

344.063 KV Special Topic:

Natural Language Processing with Deep Learning

Neural Machine Translation with Attention Networks



Navid Rekab-saz

navid.rekabsaz@jku.at

Agenda

- Attention Networks
- Machine Translation
 - Seq2seq with Attention
- Hierarchical document classification

Agenda

- **Attention Networks**
- Machine Translation
 - Seq2seq with Attention
- Hierarchical document classification

Attention Networks

- Attention is a deep learning architecture ...
 - to obtain a **composed** output embedding \mathbf{o} ...
 - from a set (matrix) of input values V ...
 - based on a given query embedding q

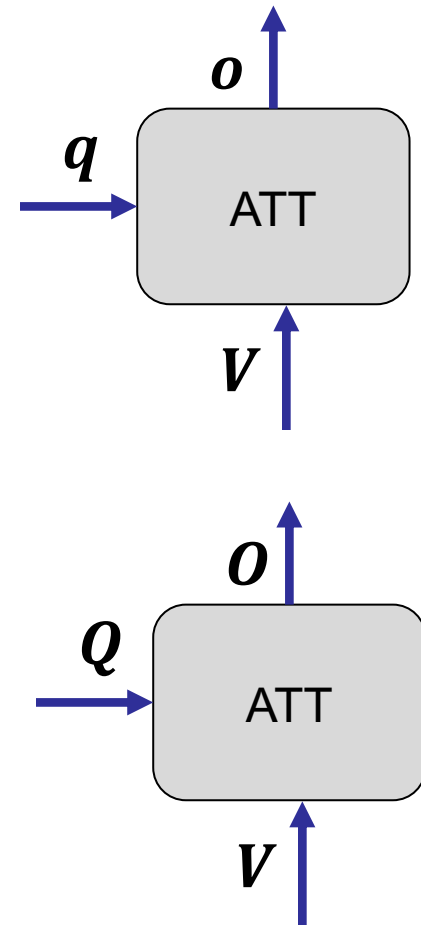
- General form of an attention network:

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

- If a set/matrix of queries Q is given, the output will become a set/matrix O :

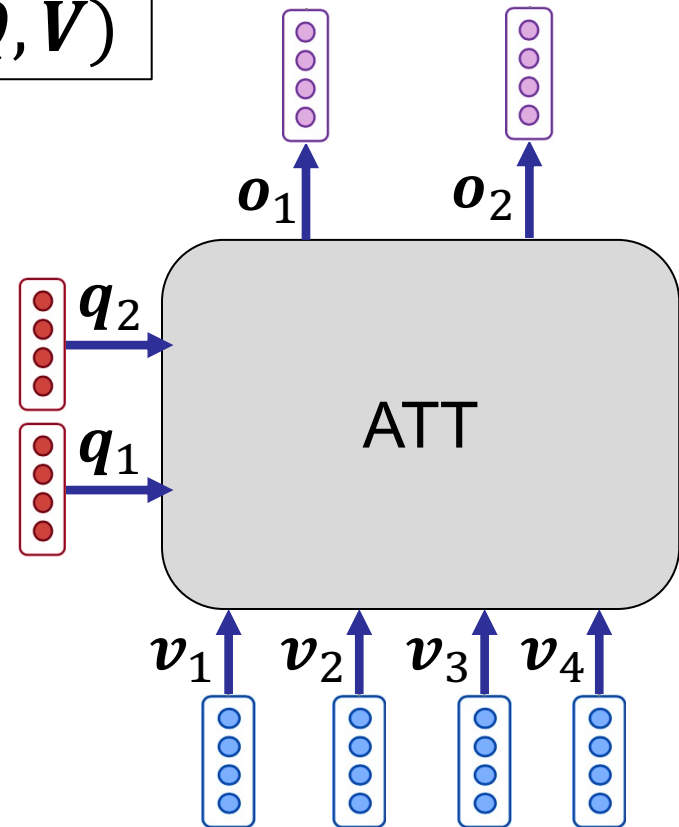
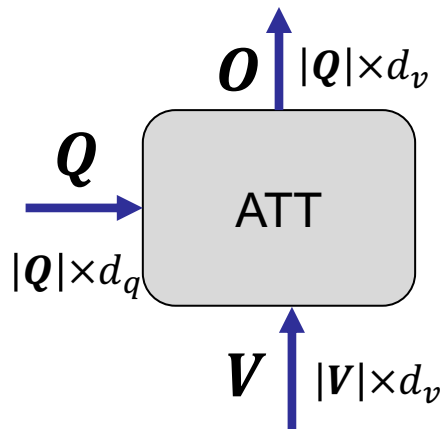
$$O = \text{ATT}(Q, V)$$

- where each output vector belongs to its respective query vector



Attention Networks

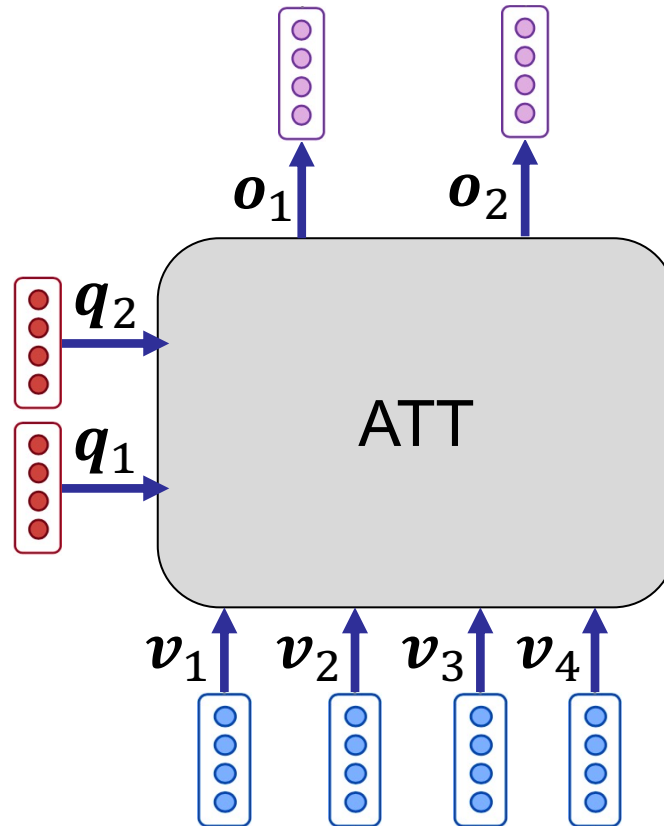
$$\mathbf{O} = \text{ATT}(\mathbf{Q}, \mathbf{V})$$



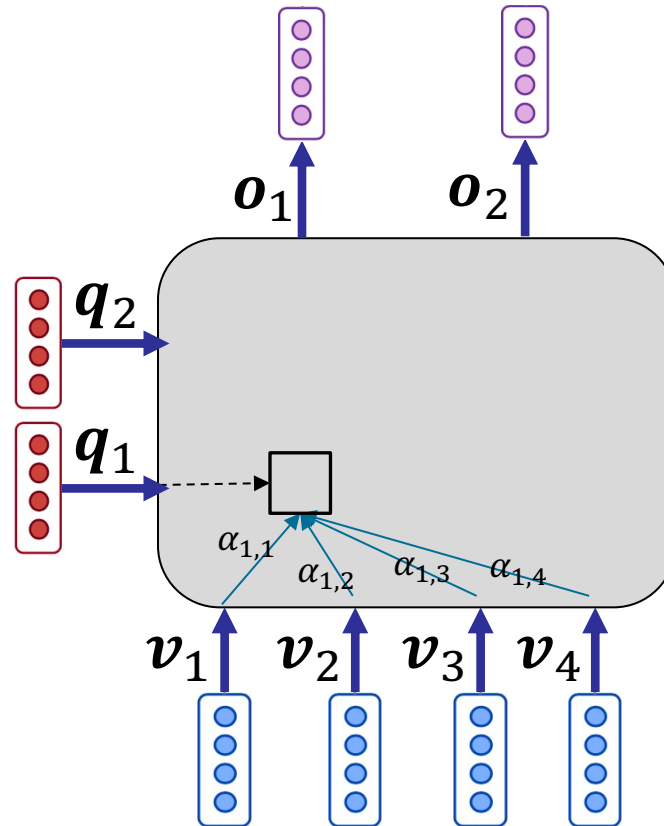
- d_q, d_v are embedding dimensions of query and value vectors, respectively

