



# Redes Transformers

Clase VI - Introducción al Aprendizaje Profundo



# ¿Qué vimos hasta ahora?

- Redes neuronales MLP.
- Redes convolucionales.
- Redes recurrentes.

# ¿Qué vamos a ver ahora?

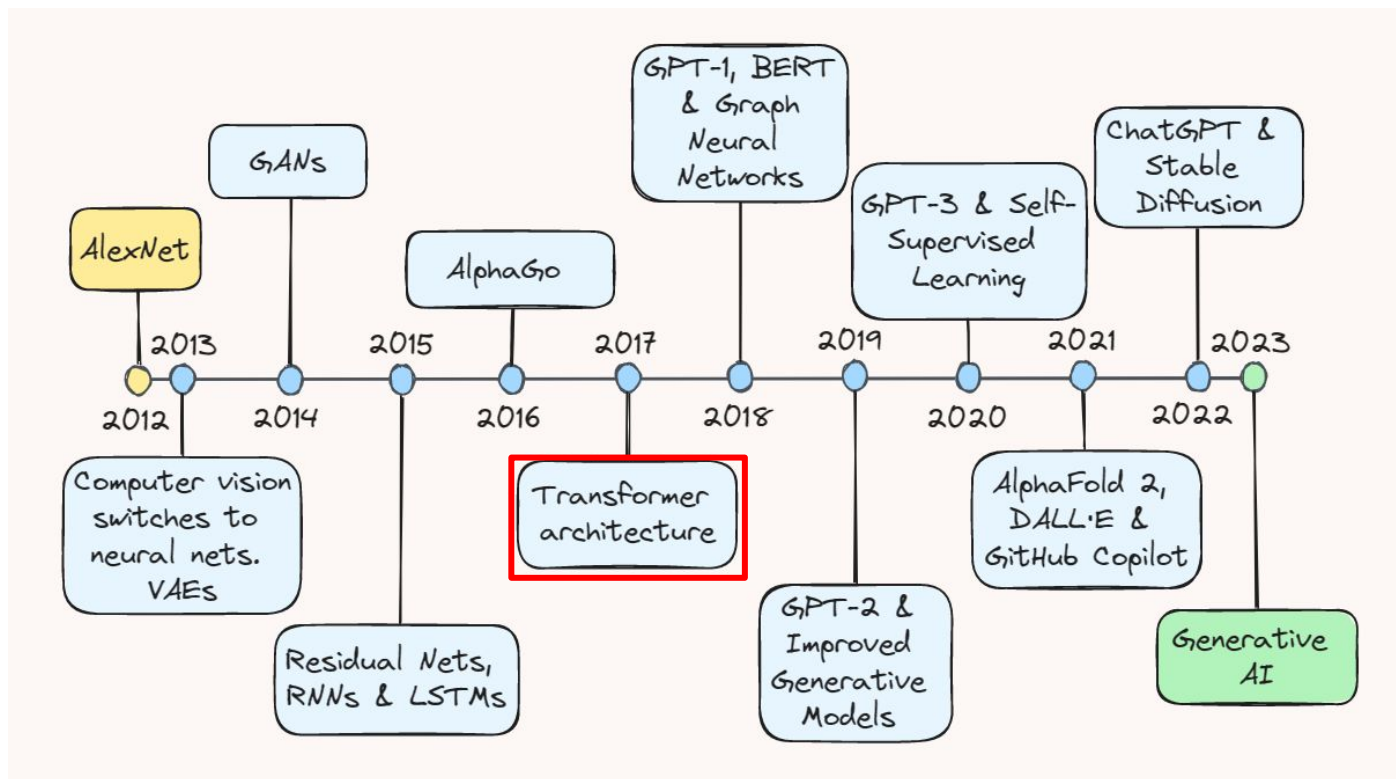
- Introducción a transformers
- Introducción a BERT y GPT



[Fuentes: del Optimus Prime: fanpop.com.](https://www.fanpop.com/optimus-prime)

[Del ícono: flaticon](#)

# Usted está aquí

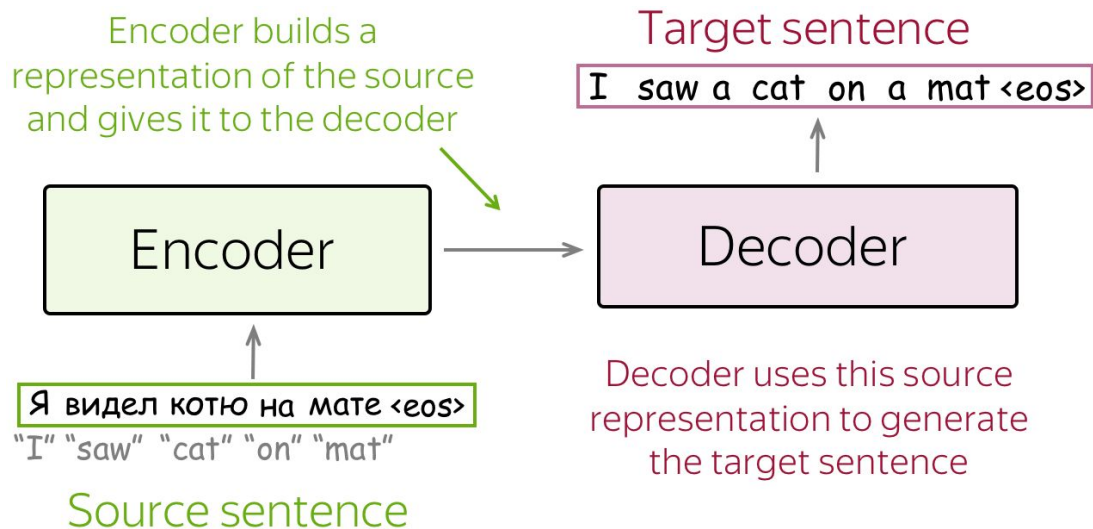


# Contexto

- El entrenamiento de las RNN es secuencial
- Esto hace que el entrenamiento no sea paralelizable como en una CNN
- Además, todo el contexto se resumía a un sólo vector
- Era un problema considerable en traducción de textos: por velocidad de entrenamiento y para mantener secuencias largas

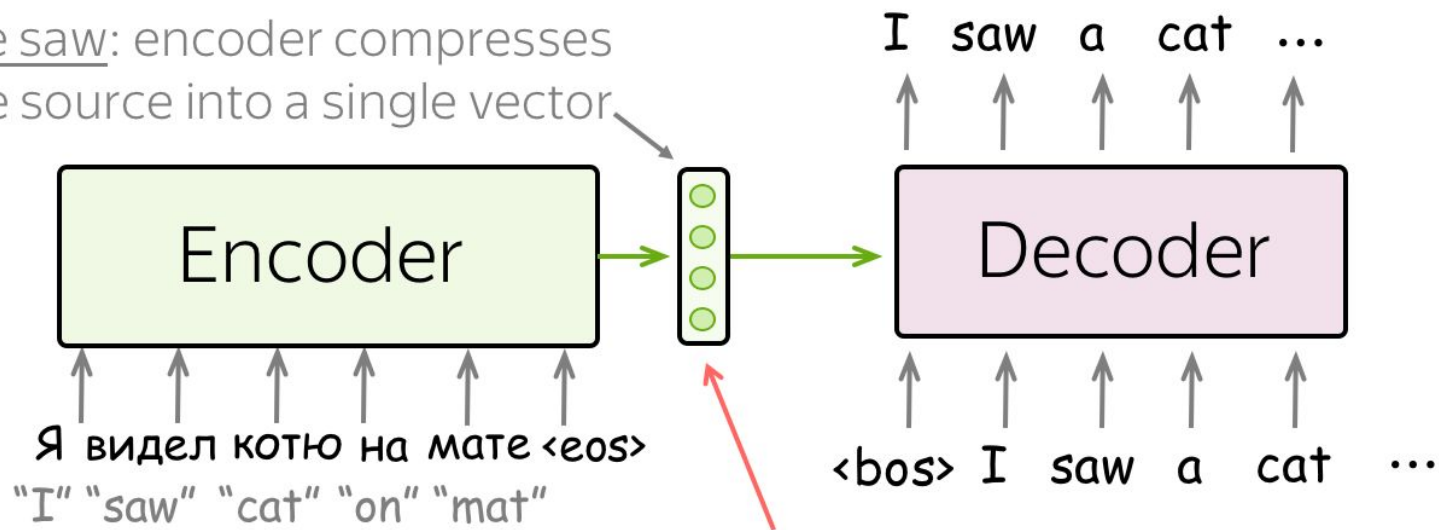
# Contexto

- Transformers surgen en 2017, originalmente para problemas de traducción
- En los mismos, se usaba una arquitectura de tipo encoder–decoder con LSTM



# Problema

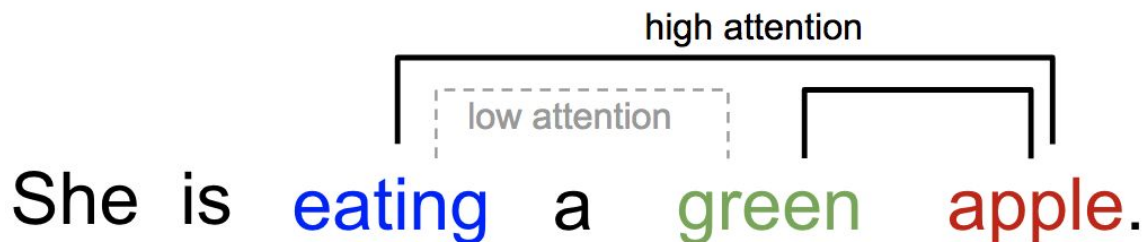
We saw: encoder compresses the source into a single vector



Problem: this is a bottleneck!

