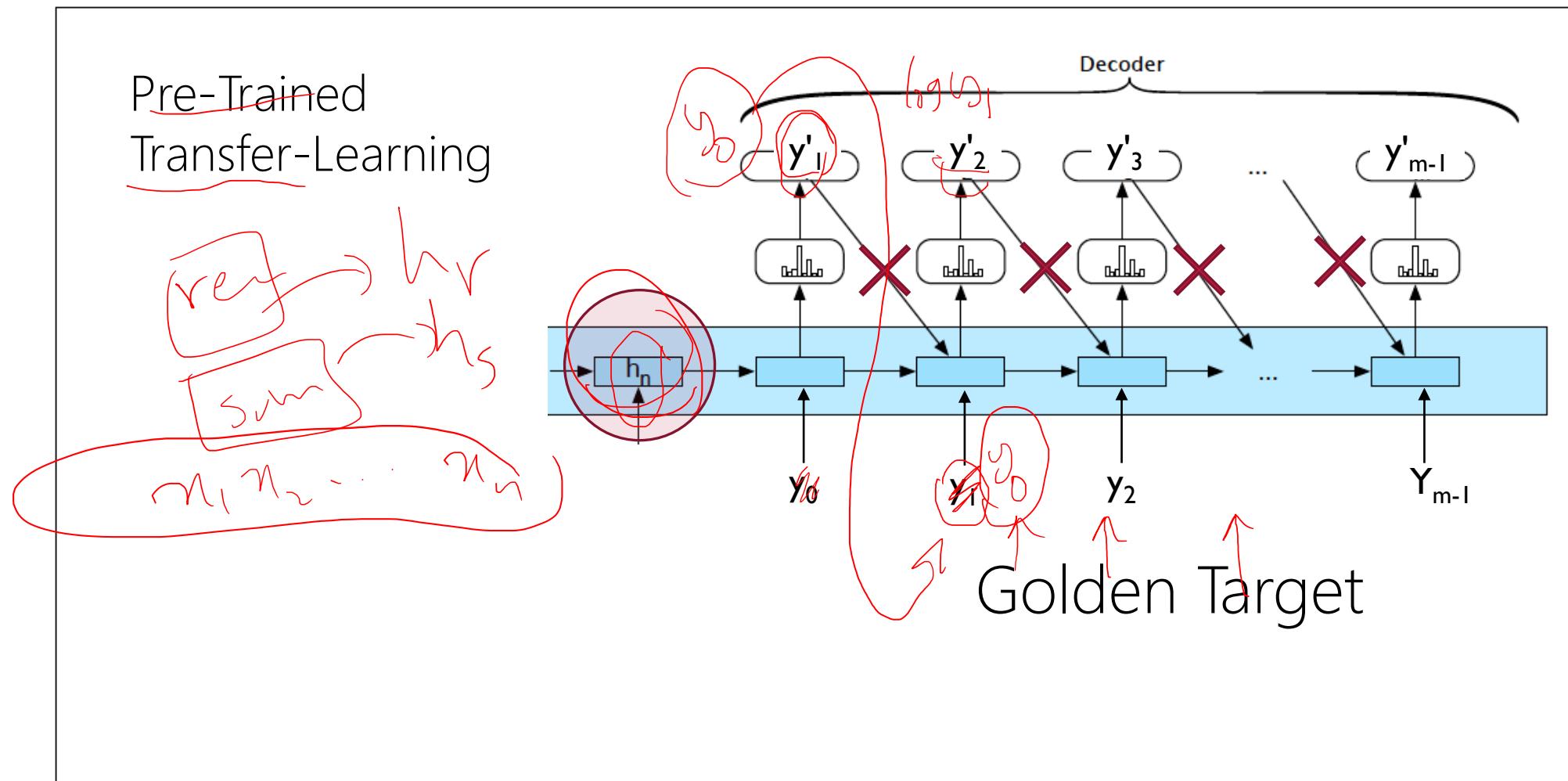


Figure 10.3 Basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN serves as the context for the decoder in its role as h_0 in the decoder RNN.

Sequence2Sequence

Sutskever, I., Vinyals, O., & Le, Q.V. (2014). Sequence to sequence learning with neural networks. In NIPS.

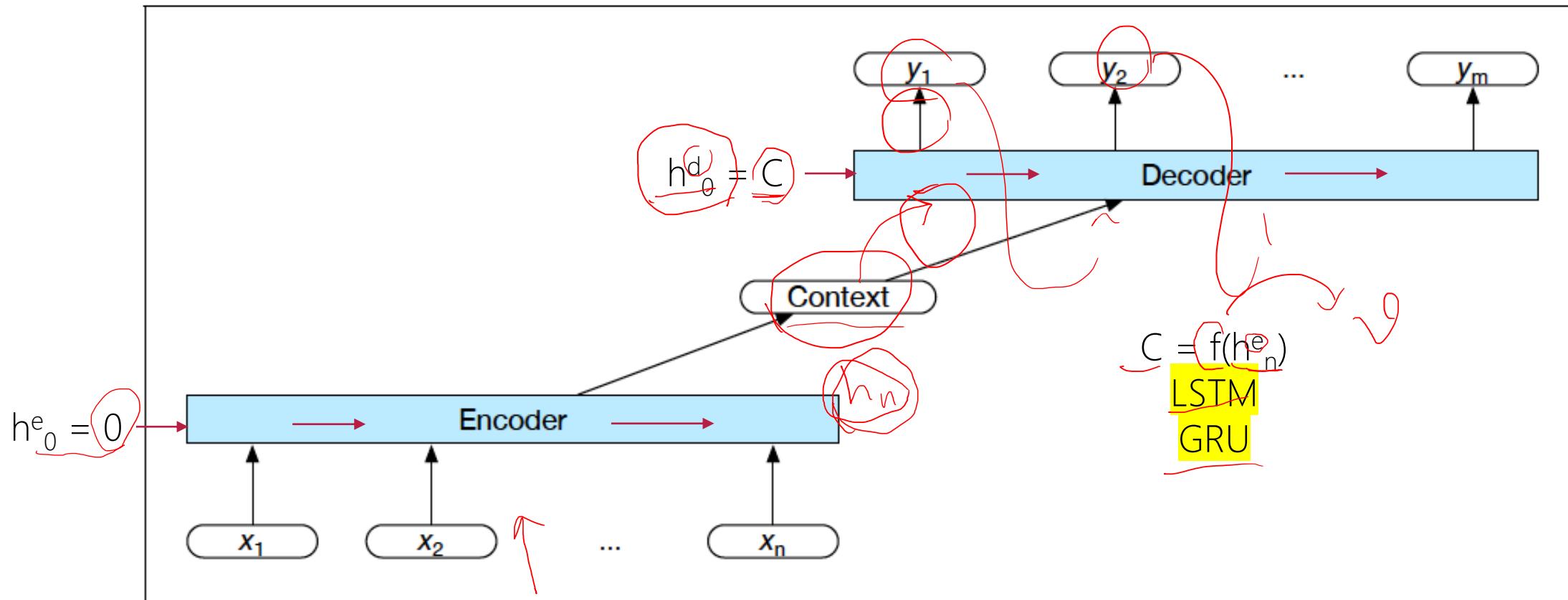
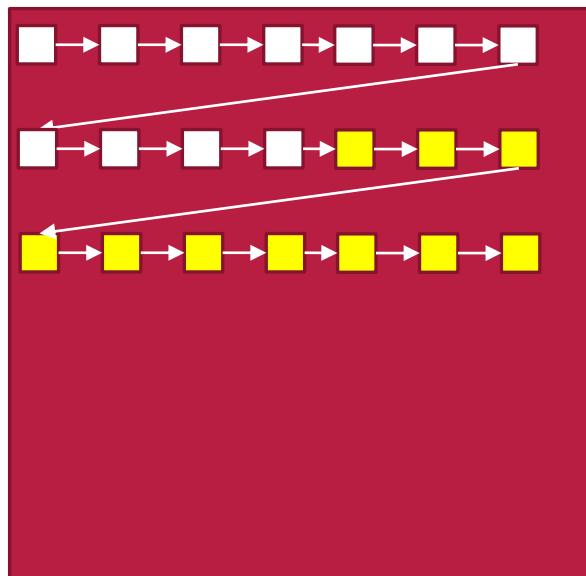


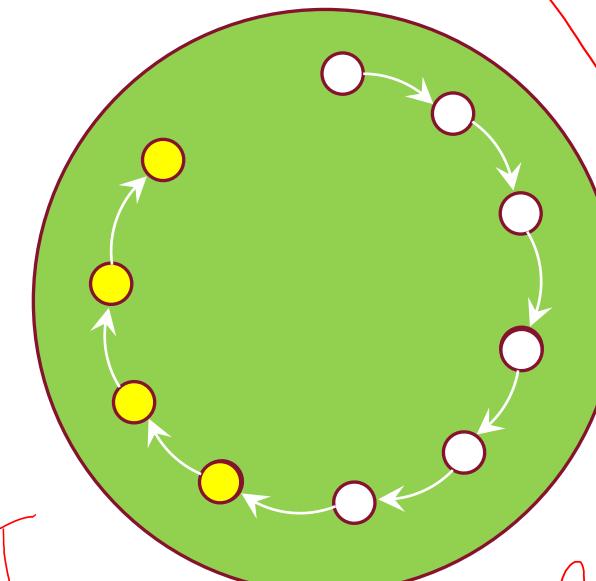
Figure 10.4 Basic architecture for an abstract encoder-decoder network. The context is a function of the vector of contextualized input representations and may be used by the decoder in a variety of ways.

Seq2Seq: Bitext Translation

For machine, they are just sequences of tokens!



$$h_i = \sigma(xW + b + h_{i-1}V)$$



< s > Prague Stock Market falls to minus by the end of the trading day < /s > < s > Die Prager Börse stürzt gegen Geschäftsschluss ins Minus < /s >

Seq2Seq: Bitext Translation

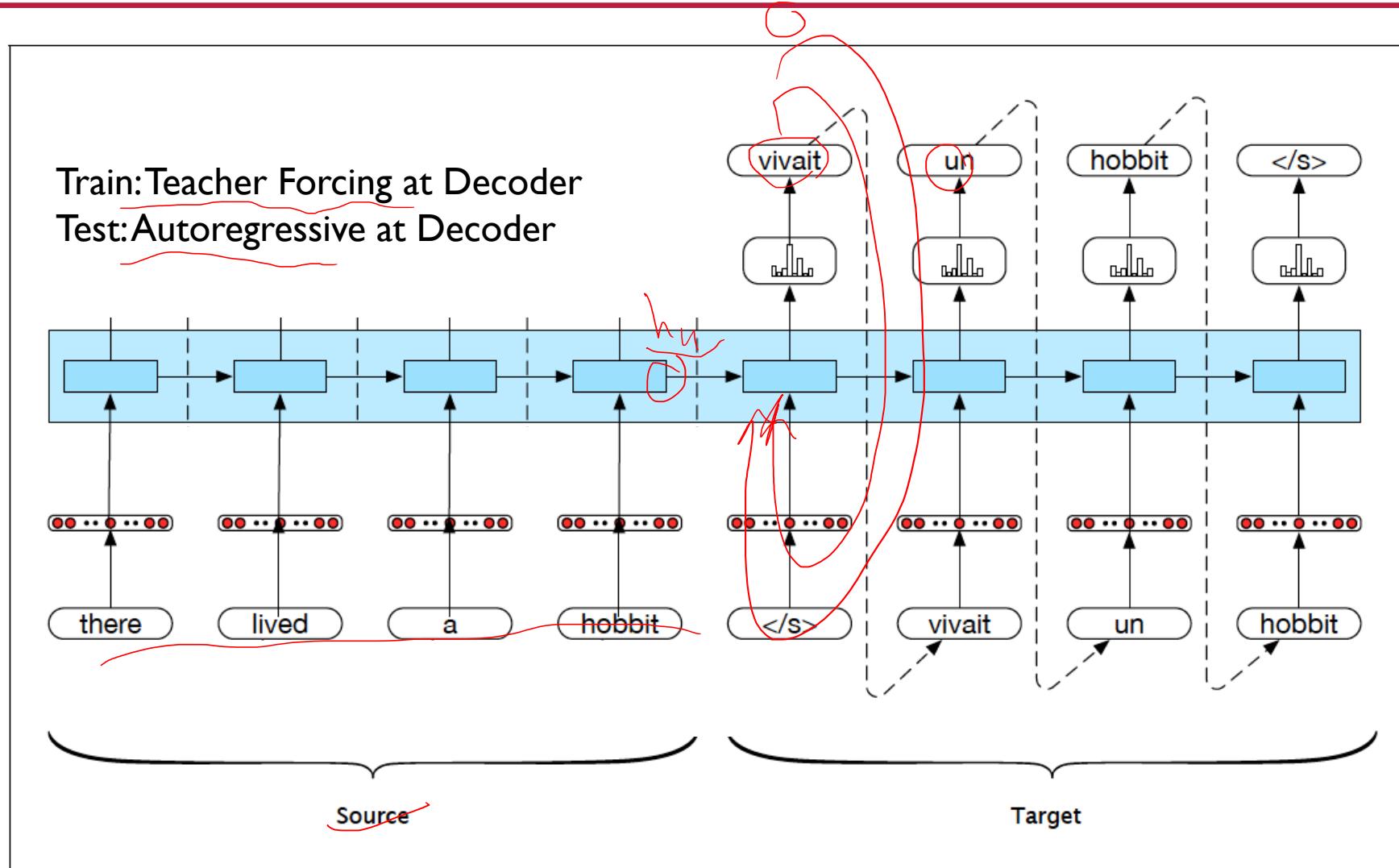
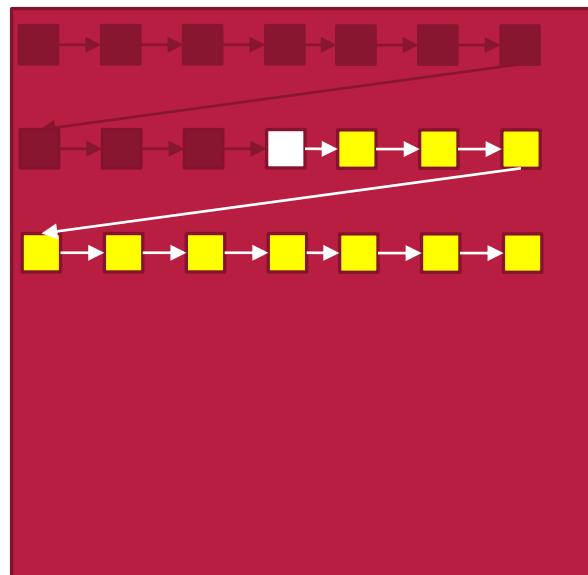