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

Attention Networks

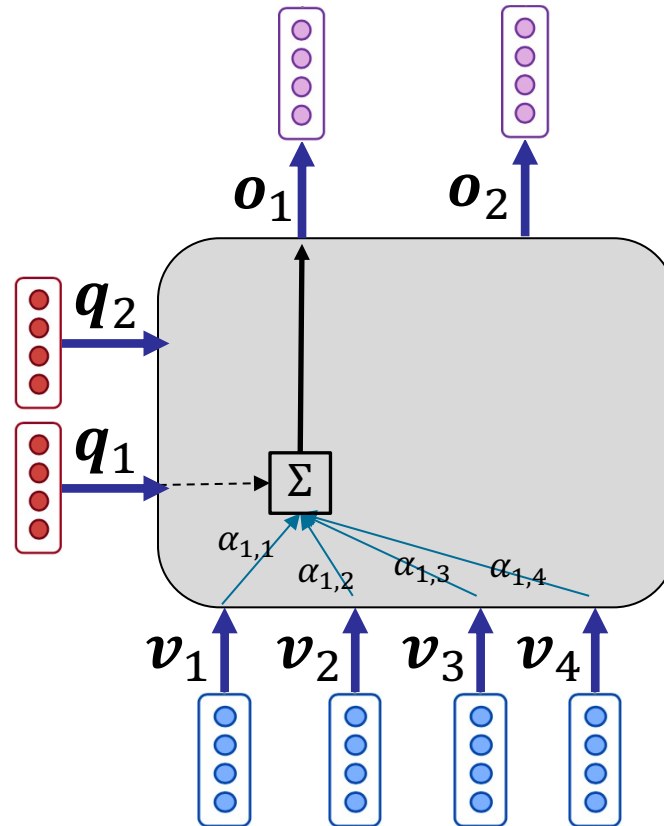


Attention Networks



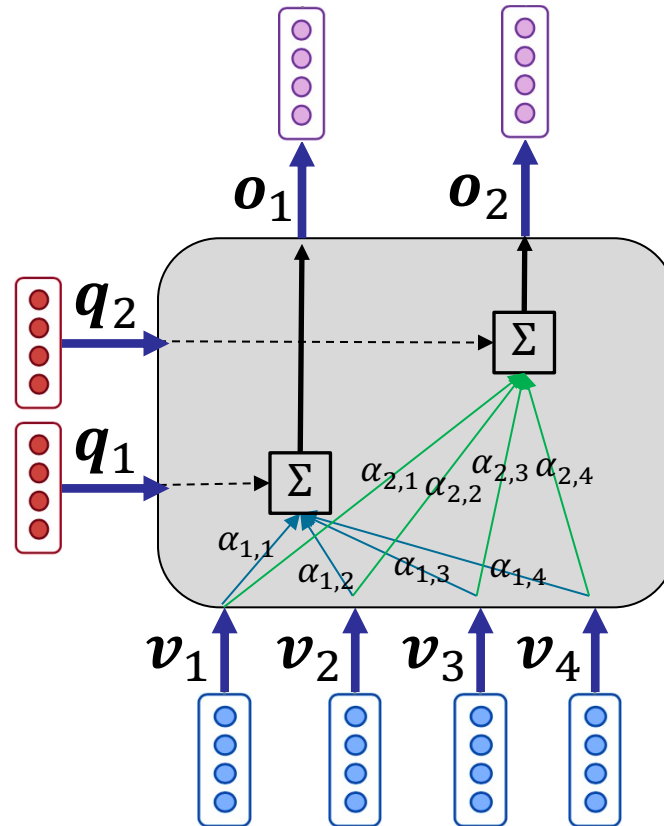
$\alpha_{i,j}$ is the attention of query q_i on value v_j

Attention Networks



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

Attention Networks



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

Attention Networks – definition

- Given a matrix of **values** V and a matrix of **queries** Q , for each query vector $q \in Q$, an **attention network** ...
 - first assigns an **attention score** to each value vector $v \in V$ based on the **similarity** of q to v ,...
 - then turns the attention scores to a **probability distribution of attentions** over value vectors, ...
 - and finally uses the **attentions** to calculate the **weighted sum** of the value vectors as the corresponding output o of the query vector q
- The output of attention networks can be viewed as a **weighted aggregation** of the value vectors, where the query (through attentions) defines the proportion of the contribution of each value vector.

Attention Networks – formulation

- Given query vector \mathbf{q}_i , an attention network uses the **attention similarity function** f to assign a **non-normalized attention score** $\tilde{\alpha}_{i,j}$ to value vector \mathbf{v}_j :

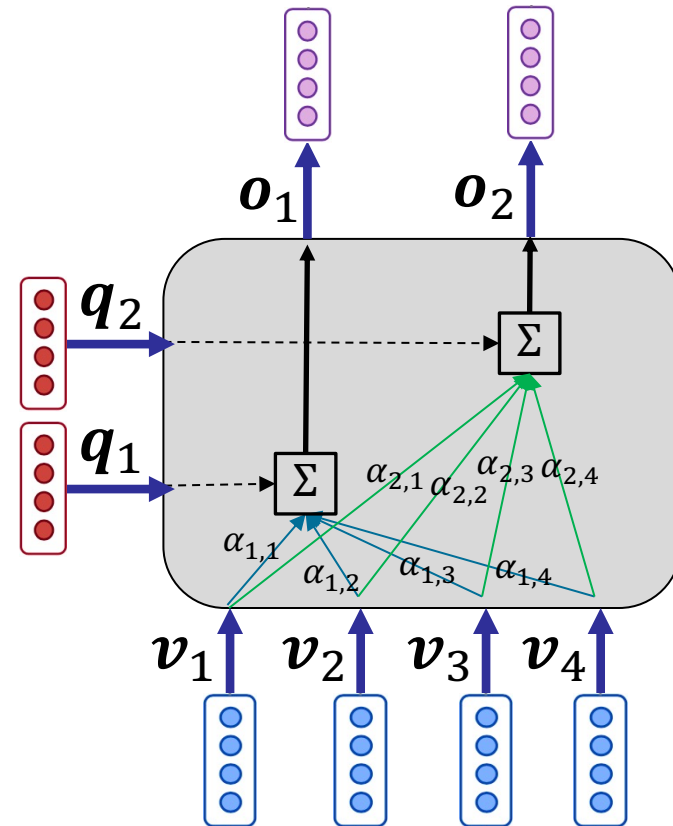
$$\tilde{\alpha}_{i,j} = f(\mathbf{q}_i, \mathbf{v}_j)$$

- Then, the attention scores over values are turned to a probability distribution using softmax:

$$\alpha_i = \text{softmax}(\tilde{\alpha}_i), \quad \sum_{j=1}^{|V|} \alpha_{i,j} = 1$$

- Finally, output vector \mathbf{o}_i regarding query \mathbf{q}_i is defined as the **sum** of the value vectors **weighted** by their corresponding attentions:

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



Attention – first implementation

Basic dot-product attention

- Non-normalized attention scores:

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

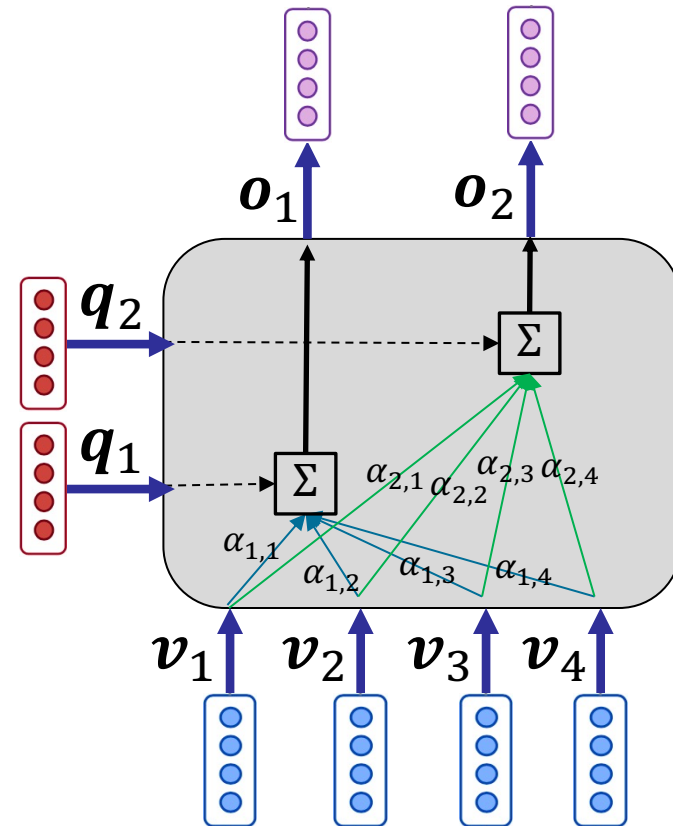
$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{v}_j^T$$

- In this case, $d_q = d_v$
- Attention network has no parameter to learn!

- Softmax over value vectors:

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

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

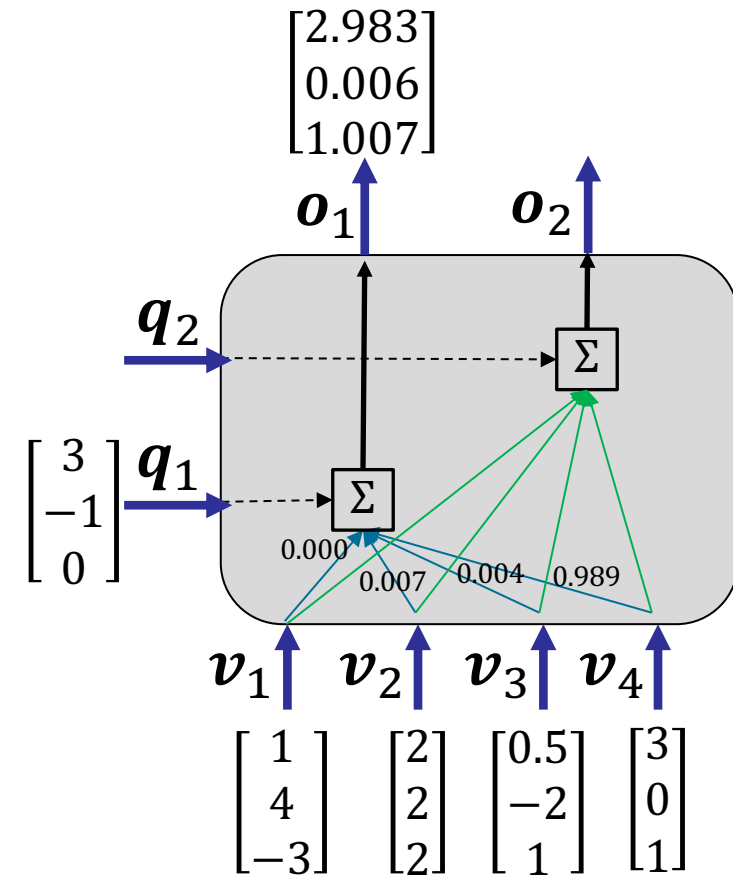


Example

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

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

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

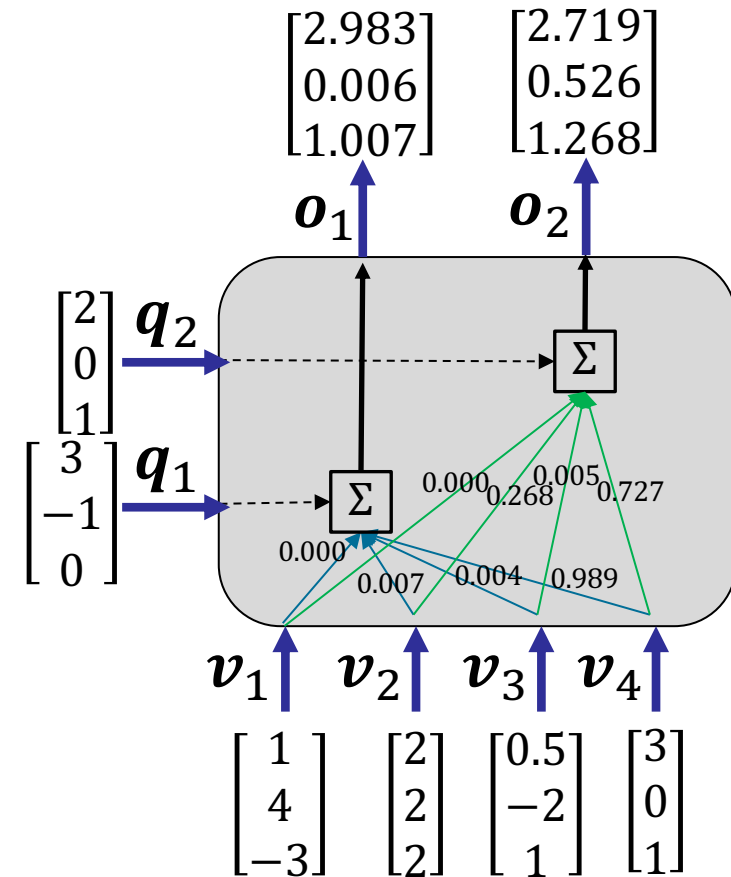


Example

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

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

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



Attention – other implementations

Multiplicative attention

- Non-normalized attention scores:

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

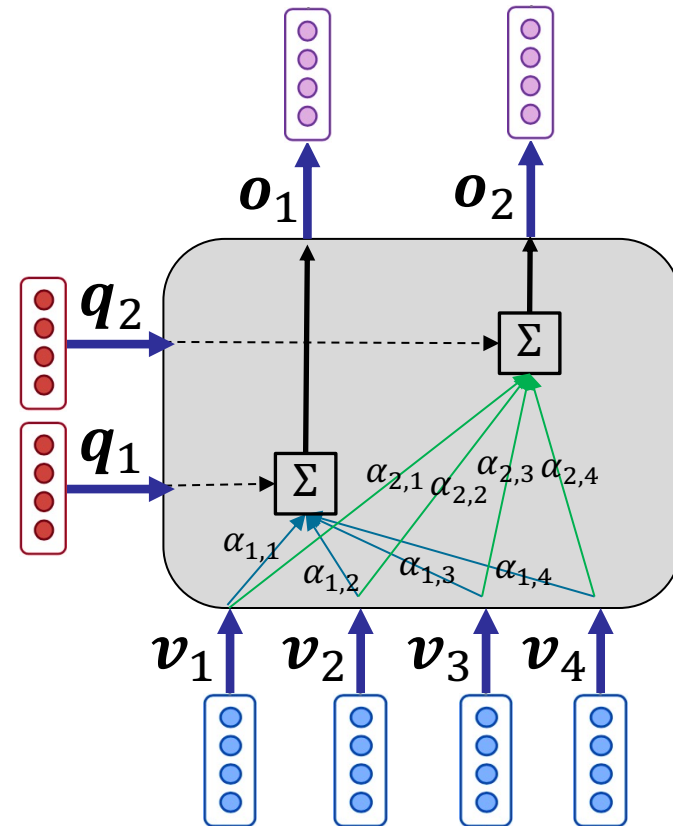
$$\tilde{\alpha}_{i,j} = \mathbf{q}_i \mathbf{W} \mathbf{v}_j^T$$

- \mathbf{W} is a matrix of parameter
- similarity of query to value is defined as a linear function

- 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 – other implementations

Additive attention

- Non-normalized attention scores:

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

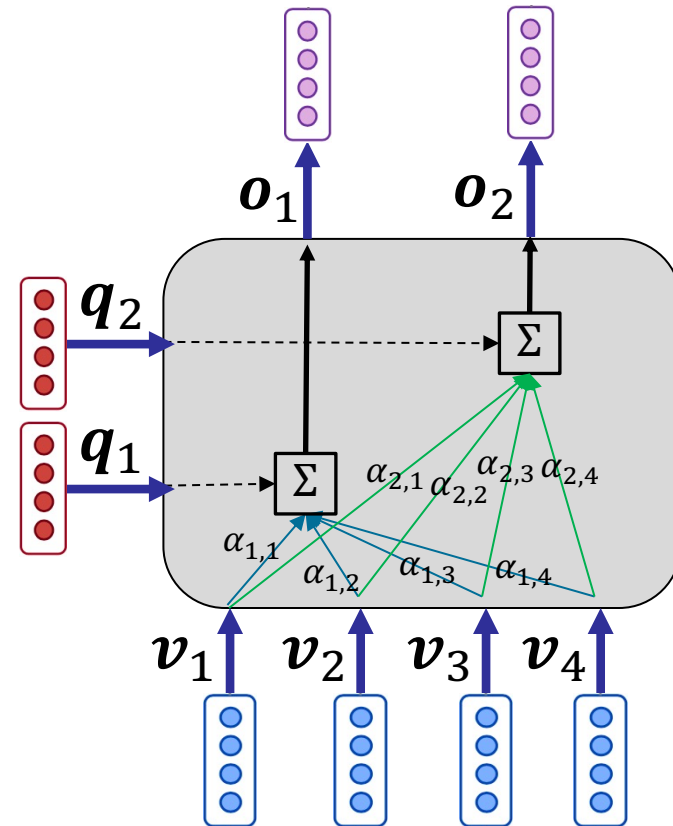
$$\tilde{\alpha}_{i,j} = \mathbf{u}^T \tanh(\mathbf{q}_i \mathbf{W}_1 + \mathbf{v}_j \mathbf{W}_2)$$

- \mathbf{W}_1 , \mathbf{W}_2 , and \mathbf{u} are model parameters
- similarity of query to value is defined as a non-linear function

- 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 – 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

- Attention Networks
- **Machine Translation**
 - 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:

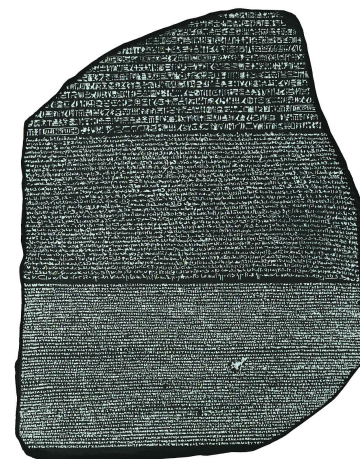
$$= \operatorname{argmax}_Y P(X|Y)P(Y)$$

Translation Model

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

Language Model

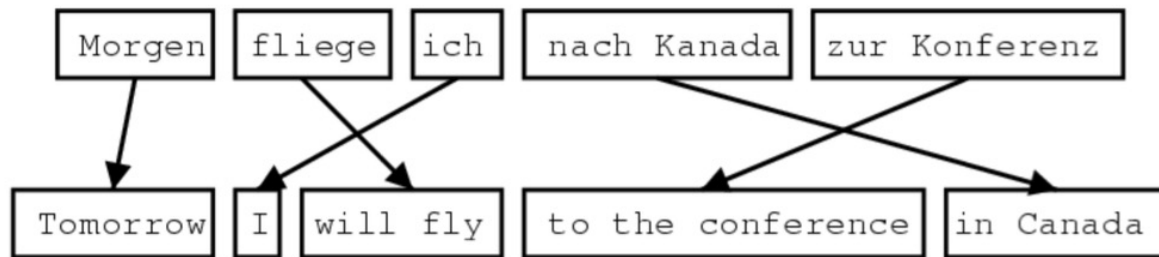
The statistical model that tells us how to write good sentences in the target language (learnt from monolingual data)



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**:

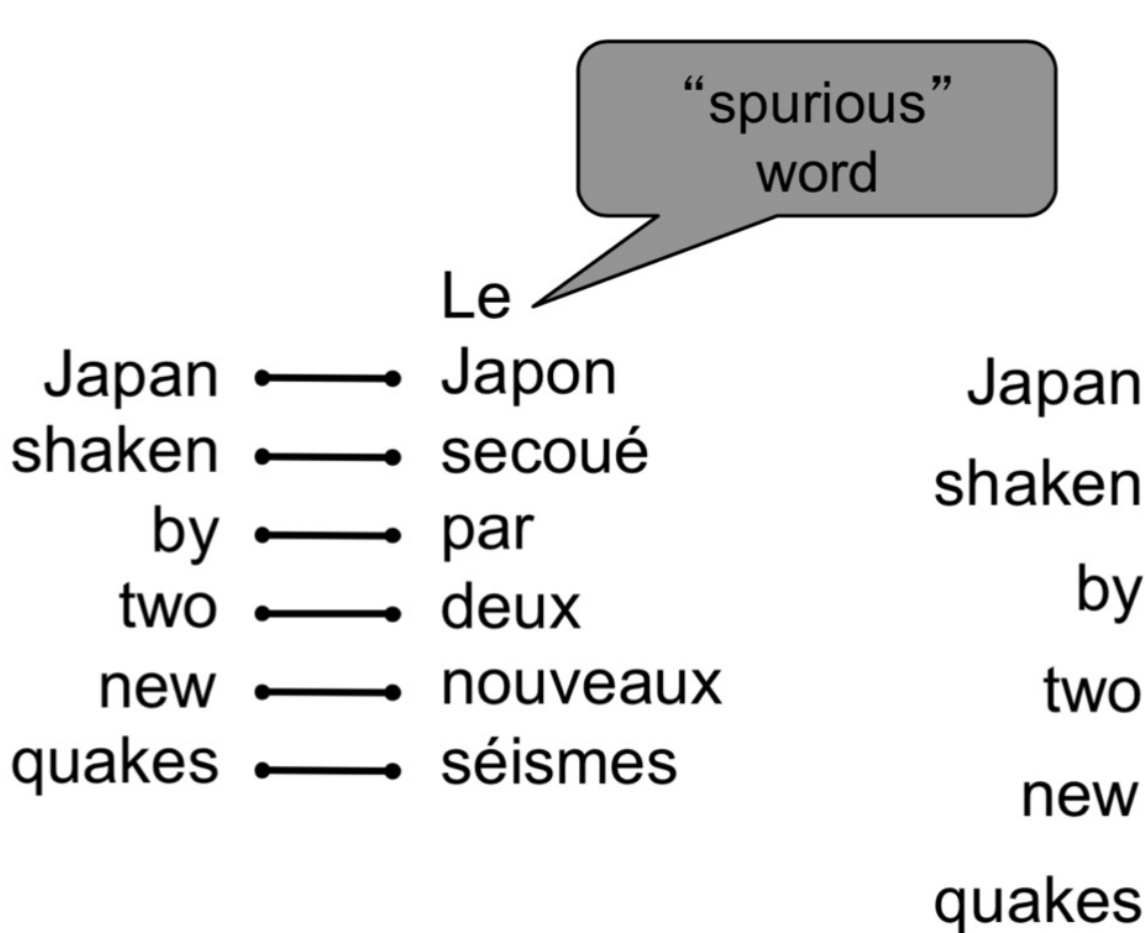


- To this end, the **alignment** latent variable a is added to the formulation of Translation model:

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

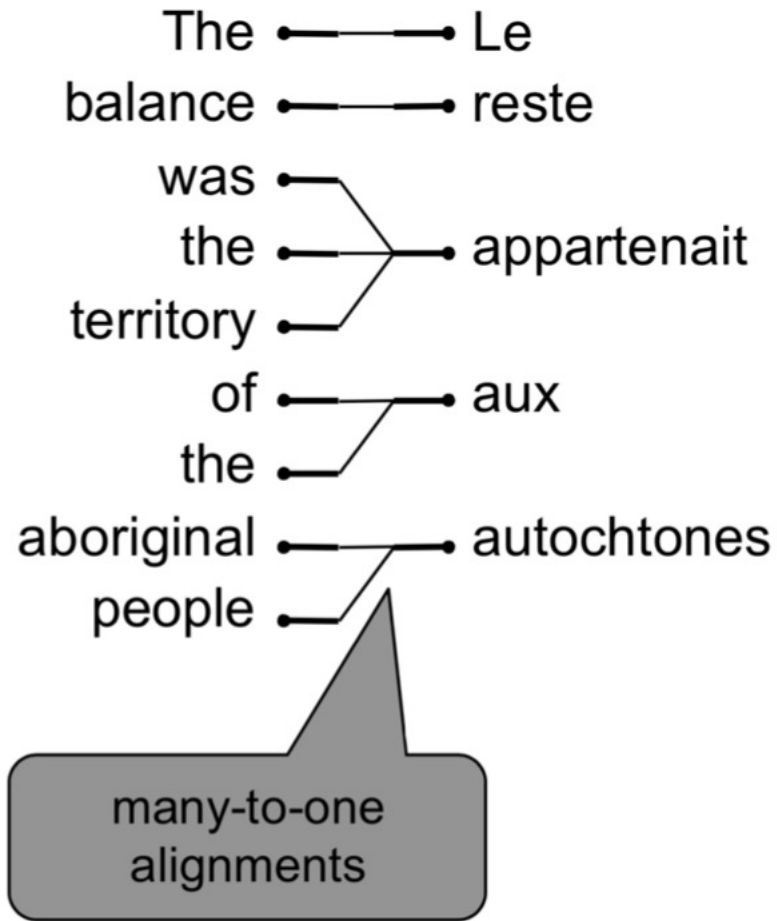
- Alignment ...
 - is a latent variable \rightarrow is not explicitly defined in the data!
 - defines the correspondence between particular words/phrases in the translation sentence pair

Alignment!



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

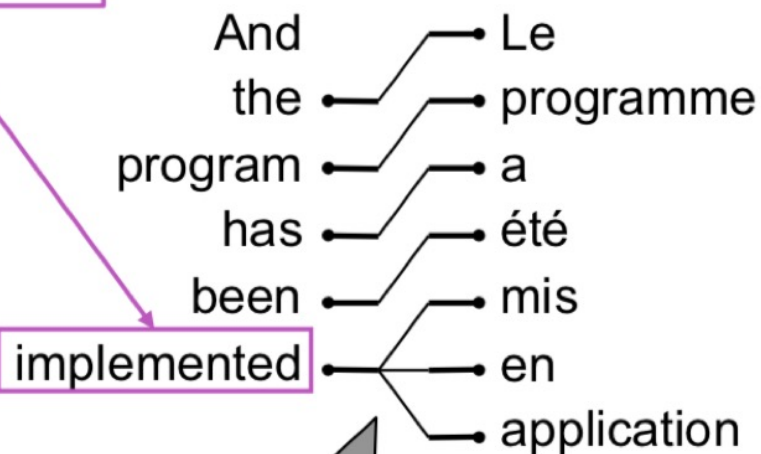
Alignment!



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

Alignment!

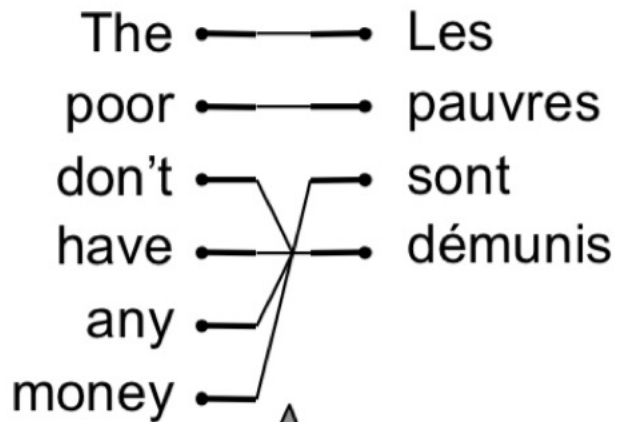
We call this a
fertile word



one-to-many
alignment

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

Alignment!



many-to-many
alignment

	Les	pauvres	sont	démunis
The				
poor				
don't				
have				
any				
money				

phrase
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

- Attention Networks
- **Machine Translation**
 - 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**:

$$P(Y|X)$$

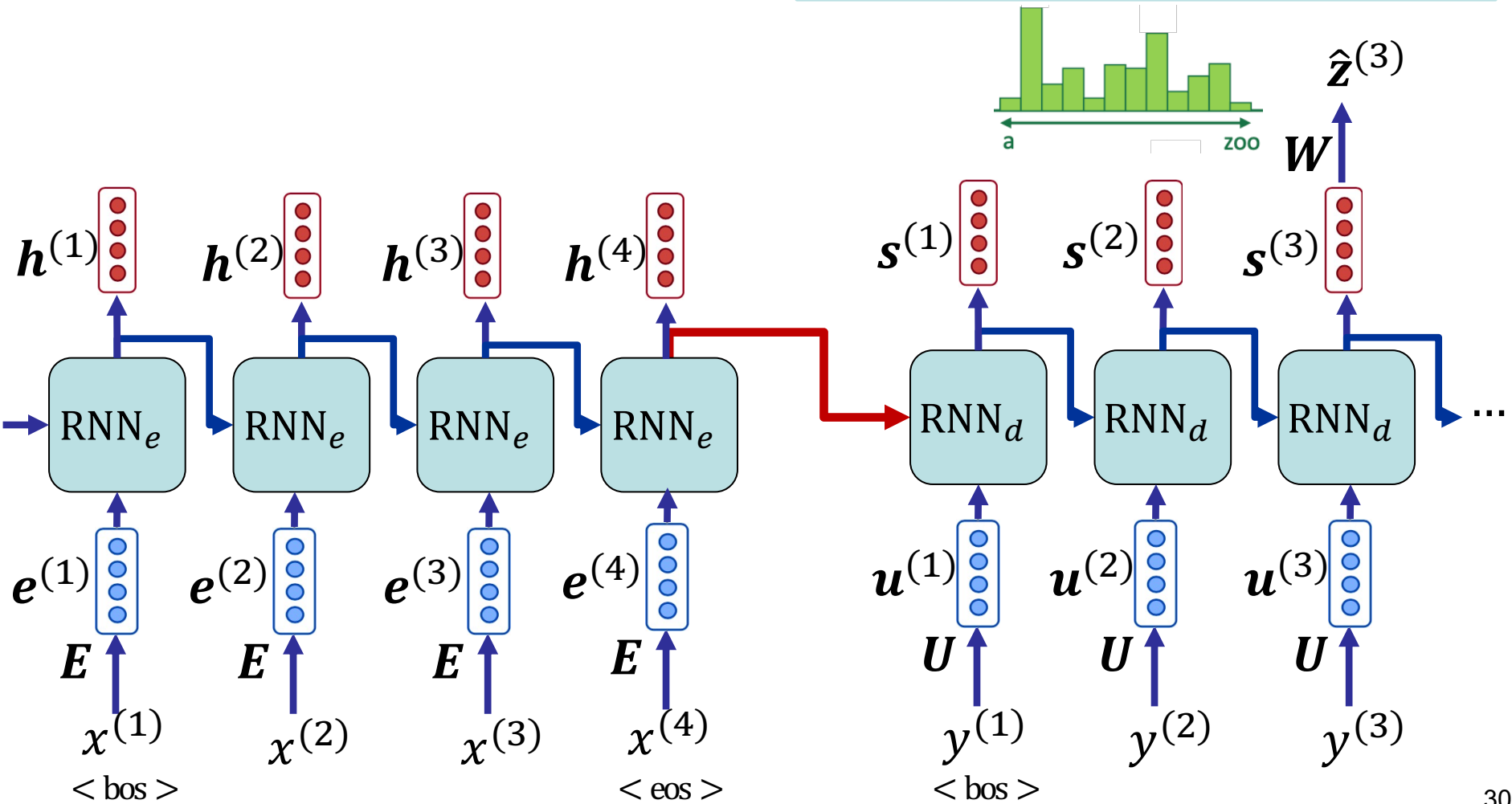
- A language model of the target language
 - Conditioned on the source language
- In contrast to SMT, no need for pre-defined alignments! 🎉
- We can simply use a seq2seq with two RNNs

Seq2seq with two RNNs (recap)

ENCODER

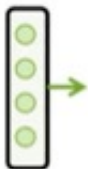
DECODER

$\hat{z}^{(i)}$: 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>

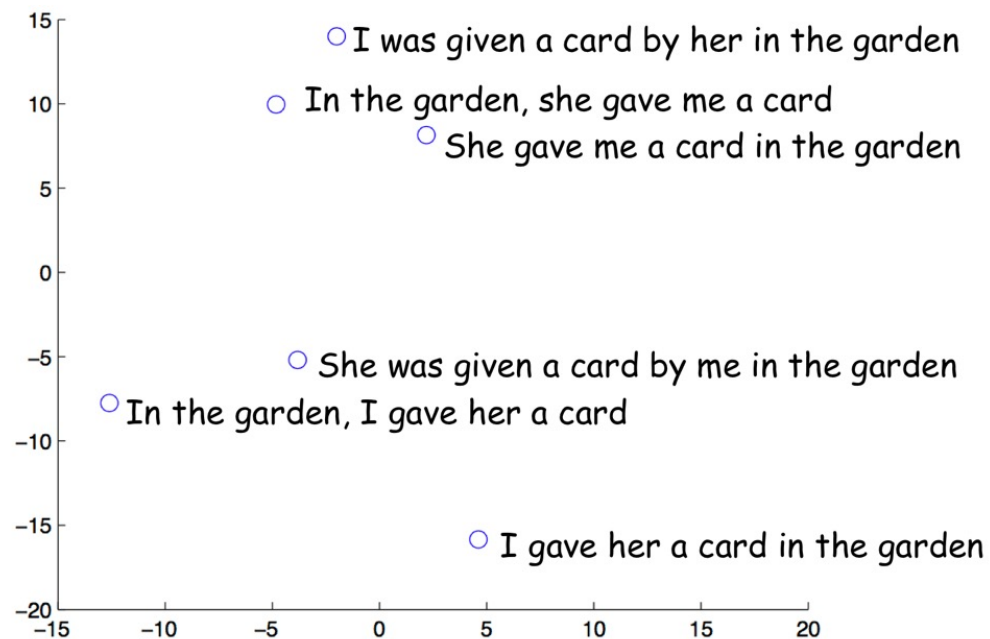
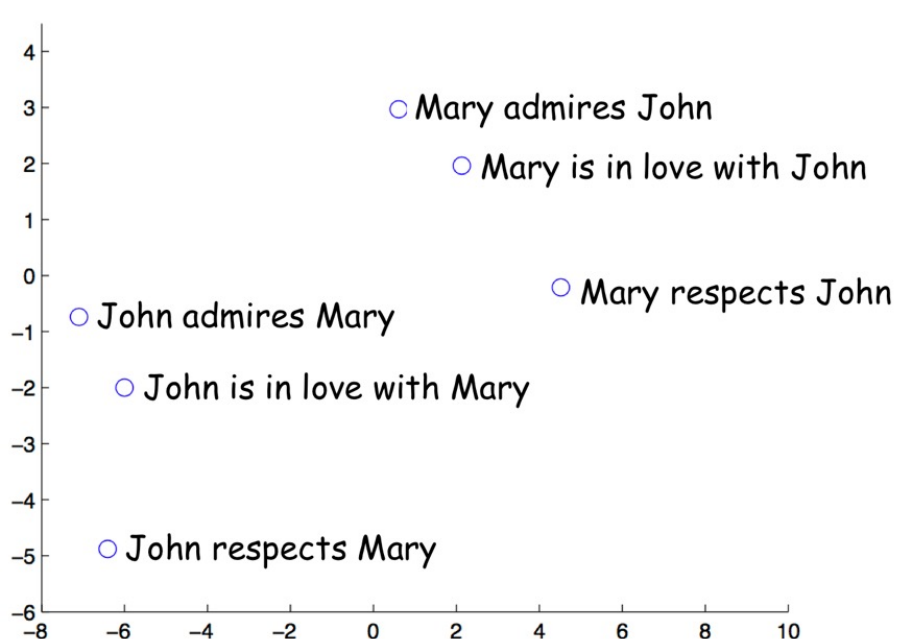
Seq2seq – decoding / beam search (recap)

<bos>

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

Sentence-level semantic representations (recap)

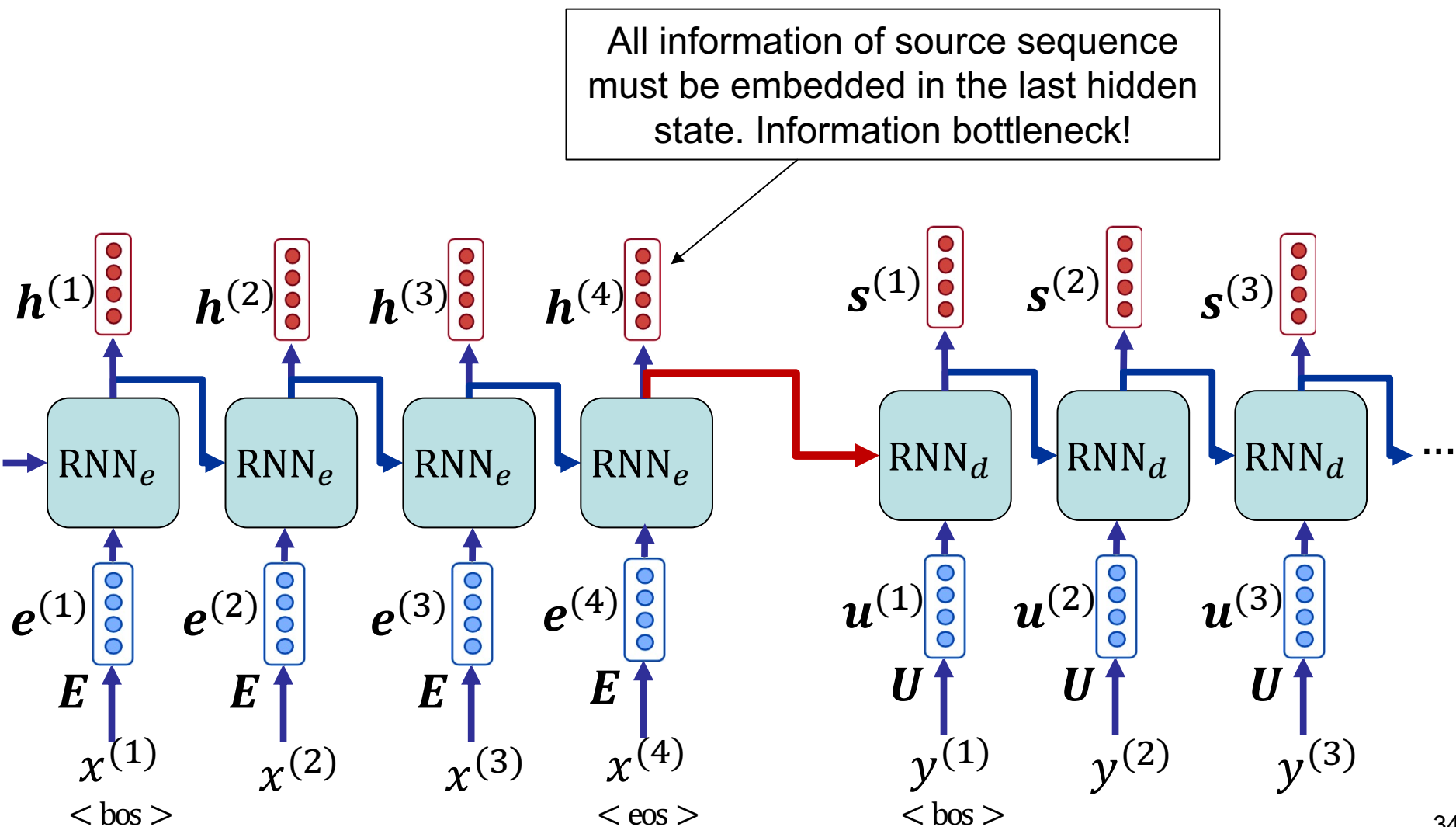
- Two-dimensional projection of the last hidden states $\mathbf{h}^{(L)}$ of RNN_e , obtained from different phrases