# Mecanismo de Atención

- Método (del 2014) por el cual el modelo, en distintos momentos, se enfoca en distintas partes de la entrada.
- Para esto se pesan los tokens para que se “preste” atención a los más relevantes en un cierto momento.
- Este pesado lo aprende el modelo en el entrenamiento.





# Transformers

- Surgen en 2017, con el paper seminal: [Attention is all you need.](#)
- Proponen una nueva arquitectura para resolver el anterior cuello de botella. Ahora pueden capturar contextos más largos y en paralelo.
- Se deja de usar recurrencia, y en su lugar va el mecanismo de atención, junto con varios ingredientes más.

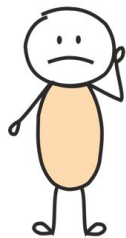
# Transformers

- En lugar de procesar token por token de una oración, procesa todos a la vez.
- Añade también otros ingredientes para poder usar tales mecanismos (referirse al material extra de la clase para un deep dive en los mismos).
- Corona el “momento computer vision” de NLP

# Transformers

I arrived at the **bank** after crossing the ... ..street? ...river?

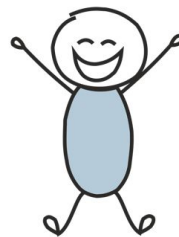
What does **bank** mean in this sentence?



I've no idea: let's wait until I read the end

RNNs

$O(N)$  steps to process a sentence with length  $N$



I don't need to wait - I see all words at once!

Transformer

Constant number of steps to process any sentence

# Desventajas vs RNN

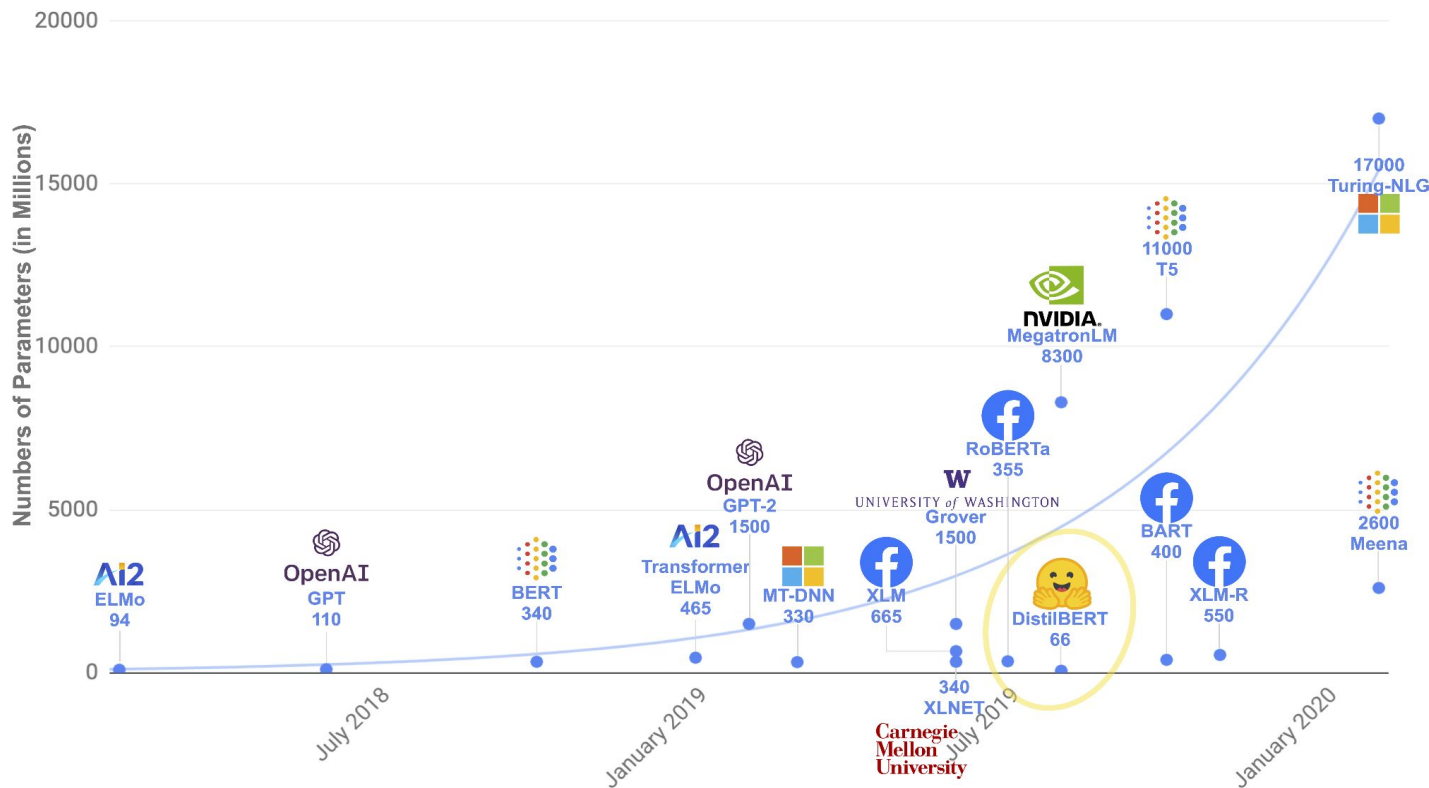
Los transformers requieren:

- Muchos (**muchos!**) más datos de entrenamiento
- Mucho más poder computacional

Además, la inferencia es más costosa y demora más

# Donde transformers se hacen LLMs...

Los transformers  
ya suponen  
modelos  
**MUY** grandes...

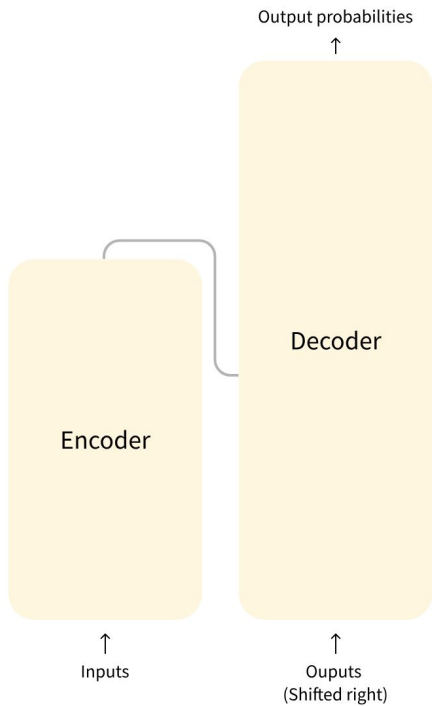


Fuente: [Huggingface](https://huggingface.co)

