

# No todo lo que necesitas es atención (Attention is not all you need)

Trabajo Fin de Grado

Grado en Matemáticas



Diego Rivera López-Brea  
20 de septiembre de 2022

## Contenidos:

1. Introducción
2. Aprendizaje Automático
3. *Transformers* y aplicación práctica
4. Causalidad
5. Conclusiones y líneas futuras

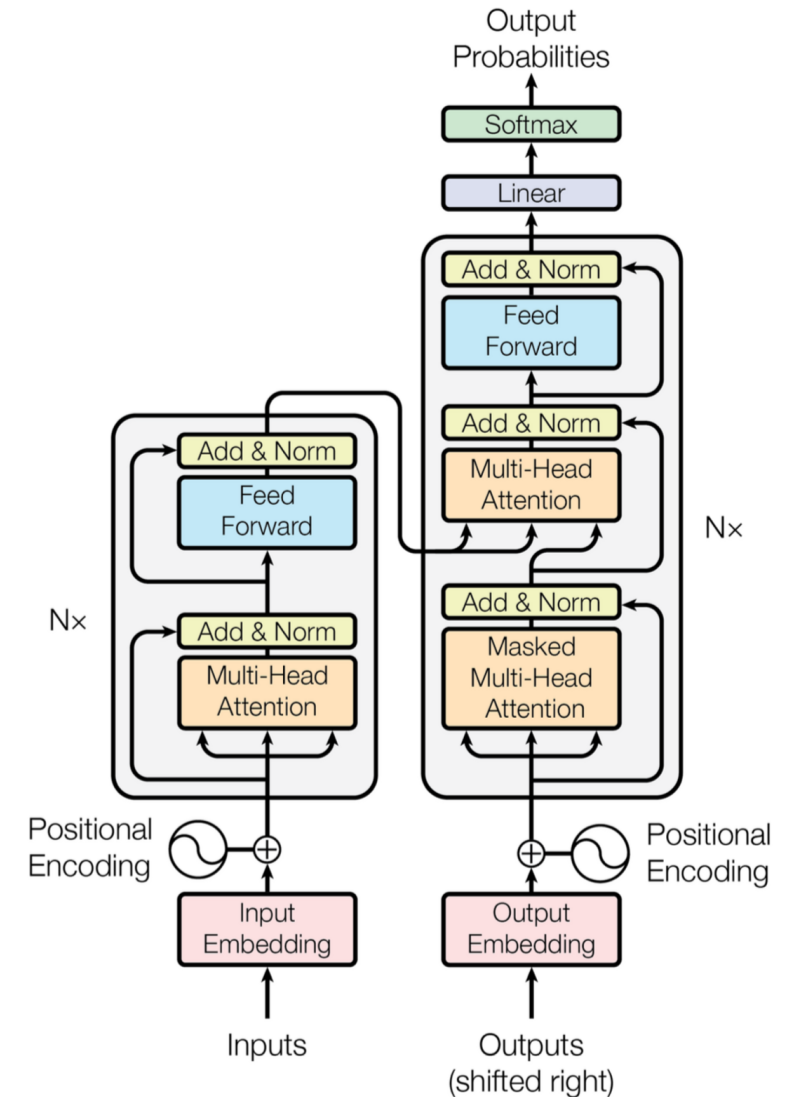
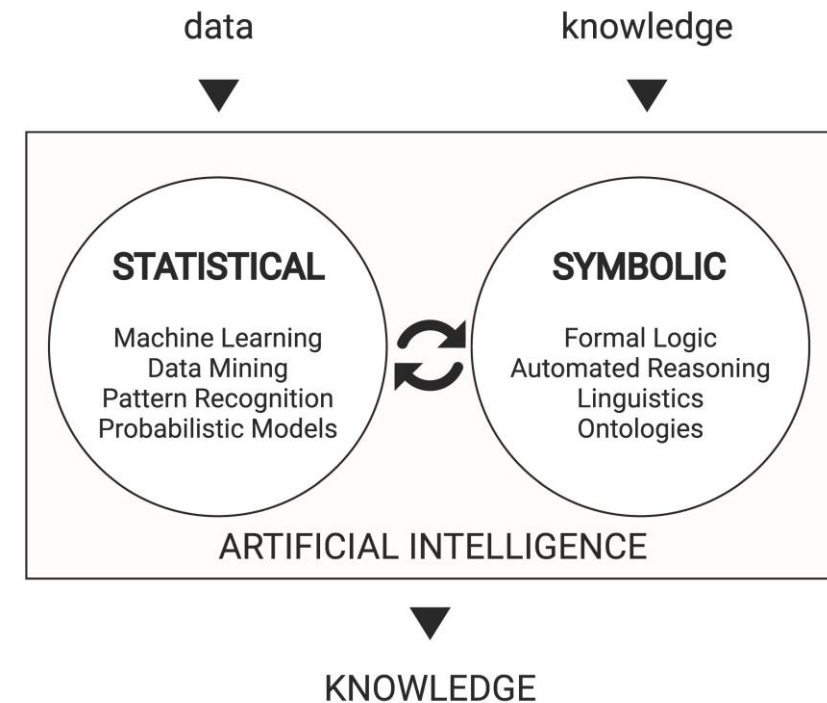


