# Attention and BERT

Alexandre Allauzen

Fall 2023

# Roadmap

# Outline

# From Embeddings to Contextualized Embeddings

Static word embeddings

$$\mathbf{x} = \begin{bmatrix} \text{this,} & \text{movie,} & \text{was,} & \text{a,} & \text{great,} & \text{experience} \end{bmatrix}$$

$\mathbf{x}_{this}$ $\mathbf{x}_{movie}$ $\mathbf{x}_{was}$ $\mathbf{x}_a$ $\mathbf{x}_{great}$ $\mathbf{x}_{experience}$

Contextualized representation with bi-recurrent encoder:



$\mathbf{x}_{this}$ $\mathbf{x}_{movie}$ $\mathbf{x}_{was}$ $\mathbf{x}_a$ $\mathbf{x}_{great}$ $\mathbf{x}_{experience}$

# Draw attention for classification



- $\mathbf{a} = (a_i)$, $\sum_{i=1}^{L} a_i = 1$ and $0 \leq a_i \leq 1$
- $\mathbf{a}$ : attention vector for the "query" $\mathbf{q}$ and the "keys" $\mathbf{X}$.
- $\mathbf{q}$ is a vector to be learnt [8, 5]

Introduction to attention

# Attention to weight inputs

- $\mathbf{a} = X\mathbf{q}$ is the attention vector

$$\mathbf{h} = \sum_{i=1}^{L} a_i \mathbf{x}_i = \mathbf{aX}$$

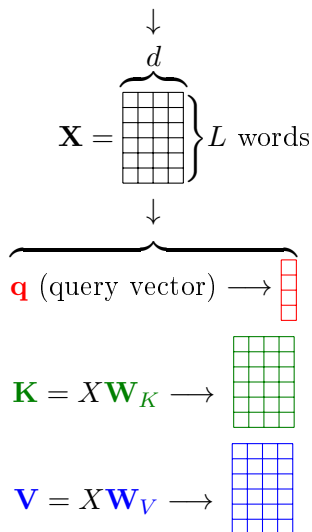- A new vector, focused on the classification task ($\mathbf{q}$)
- To summarize:

$$\mathbf{h} = \mathrm{softmax}(\mathbf{Xq})\mathbf{X} \rightarrow \text{ classification}$$

Issues:
- Scale the dot product
- X is involved everywhere !

# Basic attention mechanism for classification

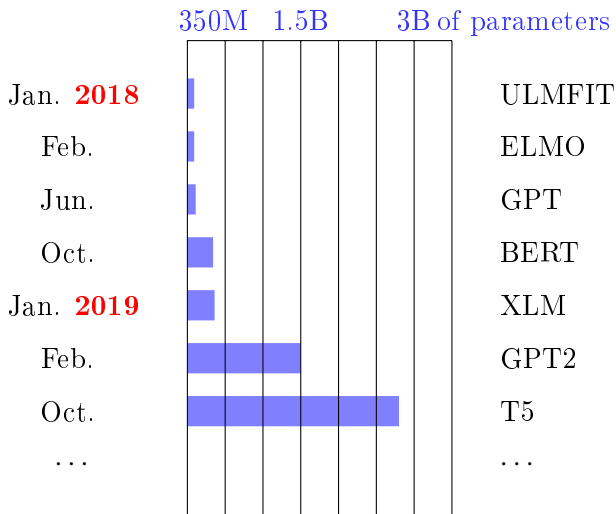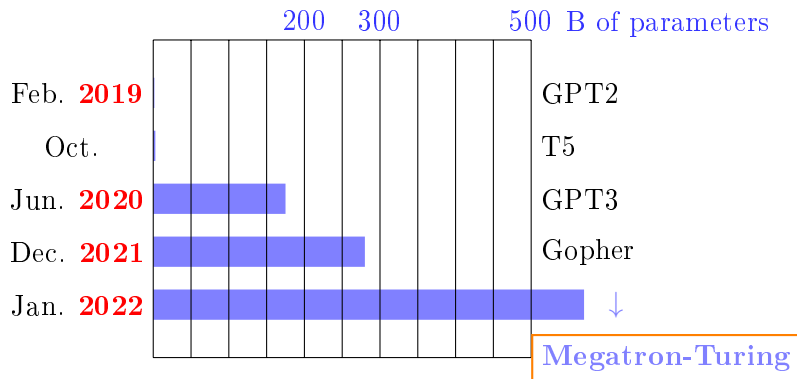this movie was a great experience

$\downarrow$



$$\mathbf{h} = \text{softmax}\left(\frac{\mathbf{Kq}}{\sqrt{d}}\right)^t \mathbf{V}$$

