

344.063/163 KV Special Topic:

Natural Language Processing with Deep Learning

Neural Machine Translation with Attention Networks



Navid Rekab-Saz

navid.rekabsaz@jku.at

Institute of Computational Perception

Agenda

- Machine Translation
- Attention Networks
- Attention in practice
 - Seq2seq with attention
 - Hierarchical document classification

Agenda

- **Machine Translation**
- Attention Networks
- Attention in practice
 - Seq2seq with attention
 - Hierarchical document classification

Machine Translation (MT)

- Machine Translation is the task of translating a sentence X from source language to sentence Y in target language
- A long-history (since 1950)
 - Early systems were mostly rule-based
- Challenges:
 - Common sense
 - Idioms!
 - Typological differences between the source and target language
 - Alignment
 - Low-resource language pairs

Statistical Machine Translation (SMT)

- Statistical Machine Translation (1990-2010) learns a **probabilistic model** using large amount of **parallel data**
- The model aims to find the best target language sentence Y^* , given the source language sentence X :

$$Y^* = \operatorname{argmax}_Y P(Y|X)$$

- SMT uses Bayes Rule to split this probability into two components that can be learnt separately:

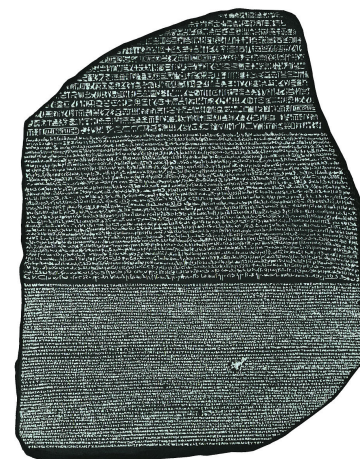
$$= \operatorname{argmax}_Y P(X|Y)P(Y)$$

Translation Model

The statistical model that defines how words and phrases should be translated (learnt from parallel data)

Language Model

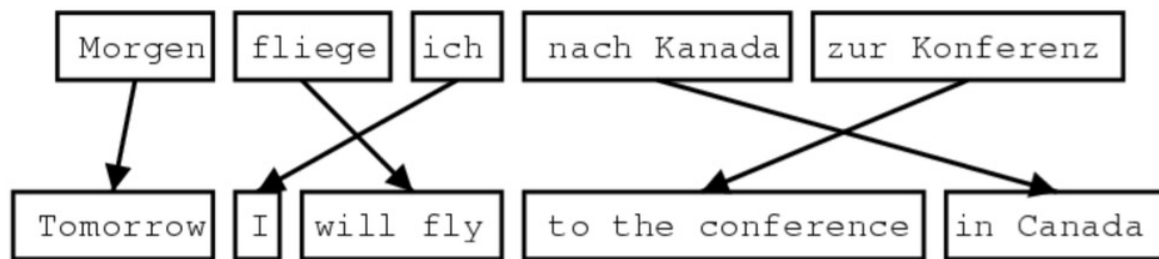
The statistical model that tells us how to write good sentences in the target language (learnt from monolingual data)



https://en.wikipedia.org/wiki/Rosetta_Stone

Learning Translation model

- To learn the Translation model $P(Y|X)$, we need to break X and Y down to **aligned words and phrases**:

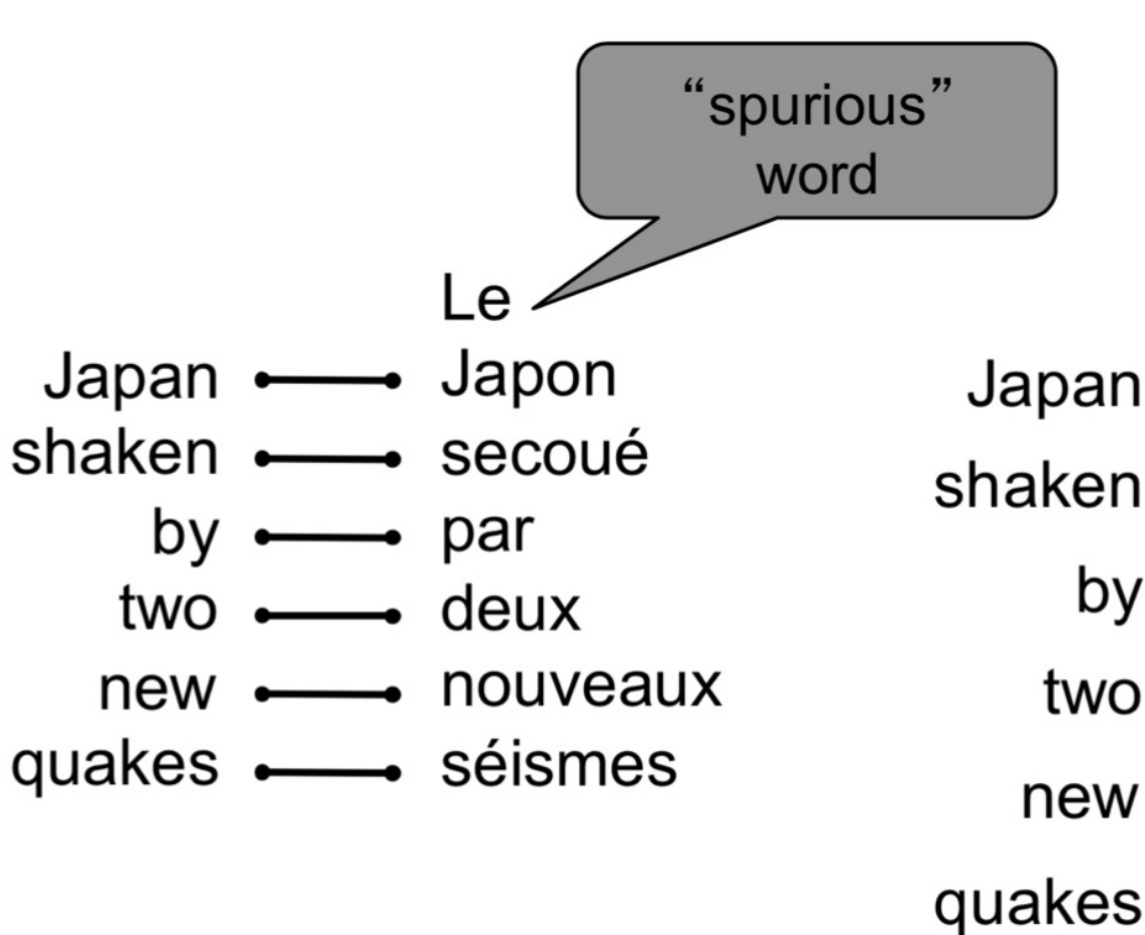


- To this end, the **alignment** latent variable a is added to the formulation of Translation model:

$$P(X, a|Y)$$

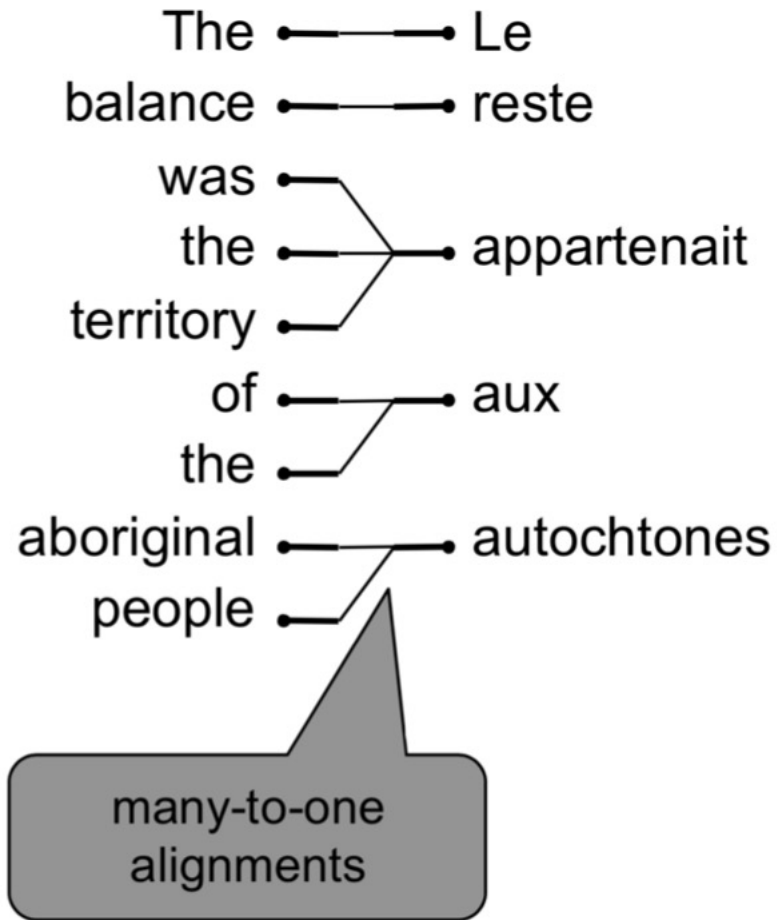
- Alignment ...
 - is a latent variable \rightarrow is not explicitly defined in the data!
 - defines the correspondence between particular words/phrases in the translation sentence pair

Alignment!



	Le	Japon	secoué	par	deux	nouveaux	séismes
Japan							
shaken							
by							
two							
new							
quakes							

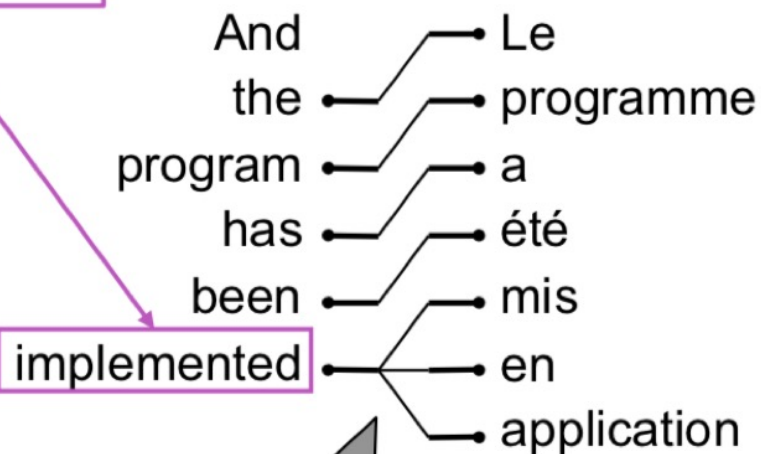
Alignment!



	Le	reste	appartenait	aux	autochtones
The					
balance					
was					
the					
territory					
of					
the					
aboriginal					
people					

Alignment!

