

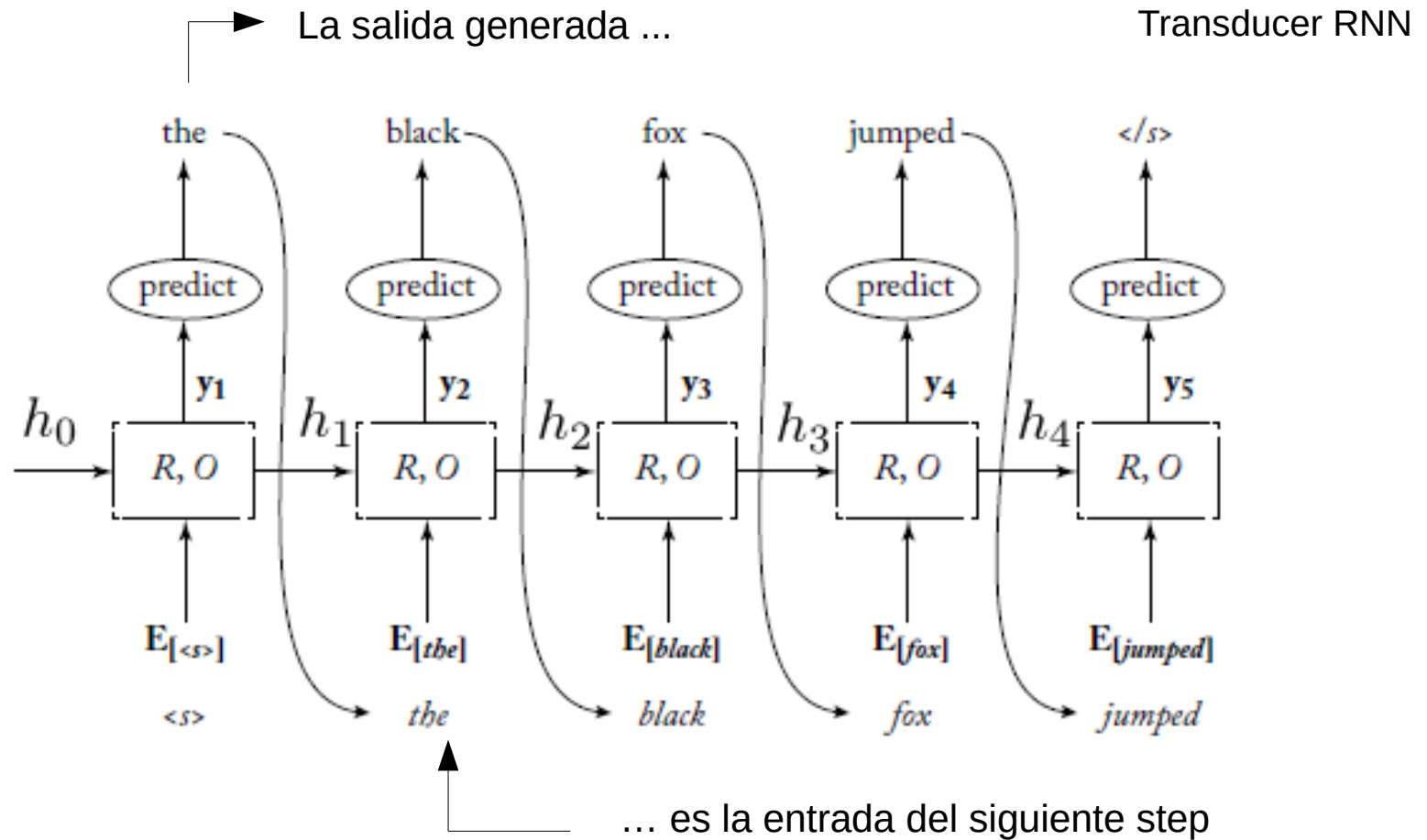
IIC3670 Procesamiento de Lenguaje Natural

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

- GENERACIÓN DE TEXTO -

Generación de texto

Generación auto-regresiva



Generación de texto

Generación condicional

$$\hat{t}_{j+1} \sim p(t_{j+1} = k \mid \hat{t}_{1:j}).$$

Transducer RNN:

$$p(t_{j+1} = k \mid \hat{t}_{1:j}) = f(\text{RNN}(\hat{t}_{1:j}))$$

$$p(t_{j+1} = k \mid \hat{t}_{1:j}) = f(O(h_{j+1}))$$

Generación de texto

Generación condicional

$$\hat{t}_{j+1} \sim p(t_{j+1} = k \mid \hat{t}_{1:j}).$$

Transducer RNN:

$$p(t_{j+1} = k \mid \hat{t}_{1:j}) = f(\text{RNN}(\hat{t}_{1:j}))$$

$$p(t_{j+1} = k \mid \hat{t}_{1:j}) = f(O(h_{j+1}))$$

└─ vector de contexto

Transducer RNN condicional:

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, c) = f(\text{RNN}(v_{1:j}))$$

