

Attention & Diffusion

Lecture 3

Alexander Schönhuth



Bielefeld University
April 25, 2023

CONTENTS

- ▶ Transformers
 - ▶ Encoder Structure (today)
 - ▶ Self Attention (today)
 - ▶ Decoder Structure (next lecture)
- ▶ Transformer Versions and Training
 - ▶ Encoder Only (next lecture)
 - ▶ Encoder-Decoder (next lecture)
 - ▶ Decoder Only (next lecture)

Transformers I

TRANSFORMERS: MOTIVATION

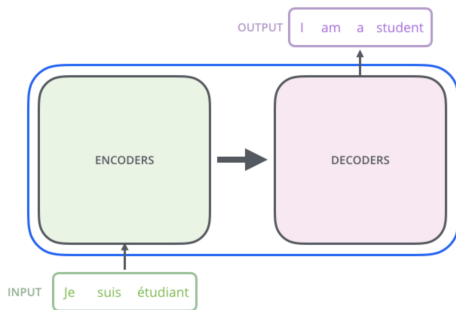


Transformers translate “languages”

From <https://jalamar.github.io>

- ▶ Inspiration for transformers: translating languages
- ▶ Transformers lend themselves to (maximum) parallelization
- ▶ Google: reference model for cloud TPU based computations

TRANSFORMERS: MOTIVATION II

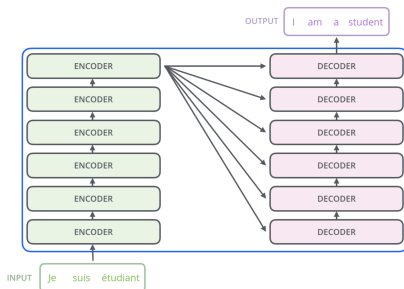


From <https://jalammar.github.io>

- ▶ Transformers employ encoder-decoder architecture
- ▶ However, neither encoder nor decoder RNN based
- ▶ *Seminal paper: "Attention is all you need"*

<https://arxiv.org/abs/1706.03762>

TRANSFORMERS: STRUCTURE I

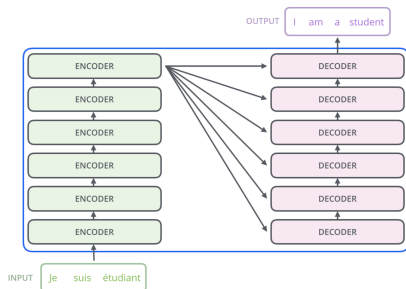


Transformers: encoders and decoders layer structured

From <https://jalammar.github.io>

- ▶ Transformers make use of stacks of encoders and decoders
 - ▶ Seminal paper: stacks are 6 layers each
 - ▶ Other numbers very well conceivable
 - ▶ Architectural design may vary by application

TRANSFORMERS: STRUCTURE II

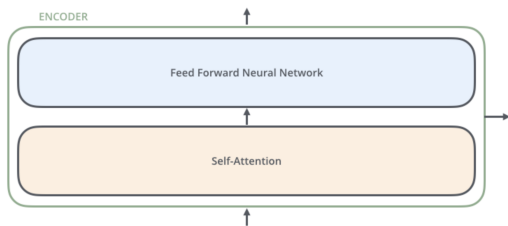


Transformers: encoders and decoders layer structured

From <https://jalamar.github.io>

- ▶ Encoders and decoders interact in different ways
 - ▶ All but last encoder provide input to next encoder
 - ▶ Last encoder provides input to all decoders
 - ▶ All but last decoder provide input to next decoder
 - ▶ Last decoder outputs translated sentence

TRANSFORMERS: ENCODER STRUCTURE I

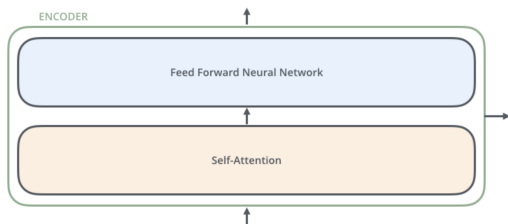


Transformers: encoders follow particular structure

From <https://jalammar.github.io>

- ▶ Encoders are identical in structure
 - ▶ But they do not share weights
- ▶ Encoders have two sublayers
 - ▶ A self-attention layer
 - ▶ A feedforward neural network layer

