

Attention & Diffusion

Lecture 4

Alexander Schönhuth



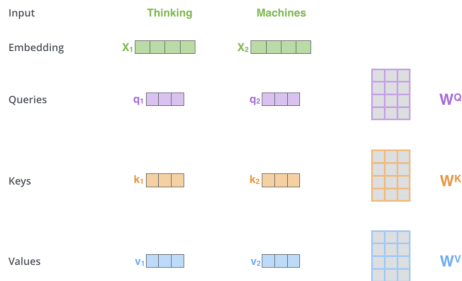
Bielefeld University
May 9, 2023

CONTENTS

- ▶ Transformers
 - ▶ Self Attention Reminder
 - ▶ Decoder Structure
- ▶ Transformer Versions and Training
 - ▶ Encoder Only
 - ▶ Encoder-Decoder
 - ▶ Decoder Only

Self Attention: Illustrated Reminder

TRANSFORMERS: SELF-ATTENTION REMINDER I

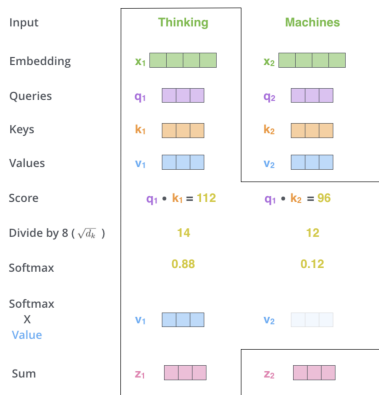


Self-attention: queries, keys and values

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 REMINDER II



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 REMINDER III

$$\mathbf{X} \times \mathbf{W}^Q = \mathbf{Q}$$

$$\mathbf{X} \times \mathbf{W}^K = \mathbf{K}$$

$$\mathbf{X} \times \mathbf{W}^V = \mathbf{V}$$

Calculating queries, keys and values

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

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

TRANSFORMERS: SELF-ATTENTION REMINDER IV

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

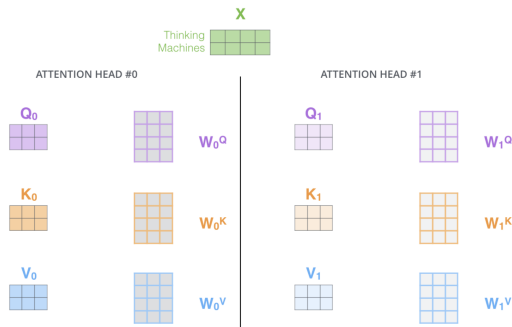
Computing values: compact matrix representation

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

1. Multiply queries with keys: $\mathbf{Q} \cdot \mathbf{K}^T$
2. Normalize relative to query/key length d_k ($= 64$ in reality)
3. Softmax across columns: $\mathbf{S} := \text{softmax}(\mathbf{Q}\mathbf{K}^T / \sqrt{d_k})$ (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

REMINDER 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

REMINDER III

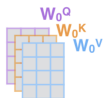
1) This is our input sentence*

Thinking
Machines

2) We embed each word*



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



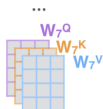
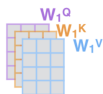
4) Calculate attention using the resulting $Q/K/V$ matrices



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



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



W^O



Z



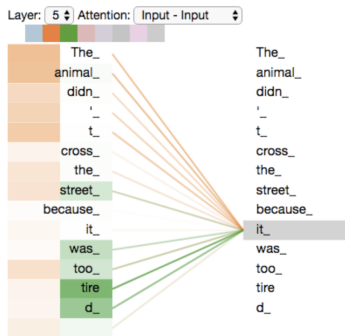
Multi-head attention: Overview / Summary

X: embedded words, input for first attention layer

R: output of earlier layer input for all but first layer

TRANSFORMERS: MULTI-HEAD ATTENTION

REMINDER IV

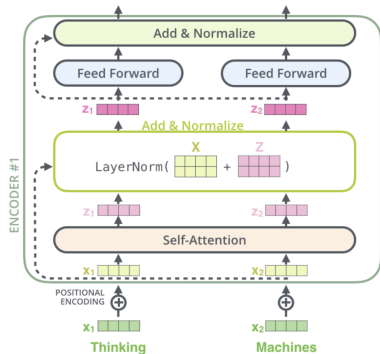


- ▶ Considering two attention heads, orange and green
- ▶ *Orange*: "it" mostly attends to "the animal"
- ▶ *Green*: "it" mostly attends to "tired"

