# Attention Is All You Need
## NeurIPS 2017

**Authors: Ashish Vaswani (Google Brain) et al.**

**Presenter: Toon Calders**

Evolutionary Tree

Open-Source
Closed-Source

transformers

Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond
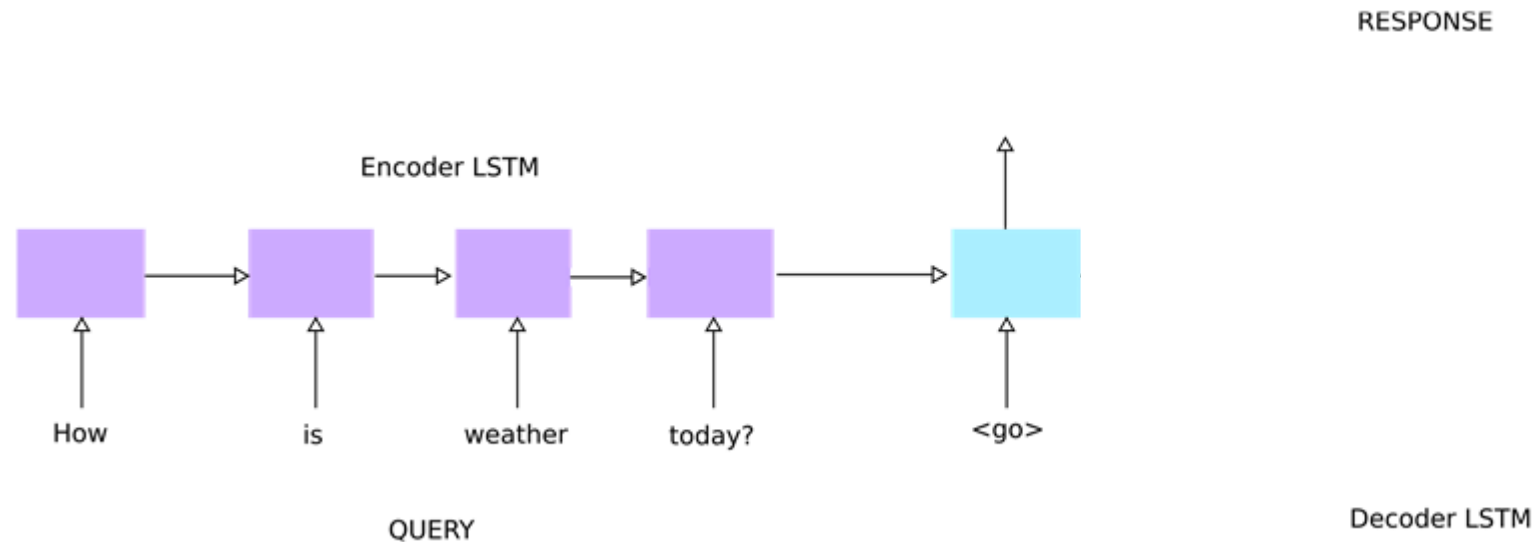
University of Antwerp
I Adrem I Adrem Data Lab

# Introduction

**Problem domain:**

- **Learning *Sequence Transduction* Models**
  - e.g. translation, Q&A

- **New model architecture: the *transformer***
  - *Encoder-decoder* architecture
  - Systematic use of *attention mechanism*
  - No convolutions nor recurrence

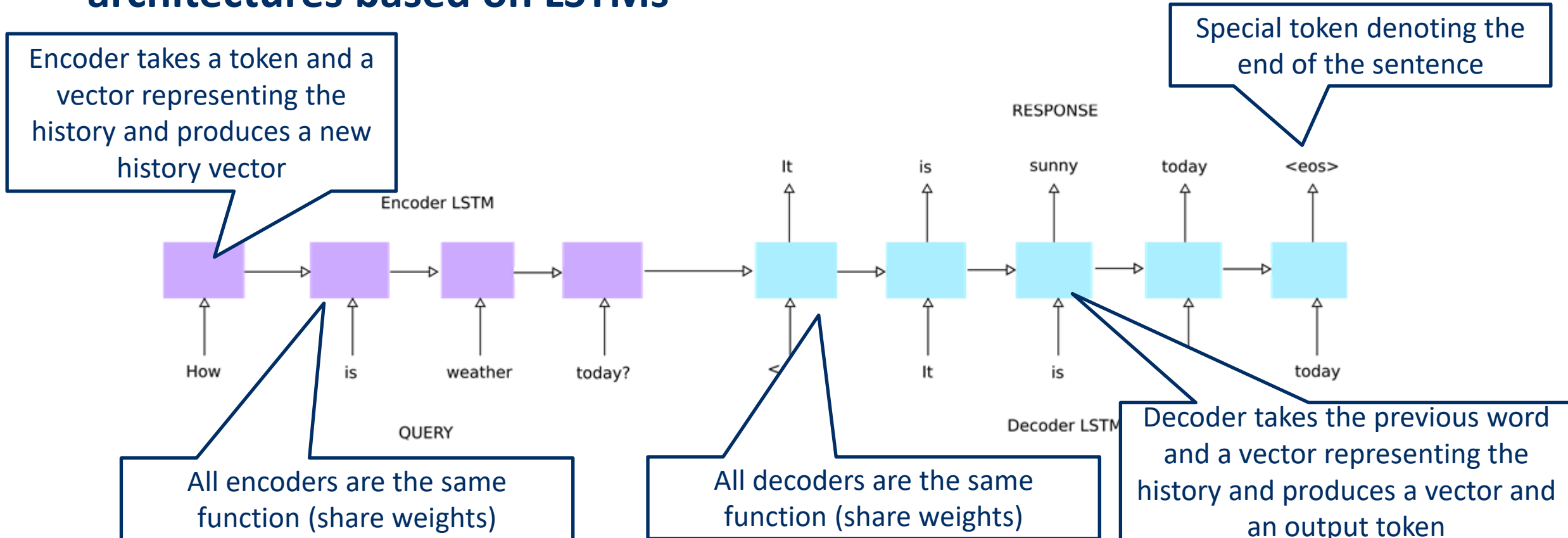- **Good performance on *Machine Translation* tasks**

# "Old" Sequence-to-Sequence Models

- **Common approach before transformers: Recurrent Neural Network architectures based on LSTMs***

RESPONSE

Encoder LSTM

How          is          weather          today?          <go>

QUERY

Decoder LSTM

\* LSTM: Long short term memory

# "Old" Sequence-to-Sequence Models

- **Common approach before transformers: Recurrent Neural Network architectures based on LSTMs***

Encoder takes a token and a vector representing the history and produces a new history vector

Special token denoting the end of the sentence

All encoders are the same function (share weights)

All decoders are the same function (share weights)

Decoder takes the previous word and a vector representing the history and produces a vector and an output token

RESPONSE

Encoder LSTM

Decoder LSTM

QUERY

It    is    sunny    today    <eos>

How    is    weather    today?

<    It    is    today

* LSTM: Long short term memory

University of Antwerp | Adrem | Adrem Data Lab

# "Old" Sequence – To – Sequence Models
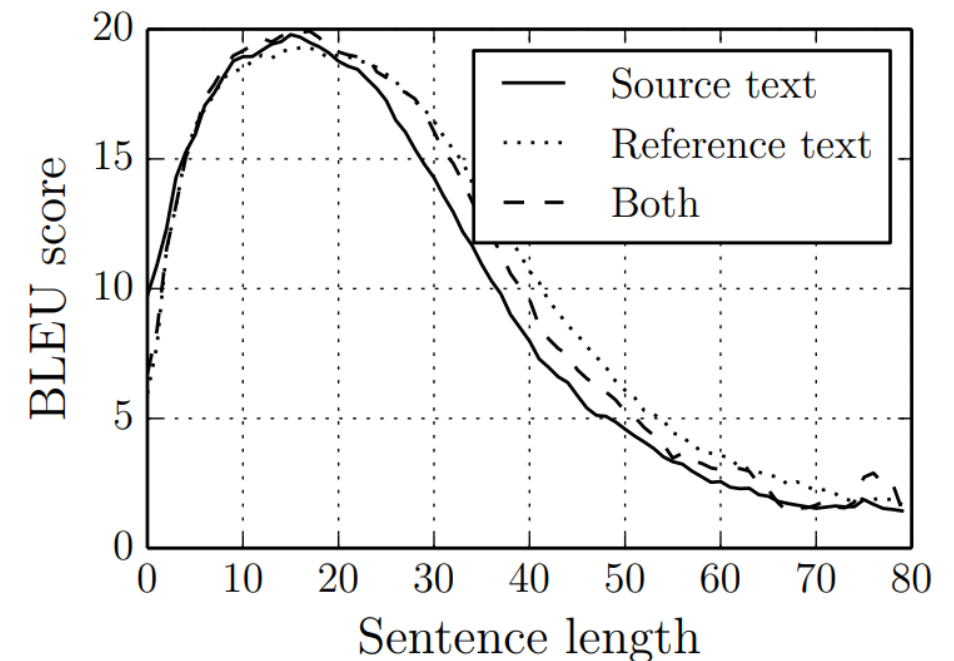
- **Common techniques:**
  - Recurrent Neural Networks based on LSTMs
  - Convolutional techniques

- **Disadvantages:**
  - Long dependencies hard to capture
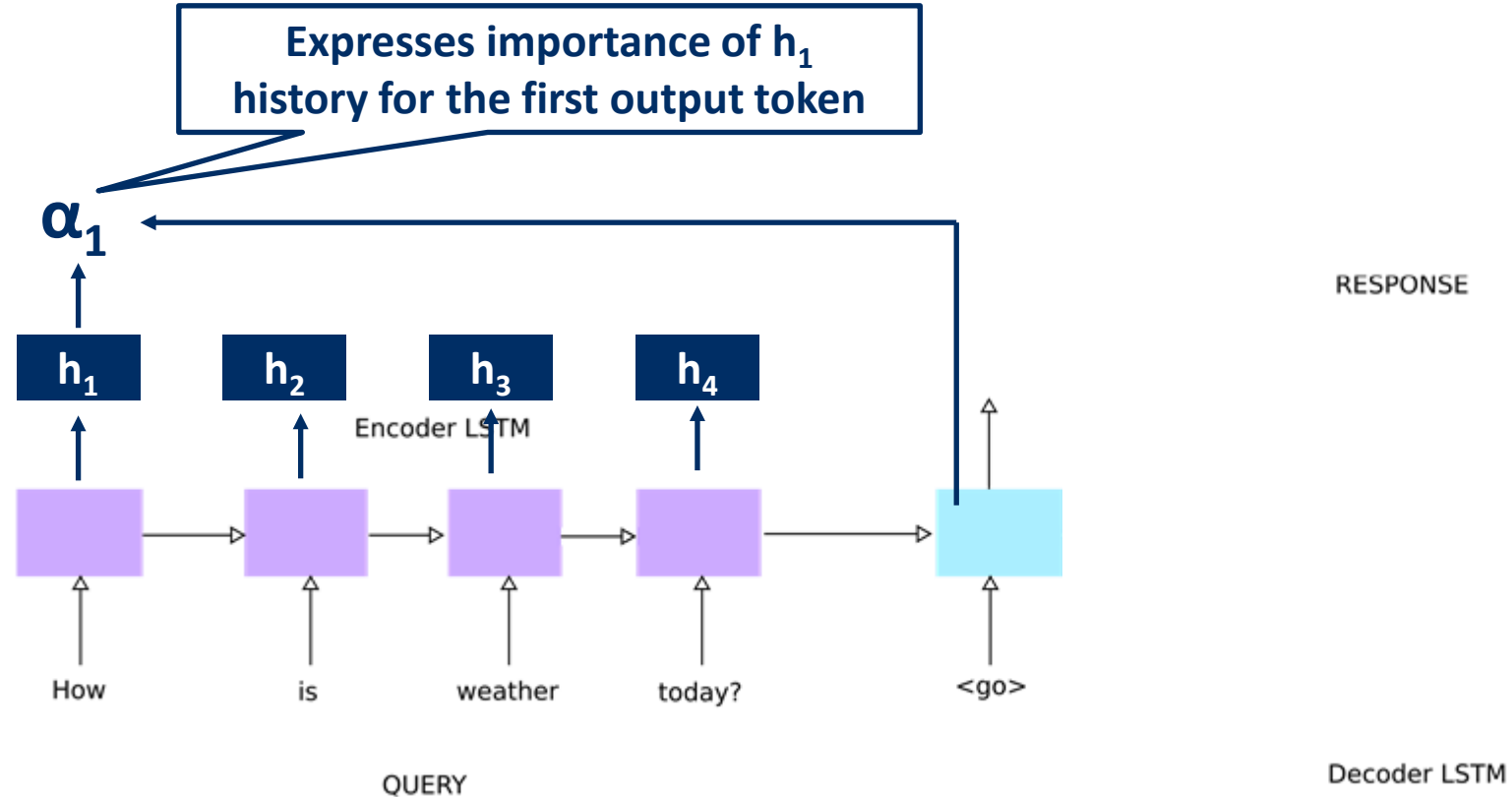  - Limited possibility for parallellism