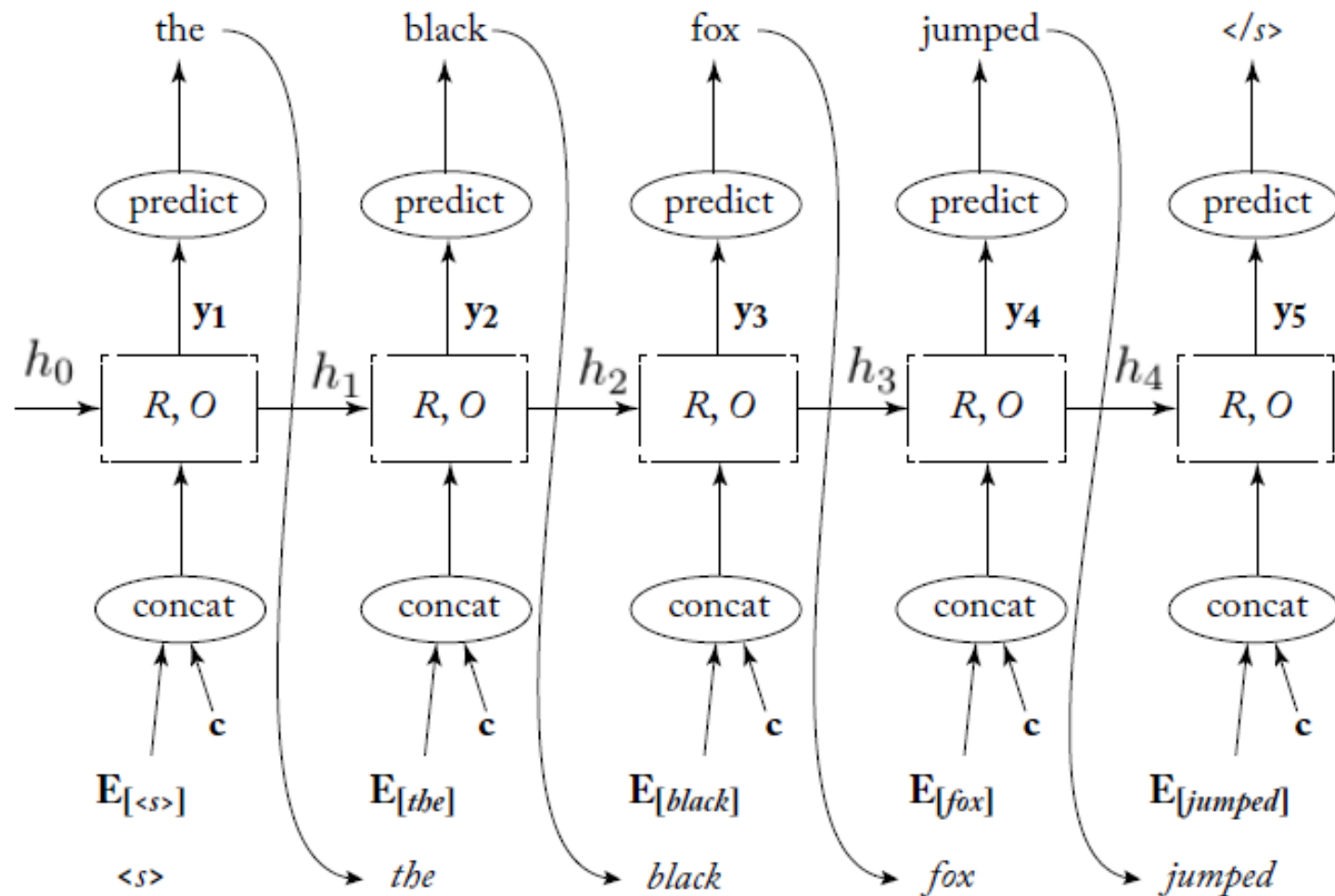
$$v_i = [\hat{t}_i; c]$$

$$\hat{t}_j \sim p(t_j \mid \hat{t}_{1:j-1}, c),$$

Generación de texto

c puede representar un tema u
otra sentencia

Transducer RNN condicional (generador):



Generación de texto

Generación condicional (encoder-decoder)

Sequence to sequence: \mathbf{c} es un vector que representa una secuencia de texto de entrada (encoder)

source sequence $\mathbf{x}_{1:n}$ \longrightarrow target output $\mathbf{t}_{1:m}$

Encoder: RNN: $\mathbf{c} = \text{RNN}^{\text{enc}}(\mathbf{x}_{1:n})$.