TRANSFORMERS: ENCODER STRUCTURE II

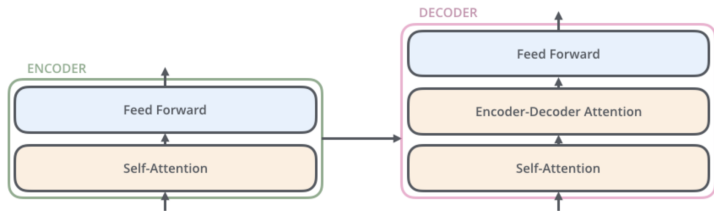


Transformers: encoders follow particular structure

From <https://jalammar.github.io>

- ▶ Self-attention layer:
 - ▶ Encoder can look at other words when encoding words
- ▶ Feedforward neural network (FFNN) layer:
 - ▶ Exact same FFNN applied for each position in sentence

TRANSFORMERS: DECODER STRUCTURE I



Transformers: encoder and decoder interact in particular way

From <https://jalamar.github.io>

- ▶ Decoder shares structure with encoder, but ...
- ▶ ... has an additional encoder-decoder attention sublayer
- ▶ Helps decoder to pay attention as guided by input

Self-Attention: Reminder

SELF-ATTENTION: DEFINITION

- ▶ Consider a sequence of tokens $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- ▶ Each token has its own query, key, and value
- ▶ Hence, each token can attend to each other token:
 - ▶ Pair the query vector with the key of the other token
 - ▶ This yields a weight for its own value
- ▶ Compute weighted sum of values as representation in next layer

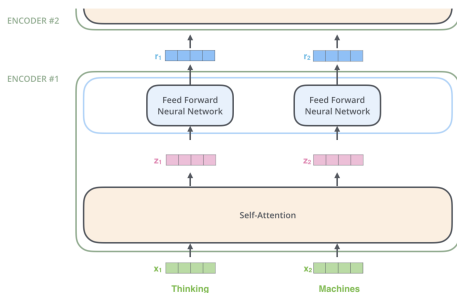
SELF-ATTENTION: FORMAL SUMMARY

- ▶ Consider a sequence of tokens $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$
- ▶ Replace \mathbf{q} with \mathbf{x} and both $\mathbf{k}_i, \mathbf{v}_i$ with \mathbf{x}_i in (1) from Lecture 2
- ▶ One obtains a new sequence $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^d$ by

$$\mathbf{z}_i := f(\mathbf{x}_i, ((\mathbf{x}_1, \mathbf{x}_1), \dots, (\mathbf{x}_n, \mathbf{x}_n))) = \sum_{j=1}^n \alpha(\mathbf{x}_i, \mathbf{x}_j) \mathbf{x}_j \in \mathbb{R}^d \quad (1)$$

Transformers continued

TRANSFORMERS: ENCODER STRUCTURE III

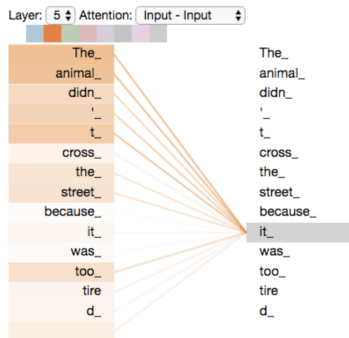


Transformers: encoders sublayer by sublayer

From <https://jalammar.github.io>

1. Words are embedded \rightarrow yields vectors x_i
2. Vectors x_i run through self-attention sublayer \rightarrow yields vectors z_i
3. Each z_i runs through exact same FFNN \rightarrow yields vectors r_i

TRANSFORMERS: SELF-ATTENTION I



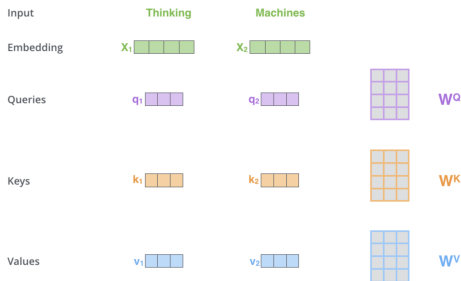
► 5th sublayer, 2nd out of 8 attention heads

- Word “it” pays most attention to “the animal”
- Word “it” pays less attention to “the street”
- Word “it” pays no attention to “because”

Words pay more/less attention to others

From <https://jalammr.github.io>

TRANSFORMERS: SELF-ATTENTION II

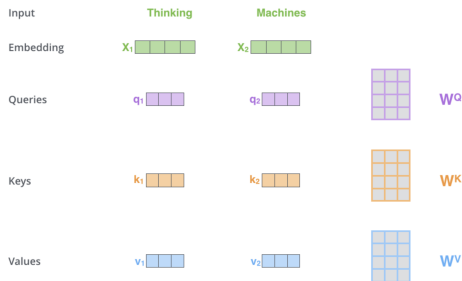


Self-attention: queries, keys and values

From <https://jalammarr.github.io>

- ▶ Input vectors x_i are transformed to
 - ▶ queries q_i , keys k_i , values v_i by
 - ▶ applying matrices W^Q , W^K , W^V to x_i from the right

