

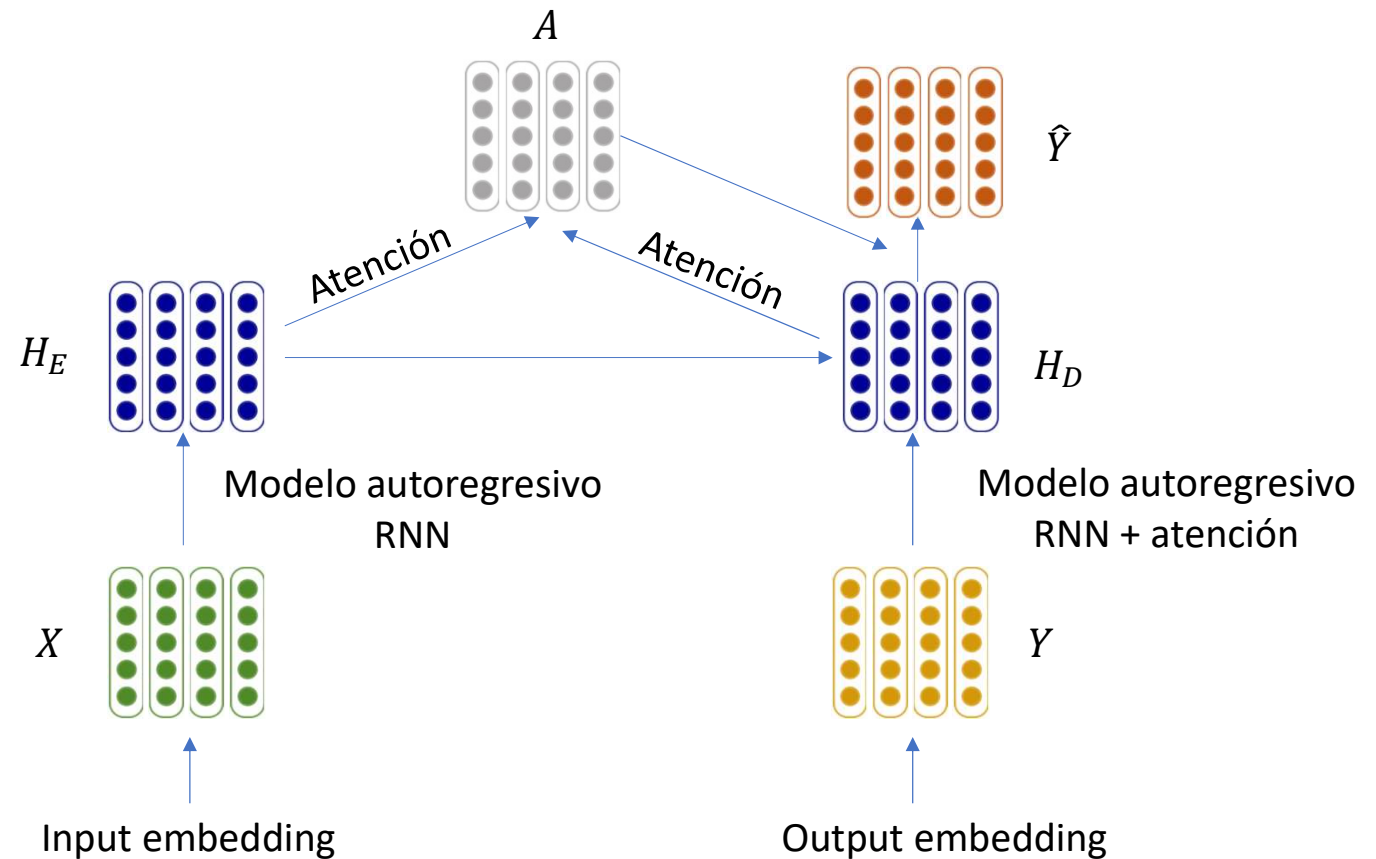
Deep learning

Transformers



Neural Machine Translation + Attention

Input: El bebé está llorando
Output: The baby is crying

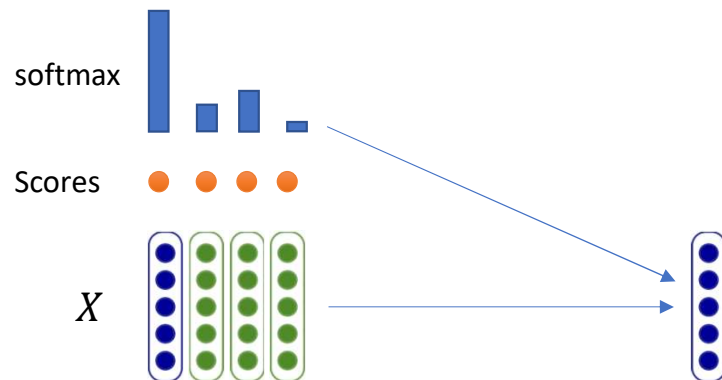


Dados dos conjuntos de vectores (que representan algo, por ejemplo una secuencia), computar nuevos vectores combinados que traten de encontrar asociaciones entre los datos originales.

Self-Attention

Se podrá aplicar atención sobre el mismo conjunto de vectores?

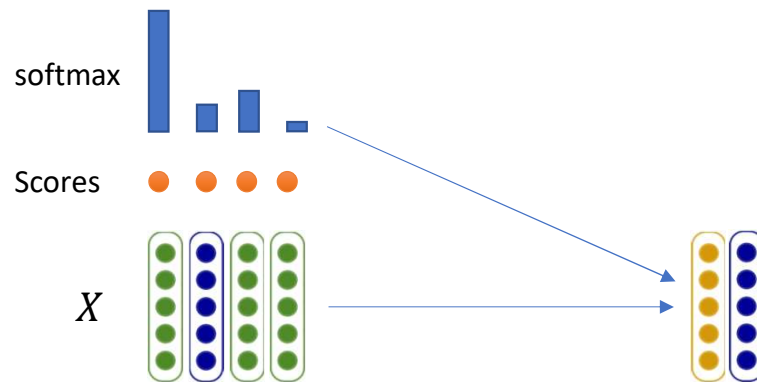