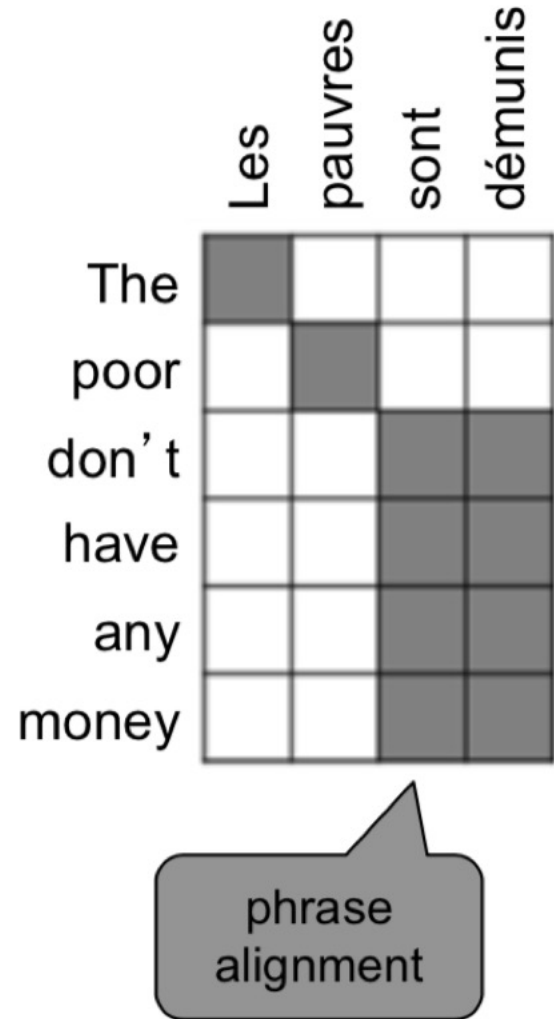
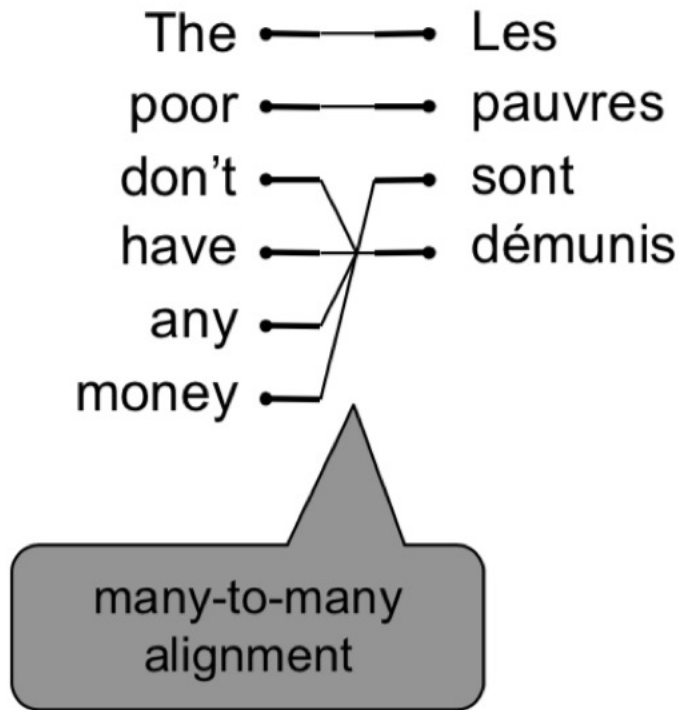
We call this a
fertile word



one-to-many
alignment

	Le	programme	a	été	mis	en	application
And							
the							
program							
has							
been							
implemented							

Alignment!



SMT – summary

- Defining alignment is complex!
 - The Translation model should jointly estimate distributions of both variables (X and a)
- SMT systems ...
 - were extremely complex with lots of features engineering
 - required extra resources like dictionaries and mapping tables between phrases and words
 - required “special attention” for each language pair and lots of human efforts

MT – Evaluation

- BLEU (Bilingual Evaluation Understudy)
- BLEU computes a **similarity score** between the machine-written translation to one or several human-written translation(s), based on:
 - ***n*-gram precision** (usually for 1, 2, 3 and 4-grams)
 - plus a penalty for too-short machine translations
- BLEU is precision-based, while ROUGE is recall-based

Details of how to calculate BLEU: <https://www.coursera.org/lecture/nlp-sequence-models/bleu-score-optional-kC2HD>

Agenda

- Machine Translation
- **Attention Networks**
- Attention in practice
 - Seq2seq with attention
 - Hierarchical document classification

Attention Networks

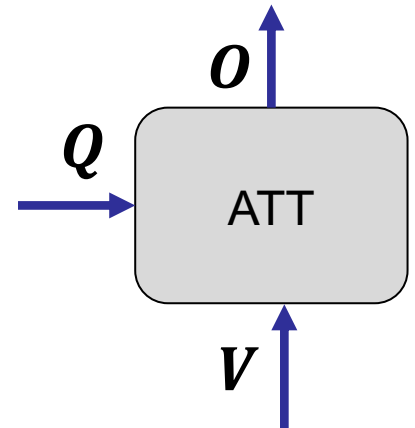
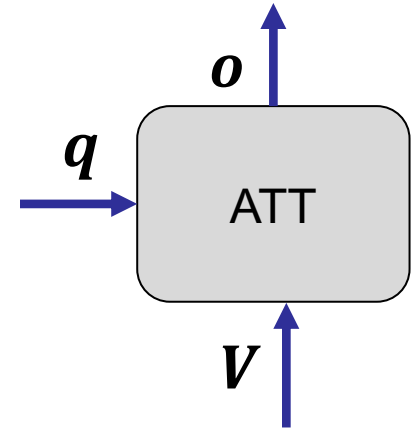
- Attention is a **generic** deep learning method ...
 - to obtain a **composed** representation (output \mathbf{o}) ...
 - from an **arbitrary size** of input representations (values \mathbf{V}) ...
 - based on another given representation (query \mathbf{q})

- General form of an attention network:

$$\mathbf{o} = \text{ATT}(\mathbf{q}, \mathbf{V})$$

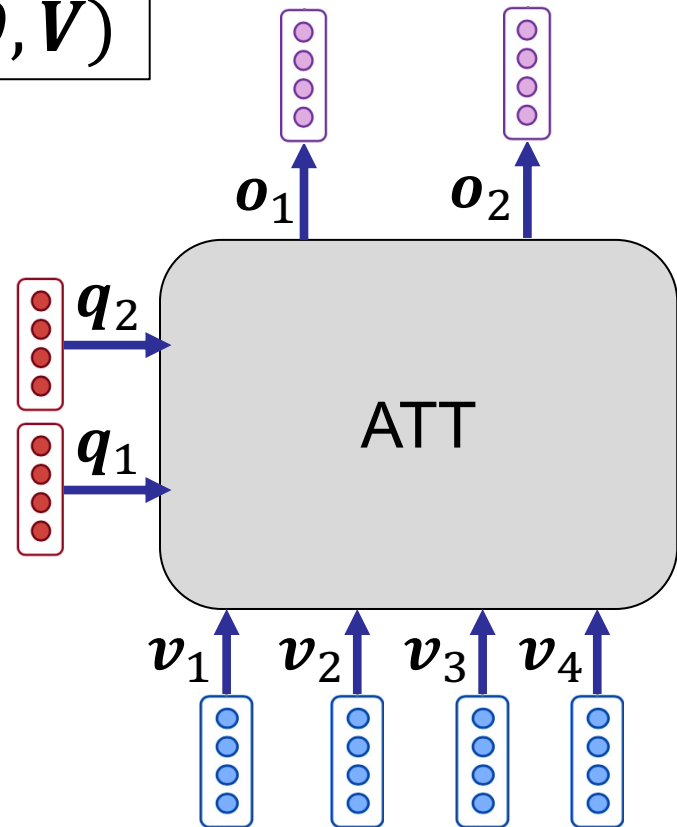
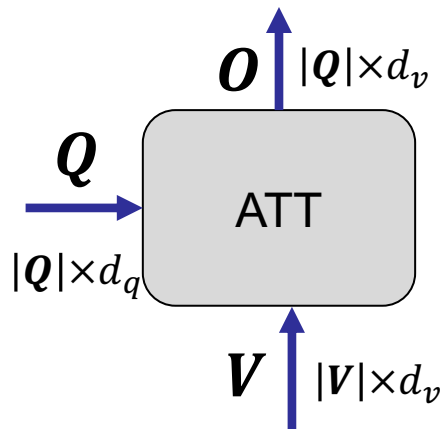
- when we have a set of queries, it will be:

$$\mathbf{O} = \text{ATT}(\mathbf{Q}, \mathbf{V})$$



Attention Networks

$$\mathbf{O} = \text{ATT}(\mathbf{Q}, \mathbf{V})$$



- d_q, d_v are embedding dimensions of query and value vectors, respectively

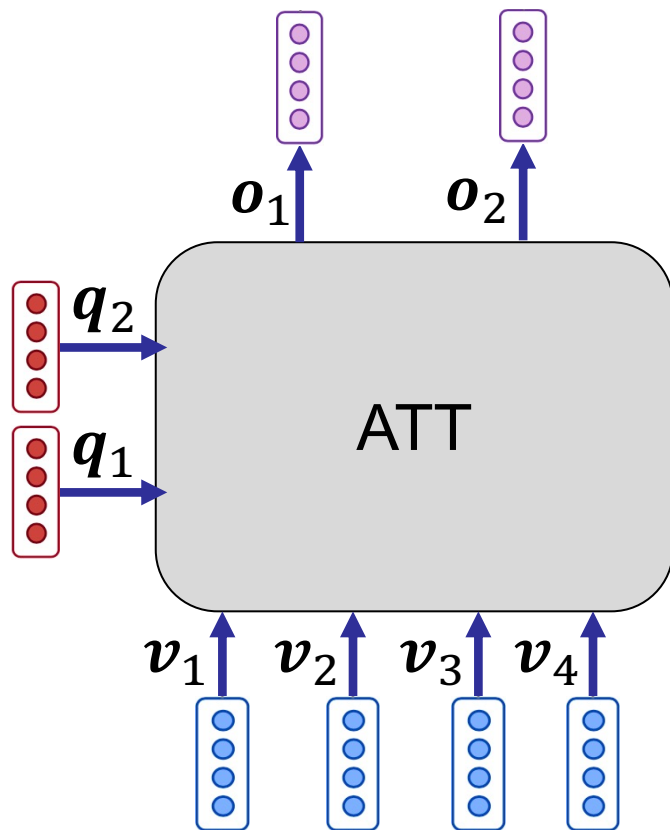
We sometime say, each **query vector** q “attends” to **value vectors**

Attention Networks – definition

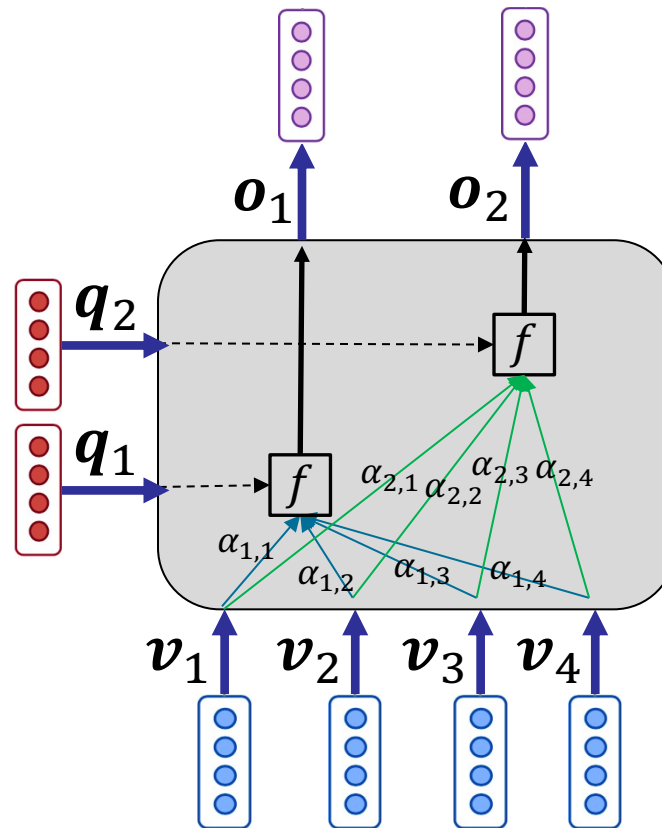
Formal definition:

- Given a set of vectors of *values* V and *queries* Q ...
 - for each query in Q , **attention** computes a *weighted sum* of the values V as output of that query
- To this end, each query vector puts some amount of *attention* on each value vector
 - This attention is used as the weight in the weighted sum
- The weighted sum in attention networks can be seen as a *selective summary* of the information of the value vectors, where the query defines what portion (attention weight) of each value should be taken

Attentions!



