

Cursos Extraordinarios

Verano 2024

**“Inteligencia Artificial y Grandes
Modelos de Lenguaje: Funcionamiento,
Componentes Clave y Aplicaciones”**

Zaragoza, del 3 al 5 de julio

Transformer

- **Arquitectura**
 - Residual layers
 - Layer normalization
 - Multi-head attention
 - Feed forward
 - Positional embedding
 - Embedding de entrada
 - Arquitectura general



Transformer

- **Attention is all you need (Vaswani 2017)** > 100k cites
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N. Kaiser L, Polosukhin, I. (2017) Attention is all you need. Advances in neural information processing systems, 30, 5998-6008..

- A new type of sequence models
- The basis model for most state of the art results nowadays:
 - Translation
 - Language modeling
 - NLP: question answering, summary
 - Language understanding,
 - Speech recognition (ASR), synthesis
 - Image recognition, segmentation, detection
 - Multimodal image+text
 - DNA folding
 - Drug discovery
 - ...

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
lukaszkaizer@google.com

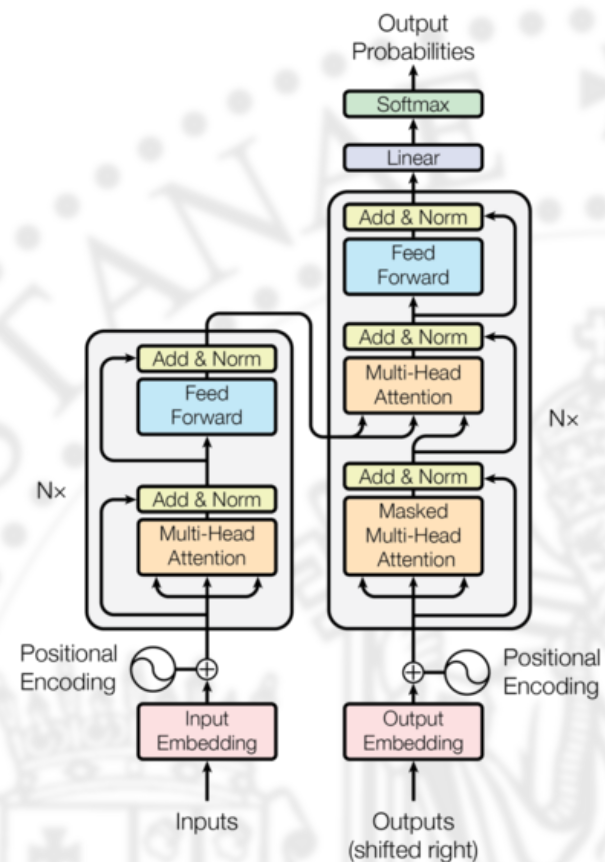
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

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.

Transformer

- Transformer
 - Architecture:
 - Encoder / Decoder**
 - Similar to **seq2seq** architecture
 - Some application only use Encoder
 - Encoder creates a processed representation for the input sentence
 - Decoder predicts next symbol/word given previously generated
 - No LSTMs or convolutions (origin of the title)
 - Self attention** (multihead) is the main computational unit:



Transformer

- **Arquitectura**
 - **Residual layers**
 - Layer normalization
 - Multi-head attention
 - Feed forward
 - Positional embedding
 - Arquitectura general

