

Any reproduction or distribution of this document, in whole or in part, is prohibited unless permission is granted by the authors

EE-559

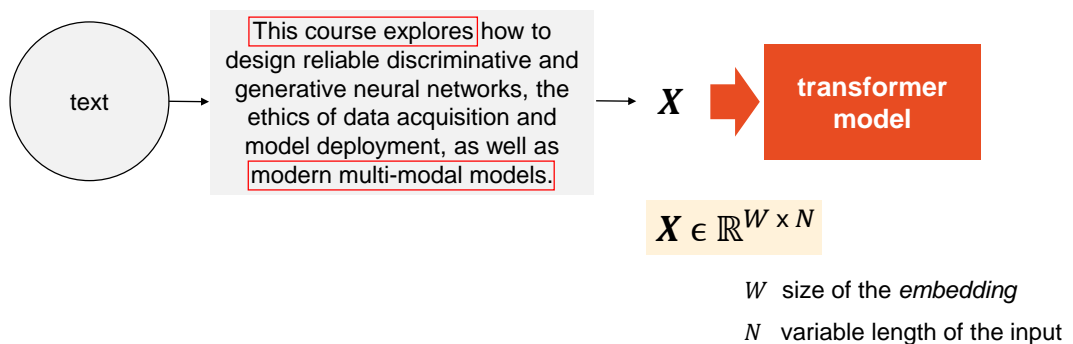
Deep Learning

What's on today?

- **Natural Language Processing**: how to analyse / synthesize language
- **Tokens**: on breaking text into small units
- **Self-attention**: how a model can focus on relevant parts of the input
- **Transformer**: how to use attention in a neural network
- **Encoder model**: how to process text to create a useful representation
- **Decoder model**: how to generate text
- **Encoder-decoder model**: how to map a sequence to a sequence
- **Exercises**: exploring attention and transformers

Natural Language Processing

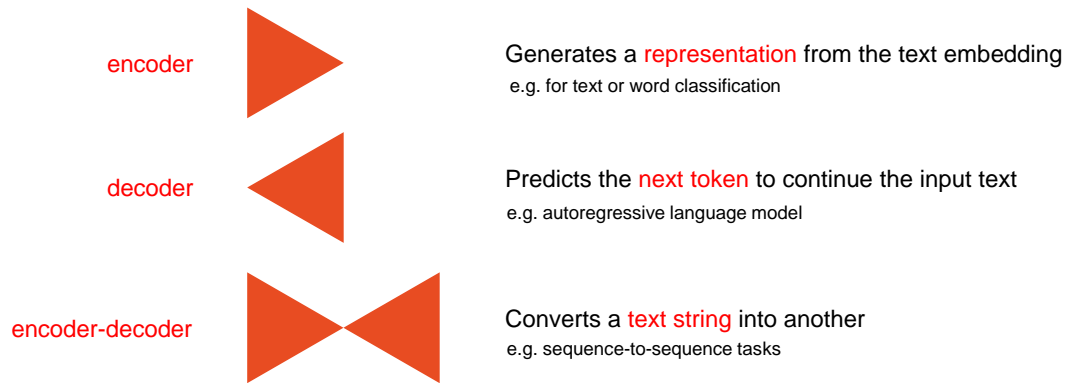
Natural Language Processing (NLP)



Concepts:

Analysis and synthesis of speech and language, syntax, pronouns, connections between words

Transformer models

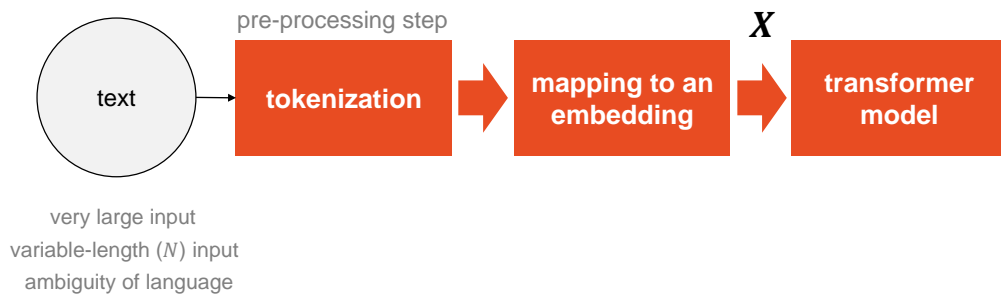


Concepts:

Sentiment analysis, name entity recognition, tokens, generative pre-trained transformers, language translation

Tokens

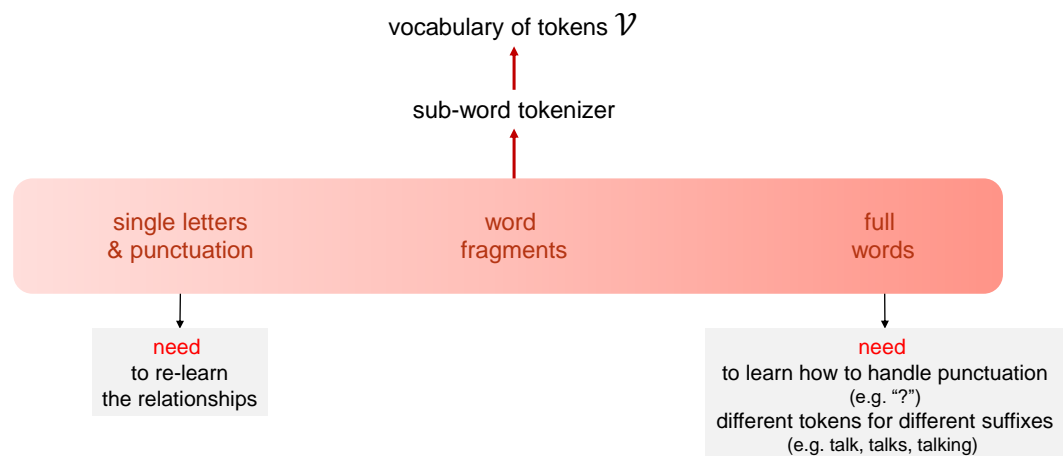
Tokenization of text



Concepts:

Words and word fragments, sub-word tokenizer, token embedding

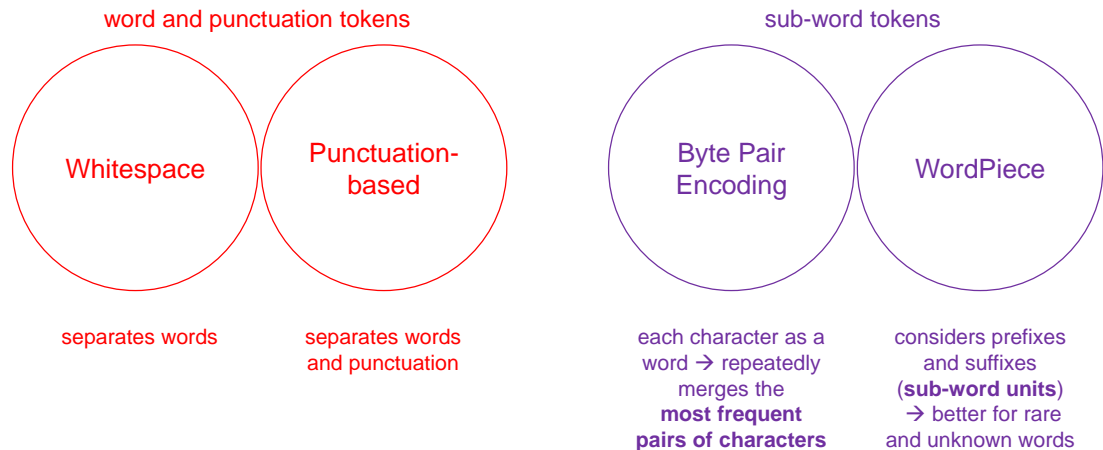
Tokenization: choices



Concepts:

Use frequency of commonly occurring sub-strings to merge them (byte pair encoding)

Tokenizers



Concepts: Counting the number of words in a text, identifying frequently occurring character combinations in training text, a tokenizer does not work equally well for all languages

More on units of text

Between words and characters: A Brief History of Open-Vocabulary Modeling and Tokenization in NLP

Sabrina J. Mielke^{1,2} Zaid Alyafeai³ Elizabeth Salesky¹
 Colin Raffel² Manan Dey⁴ Matthias Gallé⁵ Arun Raja⁶
 Chenglei Si⁷ Wilson Y. Lee⁸ Benoît Sagot^{9*} Samson Tan^{10*}

BigScience Workshop Tokenization Working Group

¹Johns Hopkins University ²HuggingFace ³King Fahd University of Petroleum and Minerals ⁴SAP

⁵Naver Labs Europe ⁶Institute for Infocomm Research, A*STAR Singapore ⁷University of Maryland

⁸BigScience Workshop ⁹Inria Paris ¹⁰Salesforce Research Asia & National University of Singapore

s.jm@sjmielke.com

Abstract

What are the units of text that we want to model? From bytes to multi-word expressions, text can be analyzed and generated at many granularities. Until recently, most natural language processing (NLP) models operated over words, treating those as discrete and atomic tokens, but starting with byte-pair encoding (BPE), subword-based approaches have become dominant in many areas, enabling small vocabularies while still allowing for fast inference. Is the end of the road character-level model or byte-level processing? In this survey, we connect several lines of work from the pre-neural and neural era, by showing how hybrid approaches

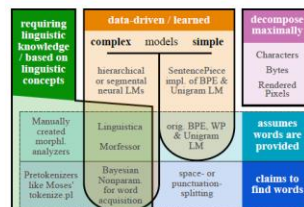
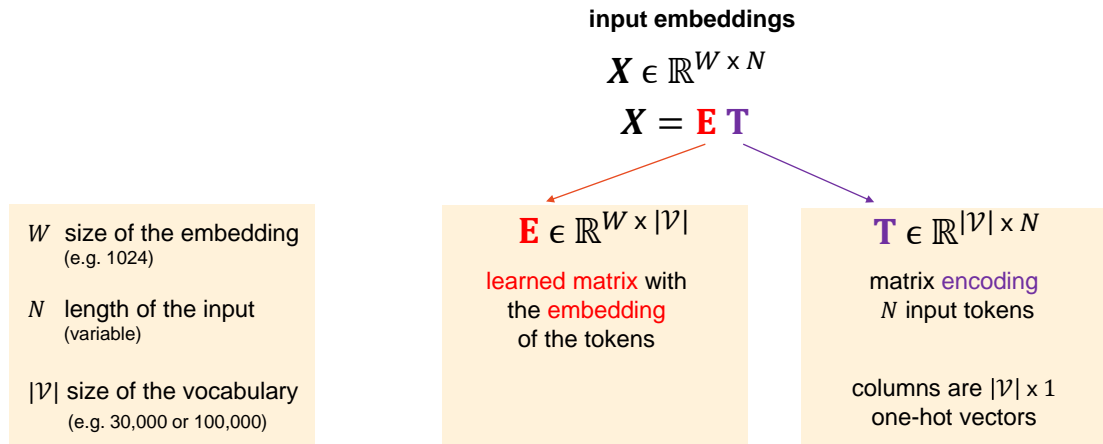