Figure 1: The Transformer - model architecture.

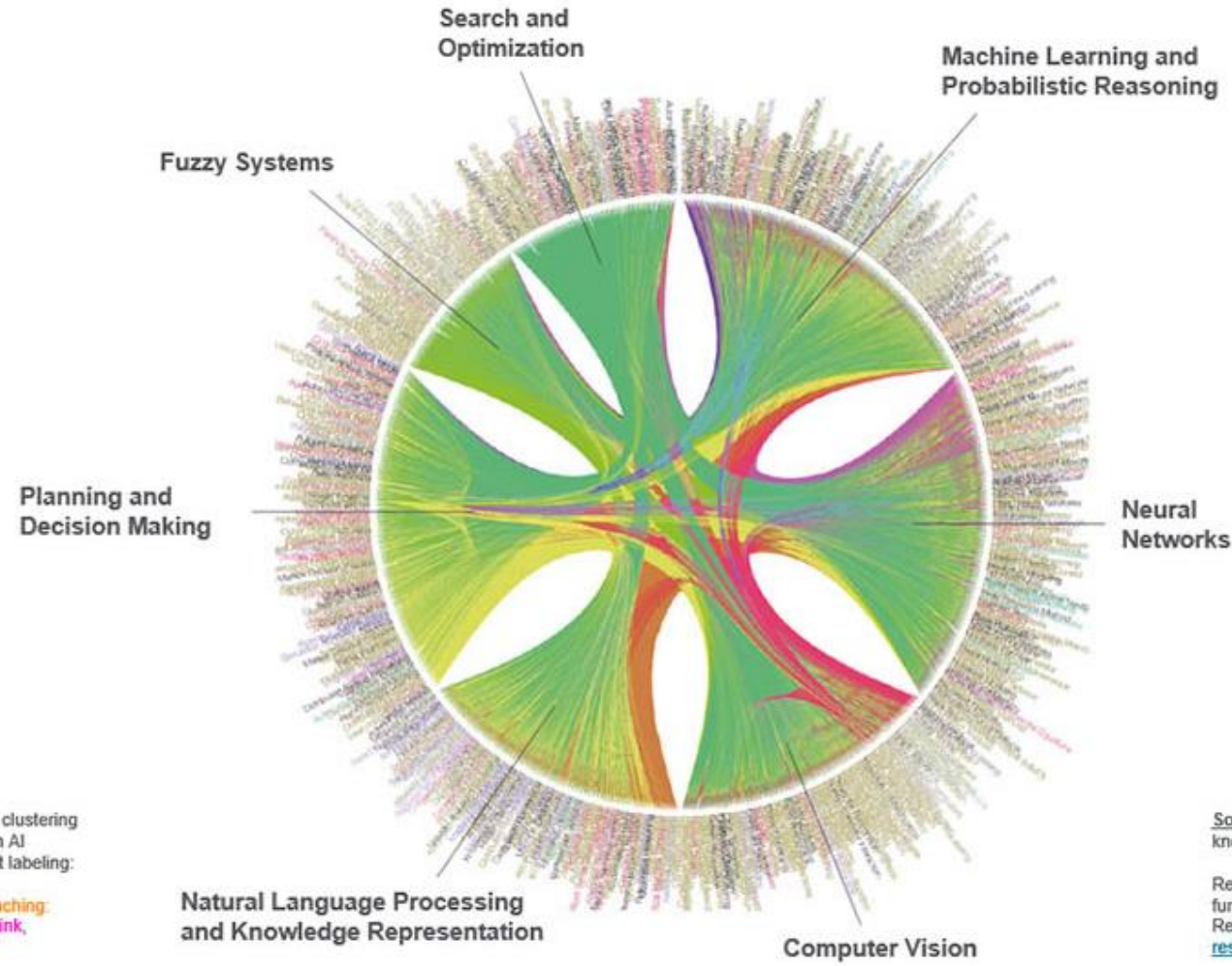
## Artificial Intelligence Composed

- The intelligence is intangible. It is composed of :-

- Reasoning
- Learning
- Problem Solving
- Perception
- Linguistic Intelligence



## The AI research field clusters around seven main research areas



Using keyword co-occurrences with unsupervised clustering on article-level, based on keywords extracted from AI documentation of different actors and using expert labeling:

The color of the keyword represents its origin: Teaching: orange, Industry: green, Research: blue, Media: pink, multiple: black

Source: Elsevier, 2018, Artificial Intelligence: how knowledge is created, transferred, and used