- **Solution for long dependencies:**
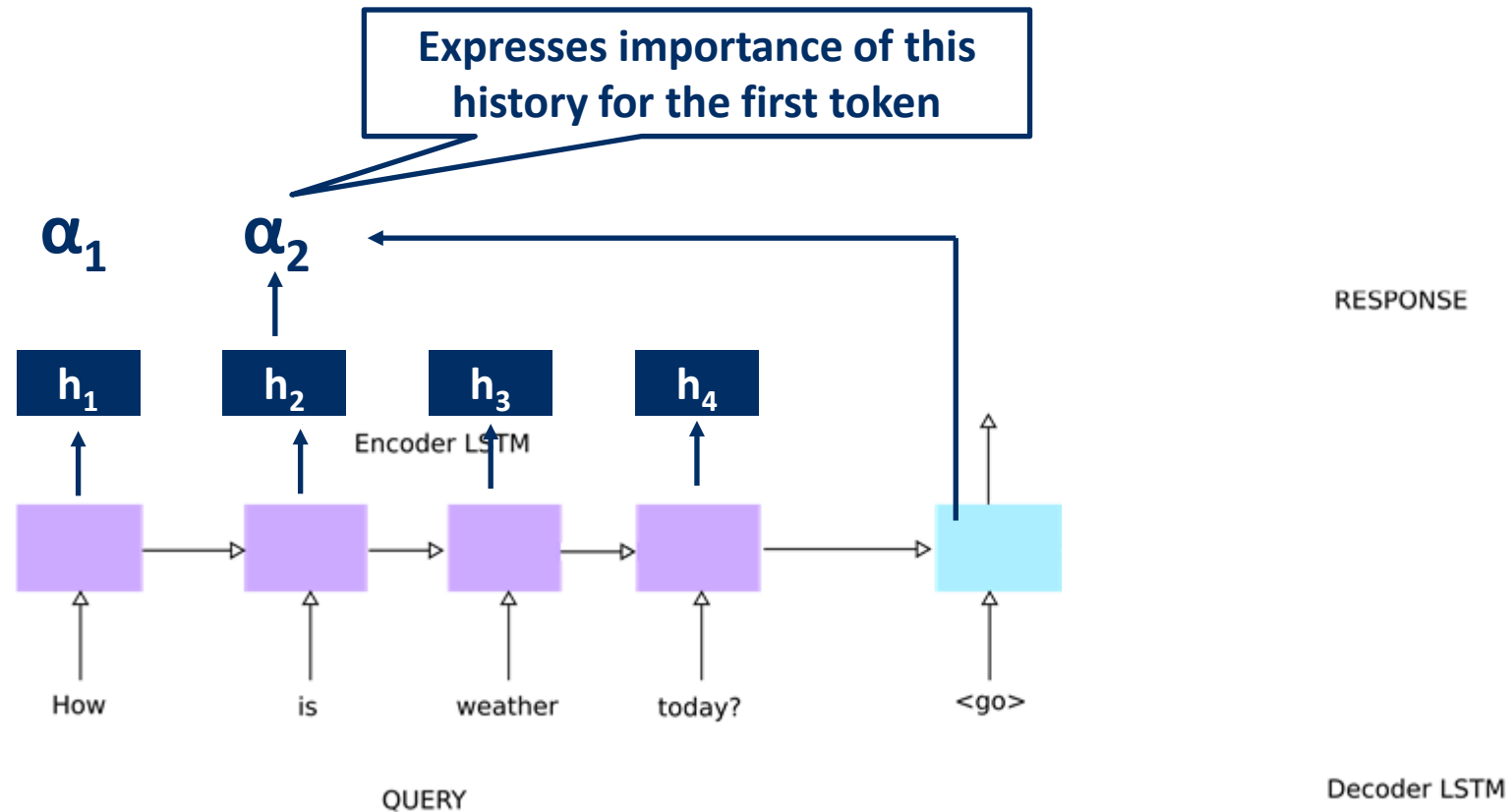  - Attention mechanisms



Cho, K., van Merrienboer, B., Bahdanau, D., and Bengio, Y. (2014b). On the properties of neural ¨ machine translation: Encoder–Decoder approaches. In Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation.
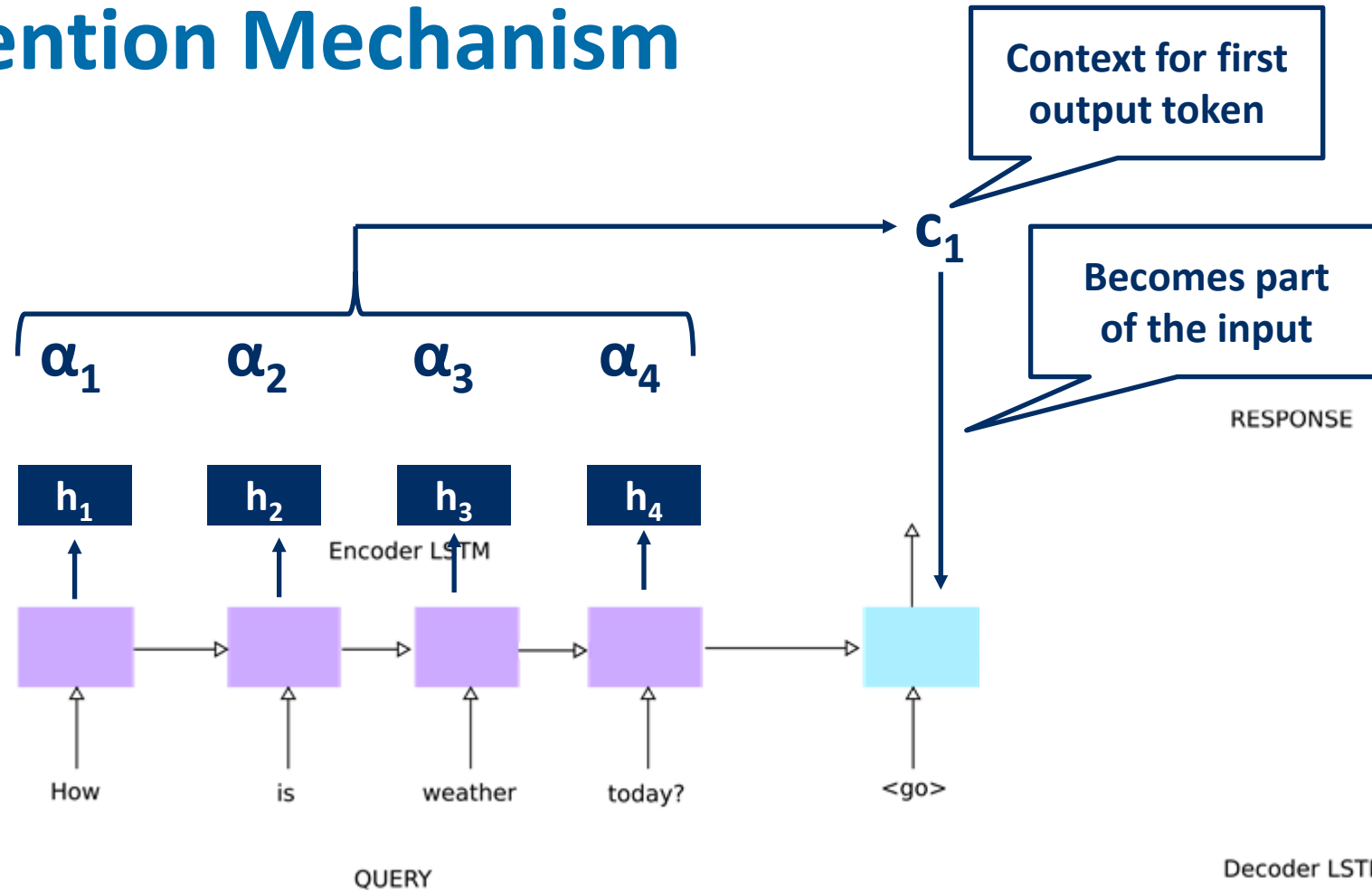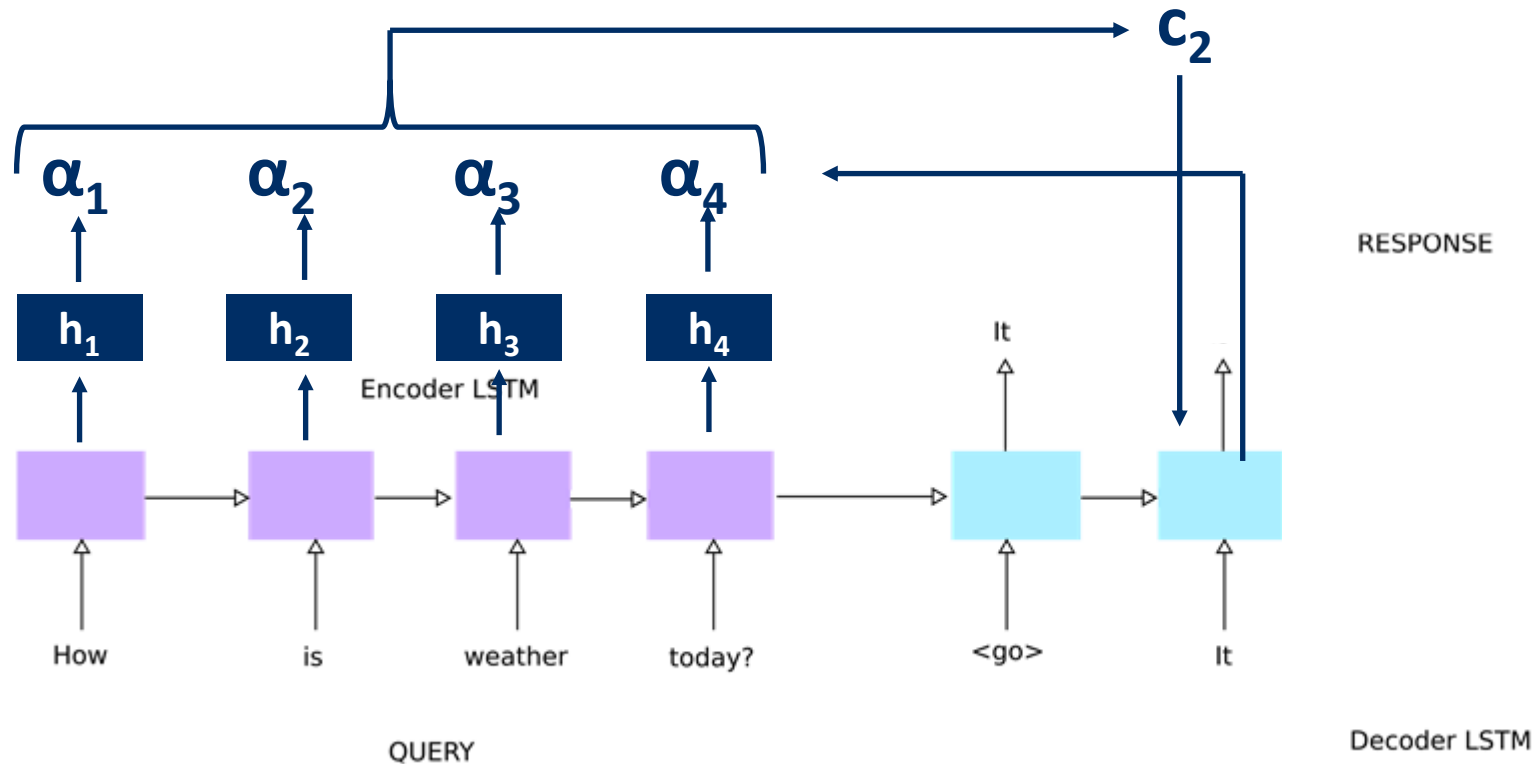
University of Antwerp
Adrem | Adrem Data Lab
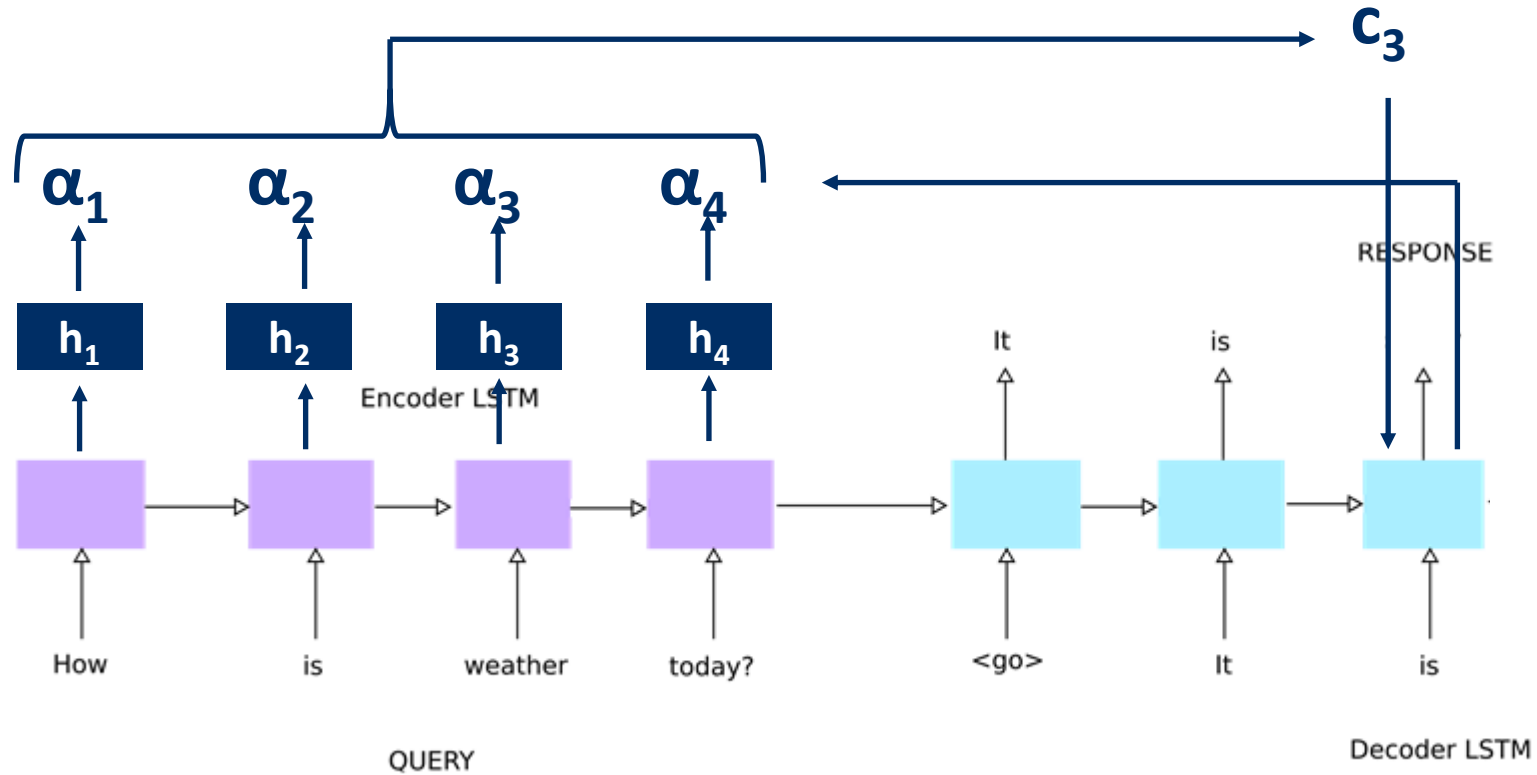
# Attention Mechanism