Figure 1: A taxonomy of segmentation and tokenization algorithms and research directions

arXiv:2112.10508

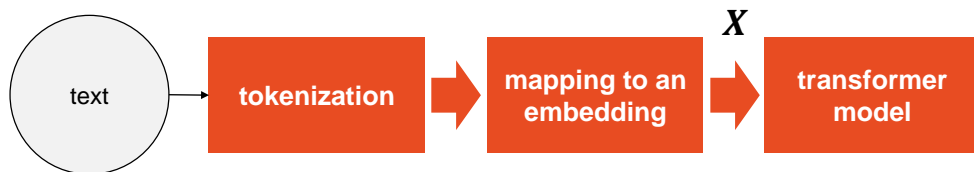
Mapping to an embedding



Concepts:

One-hot vector: the only non-zero element is the entry corresponding to the token, \mathbf{T} is sparse

Input embeddings



$$X = [x_1, x_2, \dots, x_m, \dots, x_N]$$

$$x_m \in \mathbb{R}^W$$

Concepts:

Input of length N , embedding of size W , embedding of the token (x_m)

Self-attention

Self-attention

$$\begin{aligned}
 \mathbf{s}_n[\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] &= \sum_{m=1}^N \underbrace{a[\mathbf{x}_m, \mathbf{x}_n]}_{\substack{\text{attention} \\ \text{(scalar weight)}}} (\underbrace{\boldsymbol{\theta}_0^v + \boldsymbol{\theta}_1^v \mathbf{x}_m}_{\substack{\text{value} \\ \text{(linear transformation of each } \mathbf{x}_m)}}) \\
 &= \sum_{m=1}^N a[\mathbf{x}_m, \mathbf{x}_n] \mathbf{v}_m
 \end{aligned}$$

self-attention
(connection between word representations)

$\mathbf{X} \in \mathbb{R}^{W \times N}$
 $\mathbf{x}_m \in \mathbb{R}^W$

$\mathbf{v}_m = \boldsymbol{\theta}_0^v + \boldsymbol{\theta}_1^v \mathbf{x}_m$
biases weights

$\boldsymbol{\theta}_0^v \in \mathbb{R}^W$
 $\boldsymbol{\theta}_1^v \in \mathbb{R}^{W \times W}$

$a[\cdot, \mathbf{x}_n] \geq 0$
 $\sum_{n=1}^N a[\cdot, \mathbf{x}_n] = 1$

Concepts:

Connections between the inputs, sparse attention coefficients, value (\mathbf{v}_m)

Attention

$$\mathbf{q}_n = \boldsymbol{\theta}_0^q + \boldsymbol{\theta}_1^q \mathbf{x}_n \quad \text{query}$$

$$\mathbf{k}_m = \boldsymbol{\theta}_0^k + \boldsymbol{\theta}_1^k \mathbf{x}_m \quad \text{key}$$

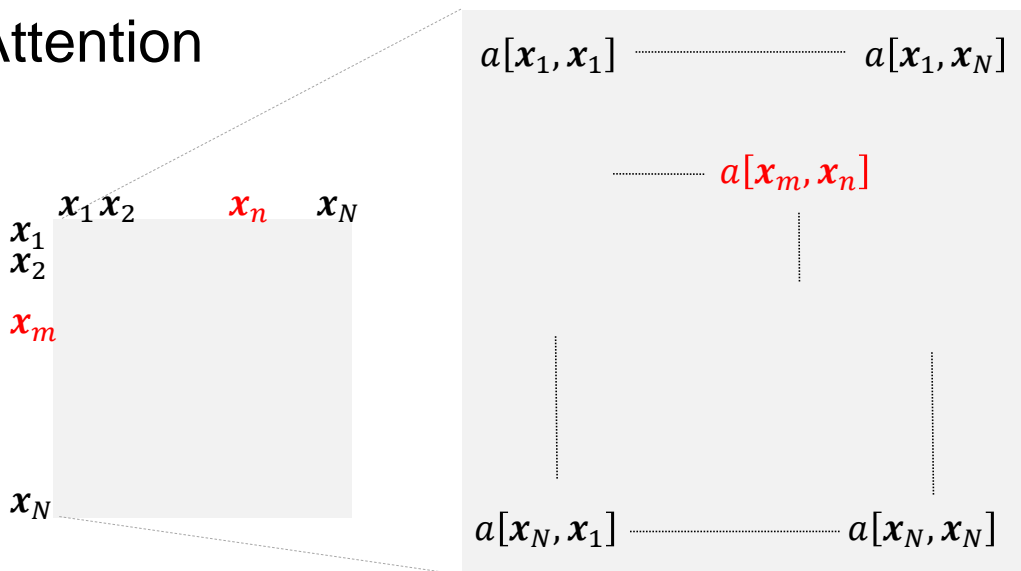
$(\mathbf{k}_m)^T \mathbf{q}_n$ measure of similarity (dot product)

$$a[\mathbf{x}_m, \mathbf{x}_n] = \text{softmax}_m \left((\mathbf{k}_*)^T \mathbf{q}_n \right) \\ = \frac{e^{(\mathbf{k}_m)^T \mathbf{q}_n}}{\sum_{m'=1}^N e^{(\mathbf{k}_{m'})^T \mathbf{q}_n}}$$

non-linear due to the dot product and the softmax

Concepts: Query and key from information retrieval, dot-product self-attention, measure of similarity between the input \mathbf{x}_n and all the other inputs

Attention



Concept:

