

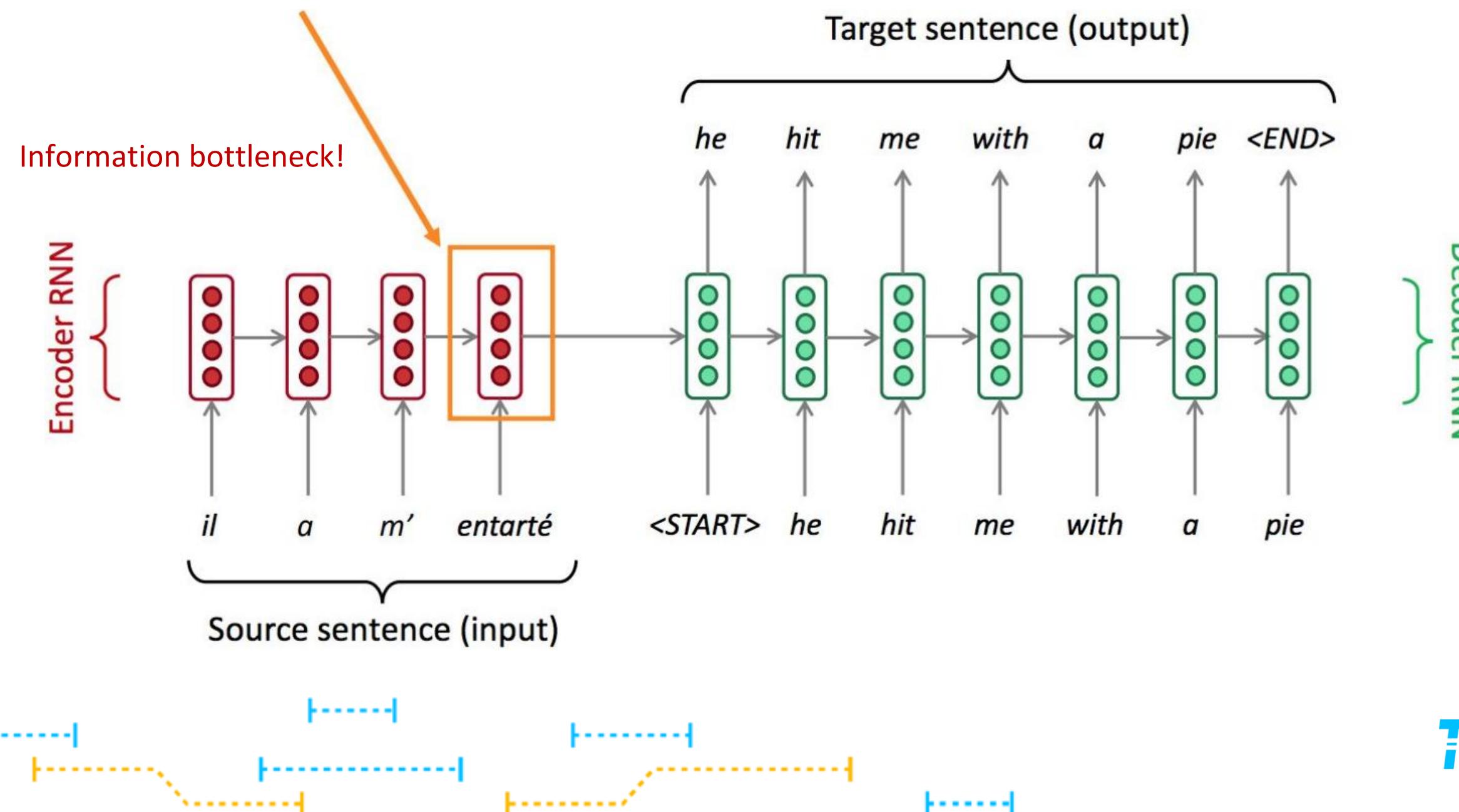
Sesión 3.0

Self-attention

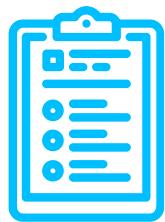
Transformers

Sequence-to-sequence models

Codificación de la oración fuente.
Debe capturar toda la información
sobre la oración fuente.



1.