# Vista general

## Tres “sabores” básicos de transformers

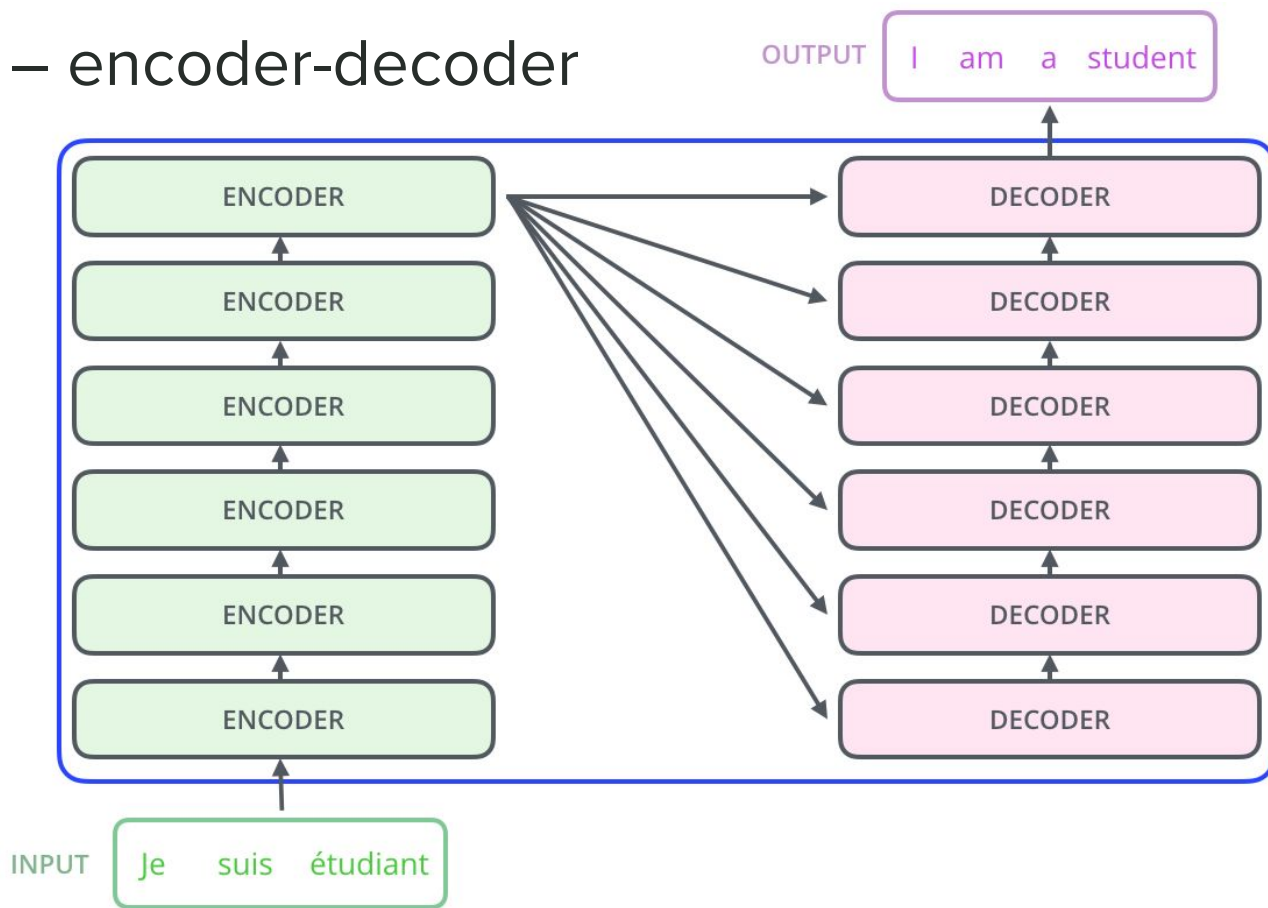
- Encoder-decoder  
(el de la imagen, el original)
- Encoder solamente
- Decoder solamente



## Vista general – encoder-decoder

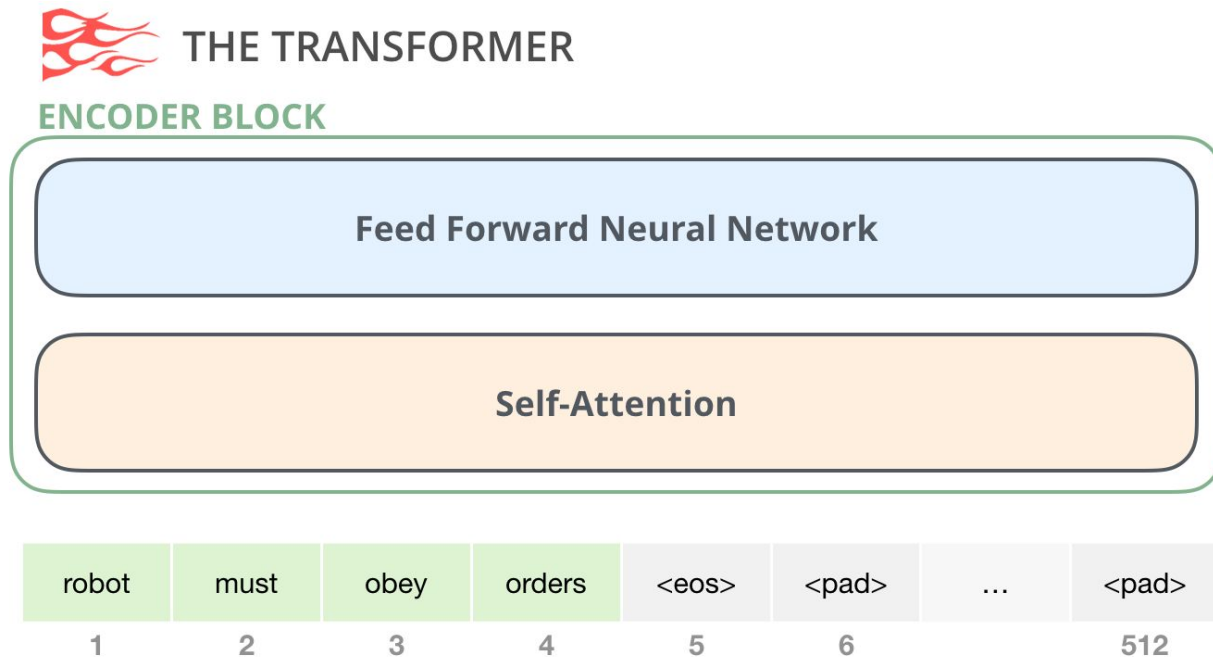
- Entrada: el texto se tokeniza, cada token se transforma en su embedding, y se le suma un encoding posicional para preservar info de la posición en la oración.
- Probabilidades de salida: se aplica softmax a los logits de salida, con lo que el modelo puede tener las probabilidades de los mismos
- Inició la arquitectura, suele usarse para tareas que involucren transformaciones como la traducción de texto.

# Vista general – encoder-decoder





# Vista general – encoder

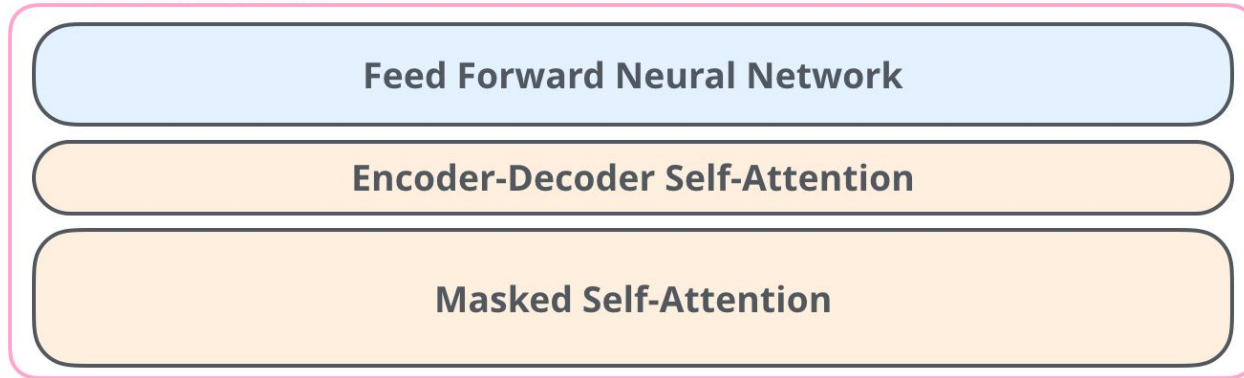


# Vista general – decoder



## THE TRANSFORMER

### DECODER BLOCK



Input



[Fuente: Jay Alammar](#)