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

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

TRANSFORMERS: ENCODER DETAILS

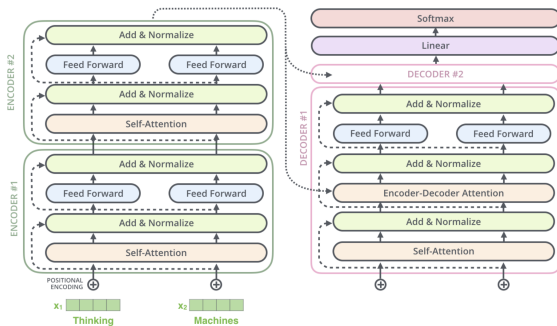


Transformer encoder block: details

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

TRANSFORMERS: ENCODER-DECODER INTERACTION



Transformer with two encoder and two decoder blocks

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

- ▶ Decoder blocks integrate encoder-decoder attention layers
 - ▶ Between decoder self attention and FFNN layer
 - ▶ Encoder output transformed into keys and values
 - ▶ Decoder output transformed into queries

TRANSFORMER: DECODER I

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

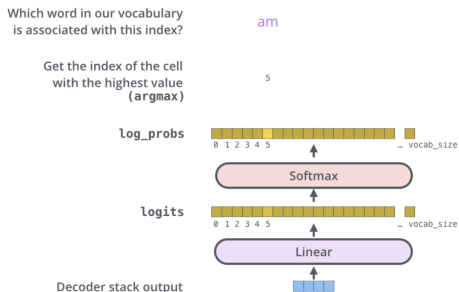
1. Encoder processes input sequence (here: with positional encoding)
2. Output of top encoder transformed into keys $\mathbf{K}_{\text{encdec}}$ and values $\mathbf{V}_{\text{encdec}}$
3. Decoder uses $\mathbf{K}_{\text{encdec}}$ and $\mathbf{V}_{\text{encdec}}$ in encoder-decoder attention layer

TRANSFORMER: DECODER II

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

1. Decoder takes in already generated tokens (words)
2. Self-attention: decoder only attends to already generated tokens
 - ▶ Achieved by masking future positions
3. Encoder-decoder attention layer generates its own queries
 - ▶ but uses keys and values from topmost encoder output

TRANSFORMERS: DECODER FINAL LAYER



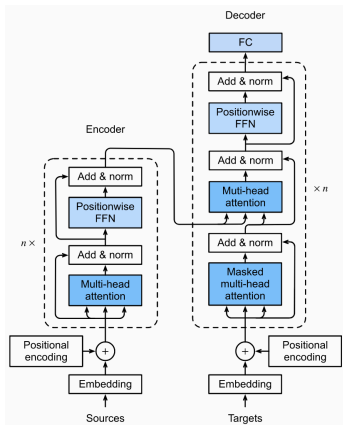
Transformer decoder: final layer consists of linear and softmax sublayer

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

- ▶ Linear layer takes decoder output, computes a value for each word
 - ▶ See *logits* layer in figure; number of words equal to size of vocabulary
- ▶ Softmax layer turns values into probabilities
 - ▶ Yields *log_probs* layer; word with greatest probability is output

Transformer Variants

TRANSFORMERS: ARCHITECTURE SUMMARY I



Encoder

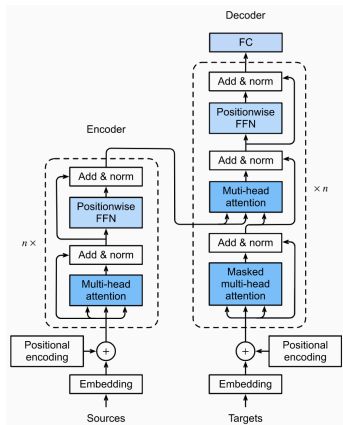
- ▶ Both encoder and decoder consist of n layers
 - ☞ original paper: $n = 6$
- ▶ Stacks identical layers
- ▶ Each layer has two sublayers
 - ▶ Multi-head attention layer
 - ▶ Positionwise feedforward neural network
- ▶ Contains skip connections
 - ☞ inspired by ResNet

Transformer: Summary.

n encoder and n decoder layers

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

TRANSFORMERS: ARCHITECTURE SUMMARY II



Transformer: Summary.
 n encoder and n decoder layers

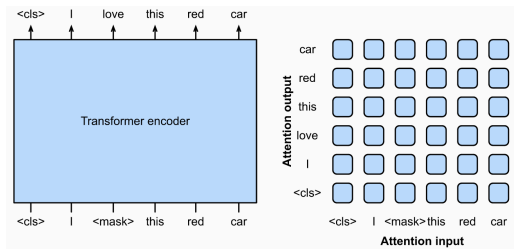
From <https://d21.ai>

Decoder

- ▶ Also stacks identical layers
- ▶ Each layer: three sublayers
 - ▶ Multi-head self attention
 - ▶ Encoder-decoder attention
 - ▶ Positionwise feedforward neural network
- ▶ Encoder-decoder attention does not exist in encoder
- ▶ Contains skip connections
↳ inspired by ResNet
- ▶ Each position only attends to earlier positions
 - ▶ *Masked* attention preserves autoregressive property