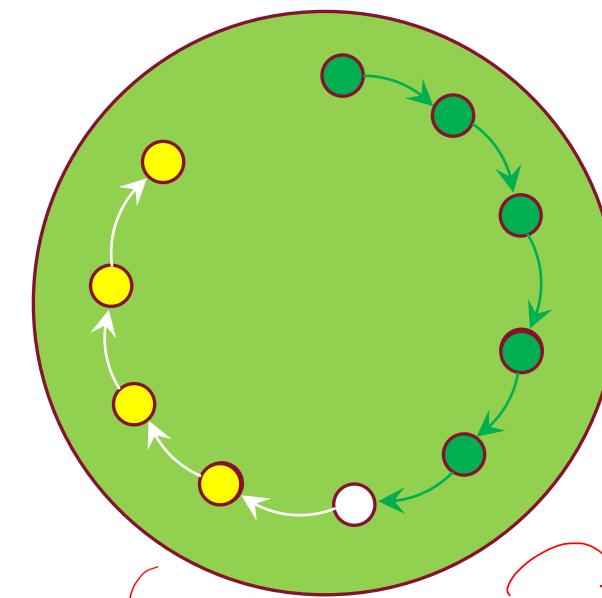
Figure 10.2 Training setup for a neural language model approach to machine translation. Source-target bitexts are concatenated and used to train a language model.

Seq2Seq: Bitext Translation

For machine, they are just sequences of tokens!



$$h_i = \sigma(xW + b + h_{i-1}V)$$



~~< /s > Prague Stock Market falls to minus by the end of the trading day < /s >~~ ~~< /s > Die Prager Börse stürzt gegen Geschäftsschluss ins Minus < /s >~~

$AU(s_1) = 2$ $AU(s_2) = 10$ $AU(s_3) = 100$



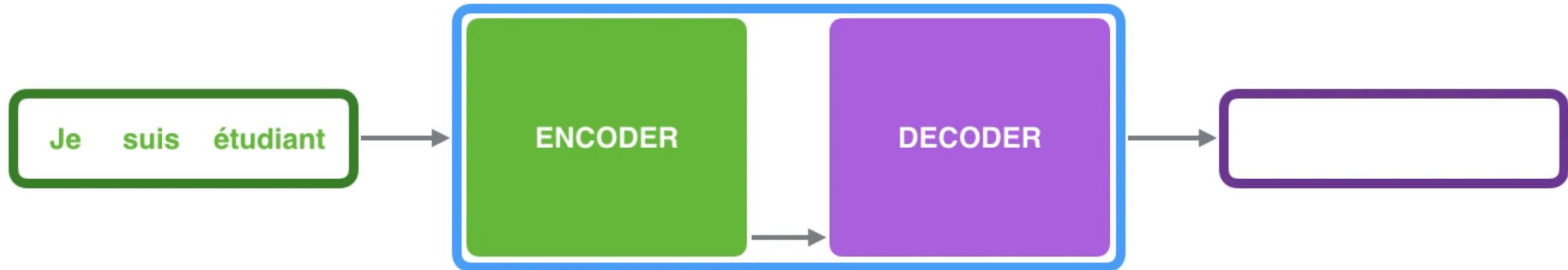
Jay Alammar

Visualizing machine learning one concept at a time.
@JayAlammar on Twitter. YouTube Channel

[Blog](#) [About](#)

Time step:

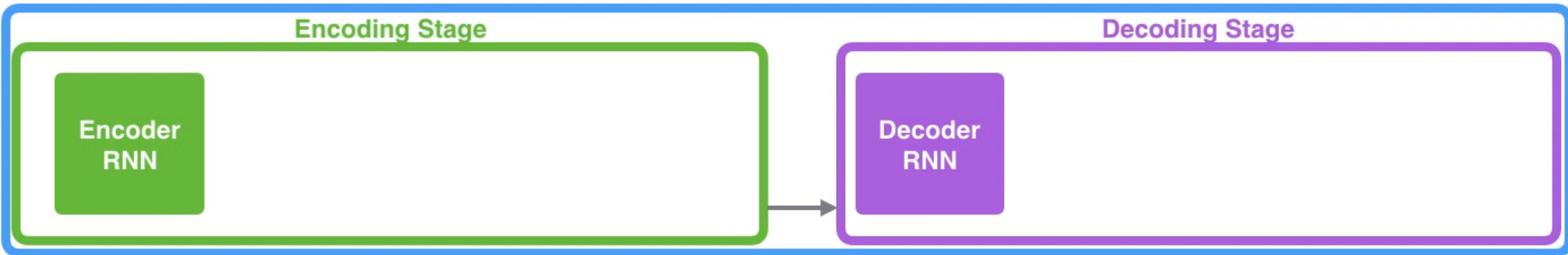
Neural Machine Translation SEQUENCE TO SEQUENCE MODEL





Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL



Je suis étudiant

Question to Answers

QA

Q A

Q ?

Anything to Sequence

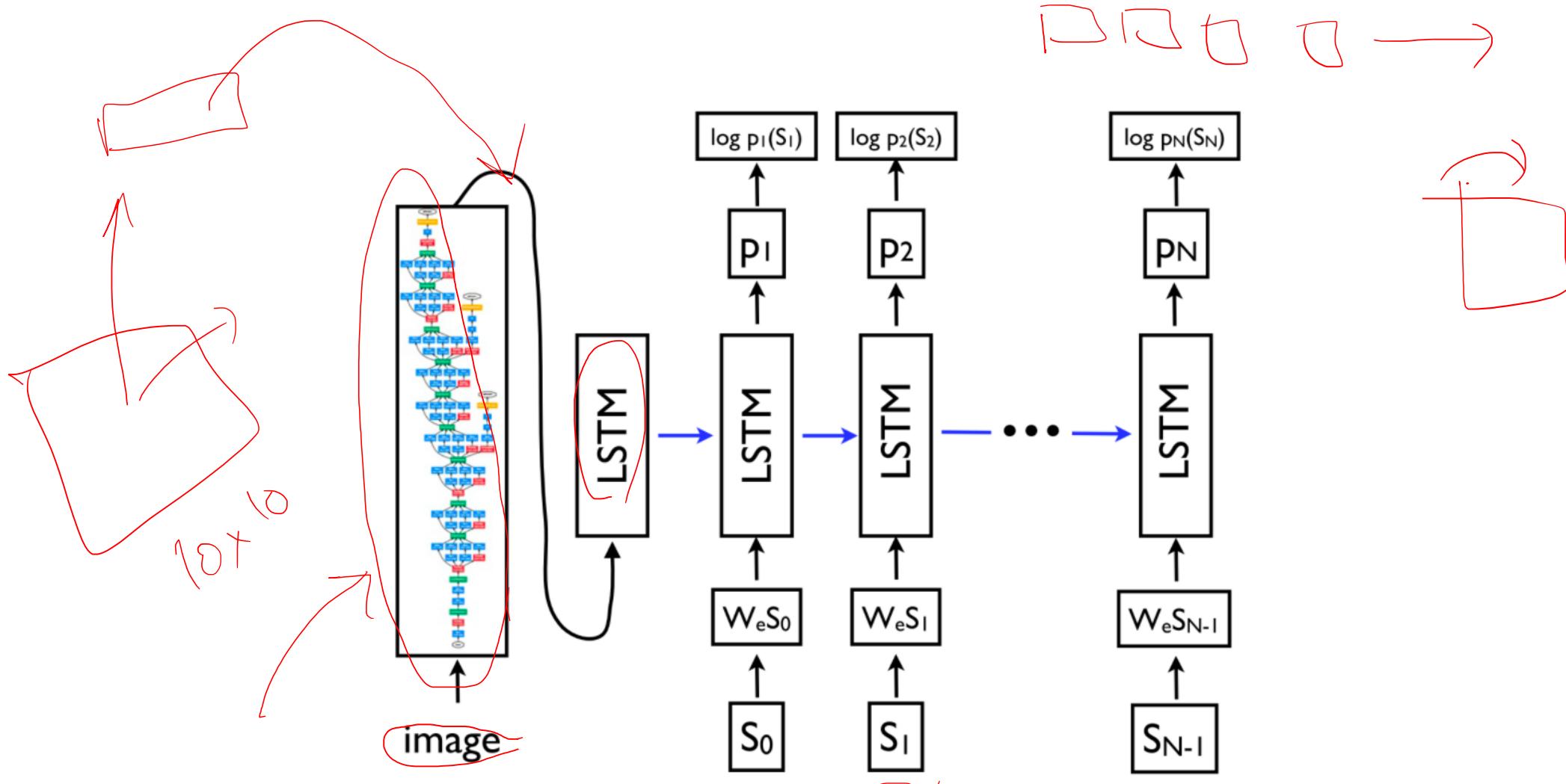


Figure 9: Generating a caption based on an image (Vinyals et al., 2015)

Seq2Seq: Information Bottleneck

Although we have LSTM, GRU, h_n falls short to encode *all* the past to decode *all* the future when the sequence becomes long!

Too much pressure to a single vector!