Attentions!



$\alpha_{i,j}$ is the attention of query q_i on value v_j

α_i is the vector of attentions of query q_i on value vectors V

α_i is a probability distribution

f is the attention function

Attention Networks – formulation

- Given the query vector \mathbf{q}_i , an attention network assigns **attention** $\alpha_{i,j}$ to each value vector \mathbf{v}_j using **attention function** f :

$$\alpha_{i,j} = f(\mathbf{q}_i, \mathbf{v}_j)$$

where α_i forms a **probability distribution** over vector values:

$$\sum_{j=1}^{|V|} \alpha_{i,j} = 1$$

- The output regarding each query is the **weighted sum** of the value vectors using attentions as weights:

$$\mathbf{o}_i = \sum_{j=1}^{|V|} \alpha_{i,j} \mathbf{v}_j$$

Attention – first formulation

Basic dot-product attention

- First, non-normalized attention scores:

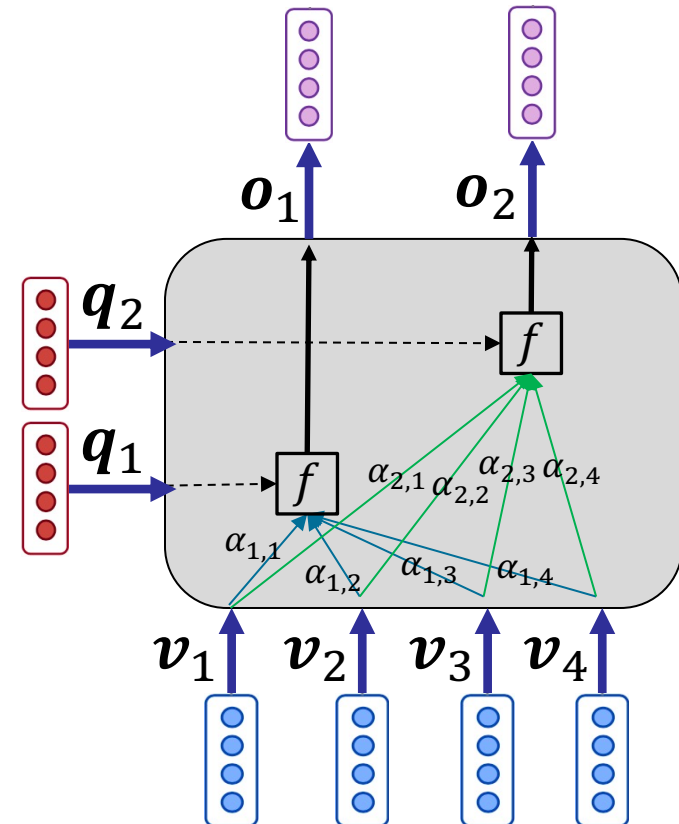
$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{v}_j^T$$

- In this variant $d_q = d_v$
- There is no parameter to learn!

- Then, softmax over values:

$$\alpha_i = \text{softmax}(\tilde{\alpha}_i)$$

- Output (weighted sum): $\mathbf{o}_i = \sum_{j=1}^{|\mathcal{V}|} \alpha_{i,j} \mathbf{v}_j$

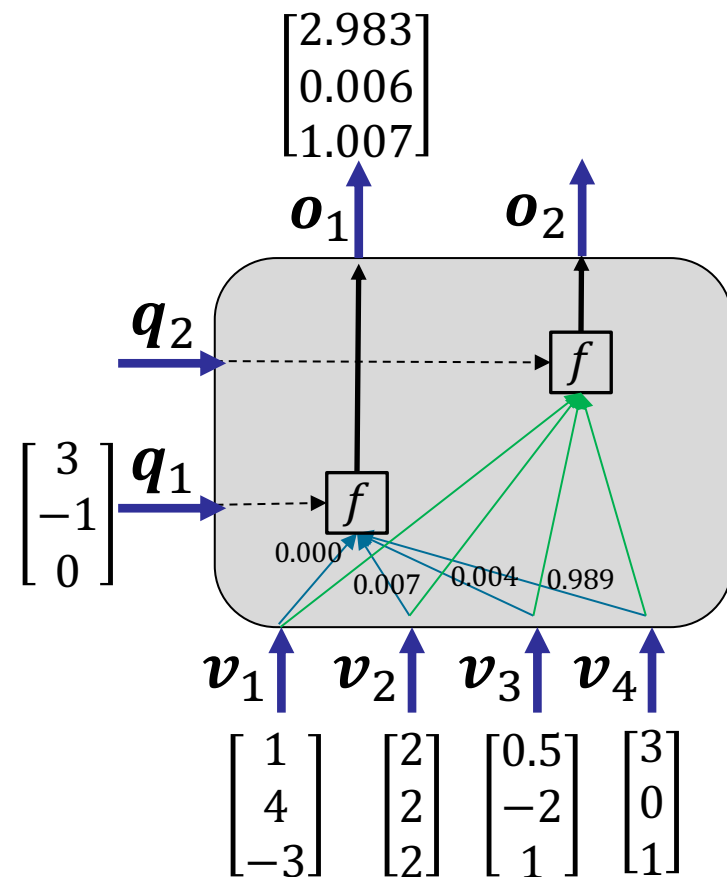


Example

$$\tilde{\alpha}_1 = \begin{bmatrix} \mathbf{q}_1 \mathbf{v}_1^T = -1 \\ \mathbf{q}_1 \mathbf{v}_2^T = 4 \\ \mathbf{q}_1 \mathbf{v}_3^T = 3.5 \\ \mathbf{q}_1 \mathbf{v}_4^T = 9 \end{bmatrix} \rightarrow \alpha_1 = \begin{bmatrix} 0.000 \\ 0.007 \\ 0.004 \\ 0.989 \end{bmatrix}$$

$$\mathbf{o}_1 = 0.000 \begin{bmatrix} 1 \\ 4 \\ -3 \end{bmatrix} + 0.007 \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + 0.004 \begin{bmatrix} 0.5 \\ -2 \\ 1 \end{bmatrix} + 0.989 \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathbf{o}_1 = \begin{bmatrix} 2.983 \\ 0.006 \\ 1.007 \end{bmatrix}$$

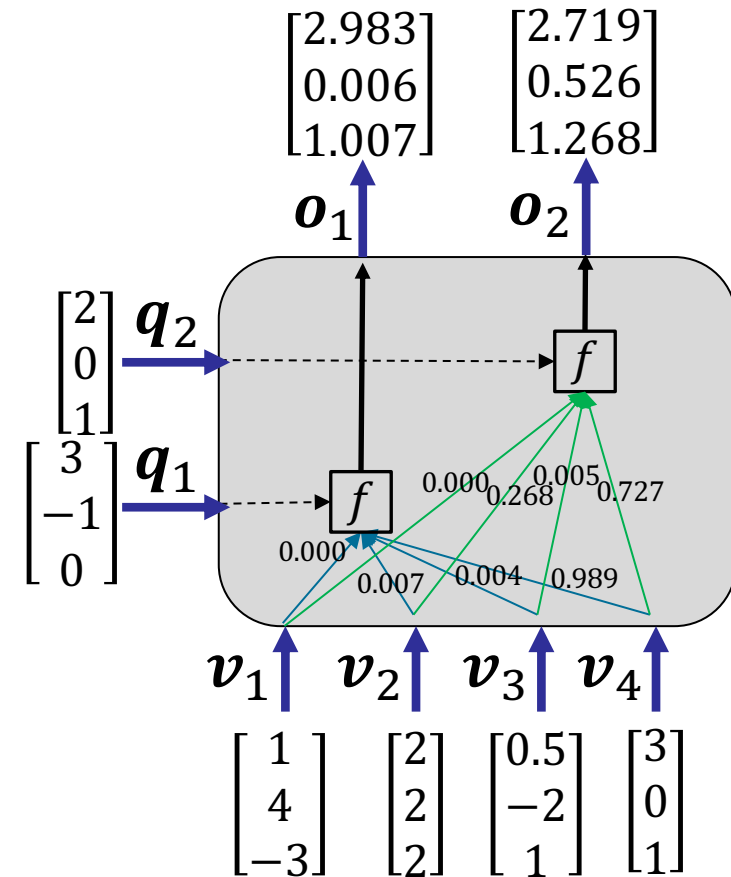


Example

$$\tilde{\alpha}_2 = \begin{bmatrix} \mathbf{q}_2 \mathbf{v}_1^T = -1 \\ \mathbf{q}_2 \mathbf{v}_2^T = 6 \\ \mathbf{q}_2 \mathbf{v}_3^T = 2 \\ \mathbf{q}_2 \mathbf{v}_4^T = 7 \end{bmatrix} \rightarrow \alpha_2 = \begin{bmatrix} 0.000 \\ 0.268 \\ 0.005 \\ 0.727 \end{bmatrix}$$

$$\mathbf{o}_2 = 0.000 \begin{bmatrix} 1 \\ 4 \\ -3 \end{bmatrix} + 0.268 \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} + 0.005 \begin{bmatrix} 0.5 \\ -2 \\ 1 \end{bmatrix} + 0.727 \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathbf{o}_2 = \begin{bmatrix} 2.719 \\ 0.526 \\ 1.268 \end{bmatrix}$$



Attention variants

Basic dot-product attention (recap)

- First, non-normalized attention scores:

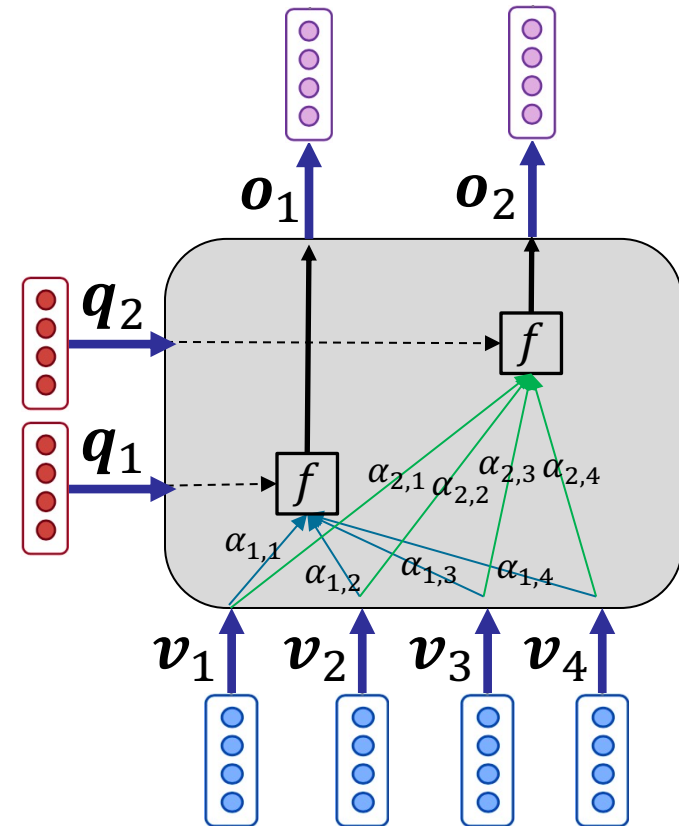
$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{v}_j^T$$

- In this variant $d_q = d_v$
- There is no parameter to learn!