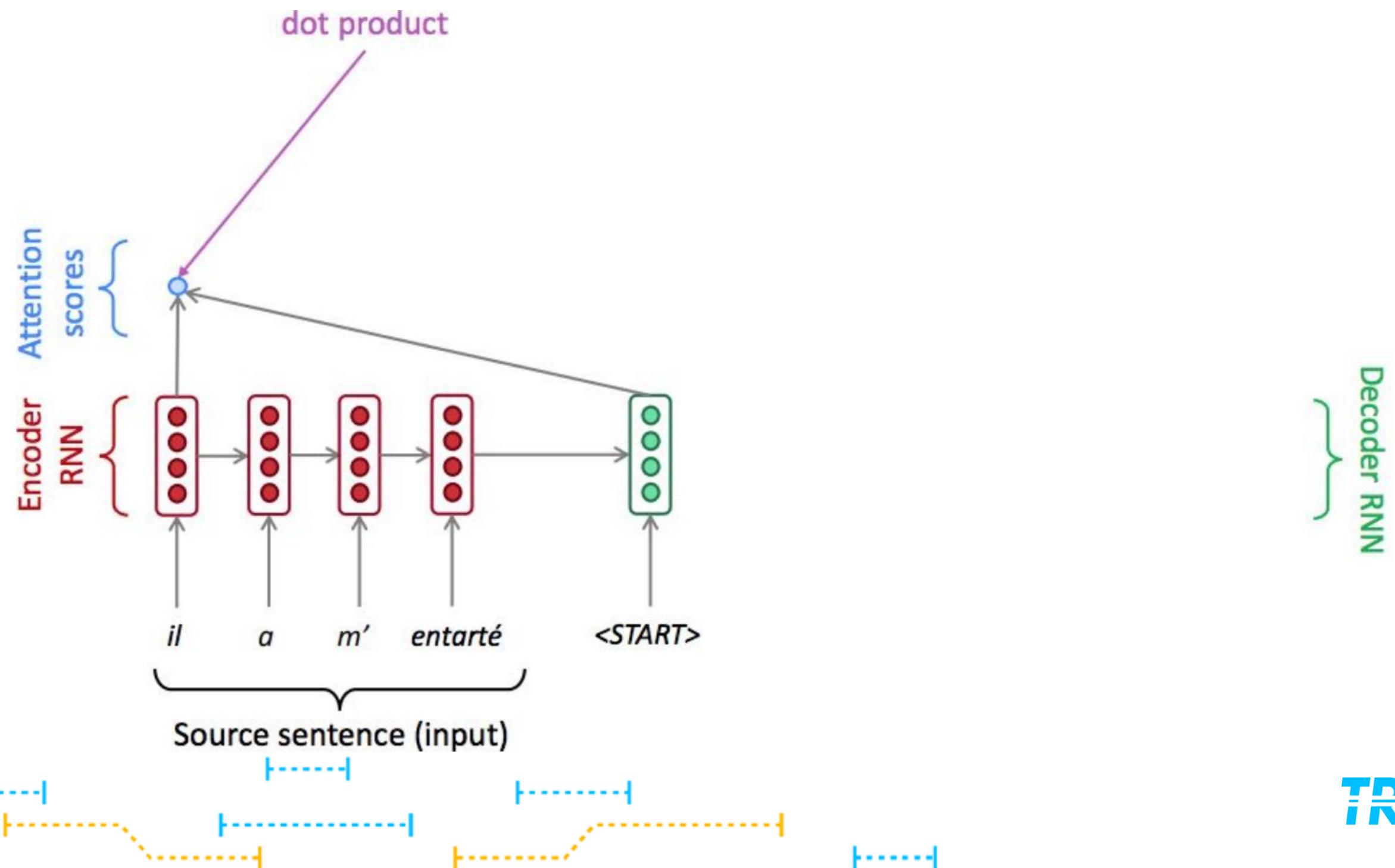


Atten*tion*

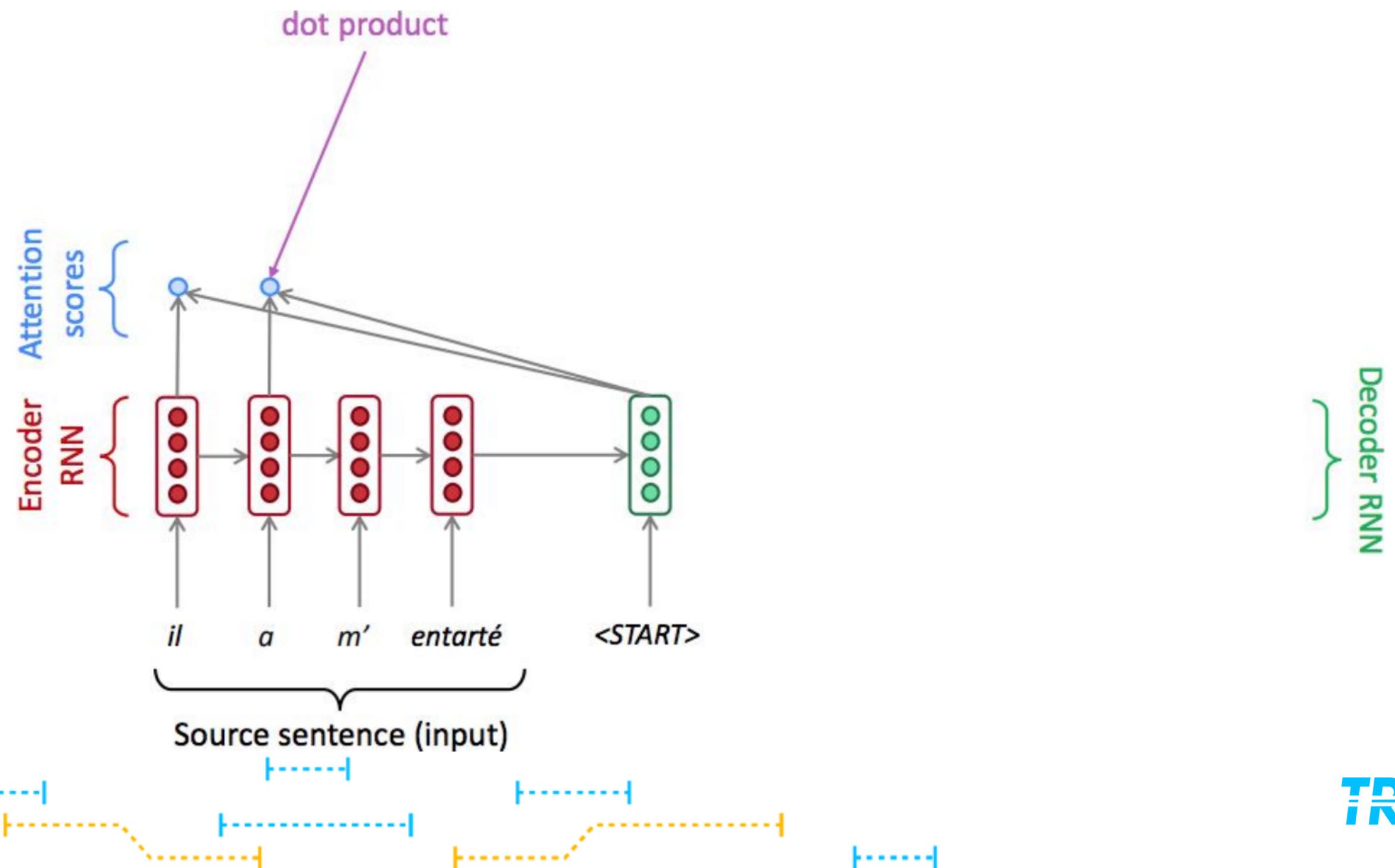
TRANSFORMATEC



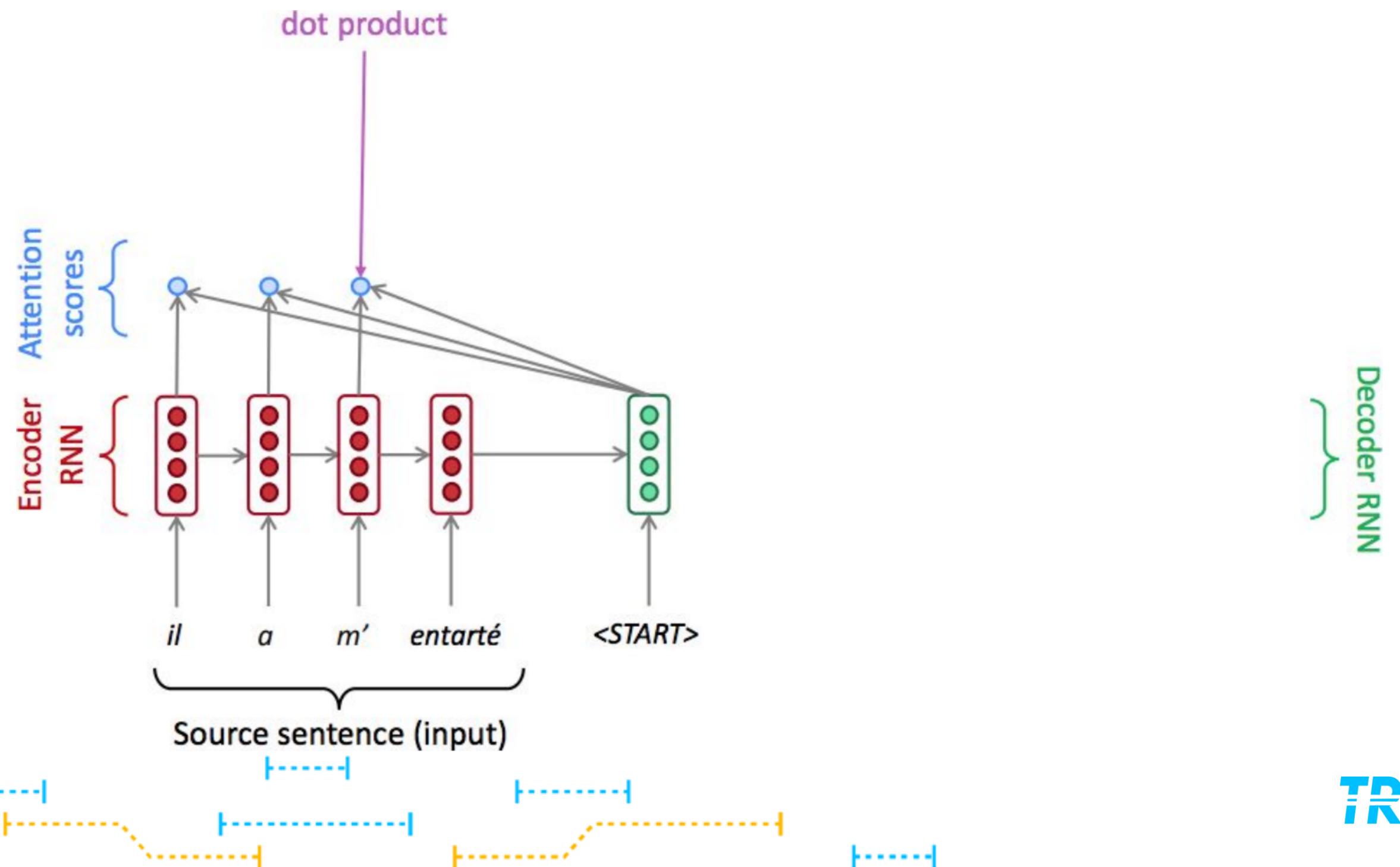
Attention Mechanism



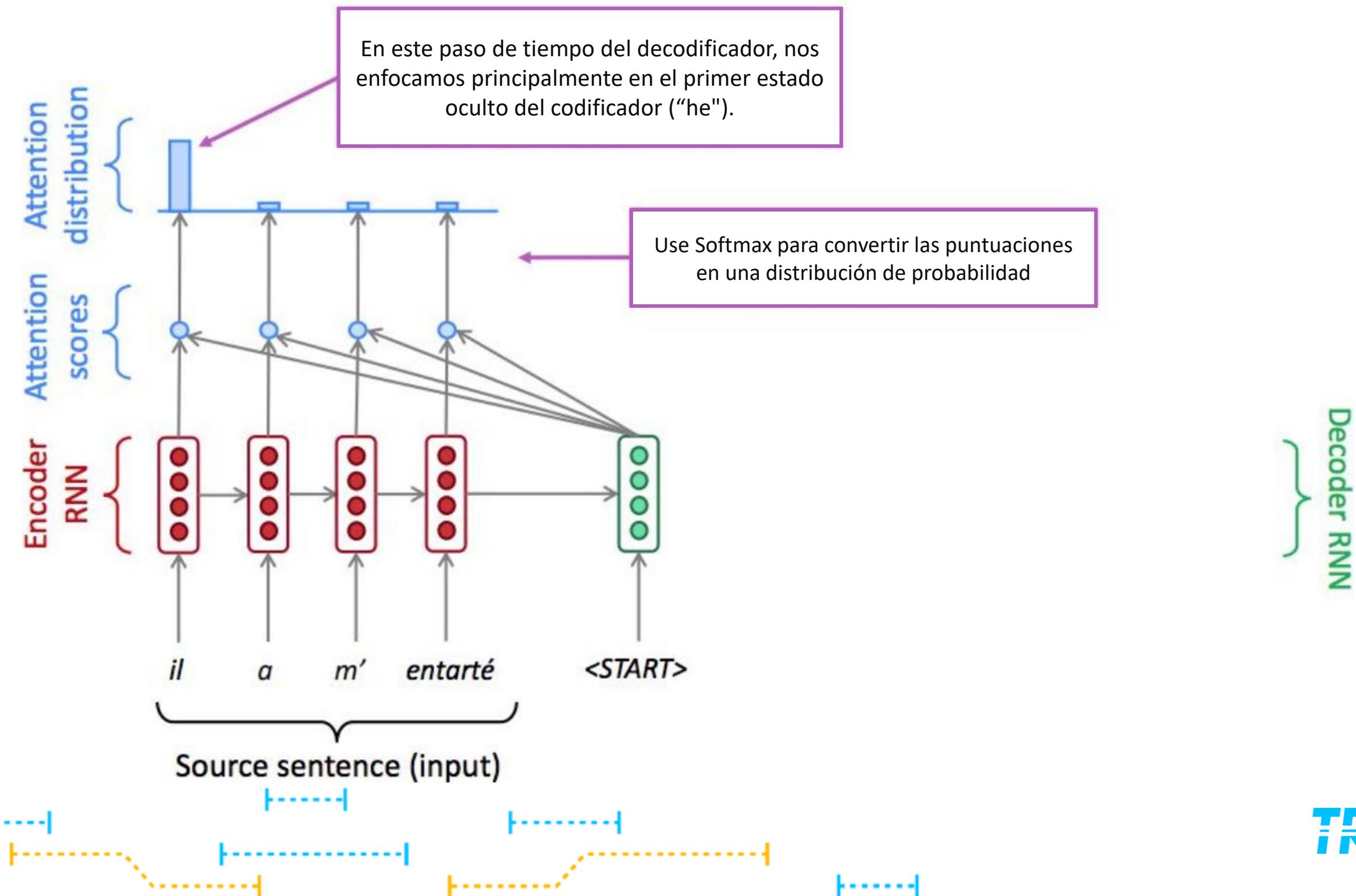
Attention Mechanism



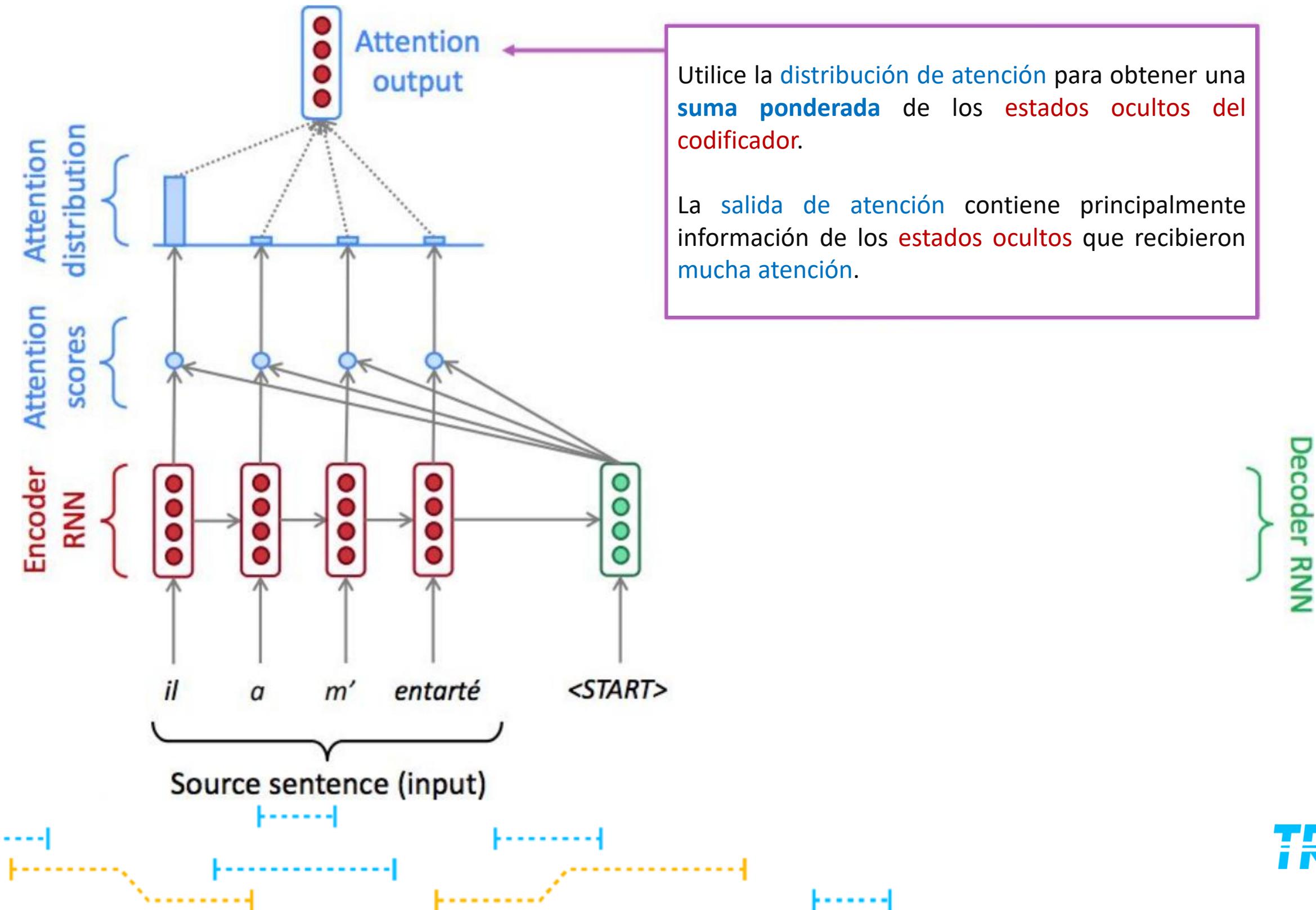
Attention Mechanism



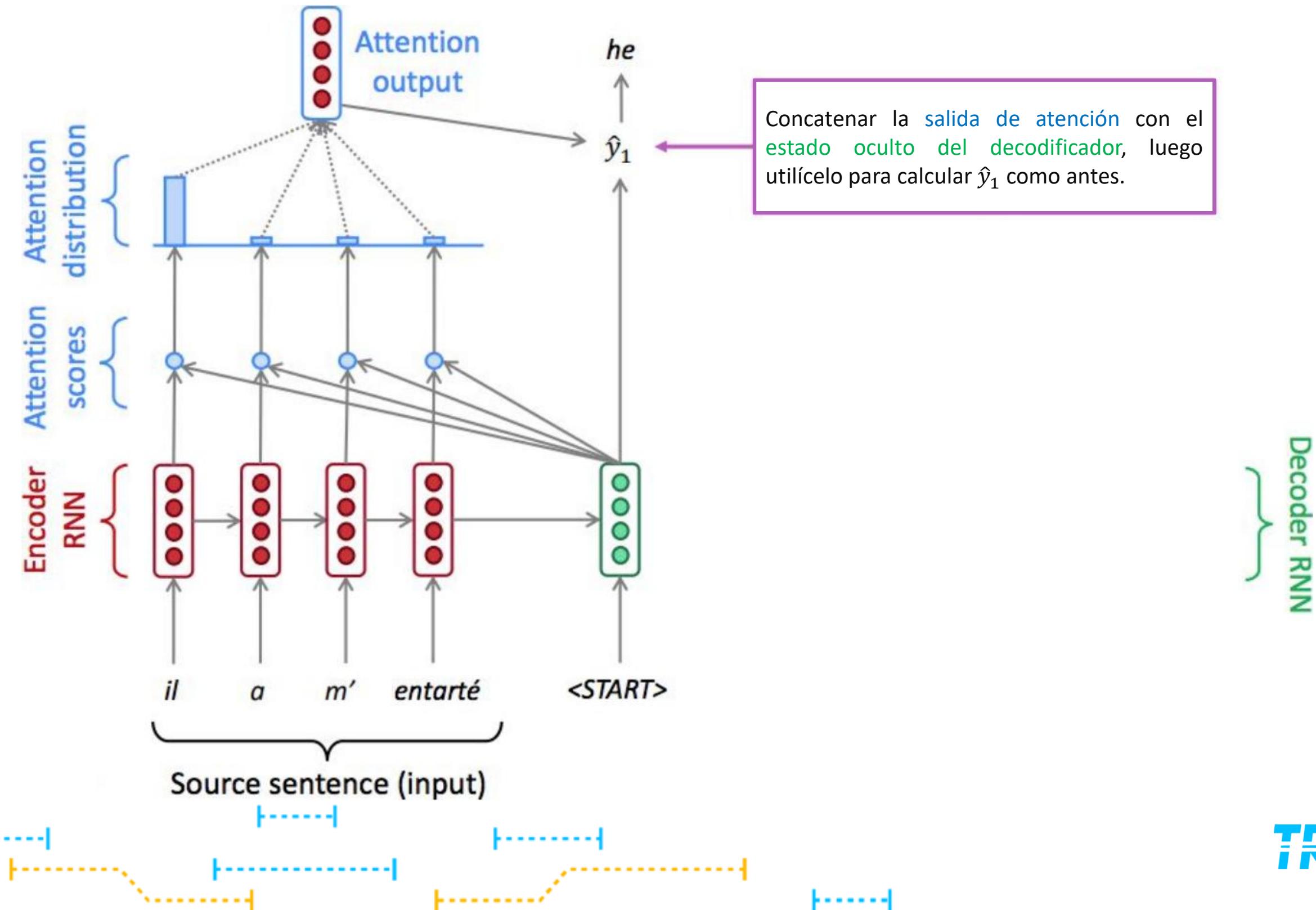
Attention Mechanism



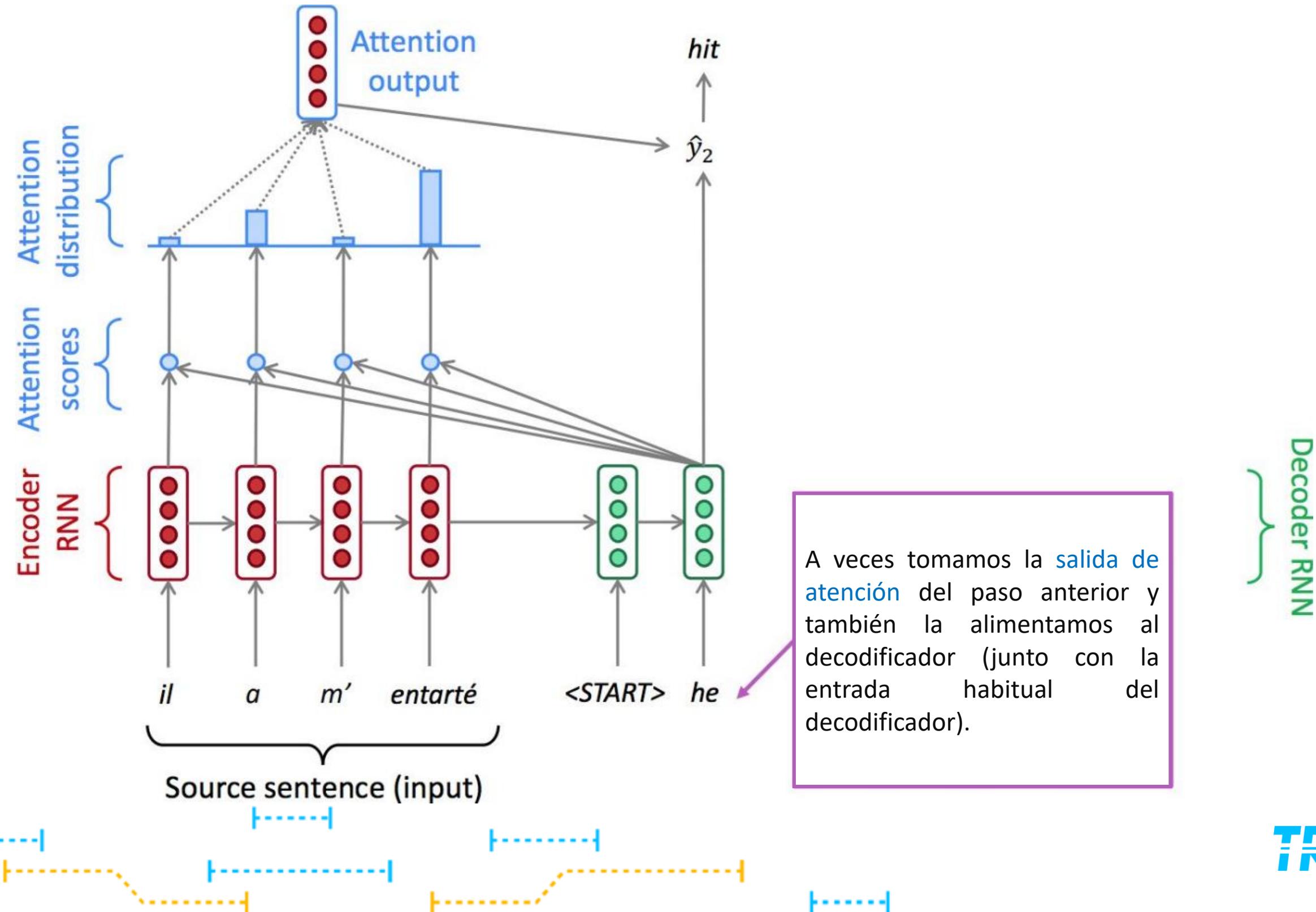
Attention Mechanism



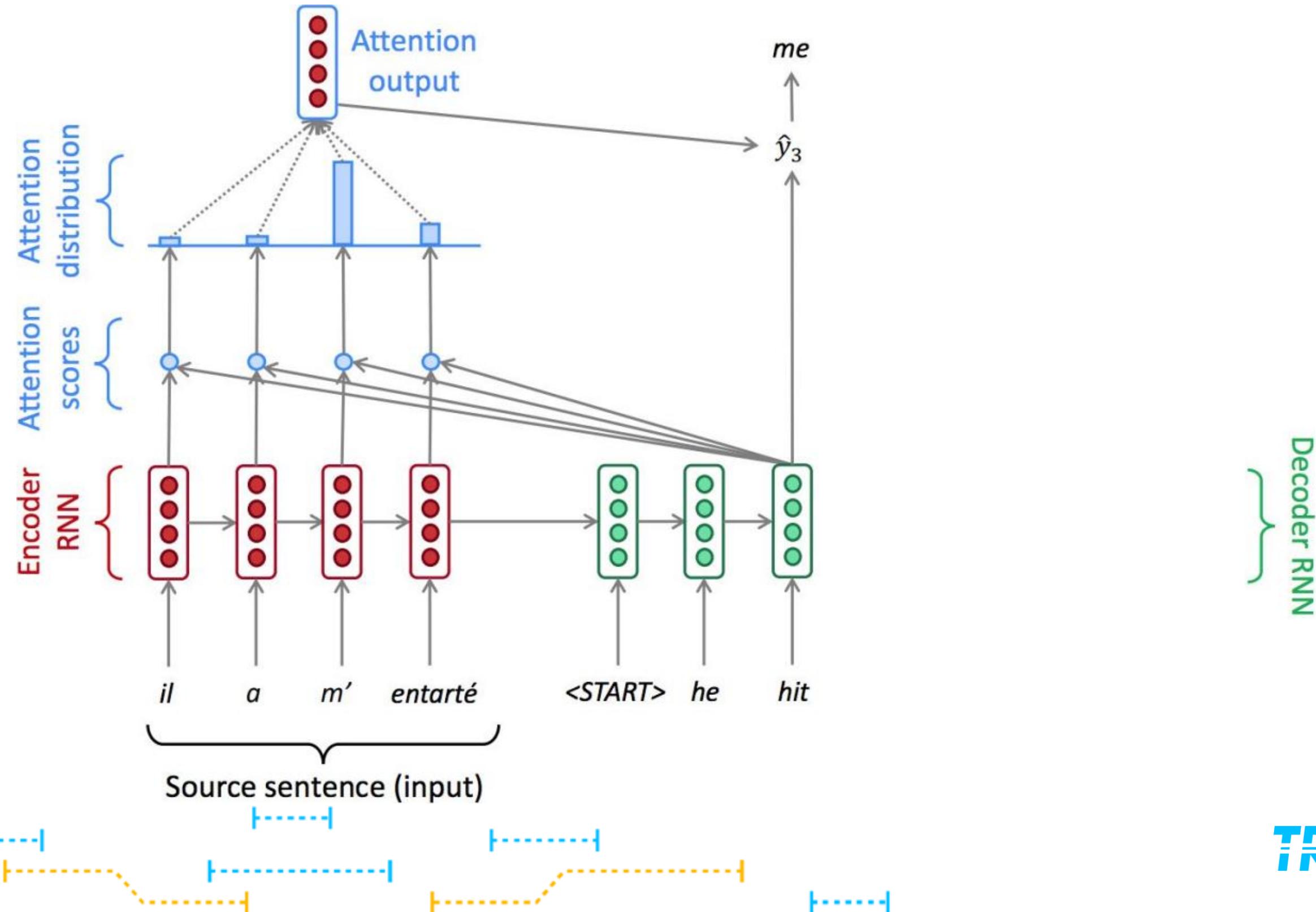
Attention Mechanism



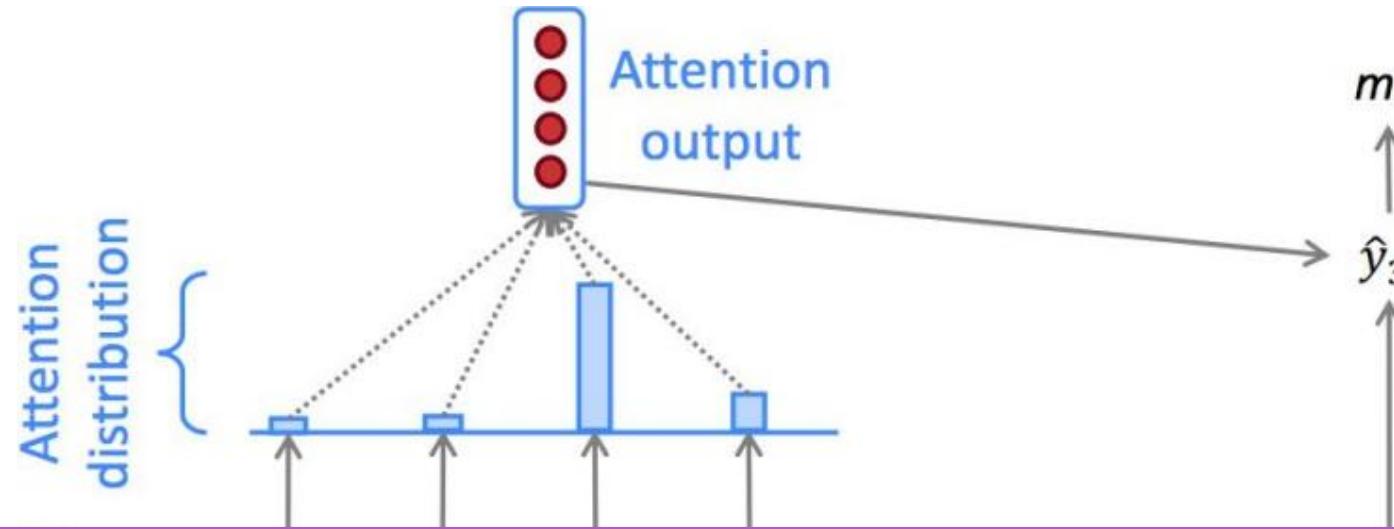
Attention Mechanism



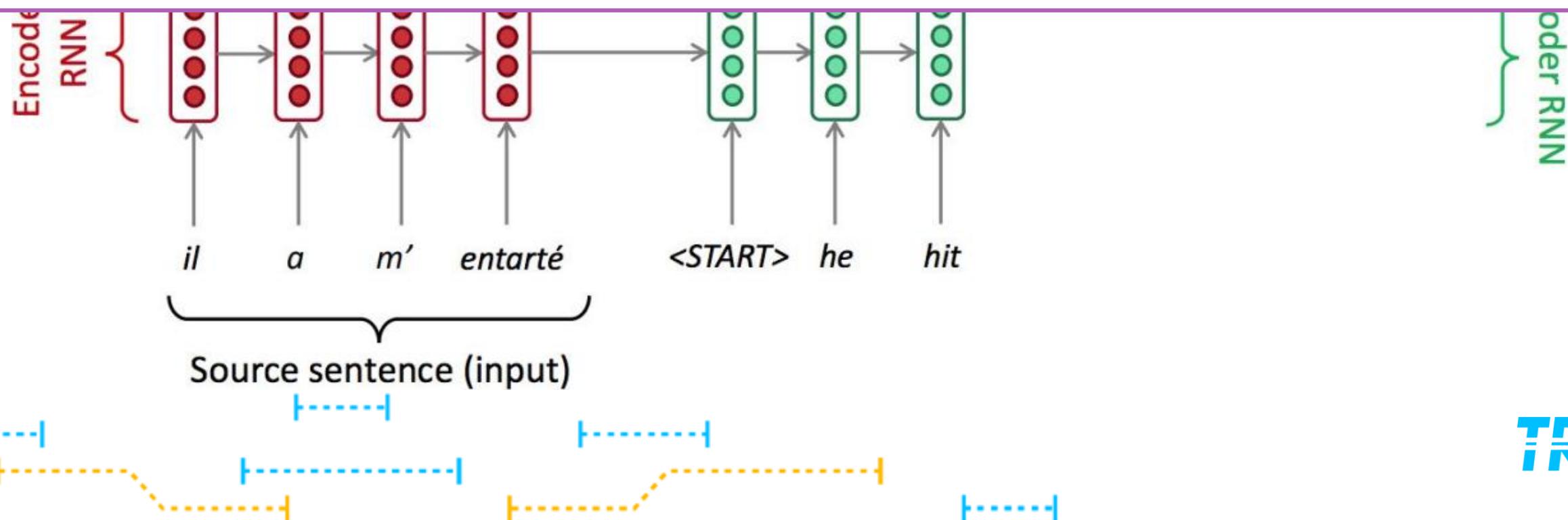
Attention Mechanism



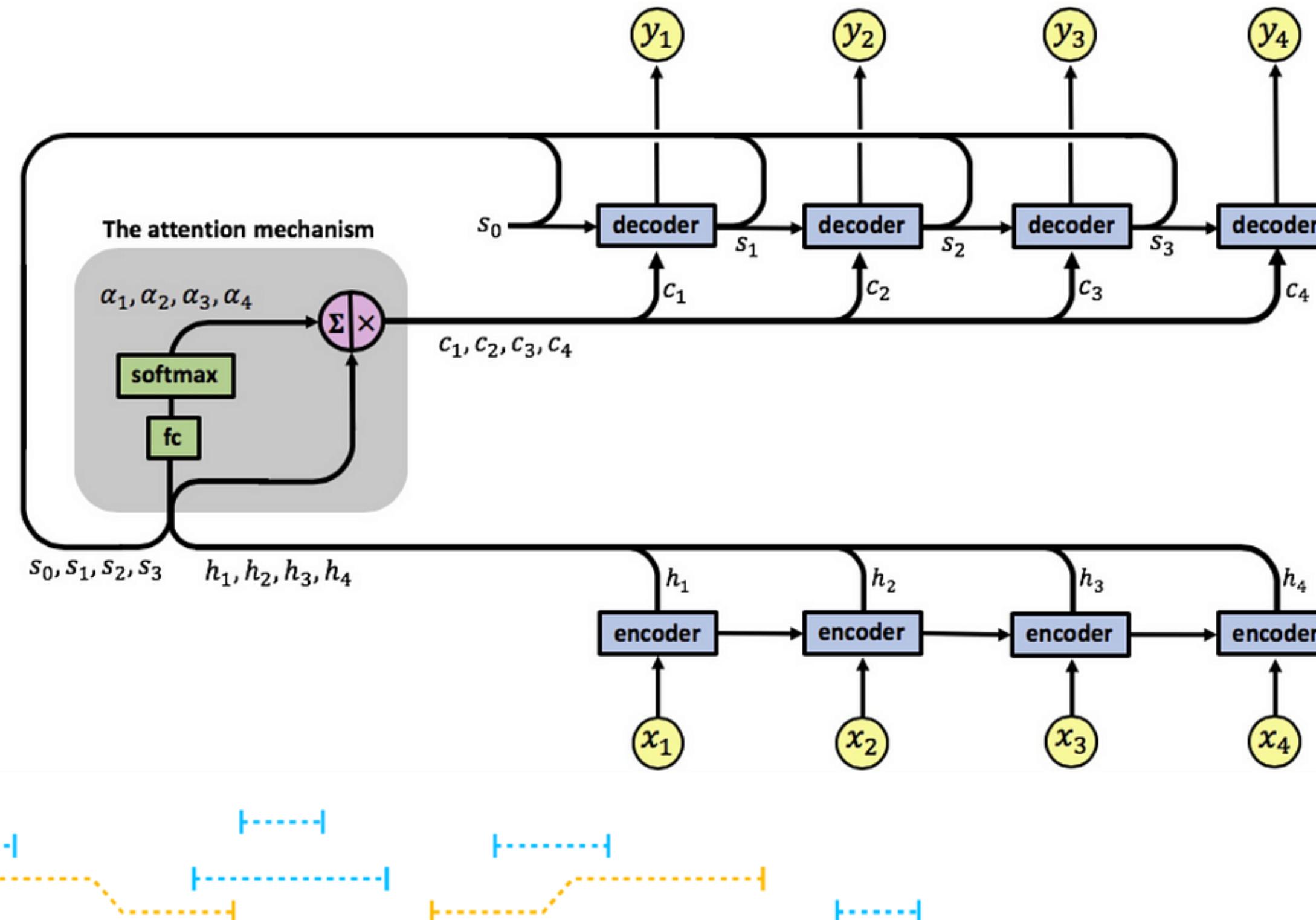
Attention Mechanism