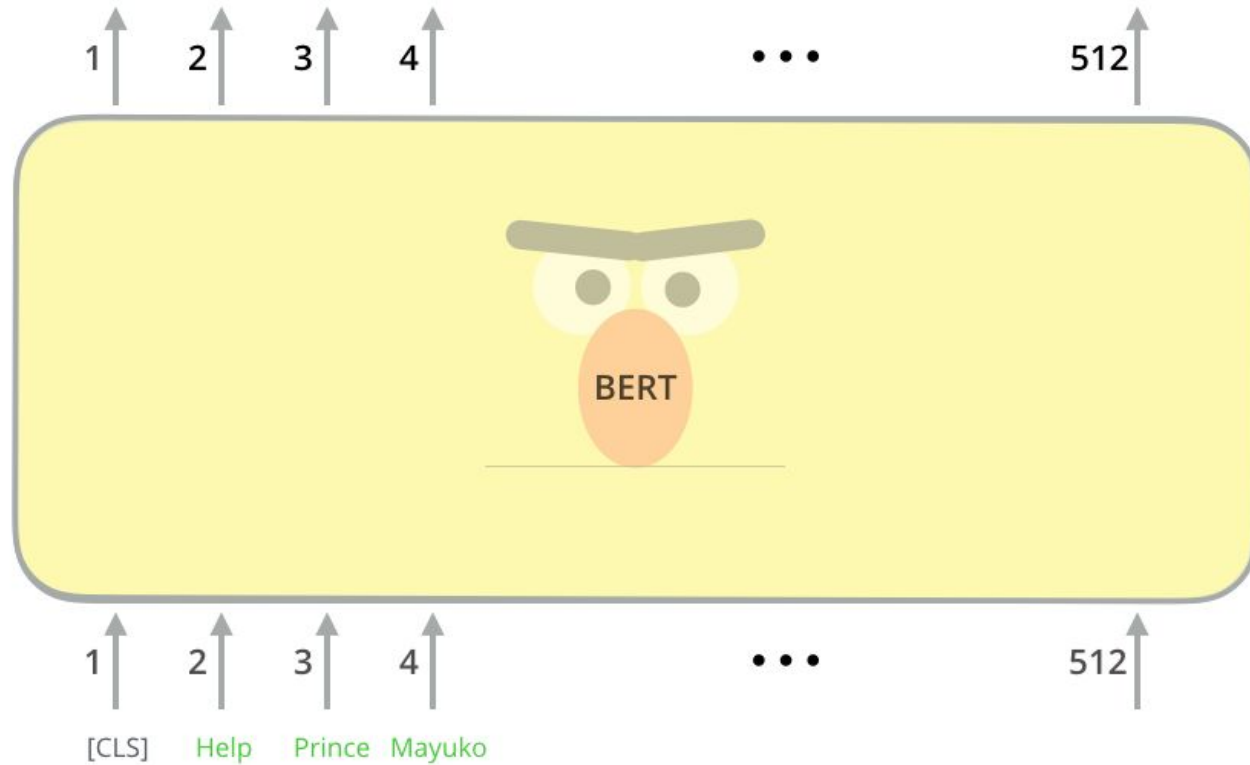
# Flavors: BERT y GPT

---

# Encoders: BERT

- Encoder-only transformer (*autoencoder*).
- Muy útiles para tareas de clasificación o generación de embeddings.
- No son modelos buenos para generar, no es su foco.
- Al no tener esta restricción, pueden ver todo el contexto de palabras.

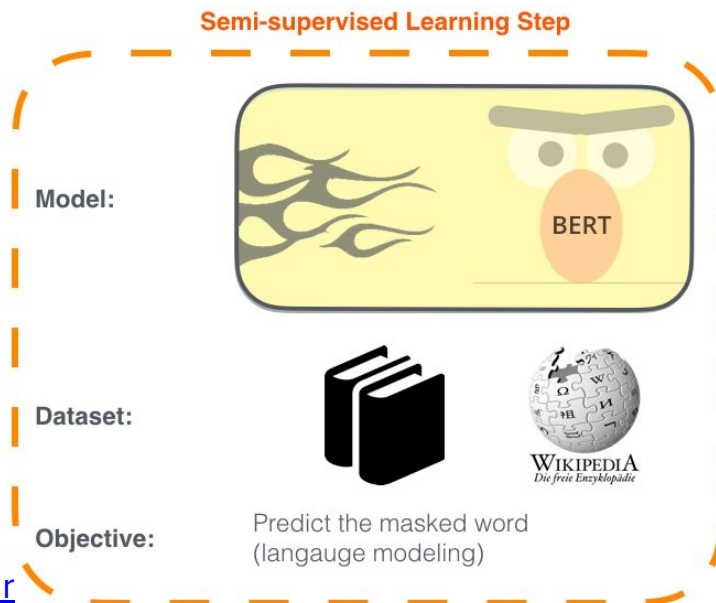
BERT



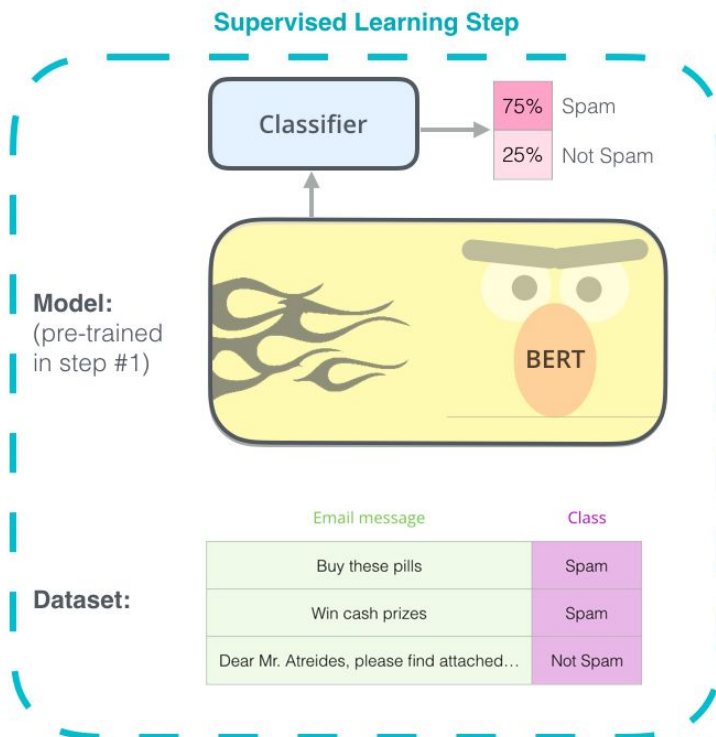
# BERT – Cómo fue entrenado – y cómo usarlo

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.



# BERT

## Cómo fue entrenado...

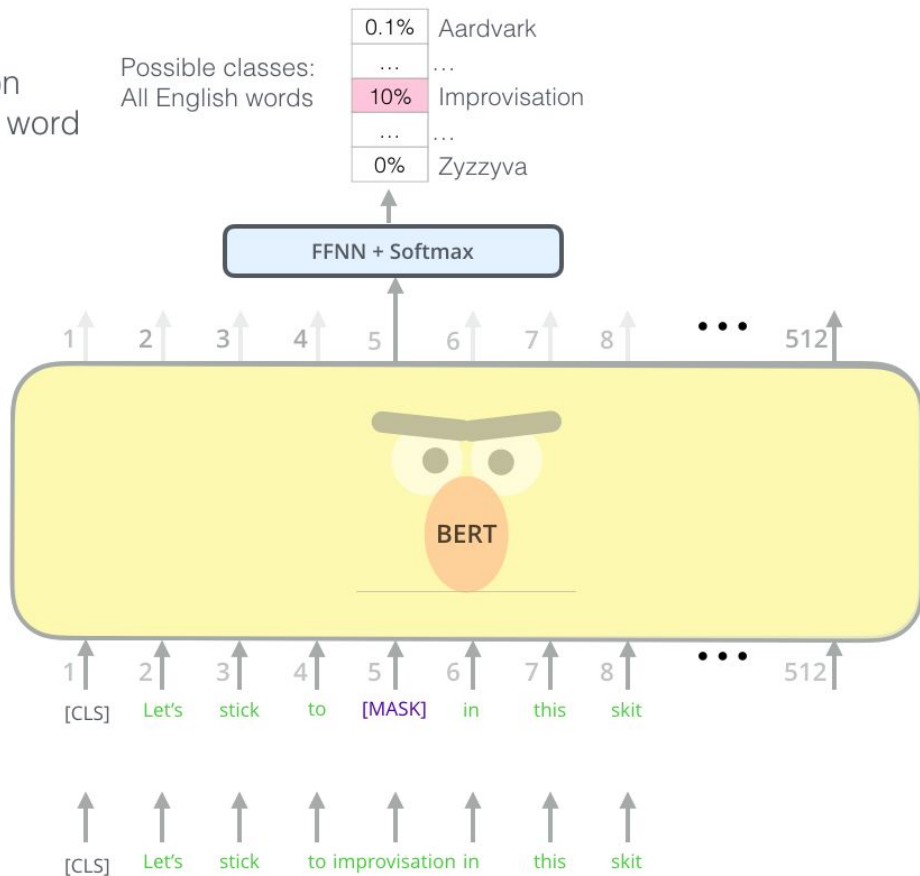
Use the output of the masked word's position to predict the masked word