Expresses importance of $h_1$ history for the first output token

$\alpha_1$

$h_1$  $h_2$  $h_3$  $h_4$

Encoder LSTM

How    is    weather    today?    <go>

QUERY

RESPONSE

Decoder LSTM

Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *3rd International Conference on Learning Representations, ICLR 2015.* 2015.

University of Antwerp
I Adrem I Adrem Data Lab

# Attention Mechanism



Expresses importance of this history for the first token

$\alpha_1$   $\alpha_2$

| $h_1$ | $h_2$ | $h_3$ | $h_4$ |

Encoder LSTM

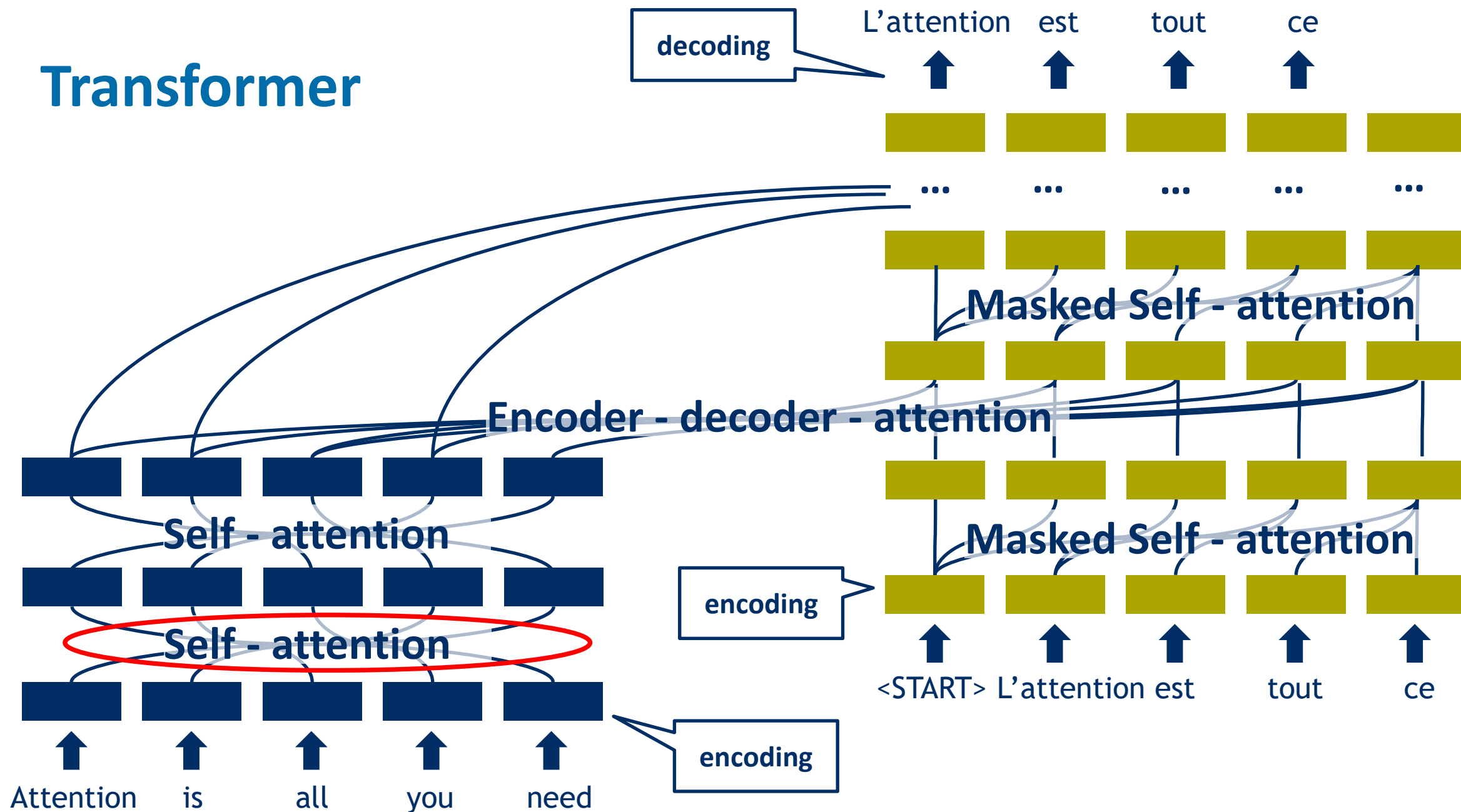How   is   weather   today?   <go>

QUERY

RESPONSE

Decoder LSTM

Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *3rd International Conference on Learning Representations, ICLR 2015*. 2015.

University of Antwerp
I Adrem I Adrem Data Lab

# Attention Mechanism



Context for first output token

Becomes part of the input

$c_1$

RESPONSE

$\alpha_1$    $\alpha_2$    $\alpha_3$    $\alpha_4$

$h_1$    $h_2$    $h_3$    $h_4$

Encoder LSTM

How    is    weather    today?    <go>

QUERY

Decoder LSTM

Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *3rd International Conference on Learning Representations, ICLR 2015*. 2015.
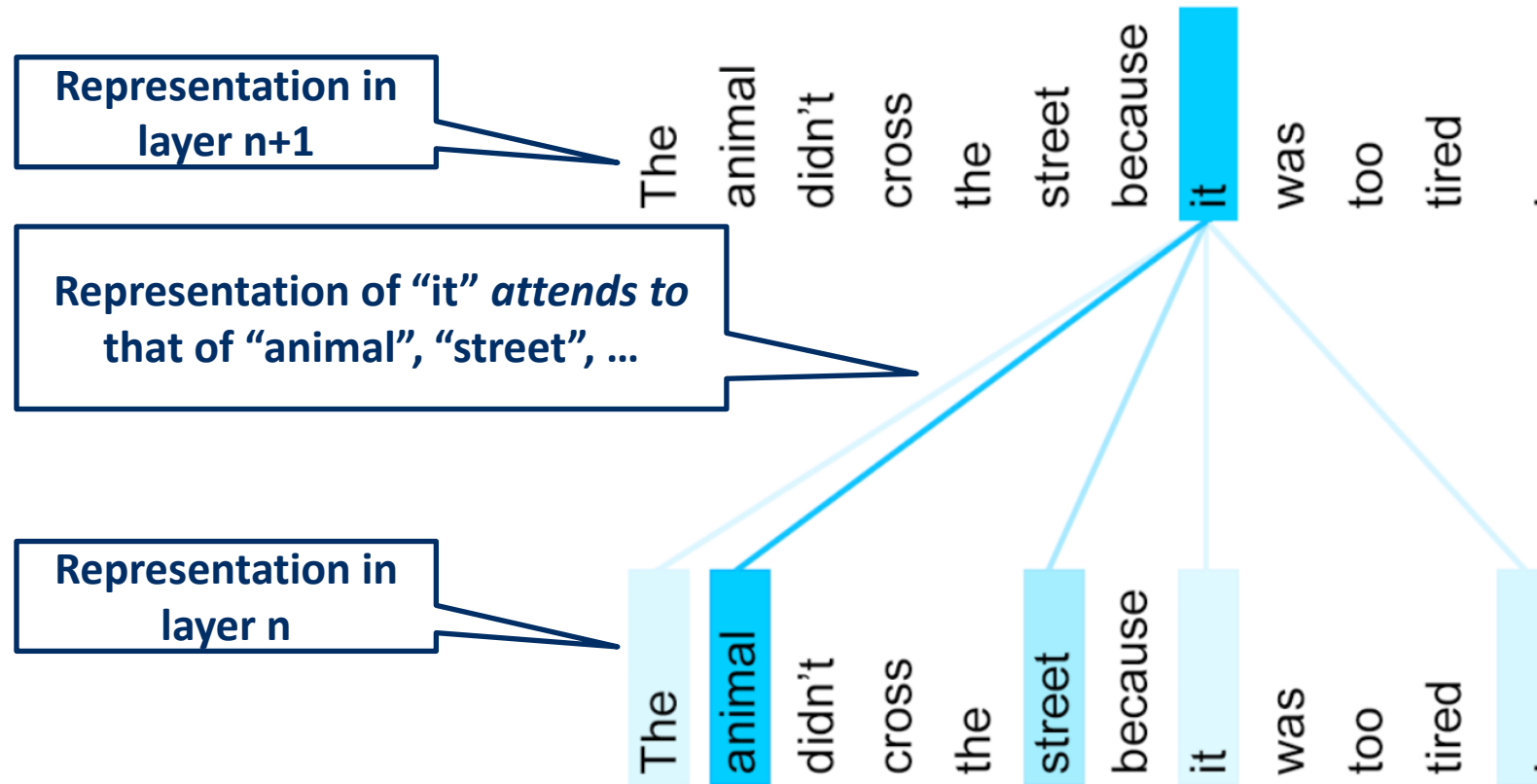
# Attention Mechanism



Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *3rd International Conference on Learning Representations, ICLR 2015*. 2015.