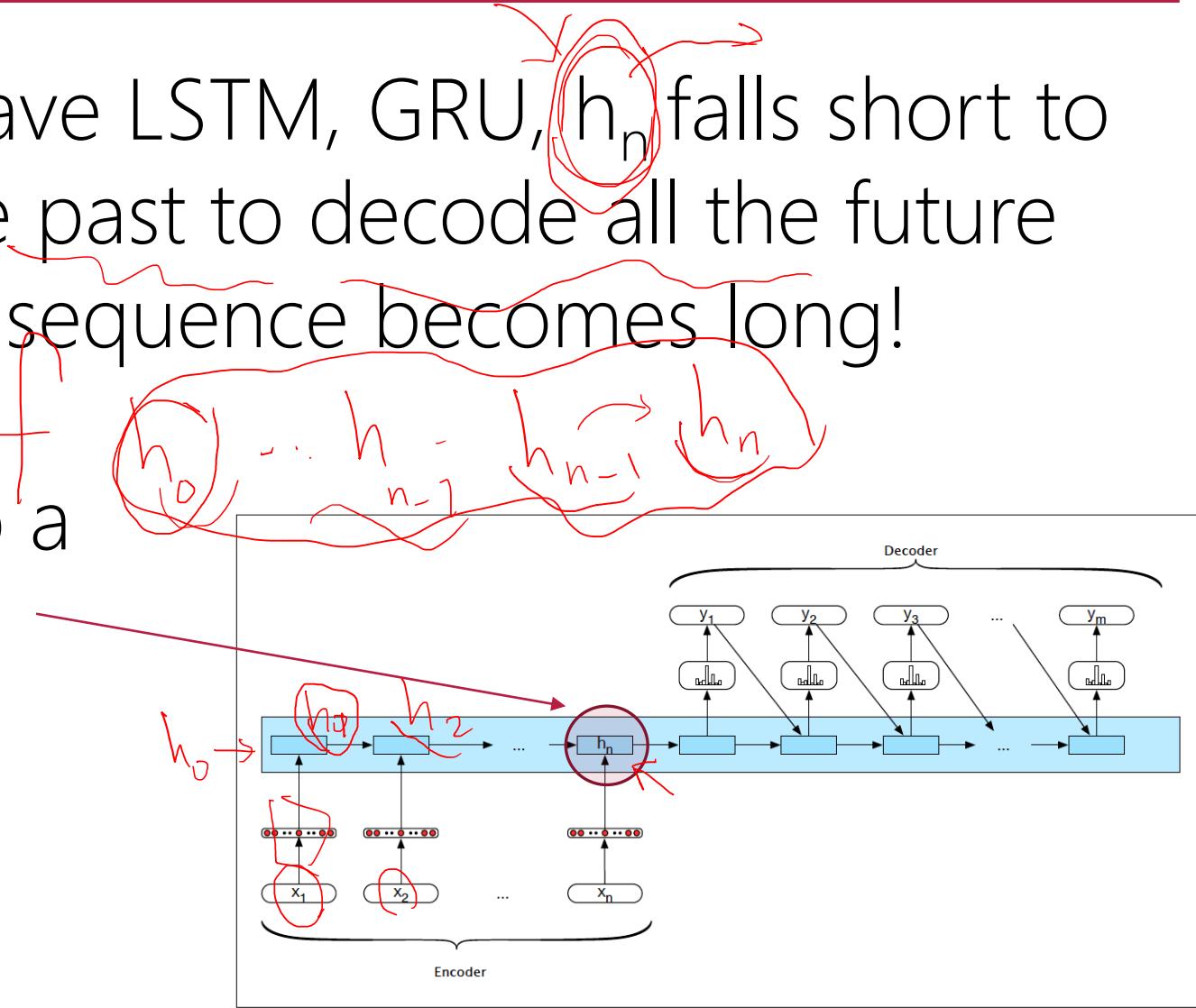


Figure 10.3 Basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN serves as the context for the decoder in its role as h_0 in the decoder RNN.

Look into the Past based on future!

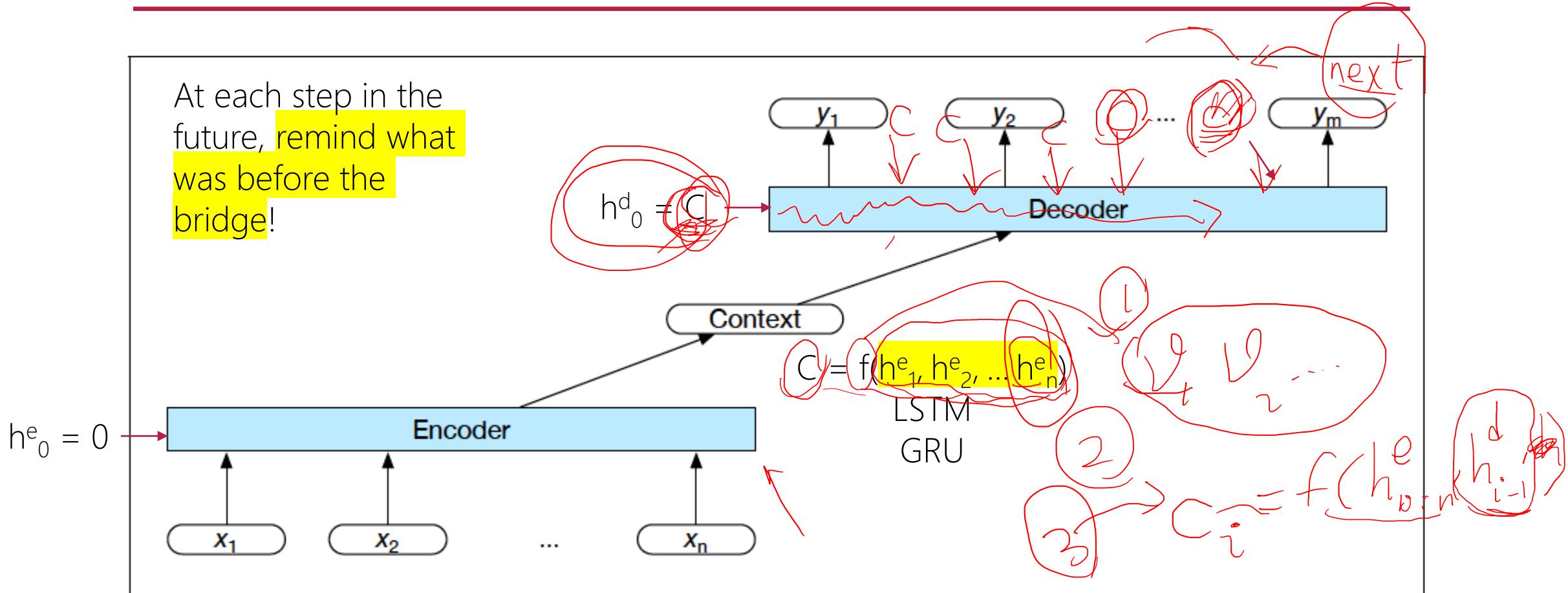


Figure 10.4 Basic architecture for an abstract encoder-decoder network. The context is a function of the vector of contextualized input representations and may be used by the decoder in a variety of ways.



Jay Alammar

Visualizing machine learning one concept at a time

[Blog](#) [About](#)

Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



Je suis étudiant

Look into the Past based on future!

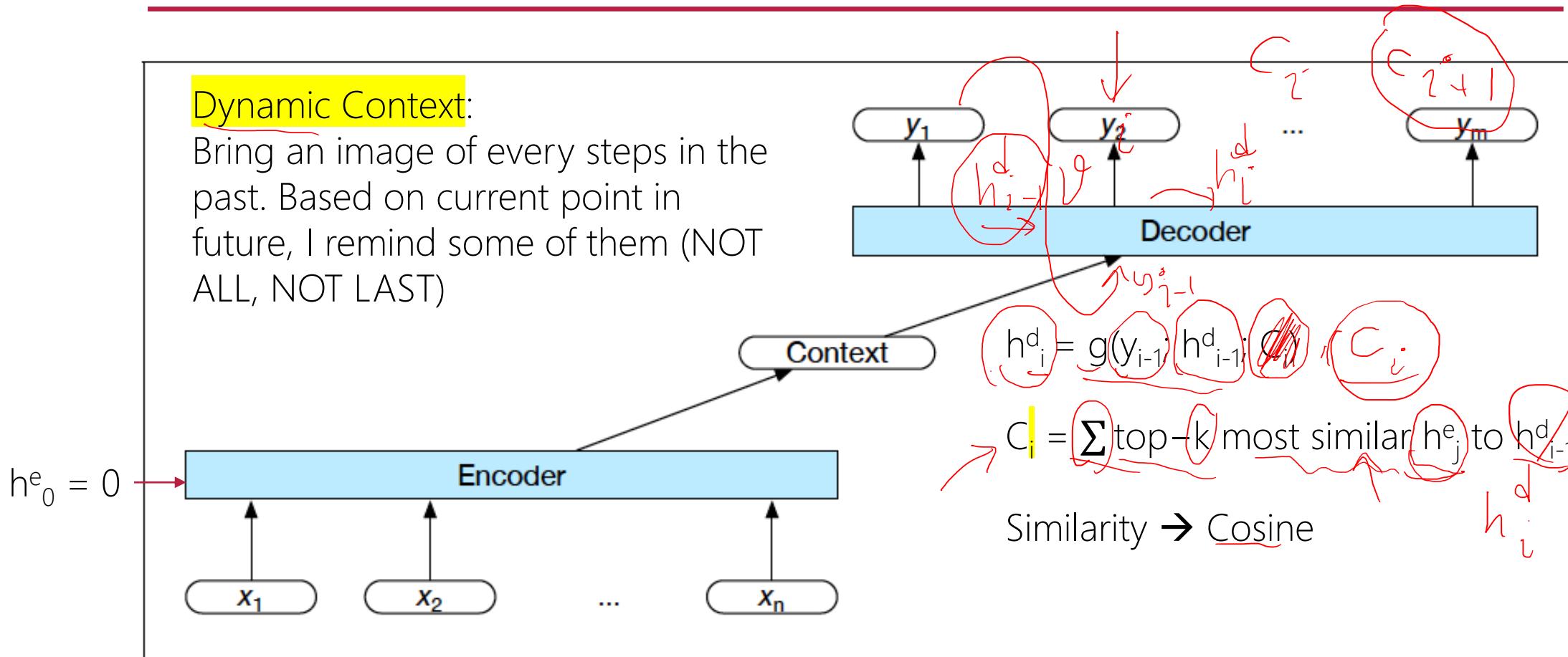


Figure 10.4 Basic architecture for an abstract encoder-decoder network. The context is a function of the vector of contextualized input representations and may be used by the decoder in a variety of ways.

Look into the Past based on future!

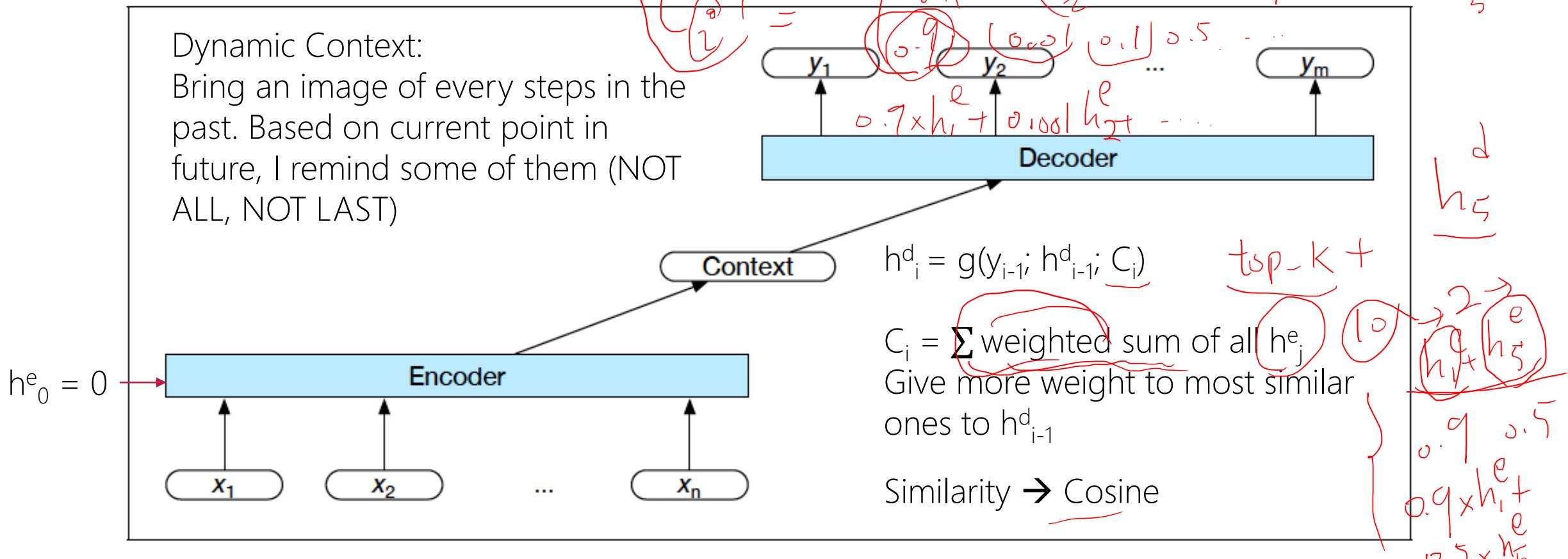
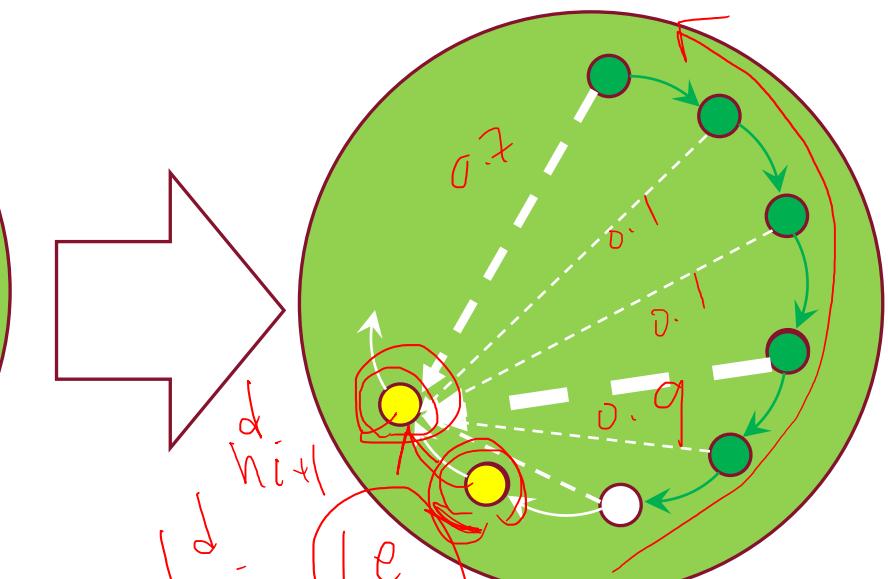
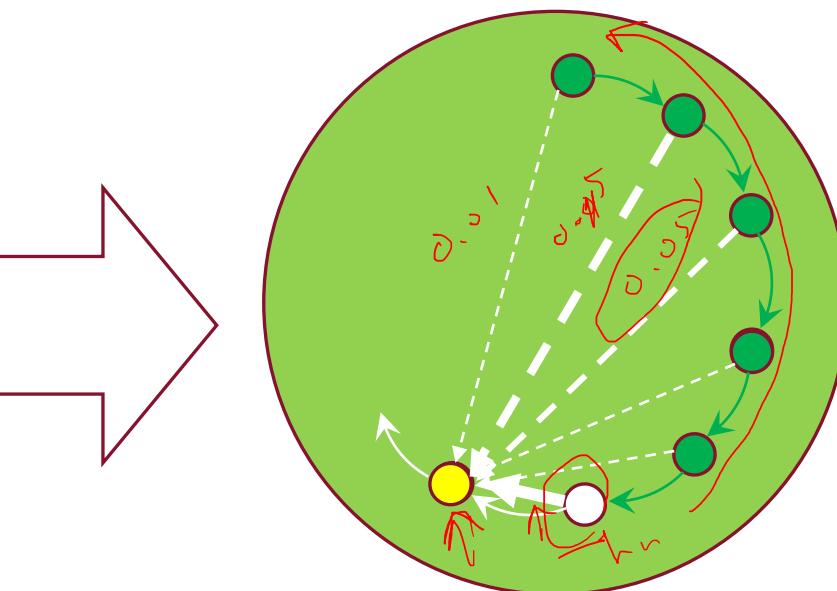
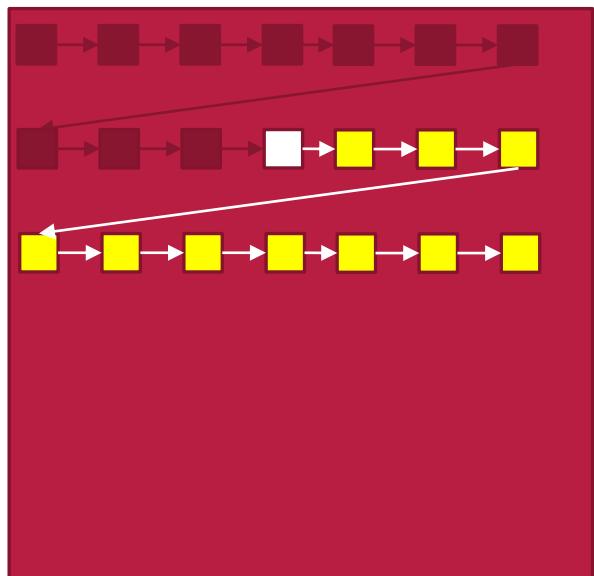