- Then, softmax over values:

$$\alpha_i = \text{softmax}(\tilde{\alpha}_i)$$

- Output (weighted sum): $\mathbf{o}_i = \sum_{j=1}^{|\mathcal{V}|} \alpha_{i,j} \mathbf{v}_j$



Attention variants

Multiplicative attention

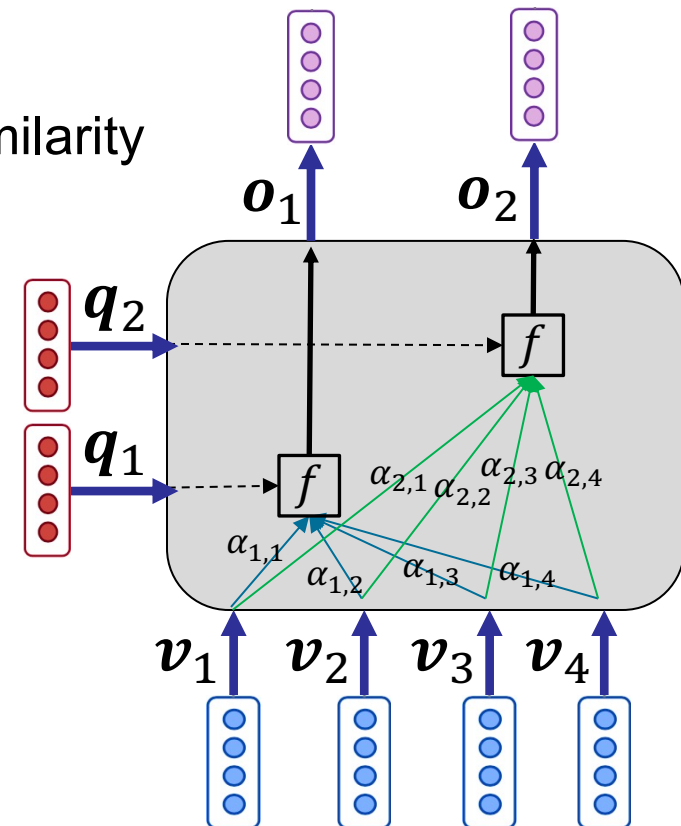
- First, non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{W} \mathbf{v}_j^T$$

- \mathbf{W} is a matrix of model parameters
 - adds a linear function to measure the similarity between query and value
- Then, softmax over values:

$$\alpha_i = \text{softmax}(\tilde{\alpha}_i)$$

- Output (weighted sum): $\mathbf{o}_i = \sum_{j=1}^{|\mathcal{V}|} \alpha_{i,j} \mathbf{v}_j$



Attention variants

Additive attention

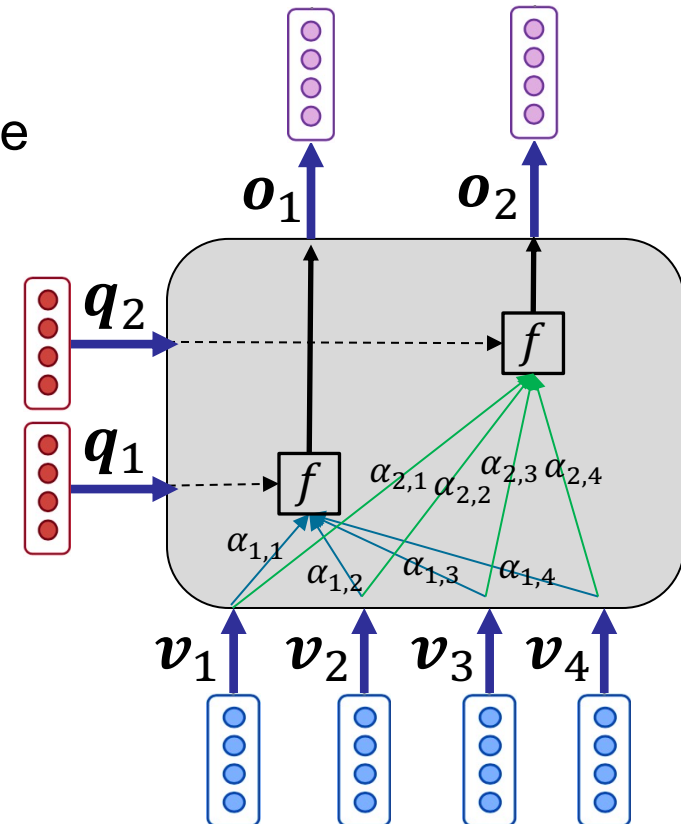
- First, non-normalized attention scores:

$$\tilde{\alpha}_{i,j} = \mathbf{u}^T \tanh(\mathbf{q}_i \mathbf{W}_1 + \mathbf{v}_j \mathbf{W}_2)$$

- \mathbf{W}_1 , \mathbf{W}_2 , and \mathbf{u} are model parameters
 - adds a non-linear function to measure the similarity between query and value
- Then, softmax over values:

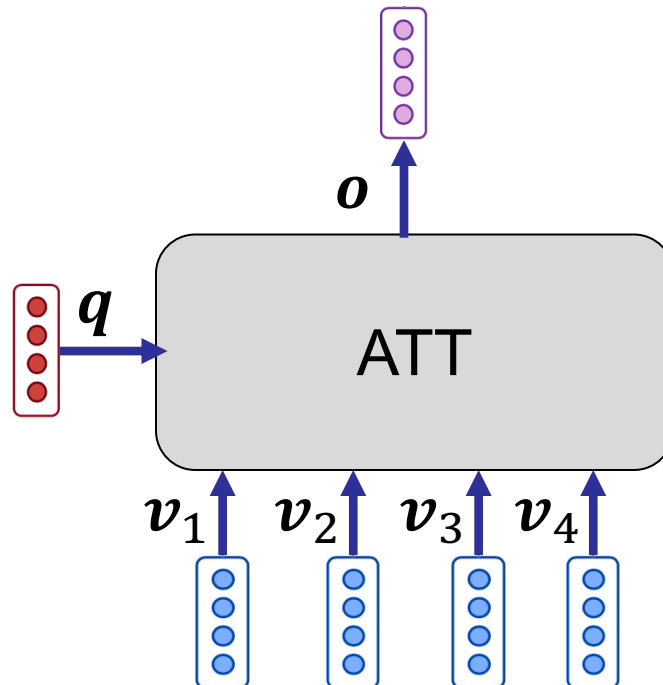
$$\alpha_i = \text{softmax}(\tilde{\alpha}_i)$$

- Output (weighted sum): $\mathbf{o}_i = \sum_{j=1}^{|\mathcal{V}|} \alpha_{i,j} \mathbf{v}_j$



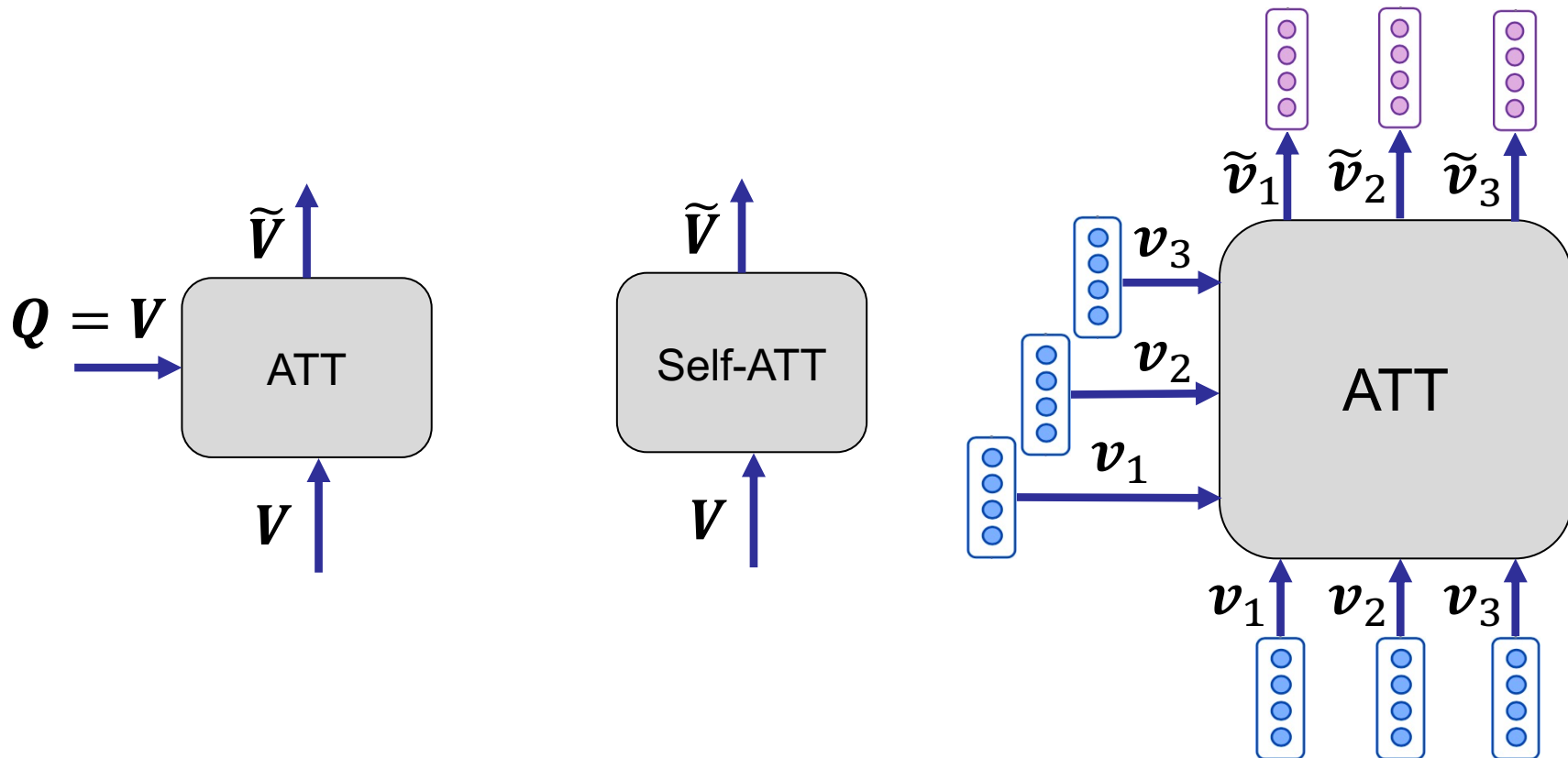
Attention in practice

- Attention is used to create a **compositional embedding** of value vectors according to a query
 - E.g., in seq2seq models or document classification (will be discussed in this lecture)