Bottleneck problem in seq2seq with two RNNs

ENCODER

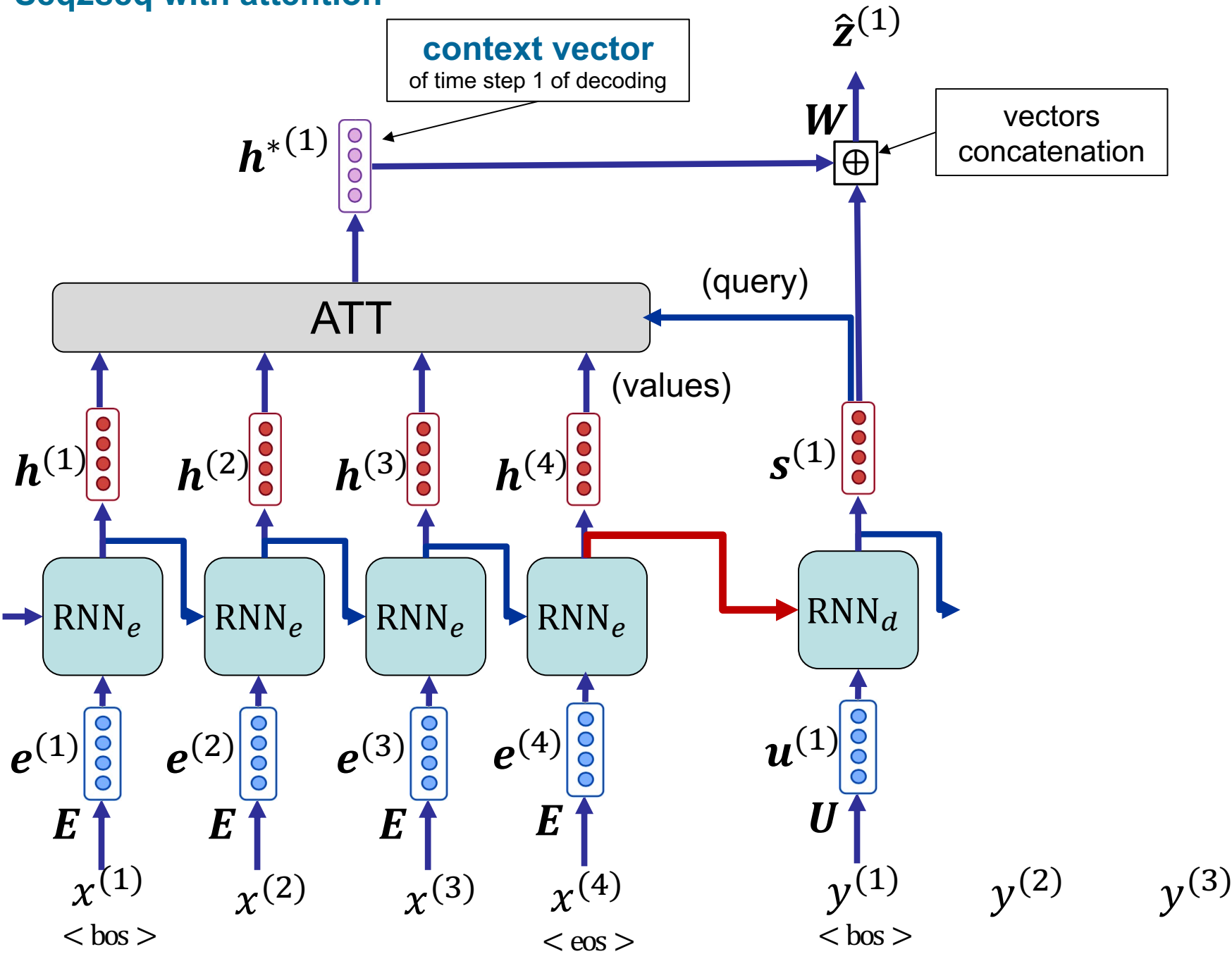
DECODER



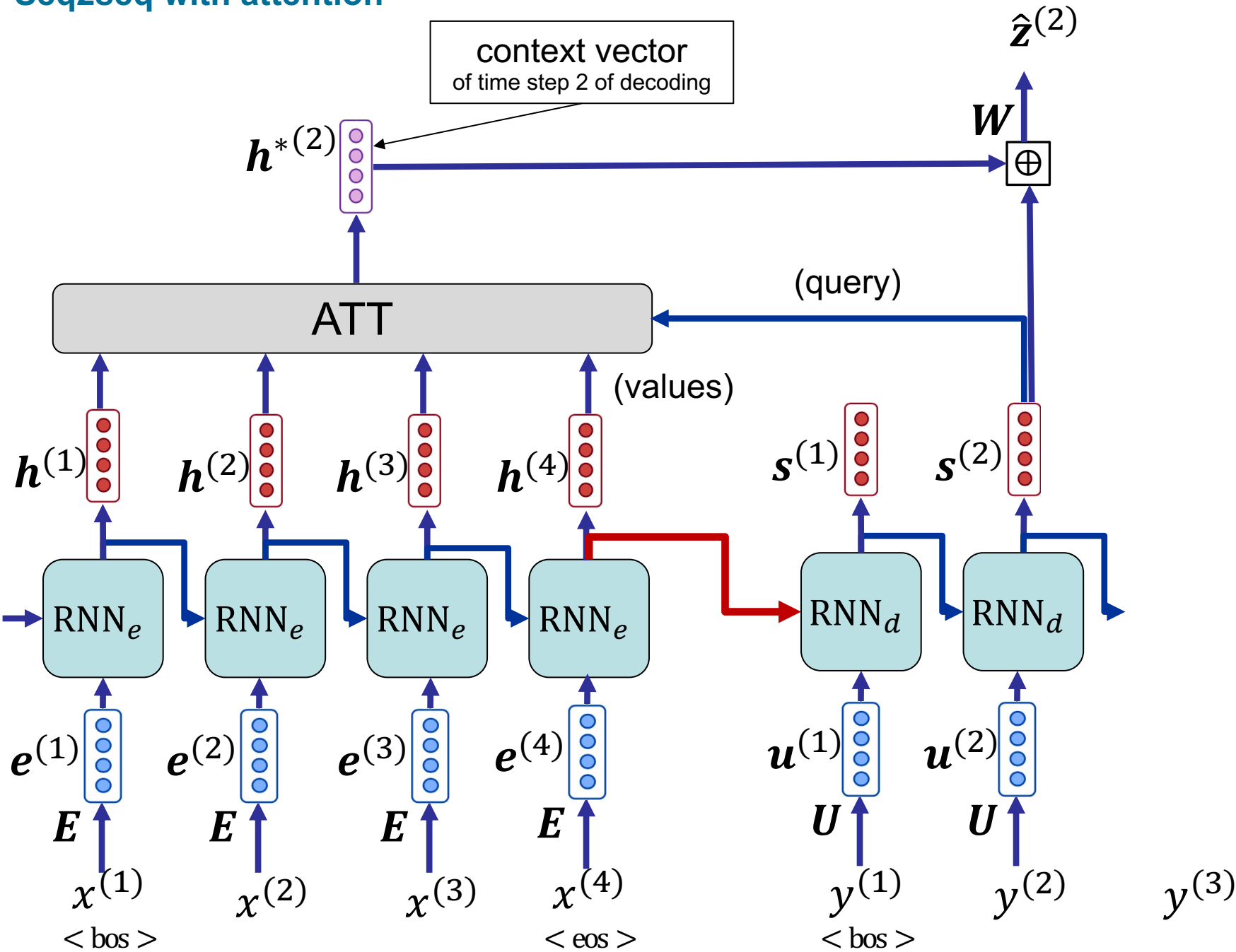
Seq2seq + Attention

- 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

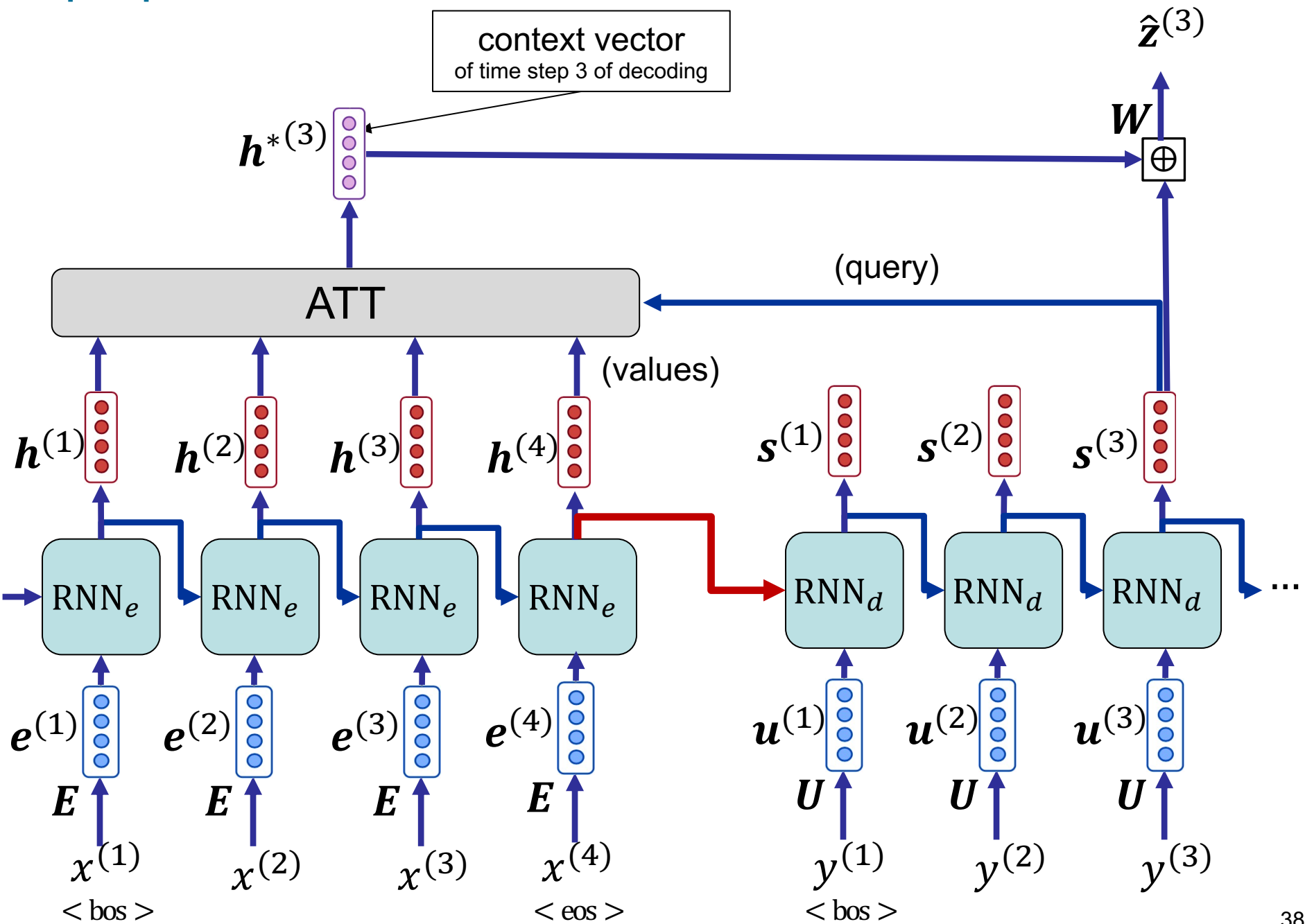
Seq2seq with attention



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

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

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

Parameters are shown in red

Seq2seq with attention – formulation

Decoder

- From words to word embeddings:
 - Decoder embeddings of target words (\mathbb{V}_d) at input $\rightarrow \textcolor{red}{U}$
 - Embedding of the target word $y^{(t)}$ at time step $t \rightarrow \textcolor{red}{u}^{(t)}$