Possible classes:  
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzzva

Randomly mask  
15% of tokens

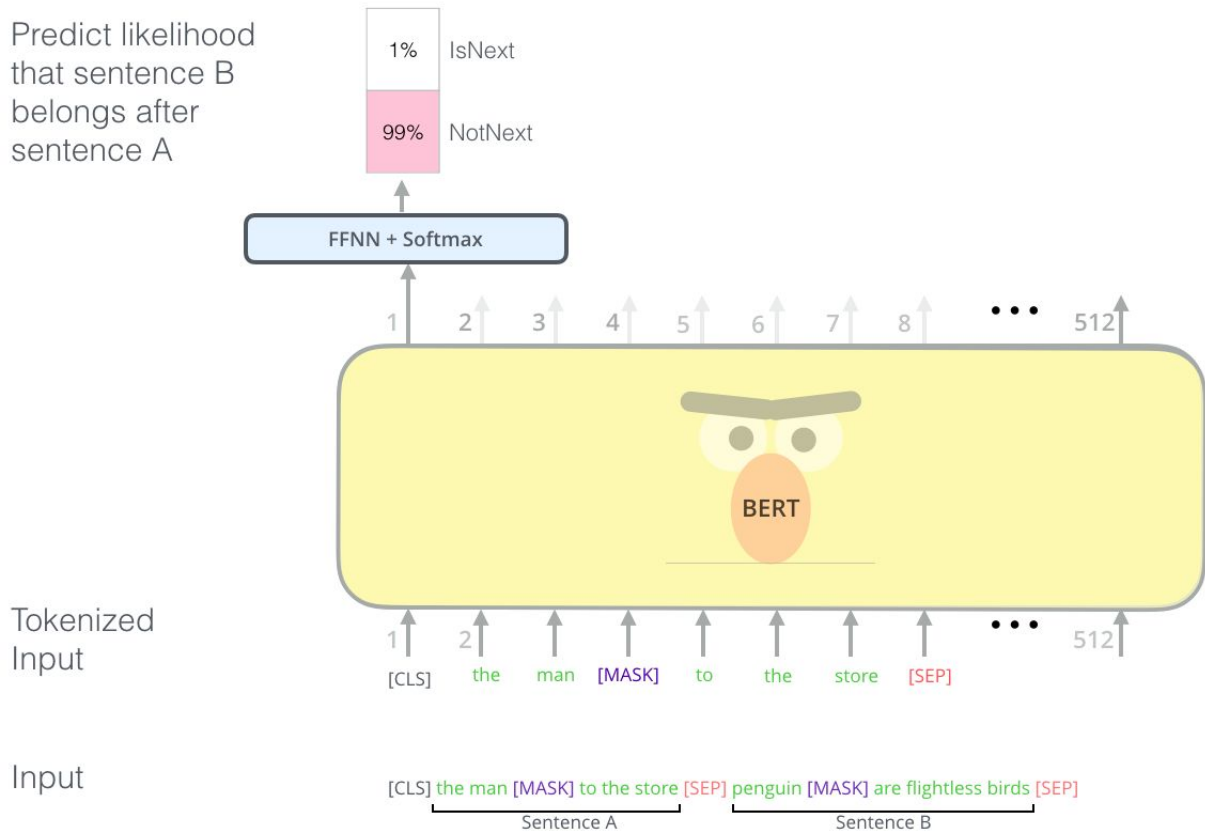
Input



# BERT

Cómo fue entrenado...

Predict likelihood  
that sentence B  
belongs after  
sentence A



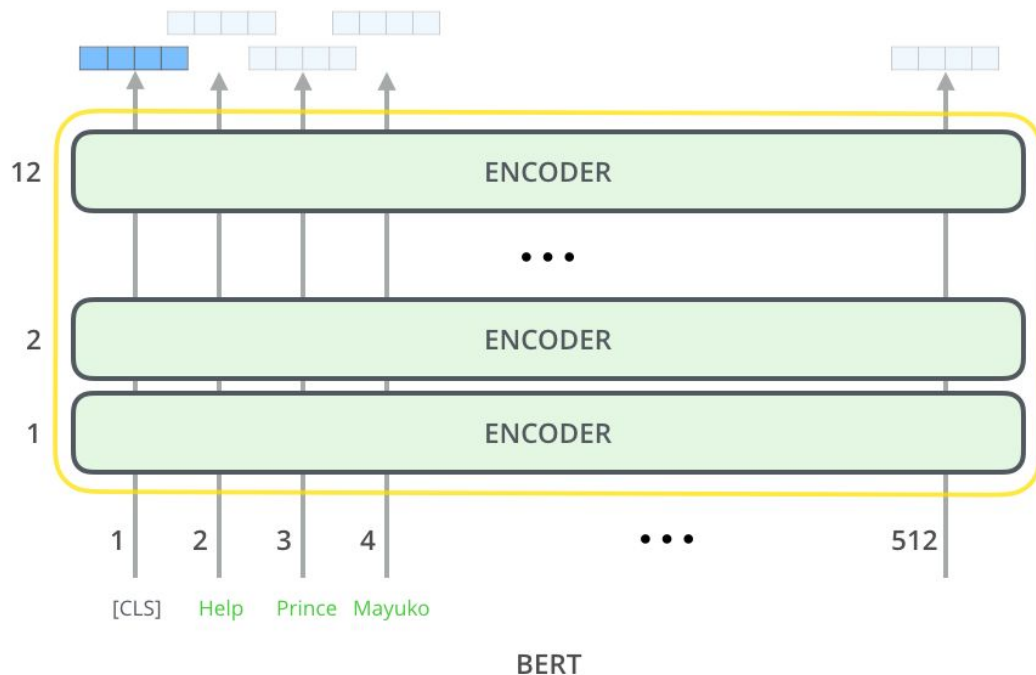
[Fuente: Jay Alammar](#)



# BERT

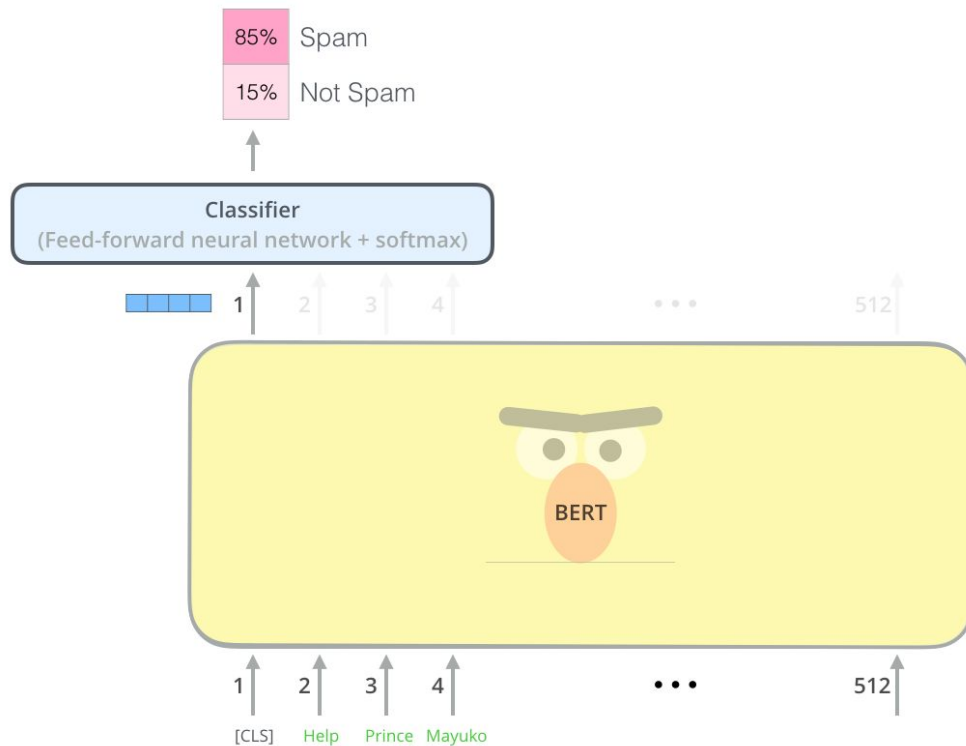
Por dentro, es un apilado de encoders que toman la concatenación de los embeddings de cada palabra.

En este caso, se puede predecir un vector por cada entrada.