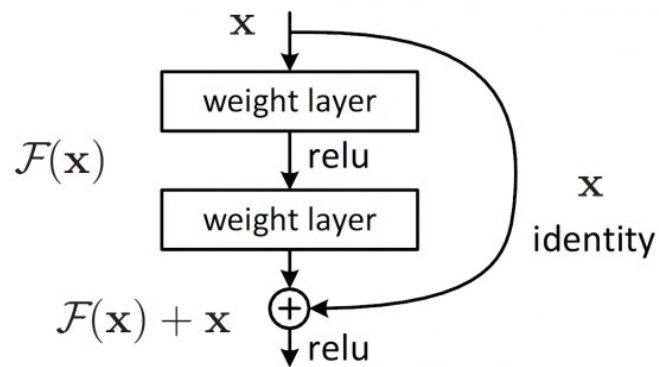


Residual layers

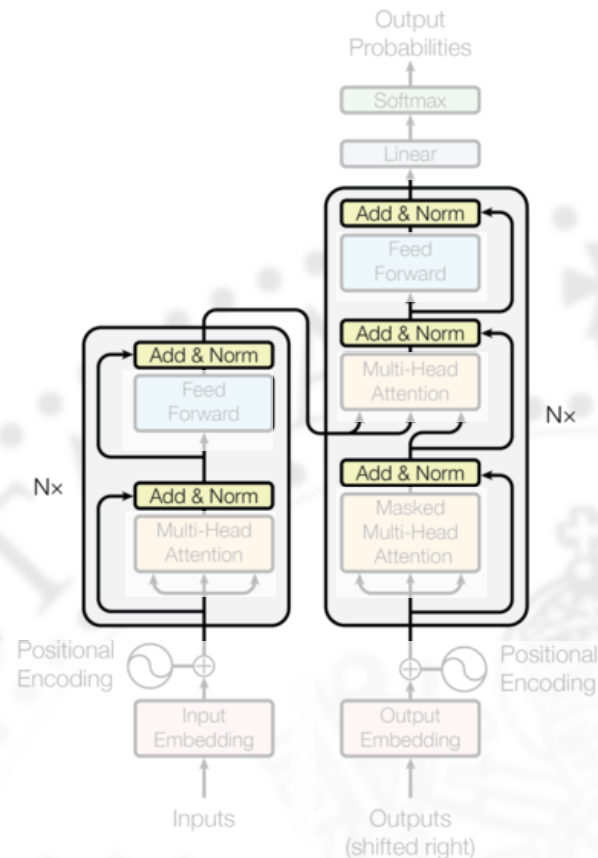
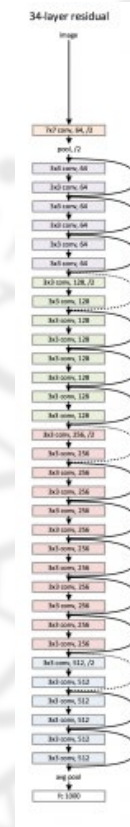
Residual layers (121k citas)

He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 770-778).

- The core idea is to modify the representation incrementally



- The vanishing gradient problem prevented deeper networks:
 - In the paper they present networks from 28 to 152 layers



Transformer

- **Arquitectura**
 - Residual layers
 - **Layer normalization**
 - Multi-head attention
 - Feed forward
 - Positional embedding
 - Embedding de entrada
 - Arquitectura general



Layer normalization

■ Normalization

■ Layer normalization (Ba 2016)

Ba, J. L., Kiros, J. R., & Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv:1607.06450*.

■ Instance Normalization(Wu 2016)

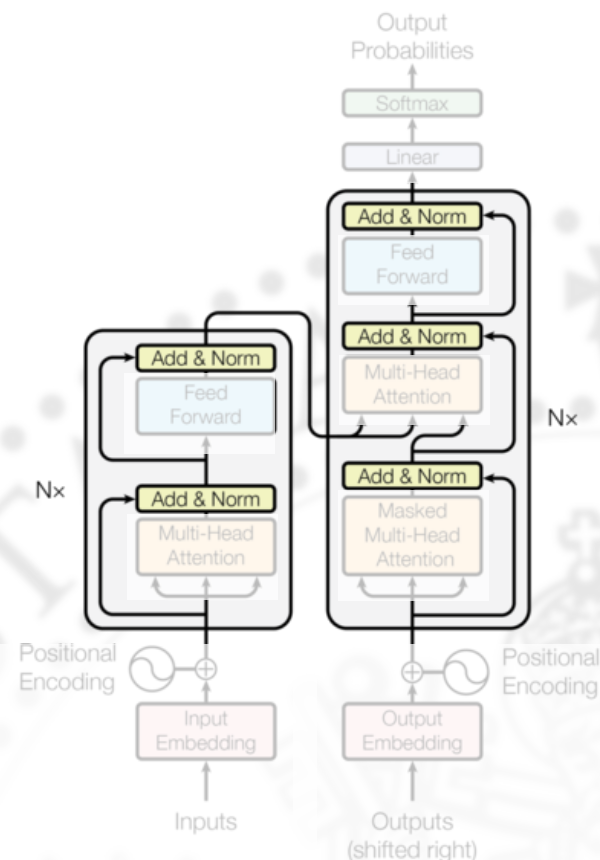
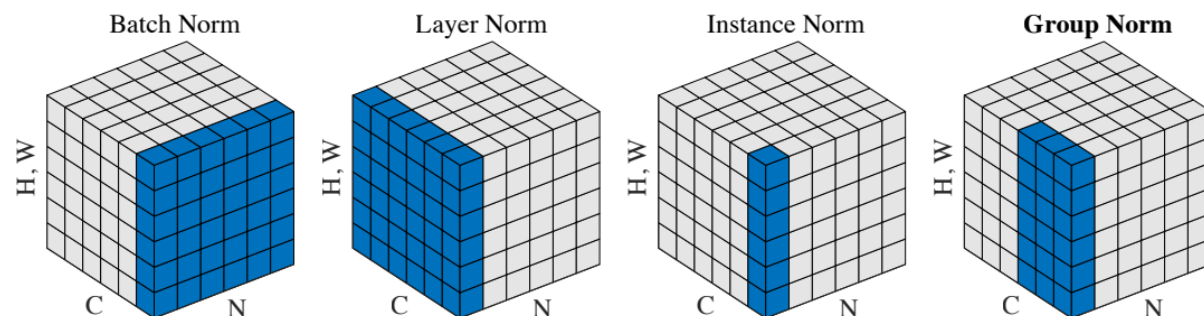
Ulyanov, D., Vedaldi, A., & Lempitsky, V. (2016). Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*.

