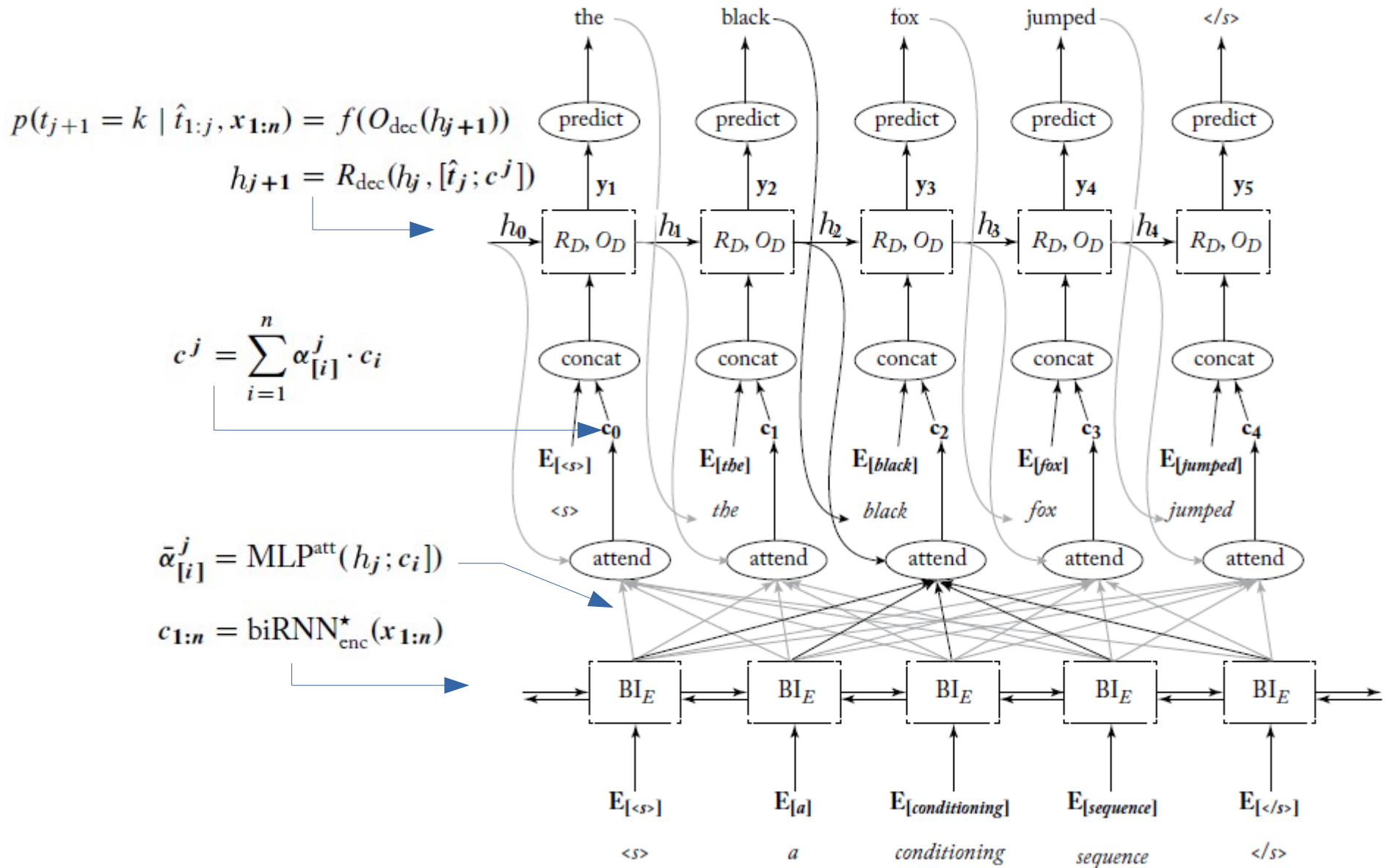




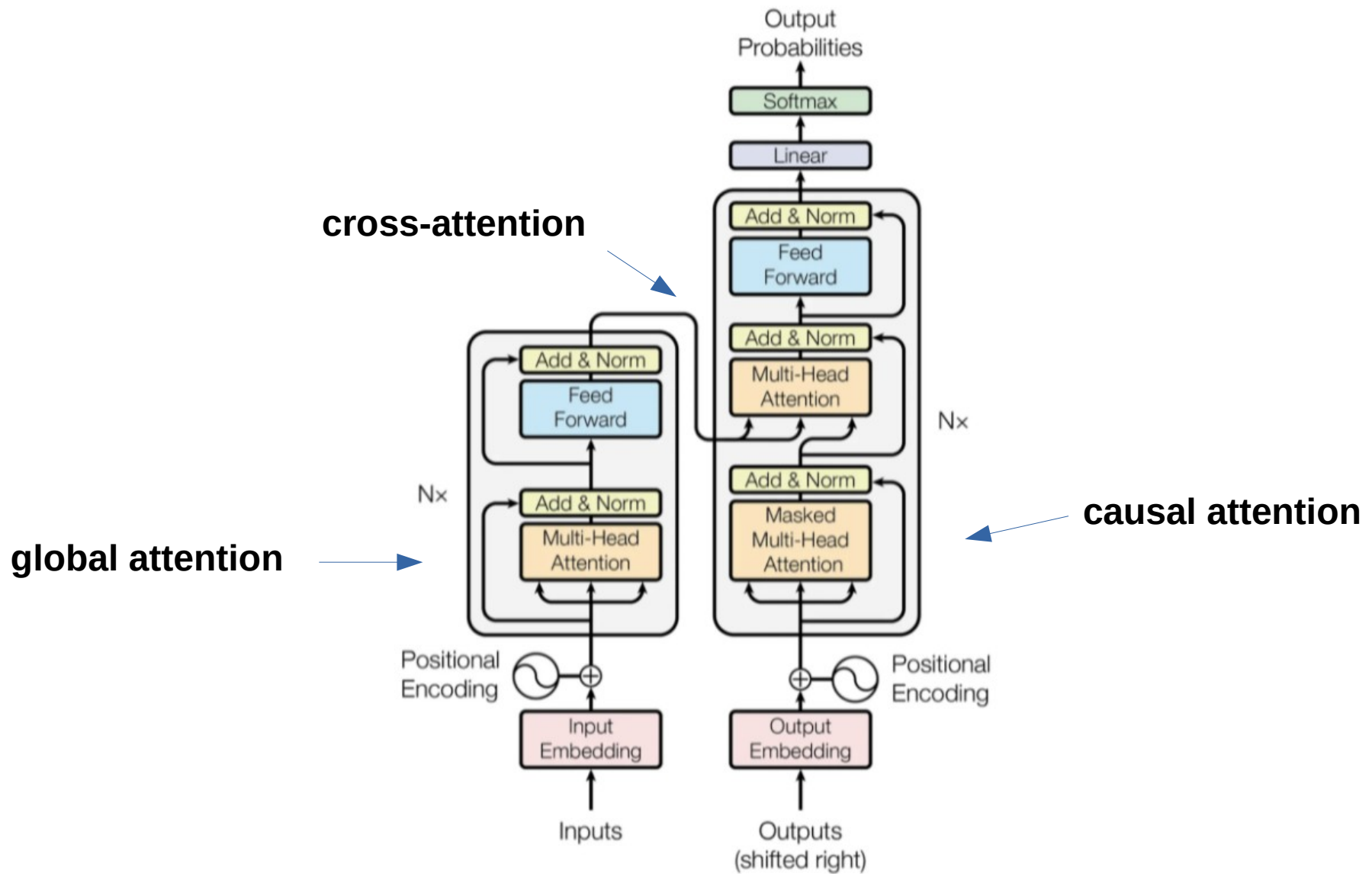
# IIC3670 Procesamiento de Lenguaje Natural