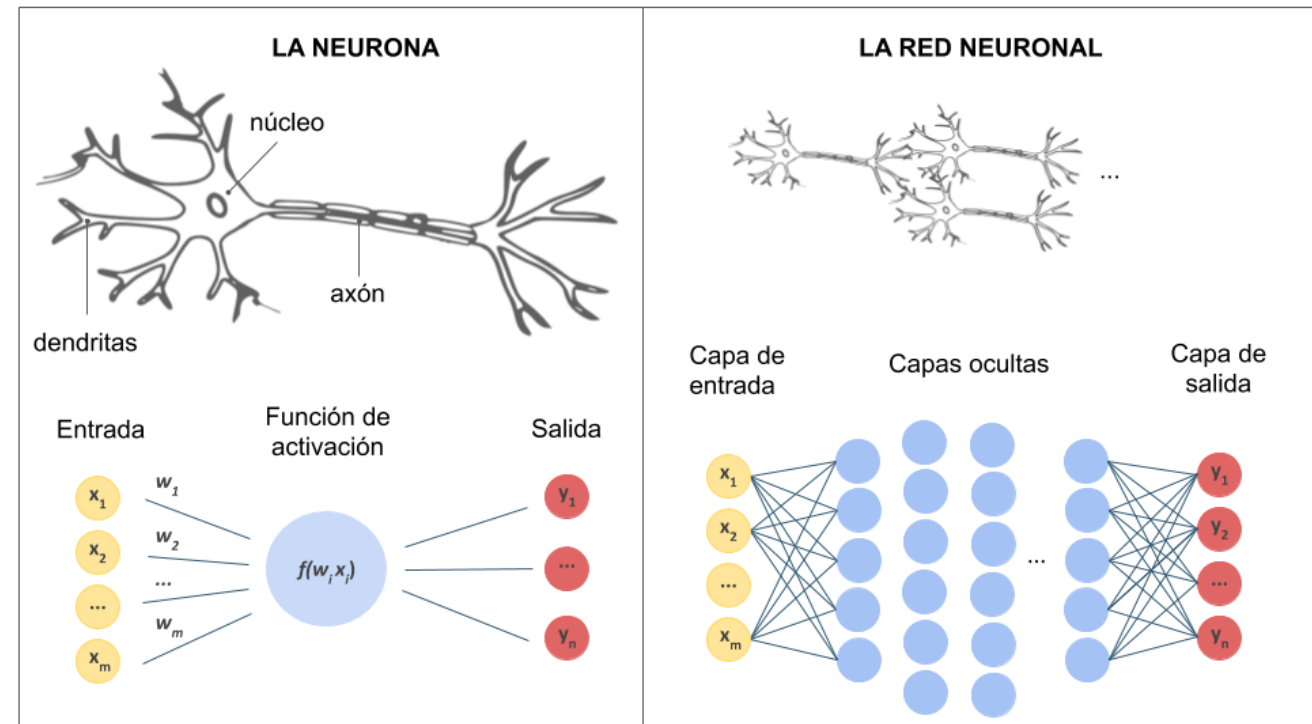
Report will be launched: Dec. 11<sup>th</sup>, 2018. Download and further information from Elsevier's Artificial Intelligence Resource Center: <https://www.elsevier.com/connect/ai-resource-center>

# *Modus operandi*

- 1) Se propone modelo hipótesis:
  - Definido por unos **hiperparámetros** (ej. Grado polinomio).
  - Acompañado por unos **parámetros** a entrenar.
- 2) **Estimador de Máxima Verosimilitud** sobre subconjunto de **entrenamiento**.
- 3) Subconjunto de **validación cruzada**: ¿mala generalización? → ajuste hiperparámetros.
- 4) Subconjunto de **test**: ¿mala generalización? → cambio de hipótesis.