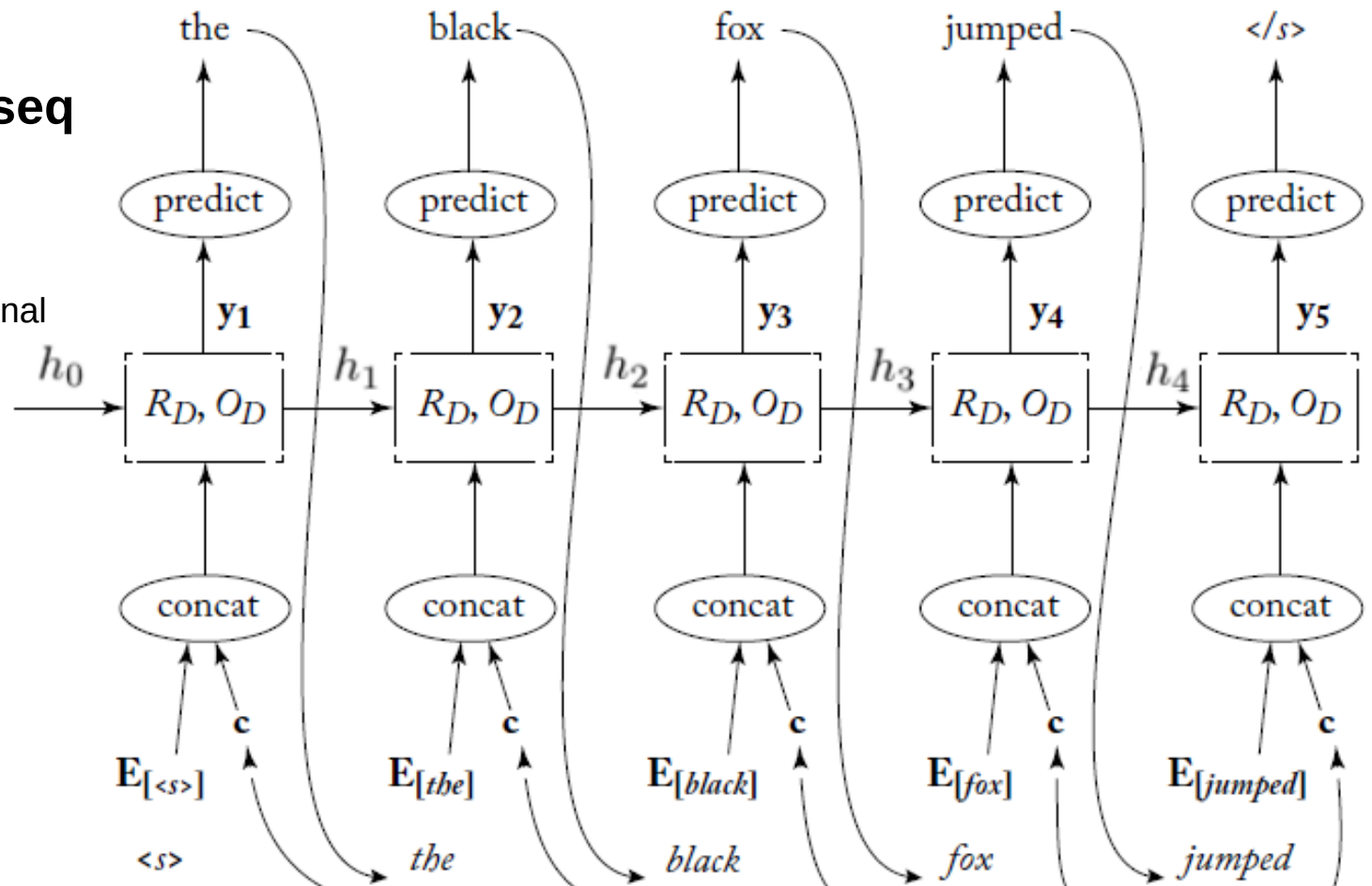
Decoder: $p(t_{j+1} = k \mid \hat{\mathbf{t}}_{1:j}, \mathbf{c}) = f(\text{RNN}(\mathbf{v}_{1:j}))$

$$\mathbf{v}_i = [\hat{t}_i; \mathbf{c}]$$

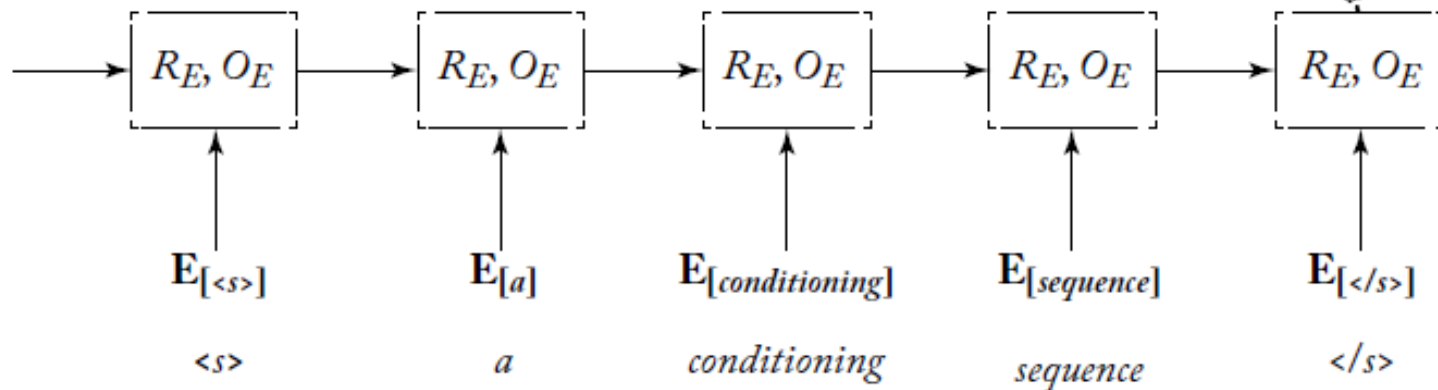
$$\hat{t}_j \sim p(t_j \mid \hat{\mathbf{t}}_{1:j-1}, \mathbf{c}),$$