Modeling data with *long-range* dependencies

Self-attention

$$\begin{array}{ccc}
 x_1 & s_1[x_1, x_2, \dots, x_N] = \sum_{m=1}^N a[x_m, x_1] v_m \\
 \vdots & \vdots \\
 x_n & s_n[x_1, x_2, \dots, x_N] = \sum_{m=1}^N a[x_m, x_n] v_m \\
 \vdots & \vdots \\
 x_N & s_N[x_1, x_2, \dots, x_N] = \sum_{m=1}^N a[x_m, x_N] v_m
 \end{array}$$

\downarrow scales linearly with N
 \downarrow quadratic dependence on N

Concepts: Value is a linear function of the input, same weights and biases applied to each input, attention weights are *sparse* and are a *non-linear* function of the input, parallelization

$$\begin{array}{l}
 X \in \mathbb{R}^{W \times N} \\
 x_m \in \mathbb{R}^W
 \end{array}$$

Self-attention: matrix form

$$\Theta = \{\Theta_0^v, \Theta_1^v, \Theta_0^q, \Theta_1^q, \Theta_0^k, \Theta_1^k\} \quad \text{shared set of parameters across inputs}$$

$$V(X) = \Theta_0^v \mathbf{1}_N^T + \Theta_1^v X$$

$$Q(X) = \Theta_0^q \mathbf{1}_N^T + \Theta_1^q X \quad \mathbf{1}_N: N\text{-dimensional vector containing only 1s}$$

$$K(X) = \Theta_0^k \mathbf{1}_N^T + \Theta_1^k X$$

$$\begin{aligned}
 S(X) &= V(X) \text{softmax} \left((K(X))^T Q(X) \right) \\
 &= V \text{softmax} \left((K)^T Q \right)
 \end{aligned}$$

scaled dot-product self-attention

$$S(X) = V \text{softmax} \left(\frac{(K)^T Q}{\sqrt{W_q}} \right)$$

\downarrow dimension of queries (and keys)

Concepts:

Scaled dot-product self-attention: facilitates training, avoids domination of the largest value

Multi-head self-attention

$$V_h(X) = \theta_{0h}^v \mathbf{1}_N^T + \theta_{1h}^v X$$

$$Q_h(X) = \theta_{0h}^q \mathbf{1}_N^T + \theta_{1h}^q X$$

$$K_h(X) = \theta_{0h}^k \mathbf{1}_N^T + \theta_{1h}^k X$$

$$\theta_h = \{\theta_{0h}^v, \theta_{1h}^v, \theta_{0h}^q, \theta_{1h}^q, \theta_{0h}^k, \theta_{1h}^k\}$$

$$S_h(X) = V_h \text{softmax}\left(\frac{(K_h)^T Q_h}{\sqrt{W_q}}\right)$$

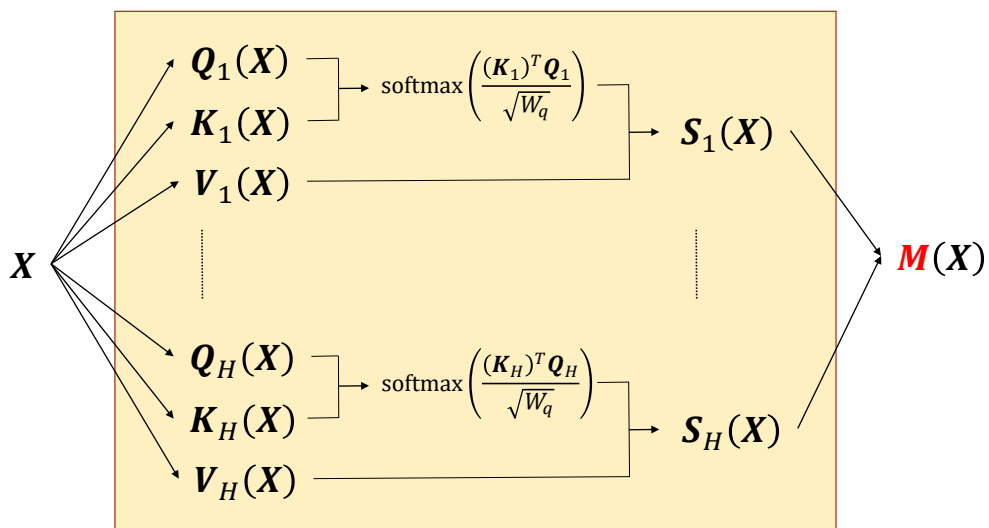
$$M(X) = \theta_M \left((S_1(X))^T, (S_2(X))^T, \dots, (S_H(X))^T \right)^T$$

θ_M linear transformation

Concepts:

Self-attention mechanisms (heads) applied in parallel, outputs are vertically concatenated