Cada vector de la colección aplica atención sobre los vectores de la colección de entrada y luego combina las atenciones para construir nuevos vectores.



Self-Attention

Se podrá aplicar atención sobre el mismo conjunto de vectores?

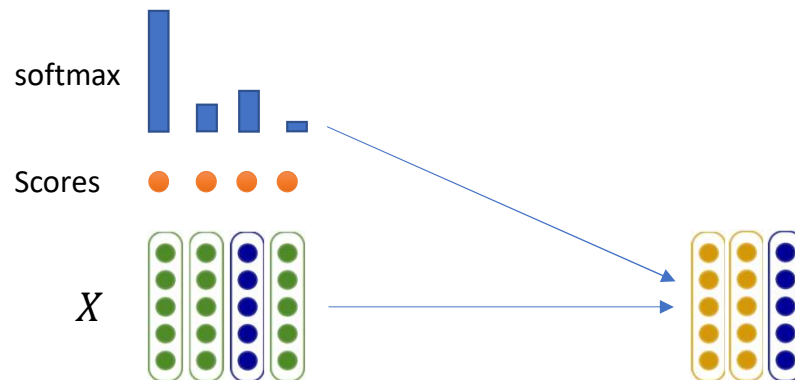
Cada vector de la colección aplica atención sobre los vectores de la colección de entrada y luego combina las atenciones para construir nuevos vectores.



Self-Attention

Se podrá aplicar atención sobre el mismo conjunto de vectores?