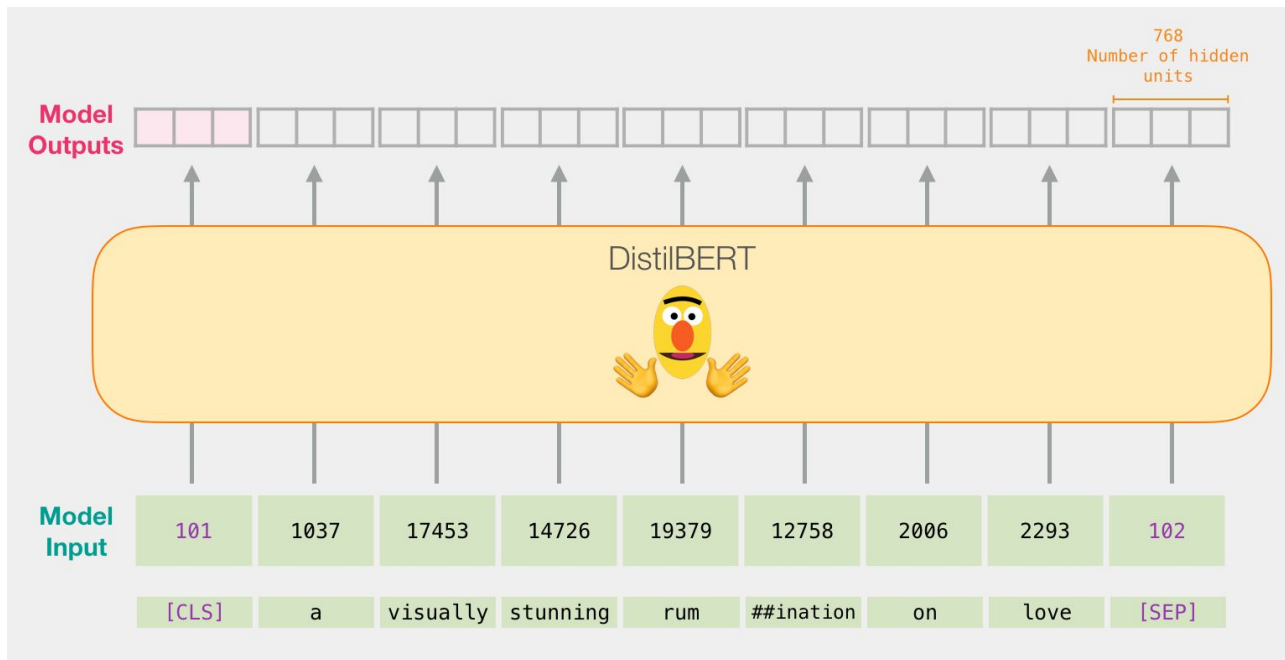


# BERT

O se puede utilizar sólo la salida del primer instante de tiempo para resolver un problema de sequence classification



## BERT – Cómo fluye una entrada



# Embeddings con BERT vs Word2Vec

## BERT

- Ofrece embeddings en el contexto de la oración.
- Ej: embedding "banco" distinto en "me sentaré en el banco" que en "depositaré en el banco".
- Toma en cuenta la posición de la palabra con los encodings posicionales.

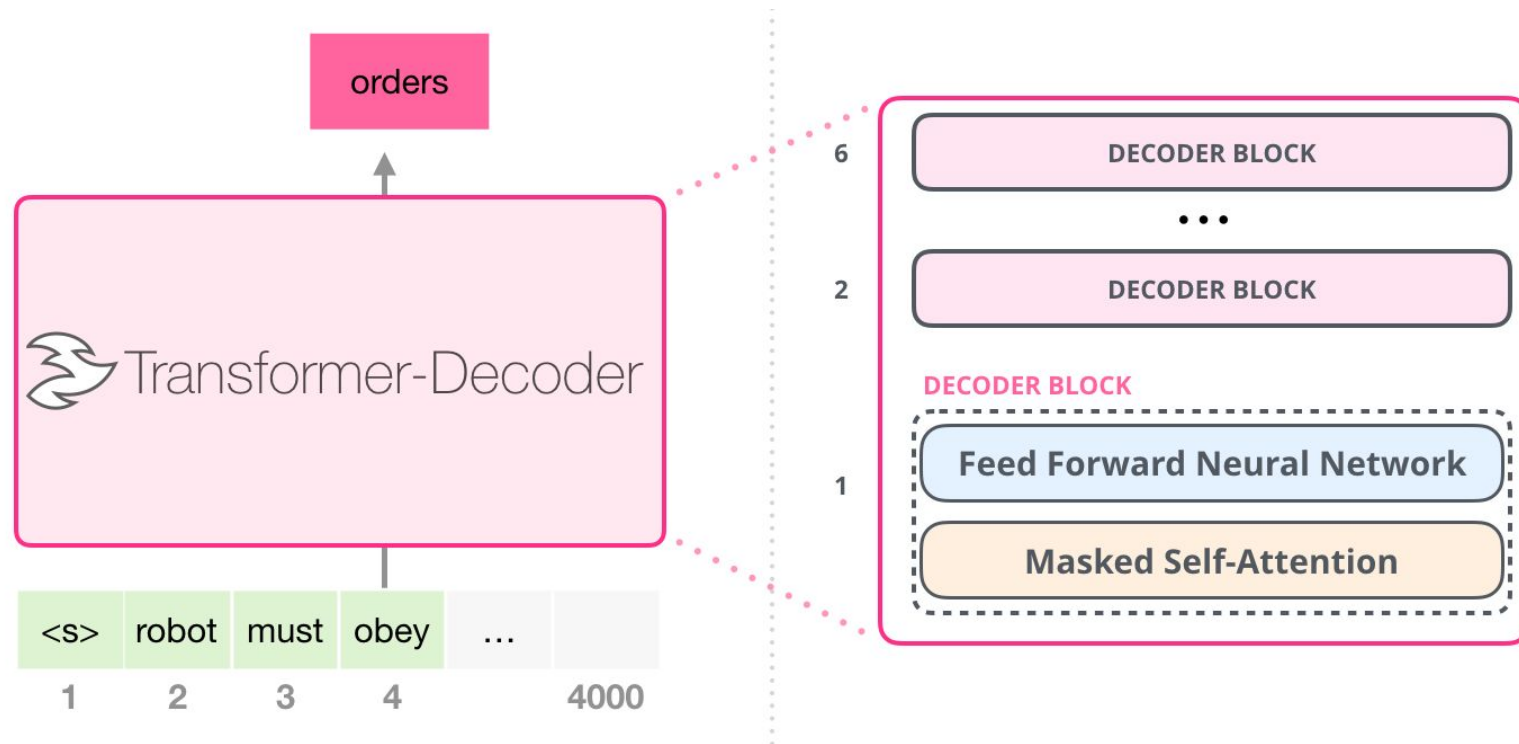
## WORD2Vec

- Mismo embedding, en el contexto general.
- Ej: mismo embedding para "banco" sin importar si es asiento o entidad bancaria.
- Es agnóstico a la posición de la palabra.

# GPT – Generative Pre-training

- Decoder-only transformer (*autoregressive*)
- Utiliza un apilado de decoders (ej. 12)
- A diferencia de encoders, tienen acceso sólo a palabras a la izquierda (o derecha, para el caso de idiomas con escritura de derecha a izquierda)
- Por esto, son modelos excelentes para generar. Pueden predecir, pero es más costoso

# GPT – Generative Pre-training



Fuente: [Jay Alammar](#)

# Hugging Face Transformers

---

# Hugging Face: modelos, datasets y librerías

Proyecto open source que contiene modelos pre-entrenados para distintas áreas de PLN, además de datasets y librerías.

La iniciativa apunta a hacer accesible modelos de estado del arte para investigadores, estudiantes y practicantes de ML, así como ayudar a reducir la huella de carbono producto de tener que re-entrenar un modelo.

Hugging Face provee implementaciones propias de dichos modelos a partir de lo expuestos en sus respectivos artículos.



## The AI community building the future.

Build, train and deploy state of the art models powered by the reference open source in natural language processing.





# Hugging Face: Transformers

Hugging face provee una librería llamada **transformers** que implementa los mecanismos de self-attention, además de habilitar el acceso a los modelos pre-entrenados desde sus servidores. Soportan **pytorch** y **tensorflow**.

```
1 from transformers import AutoTokenizer, AutoModelForMaskedLM
2 tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
3 model = AutoModelForMaskedLM.from_pretrained("bert-base-uncased")
```