Multi-head self-attention



Transformer

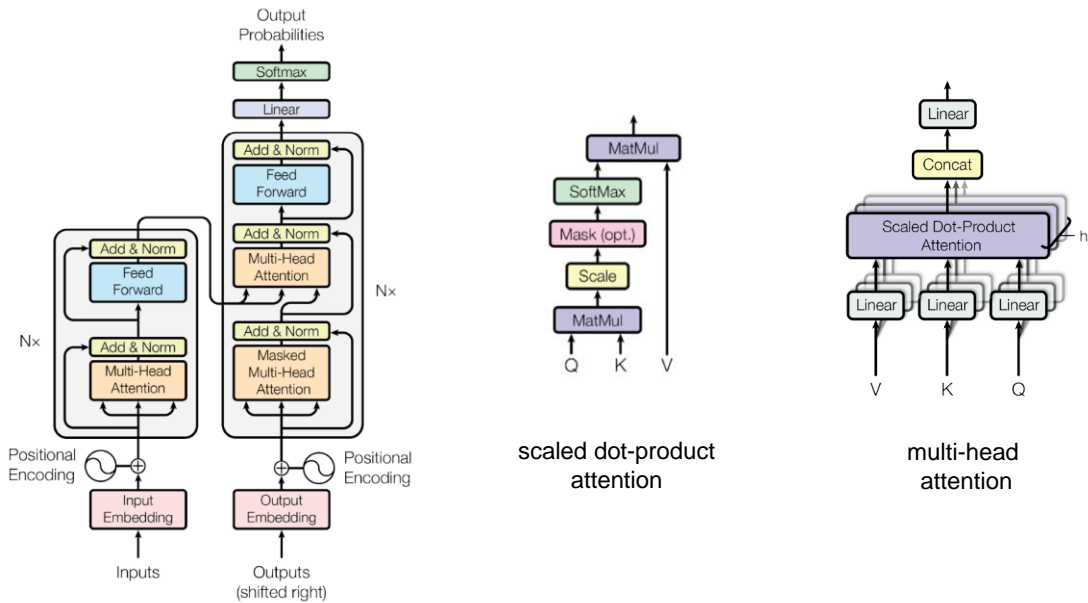
Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreit* Google Research usz@google.com
Llion Jones* Google Research llion@google.com	Aidan N. Gomez* [†] University of Toronto aidan@cs.toronto.edu	Lukasz Kaiser* Google Brain lukaszkaier@google.com	
Illia Polosukhin* [‡] illia.polosukhin@gmail.com			

Abstract

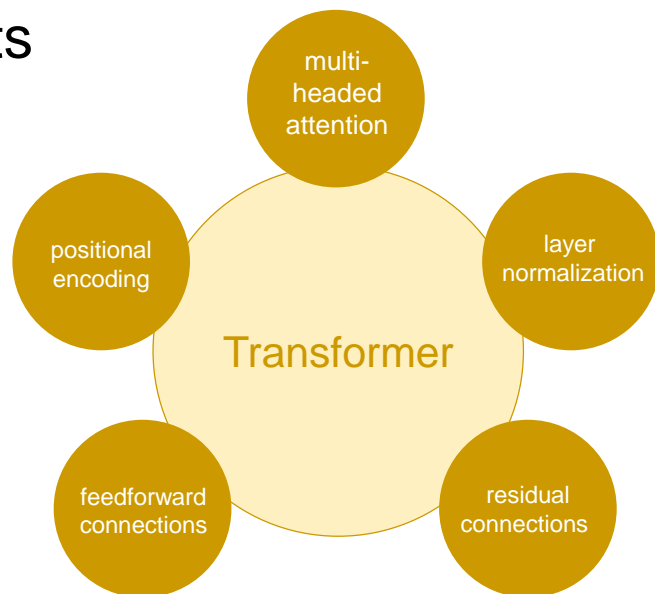
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

[arXiv:1706.03762](https://arxiv.org/abs/1706.03762)

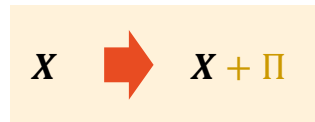


[arXiv:1706.03762](https://arxiv.org/abs/1706.03762)

Key components



Positional encoding

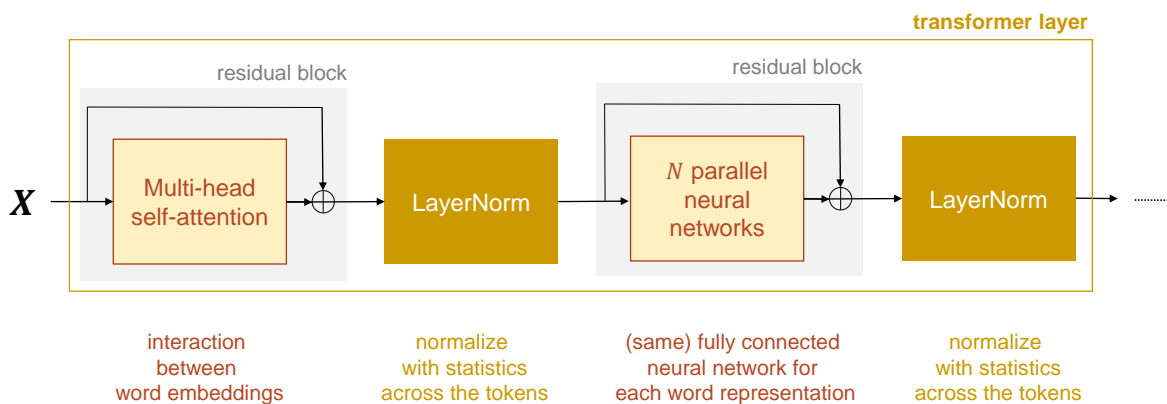


- Π → each **column** of Π is **unique** (information about the *position* in the input sequence)
- Π → **added** to the input only or at every layer
- Π → **learned** or “**hand-crafted**”

Concept:

Self-attention disregards the order of the tokens in a sentence

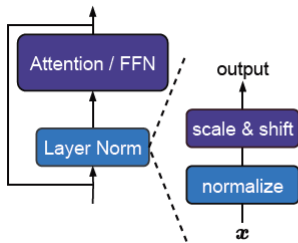
Transformer layer



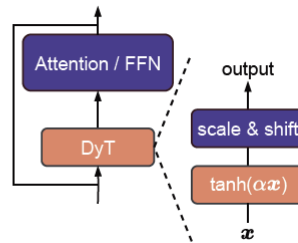
Concepts:

Output has the same size as the input ($W \times N$), residual block: output added to the original input

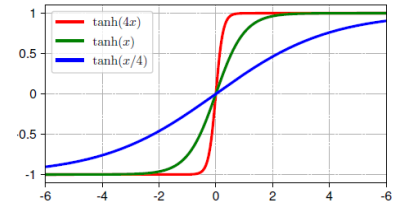
Layer normalization or Dynamic Tanh?