# Attention Mechanism



Bahdanau, Dzmitry, Kyung Hyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *3rd International Conference on Learning Representations, ICLR 2015*. 2015.

# New Model Proposed in the Paper

- Encoder – Decoder structure

- Each layer has an encoding for each token
  - "Enriched" representation in next layer computed on previous layer

- Layers are connected using attention mechanism
  - No convolutional layers, recurrent neural nets, LSTM
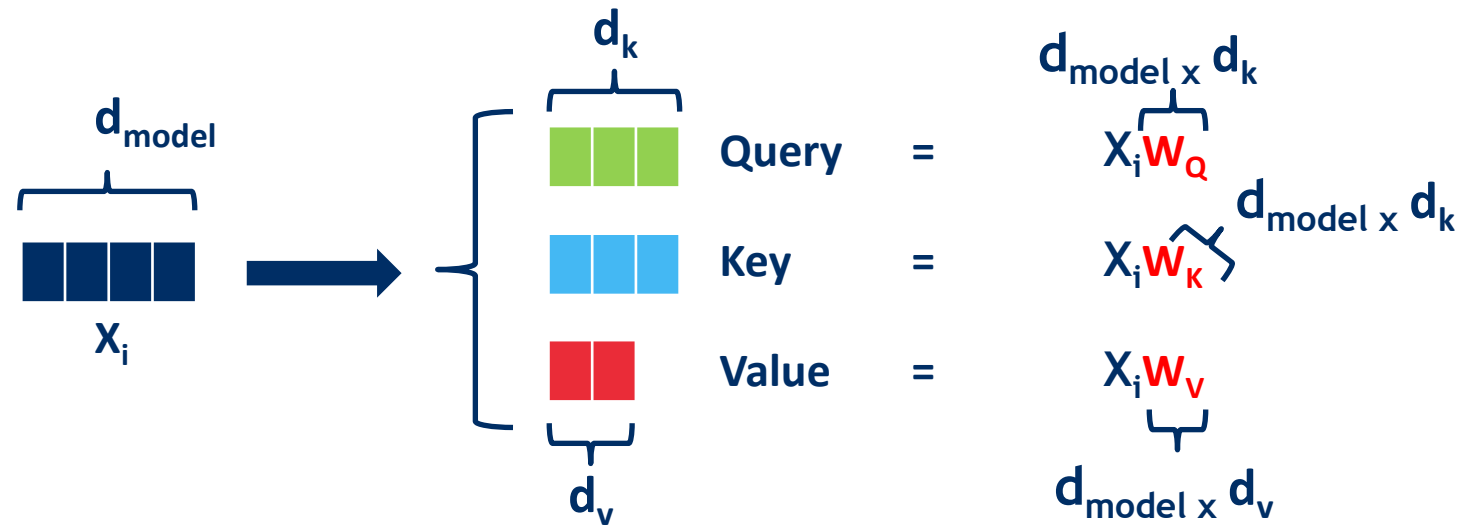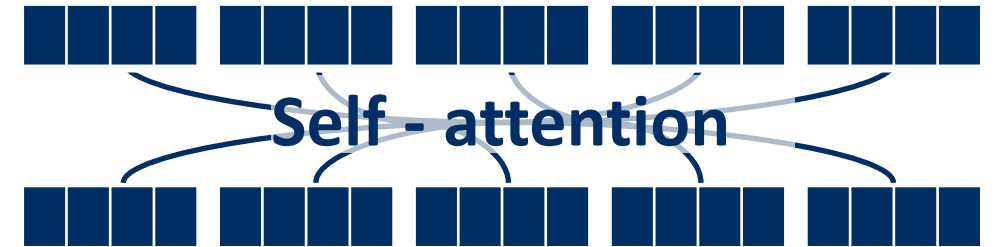  - Hence: Attention is all you need!

University of Antwerp
I Adrem I Adrem Data Lab

# Transformer

# Attention Mechanism: Intuition



**Representation in layer n+1**

**Representation of "it" *attends to* that of "animal", "street", …**

**Representation in layer n**

The animal didn't cross the street because **it** was too tired .

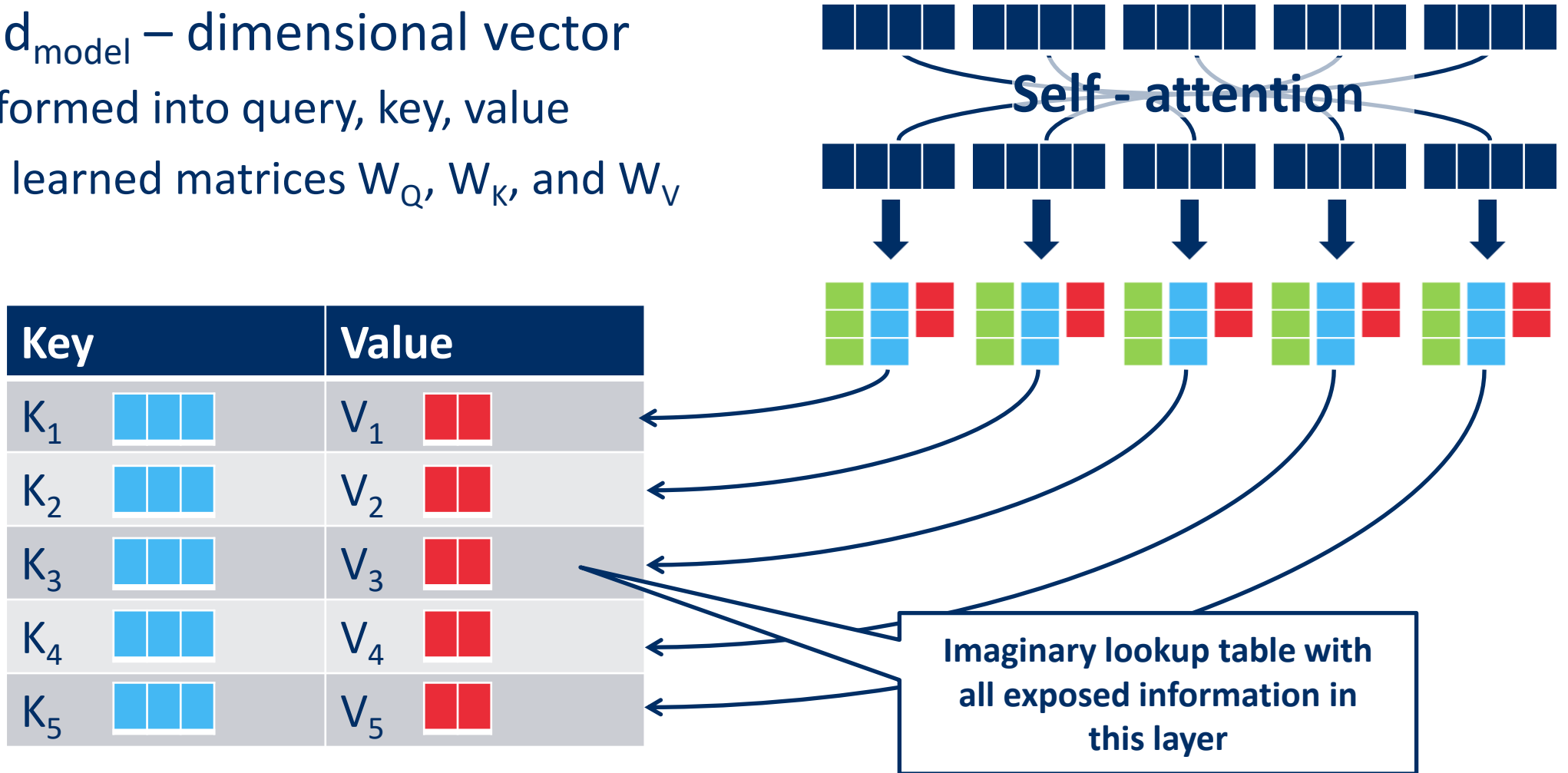The **animal** didn't cross the **street** because it was too tired .
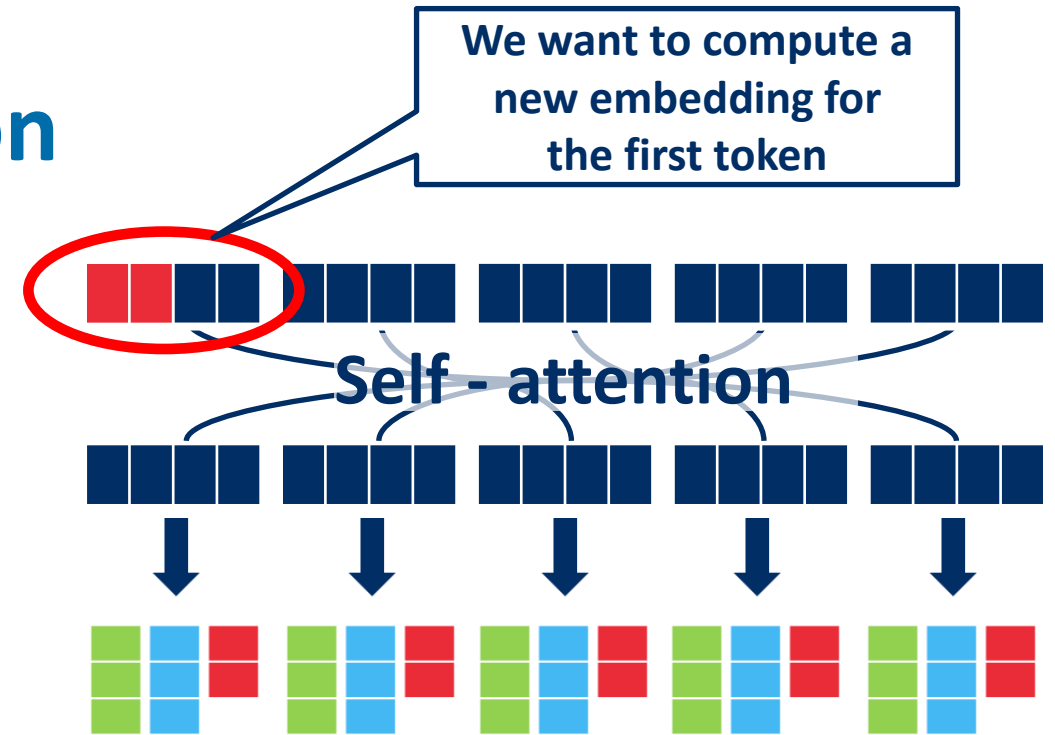
# "Scaled Dot-Product" Attention

- ▪ ■■■■ : $d_{model}$ – dimensional vector
  - ▪ Transformed into query, key, value
  - ▪ Using learned matrices $W_Q$, $W_K$, and $W_V$

**Self - attention**

$d_{model}$

$X_i$

$d_k$

Query $= X_i W_Q$ $\quad d_{model} \times d_k$

Key $= X_i W_K$ $\quad d_{model} \times d_k$

Value $= X_i W_V$

$d_v$

$d_{model} \times d_v$

# "Scaled Dot-Product" Attention

- ▪ ▪▪▪▪ : $d_{model}$ – dimensional vector
  - ▪ Transformed into query, key, value
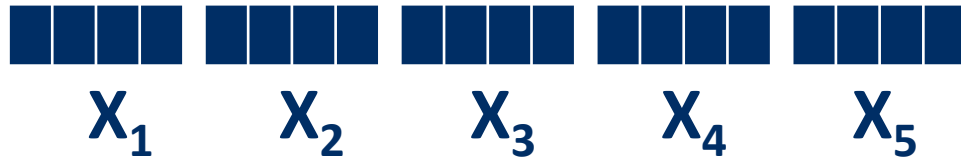  - ▪ Using learned matrices $W_Q$, $W_K$, and $W_V$

**Self - attention**

| Key | | Value | |
|-----|---|-------|---|
| $K_1$ | ▪▪▪ | $V_1$ | ▪▪ |
| $K_2$ | ▪▪▪ | $V_2$ | ▪▪ |
| $K_3$ | ▪▪▪ | $V_3$ | ▪▪ |
| $K_4$ | ▪▪▪ | $V_4$ | ▪▪ |
| $K_5$ | ▪▪▪ | $V_5$ | ▪▪ |

**Imaginary lookup table with all exposed information in this layer**

# "Scaled Dot-Product" Attention

- ▪ ▮▮▮▮ : $d_{model}$ – dimensional vector
  - ▪ Transformed into query, key, value
  - ▪ Using learned matrices $W_Q$, $W_K$, and $W_V$