- Decoder RNN: $\textcolor{red}{s}^{(t)} = \text{RNN}_d(\textcolor{red}{s}^{(t-1)}, \textcolor{red}{u}^{(t)})$
 - where the initial hidden state of the decoder RNN is set to the last hidden state of the encoder RNN: $\textcolor{red}{s}^{(0)} = \textcolor{red}{h}^{(L)}$

Seq2seq with attention – formulation

Decoder (cont.)

- 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 (cont.)

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

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

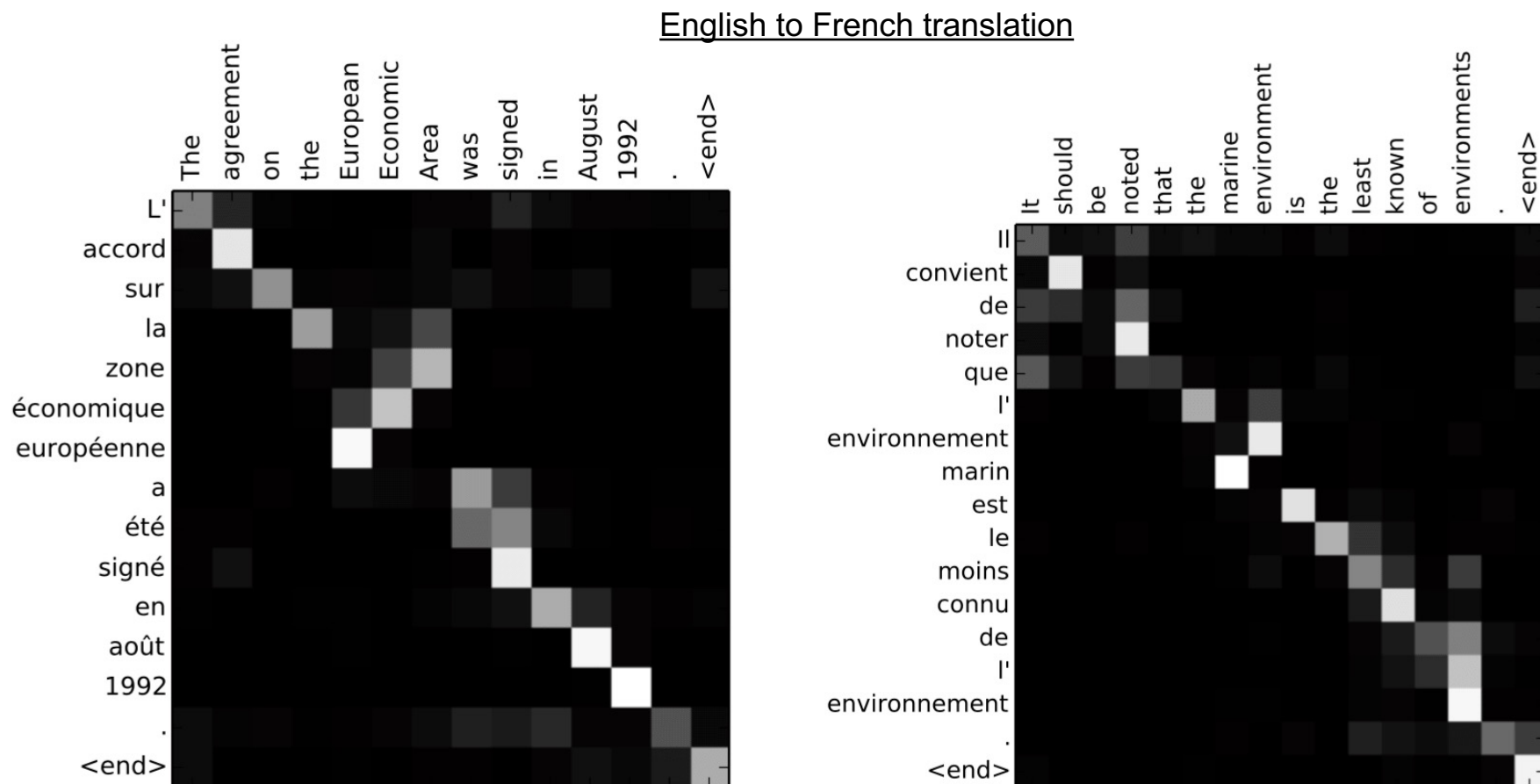
[;] denotes the concatenation of two vectors

- Probability of the next target word (at time step $t + 1$):

$$P(y^{(t+1)} | X, y^{(1)}, \dots, y^{(t-1)}, y^{(t)}) = \hat{z}_{y^{(t+1)}}^{(t)}$$

Alignment in NMT (seq2seq with attention)