# Redes neuronales (NN)

- 1) Modelo hipótesis basado (por lo menos remotamente) en la estructura del cerebro, NO en la naturaleza de los datos.
- 2) Teoremas de Aproximación Universal.
- 3) Fáciles de paralelizar → *Deep learning*.
- 4) Funcionamiento de la **memoria**.



# Procesamiento de Lenguaje Natural (NLP)

- Chomsky:
  - Distinción entre sintáctica y semántica.
  - Gramática innata vs teorías ambientalistas (marco para las NN).
- NN buenos resultados reproduciendo gramática, pero sin “sentido común” para la semántica.

# Procesamiento de Lenguaje Natural (NLP)

- **Tokenización:** palabras, caracteres, o subpalabras.
- Representaciones vectoriales de tókenes.
- *Bags of words.*
- Representaciones contextuales: modelos secuenciales.
  - *RNNs*: comprobación manual, corta memoria y difícil paralelización.
  - *Transformers.*



# *Transformers*. Mecanismo de atención

- Atención por token.
- Vector *query*  $\rightarrow$  para el token sobre el que se realiza la atención.
- Vectores *key* y *value*  $\rightarrow$  para cada token de contexto.

---

**Algoritmo 1:** Atención básica para un único token.

---

**Input:**  $e \in \mathbb{R}^{d_{in}}$ , *input embedding* del token a representar  $x_o$ .

**Input:**  $e_t \in \mathbb{R}^{d_{in}}$ , *input embedding* del token de contexto  $z_t \in T$ .

**Output:**  $\tilde{v} \in \mathbb{R}^{d_{out}}$ , la representación vectorial conjunta del token  $x_o$  y su contexto  $T$ .

**Parámetros:**  $W_q, W_k \in \mathbb{R}^{d_{attn} \times d_{in}}$ ,  $b_q, b_k \in \mathbb{R}^{d_{attn}}$ ; proyecciones lineales asociadas a *query* y a *key*.

**Parámetros:**  $W_v \in \mathbb{R}^{d_{out} \times d_{in}}$ ,  $b_v \in \mathbb{R}^{d_{out}}$ ; la proyección lineal asociada a *value*

1  $q \leftarrow W_q e + b_q$

2  $\forall t \in [N_T] : k_t \leftarrow W_k e_t + b_k$

3  $\forall t \in [N_T] : v_t \leftarrow W_v e_t + b_v$

4  $\forall t \in [N_T] : \alpha_t = \frac{\exp(q^T k_t / \sqrt{d_{attn}})}{\sum_u \exp(q^T k_u / \sqrt{d_{attn}})}$

5 **return**  $\tilde{v} = \sum_{t=1}^T \alpha_t v_t$

---

# Transformers. Mecanismo

- Atención para una secuencia de tokens.
  - Secuencia primaria  $\mathbf{X}$ , y de contexto  $\mathbf{Z}$ .
  - Softmax y Mask.
- Tipos de atención:
  - Atención unidireccional (*Unidirectional/masked self-attention*)  $\rightarrow \mathbf{X}=\mathbf{Z}$ ,  $\text{Mask}=[[t_Z \leq t_X]]$ .
  - Atención bidireccional (*Bidirectional/Unmasked self-attention*)  $\rightarrow \mathbf{X}=\mathbf{Z}$ ,  $\text{Mask}=1$ .
  - Atención cruzada (*Cross-attention*)  $\rightarrow \mathbf{X} \neq \mathbf{Z}$ ,  $\text{Mask}=1$ .
  - Atención multi-head (*Multi-head attention*).