We want to compute a new embedding for the first token

**Self - attention**

| Key | | Value | |
|-----|-----|-------|-----|
| $K_1$ | ▮▮▮ | $V_1$ | ▮▮ |
| $K_2$ | ▮▮▮ | $V_2$ | ▮▮ |
| $K_3$ | ▮▮▮ | $V_3$ | ▮▮ |
| $K_4$ | ▮▮▮ | $V_4$ | ▮▮ |
| $K_5$ | ▮▮▮ | $V_5$ | ▮▮ |

Compare the query of the first token to all keys; compute similarity and take a weighted average
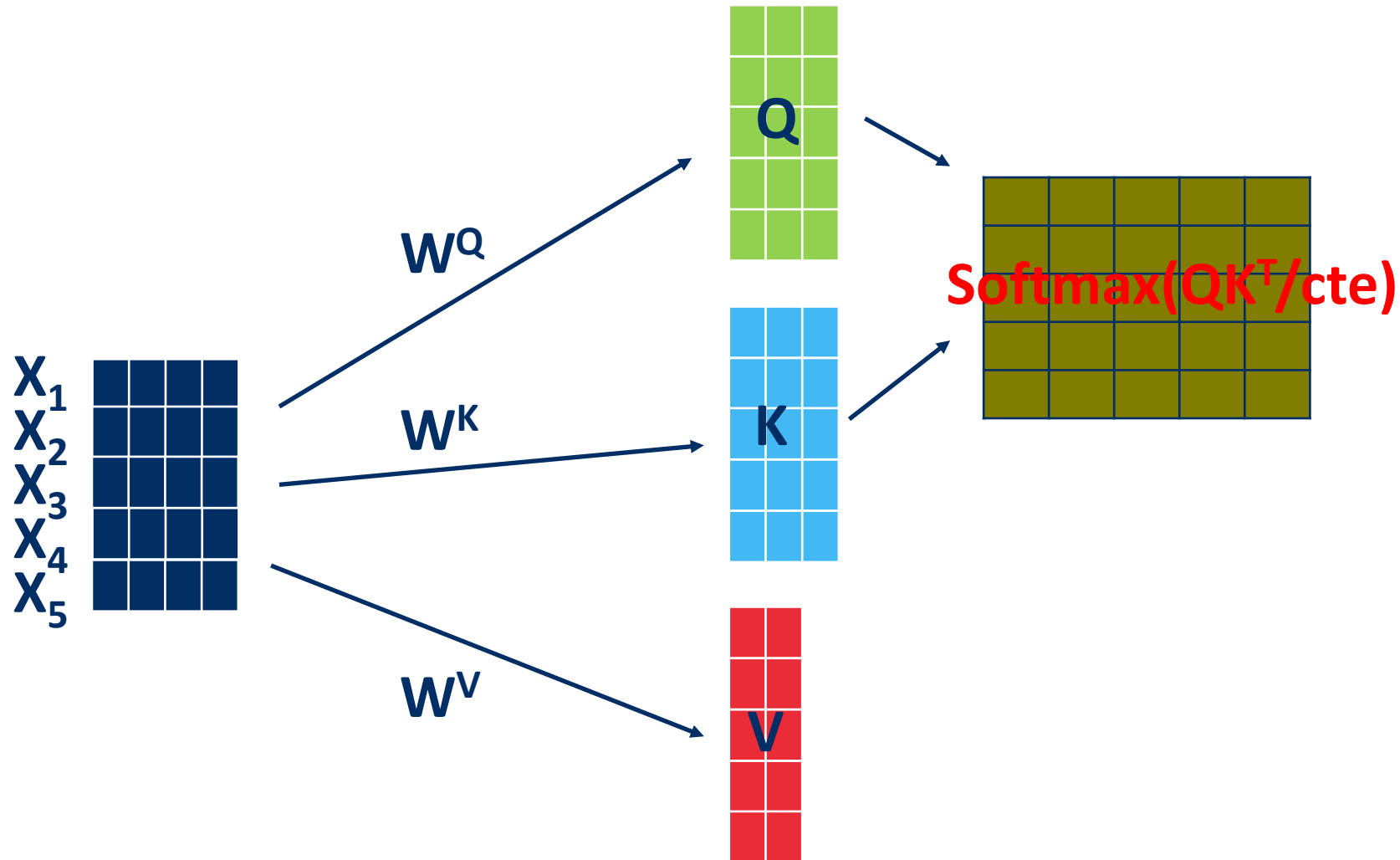
# "Scaled Dot-Product" : Implementation

$$X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5$$

# "Scaled Dot-Product" : Implementation

# "Scaled Dot-Product" : Implementation



$$Sim(Q,K) = \Sigma_i Q_i K_i = QK^T$$

# "Scaled Dot-Product" : Implementation



$W^Q$

$W^K$

$W^V$

Q

K

V

$X_1$
$X_2$
$X_3$
$X_4$
$X_5$

Softmax(QK$^T$/cte)

University of Antwerp
I Adrem I Adrem Data Lab
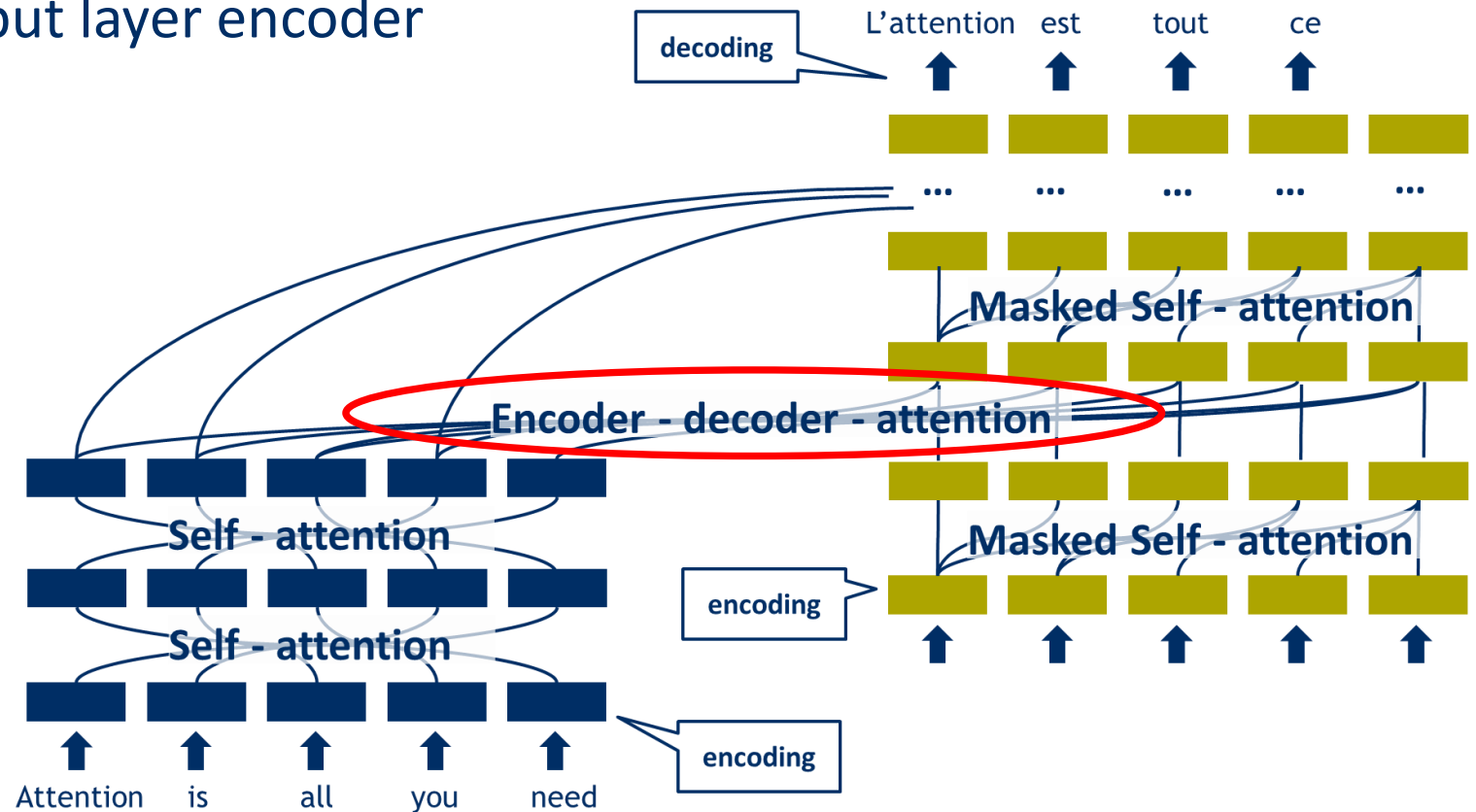
# "Scaled Dot-Product" : Implementation
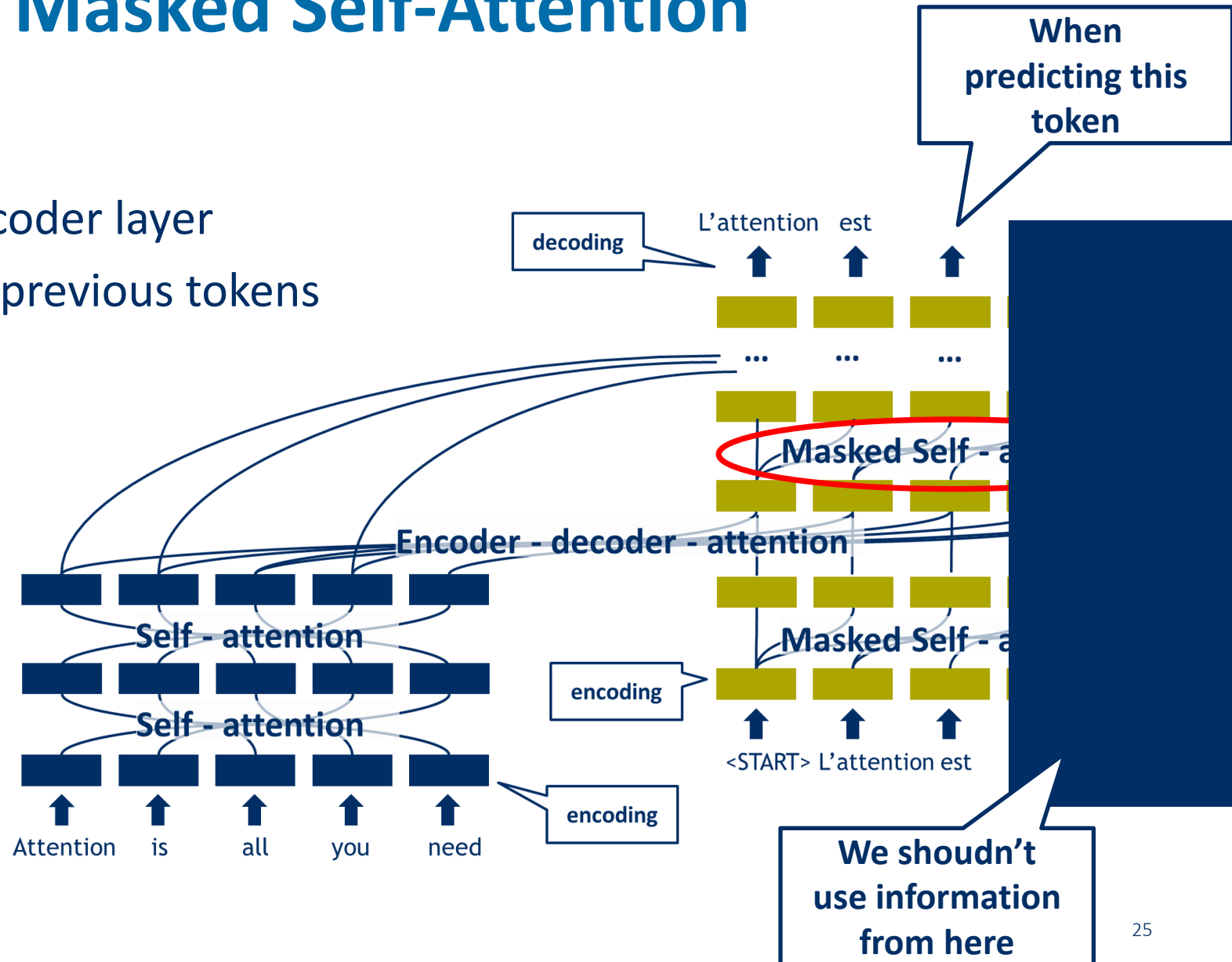
# Multi-Headed Attention