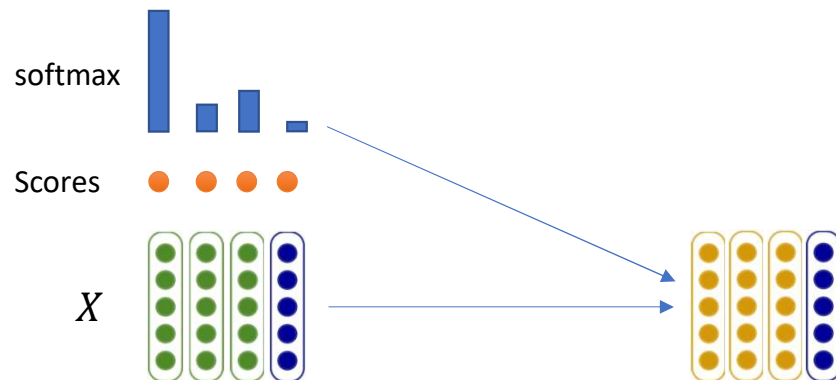
Cada vector de la colección aplica atención sobre los vectores de la colección de entrada y luego combina las atenciones para construir nuevos vectores.



Self-Attention

Se podrá aplicar atención sobre el mismo conjunto de vectores?

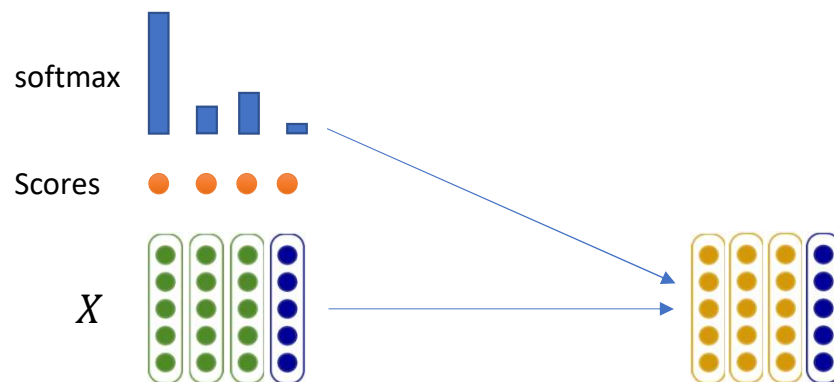
Cada vector de la colección aplica atención sobre los vectores de la colección de entrada y luego combina las atenciones para construir nuevos vectores.



Self-Attention

Se podrá aplicar atención sobre el mismo conjunto de vectores?

Cada vector de la colección aplica atención sobre los vectores de la colección de entrada y luego combina las atenciones para construir nuevos vectores.



Podemos generar nuevos datos a través de la atención

Ya no necesitamos la recurrencia!

Asumimos que los datos vienen en matrices de dimensión $seq \times dim$

Todo se puede calcular en paralelo, así:

$$H = softmax(XX^T)X$$

Pero parece no tener mucho sentido aplicar la auto-atención. Es muy posible llegar a la transformación de identidad entre X y H

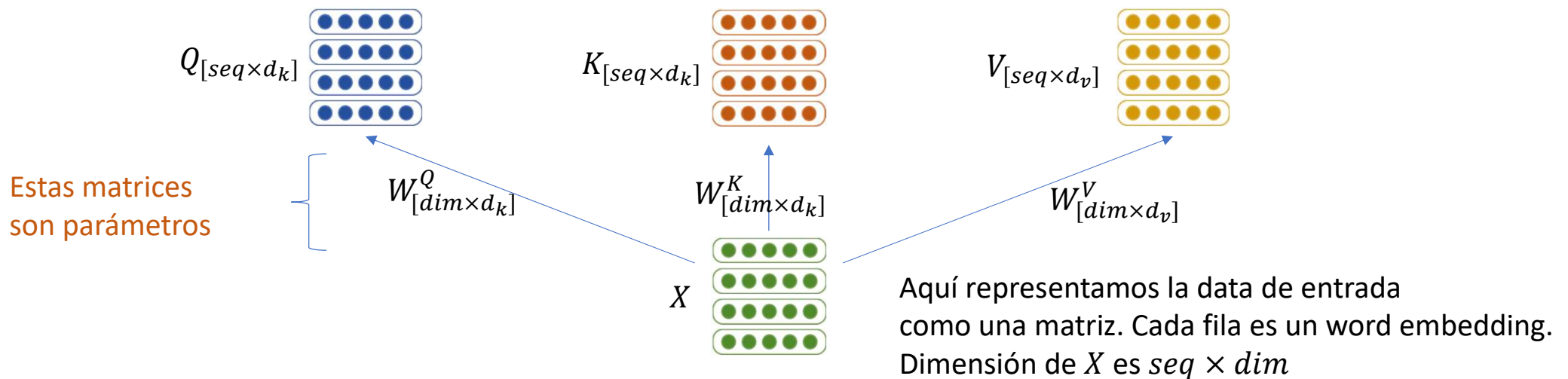
Self-Attention

Solución: Agregar un nivel de abstracción a los vectores que se usan en la auto-atención