■ Group Normalization(Wu 2018)

Wu, Y., & He, K. (2018). Group normalization. In *Proceedings of the European conference on computer vision (ECCV)* (pp. 3-19).

■ Layer normalization

- The layer normalization is an alternative to the batch normalization
- The objective and mechanism is similar to BN
- The mean and std is calculated over the signal **channel dimensions** instead



Layer normalization

■ Layer normalization

■ Layer normalization

- The mean and std is calculated over the signal **channel dimensions** instead of minibatch

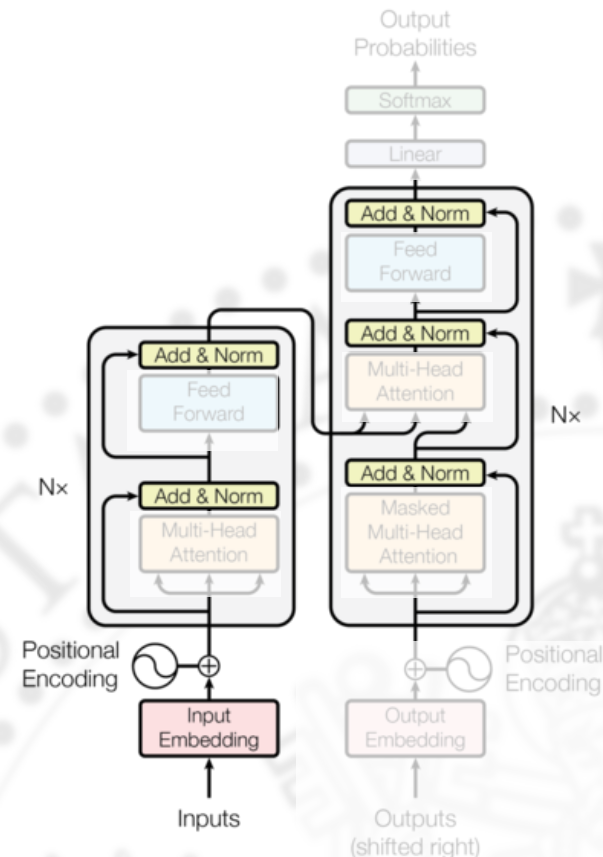
■ Steps

- 1. **Normalize** signal channel dimension
- 2. Apply **linear transformation** with learned parameters

- Special case: RMSNorm, no bias is taken into account

```
class RMSNorm(torch.nn.Module):
    def __init__(self, d_model=512, eps=1e-8, **kwargs):
        super().__init__()
        self.eps = eps
        self.scale = torch.nn.Parameter(torch.ones(1, 1, d_model))

    def forward(self, x):
        x = x / (torch.sqrt(torch.mean(x**2, dim=-1, keepdim=True)) + self.eps)
        return x * self.scale
```



Transformer

- **Arquitectura**
 - Residual layers
 - Layer normalization
 - **Multi-head attention**
 - Feed forward
 - Positional embedding
 - Embedding de entrada
 - Arquitectura general



Multi-head Attention (MHA)

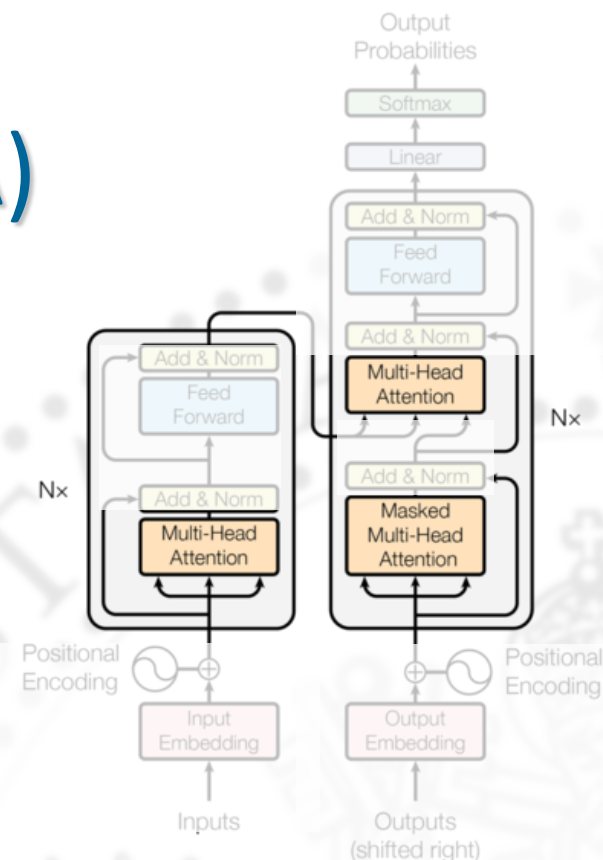
- **Multi-head Attention**
- MHA mechanism
 - Attention mechanism
 - Self attention
 - Input: sequence size ($n \times d$)
 - Output: sequence size ($n \times d$)
 - Cross attention
 - Input: sequence size ($n \times d$)
external signal ($m \times d$)
 - Output: sequence size ($n \times d$)