# LET'S CODE

## 8. Redes Transformers



# Apéndice: arquitectura en mayor detalle y material extra

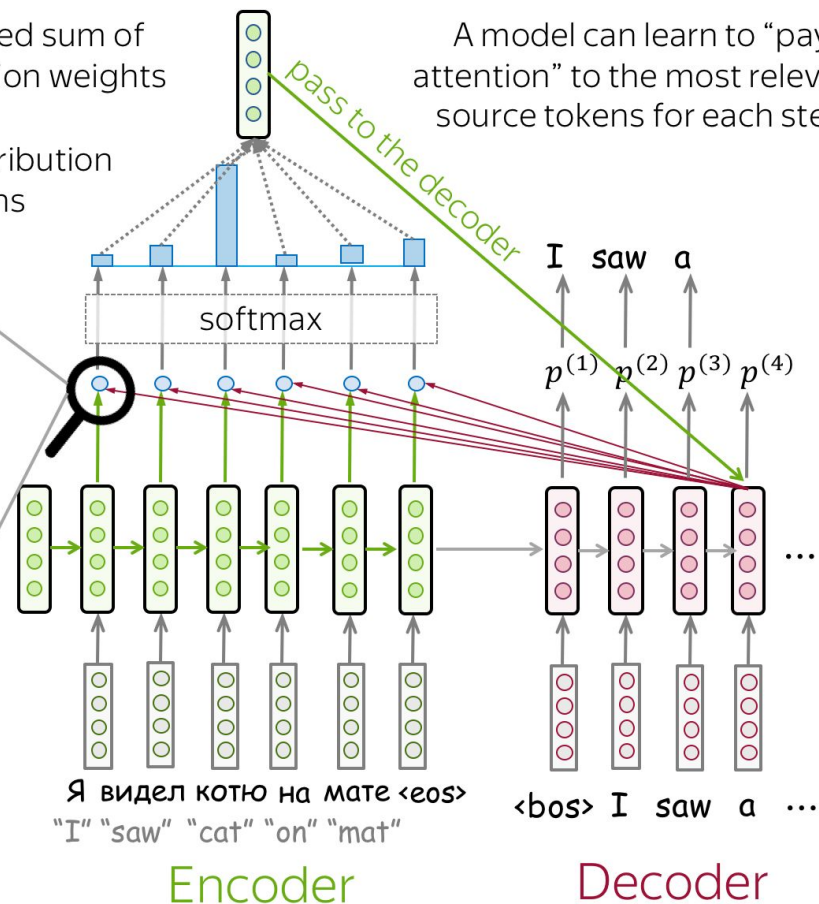
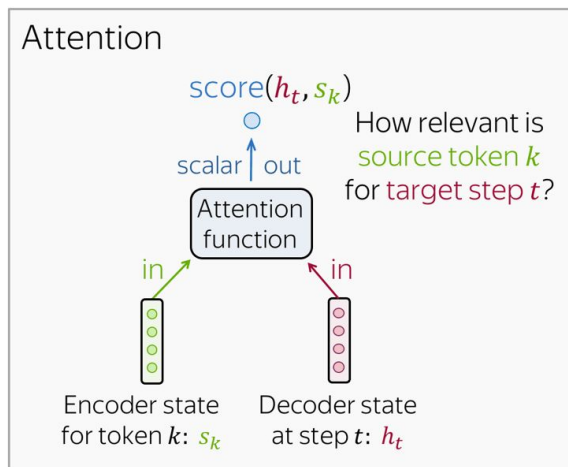
---

# Mecanismo de Atención

Attention output: weighted sum of encoder states with attention weights

Attention weights: distribution over source tokens

A model can learn to “pay attention” to the most relevant source tokens for each step



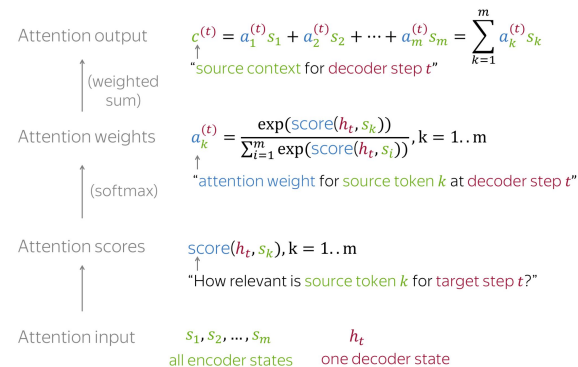
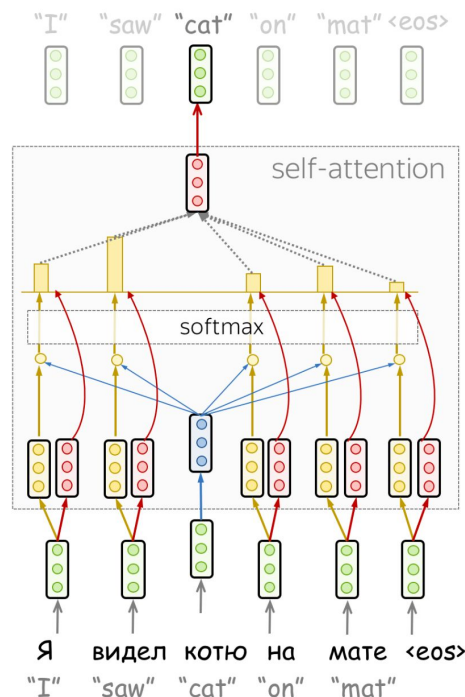
# Función de atención

Each vector receives three representations (“roles”)

$[W_Q] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{bmatrix}$  **Query:** vector **from** which the attention is looking  
 “Hey there, do you have this information?”

$[W_K] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{bmatrix}$  **Key:** vector **at** which the query looks to compute weights  
 “Hi, I have this information – give me a large weight!”

$[W_V] \times \begin{bmatrix} \text{green} \\ \text{green} \\ \text{green} \end{bmatrix} = \begin{bmatrix} \text{red} \\ \text{red} \\ \text{red} \end{bmatrix}$  **Value:** their weighted sum is attention output  
 “Here’s the information I have!”



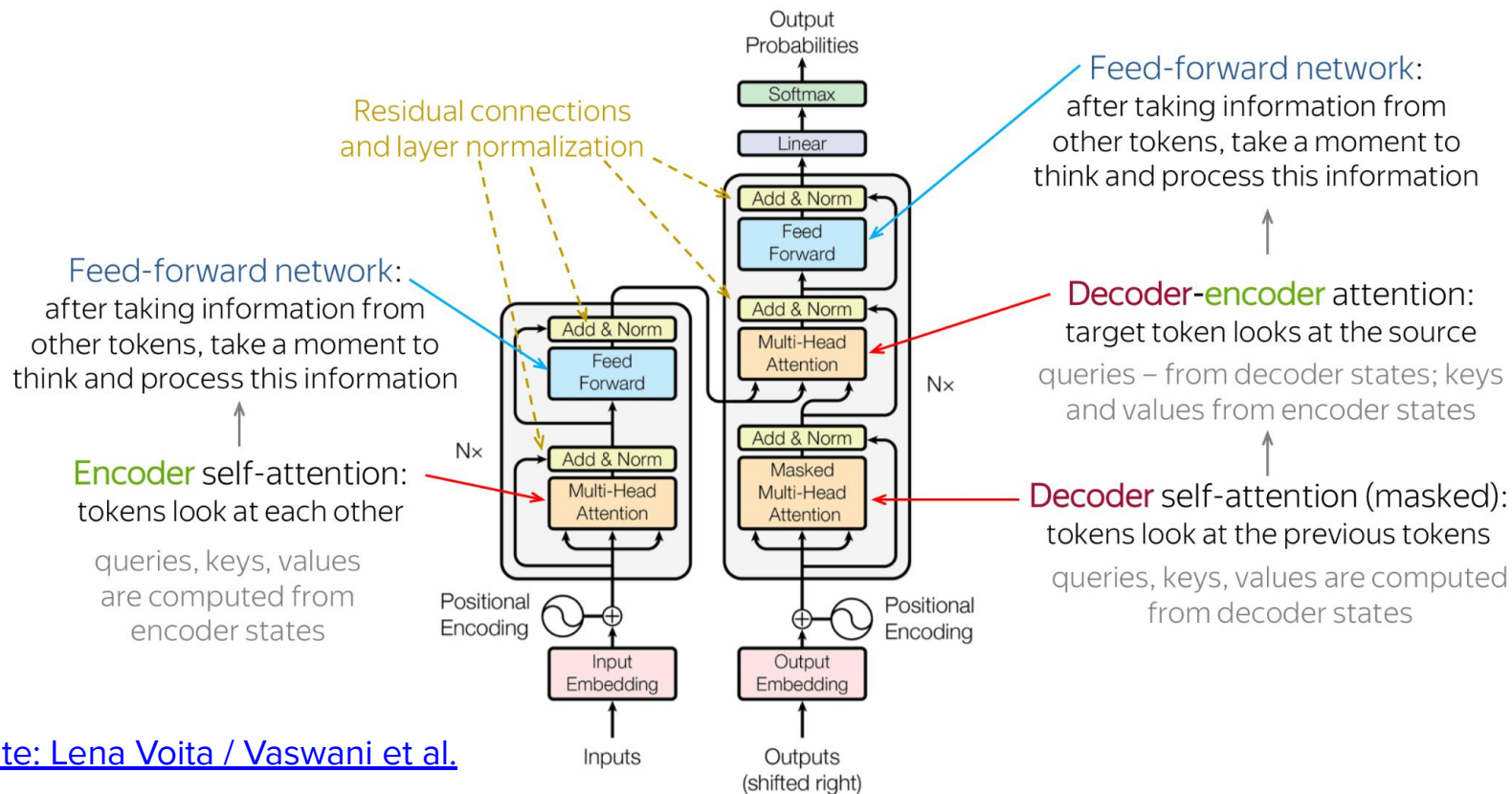
$$\text{Attention}(q, k, v) = \text{softmax}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