$$H = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

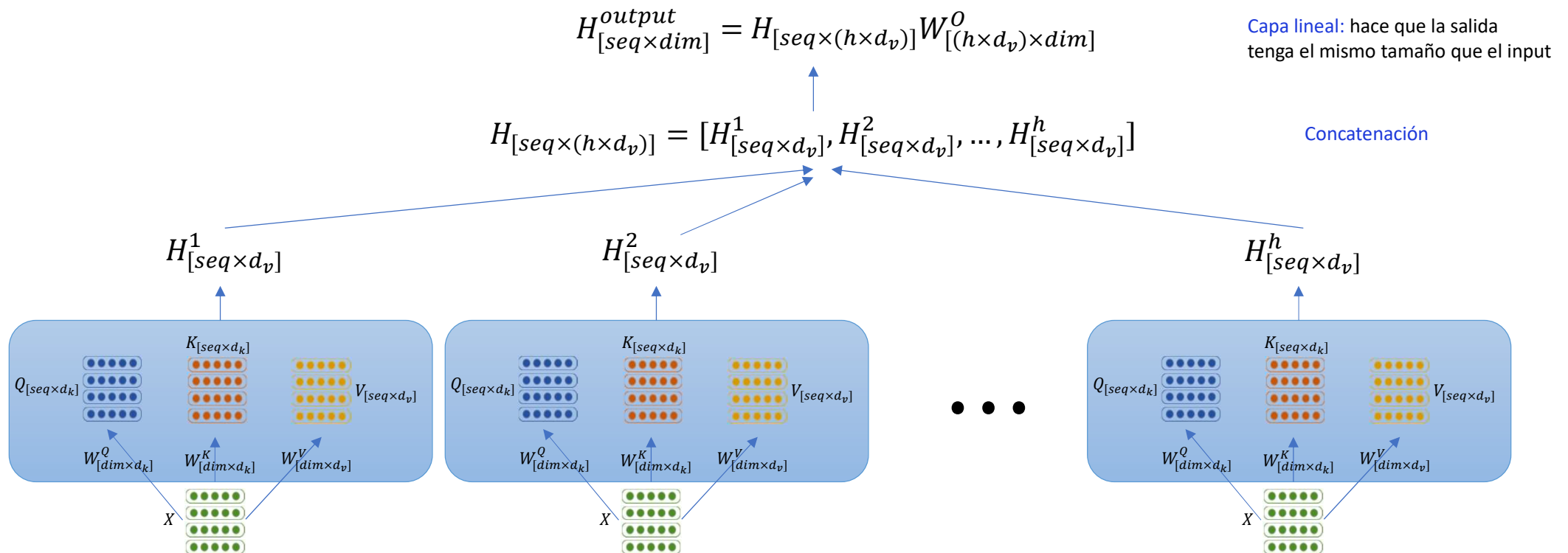
QK^T es una matriz cuadrada de scores
“similitud entre cada par de embeddings de la secuencia”