Output is a linear combination of inputs

Similar to a convolution (dynamic filter)

The decoder has a causal restriction

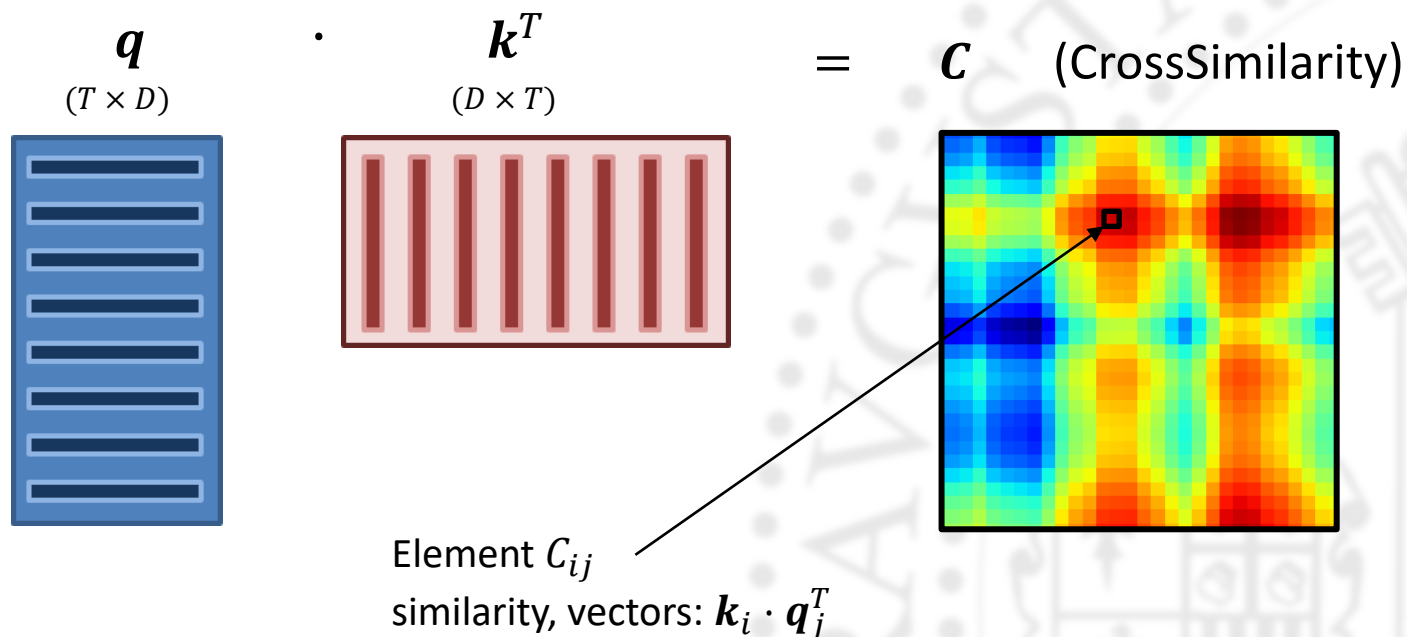
- Heads $\rightarrow h$ (from 1 to 96)
 - parallel Attention modules
 - Operation steps:
 - **Split** the vector in h parts
 - **Apply** Attention to each part
 - **Concatenate**



Multi-head Attention (MHA)

Self attention

- The core of the transformer is the self-attention operation
 - Given two sequences of vectors, k , q , of length T and dimension D
 - The product of the two matrices is used to find **similarities** between the sequences
 - In signal processing, similar to the idea of **Cross-correlation**



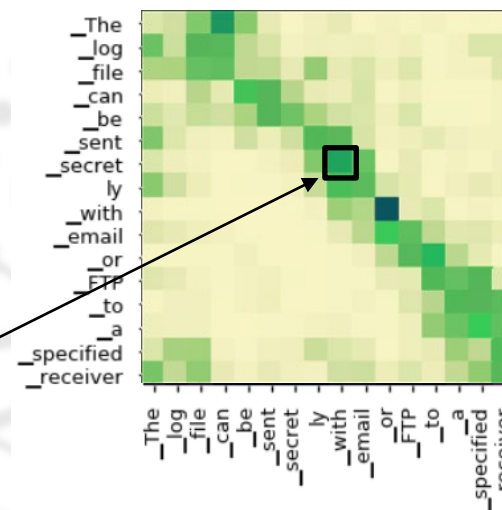
Multi-head Attention (MHA)

Self attention

- The core of the transformer is the self-attention operation
 - Given two sequences of vectors, k , q , of length T and dimension D
 - The product of the two matrices is used to find similarities between the sequences
 - A softmax operation is used to produce **mappings** of input frames to outputs (alignment matrix)

$$\text{softmax} \left(q \cdot k^T \right) = A \quad (\text{Self attention})$$

$(T \times D)$ $(D \times T)$



$$\sum_j A_{ij} = 1$$

(sum 1 cols.)