Seq2seq

Transducer RNN condicional



RNN Encoder

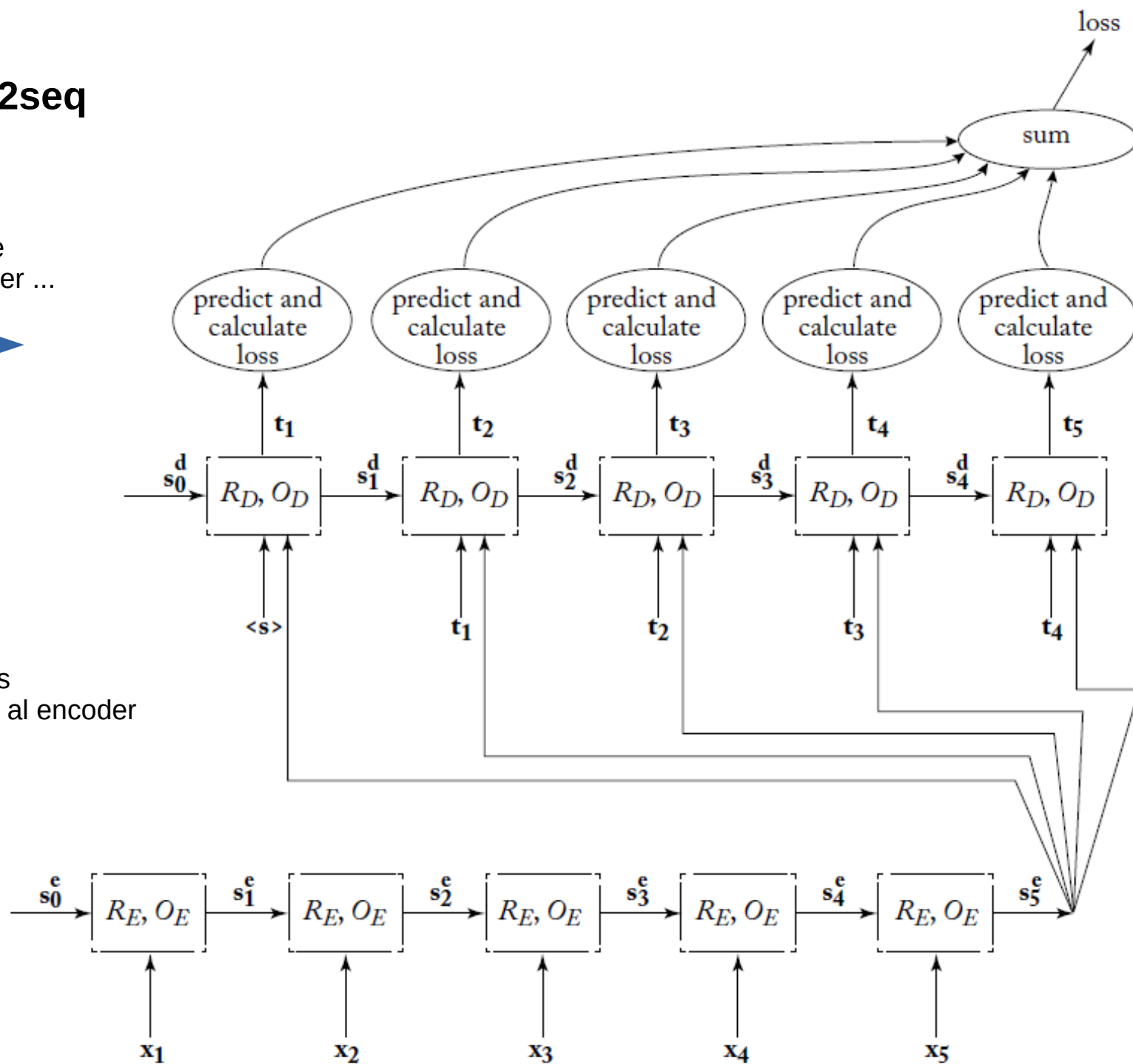


Seq2seq

La supervisión se hace en el decoder ...