El resultado es $H_{[dim \times d_v]}$

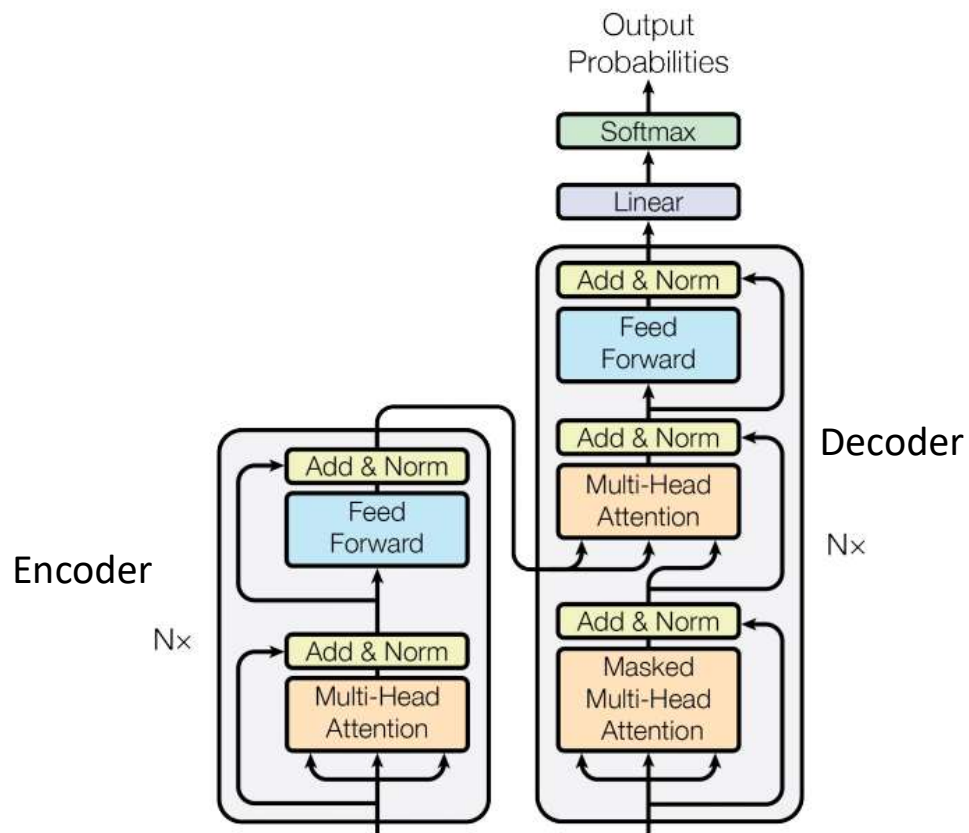


Multi-head Self-attention

Auto-atención es una transformación de los datos de entrada. Para potenciar su uso, se usa varias veces, esperando que cada aplicación extraiga diferente información.



Transformer



Encoder

- Shortcuts y normalización después de cada atención multi-cabeza
- MLP de dos capas (con ReLU)
- Este bloque se repite N veces

Decoder

- Self-attention sobre la secuencia target
- Atención sobre la salida del encoder
- Shortcuts y normalización
- MLP al finalizar el bloque

Un MLP y softmax al final

Transformer

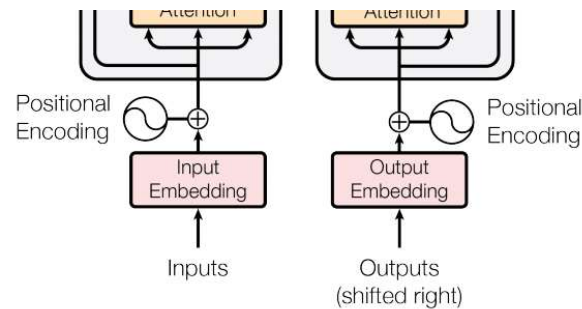
Qué sucede con el orden de la secuencia en esta arquitectura?
Hemos perdido información del orden temporal

Solución: positional encoding

- Agregar información al embedding inicial para saber en qué parte de la secuencia sucede

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{model}})$$

- La propuesta original es usar funciones periódicas en función de posición y dimensión del embedding



Transformer

Un detalle más: el decoder debería aprender a poner atención sólo a los tokens que ocurrieron hasta un determinado momento. **Al momento de hacer inferencia, no se conoce el futuro.**