<https://github.com/marcelomendoza/IIC3670>

## Generación condicional con atención

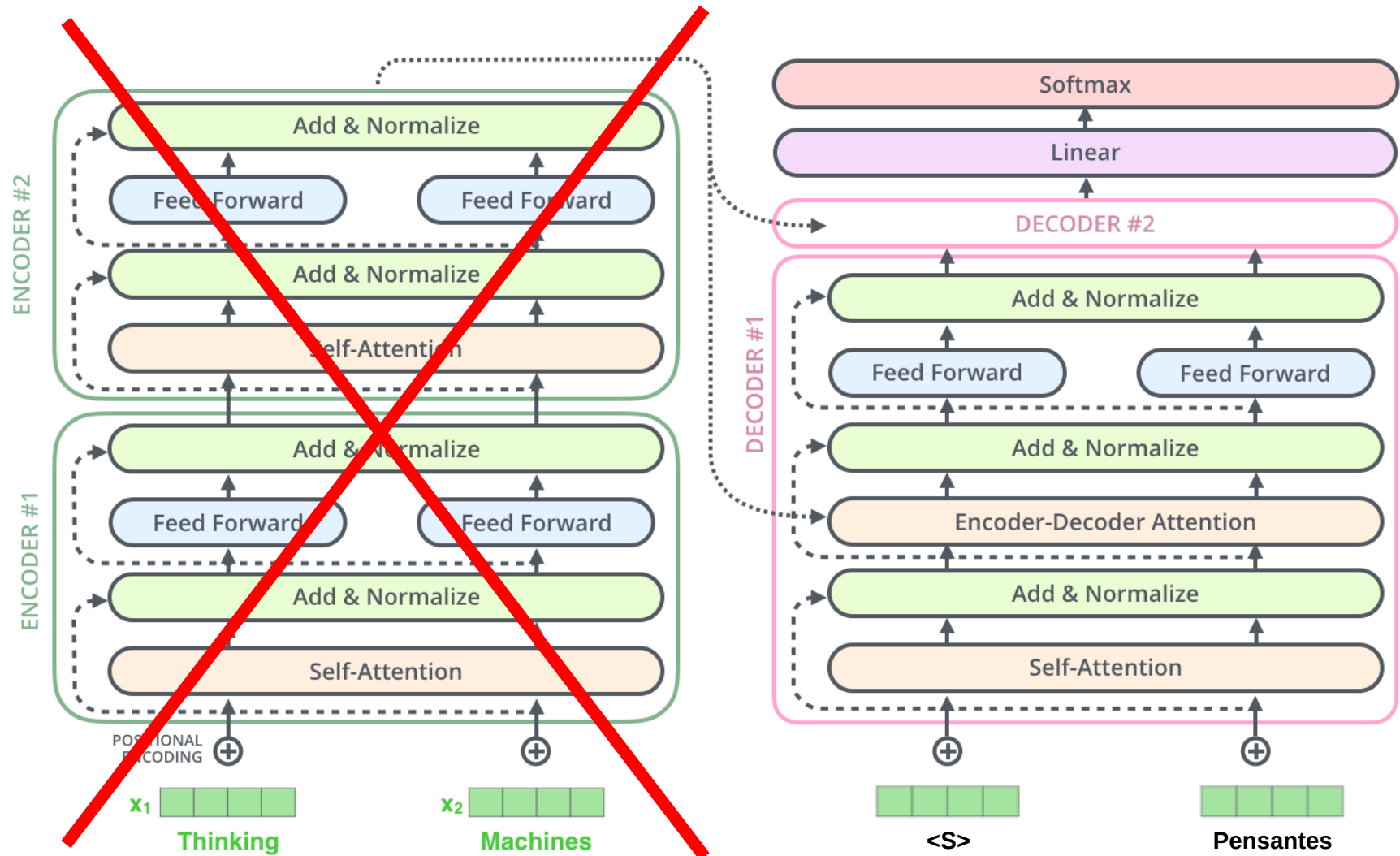


## Generación condicional con atención (transformer seq2seq)

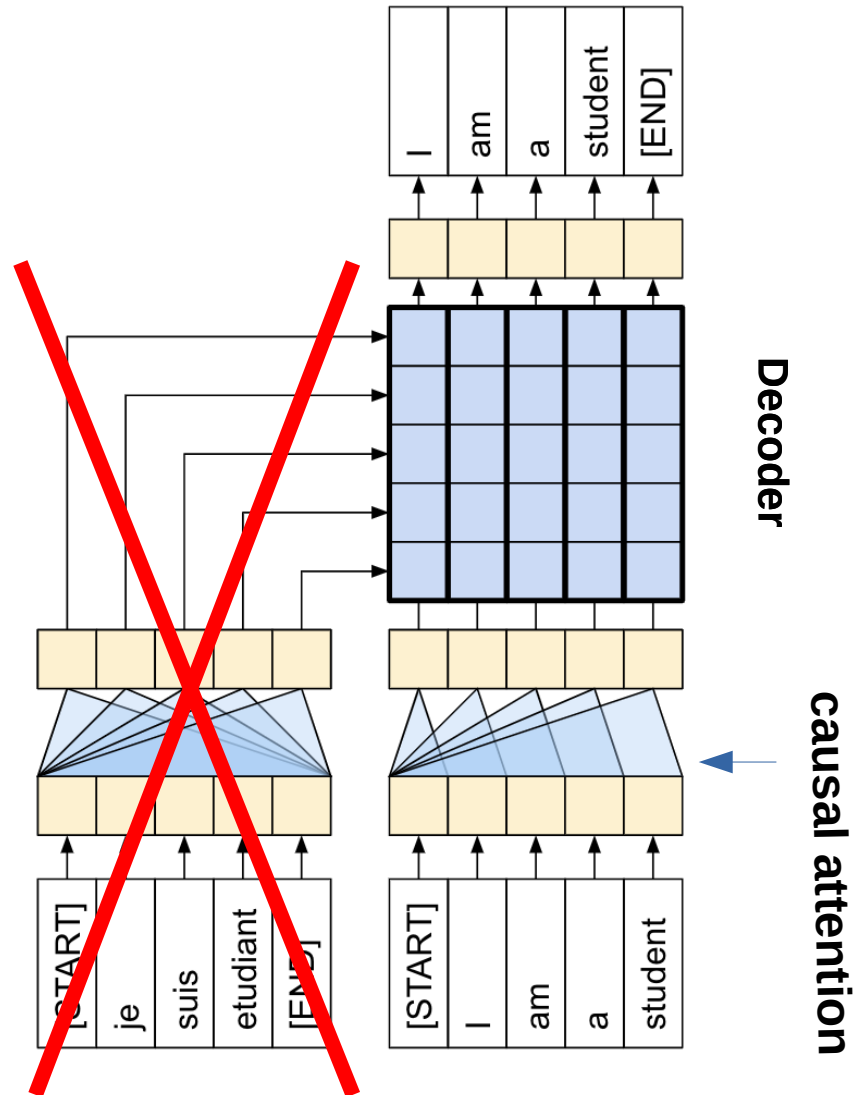


- GENERATIVE PRETRAINED TRANSFORMERS -

Se elimina el encoder!!!



Se elimina el encoder!!!



# Transformer decoder

Seq2seq

- Cada secuencia  $(m^1, \dots, m^n) \mapsto (y^1, \dots, y^\eta)$

se transforma en:

► Símbolo separador

$$(w^1, \dots, w^{n+\eta+1}) = (m^1, \dots, m^n, \delta, y^1, \dots, y^\eta)$$

- Luego, el **decoder** resuelve la siguiente tarea:

secuencias largas

$$p(w^1, \dots, w^{n+\eta}) = \prod_{j=1}^{n+\eta} p(w^j | w^1, \dots, w^{j-1})$$