Element A_{ij}
Degree of influence of input
sequence item j over output i

<https://nlp.seas.harvard.edu/2018/04/03/attention.html>

Multi-head Attention (MHA)

Self attention

- The core of the transformer is the self-attention operation
 - Given two sequences of vectors, \mathbf{k}, \mathbf{q} , of length T and dimension D
 - The product of the two matrices is used to find similarities between the sequences
 - A softmax operation is used to produce **mappings** of input frames to outputs (alignment matrix)
 - The mapping is used to produce the output of the layer after multiplying by the values matrix
 - The process can be summarized as: the application of a dynamical mapping to the matrix \mathbf{v}

$$\text{softmax}(\mathbf{q} \cdot \mathbf{k}^T) \cdot \mathbf{v} = \mathbf{o} \text{ (output)}$$


 $(D \times T)$


- The output sequence can be **reordered**
- Each output item can be dependent of **any part of the input** sequence

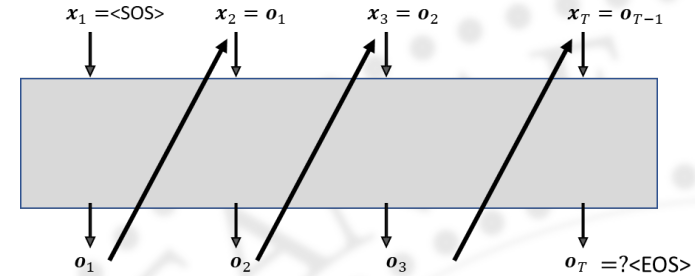

 $(T \times D)$


- The final output is the mapping/alignment multiplied by the sequence
- The mapping/alignment process can be interpreted as an **attention mechanism**

Multi-head Attention (MHA)

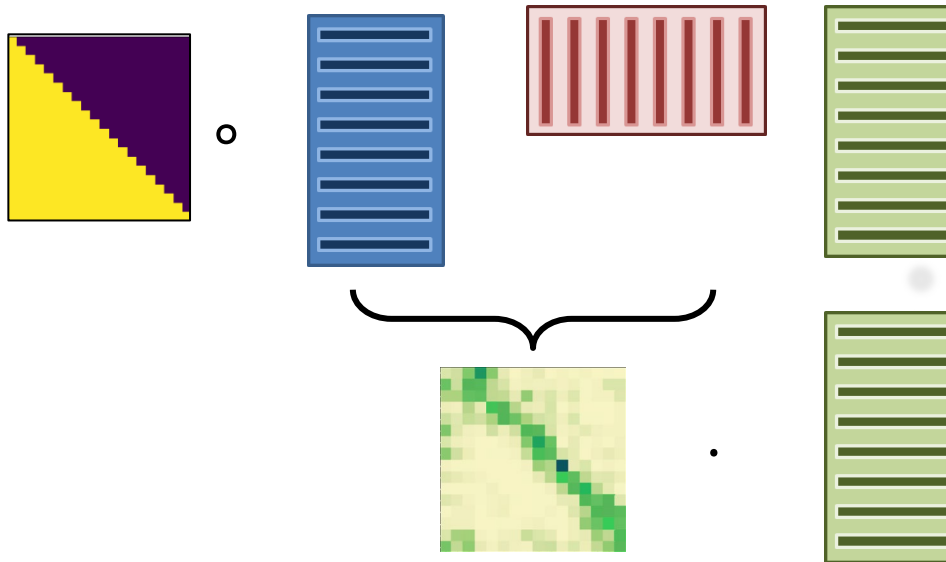
■ Causal self attention

- To operate in generators of seq2seq architectures
 - The objective is to predict next item given previous
 - The problem if the model is not modified is that self attention allows a trivial shortcut from future token present in the input to the Desired output
 - This is solved by multiplying the attention alignment with a mask matrix (upper diagonal values = 0)



$$\text{softmax}(M \circ (q \cdot k^T)) \cdot v = o \text{ (output)}$$

$(T \times D)$ $(D \times T)$



Multi-head Attention (MHA)