Figure 10.4 Basic architecture for an abstract encoder-decoder network. The context is a function of the vector of contextualized input representations and may be used by the decoder in a variety of ways.

Recurrent + Attention!

For machine, they are just sequences of tokens!



< s > Prague Stock Market falls to minus by the end of the trading day < /s > < s > Die Prager Börse stürzt gegen Geschäftsschluss ins Minus < /s >



Jay Alammar

Visualizing machine learning one concept at a time.

[@JayAlammar](#) on Twitter. [YouTube Channel](#)

[Blog](#) [About](#)

Attention at time step 4



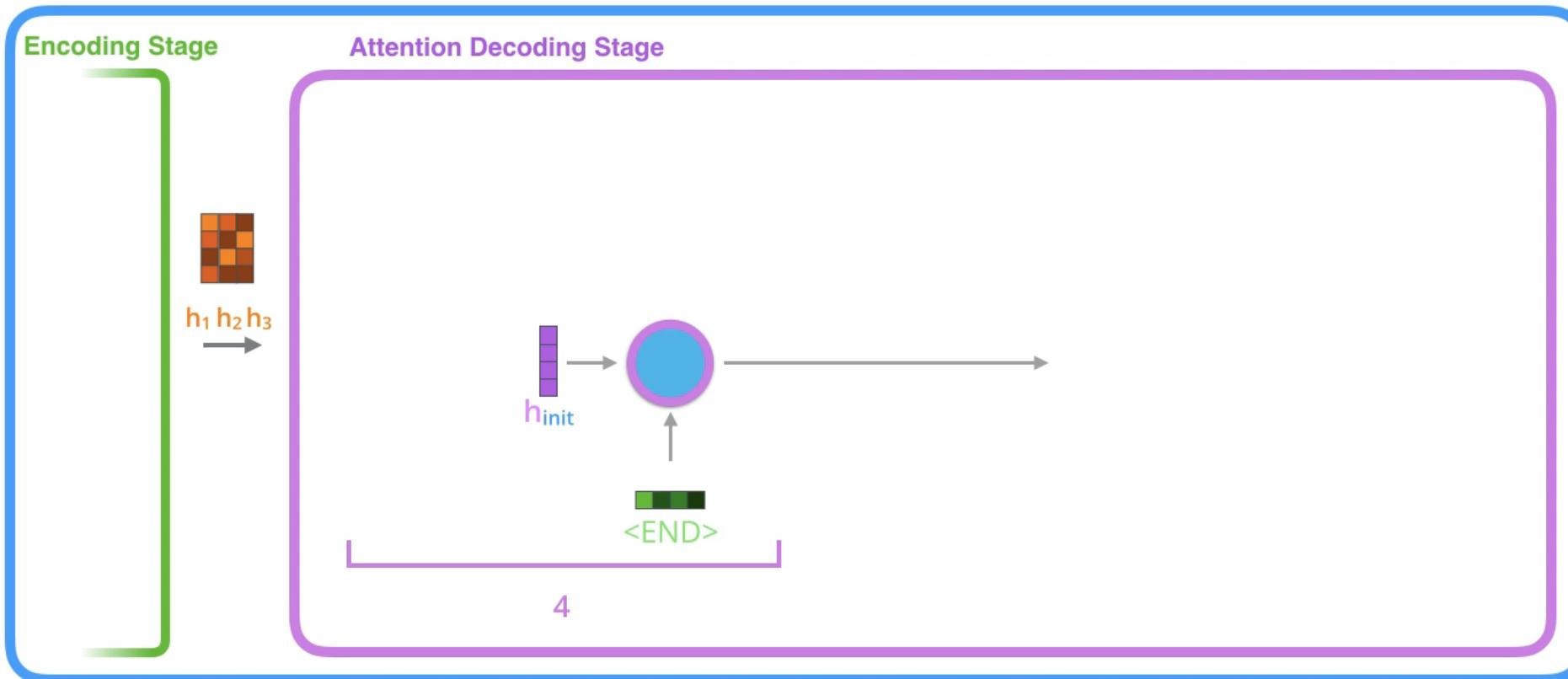


Jay Alammar

Visualizing machine learning one concept at a time.
@JayAlammar on Twitter. YouTube Channel

[Blog](#) [About](#)