# Encoder – Decoder attention

- Very similar construction
  - key, value come from output layer encoder
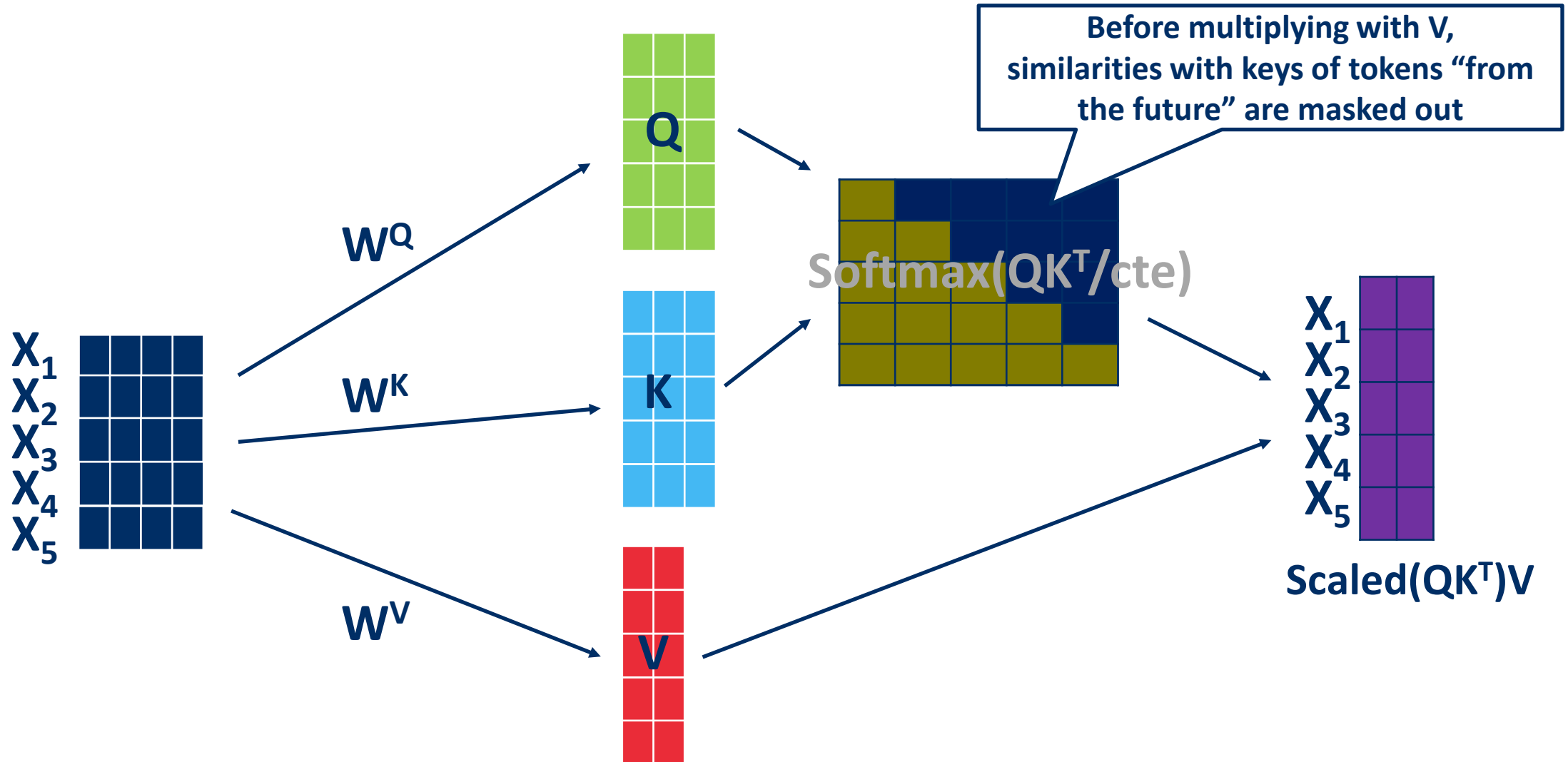  - query from decoder layer
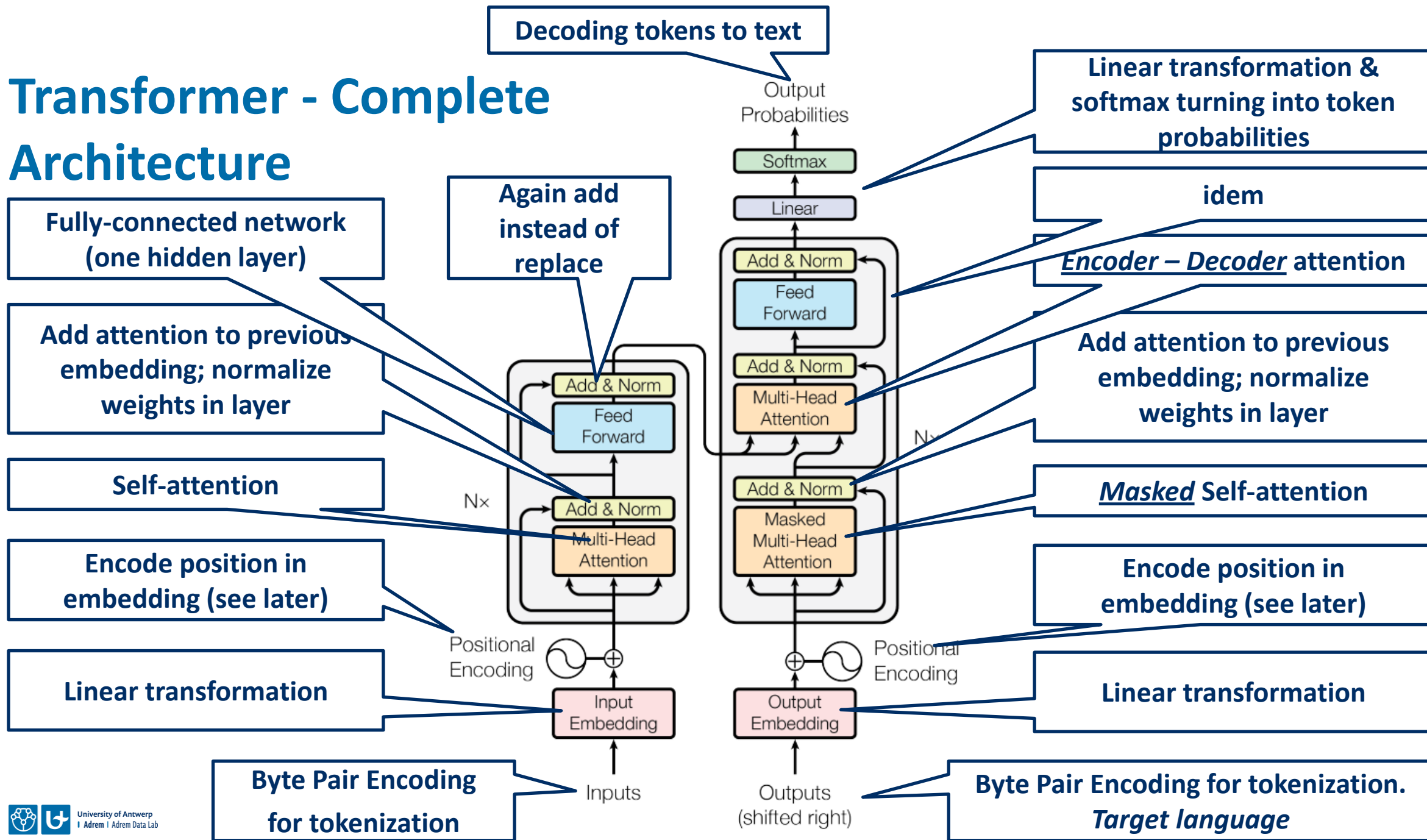
# Decoder Attention: Masked Self-Attention

- Again, very similar
  - Key, value, query from decoder layer
  - Tokens can only attend to previous tokens



When predicting this token

decoding

L'attention est

...    ...    ...

Masked Self - a

Encoder - decoder - attention

Self - attention

Self - attention

encoding

Masked Self - a

encoding

<START> L'attention est

Attention  is  all  you  need

We shoudn't use information from here

University of Antwerp
Adrem | Adrem Data Lab

# Masked Self-Attention

# Transformer - Complete Architecture

Decoding tokens to text

Linear transformation & softmax turning into token probabilities

idem

*Encoder – Decoder* attention

Again add instead of replace

Fully-connected network (one hidden layer)

Add attention to previous embedding; normalize weights in layer

Add attention to previous embedding; normalize weights in layer

Self-attention

*Masked* Self-attention

Encode position in embedding (see later)

Encode position in embedding (see later)

Linear transformation

Linear transformation

Byte Pair Encoding for tokenization

Byte Pair Encoding for tokenization. *Target language*

# Byte Pair Encoding

- Turn text into numbers
  - fixed (subword-)tokenization scheme
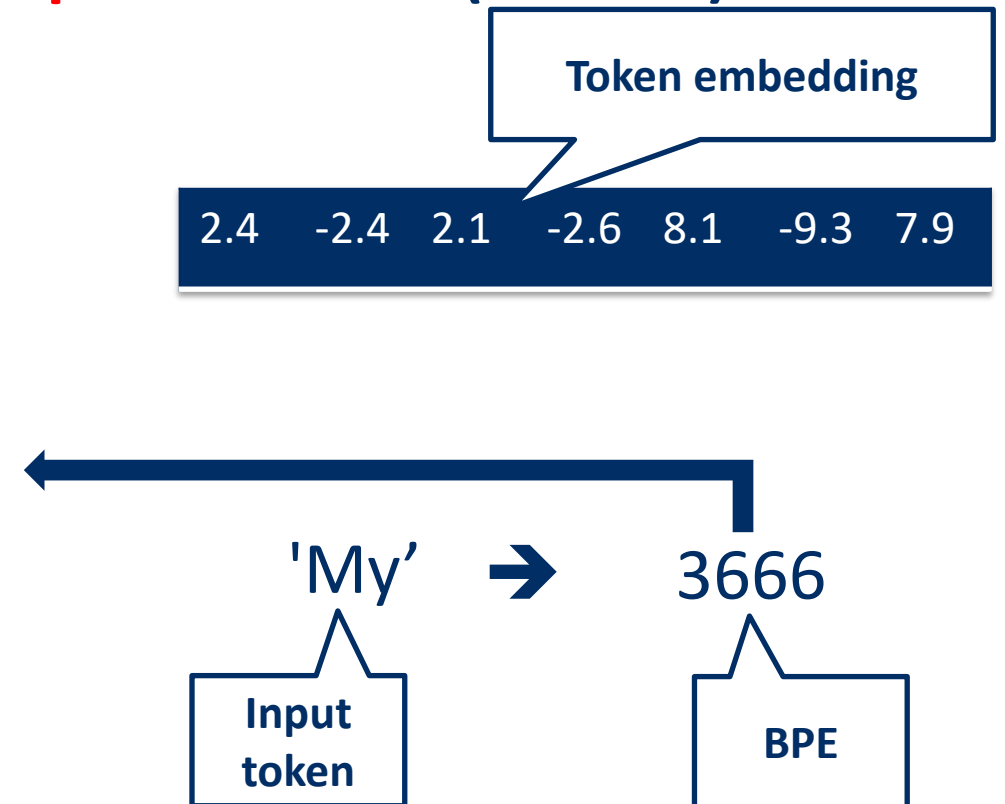  - Every sequence of characters can be tokenized

**Example:**

My name is Toon. ➔ ['My', ' name', ' is', ' To', 'on', '.']

[3666, 1438, 318, 1675, 261, 13]

My name  is  Toon! ➔ ['My', ' name', ' ', ' is', ' ', ' To', 'on', '!']