---

**Algoritmo 2:**  $\tilde{\mathbf{V}} \leftarrow \text{Atención}(\mathbf{X}, \mathbf{Z} | \mathcal{W}_{qkv}, \text{Mask})$

---

*/\* Calcula una única head de self- o cross-attention (con máscara). \*/*

**Input:**  $\mathbf{X} \in \mathbb{R}^{d_X \times l_X}$ ,  $\mathbf{Z} \in \mathbb{R}^{d_Z \times l_Z}$  los *input embeddings* de la secuencias primaria y de contexto.

**Output:**  $\tilde{\mathbf{V}} \in \mathbb{R}^{d_{out} \times l_X}$  la representación vectorial conjunta de los tokens en  $\mathbf{X}$  y su contexto  $\mathbf{Z}$ .

**Parámetros:**  $\mathcal{W}_{qkv}$  formada por:  $\mathbf{W}_q \in \mathbb{R}^{d_{attn} \times d_X}$ ,  $\mathbf{b}_q \in \mathbb{R}^{d_{attn}}$ ,  
 $\mathbf{W}_k \in \mathbb{R}^{d_{attn} \times d_Z}$ ,  $\mathbf{b}_k \in \mathbb{R}^{d_{attn}}$ ,  
 $\mathbf{W}_v \in \mathbb{R}^{d_{out} \times d_Z}$ ,  $\mathbf{b}_v \in \mathbb{R}^{d_{out}}$ .

**Hiperparámetros:**  $\text{Mask} \in \{0, 1\}^{l_Z \times l_X}$

1  $\mathbf{Q} \leftarrow \mathbf{W}_q \mathbf{X} + \mathbf{b}_q \mathbf{1}^T$ ,  $\triangleright \mathbf{Query} \in \mathbb{R}^{d_{attn} \times l_X}$

2  $\mathbf{K} \leftarrow \mathbf{W}_k \mathbf{Z} + \mathbf{b}_k \mathbf{1}^T$ ,  $\triangleright \mathbf{Key} \in \mathbb{R}^{d_{attn} \times l_Z}$

3  $\mathbf{V} \leftarrow \mathbf{W}_v \mathbf{Z} + \mathbf{b}_v \mathbf{1}^T$ ,  $\triangleright \mathbf{Value} \in \mathbb{R}^{d_{out} \times l_Z}$

4  $\mathbf{S} \leftarrow \mathbf{K}^T \mathbf{Q}$ ,  $\triangleright \mathbf{Score} \in \mathbb{R}^{l_Z \times l_X}$

5  $\forall t_Z, t_X$  si  $\text{Mask}[t_X, t_Z] \neq 1$  entonces  $\mathbf{S}[t_X, t_Z] \leftarrow -\infty$

6 **return**  $\tilde{\mathbf{V}} = \mathbf{V} \cdot \text{softmax}(\mathbf{S} / \sqrt{d_{attn}})$

---

# Transformers: Encoder-Decoder original

- **Objetivo:** entrenar en predecir secuencia de tokens  $x$  dada otra secuencia de tokens original  $z$ .
- **Núcleo:** varios bloques de atención distribuidos en dos módulos: un *encoder* y un *decoder*.