Neural Machine Translation SEQUENCE TO SEQUENCE MODEL WITH ATTENTION



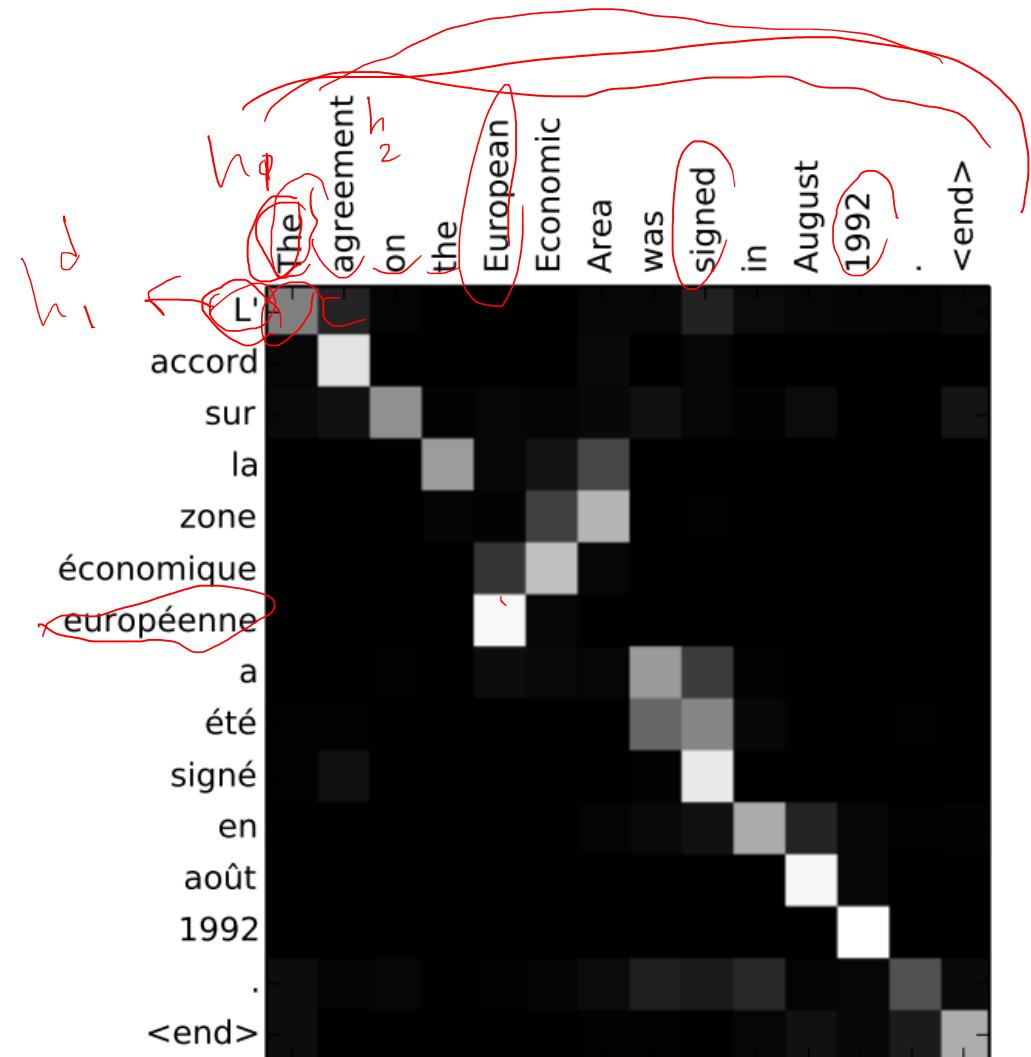


Jay Alammar

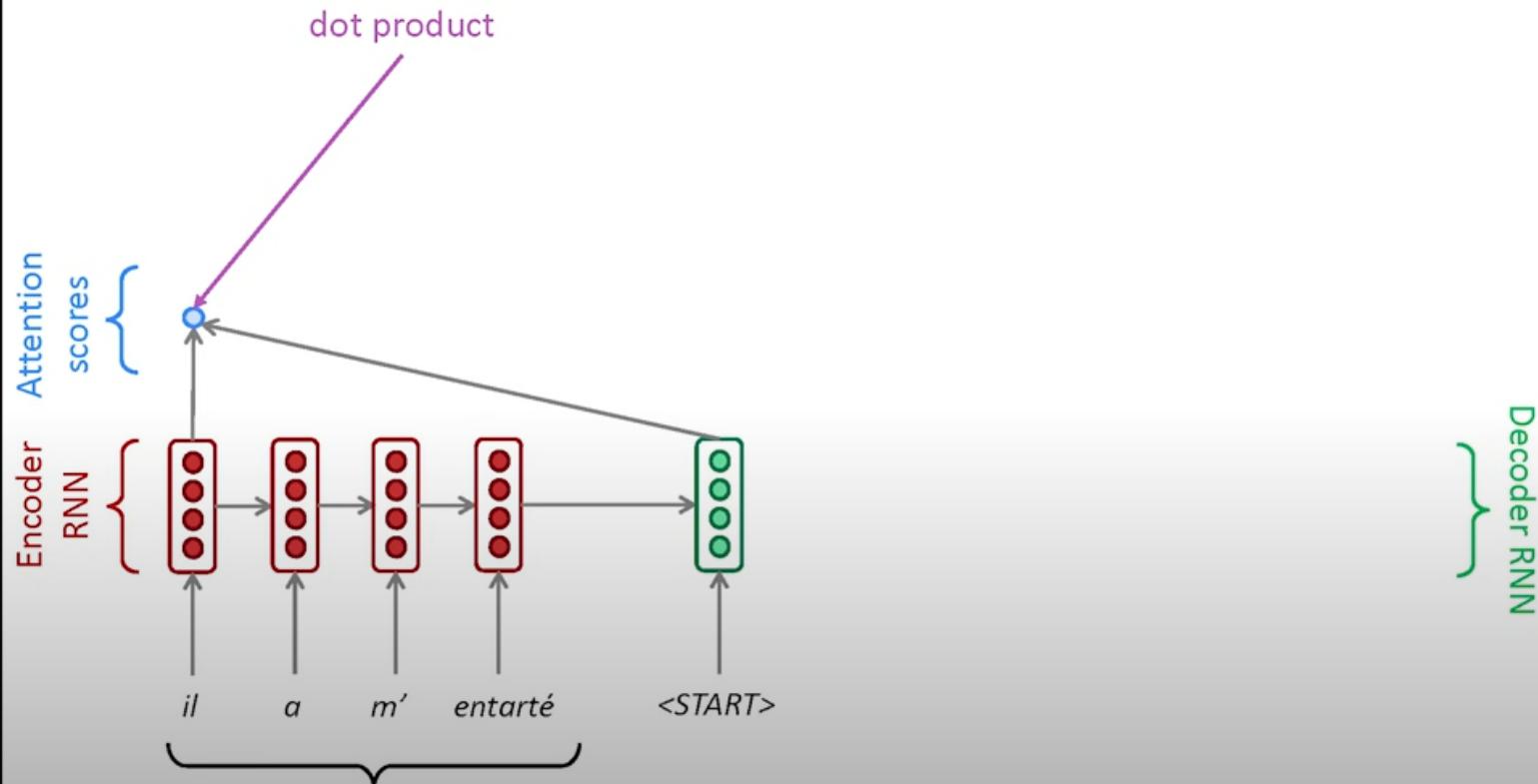
Visualizing machine learning one concept at a time.
@JayAlammar on Twitter. YouTube Channel

Blog About

RNN



Sequence-to-sequence with attention



61

1:03:04 / 1:16:56 • So is Machine Translation solved? >



Stanford CS224N: NLP with Deep Learning | Winter 2019 | Lecture 8 – Translation, Seq2Seq, Attention

100,227 views • May 3, 2019

1.3K DISLIKE SHARE CLIP SAVE ...



Stanford Online ✓
346K subscribers

SUBSCRIBE

Recurrent + Attention!

Published as a conference paper at ICLR 2015

$$\text{Seq2Seq} \Rightarrow \text{MT}$$

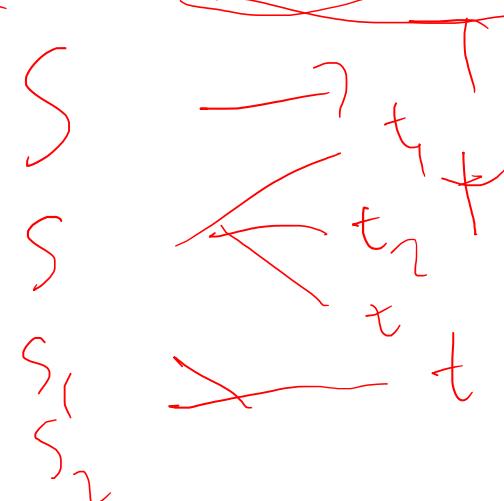
NEURAL MACHINE TRANSLATION BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry Bahdanau

Jacobs University Bremen, Germany

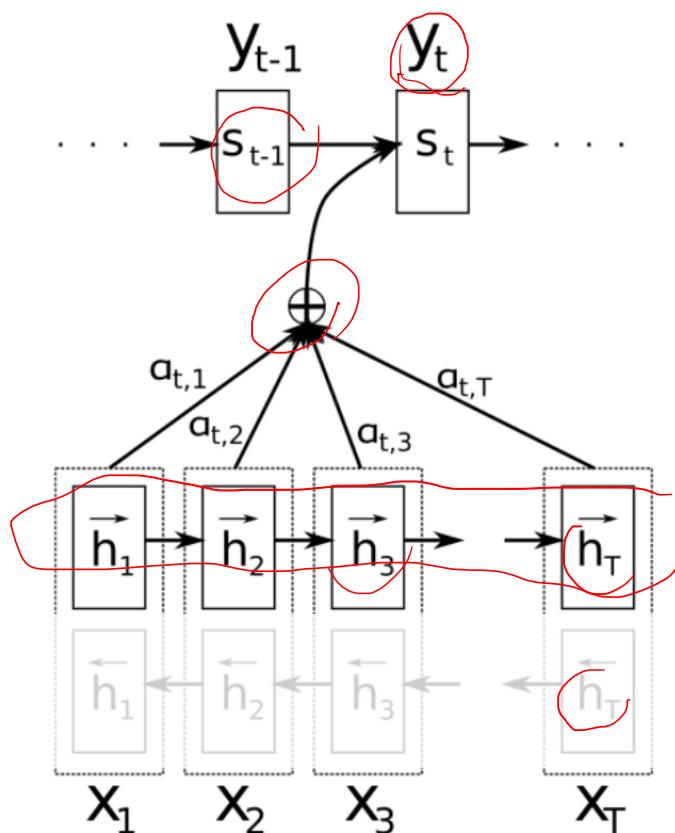
KyungHyun Cho Yoshua Bengio*

Université de Montréal



Recurrent + Attention!

Bahdanau, D., Cho, K., & Bengio, Y. (2015). Neural Machine Translation by Jointly Learning to Align and Translate. In ICLR 2015



Similarity → Cosine → Learn Similarity as a functions of inputs!!

$$\text{score}(h_{d_{t-1}}^d; h_j^e) = h_{d_{t-1}}^d W_s h_j^e$$

$$a_{ij} = \text{softmax}(\text{score}(h_{d_{t-1}}^d, h_j^e)) \text{ for all } j \text{ in encoder}$$

$$c_i = \sum_{j=1}^n a_{ij} h_j^e$$

The point in future make a query to find out how much attention should put to what part of past (values)!

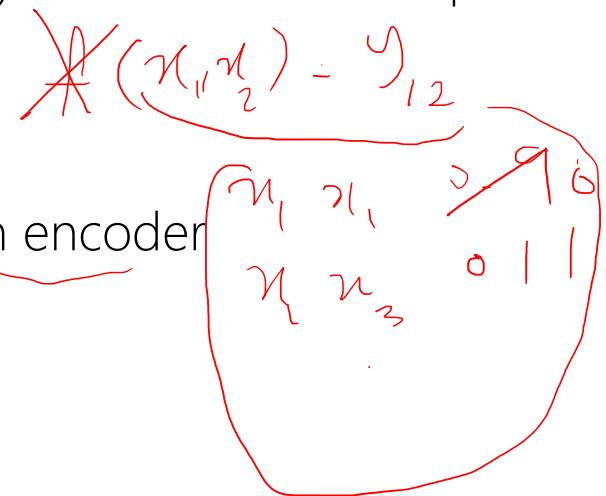


Figure 11: Attention (Bahdanau et al., 2015)

Recurrent + Attention!



A woman is throwing a frisbee in a park.

Figure 12: Visual attention in an image captioning model indicating what the model is attending to when generating the word "frisbee". (Xu et al., 2015)

~~Recurrent~~ + Attention!

Attention Is All You Need

Encoder

WMT

Sys Seq

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

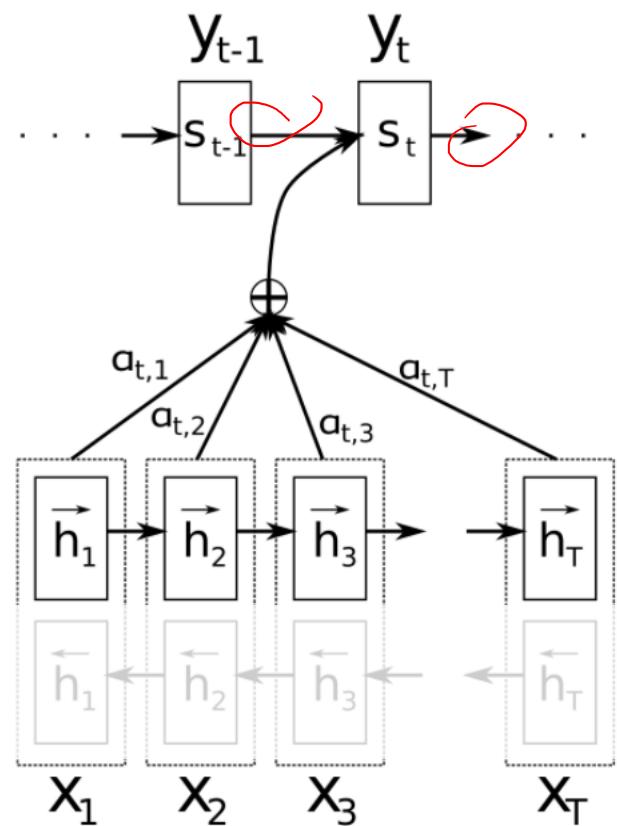
Lukasz Kaiser*
Google Brain
lukasz.kaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

Recurrent + Attention!

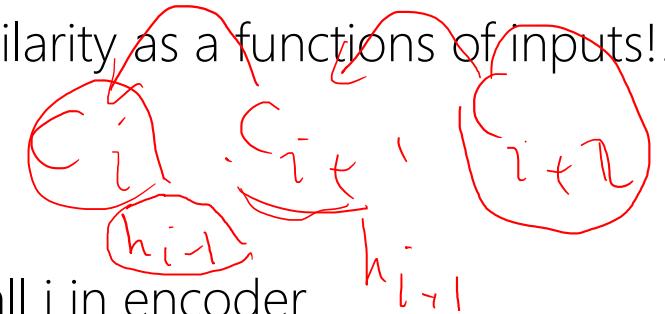


Similarity → Cosine → Learn Similarity as a function of inputs!!

$$\text{score}(h_{i-1}^d; h_j^e) = h_{i-1}^d W_s h_j^e$$

$$a_{ij} = \text{softmax}(\text{score}(h_{i-1}^d; h_j^e)) \text{ for all } j \text{ in encoder}$$

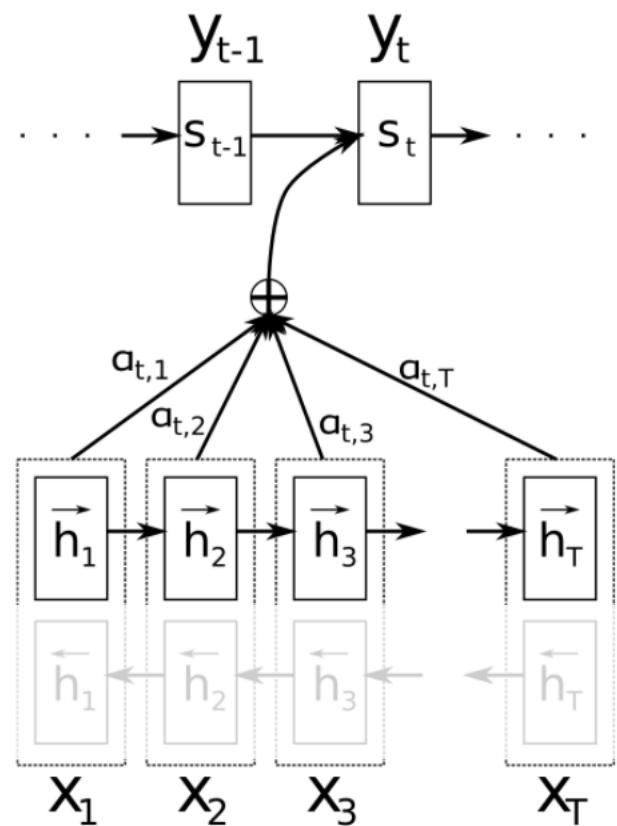
$$c_i = \sum_{j=1}^n a_{ij} h_j^e$$



This is not parallelizable for all $[h_{i-1}^d]$
Why?

Figure 11: Attention (Bahdanau et al., 2015)

Recurrent + Attention!