original transformer block
with layer normalization



block with Dynamic Tanh
(DyT) layer

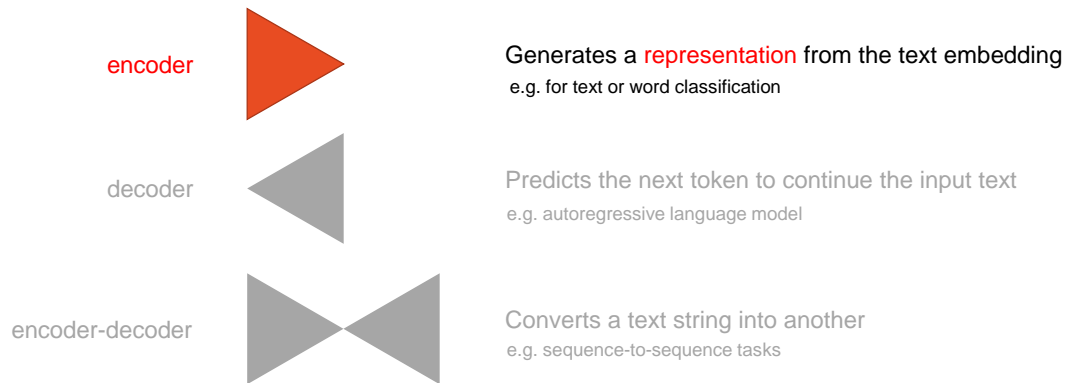


[arXiv:2503.10622](https://arxiv.org/abs/2503.10622)

Concept:
Transformers without normalization

Encoder model

Transformer models



Encoder model example: BERT

$P = 340,000,000$ parameters

$W = 1,024$ size of the embedding

$|\mathcal{V}| = 30,000$ size of the vocabulary (number of tokens)

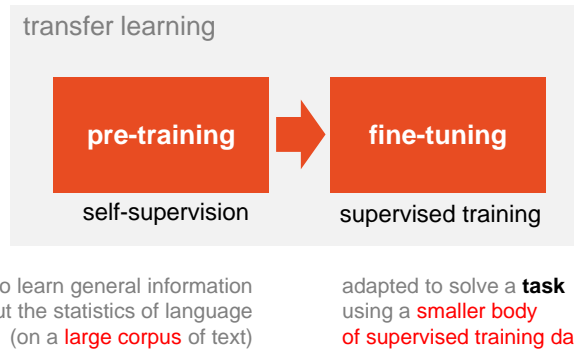
24 transformers (transformer layers)

$H = 16$ heads for the self-attention of each transformer

$W_v = W_q = W_k = 64$ dimension of the value, query, key

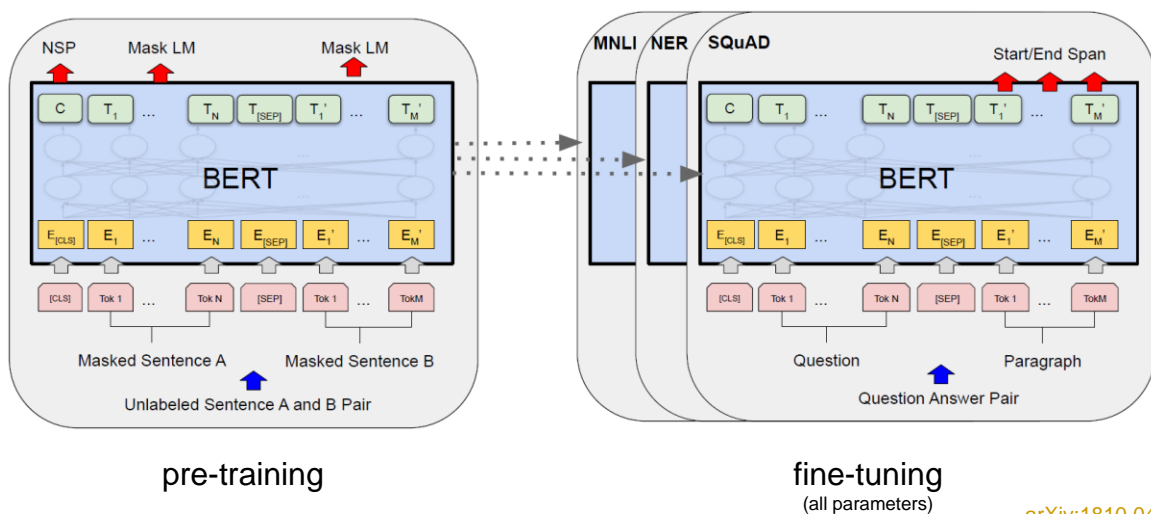
Concept: Bidirectional Encoder Representations from Transformers (BERT): jointly conditioning on both left and right context (in all layers)

Training BERT

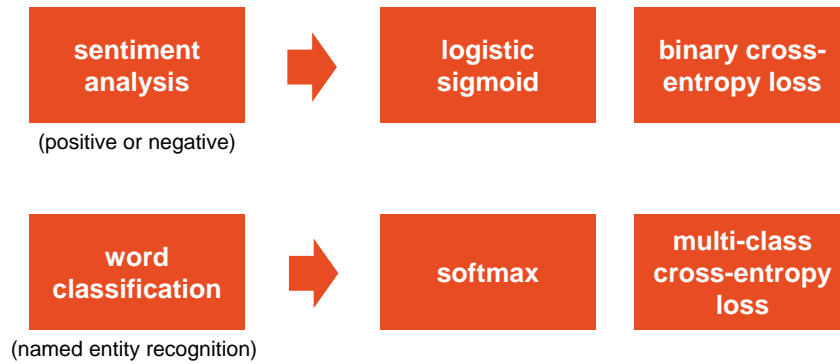


Concepts: Missing word prediction, understanding syntax, (statistical) understanding about the world, fine-tune with extra layer after the transformer network