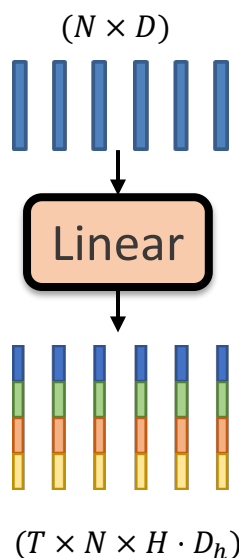
Multi-Head attention

- The representation dimension: split in h subvectors: k_h, q_h, v_h
- The self-attention is then used with each subvector

$$\text{Attention}(q_h, k_h, v_h) = \text{softmax}\left(\frac{q_h \cdot k_h^T}{\text{scale}}\right) \cdot v_h$$

scale constant (scale = $1 / \sqrt{D_h}$) plays the role of temperature

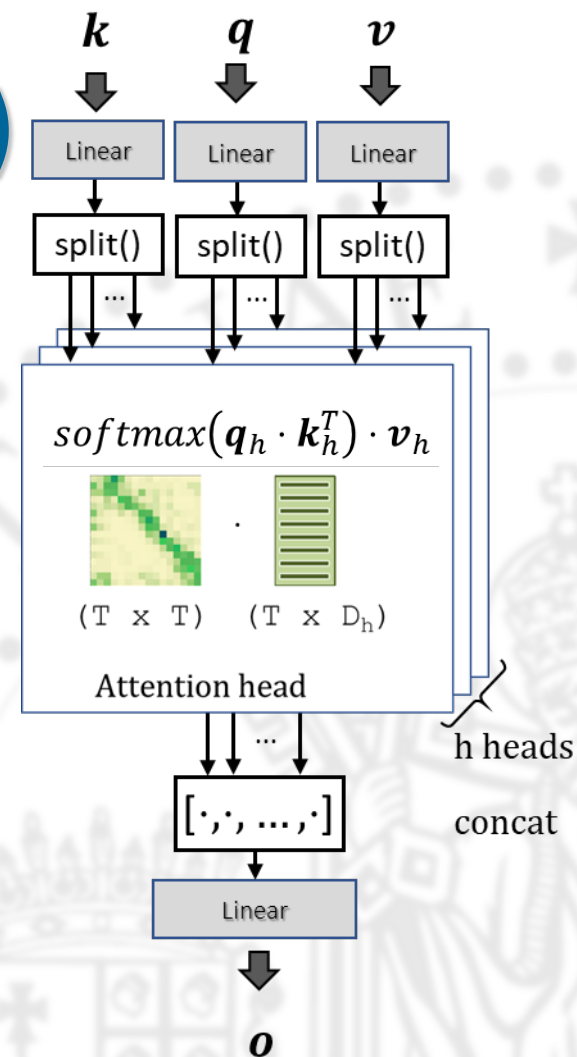
- The resulting sequences are **concatenated**
- Finally a linear layer is applied to obtain the output



```
self.linear = nn.Linear(d, h*dh)
```

```
q = self.linear(q).reshape(b, n, h, dh)
```

split the output in h subvectors of dimension dh

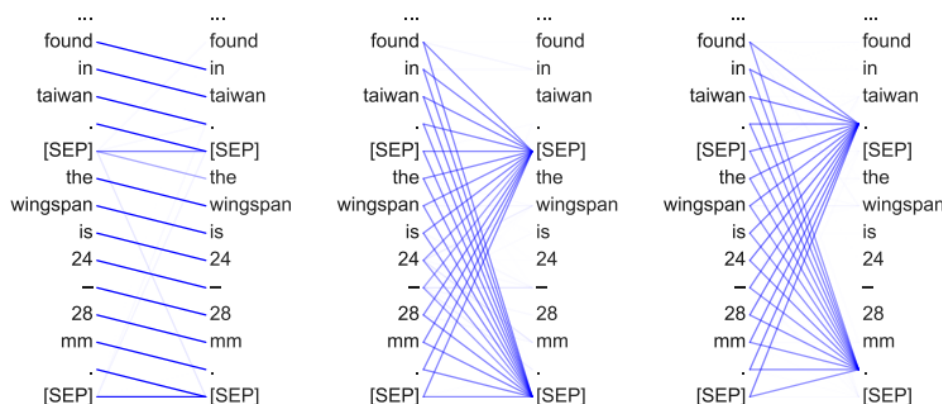
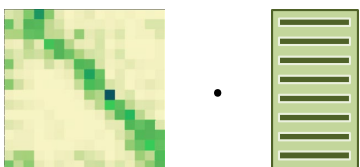


Multi-head Attention (MHA)