from to

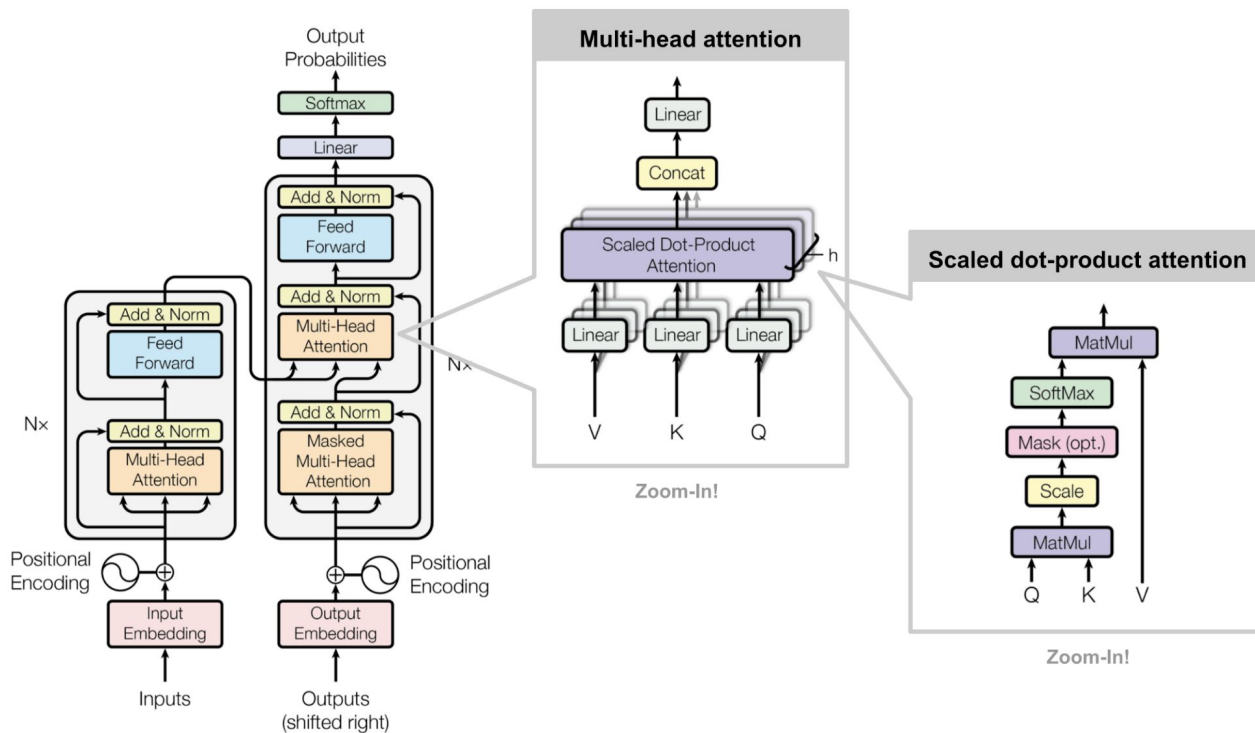
Attention weights

vector dimensionality of K, V

# Transformer – arquitectura completa



# Transformer – arquitectura completa



Fuente: [Lilian Weng / Vaswani et al.](#)

# GPT – Visualizado

Vemos desde  
un visualizador

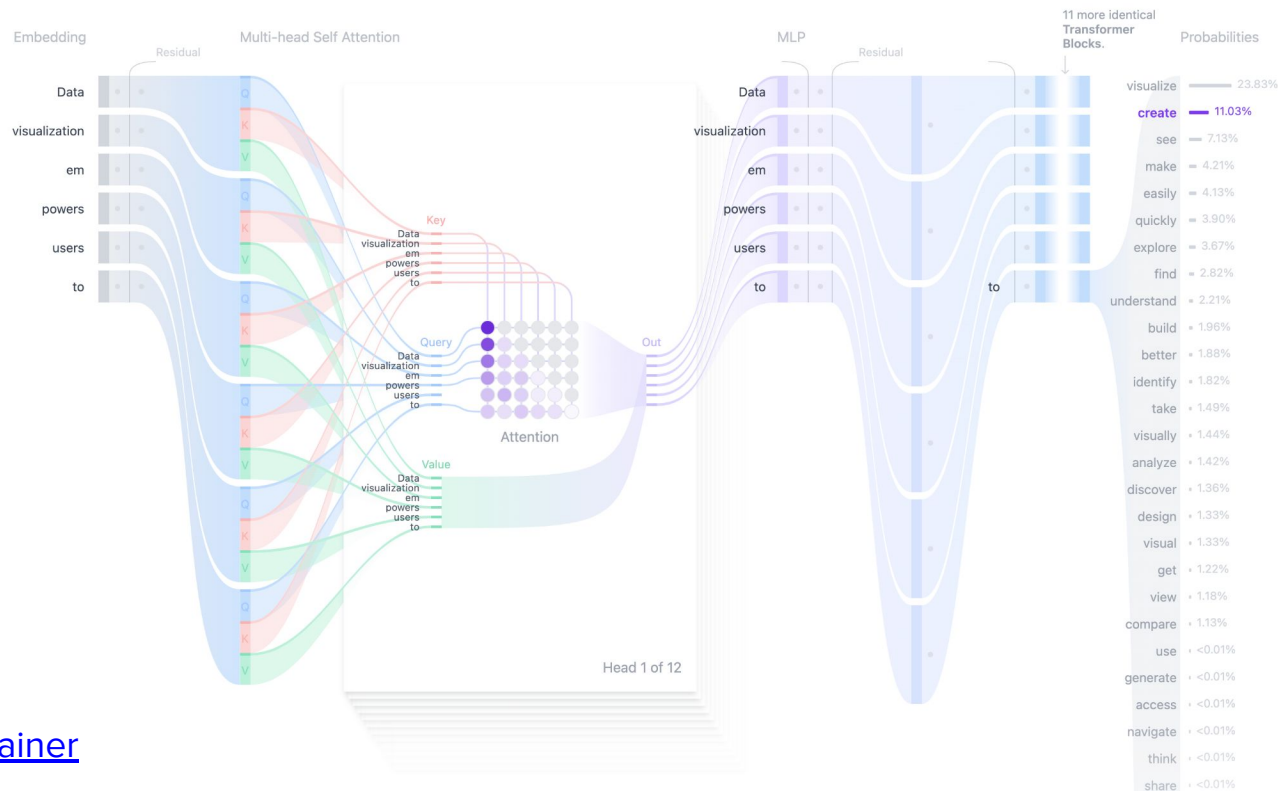
TRANSFORMER EXPLAINER

Examples ▾

Data visualization empowers users to create

Generate

Temperature 1



[Fuente: Transformer Explainer](#)



# Vision Transformers

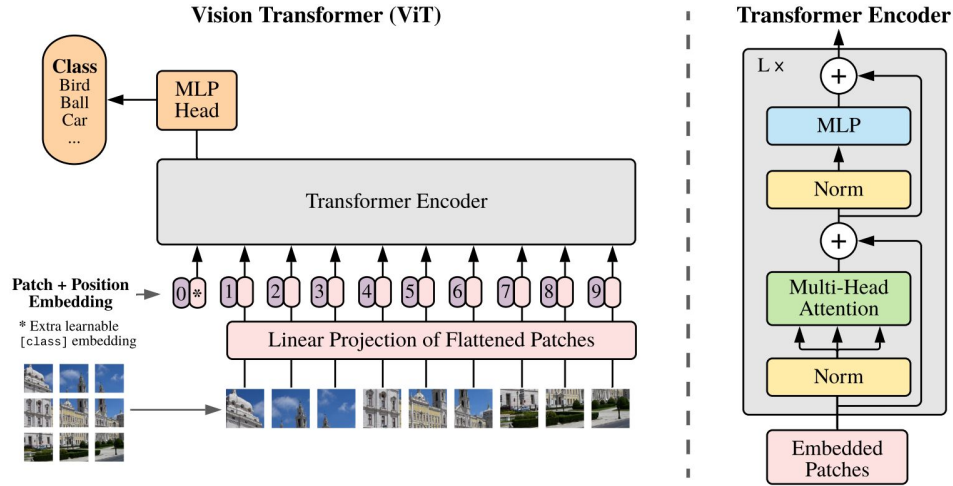


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).

## Material extra

- [Curso de NLP de Lena Voita: Transformers.](#)
- Blogs ilustrados de Jay Alammar: [Seq2seq](#), [Transformers](#), [BERT 1](#), [2](#), [GPT-2](#) y [GPT-3](#).
- [Explicador de transformers: GPT-2.](#)
- [The Vision Transformer Model.](#)