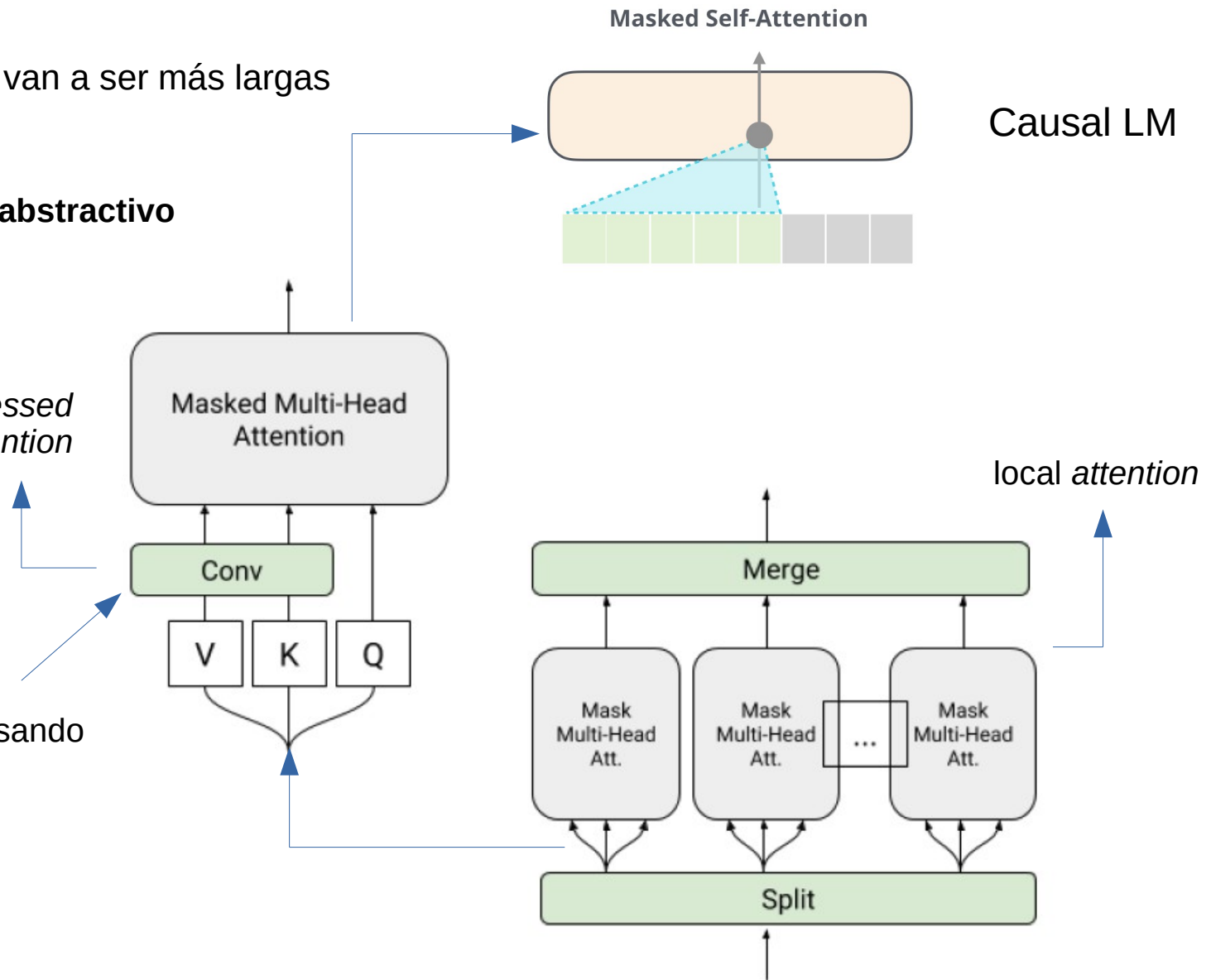
# Transformer decoder

Las oraciones van a ser más largas

Modelo **abstractivo**

*Memory-compressed attention*

Reduce las V y K usando kernels de 3x3

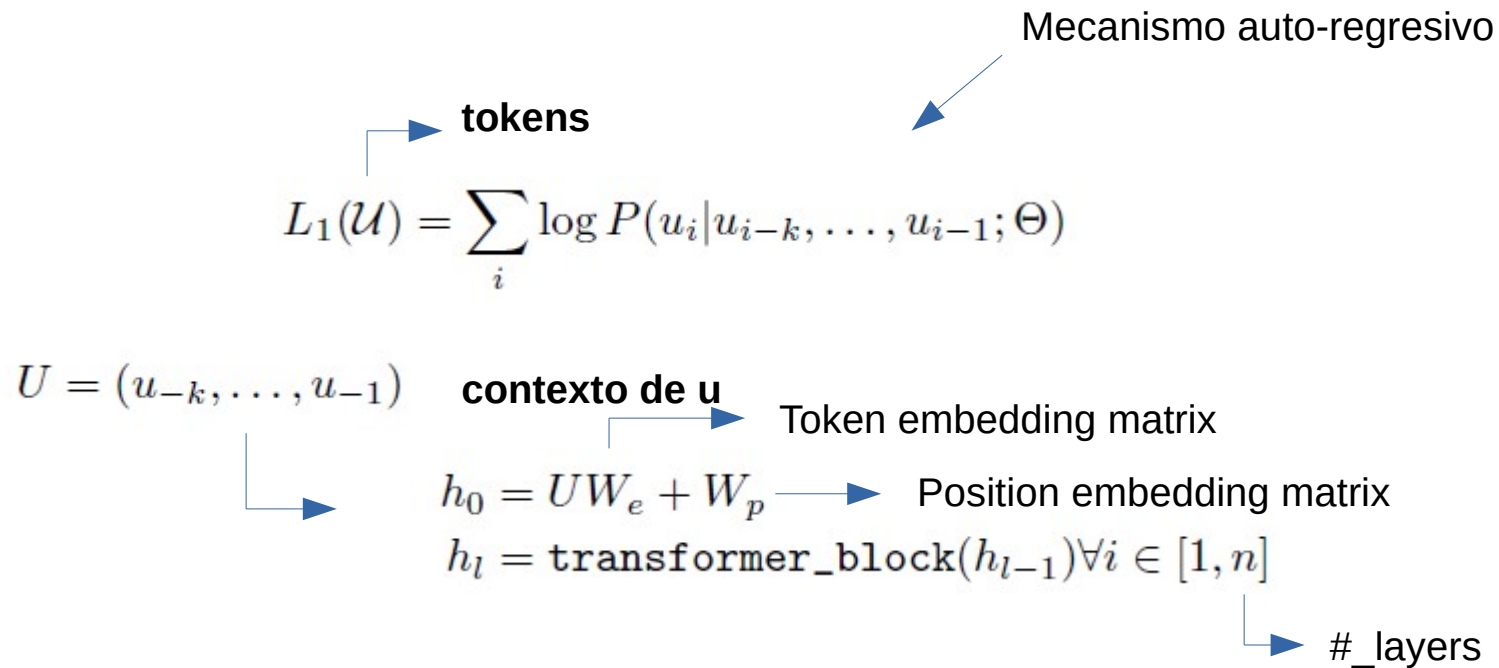


Divide la oración en *batches* del mismo tamaño



# Transformer decoder

LLM



task

$$P(u) = \text{softmax}(h_n W_e^T)$$

Idea propuesta en:



GENERATING WIKIPEDIA BY SUMMARIZING LONG SEQUENCES

Peter J. Liu, Mohammad Saleh,  
Etienne Pot, Ben Goodrich, Ryan Sepassi, Łukasz Kaiser, Noam  
Shazeer  
ICLR 2018

# GPT 1

Transformer decoder + supervised fine tuning

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta) \quad \text{Transformer decoder}$$

$$P(y | x^1, \dots, x^m) = \text{softmax}(h_l^m W_y). \quad \text{Supervised fine-tuning}$$

label  $\leftarrow$   $\rightarrow$  Última capa del transformer decoder

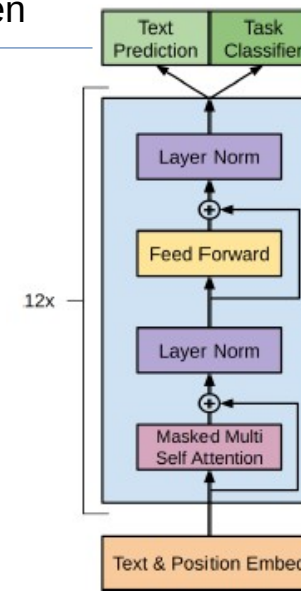
en LM y es un token

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y | x^1, \dots, x^m).$$

La pérdida combinada es:

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

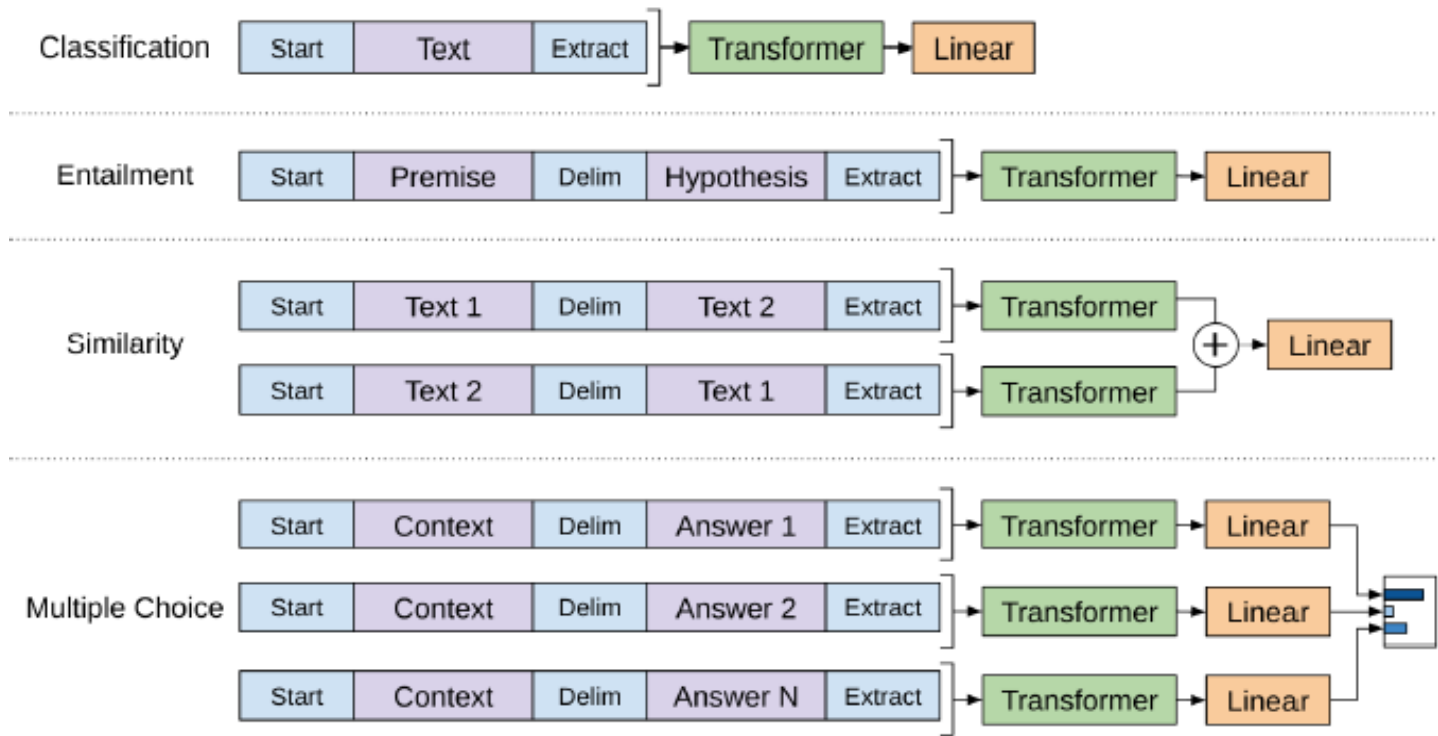
GPT 1 resuelve ambas tareas simultáneamente (LM y specific task)  $\rightarrow$



# GPT 1

## Supervised fine-tuning

La secuencia de entrada depende de la tarea



GPT 1



Radford, A., Narasimhan, K., Salimans, T., Sutskever, I. Improving language understanding by generative pre-training, 2018.