Self attention

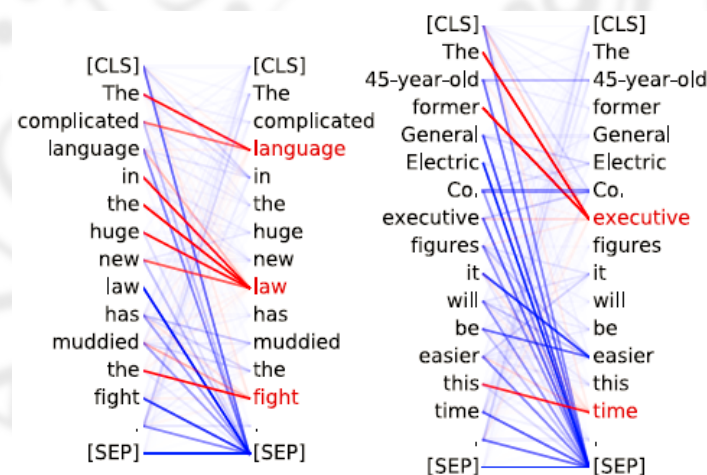
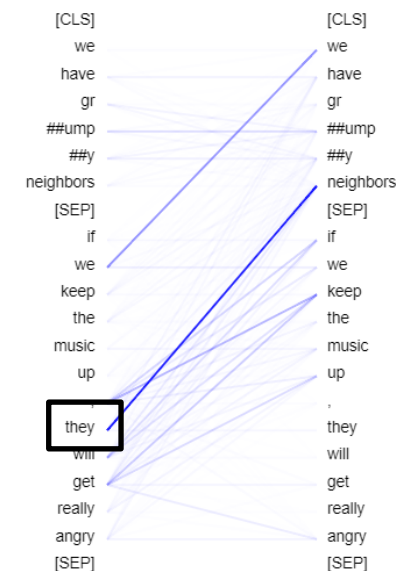
- There have been many studies trying to analyze the meaning of the mappings
 - One of the classic problems of NLP has been to determine automatically the relationship of a pronoun with previous text
 - e. g. **they** in the example on the left

$$\text{softmax}(q \cdot k^T) \cdot v$$



Attention to previous item

Attention to end of sentence



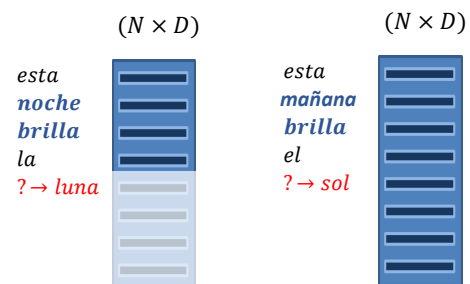
Noun modifiers, e.g. determiners attend to their noun

<https://medium.com/synapse-dev/understanding-bert-transformer-attention-isnt-all-you-need-5839ebd396db>

<https://medium.com/dair-ai/aspects-of-language-captured-by-bert-32bc3c54016f>

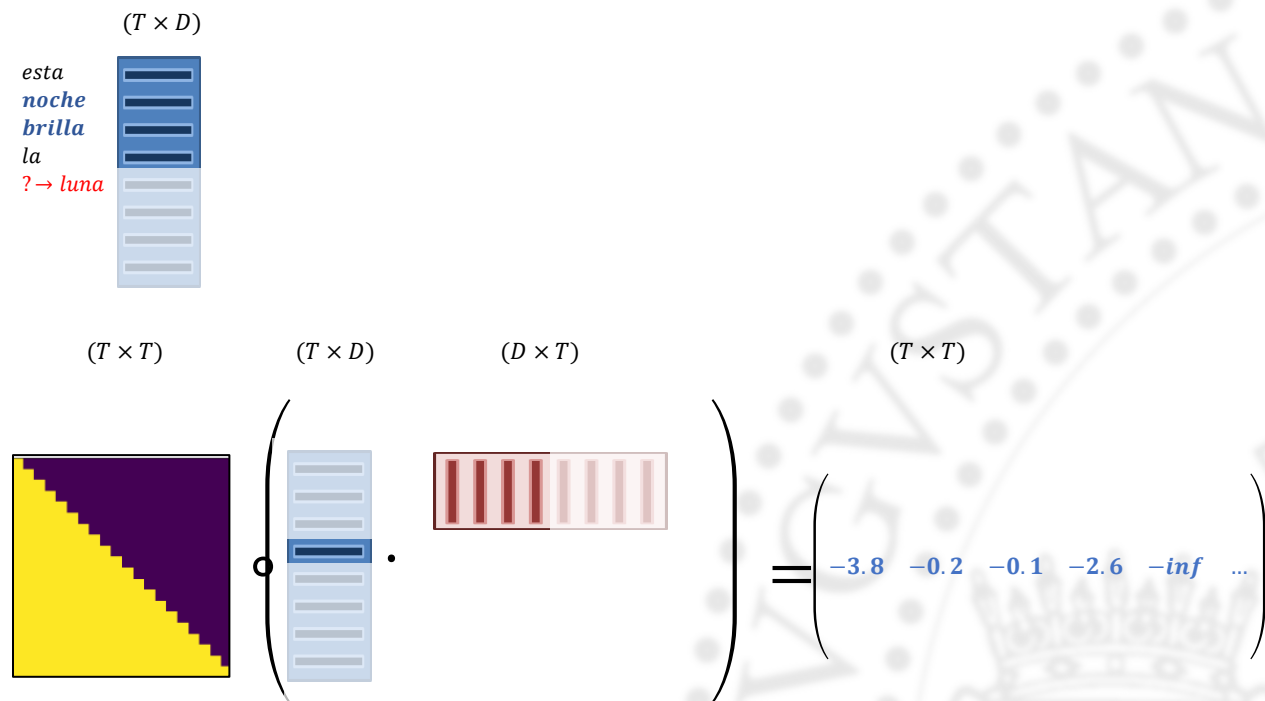
Multi-head Attention (MHA)