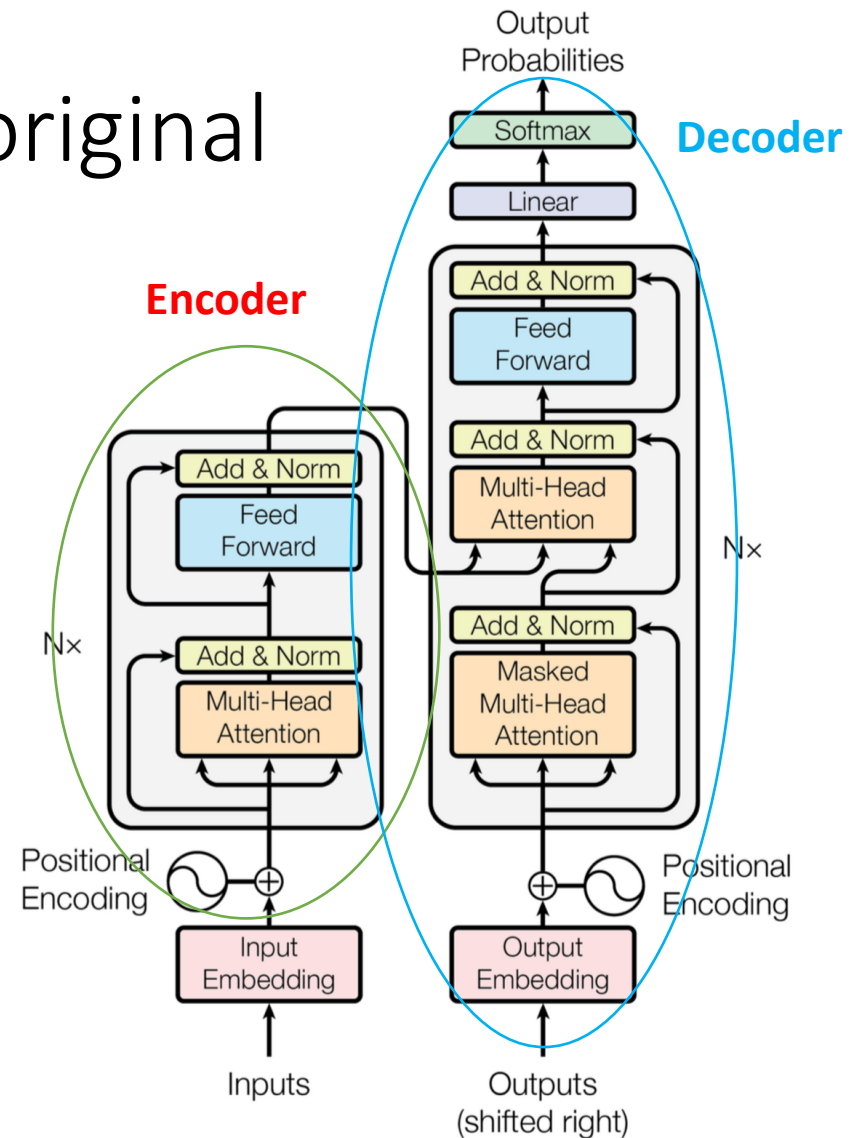


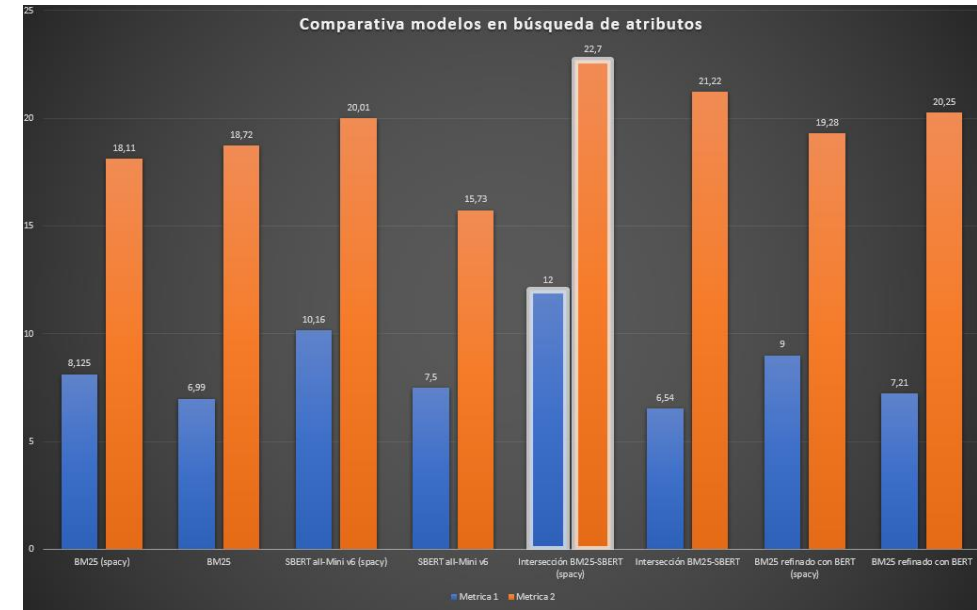
Figure 1: The Transformer - model architecture.

# BERT y GPT

- El entrenamiento de un *transformer* se centra en predecir un solo token.
- *Bidirectional encoders* → BERT, RoBERTa
- *Autoregressive decoders* → GPT, Gopher

# SBERT y aplicación práctica

- *Sentence similarity: cross-encoders vs bi-encoders.*
- Búsqueda de atributos más relevantes para una lista de productos.



# Valoración

- 1) Muchas más información lingüística que cualquier ser humano pero peores resultados.
- 2) Los *transformers* Sí que están basados en la naturaleza del lenguaje, además de la estructura neuronal del cerebro. Pero no es suficiente.
- 3) Posible vía para reflejar “sentido común” → Causalidad.