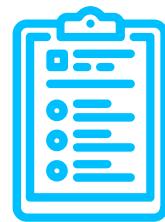
La atención es una forma de obtener una representación de tamaño fijo de un conjunto arbitrario de representaciones (**values**), dependiendo de alguna otra representación (**query**)



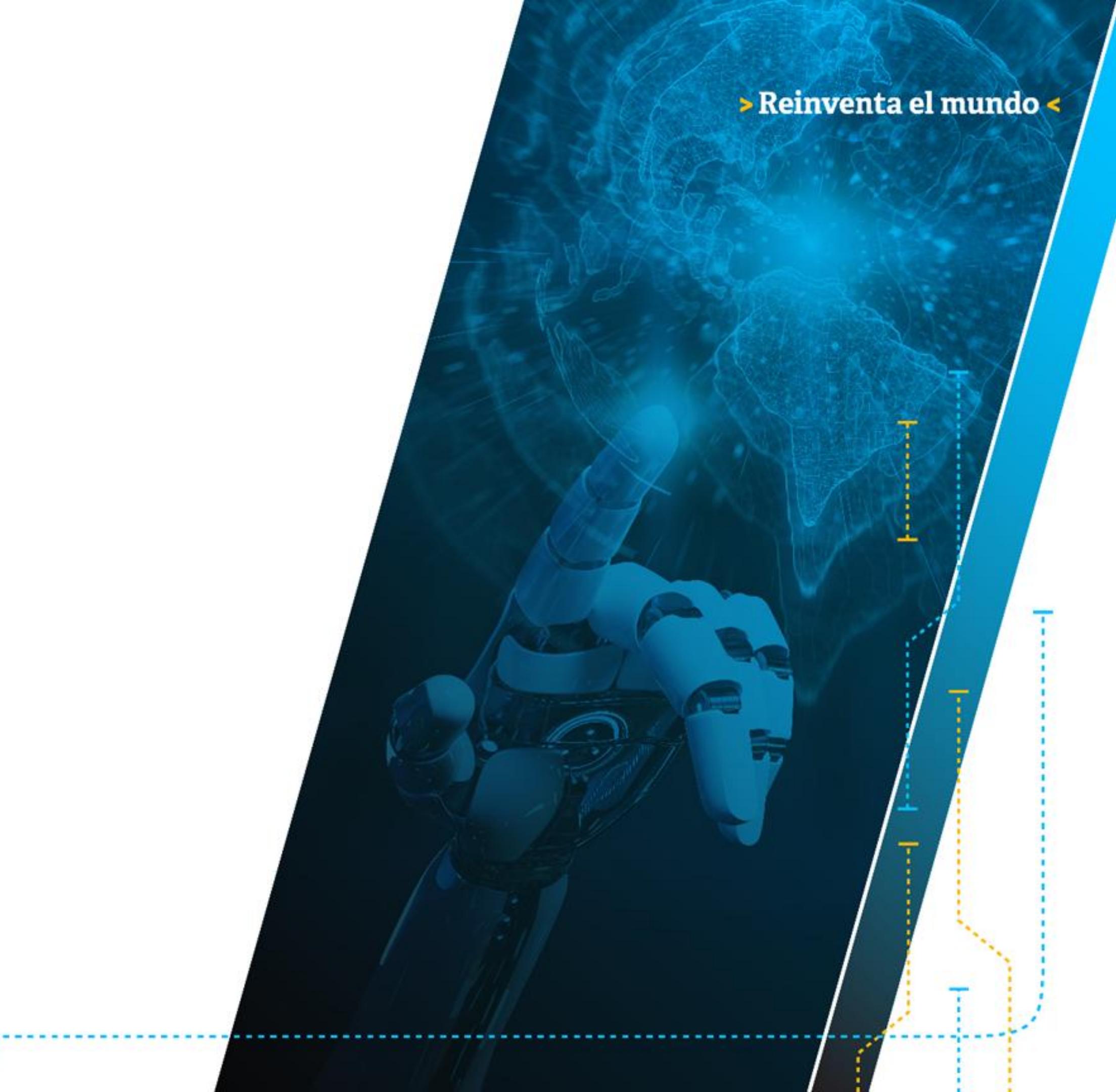
Attention Mechanism



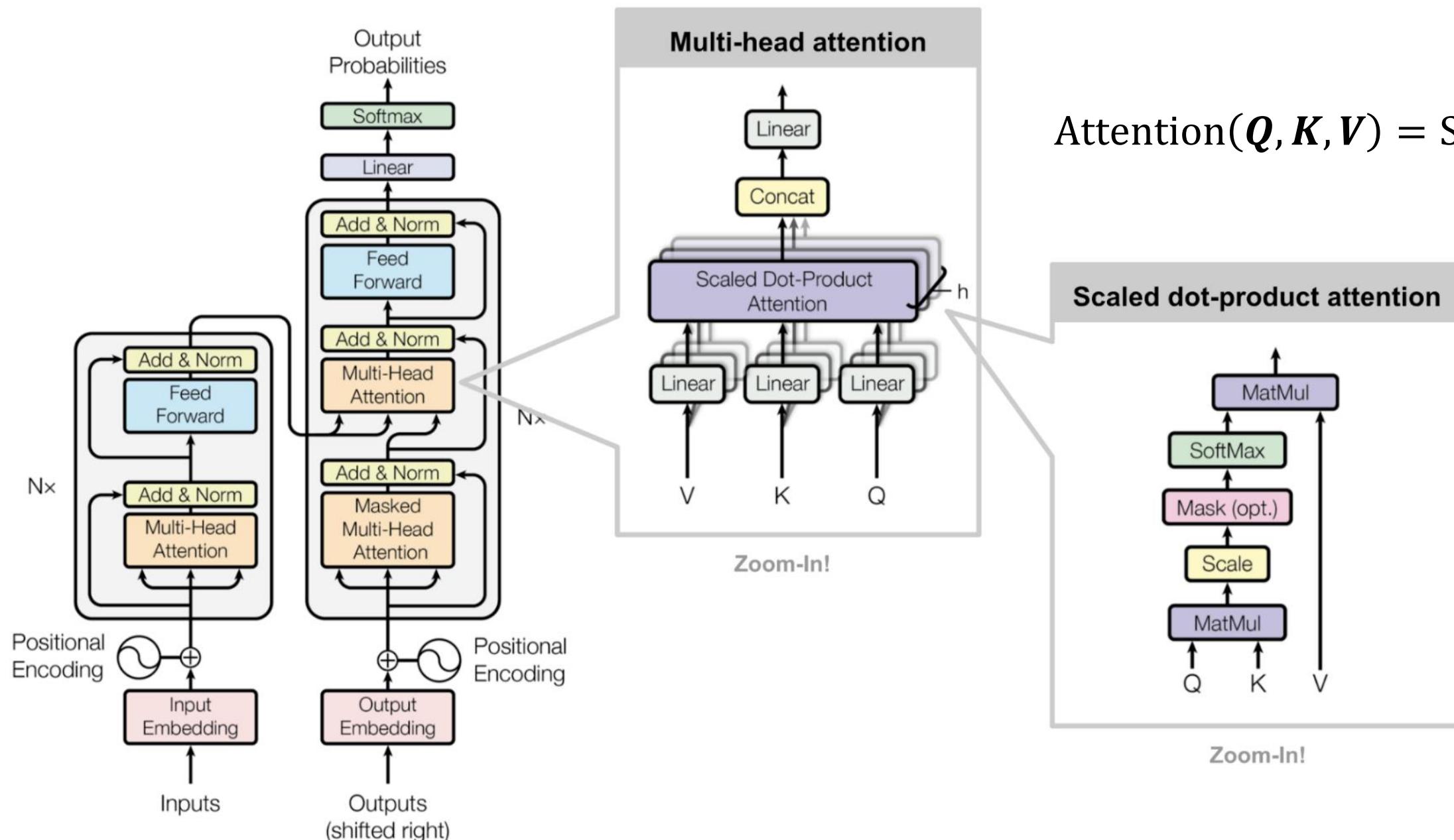
2.



Self-Attention

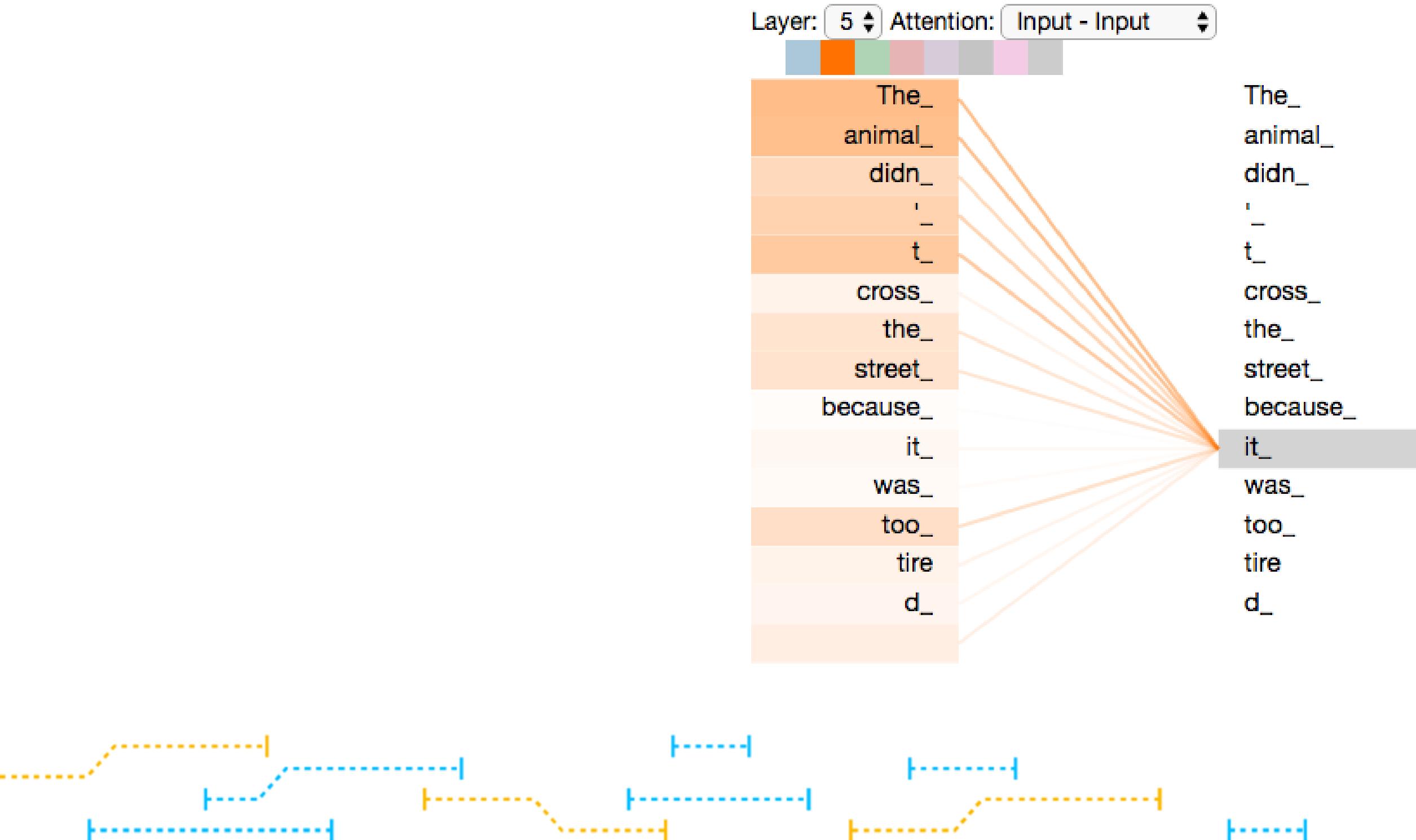


Attention is *all you need*



$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

Attention is *all you need*



Attention is *all you need*

Key (K)

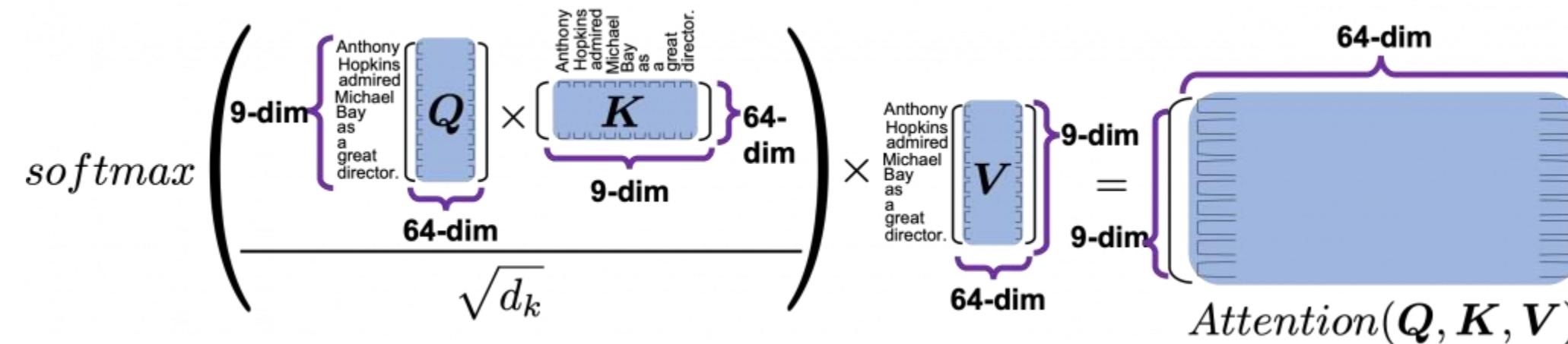
Actúa como un identificador de cada entrada. Es lo que el query usa para decidir qué tan relevante es una entrada en específico.