- Self attention: example



Multi-head Attention (MHA)

- Self attention: example



Multi-head Attention (MHA)

- Self attention: example

$$\begin{array}{c}
 \begin{array}{c} (N \times D) \qquad (N \times D) \\
 \text{softmax} \left(\begin{array}{cccccc} -2.8 & -0.2 & -0.1 & -2.6 & -inf & \dots \end{array} \right) \cdot \begin{array}{|c|} \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \end{array} = \\
 \\
 \begin{array}{c} (N \times N) \qquad (N \times D) \qquad (N \times D) \\
 \left(\begin{array}{cccccc} 0.03 & 0.45 & 0.48 & 0.04 & 0 & \dots \end{array} \right) \cdot \begin{array}{|c|} \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \end{array} \xrightarrow[\text{=}]{} \begin{array}{|c|} \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \hline \end{array} \\
 \end{array}
 \end{array}$$

Diagram illustrating the Self Attention mechanism (MHA) for a sequence of length N and dimension D .

The input sequence is processed through a softmax operation to generate attention weights (represented by the first matrix, $(N \times D)$).

The attention weights are then multiplied by the input sequence (represented by the second matrix, $(N \times D)$) to produce the output sequence (represented by the third matrix, $(N \times D)$).

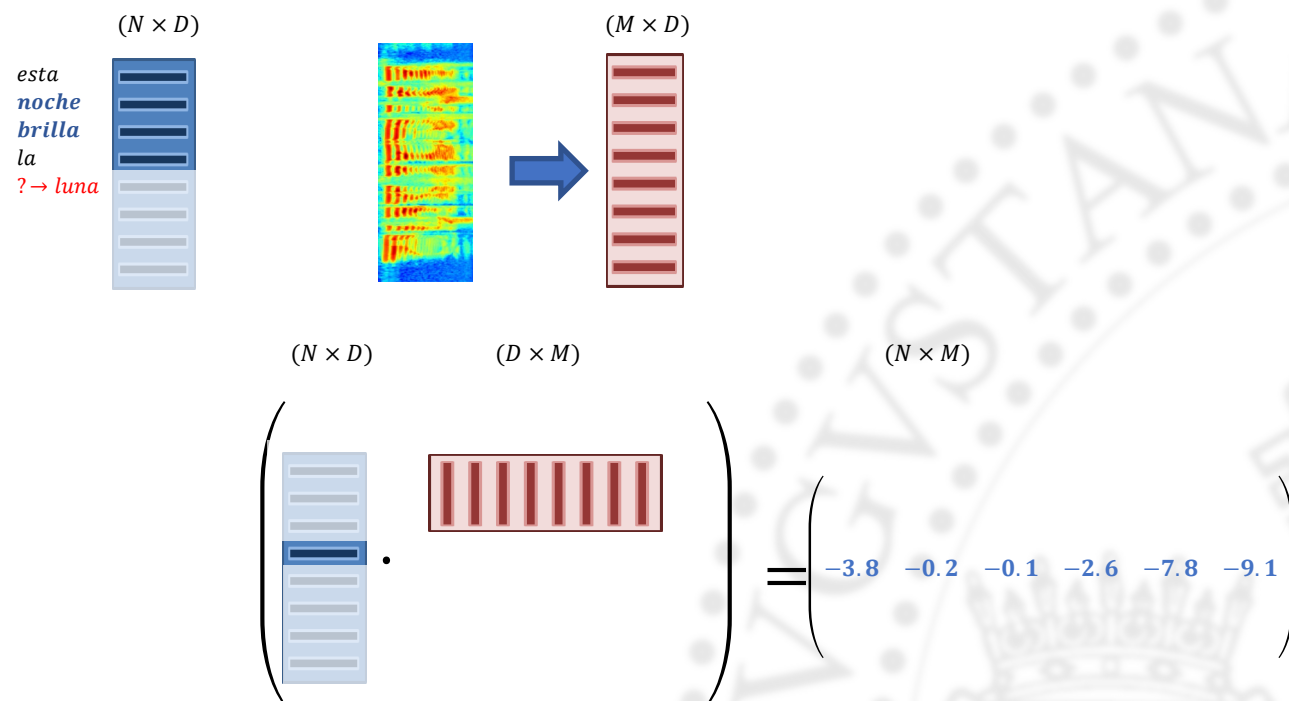
The diagram shows the calculation of the attention weights for the first row of the input sequence, where the weights are $-2.8, -0.2, -0.1, -2.6, -inf, \dots$.

The resulting attention weights are then used to calculate the output sequence, where the weights are $0.03, 0.45, 0.48, 0.04, 0, \dots$.

The output sequence is calculated as the weighted sum of the input sequence, where the weights are the attention weights.

Multi-head Attention (MHA)

- Cross attention: example



Multi-head Attention (MHA)

- Cross attention: example