BERT: pre-training and fine-tuning

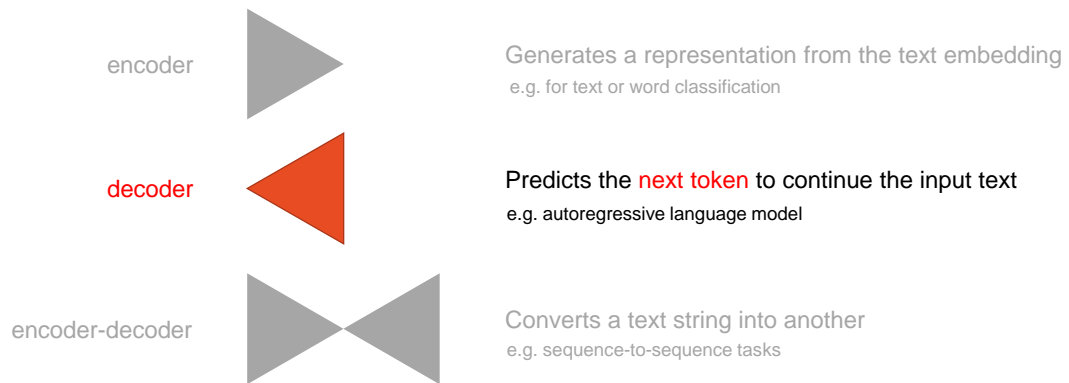


BERT: tasks



Decoder model

Transformer models



Decoder model example: GPT3

$P = 175,000,000,000$ parameters

$W = 12,288$ size of the embedding

$|\mathcal{V}| = 50,257$ size of the vocabulary (number of tokens)

96 transformers

$H = 96$ heads for the self-attention of each transformer

$W_v = W_q = W_k = 128$ dimension of the value, query, key

Context:

GPT4 has got $P = 1,800,000,000,000$ parameters

GPT3: next token prediction

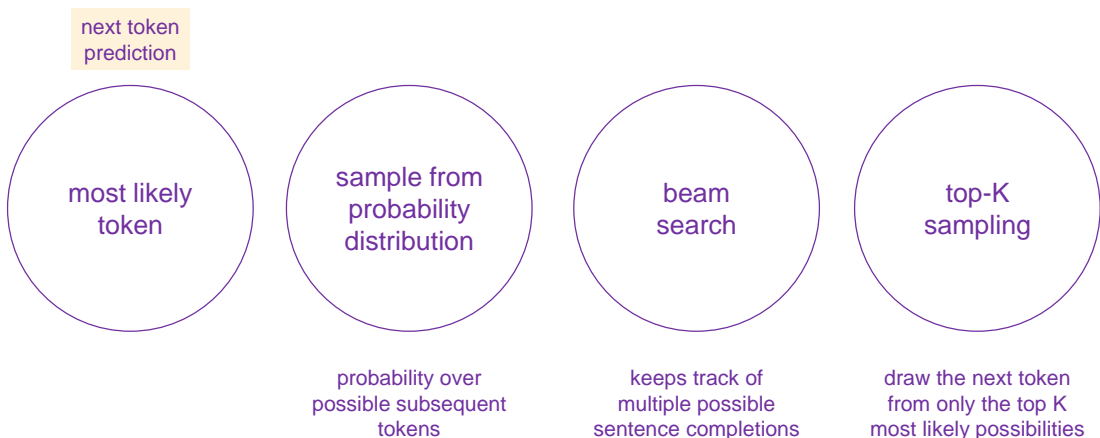
Each token t_n only interacts with previous tokens (**masked self-attention**)

$$P(\underbrace{t_1, t_2, \dots, t_{N-1}, t_N}_{\text{sentence}}) = P(t_1) \prod_{n=2}^N P(t_n | t_1, t_2, \dots, t_{n-1})$$

New extended sequence: fed back to the encoder for the **next token prediction**

Concepts: Autoregressive language model, predict the next token w/o access to the future, connection between maximizing the log probability of the tokens (loss function) and the next token prediction task

GPT3: generative model



slido



To understand transformers, was today's lecture more useful than prompting a large language model?

① Start presenting to display the poll results on this slide.

slido

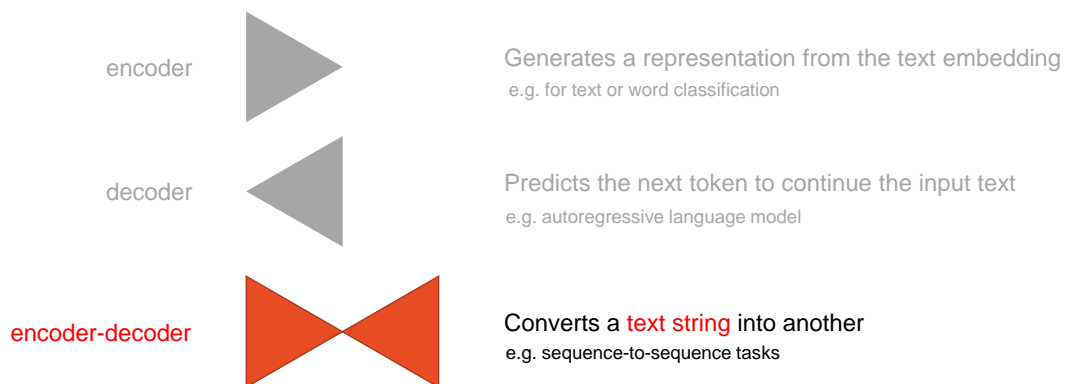


In terms of learning about transformers, how did today's lecture compare to using a large language model?