Query (Q)

Representa lo que esta buscando. Se genera para cada entrada y se contra los keys de los otros elementos de la secuencia.

Value (V)

Contiene la información que será usada en la salida. Una vez que se decida a quién prestar atención (a partir de queries y keys), ponderamos los values para generar la representación final.



Attention is *all you need*

Key (K)	Actúa como un identificador de cada entrada. Es lo que el query usa para decidir qué tan relevante es una entrada en específico.
Query (Q)	Representa lo que esta buscando. Se genera para cada entrada y se contra los keys de los otros elementos de la secuencia.
Value (V)	Contiene la información que será usada en la salida. Una vez que se decida a quién prestar atención (a partir de los queries y keys), ponderamos los values para generar la representación final.

Query (Q):

¿Qué acción está asociada con 'gato'?

El **gato** persigue al ratón

Key (K):

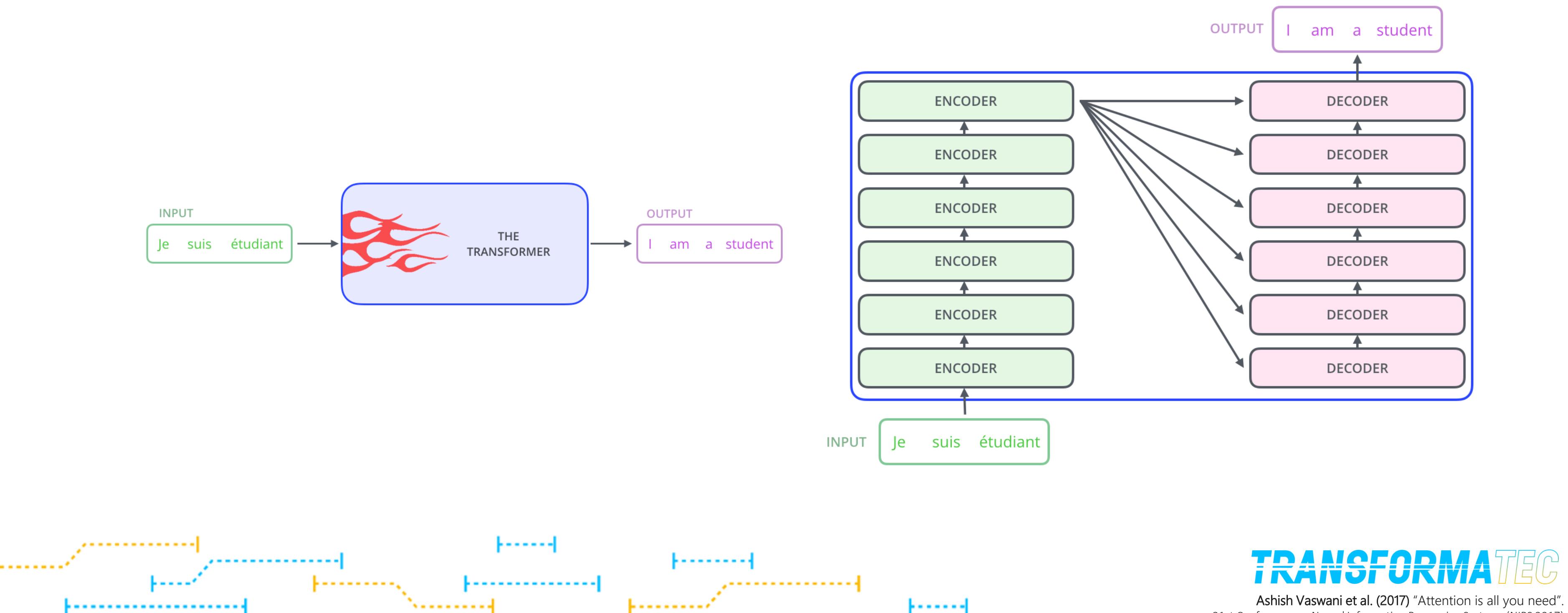
Es un verbo de acción

Value (V):

Información semántica de la palabra



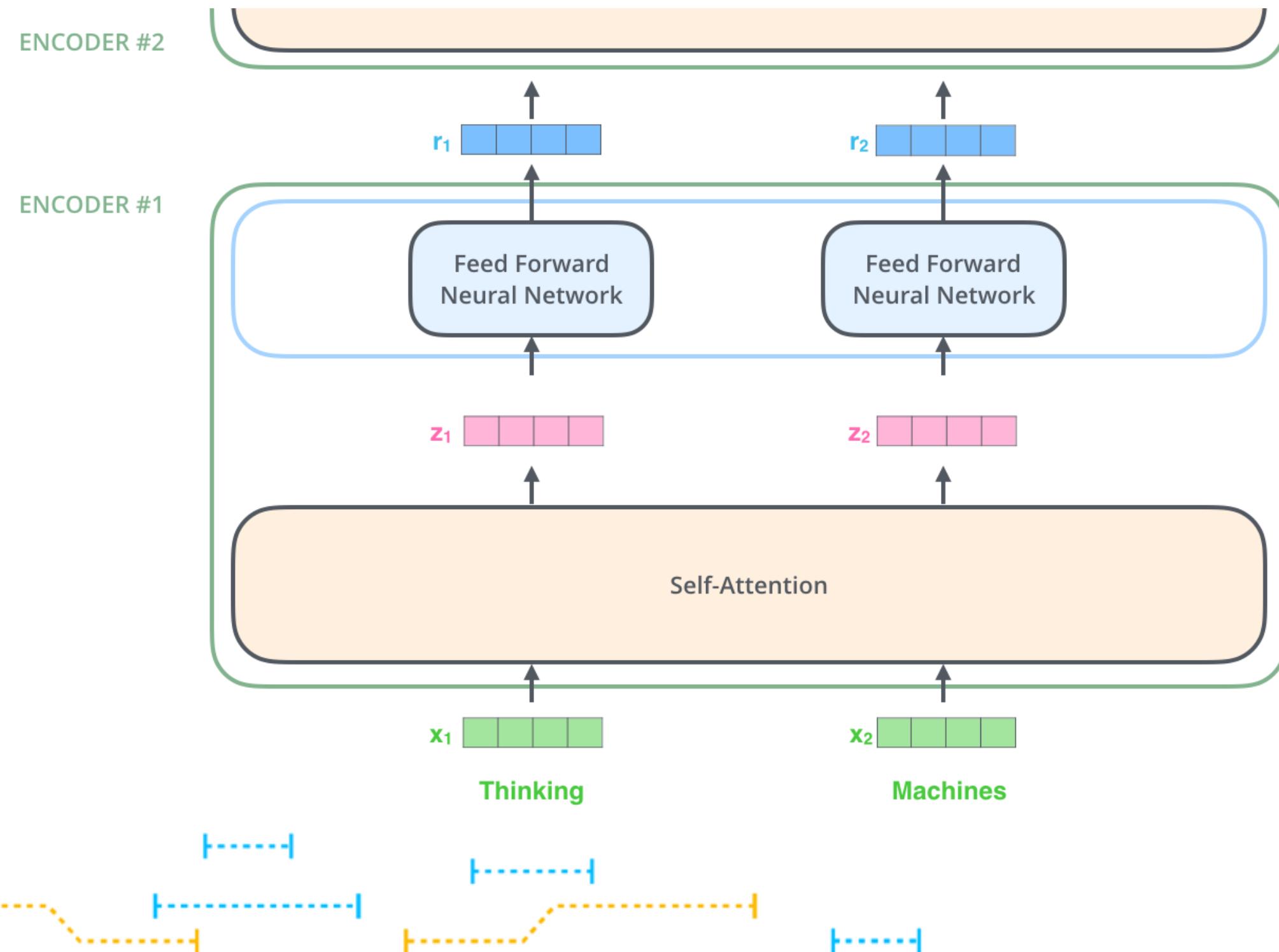
Attention is *all you need*



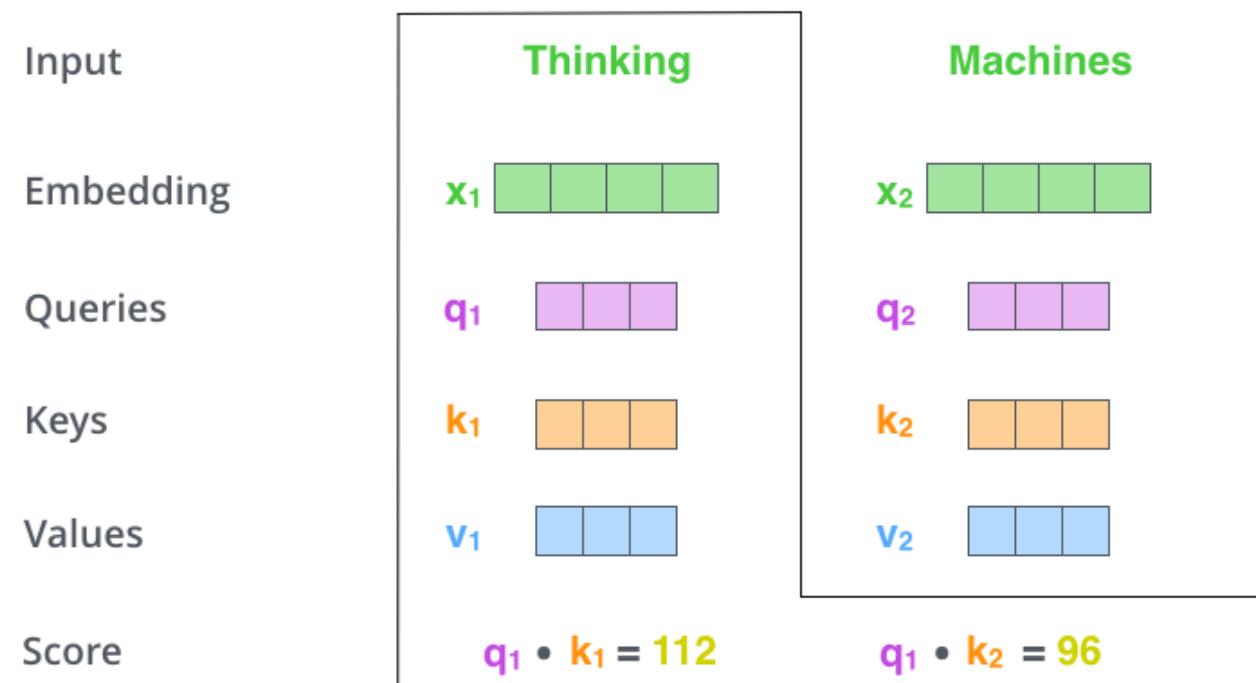
TRANSFORMATEC

Ashish Vaswani et al. (2017) "Attention is all you need".
31st Conference on Neural Information Processing Systems (NIPS 2017)

Attention is *all you need*



Attention is *all you need*

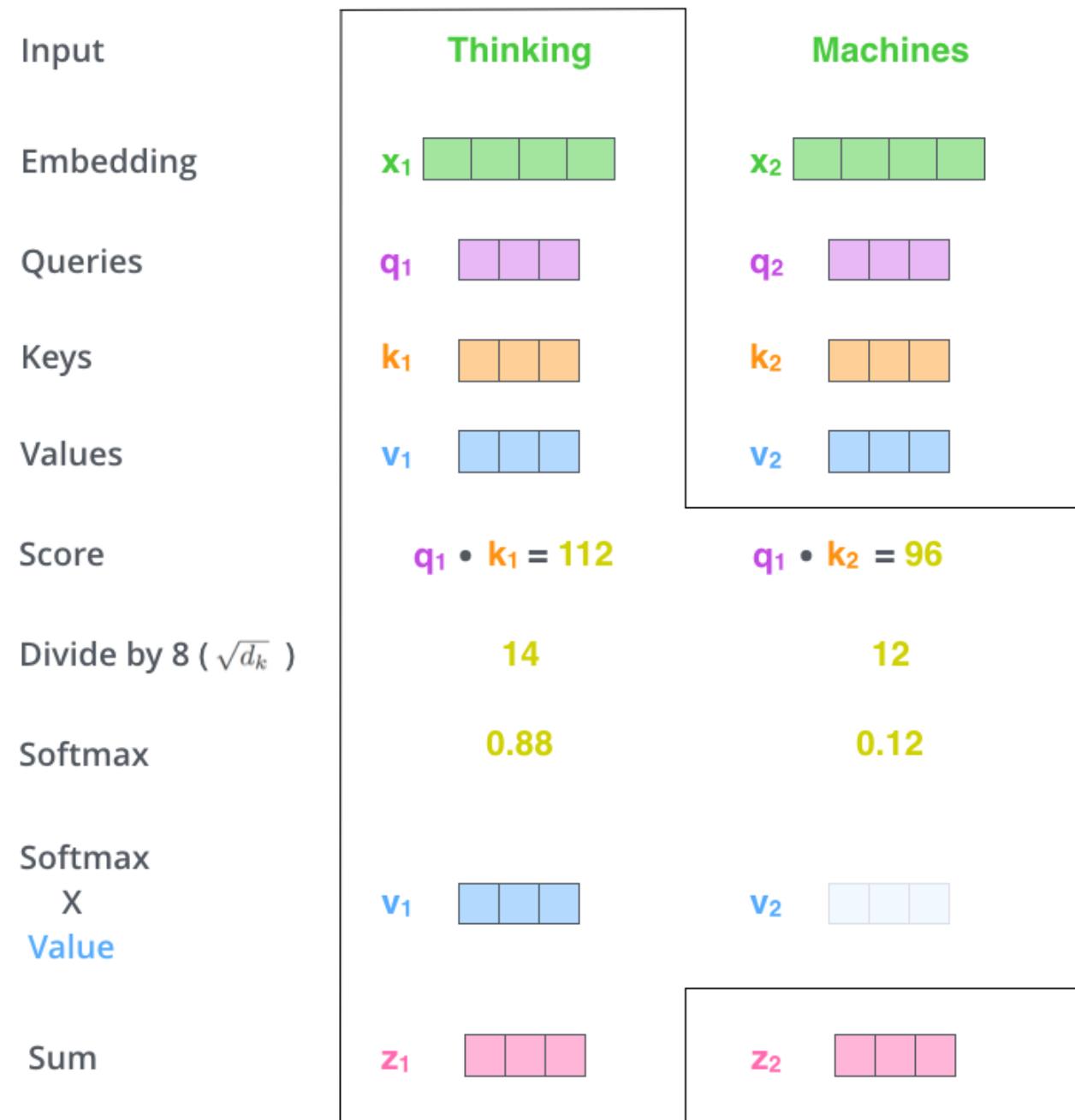


Attention is *all you need*

Input	Thinking		Machines	
Embedding	x_1	[4 green boxes]	x_2	[4 green boxes]
Queries	q_1	[3 purple boxes]	q_2	[3 purple boxes]
Keys	k_1	[3 orange boxes]	k_2	[3 orange boxes]
Values	v_1	[3 blue boxes]	v_2	[3 blue boxes]
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by 8 ($\sqrt{d_k}$)	14		12	
Softmax	0.88		0.12	