Self-attention (next lectures)

- No query is given; queries are the same as values: $Q = V$
- Self-attention is used to **encode** a sequence V to a **contextualized sequence** \tilde{V}
 - Each encoded vector \tilde{v}_i is a **contextual embedding** of the corresponding input vector v_i



Attention – summary

- Attention is a way to define the **distribution of focus** on inputs based on a query, and create a compositional embedding of inputs
- Attention networks define an attention distribution over inputs and calculate their weighted sum
- The original definition of attention network has two inputs: **key vectors K** , and **value vectors V**
 - Key vectors are used to calculate attentions
 - and output is the weighted sum of value vectors
 - In practice, in most cases $K = V$.
 - In this course, we use our slightly simplified definition

Agenda

- Machine Translation
- Attention Networks
- **Attention in practice**
 - Seq2seq with attention
 - Hierarchical document classification

Neural Machine Translation (NMT)

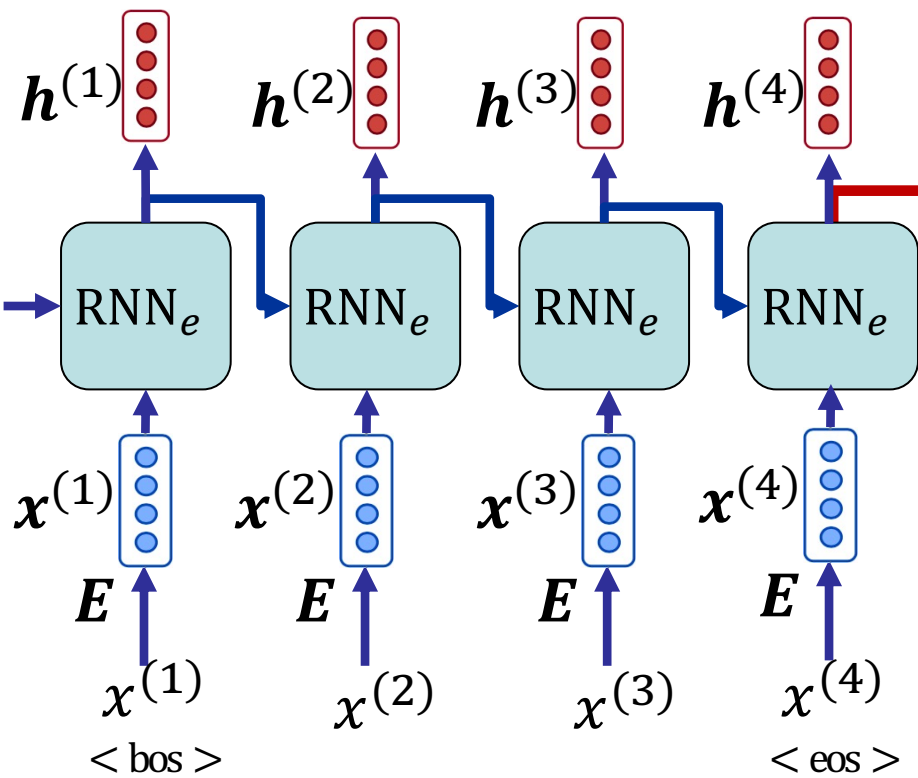
- Given the source language sentence X and target language sentence Y , NMT uses seq2seq models to calculate the **conditional language model**:

$$P(Y|X)$$

- A language model of the target language
 - Conditioned on the source language
- In contrast to SMT, no need for pre-defined alignments! 🎉
- We can simply use a seq2seq with two RNNs

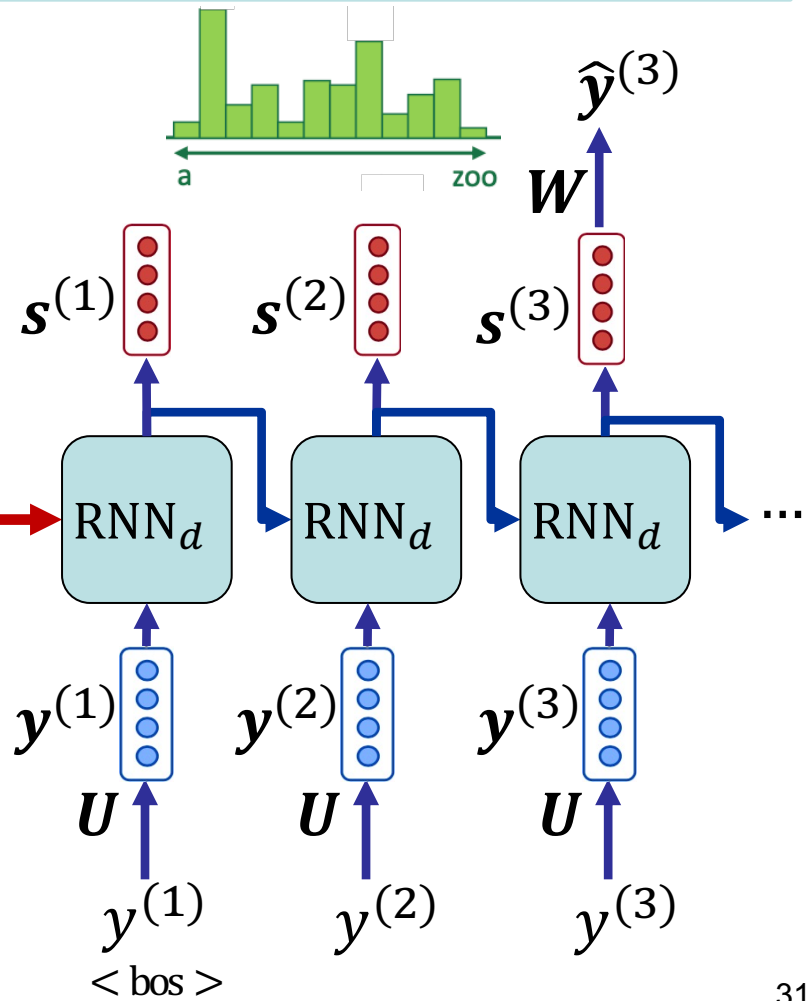
Seq2seq with two RNNs (recap)

ENCODER



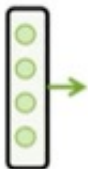
DECODER

$\hat{y}^{(i)}$: predicted probability distribution of the next target word, given the source sequence and previous target words



Seq2seq with two RNNs – training (recap)

Encoder: read source



we are here

Source: Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

Target: I saw a cat on a mat <eos>

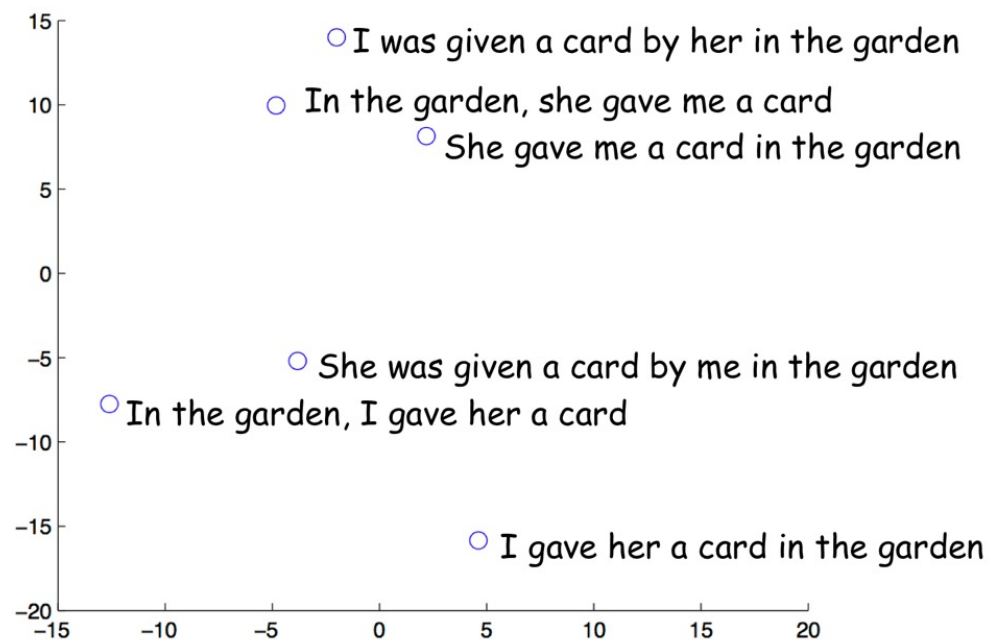
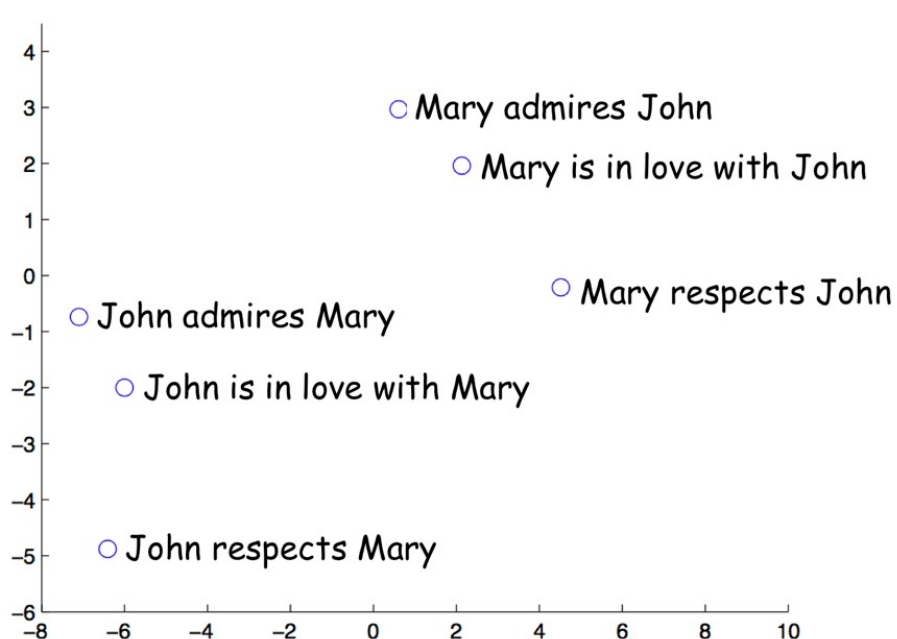
Seq2seq with two RNNs – decoding / beam search (recap)