Transformer Variants: Encoder Only

TRANSFORMER VARIANTS: ENCODER ONLY I

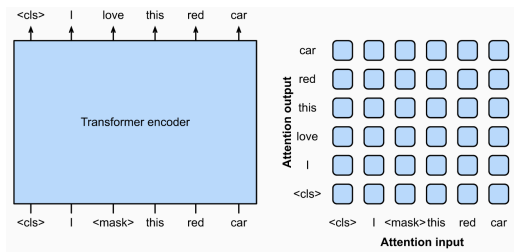


Transformer encoder only: pretraining

From <https://d21.ai>

- ▶ *Attention input* \leftrightarrow embedded words x_i
- ▶ *Attention output* \leftrightarrow new "words" z_i
- ▶ Meaning right panel: each x_i contributes to each z_i

TRANSFORMER VARIANTS: ENCODER ONLY II

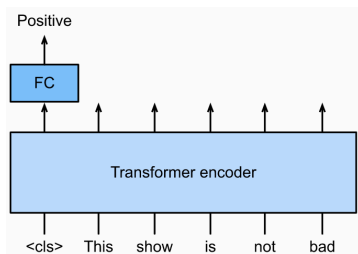


Transformer encoder only: pretraining

From <https://d21.ai>

- ▶ Prominent example: *Bidirectional Encoder Representations from Transformers (BERT)*, see <https://arxiv.org/abs/1810.04805>
- ▶ *Pretraining* supposed to pick up basic language structure
- ▶ *Principle*: Learn masked words in sentences

TRANSFORMER VARIANTS: ENCODER ONLY III



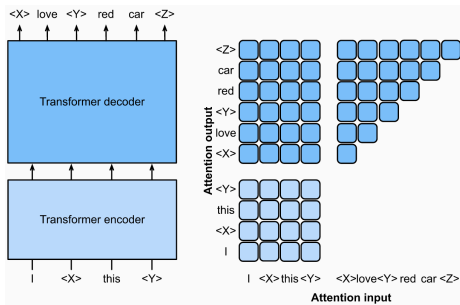
Transformer encoder only: finetuning for sentiment analysis

From <https://d21.ai>

- ▶ After pretraining, encoder-only transformer is *finetuned*
 - ▶ Involves different kind of training
- ▶ *Example*: Sentiment analysis
 - ▶ Predicting sentiments inherent to sentences
- ▶ *Principle*: Use final representation of special token < cls >

Transformer Variants: Encoder-Decoder

TRANSFORMER VARIANTS: ENCODER-DECODER I

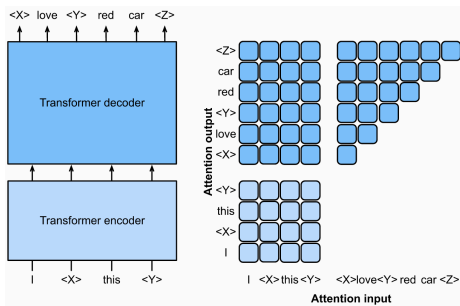


Transformer encoder-decoder: pretraining

From <https://d21.ai>

- ▶ *Advantage:* Output can vary in length
- ▶ Prominent example: *T5*, see <https://arxiv.org/abs/1910.10683>

TRANSFORMER VARIANTS: ENCODER-DECODER II

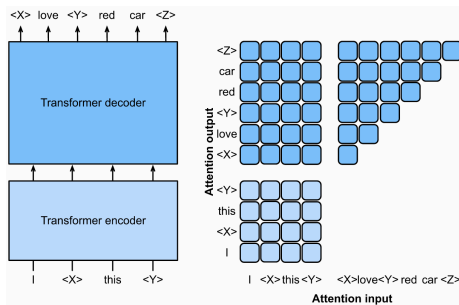


Transformer encoder-decoder: pretraining

From <https://d21.ai>

- ▶ *Pretraining Example:* Predict consecutive spans
- ▶ *Here:* Replace “<X>” with “<X> love” and “<Y>” with “<Y> red car”