TRANSFORMERS: SELF-ATTENTION III

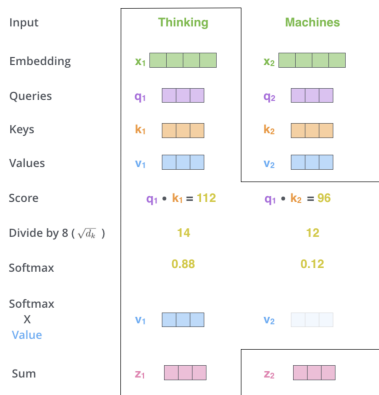


Self-attention: queries, keys and values

From <https://jalammarm.github.io>

- ▶ Seminal paper: dimension of $x_i = 512$, of $q_i, k_i, v_i = 64$
 - ▶ So, $W^Q, W^K, W^V \in \mathbb{R}^{512 \times 64}$
 - ▶ Recall: $q_i = x_i W^Q, k_i = x_i W^K, v_i = x_i W^V$

TRANSFORMERS: SELF-ATTENTION IV



Self-attention: from input to output

From <https://jalamar.github.io>

- ▶ Scores for x_1 w.r.t. v_1, v_2
 - ▶ v_1 : Compute $q_1 \cdot k_1$, divide by 8, yields 112
 - ▶ v_2 : Compute $q_1 \cdot k_2$, divide by 8, yields 96
- ▶ Softmax'ing: Probabilities 0.88, 0.12 for v_1, v_2
- ▶ Final output for x_1 :

$$0.88 \cdot v_1 + 0.12 \cdot v_2$$

TRANSFORMERS: SELF-ATTENTION V



Calculating queries, keys and values

From <https://jalammar.github.io>

- ▶ Pack embedded words into matrix X
 - ▶ Each row corresponds to one word
- ▶ Multiply X with trained matrices W^Q, W^K, W^V
- ▶ Recall real dimensions:
 - ▶ Words: 512 (here: 4);
 Q, K, V : 64 (here: 3)

TRANSFORMERS: SELF-ATTENTION VI

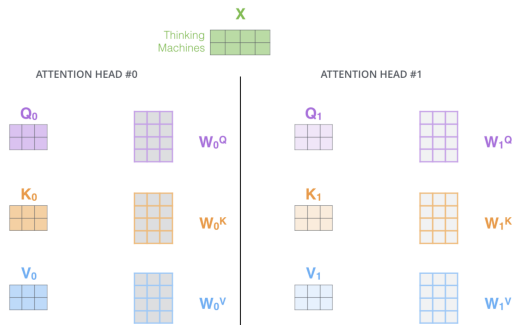
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \mathbf{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}} \right) \begin{matrix} \mathbf{V} \\ \begin{matrix} \square & \square \\ \square & \square \end{matrix} \end{matrix} \\ = \begin{matrix} \mathbf{Z} \\ \begin{matrix} \square & \square & \square \\ \square & \square & \square \end{matrix} \end{matrix}$$

Computing values: compact matrix representation

From <https://jalammar.github.io>

1. Multiply queries with keys: $\mathbf{Q} \cdot \mathbf{S}^T$
2. Normalize relative to query/key length d_k ($= 64$ in reality)
3. Softmax across columns: $\mathbf{S} := \text{softmax}(\mathbf{Q}\mathbf{S}^T / \sqrt{d_k})$ (here: $\in \mathbb{R}^{2 \times 2}$)
4. Compute weighted sum for each word: $\mathbf{Z} = \mathbf{S} \cdot \mathbf{V}$

TRANSFORMERS: MULTI-HEAD ATTENTION I

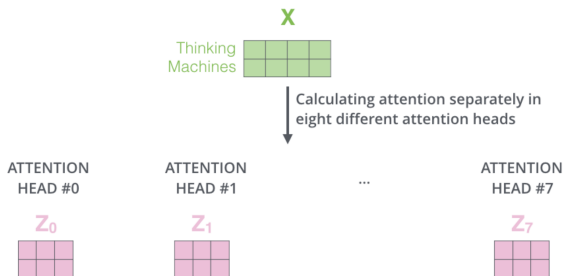


Multi-head attention with 2 heads

From <https://jalammar.github.io>

- ▶ Learn several “heads”, attending to different interactions
 - ▶ Learn several different W_i^Q, W_i^K, W_i^V (here: $i = 2$)
 - ▶ Establish differences by randomized initialization

TRANSFORMERS: MULTI-HEAD ATTENTION II



Original paper: multi-head attention with 8 heads

From <https://jalammar.github.io>

- ▶ Seminal paper uses 8 different attention heads
- ▶ How to summarize / combine the 8 resulting outputs?