Similarity → Cosine → Learn Similarity as a functions of inputs!!

$$\text{score}(h_{i-1}^d; h_j^e) = h_{i-1}^d W_s h_j^e$$

$a_{ij} = \text{softmax}(\text{score}(h_{i-1}^d; h_j^e))$ for all j in encoder

$$C_i = \sum_{j=1}^n a_{ij} h_j^e$$

This is not parallelizable for all $[h_{i-1}^d]$

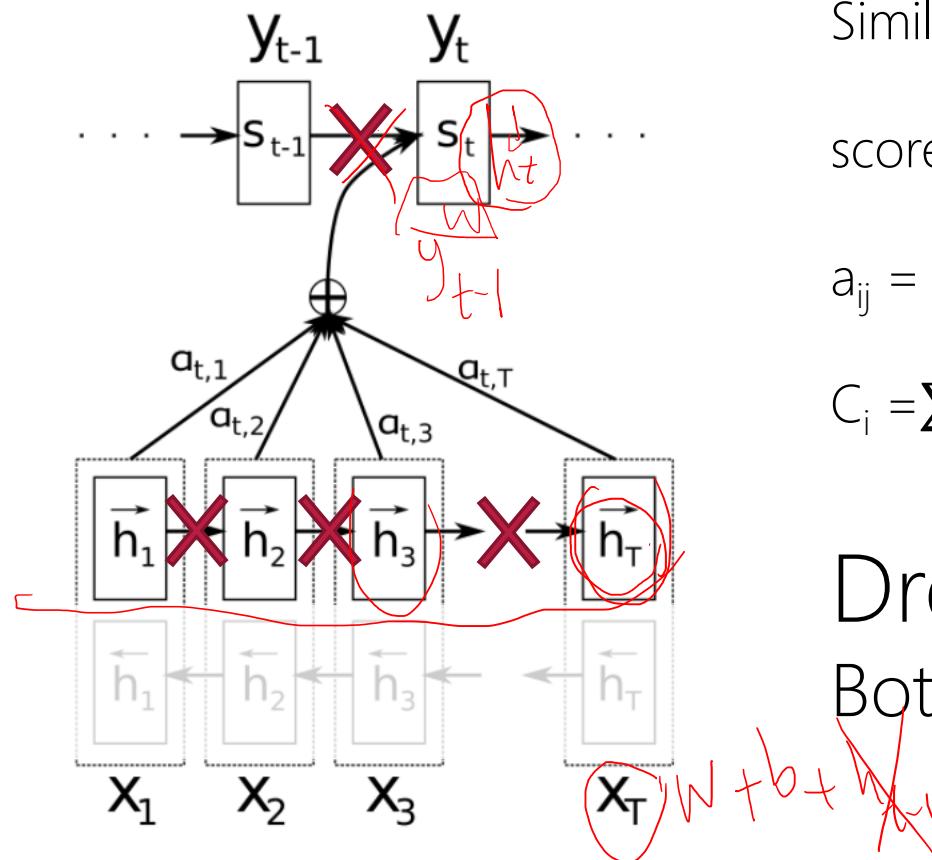
For $[h_{i-1}^d]$, we need to know $[h_{i-1}^d]$

For $[h_{i-1}^d]$, we need to know $[h_{i-2}^d]$

...

Figure 11: Attention (Bahdanau et al., 2015)

Recurrent + Attention!



Similarity → Cosine → Learn Similarity as a functions of inputs!!

$$\text{score}(h_{i-1}^d; h_j^e) = h_{i-1}^d W_s h_j^e$$

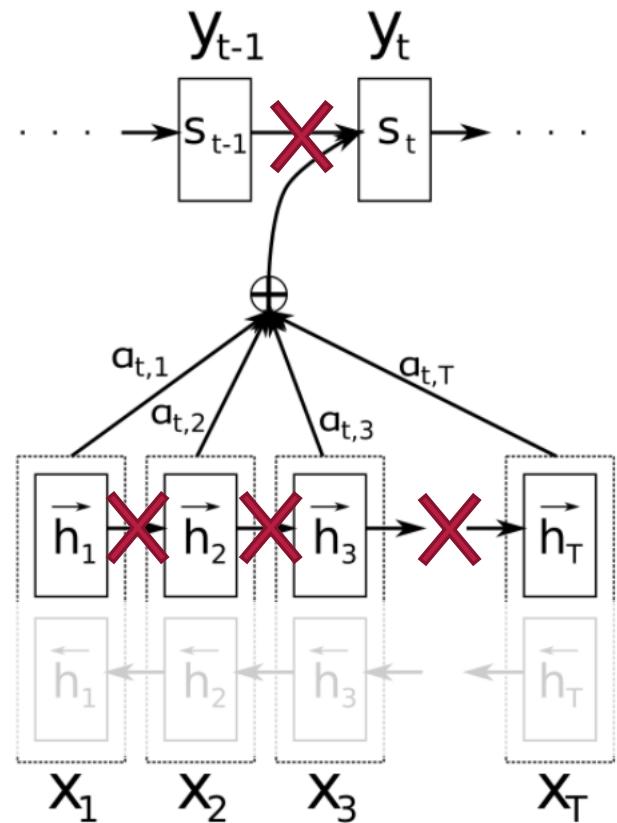
$a_{ij} = \text{softmax}(\text{score}(h_{i-1}^d; h_j^e))$ for all j in encoder

$$C_i = \sum_{j=1}^n a_{ij} h_j^e$$

Drop the recurrence!
Both at encode and decoder.

Figure 11: Attention (Bahdanau et al., 2015)

Recurrent + Attention!



Similarity → Cosine

$$\text{score}(h_i^d; h_j^e) = \underline{h_i^d} \cdot \underline{h_j^e}$$

$a_{ij} = \text{softmax}(\text{score}(h_i^d; h_j^e))$ for all j in encoder

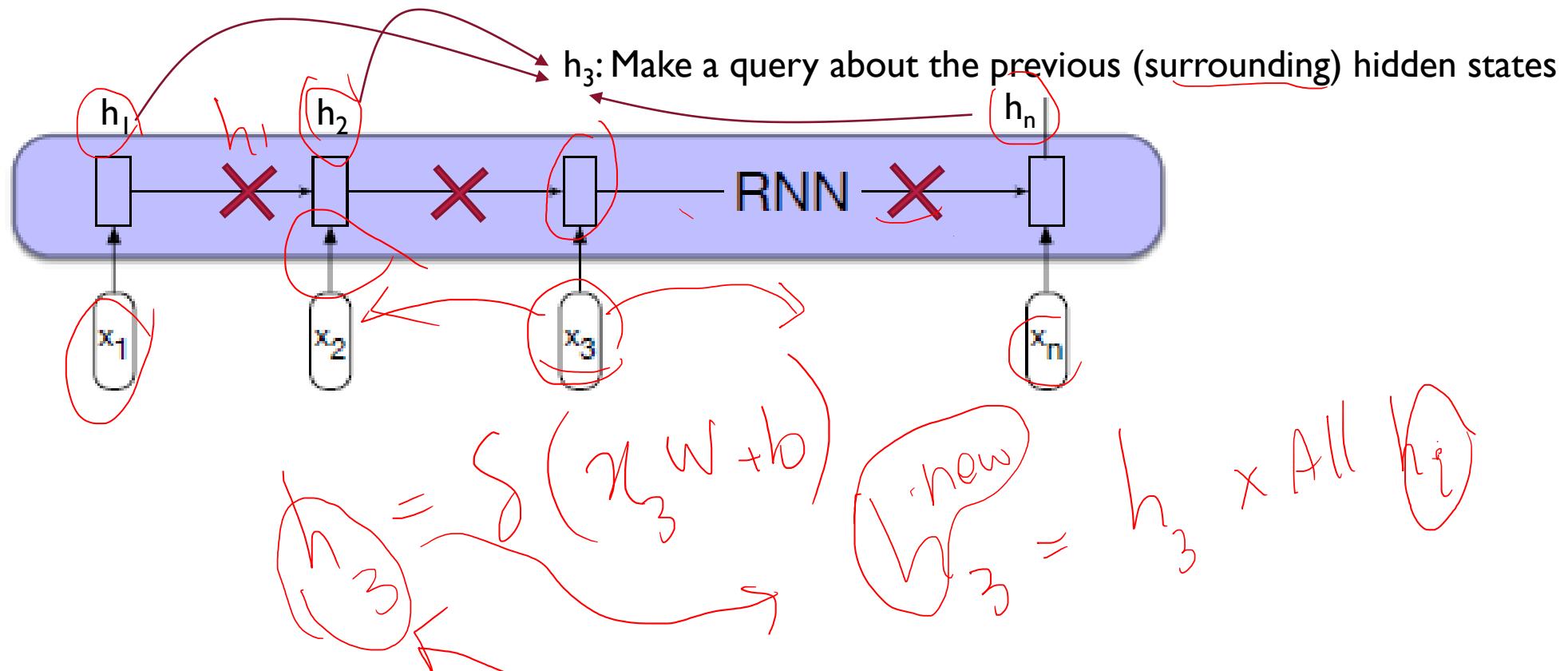
$$C_i = \sum_{j=1}^n a_{ij} h_j^e$$

Back to Cosine!

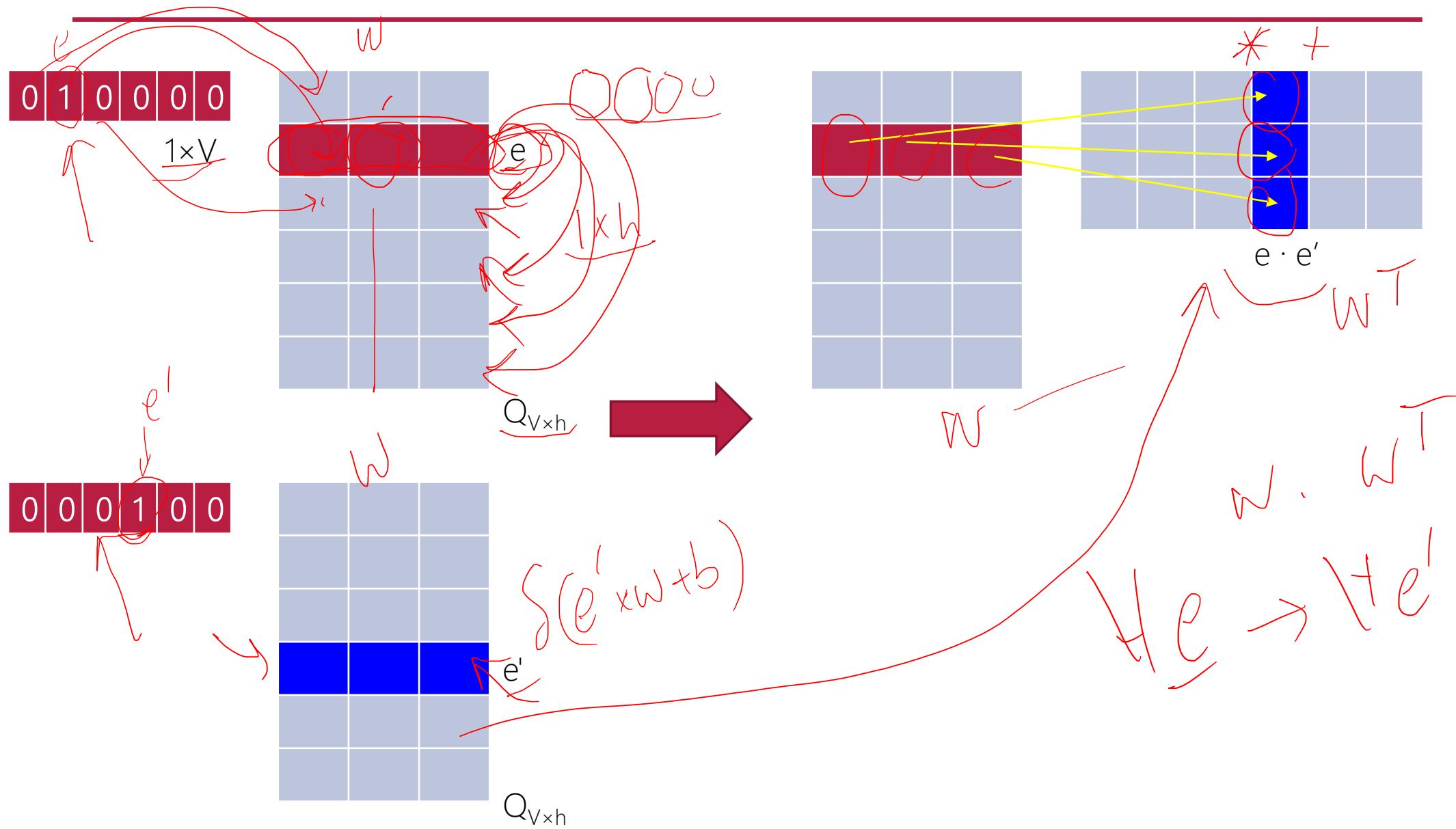
Figure 11: Attention (Bahdanau et al., 2015)