- X can be static emb.
- Derived from bi-LSTM
- $\mathbf{q}$ is learnt as a target for selection
- $pa = \mathbf{Kq}$: selection in $\mathbf{V}$

$\mathbf{X} =$ $\left.\begin{array}{}\\\\\\\end{array}\right\} L$ words

$\mathbf{q}$ (query vector) $\longrightarrow$

$\mathbf{K} = X\mathbf{W}_K \longrightarrow$

$\mathbf{V} = X\mathbf{W}_V \longrightarrow$

# In a few dates



| | 350M | 1.5B | | 3B of parameters | |
|---|---|---|---|---|---|
| Jan. **2018** | | | | | ULMFIT |
| Feb. | | | | | ELMO |
| Jun. | | | | | GPT |
| Oct. | | | | | BERT |
| Jan. **2019** | | | | | XLM |
| Feb. | | | | | GPT2 |
| Oct. | | | | | T5 |
| . . . | | | | | . . . |

# Bigger is . . .



200    300                500 B of parameters

| | |
|---|---|
| Feb. **2019** | GPT2 |
| Oct. | T5 |
| Jun. **2020** | GPT3 |
| Dec. **2021** | Gopher |
| Jan. **2022** | ↓ |
| | **Megatron-Turing** |

# Outline

# Contextualized word embeddings

Consider the word <span style="color:red">driver</span>:

| the | audio | driver | is | really | outdated |
|-----|-------|--------|-----|--------|----------|
| the | driver | exceeded | the | speed | limit |

# Contextualized word embeddings

Consider the word <span style="color:red">driver</span>:

| the | audio | <span style="color:red">driver</span> | is | really | outdated |
|-----|-------|--------|-----|--------|----------|
| the | <span style="color:red">driver</span> | exceeded | the | speed | limit |

## The context

| The | | | The | | $\lambda_{1,2}$ |
|-----|--|--|-----|--|------|
| audio | | | <span style="color:red">driver</span> | | $\lambda_{2,2}$ |
| <span style="color:red">driver</span> | | | exceeded | | $\lambda_{3,2}$ |
| is | | | the | | $\lambda_{4,2}$ |
| really | | | speed | | $\lambda_{5,2}$ |
| outdated | | | limit | | $\lambda_{6,2}$ |

# Self attention

Consider the word driver:



the   driver   exceeded   the   speed   limit

Embeddings: $\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{x}_3$   $\mathbf{x}_4$   $\mathbf{x}_5$   $\mathbf{x}_6$

Attention: $\lambda_{1,2}$   $\lambda_{2,2}$   $\lambda_{3,2}$   $\lambda_{4,2}$   $\lambda_{5,2}$   $\lambda_{6,2}$

Output: $\mathbf{z}_2 = f\left(\sum_i [\lambda_{i,2} \times g(\mathbf{x}_i)]\right)$

- $(\lambda_{i,j})$ are the attention coefficients, $\sum_i \lambda_{i,j} = 1$, and
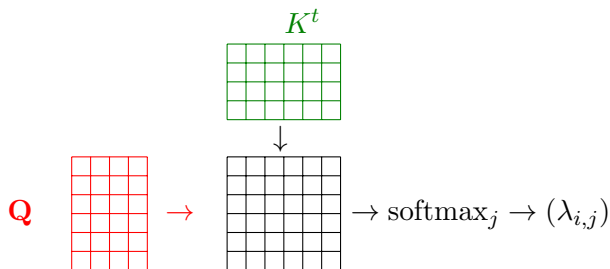- Reflects the influence of $\mathbf{x_i}$ on $\mathbf{x_j}$ (transformed version)

# Transformer : Queries, Keys, Values

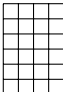# Tranformer : Attention matrix

The distance matrix between $Q$ and $K$



Scaled Dot-Product Attention

$$\mathbf{Z} = \text{softmax}\left(\frac{\textcolor{red}{\mathbf{Q}}\textcolor{green}{\mathbf{K^t}}}{\sqrt{d}}\right)\textcolor{blue}{\mathbf{V}} =$$
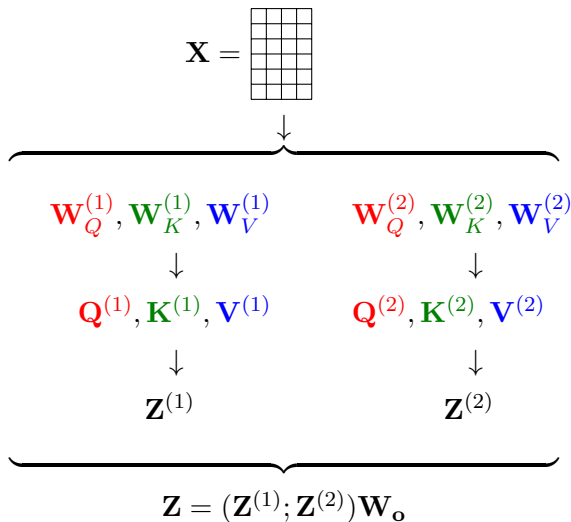
# QKV and Metric Learning