TRANSFORMERS: MULTI-HEAD ATTENTION III

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^O that was trained jointly with the model

\times



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



Combining outputs of different attention heads

From <https://jalammar.github.io>

► Combining attention head outputs:

1. Concatenate all outputs
2. Multiply resulting matrix with learned matrix W^O
3. Yields output being equal to input in dimension

👉 Remark: Need to learn W^O also for single head

TRANSFORMERS: MULTI-HEAD ATTENTION IV

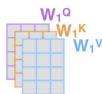
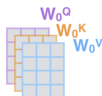
1) This is our input sentence*

Thinking
Machines

2) We embed each word*



3) Split into 8 heads. We multiply X or R with weight matrices



4) Calculate attention using the resulting $Q/K/V$ matrices



5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer



W^O



Z



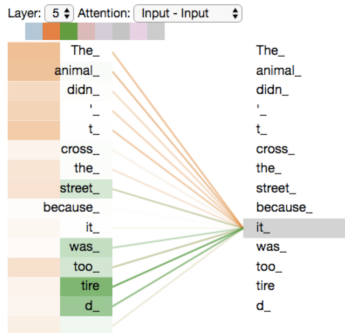
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



Multi-head attention: Overview / Summary
X: embedded words, input for first attention layer
R: output of earlier layer input for all but first layer

From <https://jalammar.github.io>

TRANSFORMERS: MULTI-HEAD ATTENTION V

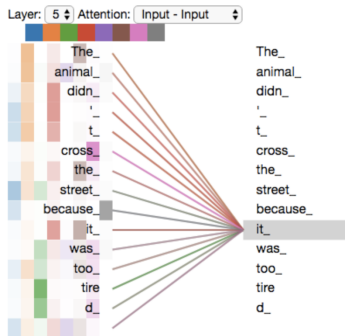


- ▶ Considering two attention heads, orange and green
- ▶ *Orange*: "it" mostly attends to "the animal"
- ▶ *Green*: "it" mostly attends to "tired"

Multi-head attention:
Considering two (of eight) heads

From <https://jalammr.github.io>

TRANSFORMERS: MULTI-HEAD ATTENTION VI



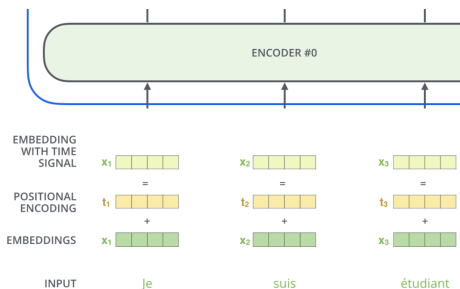
- ▶ Considering all eight attention heads
- ▶ Things are more difficult to interpret
- ▶ Each head reflects different relationships

Multi-head attention:

Considering all (eight) heads

From <https://jalamar.github.io>

TRANSFORMERS: POSITIONAL ENCODING

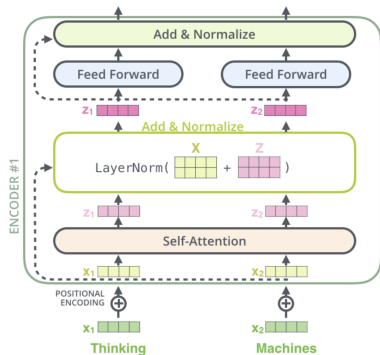


Integrating positional encodings

From <https://jalammr.github.io>

- ▶ *Problem:* Self attention unaware of order
- ▶ *Solution:* Consider vectors t_i that code order of word x_i
 - ▶ Add t_i to x_i → position i of x_i can be determined
 - ▶ Details of generation of t_i not discussed here

TRANSFORMERS: ENCODER DETAILS



Transformer encoder block: details

From <https://jalammr.github.io>

1. Embedded words are equipped with positional encodings
2. Self attention is applied
 - 2.1 Original x_i is added to z_i
 - ↳ Residual skip connection
 - 2.2 Layer norm is applied
 - ↳ Normalizes values across layer
3. Each resulting z_i passed through identical feedforward NN (FFNN)
 - 3.1 Original z_i added to FFNN output
 - ↳ Residual skip connection
 - 3.2 Layer norm is applied
 - ↳ Normalizes values across layer

REFERENCES

- ▶ Jay Alammar's ML blog:
 - ▶ "Visualizing A Neural Machine Translation Model", see <https://jalammar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention>
 - ▶ "The Illustrated Transformer", see <https://jalammar.github.io/illustrated-transformer/>
- ▶ <http://d2l.ai>, 11.4–11.7, 11.9

Thanks for your attention!!