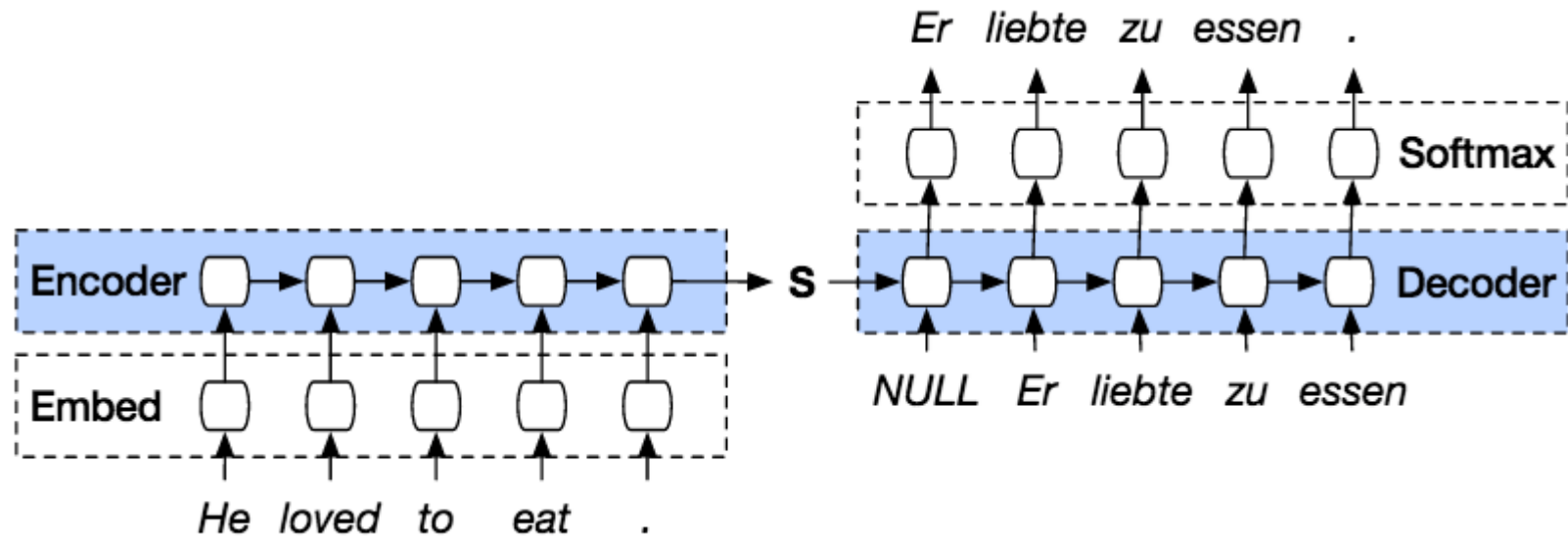


... pero los gradientes también se propagan al encoder



Seq2seq

- Aplicación más conocida de la arquitectura: **Machine translation**

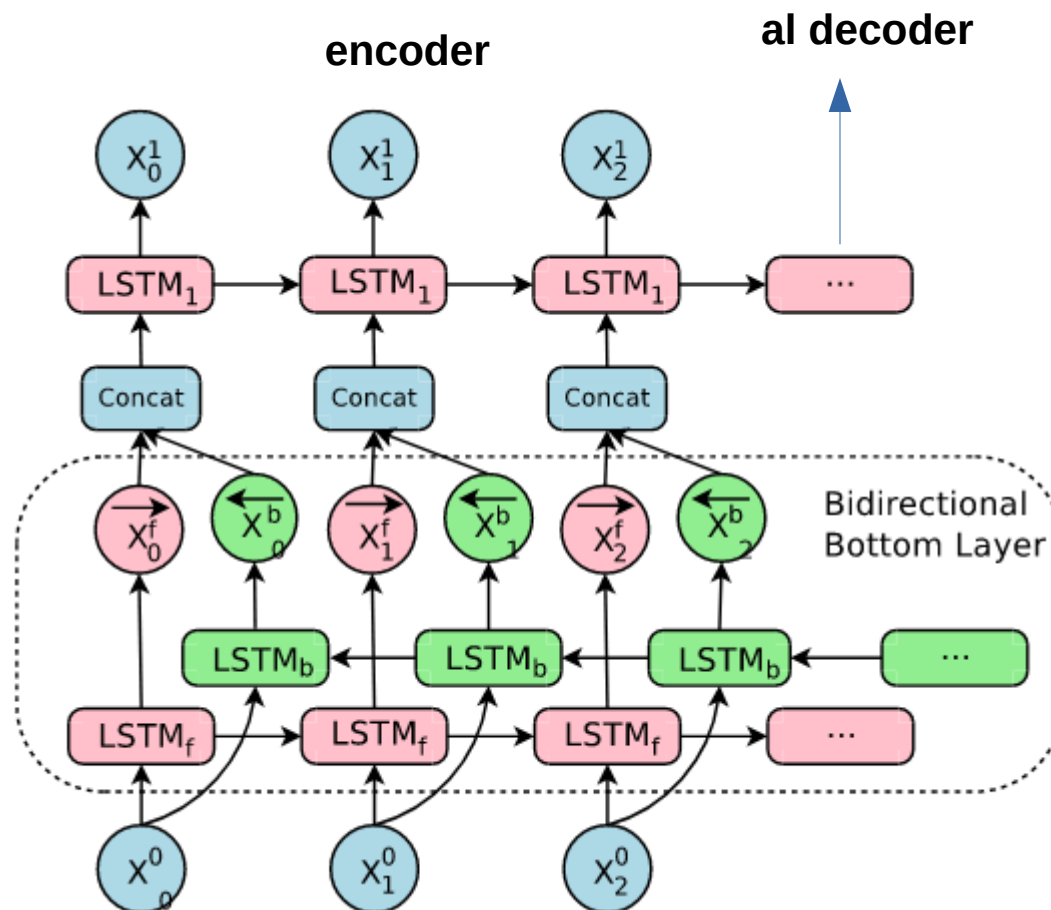


- Otras: **QA (closed)**

Seq2seq

Google Machine translation:

Usan un **encoder** más complejo



Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation, <https://arxiv.org/pdf/1609.08144.pdf> , 2016

- GENERACIÓN CONDICIONAL CON ATENCIÓN -

Generación condicional con atención

Idea: usar un mecanismo de atención en el **decoder**

Largo de la sentencia de entrada (encoder) \leftarrow $\text{attend}(c_{1:n}, \hat{t}_{1:j}) = c^j$ \rightarrow Posición en la sentencia de salida (decoder)

\rightarrow factores de atención

\rightarrow vectores de contexto (encoder)