$$\mathbf{Q}\mathbf{K}^t = \mathbf{X}\mathbf{W}_K \times (\mathbf{X}\mathbf{W}_K)^t = \mathbf{X}\mathbf{W}_Q \times (\mathbf{W}_K^t \mathbf{X}^t)$$
$$= \mathbf{X}\mathbf{M}\mathbf{X}^t$$

- If $\mathbf{M}$ would be PSD, it is a metric.
- Otherwise, it is a transformed similarity (bilinear similarity)

$\mathbf{M}$ is learnt: a transformer block learns its own similarity.
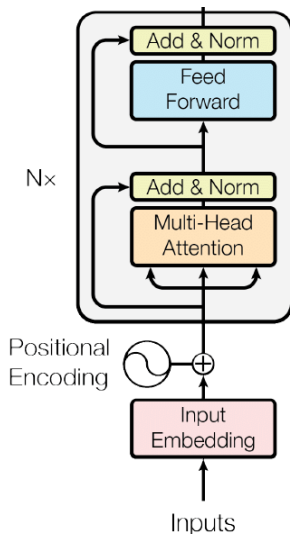
# Multi-head attention (with 2 heads)

$$\mathbf{X} = \begin{array}{|c|c|c|c|}\hline \phantom{x} & & & \\\hline & & & \\\hline & & & \\\hline & & & \\\hline & & & \\\hline\end{array}$$

$$\downarrow$$

$$\mathbf{W}_Q^{(1)}, \mathbf{W}_K^{(1)}, \mathbf{W}_V^{(1)} \qquad \mathbf{W}_Q^{(2)}, \mathbf{W}_K^{(2)}, \mathbf{W}_V^{(2)}$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$\mathbf{Q}^{(1)}, \mathbf{K}^{(1)}, \mathbf{V}^{(1)} \qquad \mathbf{Q}^{(2)}, \mathbf{K}^{(2)}, \mathbf{V}^{(2)}$$

$$\downarrow \qquad\qquad\qquad \downarrow$$

$$\mathbf{Z}^{(1)} \qquad\qquad\qquad \mathbf{Z}^{(2)}$$

$$\mathbf{Z} = (\mathbf{Z}^{(1)}; \mathbf{Z}^{(2)})\mathbf{W_o}$$

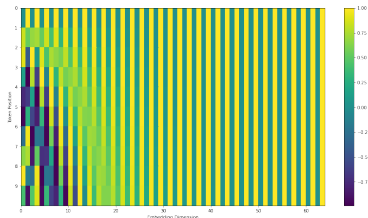# Putting all together (with more tricks)

**Transformer block**
From [7]

- Inputs is $\mathbf{X}$
- Positional embeddings
- Multihead attention
- Residual connections [4]
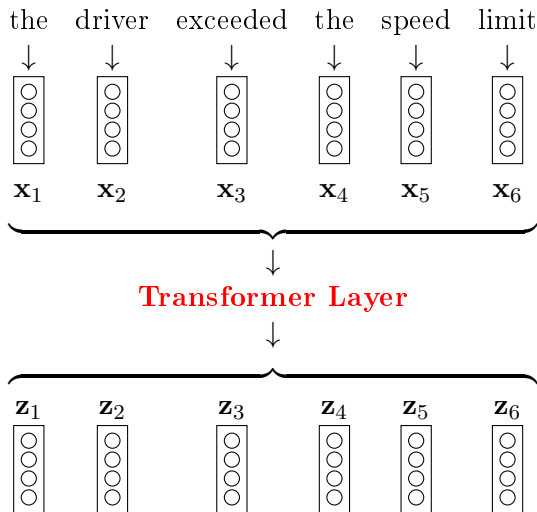- Layer Normalization [2]
- Final filtering