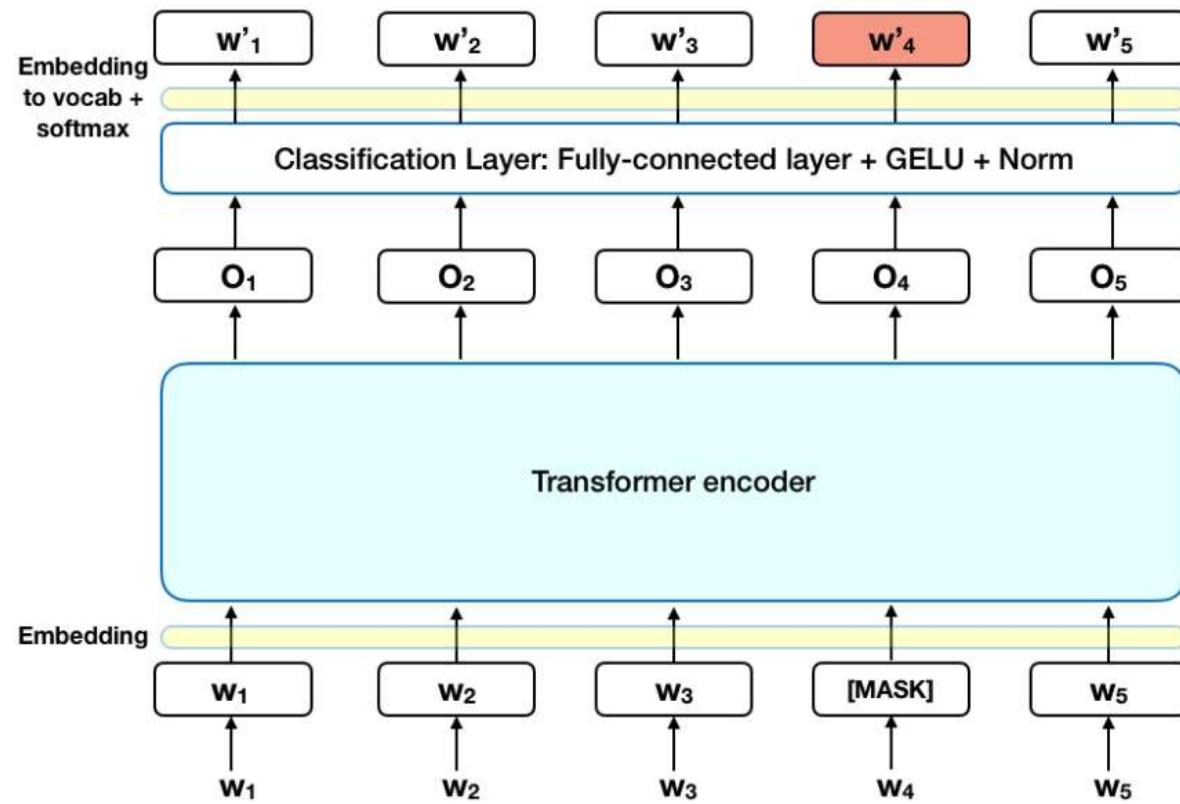
Esto se puede resolver fácil en el mismo cómputo de la auto-atención.

Usar a propósito un score de atención $-\infty$ para posiciones futuras. Softmax se hará cargo de poner en cero los scores de atención.

The diagram illustrates the process of applying a Softmax function to an attention matrix. On the left, a 4x4 matrix is shown with values: $\begin{bmatrix} 0.7 & -\infty & -\infty & -\infty \\ 0.1 & 0.6 & -\infty & -\infty \\ 0.1 & 0.3 & 0.6 & -\infty \\ 0.1 & 0.3 & 0.3 & 0.3 \end{bmatrix}$. This matrix is enclosed in a red border. To its left is the text "Softmax(" and to its right is a closing parenthesis ")", with an equals sign "=" following. To the right of the equals sign is another 4x4 matrix, enclosed in a blue border, representing the result of the Softmax operation. The resulting matrix has values: $\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0.37 & 0.62 & 0 & 0 \\ 0.26 & 0.31 & 0.43 & 0 \\ 0.21 & 0.26 & 0.26 & 0.26 \end{bmatrix}$. Above the resulting matrix, the column headers are "<start>", "I", "am", and "fine". To the left of the resulting matrix, the row headers are "<start>", "I", "am", and "fine".

	<start>	I	am	fine
<start>	1	0	0	0
I	0.37	0.62	0	0
am	0.26	0.31	0.43	0
fine	0.21	0.26	0.26	0.26

BERT



ViT