# Causalidad

- Una causa no es una “correlación fuerte”.
- Los datos son resultados de una jerarquía causal invisible, pero necesaria para una descripción correcta de la realidad.
- Jerarquía causal representada por grafos acíclicos: diagramas causales.

# Escalera de causalidad

- 1) Estadística: observación de datos.
- 2) Intervención: nuevo operador en probabilidad: *do()*
- 3) Razonamiento contrafáctico: ¿qué pasaría si en vez de A fuese B? → Imaginación (información desde datos no existentes)

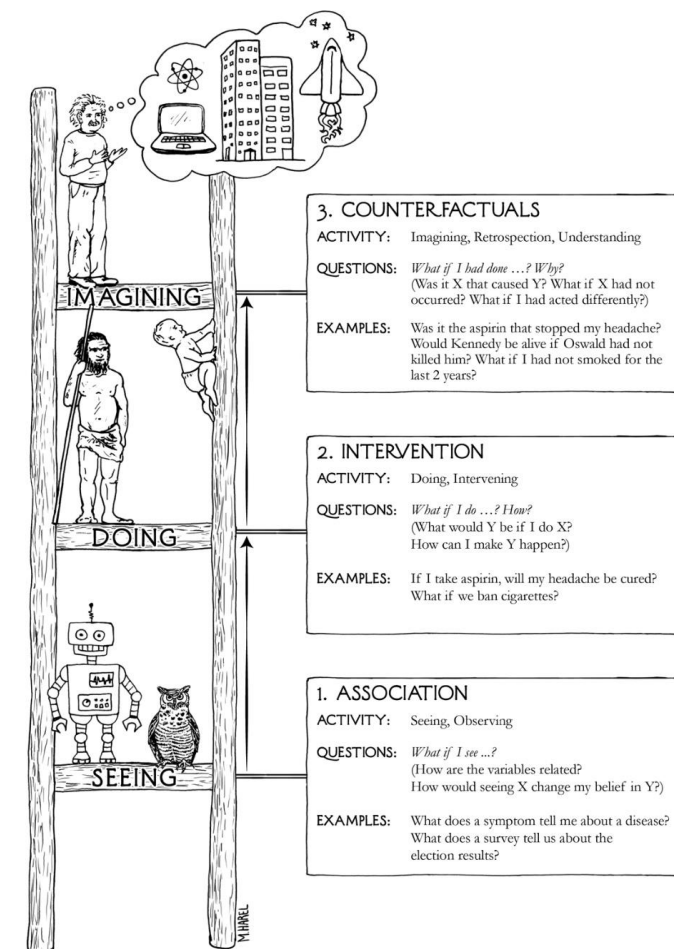


FIGURE 1.2. The Ladder of Causation, with representative organisms at each level. Most animals, as well as present-day learning machines, are on the first rung, learning from association. Tool users, such as early humans, are on the second rung if they act by planning and not merely by imitation. We can also use experiments to learn the effects of interventions, and presumably this is how babies acquire much of their causal knowledge. Counterfactual learners, on the top rung, can imagine worlds that do not exist and infer reasons for observed phenomena. (Source: Drawing by Maayan Harel.)



# Modelo de razonamiento contrafáctico<sup>[1]</sup>

- No construye estadística nueva.
- Aprende razonamiento contrafáctico de ejemplos contrafácticos.
- Como ocurría con los *transformers*, su estructura responde a la naturaleza del problema que intenta resolver.

[1] F. Feng, J. Zhang, X. He, H. Zhang, and T.-S. Chua, “Empowering Language Understanding with Counterfactual Reasoning,” vol. abs/2106.03046, 2021.

# Conclusiones y líneas futuras.

- Distinción entre acercamiento simbólico y estadístico es ambigua con los *Transformers*.
- Obtienen mejores resultados que los métodos tradicionales en un problema de búsqueda y recuperación de la información.
- Consiguen reproducir de forma adecuada la gramática del lenguaje desde un marco ambientalista, pero no así su semántica.
- La causalidad es una buena baza para línea de investigación.
  - Construcción con inferencia causal en vez de con estadística clásica.
  - Aprendizaje estadístico a partir de ejemplos contrafácticos.

Gracias por su “atención”