$$c^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot c_i$$

los factores se aprenden con una **softmax** \leftarrow

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

$$\bar{\alpha}_{[i]}^j = \text{MLP}^{\text{att}}([h_j; c_i])$$

\rightarrow la **softmax** opera a la salida de una MLP que opera sobre h_j y c_i

Ojo, h del decoder!!! \leftarrow

Generación condicional con atención

encoder

$$c_{1:n} = \text{biRNN}_{\text{enc}}^*(x_{1:n}) \quad \text{decoder}$$

$$\bar{\alpha}_{[i]}^j = \text{MLP}^{\text{att}}(h_j; c_i) \rightarrow \text{MLP}^{\text{att}}([h_j; c_i]) = v \tanh([h_j; c_i]U + b)$$

soft attention

$$\alpha^j = \text{softmax}(\bar{\alpha}_{[1]}^j, \dots, \bar{\alpha}_{[n]}^j)$$

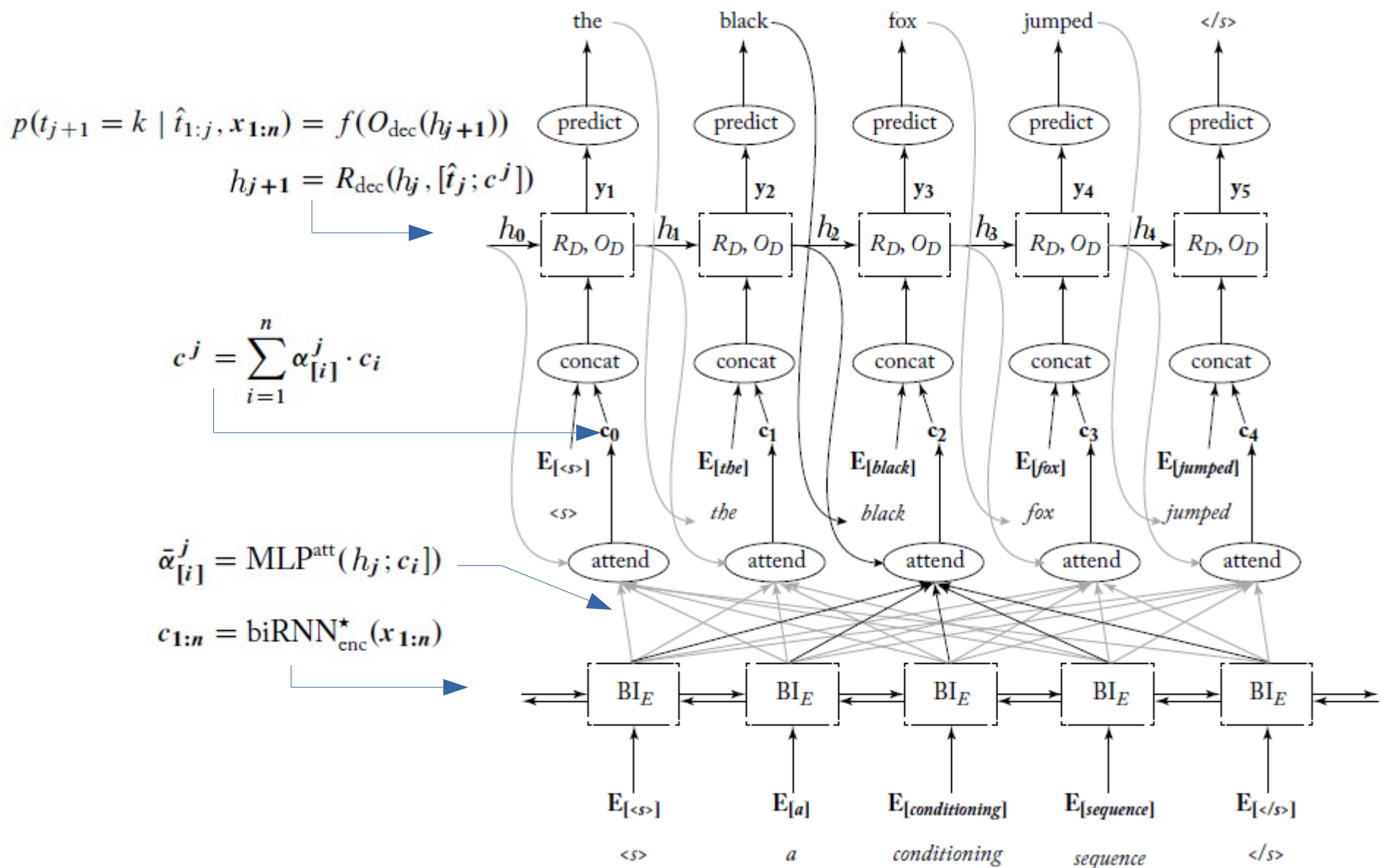
$$c^j = \sum_{i=1}^n \alpha_{[i]}^j \cdot c_i$$