TRANSFORMER VARIANTS: ENCODER-DECODER III

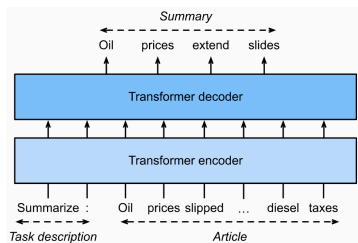


Transformer encoder-decoder: pretraining

From <https://d21.ai>

- ▶ *Encoder*: Each input token attends to each other
- ▶ *Decoder*: Target tokens attend to
 - ▶ all input tokens (*encoder-decoder attention*)
 - ▶ only past and present target tokens (*causal attention*)

TRANSFORMER VARIANTS: ENCODER-DECODER IV



Transformer encoder-decoder: Finetuning for generating text summaries

From <https://d21.ai>

- ▶ After pretraining, encoder-decoder transformer is *finetuned*
 - ▶ Involves different training principle
- ▶ *Example*: Summarization of large texts
 - ▶ *Input*: Task description and large text
 - ▶ *Output*: Brief summary of large text

TRANSFORMER VARIANTS: ENCODER-DECODER V



Teddy bears swimming at the Olympics 400m Butterfly event.



A cute corgi lives in a house made out of sushi.



A cute sloth holding a small treasure chest. A bright golden glow is coming from the chest.

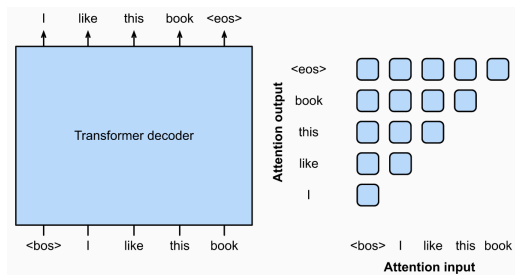
Imagen, based on T5 encoder: Turning texts into images

From <https://d21.ai>

- ▶ Generate image that reflects text contents
- ▶ Text-to-image model “Imagen”, see <https://arxiv.org/abs/2205.11487>
- ▶ Imagen based on “frozen” T5 encoder

Transformer Variants: Decoder Only

TRANSFORMER VARIANTS: DECODER ONLY I

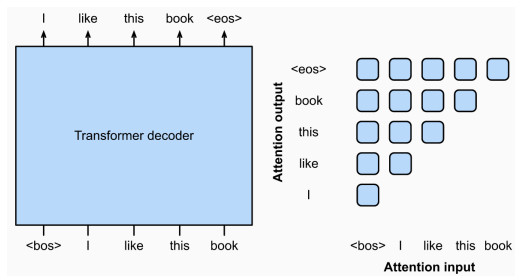


Transformer decoder only: pretraining

From <https://d21.ai>

- ▶ De facto architecture in large-scale language modeling
- ▶ Encoder-decoder attention sublayers removed
- ▶ *Pretraining*: Teacher forcing
 - ▶ Target sequence is input sequence shifted by one token

TRANSFORMER VARIANTS: DECODER ONLY II

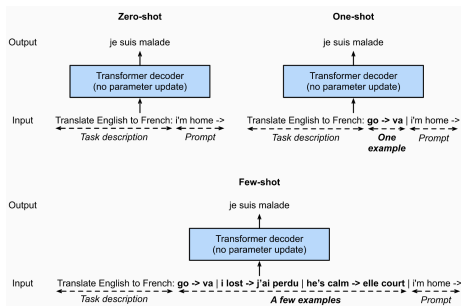


Transformer decoder only: pretraining

From <https://d21.ai>

- ▶ *Self-supervised learning*: Learns structures in unlabeled data
 - ▶ Leverages abundantly existing, unlabeled text corpora
- ▶ Prominent example: *GPT-3*, see <https://arxiv.org/abs/2005.14165>
 - ▶ Basis of *ChatGPT*, for example

TRANSFORMER VARIANTS: DECODER ONLY III

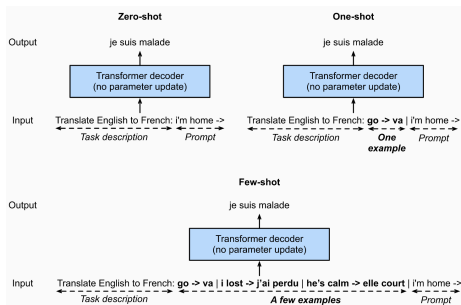


Transformer decoder only: finetuning

From <https://d21.ai>

- ▶ GPT-2 demonstrated that model can be re-used for other tasks
 - ▶ without parameter re-training / updating (!), so no finetuning
- ▶ GPT-3 exploits the *in-context learning* principle further

TRANSFORMER VARIANTS: DECODER ONLY IV



Transformer decoder only: finetuning

From <https://d2l.ai>

- ▶ In-context learning requires task description and prompt, as task input
- ▶ In addition, in-context learning may involve no examples (*zero-shot*), one example (*one-shot*) or few examples: *few-shot*

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!!