[3666, 1438, 220, 318, 220, 1675, 261, 0]

University of Antwerp
Adrem | Adrem Data Lab

# Input Embedding

- **BPE has 37K tokens**
- **37K dim. 1-hot encoding ➔ lower-dim. representation (trained)**

| | | | | | | |
|---|---|---|---|---|---|---|
| -1.5 | 3.3 | -2.1 | 8.5 | 2.1 | -2.6 | 8.1 |
| 0.3 | -0.2 | 2.4 | -2.1 | 6.6 | 2.1 | 7.2 |

**Token embedding**

| 2.4 | -2.4 | 2.1 | -2.6 | 8.1 | -9.3 | 7.9 |
|---|---|---|---|---|---|---|

| | | | | | | |
|---|---|---|---|---|---|---|
| 2.1 | -2.6 | 8.1 | 3.3 | -2.1 | 8.6 | 2.1 |
| 2.4 | -2.4 | 2.1 | -2.6 | 8.1 | -9.3 | 7.9 |
| 3.4 | -2.1 | 6.6 | 2.1 | -4.2 | 1.3 | 9.2 |

| | | | | | | |
|---|---|---|---|---|---|---|
| -3.4 | 2.1 | -2.6 | 8.1 | -2.3 | -2.2 | 3.2 |
| 3.4 | -2.1 | 8.6 | 2.1 | 2.1 | -2.6 | 8.1 |

1
2
...
3,665
3,666
3,667
...
36,999
37,000

'My' ➔ 3666

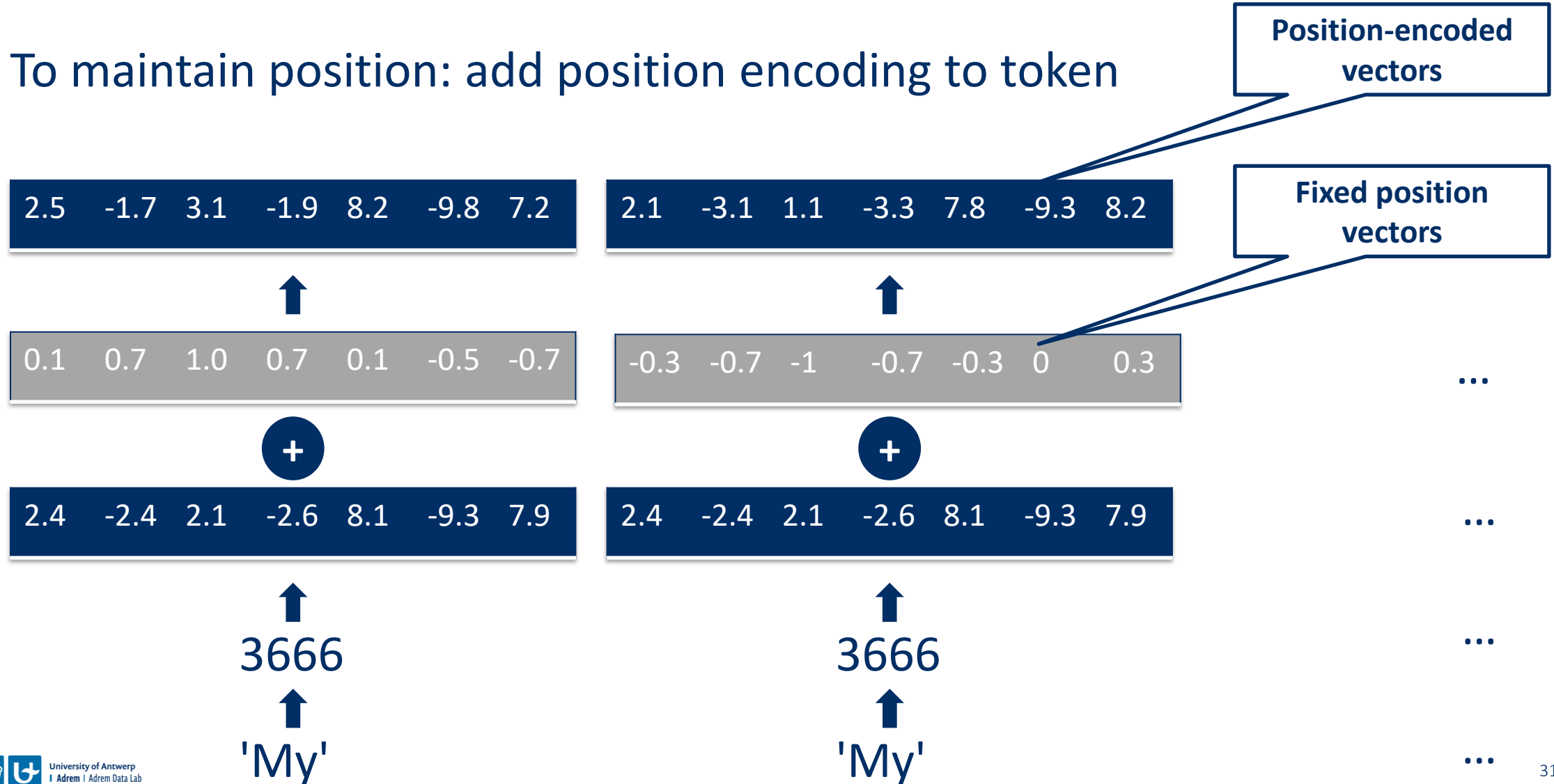**Input token**

**BPE**

# Positional Encoding

- Order of tokens is important:
  - Woman, without her **man,** is nothing.
  - Woman, without her**, man** is nothing.
- BPE:
  - [48081, 11, 1231, 607, **582**, **11**, 318, 2147, 13]
  - [48081, 11, 1231, 607, **11**, **582**, 318, 2147, 13]
- Gives *the same* key-value pairs
  - Is token "Woman" subject of the sentence or not?
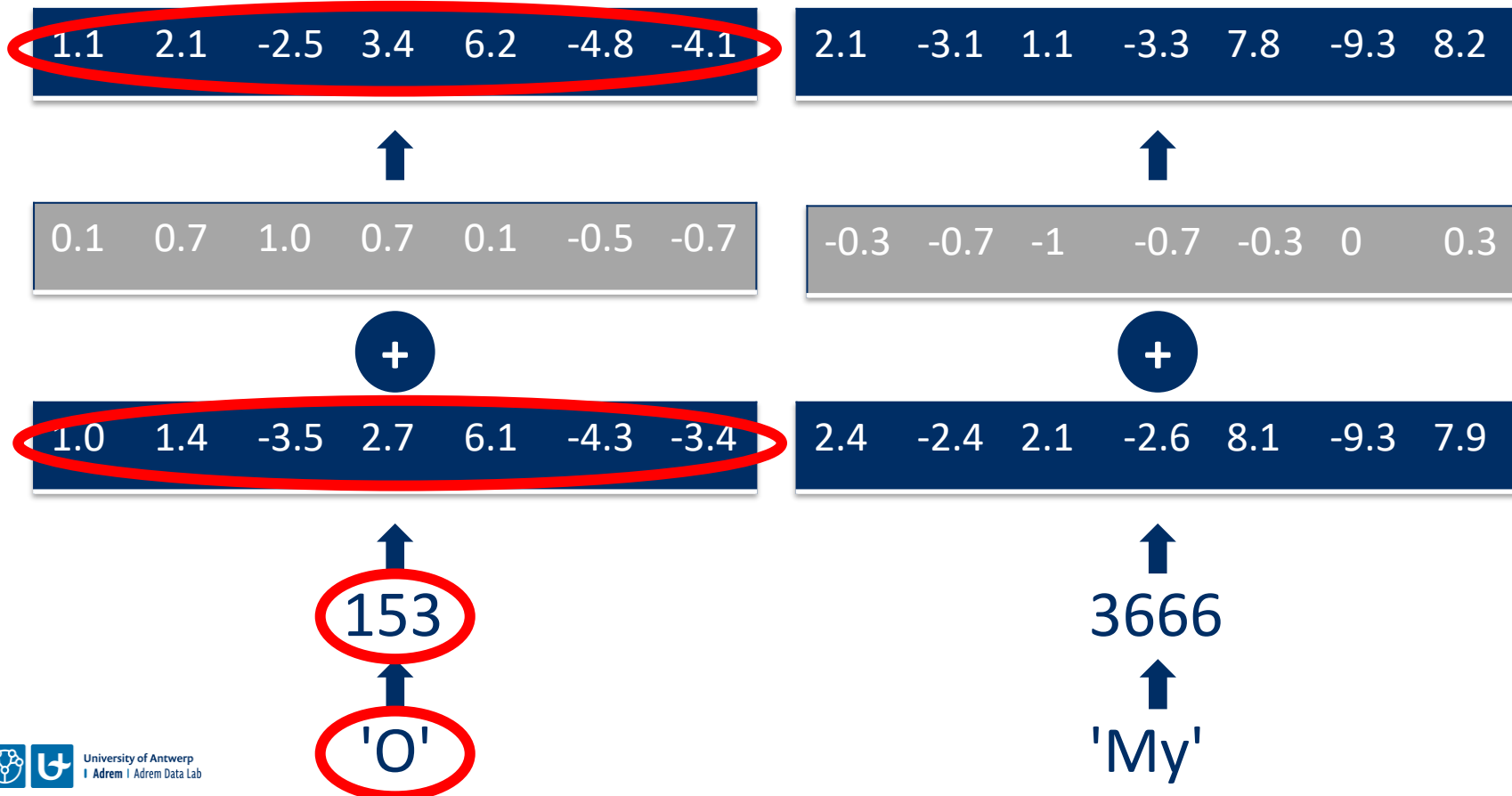  - Position of comma is essential!

University of Antwerp
I Adrem I Adrem Data Lab

# Positional Encoding

- To maintain position: add position encoding to token

| Position-encoded vectors |
| --- |

| Fixed position vectors |
| --- |

| 2.5 | -1.7 | 3.1 | -1.9 | 8.2 | -9.8 | 7.2 |
|---|---|---|---|---|---|---|

| 2.1 | -3.1 | 1.1 | -3.3 | 7.8 | -9.3 | 8.2 |
|---|---|---|---|---|---|---|

↑

| 0.1 | 0.7 | 1.0 | 0.7 | 0.1 | -0.5 | -0.7 |
|---|---|---|---|---|---|---|

| -0.3 | -0.7 | -1 | -0.7 | -0.3 | 0 | 0.3 |
|---|---|---|---|---|---|---|

...

**+**

**+**

| 2.4 | -2.4 | 2.1 | -2.6 | 8.1 | -9.3 | 7.9 |
|---|---|---|---|---|---|---|

| 2.4 | -2.4 | 2.1 | -2.6 | 8.1 | -9.3 | 7.9 |
|---|---|---|---|---|---|---|

...

↑

↑

3666

3666

...

↑

↑

'My'

'My'

...