Self-Attention

Back to Recurrent LM



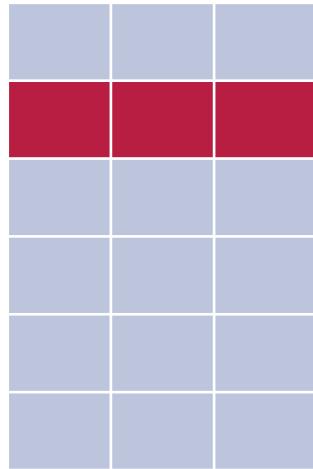
Self-Attention



Self-Attention

0|1|0|0|0|0

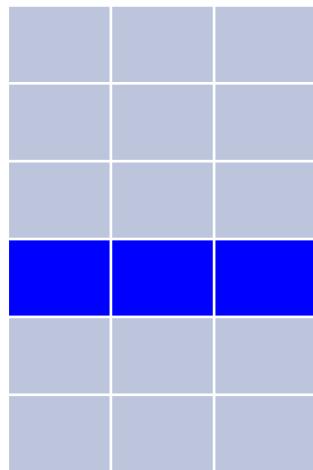
$1 \times V$



e

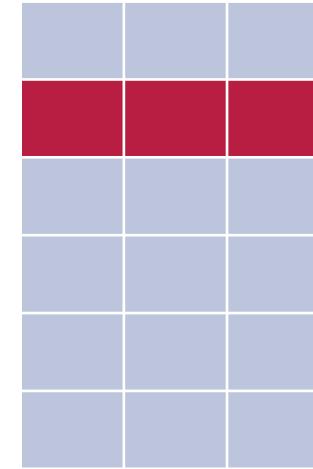
$Q_{V \times h}$

0|0|0|1|0|0



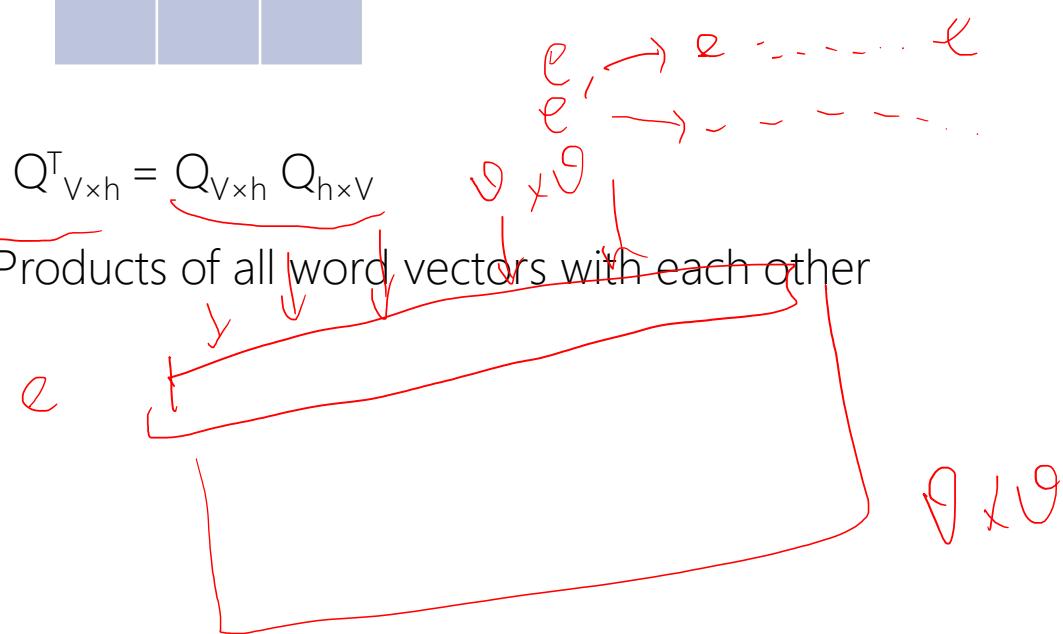
e'

$Q_{V \times h}$



$$Q_{V \times h} Q_{V \times h}^T = Q_{V \times h} Q_{h \times V}$$

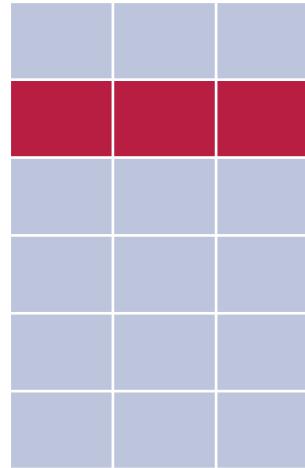
Dot Products of all word vectors with each other



Self-Attention

0|1|0|0|0|0

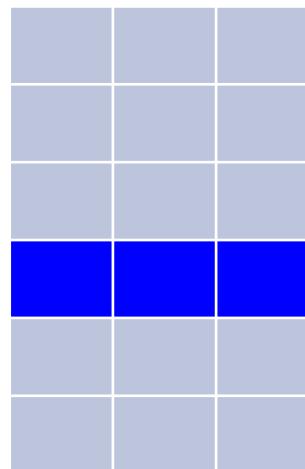
$1 \times V$



e

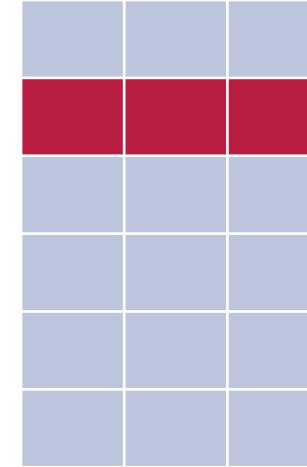
$Q_{V \times h}$ 

0|0|0|1|0|0



e'

$Q_{V \times h}$



$e \cdot e^0$

$e \cdot e^1$

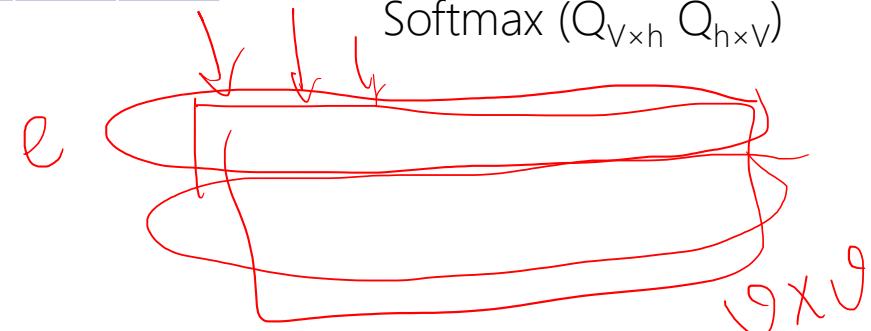
$e \cdot e^2$

$e \cdot e^3$

$e \cdot e^4$

$e \cdot e^5$

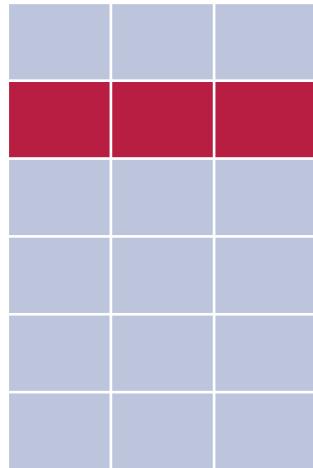
Softmax ($Q_{V \times h} Q_{h \times V}$)



Self-Attention

$0|1|0|0|0|0|0$

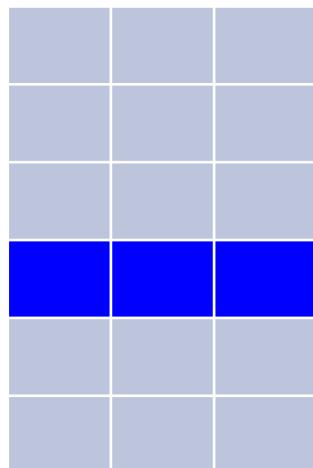
$1 \times V$



e

$Q_{V \times h}$

$0|0|0|1|0|0|0$

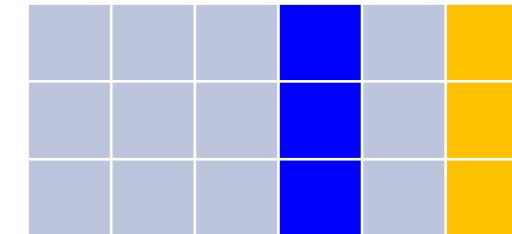
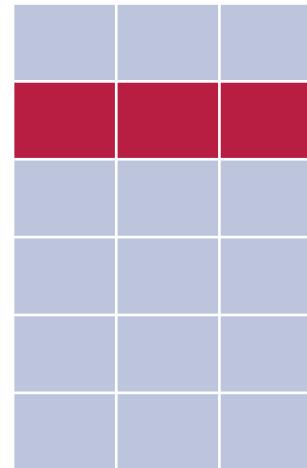


e'

$Q_{V \times h}$



$$(s_{e0} \times e_0) + (s_{e1} \times e_1) + \dots$$

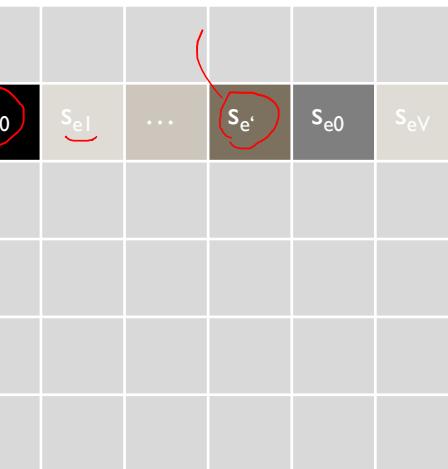


$e \cdot e_0$

$e \cdot e_1$

$e \cdot e'$

$e \cdot e_V$



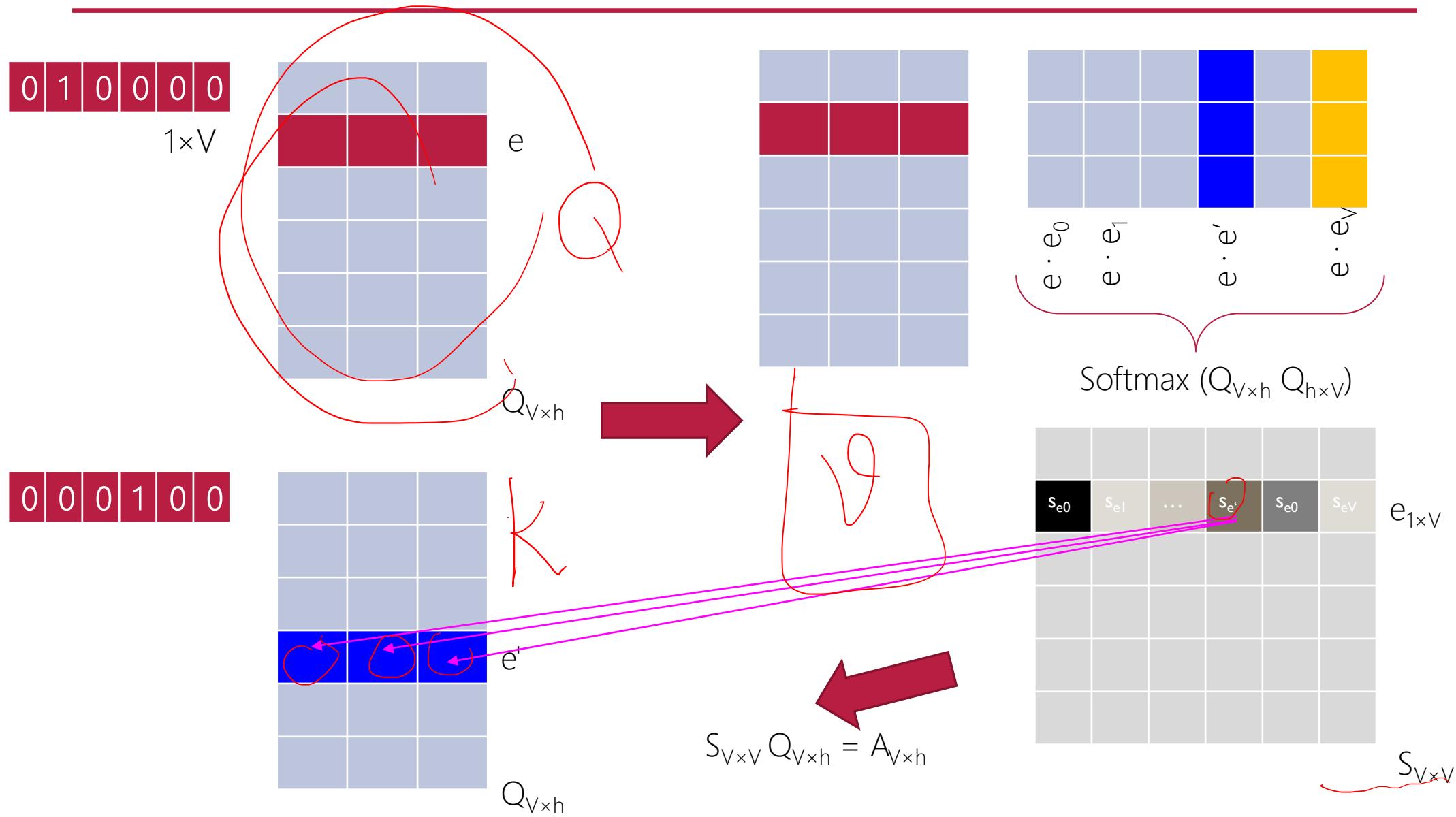
Softmax ($Q_{V \times h}, Q_{h \times V}$)