$$p(t_{j+1} = k \mid \hat{t}_{1:j}, x_{1:n}) = f(O_{\text{dec}}(h_{j+1}))$$

decoder

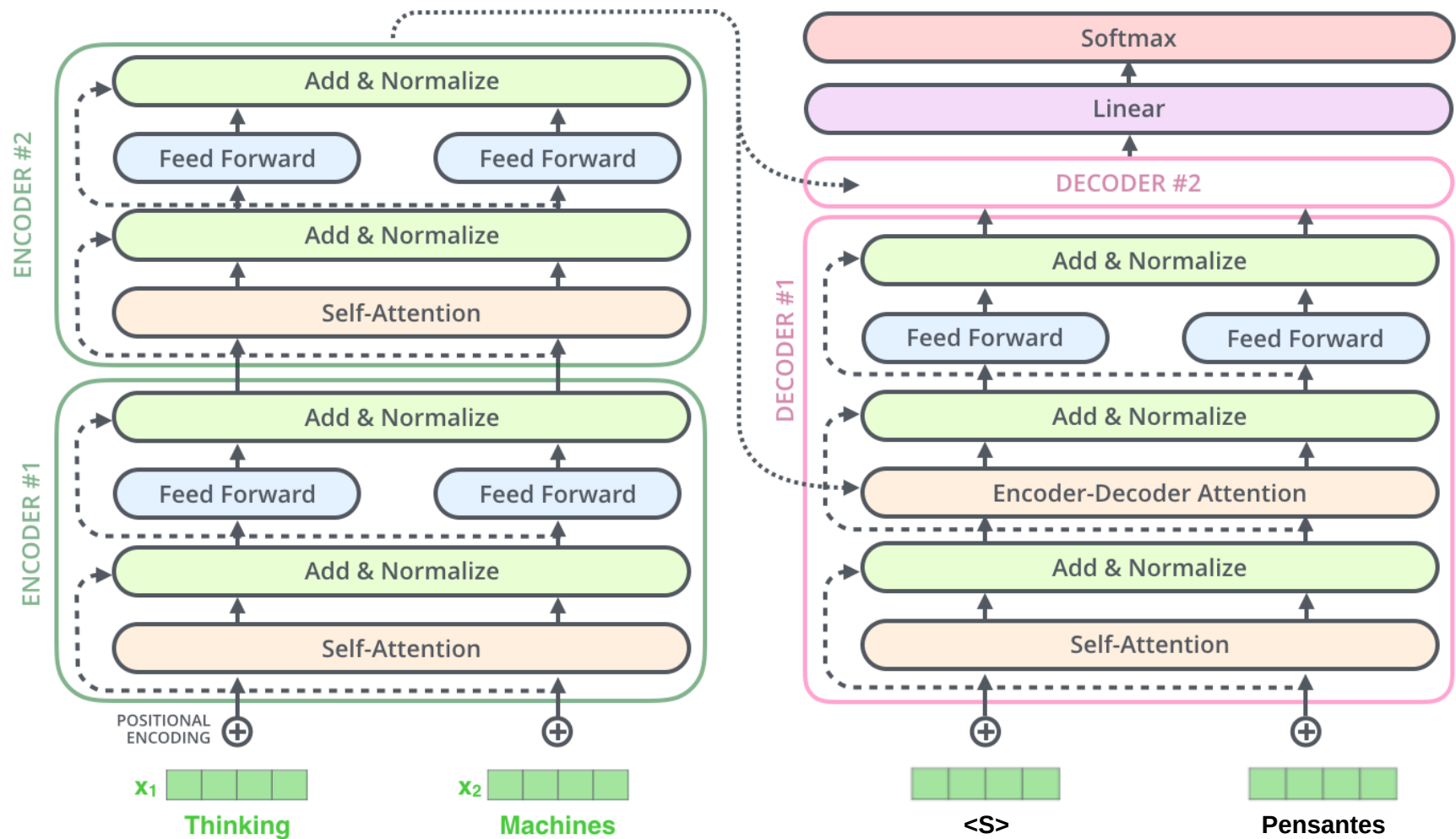
$$h_{j+1} = R_{\text{dec}}(h_j, [\hat{t}_j; c^j])$$