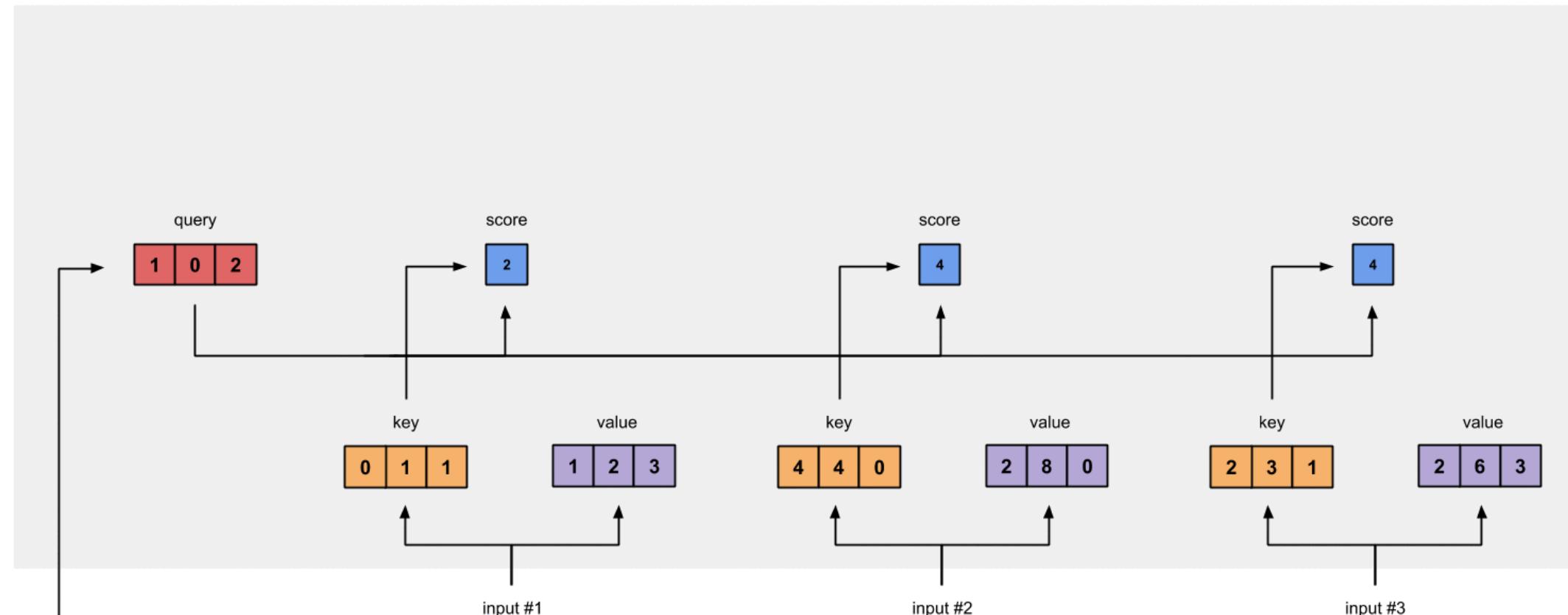


Attention is *all you need*



Attention is *all you need*

Self-attention

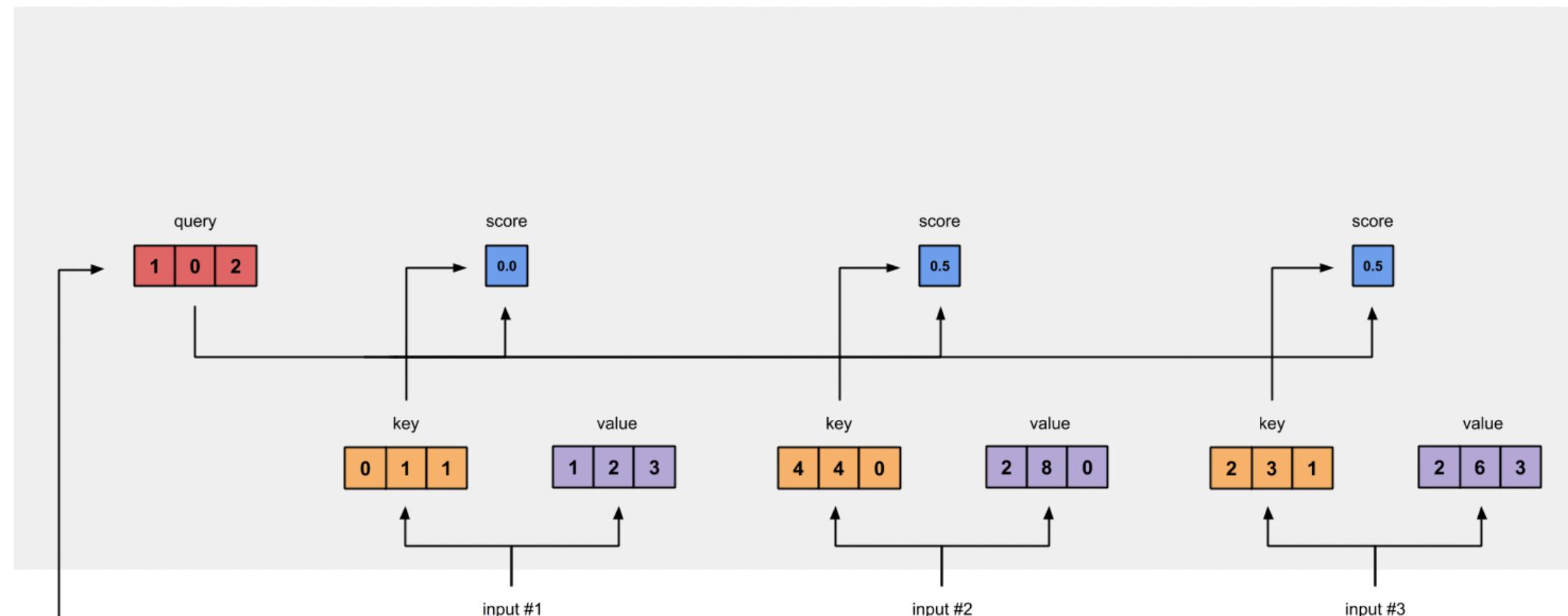


TRANSFORMATEC

Ashish Vaswani et al. (2017) "Attention is all you need".
31st Conference on Neural Information Processing Systems (NIPS 2017)

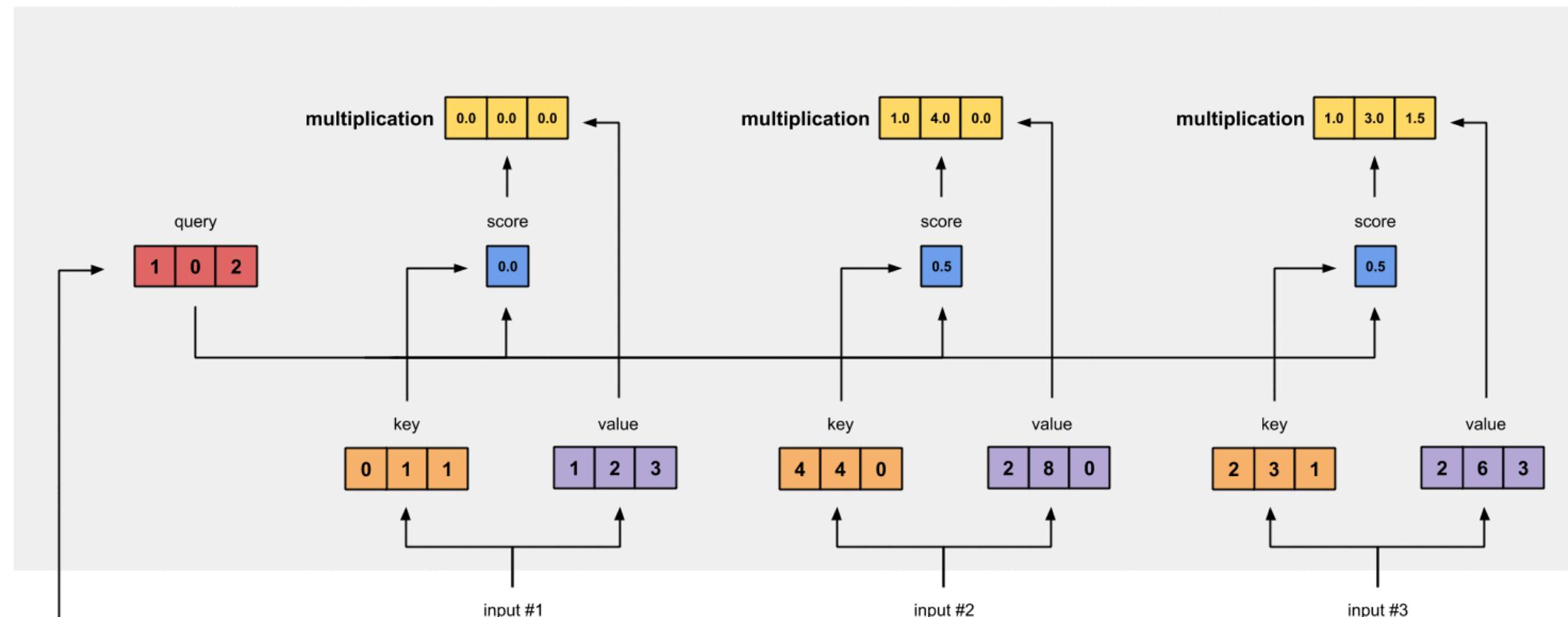
Attention is *all you need*

Self-attention



Attention is *all you need*

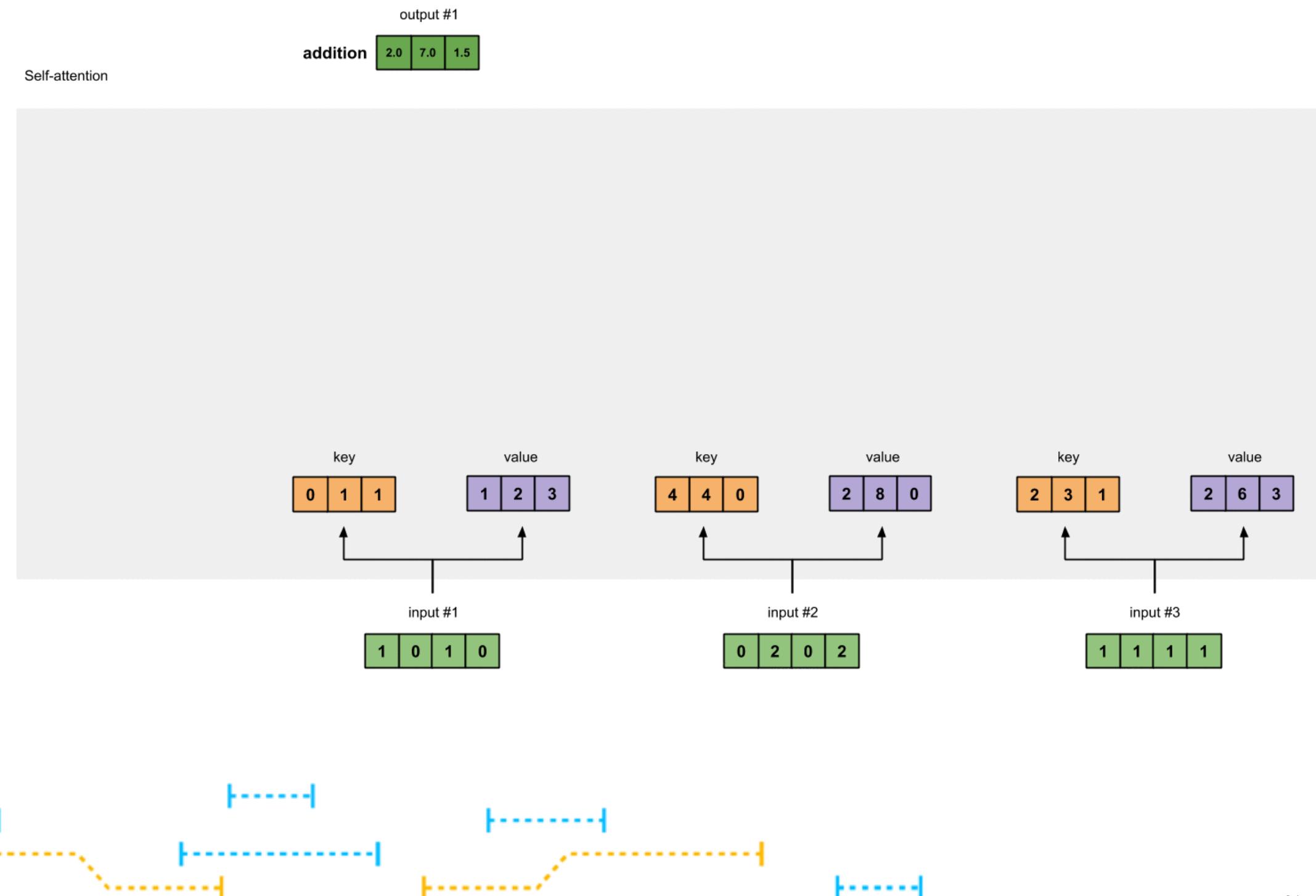
Self-attention



TRANSFORMATEC

Ashish Vaswani et al. (2017) "Attention is all you need".
31st Conference on Neural Information Processing Systems (NIPS 2017)

Attention is *all you need*



Attention is *all you need*

$$\begin{matrix} \mathbf{X} \\ \begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} \end{matrix} \times \begin{matrix} \mathbf{W}^Q \\ \begin{matrix} \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \end{matrix} \end{matrix} = \begin{matrix} \mathbf{Q} \\ \begin{matrix} \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \end{matrix} \end{matrix}$$

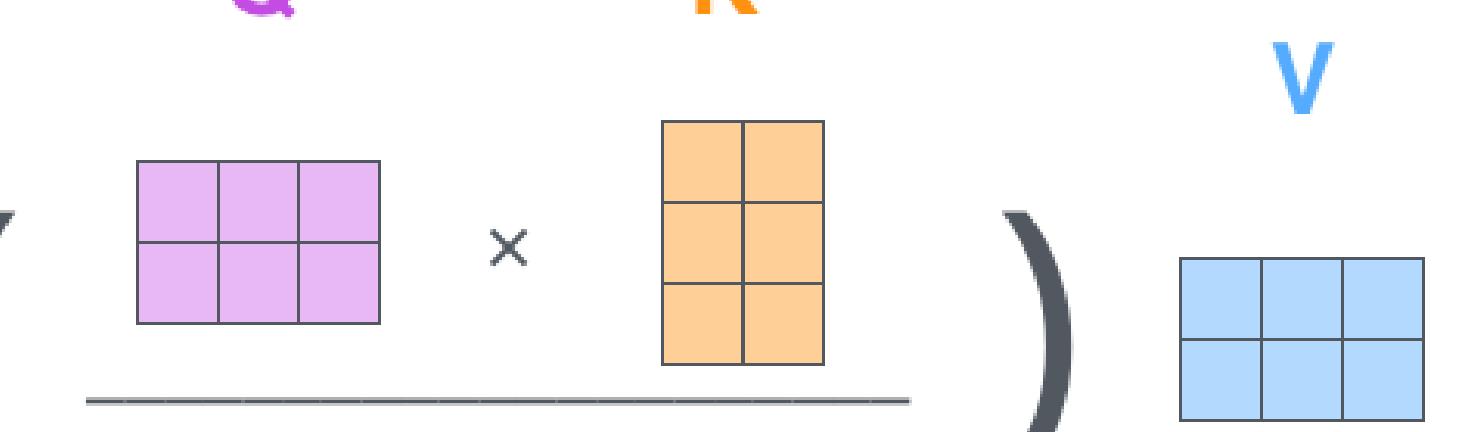
$$\begin{matrix} \mathbf{X} \\ \begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} \end{matrix} \times \begin{matrix} \mathbf{W}^K \\ \begin{matrix} \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} \\ \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} \\ \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} \end{matrix} \end{matrix} = \begin{matrix} \mathbf{K} \\ \begin{matrix} \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} \\ \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} \\ \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} & \textcolor{lightorange}{\square} \end{matrix} \end{matrix}$$

$$\begin{matrix} \mathbf{X} \\ \begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} \end{matrix} \times \begin{matrix} \mathbf{W}^V \\ \begin{matrix} \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \end{matrix} \end{matrix} = \begin{matrix} \mathbf{V} \\ \begin{matrix} \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \\ \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} & \textcolor{lightblue}{\square} \end{matrix} \end{matrix}$$