- Attention automatically learns (nearly) alignment



Seq2seq with attention – summary

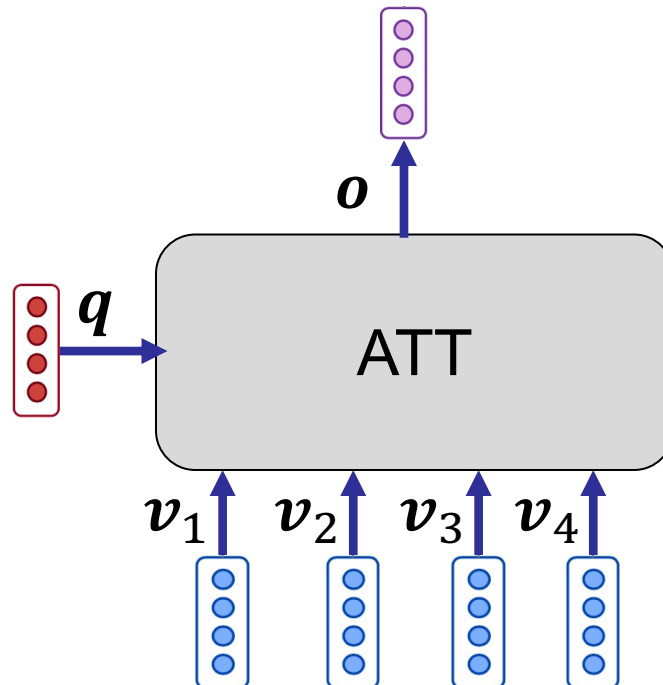
- Attention on source sequence facilitates the **focus on relevant words** and a better **flow of information**
- Adding the attention network also helps avoiding **vanishing gradient** problem by providing a shortcut to faraway states

Agenda

- Attention Networks
- Machine Translation
 - Seq2seq with Attention
- **Hierarchical document classification**

Attention in practice

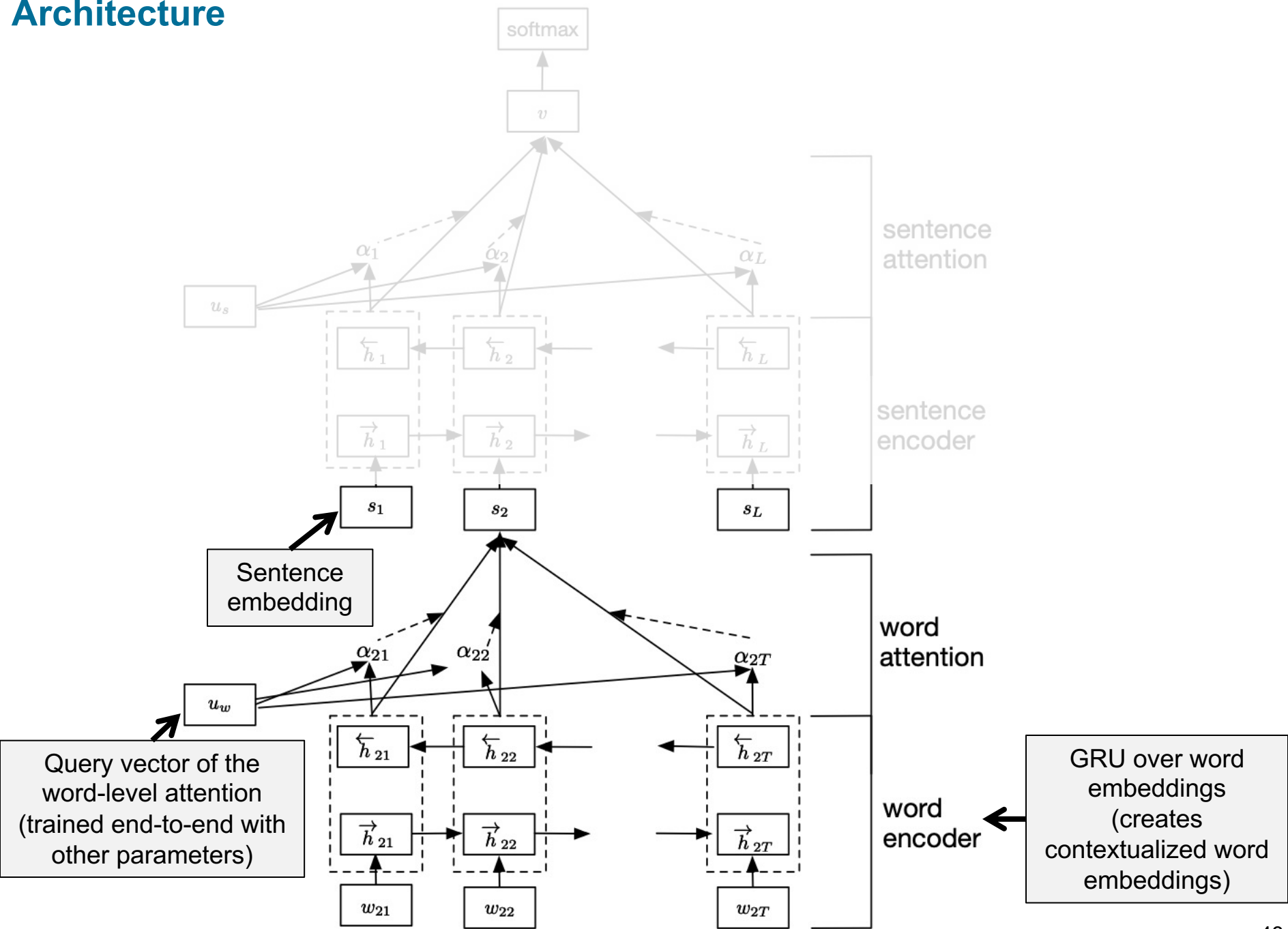
- Attention is used to create a **compositional embedding** of value vectors according to a query
 - as we already saw in seq2seq models ...
 - but it can also in tasks like sequence classification



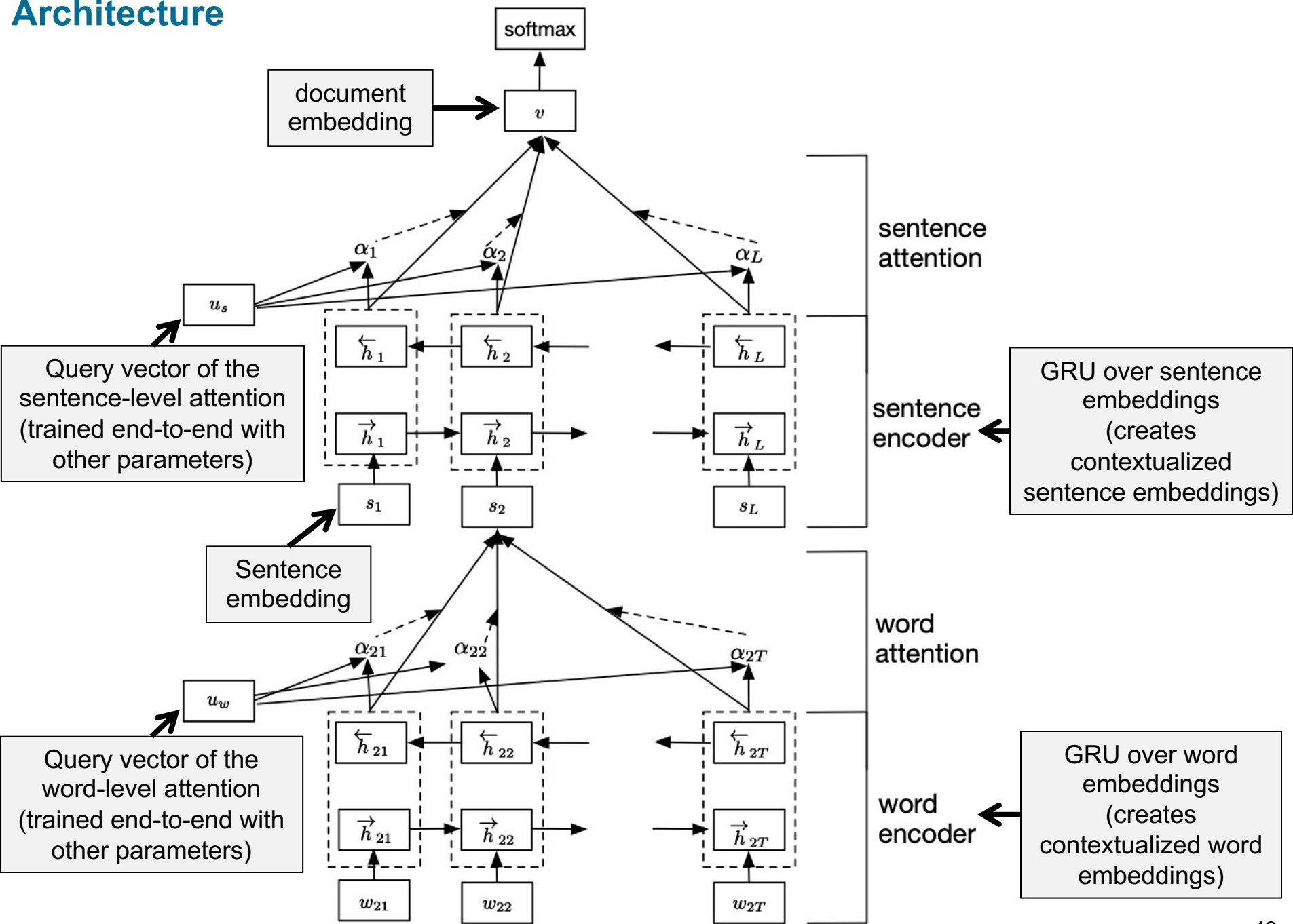
Hierarchical document classification with attention

- 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.

Sequence classification with attention – summary

- Attention can be used to compose a sequence vector from its token vectors
 - In this case, the query vector is a set of parameters that will be trained with other model parameters
 - The composed vector is in fact the weighted average of the token vectors based on attention weights
- Attention provides **some interpretability**
 - Looking at attention distributions, one may assume what the model is focusing on
 - We should however be careful about directly taking attention distributions 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/>