

# Transformers

PA154 Language Modeling (10.1)

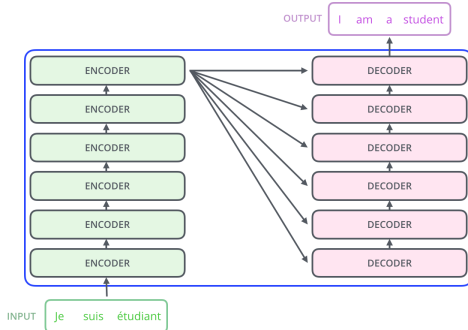
Pavel Rychlý

pary@fi.muni.cz

April 20, 2023

## Attention

- each decoder layer has access to all hidden states from the last encoder
- use attention to extract important parts (vector)

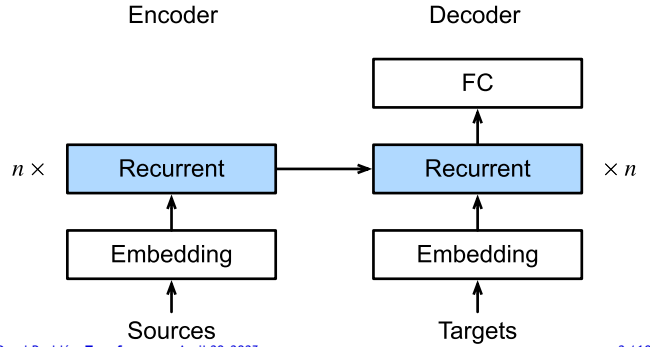


## Self-Attention

- instead of sequential processing
- attention to previous (and following) tokens
- fully parallel processing during training

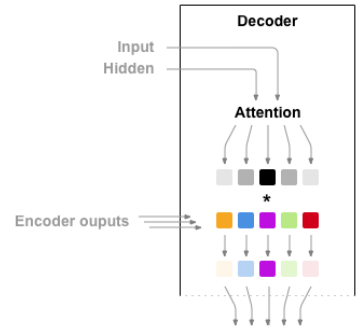
## Multi-layer encoder/decoder

- Encoder: Input sequence → state
- Decoder: state + sentence delimiter → output
- Problem: fix size state



## Attention

- use attention to extract important parts (vector)
- important = similar to "me"



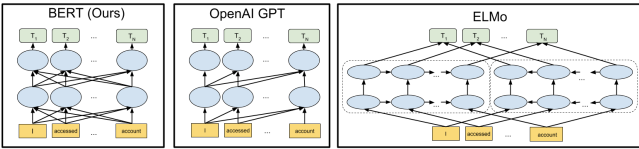
## Transformes

- Attention is All You Need
- self-attention in both encoder and decoder
- masked cross-attention in decoder

http://jalamar.github.io/illustrated-transformer/

## Transformers variants

- using context to compute token/sentence/document embedding
- BERT = Bidirectional Encoder Representations from Transformers
- GPT = Generative Pre-trained Transformer
- many variants: tokenization, attention, encoder/decoder connections



## Using pre-trained models

- (BERT) trained on huge amount of data
- finetuned on task specific data
- using output of BERT as an input to task specific model (without modification of BERT)

## BERT

- Google
- encoder only
- pre-training on raw text
- masking tokens, is-next-sentence
- big pre-trained models available
- domain (task) adaptation

**Input:** The man went to the [MASK]<sub>1</sub>. He bought a [MASK]<sub>2</sub> of milk.  
**Labels:** [MASK]<sub>1</sub> = store; [MASK]<sub>2</sub> = gallon

**Sentence A** = The man went to the store.  
**Sentence B** = He bought a gallon of milk.  
**Label** = IsNextSentence

**Sentence A** = The man went to the store.  
**Sentence B** = Penguins are flightless.  
**Label** = NotNextSentence

## GPT

- Open AI
- decoder only
- pre-training on raw text
- trained on prediction of next token