① Start presenting to display the poll results on this slide.

Encoder-decoder model

Transformer models



Example: machine translation

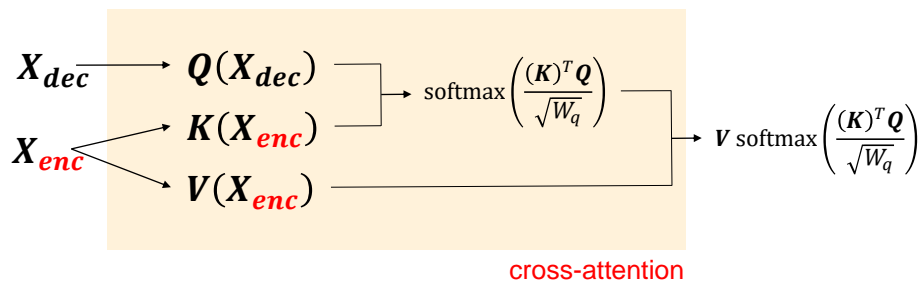
decoder layers also attend to the output of the **encoder**



each output token is *conditioned*
on the previous output tokens
and the source sentence

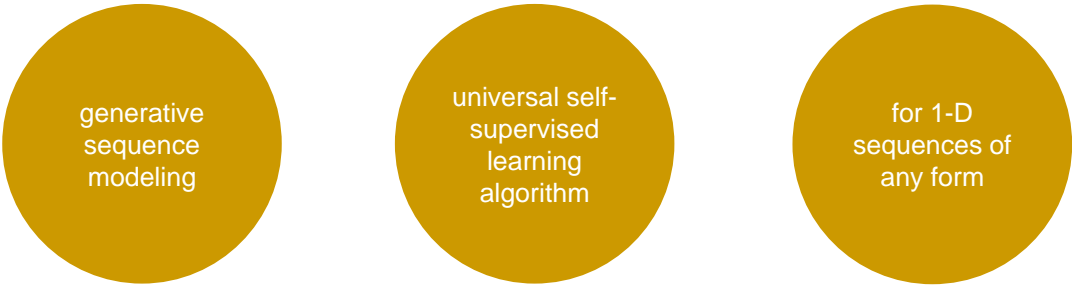
Concepts: Sequence-to-sequence task, information about the source language (encoder) and the target language statistics (decoder)

Encoder-decoder attention: cross attention



Concepts: Cross-attention (encoder-decoder attention), queries computed from the decoder embeddings, keys and values computed from the encoder embeddings

Transformers



generative
sequence
modeling

universal self-
supervised
learning
algorithm

for 1-D
sequences of
any form

Concepts: Transformer as *self-supervised* learning algorithm, sequences of *bytes*, trained to maximize the likelihood (mode covering), can be applied to any data type

Transformers for visual tasks

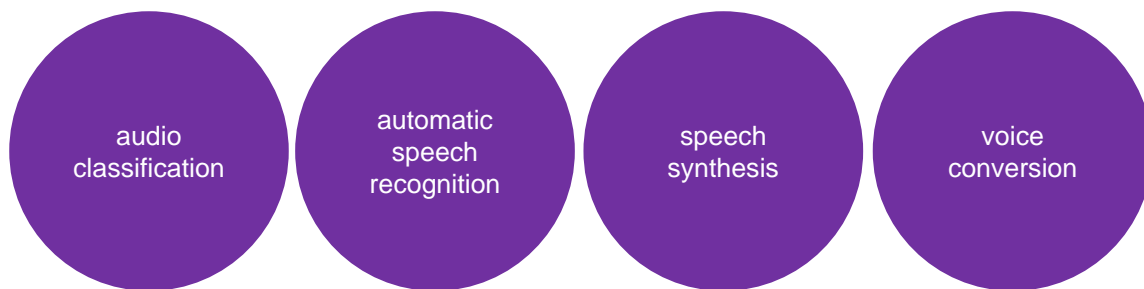


image
generation

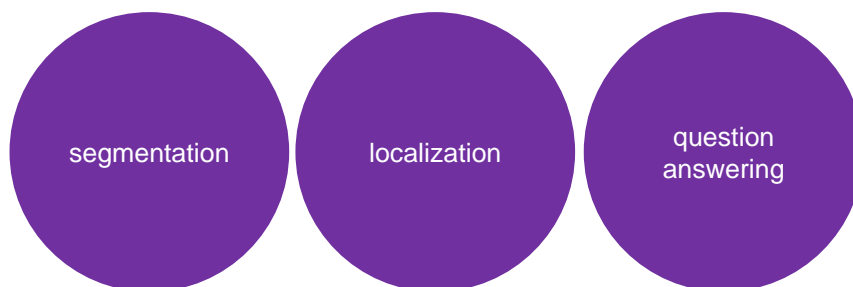
image
completion

image
classification

Transformers for audio tasks



Transformers for audio-visual tasks



Exercises

Today's exercises

Practice

You will become familiar with fine-tuning a **transformer** for binary text classification

Marked *(last one for Submission 1!)*

Includes **positional encoding**, scaled dot-product **attention** and multi-head **attention**

slido



What questions do you have regarding the mini-project?

① Start presenting to display the poll results on this slide.

What did we learn today?

- Natural Language Processing
- Tokens
- Self-attention
- Transformer
- Encoder model
- Decoder model
- Encoder-decoder model
- Exercises

EE-559

Deep Learning

andrea.cavallaro@epfl.ch