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V}$$

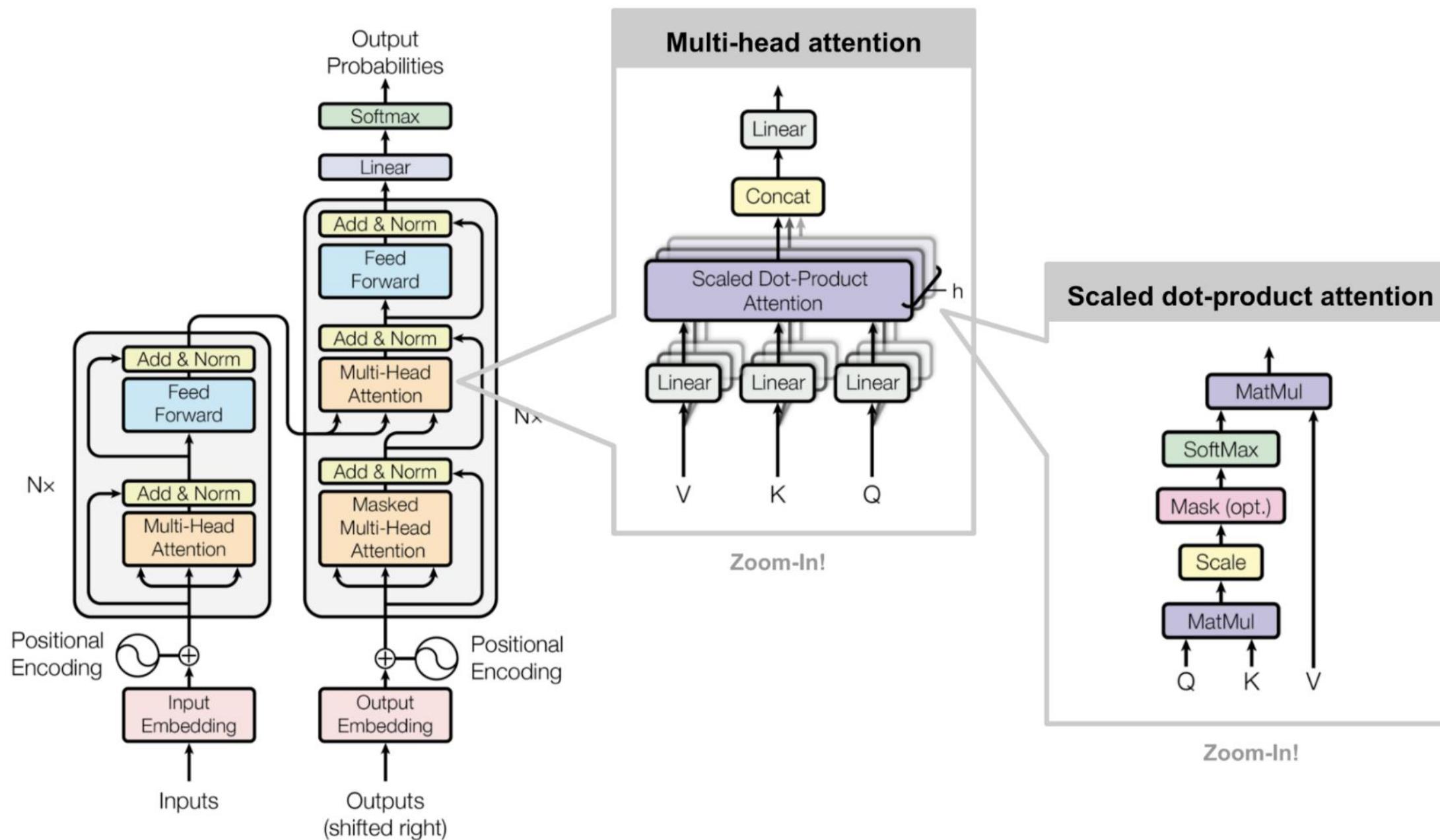


Attention is *all you need*

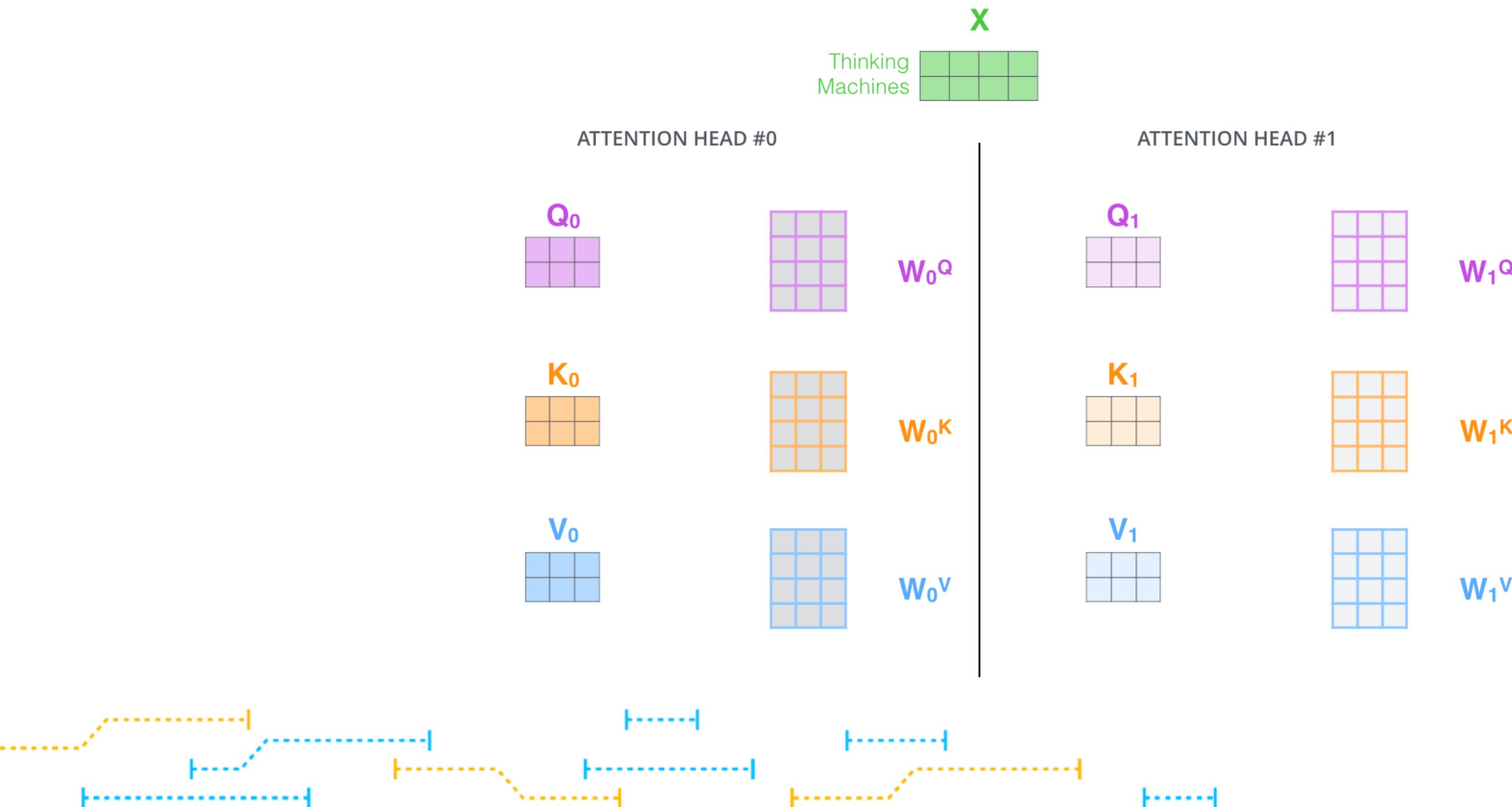
$$\text{softmax} \left(\frac{\begin{matrix} \mathbf{Q} \\ \times \\ \mathbf{K}^T \end{matrix}}{\sqrt{d_k}} \right) \mathbf{V} = \mathbf{z}$$




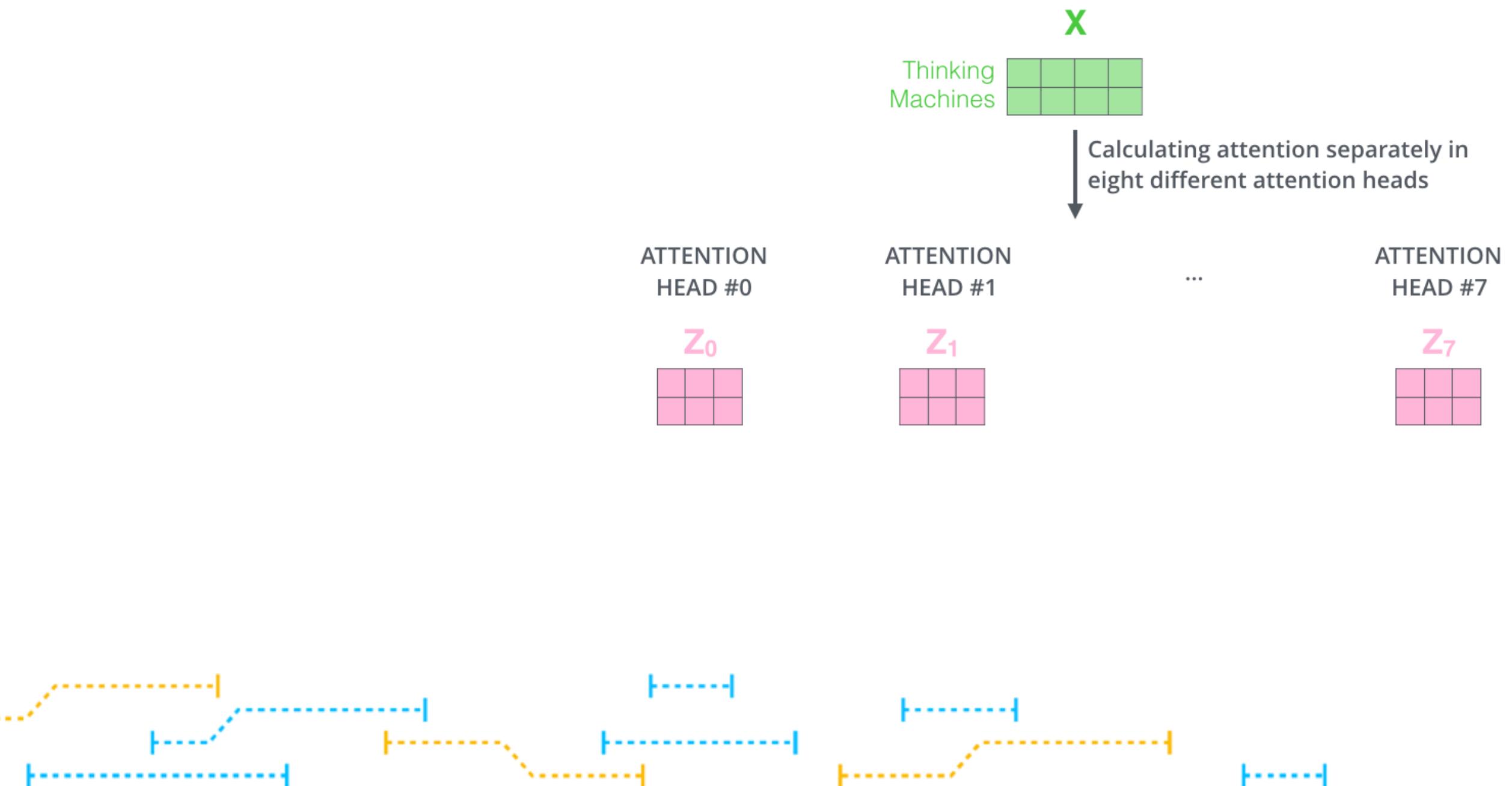
Attention is *all you need*



Attention is *all you need*



Attention is *all you need*



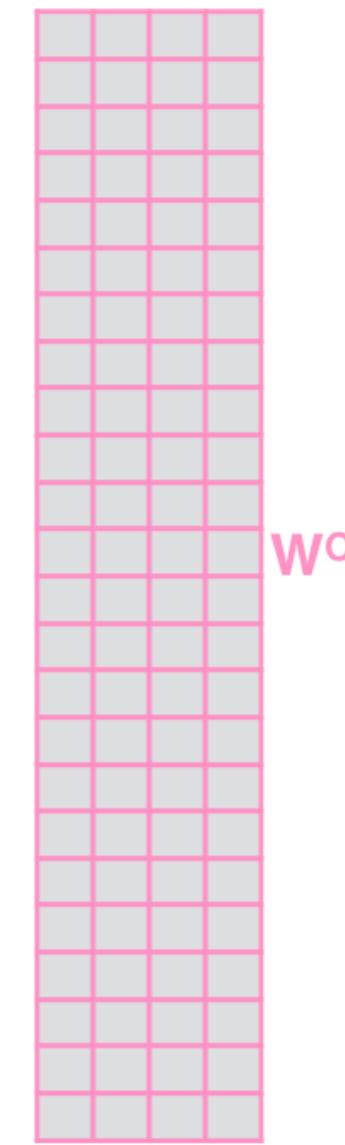
Attention is *all you need*

1) Concatenate all the attention heads



2) Multiply with a weight matrix W^o that was trained jointly with the model

X

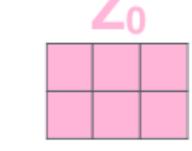


3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN

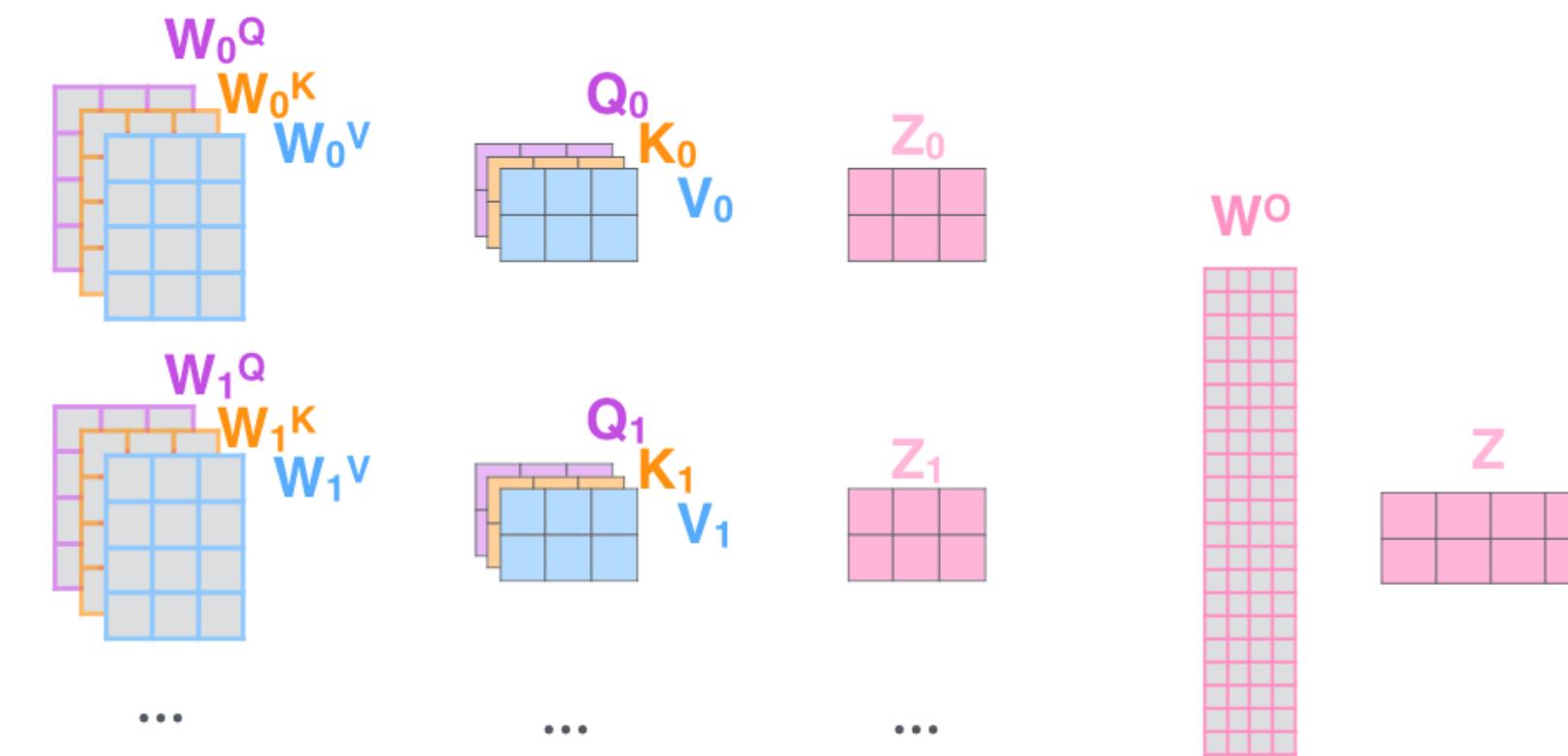
$$= \begin{matrix} Z \\ \vdots \end{matrix}$$



Attention is *all you need*

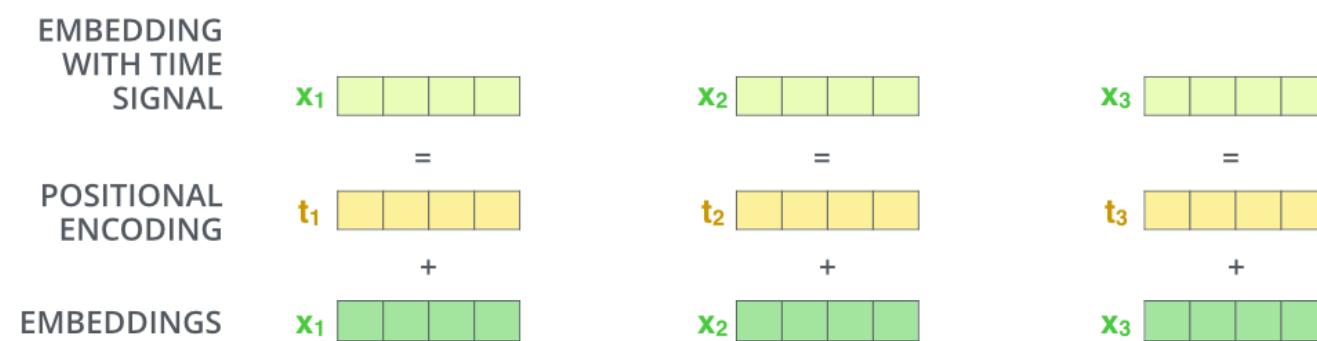
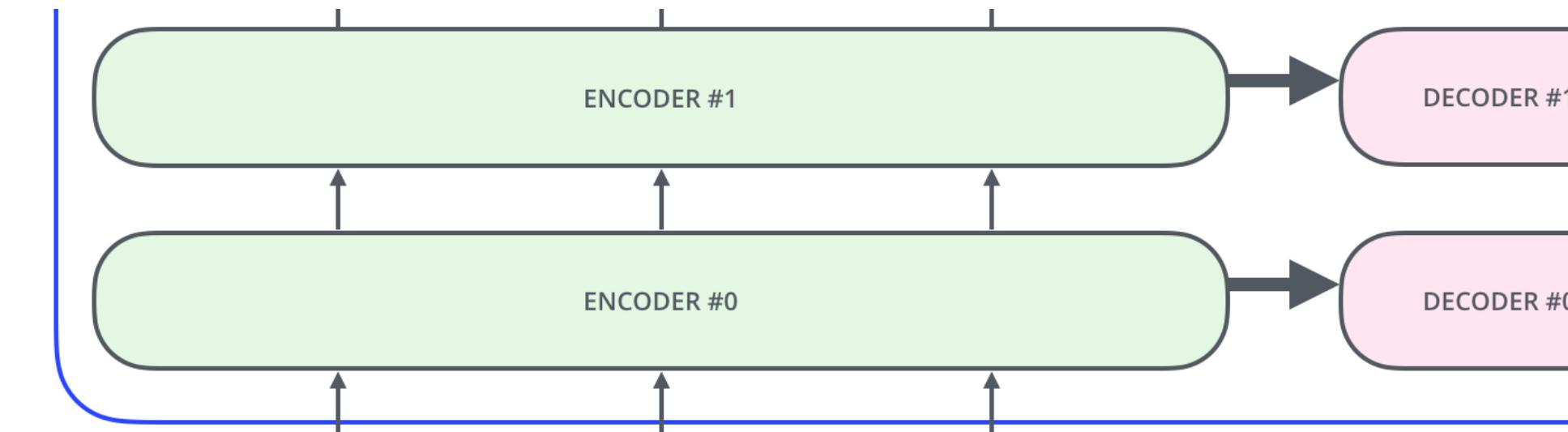
- 1) This is our input sentence* 
- 2) We embed each word* 
- 3) Split into 8 heads. We multiply  with weight matrices W_0^Q , W_0^K , W_0^V , W_1^Q , W_1^K , W_1^V , ..., W_7^Q , W_7^K , W_7^V
- 4) Calculate attention using the resulting $Q/K/V$ matrices Q_0, K_0, V_0 , Q_1, K_1, V_1 , ..., Q_7, K_7, V_7
- 5) Concatenate the resulting  matrices, then multiply with weight matrix W^O to produce the output of the layer

* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



Attention is *all you need*

Positional encoding

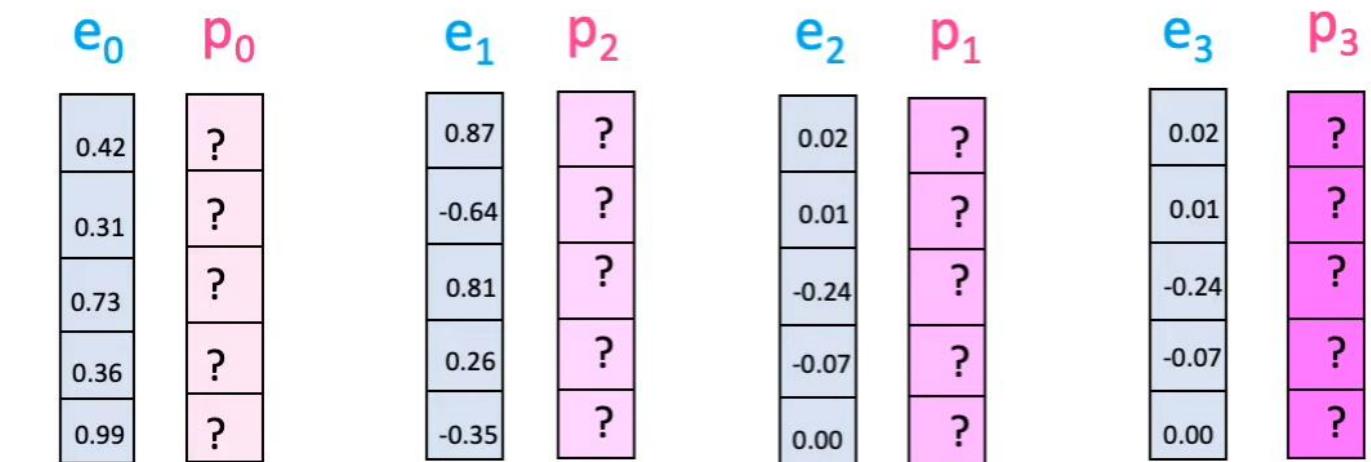
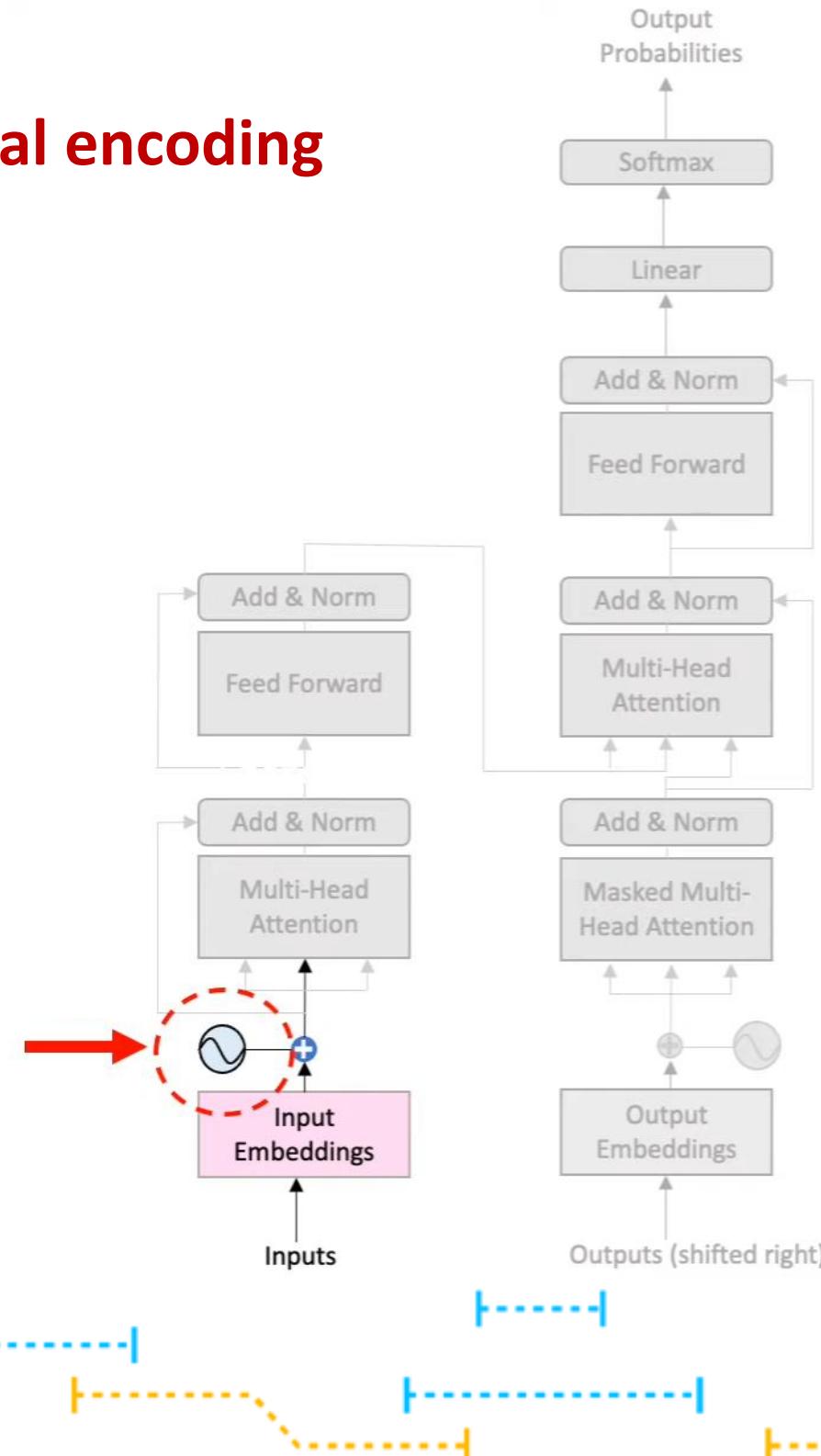


INPUT Je suis étudiant



Attention is *all you need*

Positional encoding



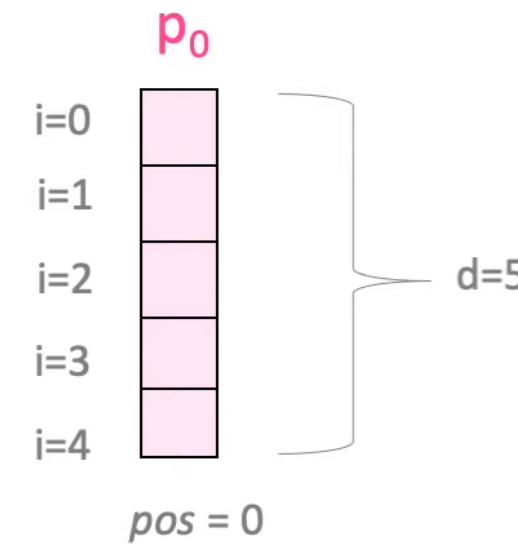
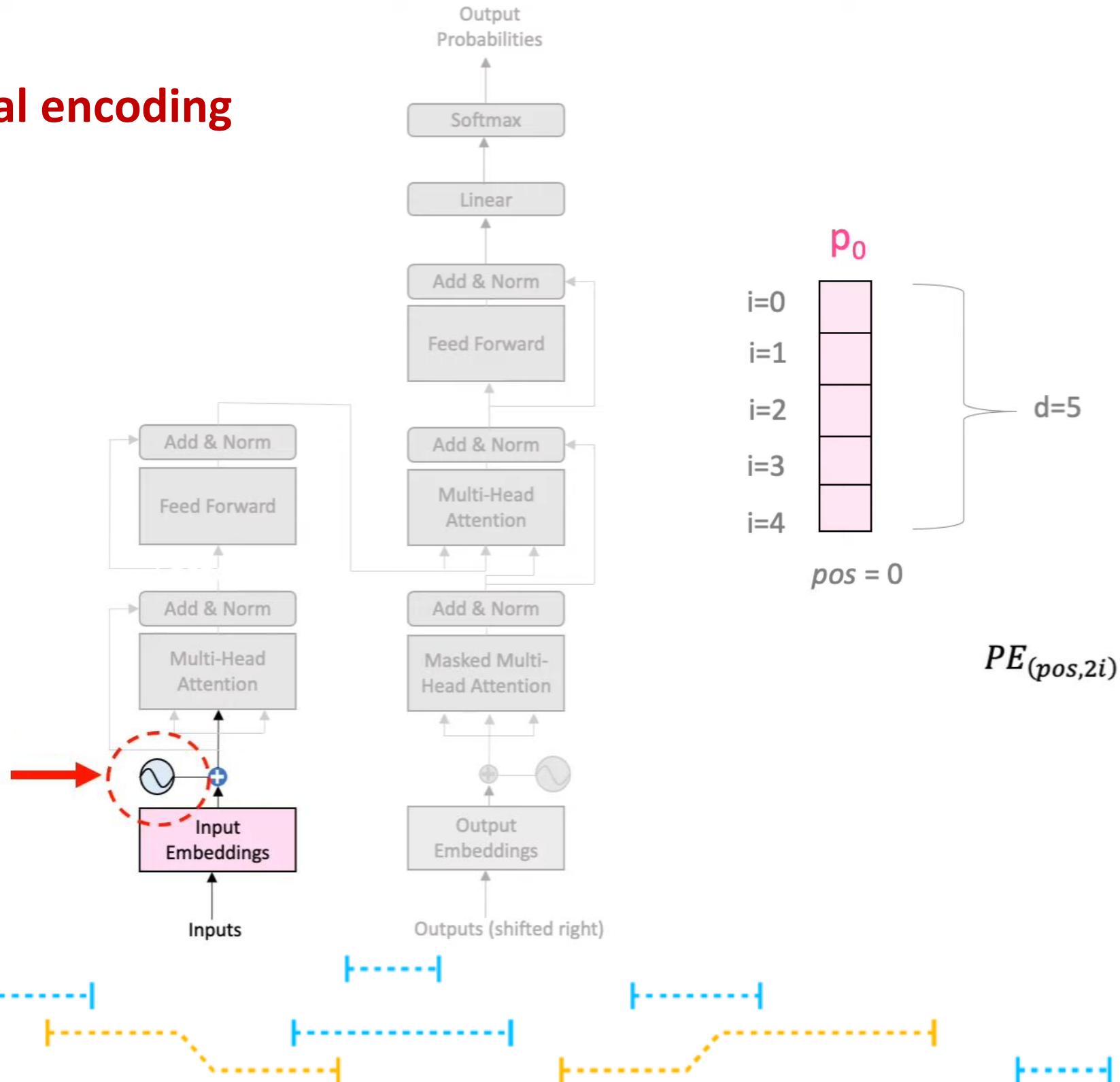
Vaswani et al 2017

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

Attention is *all you need*

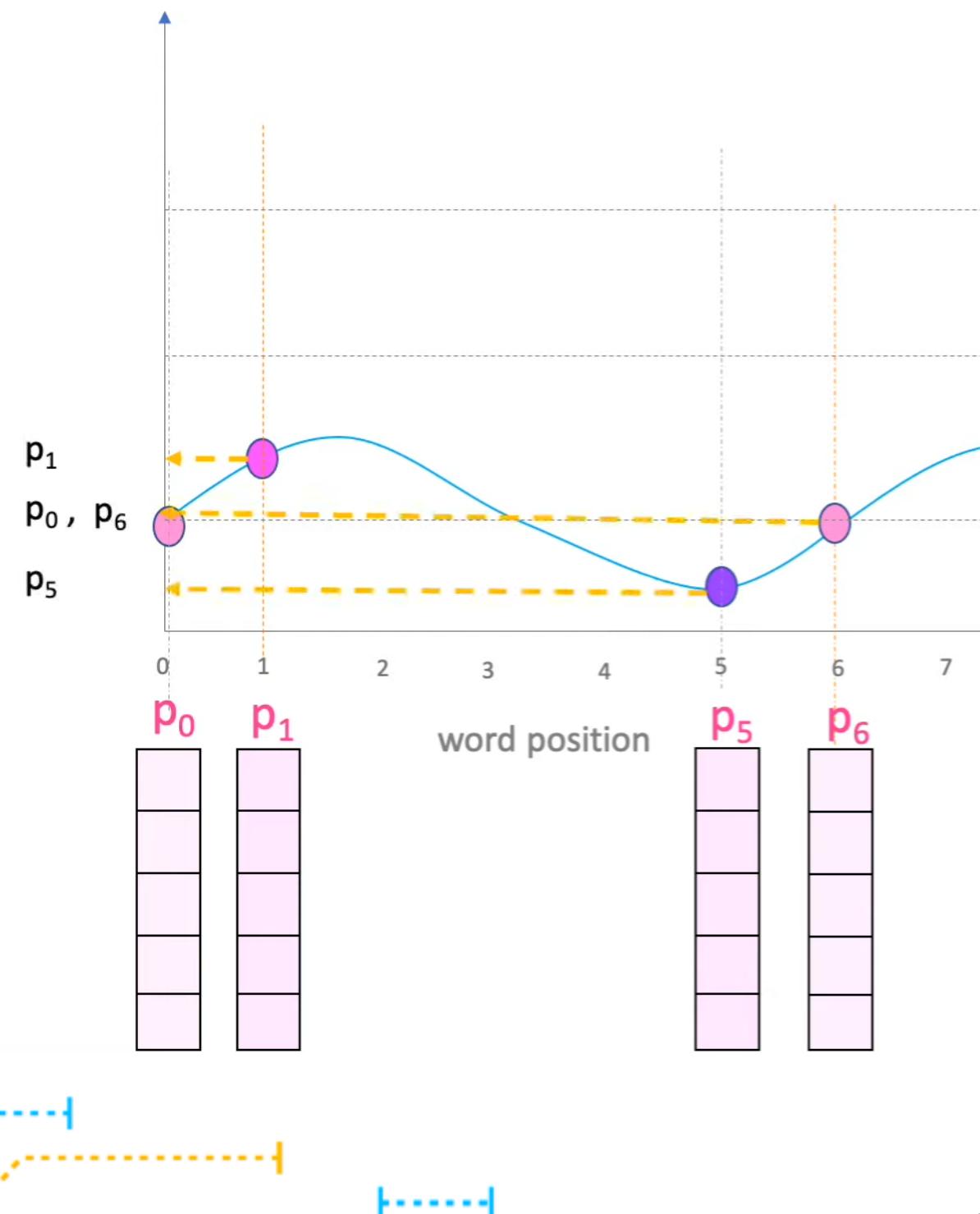
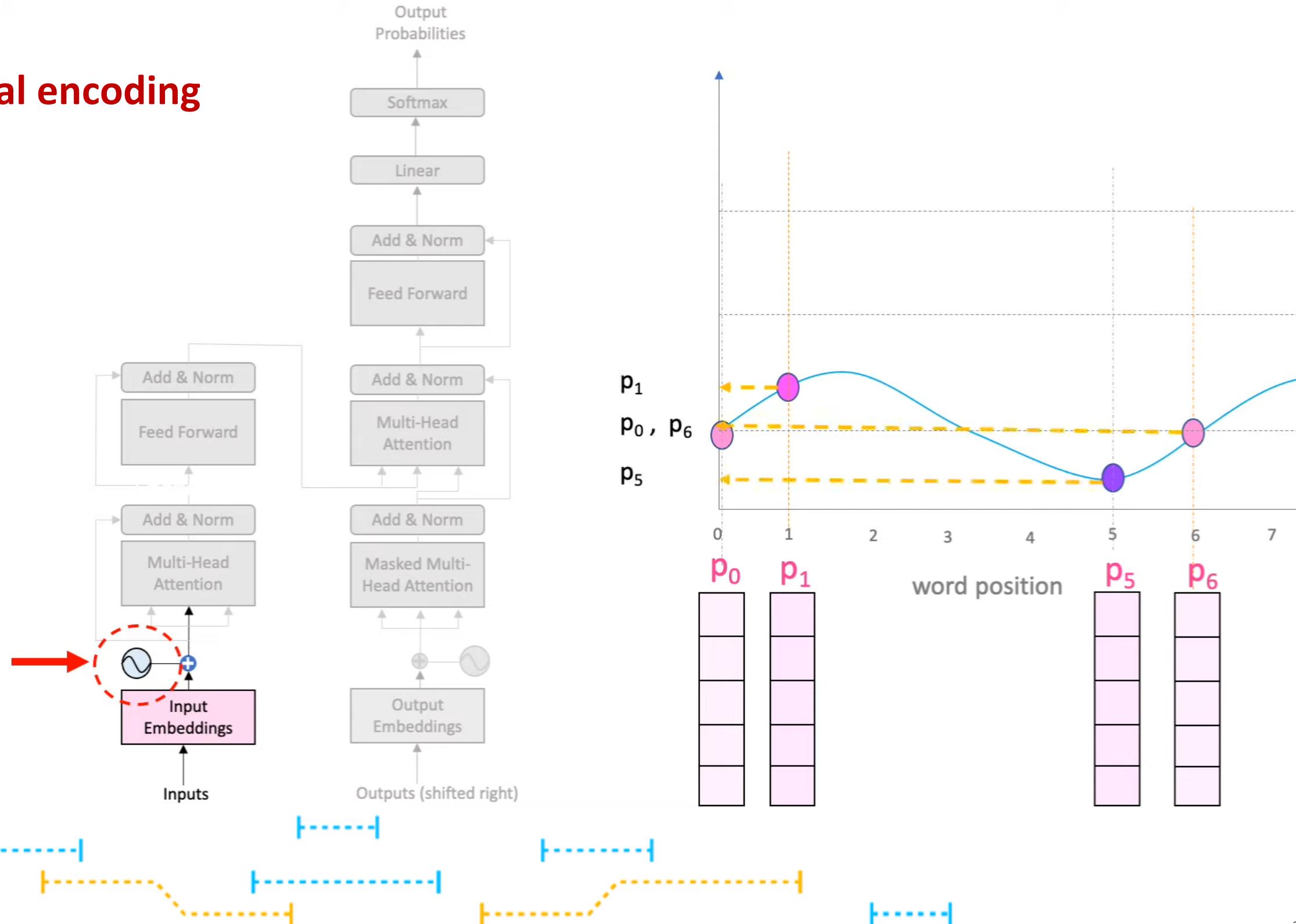
Positional encoding