University of Antwerp
Adrem | Adrem Data Lab

# Positional Encoding

- To maintain position: add position encoding to token

# Number of Parameters? Base Model

- **N=6, $d_{model}$=512, $d_{ff}$=2048, h=8, $d_k$=64, $d_v$=64**

- **Input & output embedding: ≈ 19M**
  - They are shared
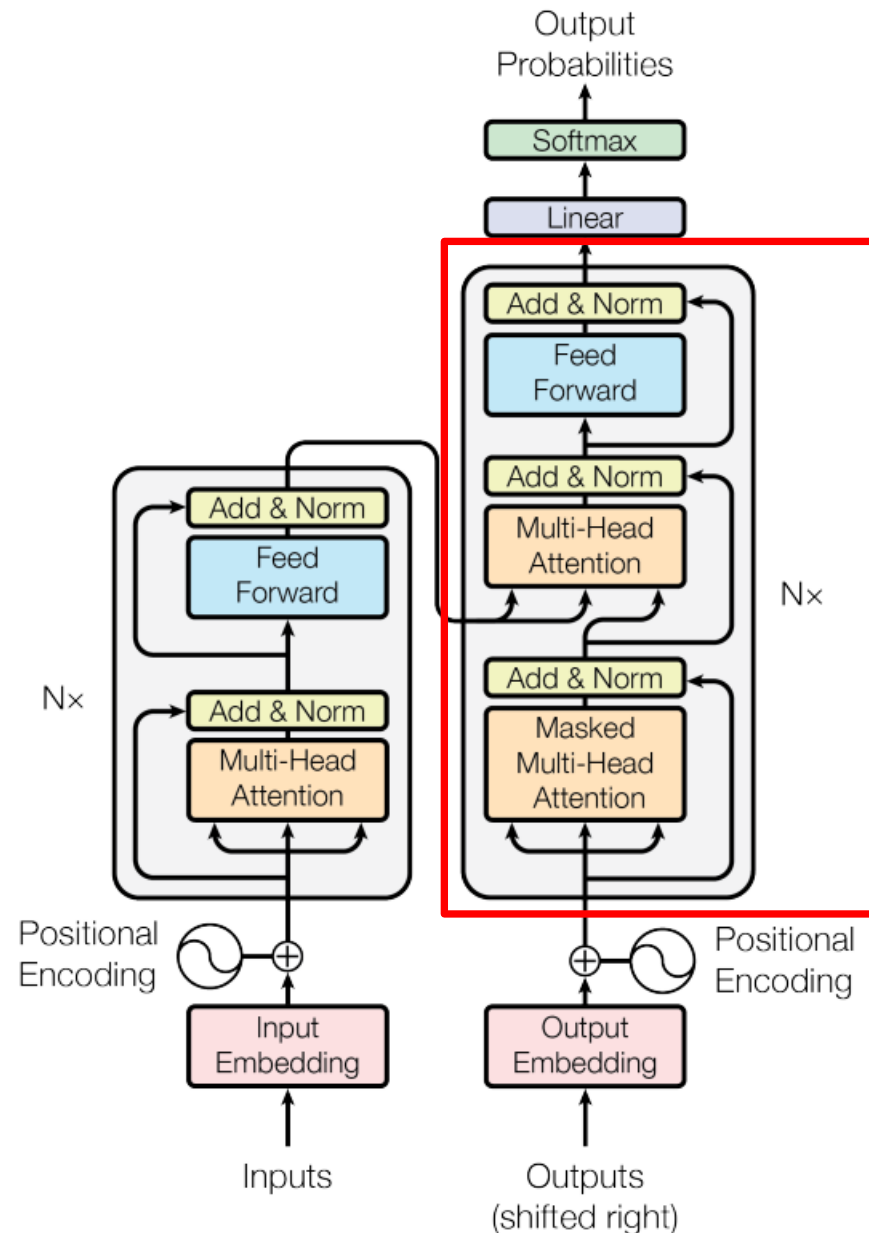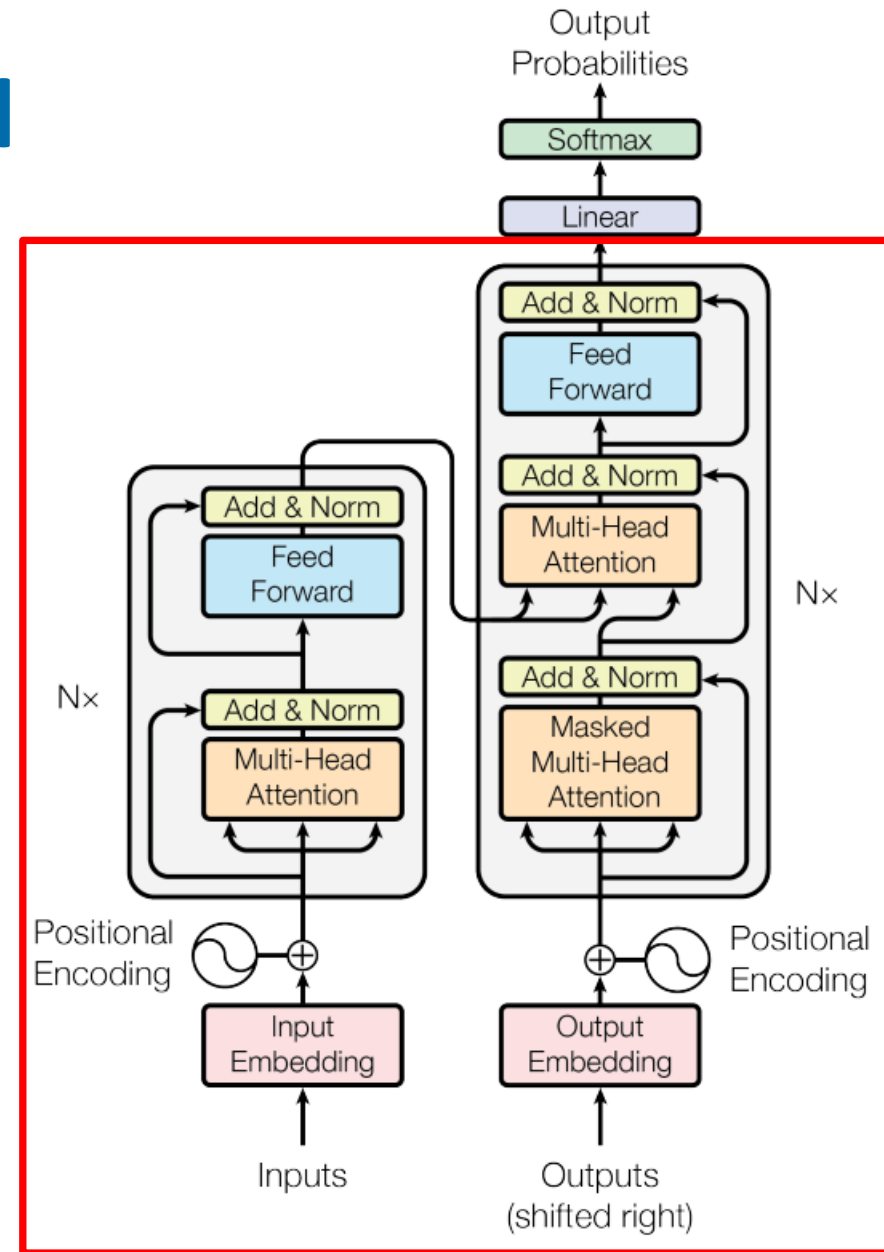  - Shared vocabulary of 37,000 tokens
  - 512-dimensional representation

# Number of Parameters? Base Model

- **$N=6$, $d_{model}=512$, $d_{ff}=2048$, $h=8$, $d_k=64$, $d_v=64$**

- **Attention mechanism: $\approx 6 \times 1M$**
  - 8 heads
    - Each matrix $W^Q$, $W^K$, $W^V$ : 512 x 64
  - Matrix $W^O$ : 512 x 512 = 262,144
- **Feed Forward: $\approx 6 \times 2.1M$**
  - (512+1) x 2048 + (2048+1) x 512
- **Total: $\approx 18.3M$**

# Number of Parameters? Base Model

- **N=6, $d_{model}$=512, $d_{ff}$=2048, h=8, $d_k$=64, $d_v$=64**

- **Attention mechanism: ≈ 2 x 6 x 1M**
  - Encoder – decoder attention
  - Masked self-attention
- **Feed Forward: ≈ 6 x 2.1M**
  - (512+1) x 2048 + (2048+1) x 512
- **Total: ≈ 24.6M**

# Number of Parameters? Base Model

- **$N=6$, $d_{model}=512$, $d_{ff}=2048$, $h=8$, $d_k=64$, $d_v=64$**

- **Total:**
  - Input/output embedding: ≈ 19M
  - Encoder: ≈ 18.3M
  - Decoder: ≈ 24.6M

- **≈ 61.9M parameters ***



* Paper states 65M Parameters ; this is due to estimations and some details we omitted (Norm)

# Number of Parameters?

- **Base Model**
  - N=6, $d_{model}$=512, $d_{ff}$=2048, h=8, $d_k$=64, $d_v$=64
  - ≈ 65M parameters

- **Big Model**
  - N=6, $d_{model}$=1024, $d_{ff}$=4096, h=16, $d_k$=64, $d_v$=64
  - ≈ 213M parameters

# Training Regime

- **Dataset consists of pairs:**
  - Text A = $[a_1, a_2, ..., a_n]$
  - Text B = $[b_1, b_2, ..., b_m]$

- **Input to the network:**
  - $[a_1, a_2, ..., a_n, <BOS>, b_1, b_2, ..., b_m]$

- **Expected output:**
  - $[b_1, b_2, ..., b_m]$

# Training Regime

- **Objective function: logloss**
  - Model output = 1 distribution over tokens per output slot
  - Reward high probabilities for the correct token

- **Adam optimizer**

- **Different regularization techniques were used**
  - Dropout, Label smoothing

- **Learning rate varied during training**
  - First increase, then decrease

University of Antwerp
I Adrem I Adrem Data Lab

# Testing: Generating All Answers

- Generate response tokens one by one
  - "autoregression"

Attention
is
all
you
need

**Encoder**

**Decoder**

<START>

L'attention

# Testing: Generating All Answers

- Generate response tokens one by one
  - "autoregression"

Attention
is
all
you
need

**Encoder**

**Decoder**

<START>
L'attention

L'attention
**est**

# Testing: Generating All Answers

- Generate response tokens one by one
  - "autoregression"

Attention
is
all
you
need

**Encoder**

<START>
L'attention
est

**Decoder**

L'attention
est
**tout**

University of Antwerp
I Adrem I Adrem Data Lab

# Testing: Generating All Answers

- Generate response tokens one by one
  - "autoregression"



Attention
is
all
you
need

**Encoder**

<START>
L'attention
est
tout

**Decoder**

L'attention
est
tout
**ce**

# Experimental Results

- **2 translation tasks: EN-DE and EN-FR**

- **BLEU score used to assess quality of result (higher is better)**

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [18] | 23.75 | | | |
| Deep-Att + PosUnk [39] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [38] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [9] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [32] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [38] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9] | 26.36 | **41.29** | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.8** | $2.3 \cdot 10^{19}$ | |

# Conclusion

- **New architecture Transformer proposed**
  - Based on attention mechanism
  - No recurrent neural nets, convolutions needed

- **Experiments show promising behavior**

- **Translation task:**
  - Outperforms state-of-the art in accuracy
  - For lower or comparable training costs

- **... and the rest is history ...**

# Machine Translation on WMT2014 English-German
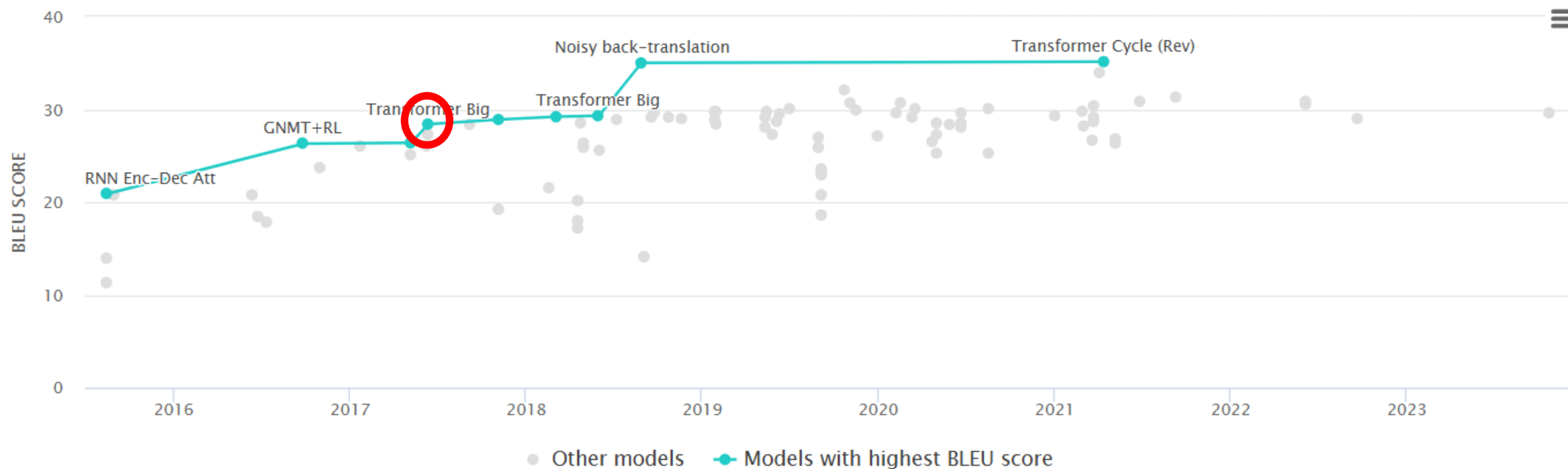
# Sources

- Vaswani, A. et al. (2017). Attention is all you need. Advances in Neural Information Processing Systems.
- Cho, K. et al. (2014). On the properties of neural machine translation: Encoder–Decoder approaches. In Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation.
- Bahdanau, D. et al. (2015). "Neural machine translation by jointly learning to align and translate." 3rd International Conference on Learning Representations
- Papineni, K. et al. (2002). Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th annual meeting of the Association for Computational Linguistics (pp. 311-318).
- Yang, J. et al. (2024) "Harnessing the power of llms in practice: A survey on chatgpt and beyond." ACM Transactions on Knowledge Discovery from Data 18.6: 1-32.
- Raschka, S. (2024). Build a Large Language Model (From Scratch). Simon and Schuster.

- https://www.kaggle.com/datasets/mohamedlotfy50/wmt-2014-english-german
- https://paperswithcode.com/task/machine-translation
- https://jalammar.github.io/illustrated-transformer/