Generación condicional con atención

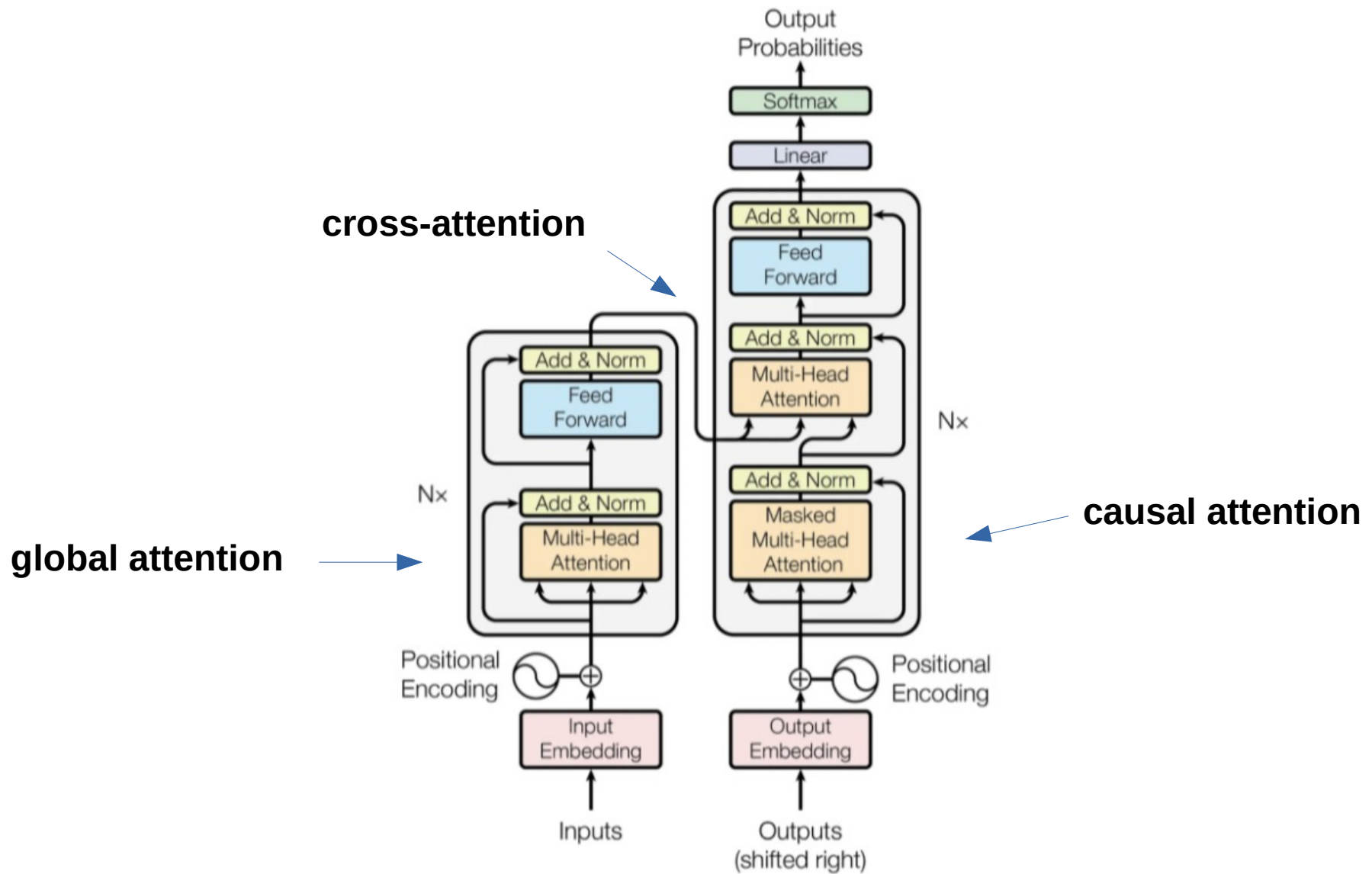


- GENERACIÓN CONDICIONAL CON TRANSFORMERS -

Generación condicional con atención (transformer seq2seq)



Generación condicional con atención (transformer seq2seq)



Generación condicional con atención (transformer seq2seq)