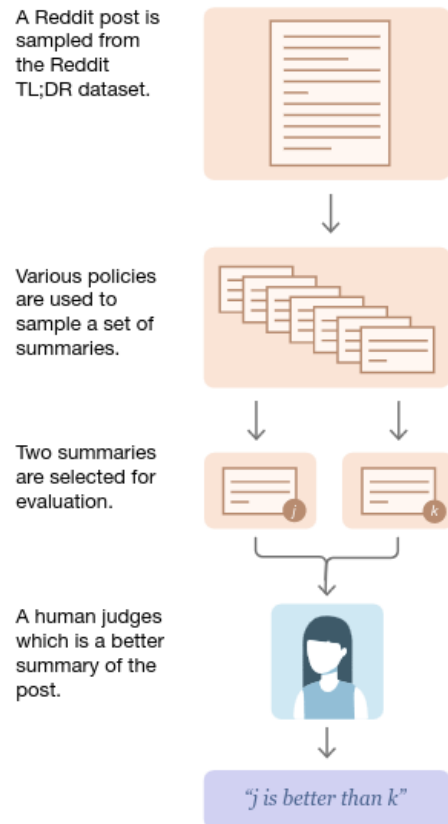


- INSTRUCT GPT -

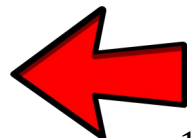
# Human in the loop

Inicialmente este tema se abordó para construcción de resúmenes

## 1 Collect human feedback



Stiennon et al. Learning to summarize from human feedback, NeurIPS 2020



# Human in the loop

Inicialmente este tema se abordó para construcción de resúmenes

## 1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



"j is better than k"

## 2 Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward  $r$  for each summary.



$r_j$

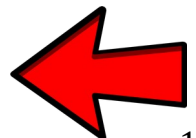
$r_k$

The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$\text{loss} = \log(\sigma(r_j - r_k))$$

"j is better than k"

Stiennon et al. Learning to summarize from human feedback, NeurIPS 2020



# Human in the loop

Inicialmente este tema se abordó para construcción de resúmenes

## 1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



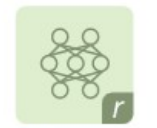
"j is better than k"

## 2 Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward  $r$  for each summary.



$r_j$

$r_k$

The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$\text{loss} = \log(\sigma(r_j - r_k))$$

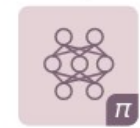
"j is better than k"

## 3 Train policy with PPO

A new post is sampled from the dataset.



The policy  $\pi$  generates a summary for the post.



The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.

$r$

Stiennon et al. Learning to summarize from human feedback, NeurIPS 2020

## Reward Model (RM)

- El RM es un modelo lineal que entrega un escalar (regresión).
- Se entrena el modelo para predecir a partir de resúmenes cuál es mejor:

$$\text{loss}(r_\theta) = -E_{(x, y_0, y_1, i) \sim D} [\log(\sigma(r_\theta(x, y_i) - r_\theta(x, y_{1-i})))]$$

resumen preferido

$y \in \{y_0, y_1\}$

- Una vez que termina el entrenamiento, se normaliza el RM para que el reward tenga un score medio en 0.



# Proximal Policy Optimization (PPO)

## 3 Train policy with PPO

A new post is sampled from the dataset.



The policy  $\pi$  generates a summary for the post.



The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.



- La salida del RM es un reward para el sistema.
- Se penaliza la divergencia entre el modelo mejorado (RL) y el original SFT para evitar que las salidas de los modelo sean muy distintas a las vistas por el RM durante el entrenamiento.
- El full reward queda dado por:

$$R(x, y) = r_{\theta}(x, y) - \beta \log[\pi_{\phi}^{\text{RL}}(y|x) / \pi^{\text{SFT}}(y|x)]$$