# Positional embeddings



- Originally "absolute"
- Can be learnt [3, 1]
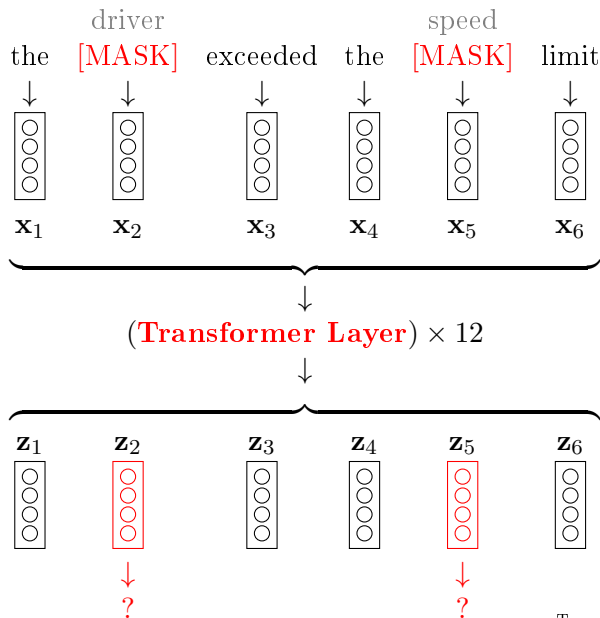- Or relative [6]

(figure generated by the following code
`https://github.com/jalammar/jalammar.github.io/blob/master/notebookes/`
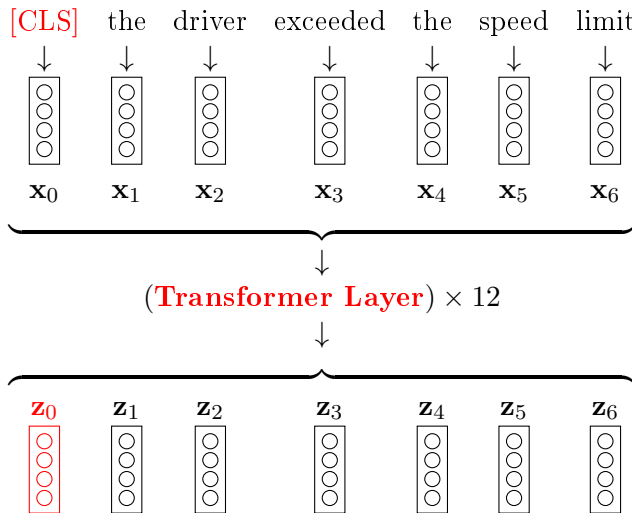`transformer/transformer_positional_encoding_graph.ipynb`)

# A Transformer layer



Transformer layers can be stacked !
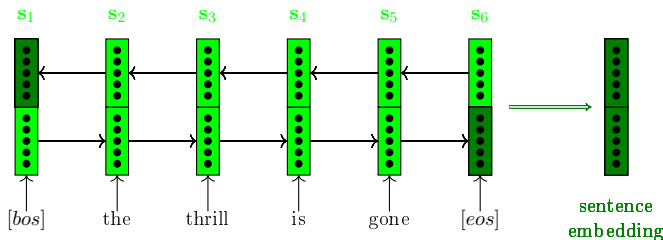
# Pre-training as a (Masked) language model

# BERT Encoder for text classification

# Transformers / bi-lstm encoders

## Reminder of bi-recurrent encoder



$\mathbf{s}_1$  $\mathbf{s}_2$  $\mathbf{s}_3$  $\mathbf{s}_4$  $\mathbf{s}_5$  $\mathbf{s}_6$

[*bos*]  the  thrill  is  gone  [*eos*]

sentence embedding

## The difference

- Two different ways to encode the dependence
- Richer for attention since we stack transformers
- all the deep-learning tricks ⇒ over-parametrization

# Outline

# Transformers are everywhere

### State of the art encoder
- For text !
- And also for speech, DNA, vision, . . .

### Also a powerful generator
- For text (GPT, . . .)
- Speech, . . .  sequences

# Outline

[1] Rami Al-Rfou et al. *Character-Level Language Modeling with Deeper Self-Attention*. 2018. arXiv: 1808.04444 [cs.CL].

[2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. *Layer Normalization*. 2016. arXiv: 1607.06450 [stat.ML].

[3] Jonas Gehring et al. "Convolutional Sequence to Sequence Learning". In: *CoRR* abs/1705.03122 (2017). arXiv: 1705.03122. URL: http://arxiv.org/abs/1705.03122.

[4] Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 770–778. arXiv: 1512.03385 [cs.CV].

[5] Zhouhan Lin et al. "A STRUCTURED SELF-ATTENTIVE SENTENCE EMBEDDING". In: *International Conference on Learning Representations*. 2017. URL: https://openreview.net/forum?id=BJC_jUqxe.

[6] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. *Self-Attention with Relative Position Representations*. 2018. arXiv: 1803.02155 [cs.CL].

[7] Ashish Vaswani et al. "Attention is All you Need". In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. Curran Associates, Inc., 2017, pp. 6000–6010. URL: http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf.

[8] Zichao Yang et al. "Hierarchical Attention Networks for Document Classification". In: *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*.