<bos>

Start with the begin of sentence token or with an empty sequence

Sentence-level semantic representations (recap)

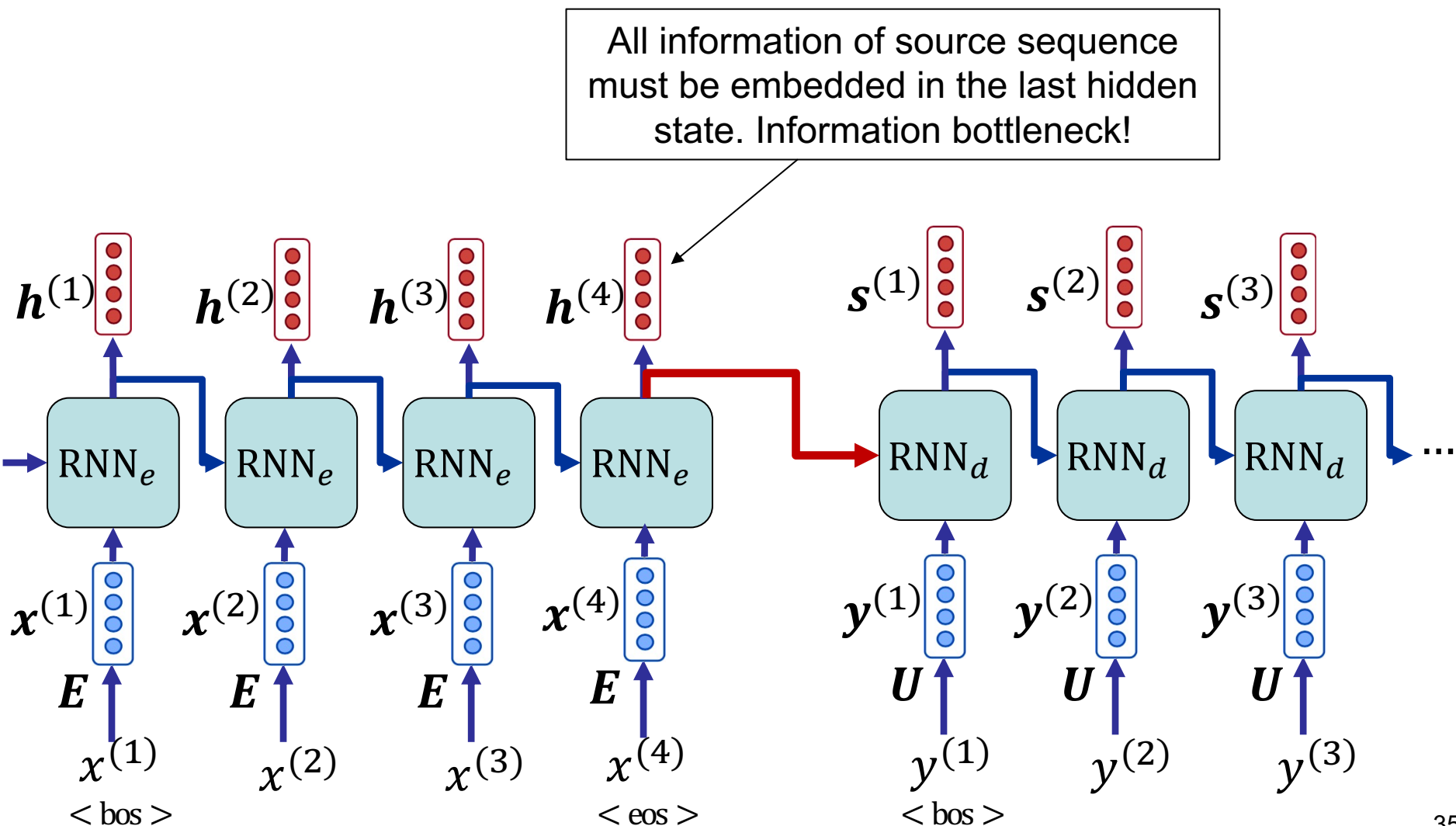
- 2-dimensional projection of the last hidden states ($h^{(L)}$) of RNN_e that are obtained from different phrases



Bottleneck problem in seq2seq with two RNNs

ENCODER

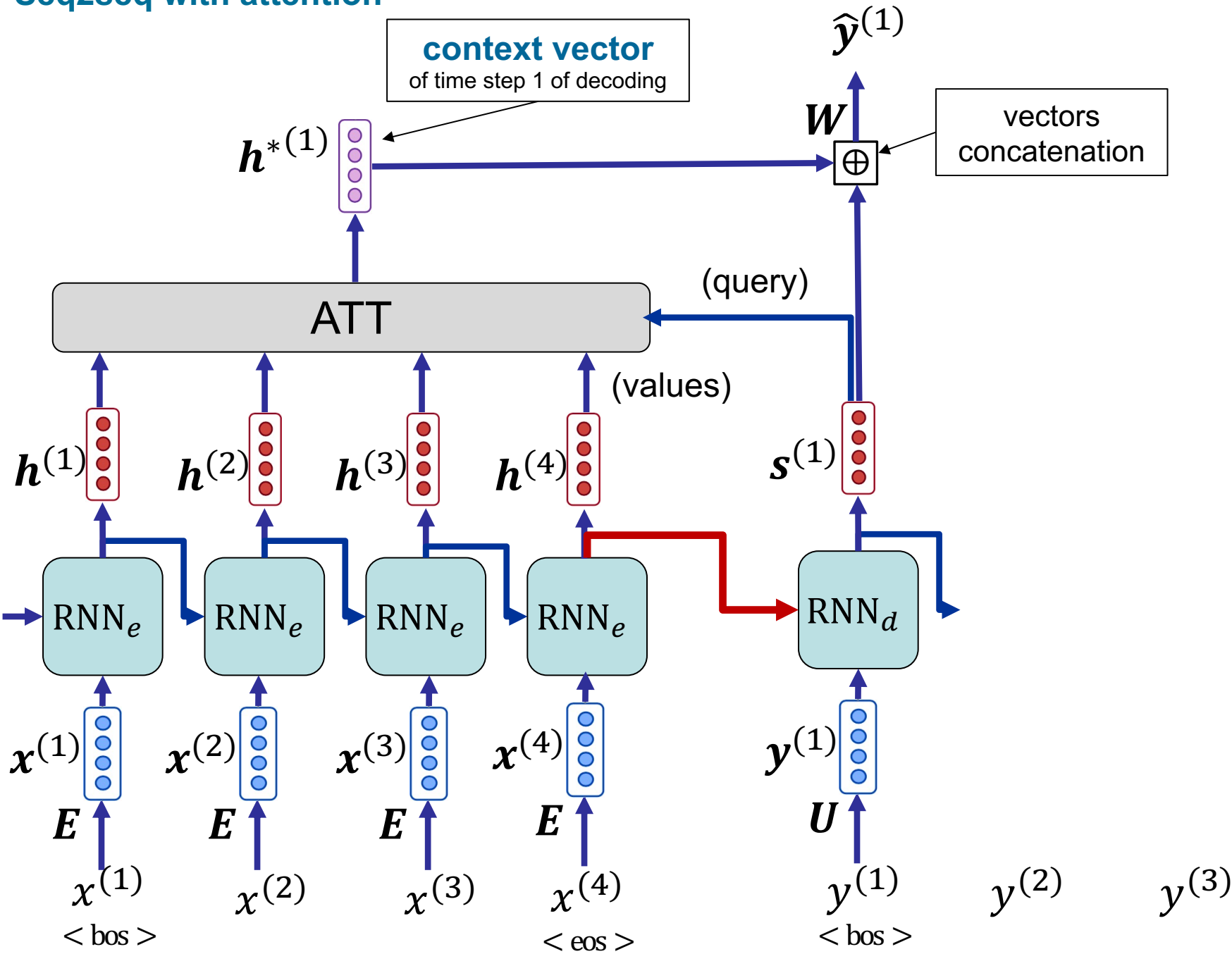
DECODER



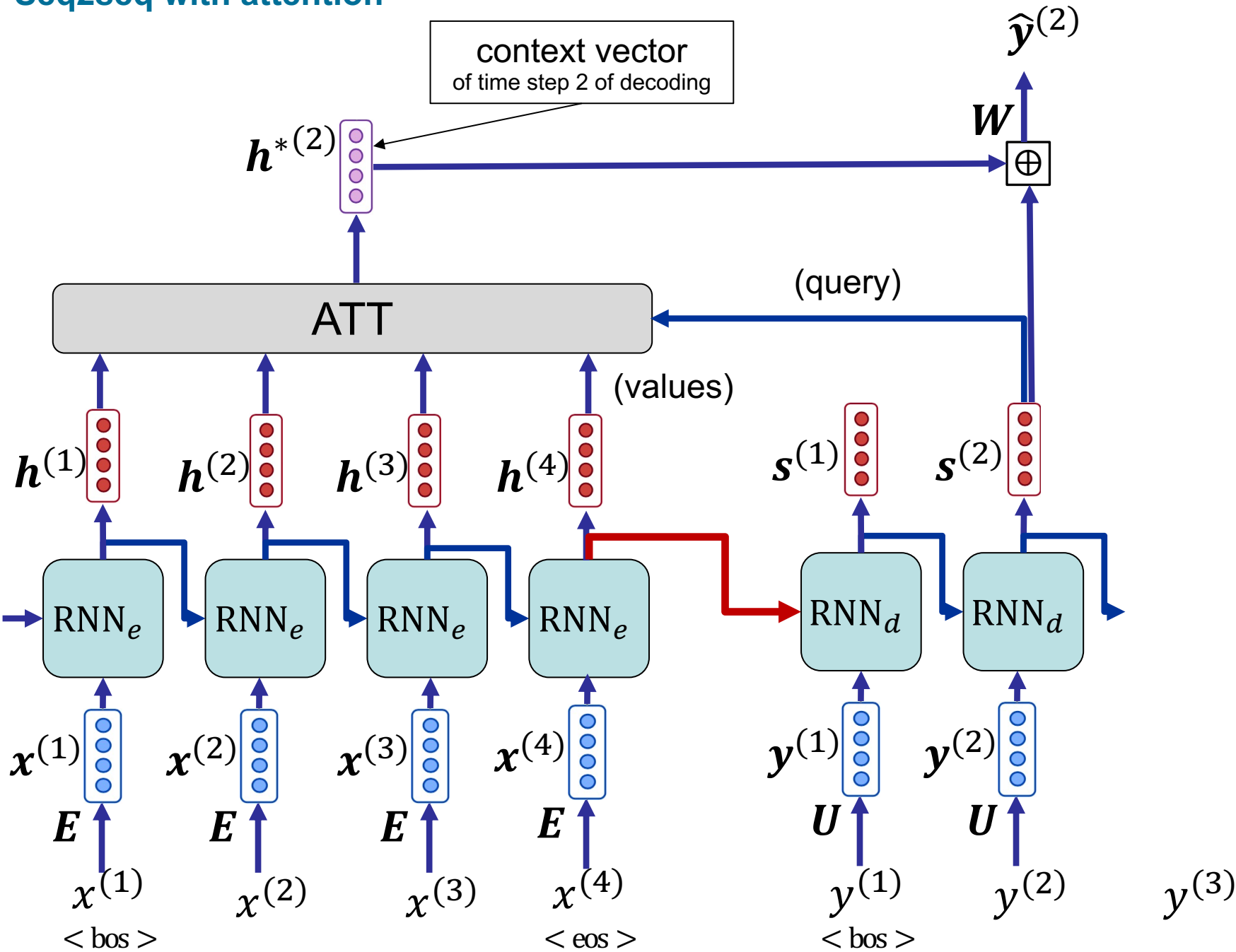
Seq2seq + Attention

- It can be useful, if we allow decoder the **direct access to all elements of source sequence**,
 - Decoder can decide on which element of source sequence, it wants to put attention
- Attention is a solution to the bottleneck problem
- Seq2seq with attention
 - adds an attention network to the architecture of basic seq2seq (two RNNs)
 - At each time step, decoder uses the attention network to **attend to all contextualized vectors** of the source sequence
 - Training and inference (decoding) processes are the same as basic seq2seq