$S_{V \times V}$

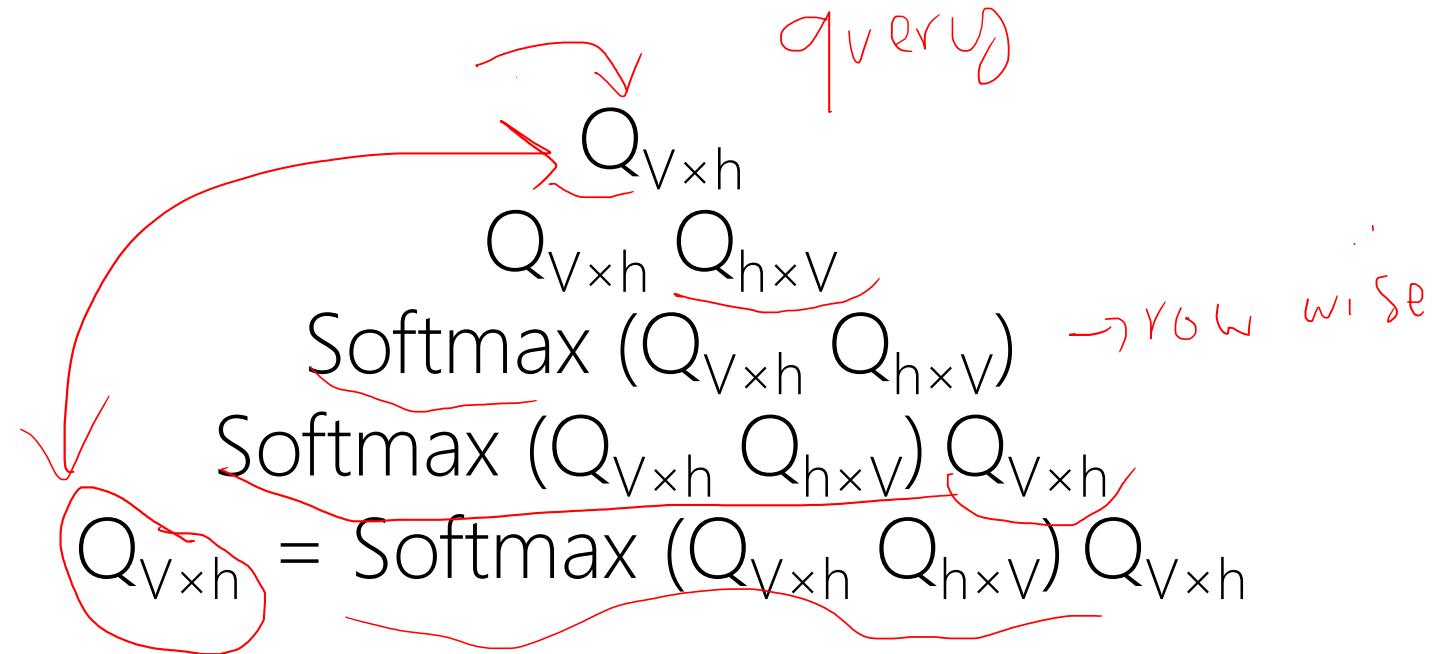
$e_{1 \times V}$

$$S = \sum_{i=0}^{V-1} s_{ei} (W^T)^i W$$

Self-Attention

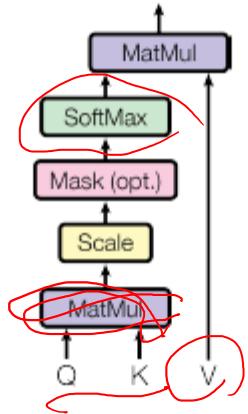


Self-Attention! Renew Representation



Self-Attention! Renew Representation

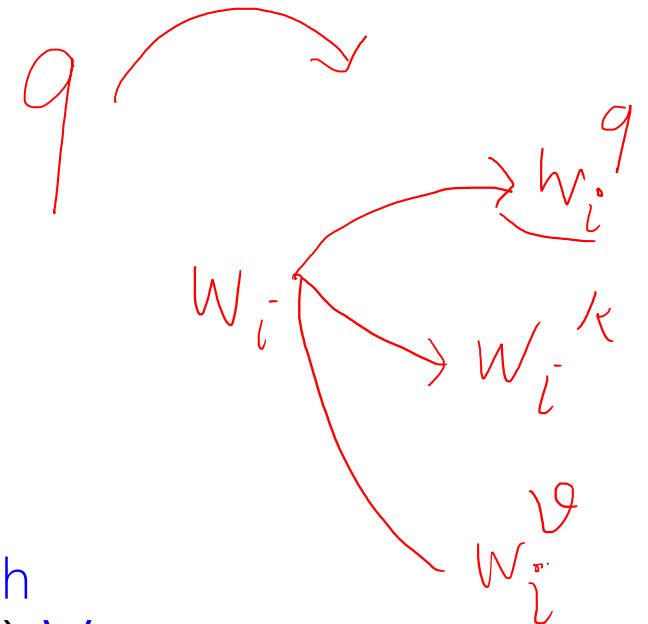
Scaled Dot-Product Attention



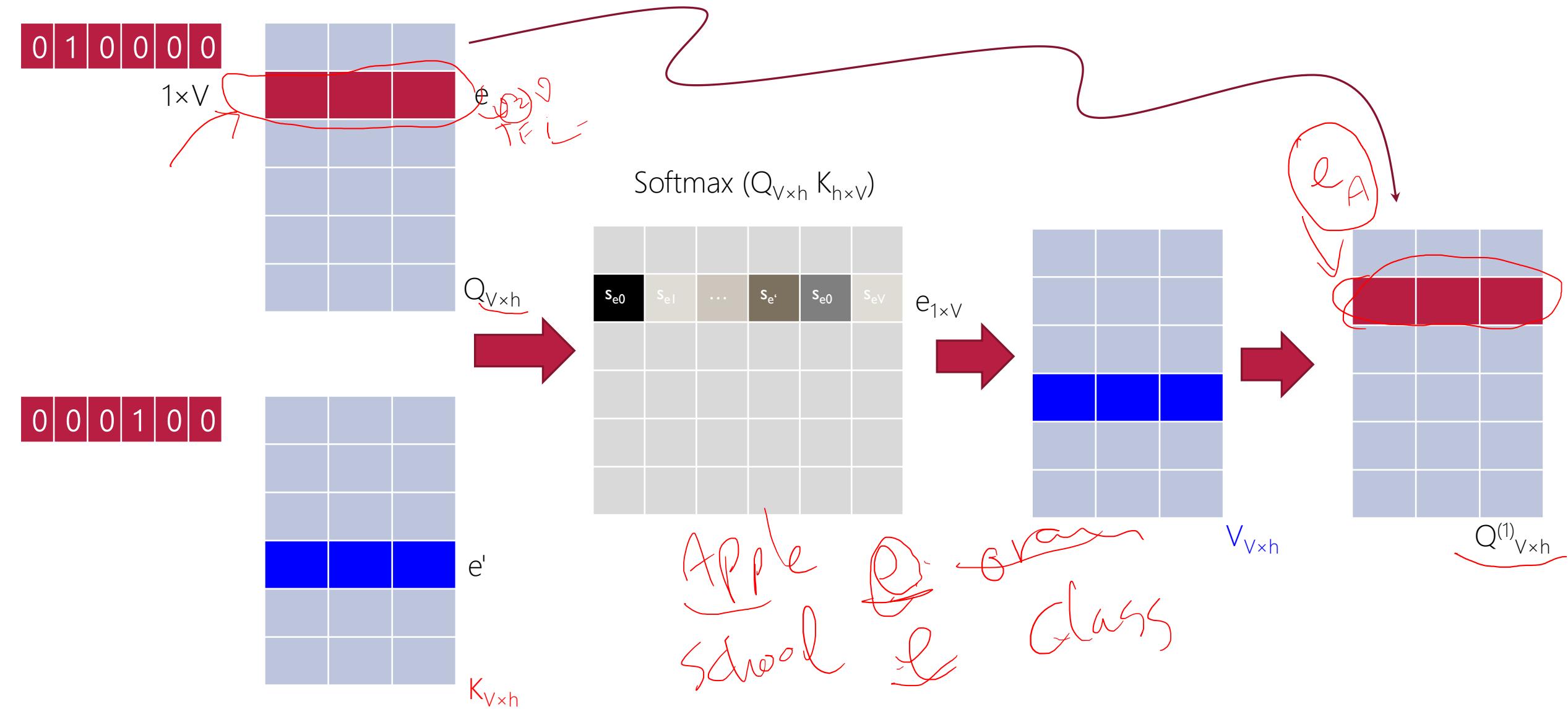
$$Q^{(1)}_{V \times h} = \text{Softmax} \left(Q_{V \times h} K_{h \times V} \right) V_{V \times h}$$

Annotations in red:

- $Q_{V \times h}$ is circled.
- $K_{h \times V}$ is circled.
- $Q_{V \times h} K_{h \times V}$ is circled.
- $Q_{V \times h} K_{h \times V}$ is circled.
- $V_{V \times h}$ is circled.
- key is written next to $K_{h \times V}$.



Self-Attention



The Transformer

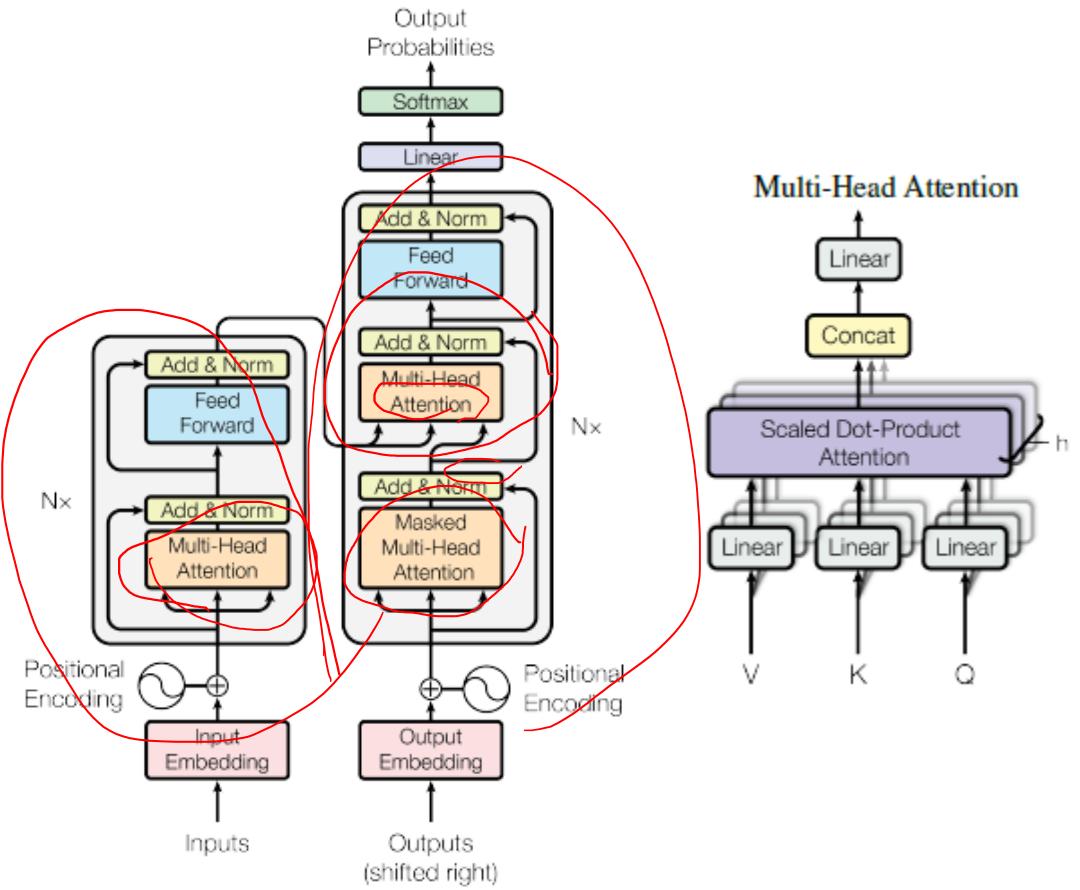


Figure 1: The Transformer - model architecture.

Contextualized Word Embedding

Positional Embedding

BERT: Bidirectional Encoder Representations from Transformers