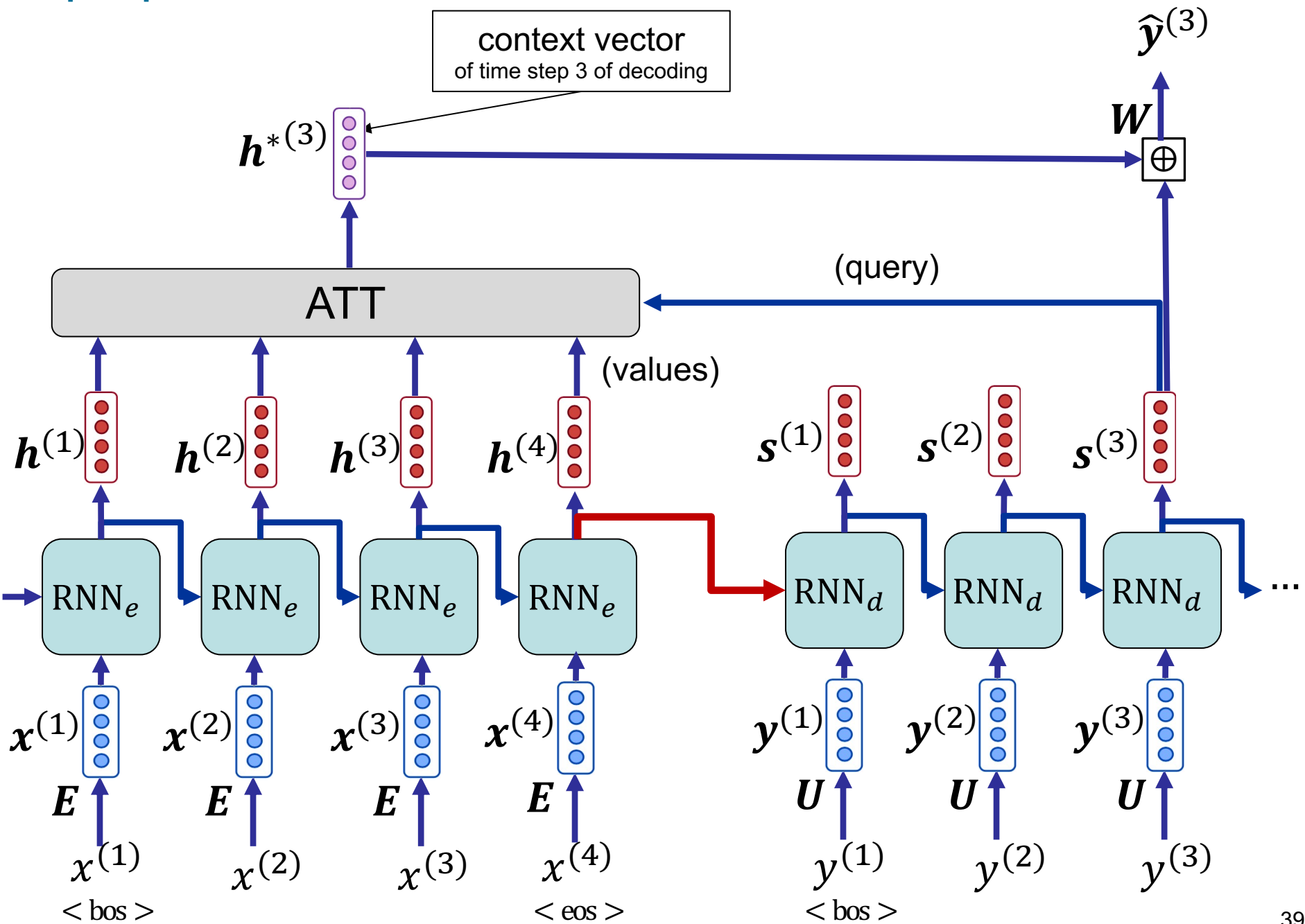
Seq2seq with attention



Seq2seq with attention



Seq2seq with attention



Seq2seq with attention – formulation

- Two sets of vocabularies
 - \mathbb{V}_e is the set of vocabularies for source sequences
 - \mathbb{V}_d is the set of vocabularies for target sequences

ENCODER is the same as seq2seq with two RNNs

- Encoder embedding
 - Encoder embeddings for source words (\mathbb{V}_e) \rightarrow **E**
 - Embedding of the source word $x^{(l)}$ (at time step l) $\rightarrow \mathbf{x}^{(l)}$

- Encoder RNN:

$$\mathbf{h}^{(l)} = \text{RNN}_e(\mathbf{h}^{(l-1)}, \mathbf{x}^{(l)})$$

Parameters are shown in red

Seq2seq with attention – formulation

DECODER – input

- Decoder embedding

- Decoder embeddings *at* input for target words (\mathbb{V}_d) \rightarrow **U**
- Embedding of the target word $y^{(t)}$ (at time step t) $\rightarrow \mathbf{y}^{(t)}$

- Decoder RNN

$$\mathbf{s}^{(t)} = \text{RNN}_d(\mathbf{s}^{(t-1)}, \mathbf{y}^{(t)})$$

- The values of the **last hidden state** of the encoder RNN are passed to the **initial hidden state** of the decoder RNN:

$$\mathbf{s}^{(0)} = \mathbf{h}^{(L)}$$

Seq2seq with attention – formulation

DECODER - attention

- Attention context vector

$$\mathbf{h}^{*(t)} = \text{ATT}(\mathbf{s}^{(t)}, \{\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)}\})$$

For instance, if ATT is a “basic dot-product attention”, this is done by:

- First calculating non-normalized attentions:

$$\tilde{\alpha}_l^{(t)} = \mathbf{s}^{(t)\top} \mathbf{h}_l$$

- Then, normalizing the attentions:

$$\alpha^{(t)} = \text{softmax}(\tilde{\alpha}^{(t)})$$

- and finally calculating the weighted sum of encoder hidden states

$$\mathbf{h}^{*(t)} = \sum_{l=1}^L \alpha_l^{(t)} \mathbf{h}_l$$

Seq2seq with attention – formulation

DECODER - output

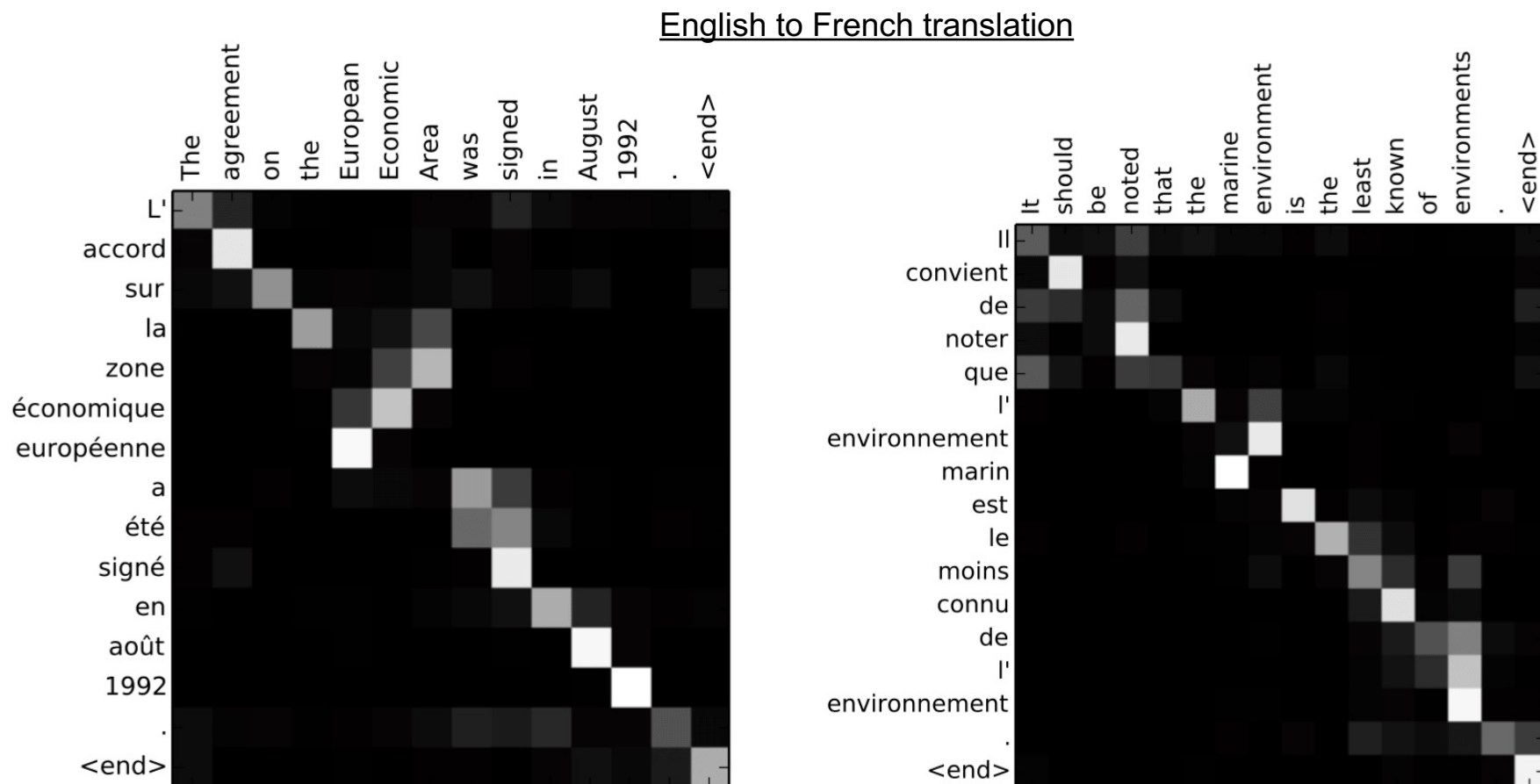
- Decoder output prediction
 - Predicted probability distribution of words at the next time step:

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{W}[\mathbf{s}^{(t)}; \mathbf{h}^{*(t)}] + \mathbf{b}) \in \mathbb{R}^{|\mathcal{V}|}$$

[;] denotes the concatenation of two vectors

Alignment in NMT (seq2seq with attention)

- Attention automatically learns (nearly) alignment



Seq2seq with attention – summary

- Attention on source sequence facilitates the **focus on relevant words** and better **flow of information**
 - It also helps avoiding **vanishing gradient** problem by providing a shortcut to faraway states
- Attention provides **some interpretability**
 - Looking at attention distributions, we can assume what the decoder is focusing on
 - It is however disputable whether attention distributions should be taken as model explanations (particularly in Transformers)!
 - Jain, Sarthak, and Byron C. Wallace. "Attention is not Explanation." *In proc. of NAACL-HTL* 2019. <https://www.aclweb.org/anthology/N19-1357.pdf>
 - Wiegrefe, Sarah, and Yuval Pinter. "Attention is not not Explanation." *In proc. of EMNLP-IJCNLP*. 2019. <https://www.aclweb.org/anthology/D19-1002/>

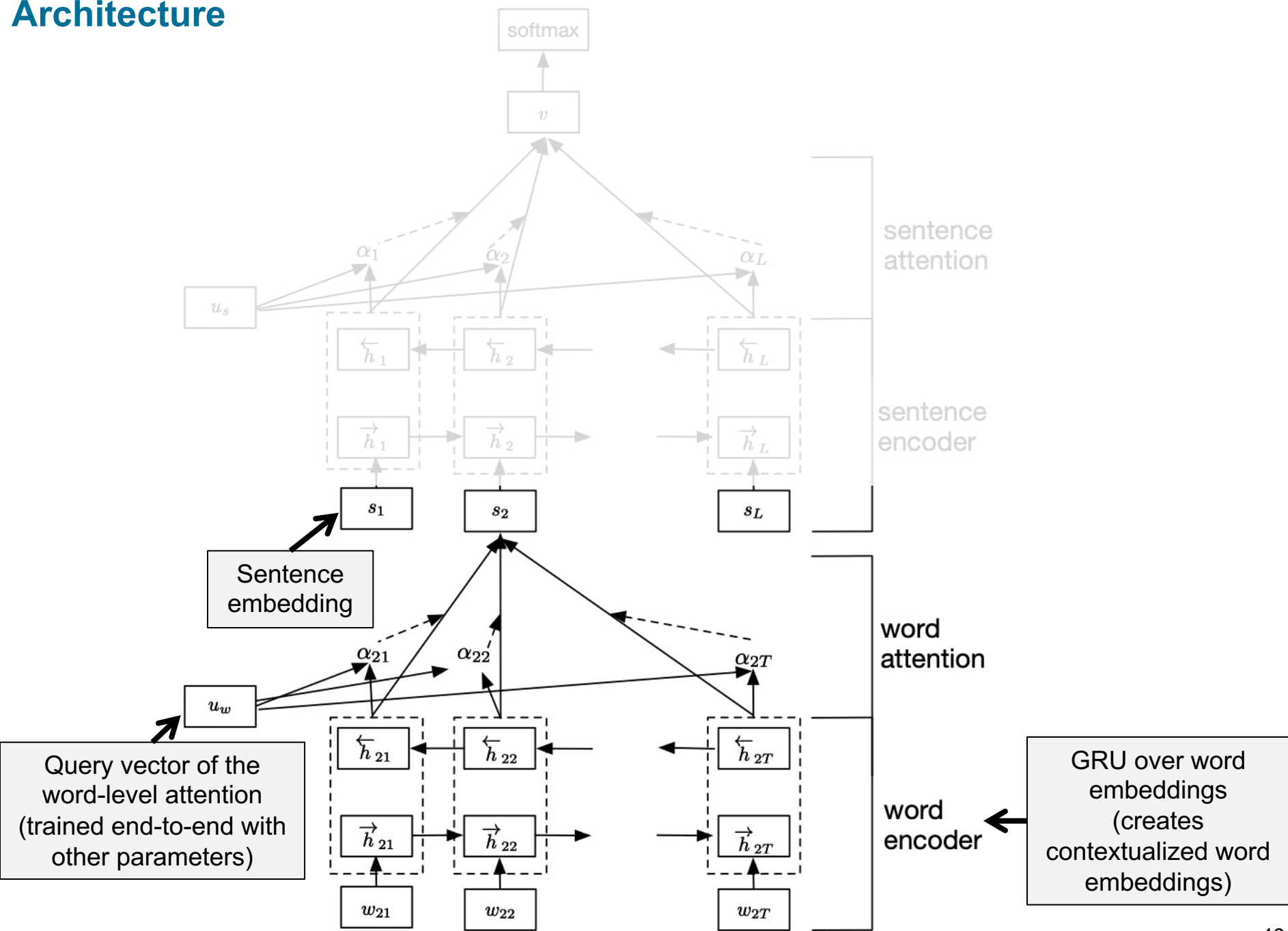
Agenda

- Machine Translation
- Attention Networks
- **Attention in practice**
 - Seq2seq with attention
 - **Hierarchical document classification**

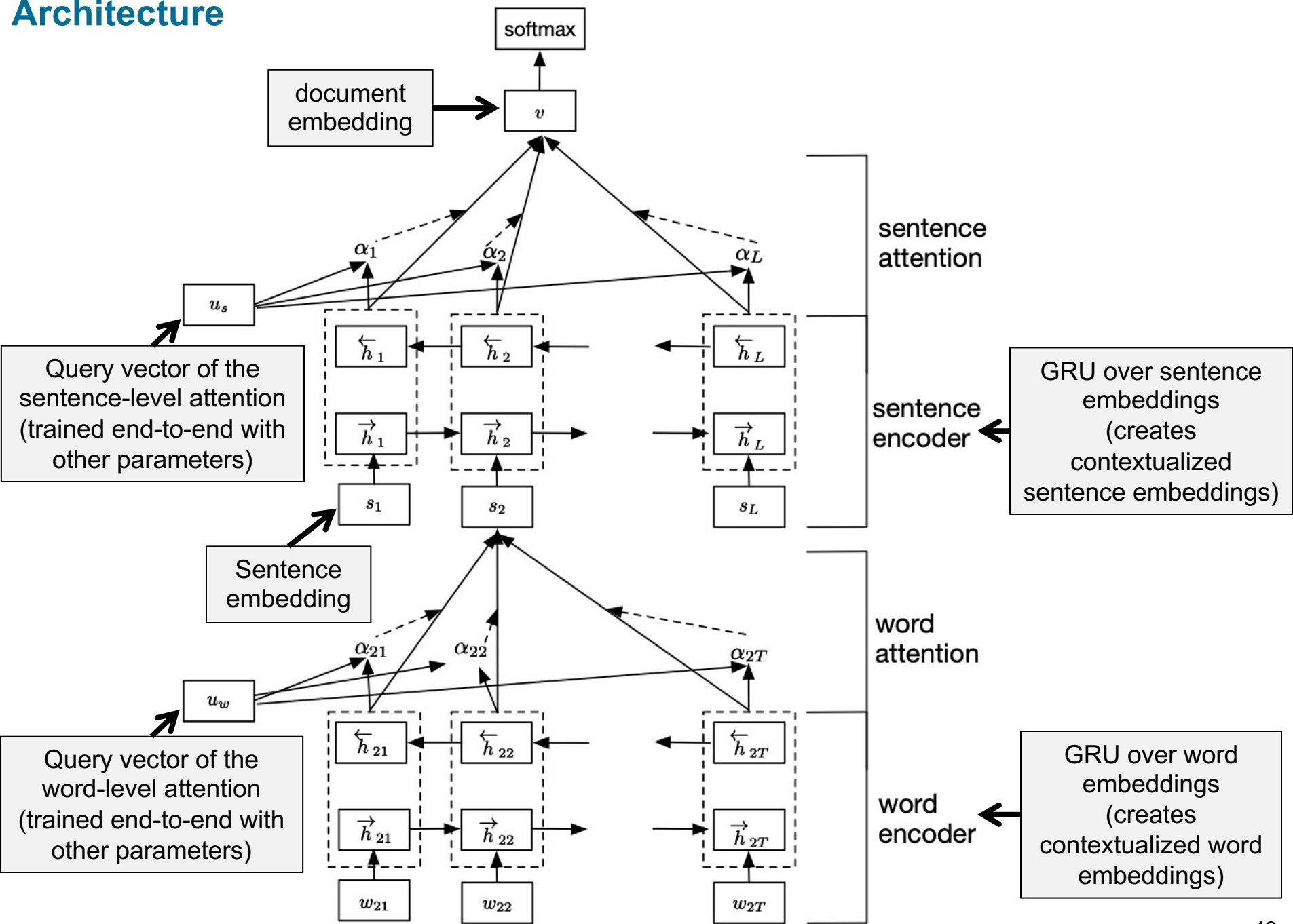
Hierarchical document classification with attention

- Document classification with attention
 - An attention network is applied to word embeddings as **values** (inputs) to compose a document vector (output)
 - Document embedding is then used as features for classification
 - The **query** of the attention network is a randomly initialized parameter vector, whose weights are trained end-to-end with the model
- Hierarchical document classification
 - Split the document into sentences
 - Use a word-level attention to create a sentence embedding from the word embeddings of each sentence
 - Use a sentence-level attention to create the document embedding from the sentence embeddings

Architecture



Architecture



Examples

GT: 4 Prediction: 4

pork belly = delicious .
scallops ?
i do n't .
even .
like .
scallops , and these were a-m-a-z-i-n-g .
fun and tasty cocktails .
next time i 'm in phoenix , i will go
back here .
highly recommend .


GT: 0 Prediction: 0

terrible value .
ordered pasta entree .
.
\$ 16.95 good taste but size was an
appetizer size .
.
no salad , no bread no vegetable .
this was .
our and tasty cocktails .
our second visit .
i will not go back .

Figure 5: Documents from Yelp 2013. Label 4 means star 5, label 0 means star 1.


Example

GT: 1 Prediction: 1



why does zebras have stripes ?
what is the purpose or those stripes ?
who do they serve the zebras in the
wild life ?
this provides camouflage - predator
vision is such that it is usually difficult
for them to see complex patterns

GT: 4 Prediction: 4



how do i get rid of all the old web
searches i have on my web browser ?
i want to clean up my web browser
go to tools > options .
then click “ delete history ” and “
clean up temporary internet files . ”

Figure 6: Documents from Yahoo Answers. Label 1 denotes Science and Mathematics and label 4 denotes Computers and Internet.

Recap

- Attention is a general deep learning approach to learn to distribute the focus on certain parts, and compose outputs
- Attention significantly helps seq2seq models in machine translation!
- Attention can also be used for encoding text, e.g., in document classification