$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

Attention is *all you need*

Positional encoding

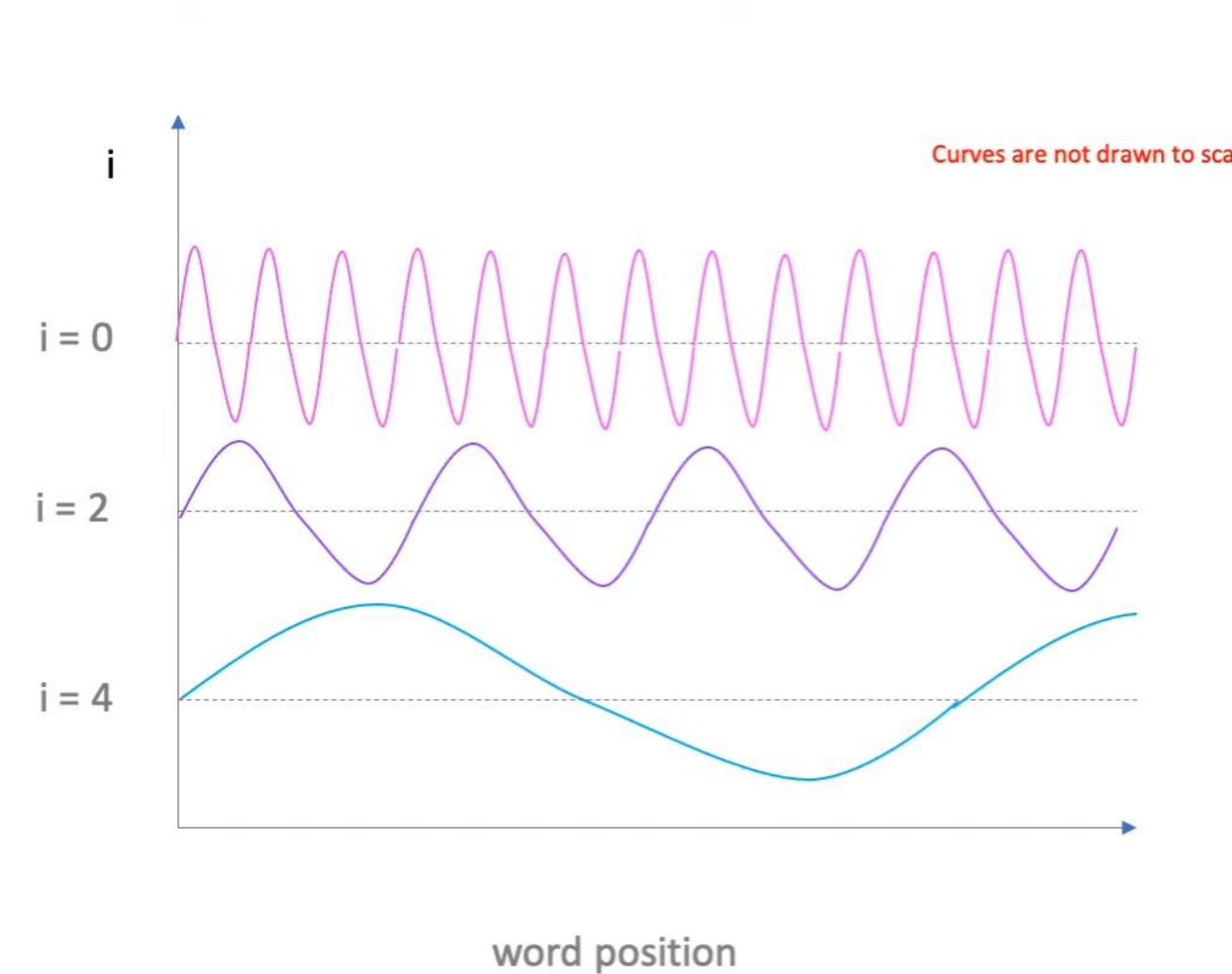
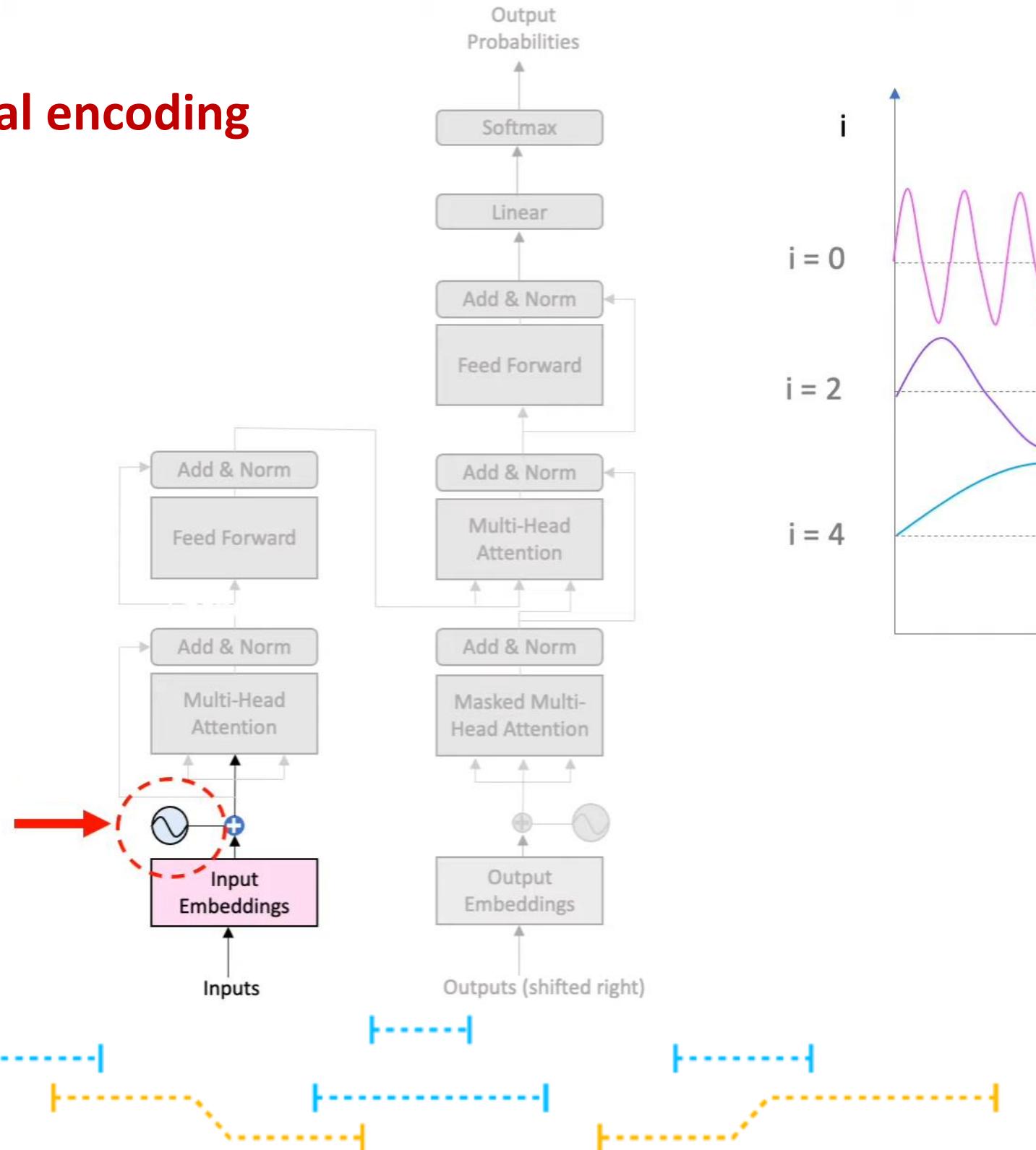


TRANSFORMATEC

Ashish Vaswani et al. (2017) "Attention is all you need".
31st Conference on Neural Information Processing Systems (NIPS 2017)

Attention is *all you need*

Positional encoding



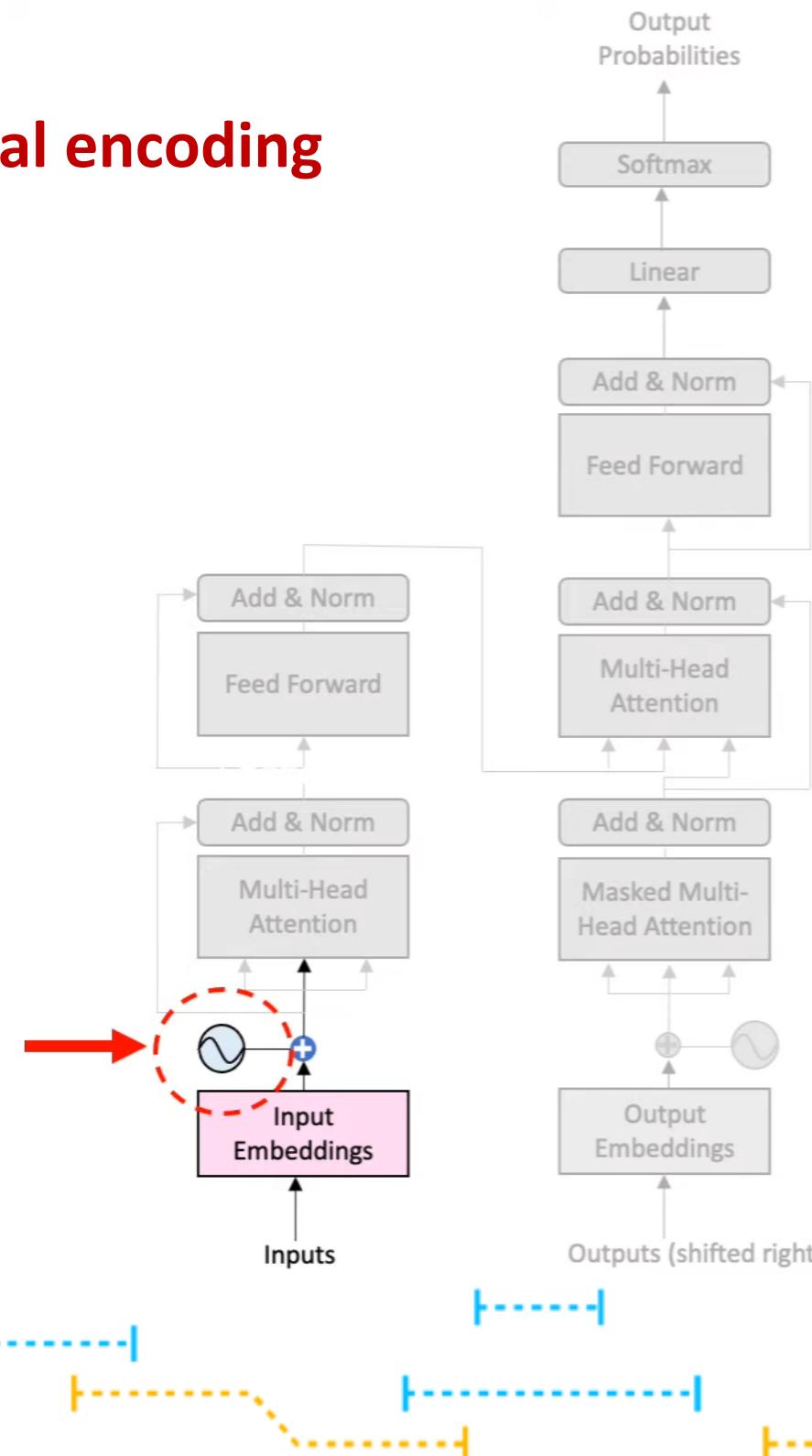
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d}}}\right)$$

TRANSFORMATEC

Ashish Vaswani et al. (2017) "Attention is all you need".
31st Conference on Neural Information Processing Systems (NIPS 2017)

Attention is *all you need*

Positional encoding



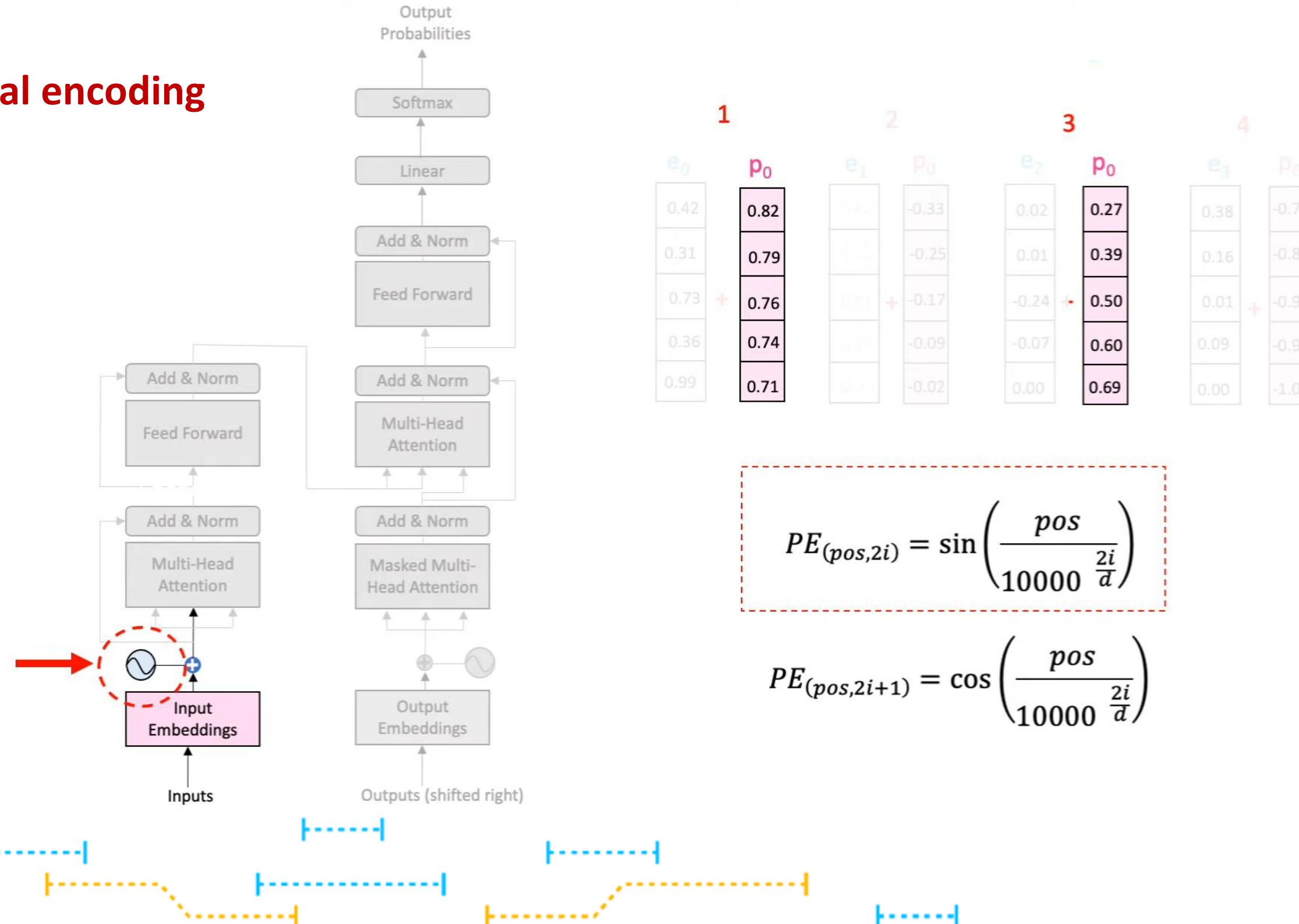
When	pos = 1	you	play	the	...
	e_0	e_1	e_2	e_3	
	0.42	0.87	0.02	0.38	
	0.31	-0.64	0.01	0.16	
	0.73	0.76	-0.24	0.01	
	0.36	0.74	-0.09	-0.07	
	0.99	-0.35	0.00	0.00	
	0.71	-0.02	0.69	-1.00	
	p_0	p_0	p_0	p_0	
	0.82	-0.33	0.27	-0.78	
	0.79	-0.25	0.39	-0.87	
	-0.64	-0.17	0.50	-0.94	
	0.81	0.09	0.60	-0.98	
	0.26	-0.02	0.69		
	-0.35				

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

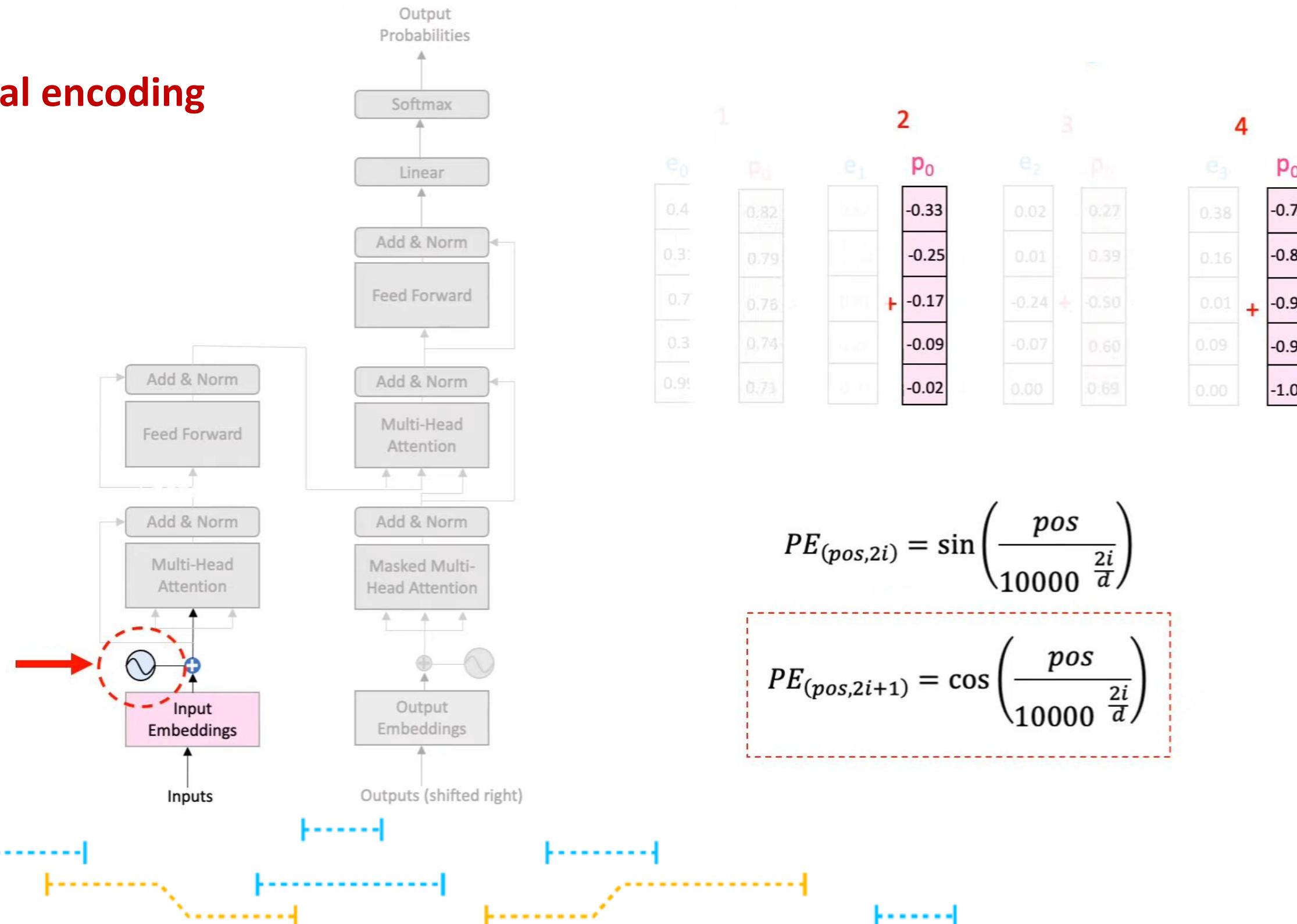
Attention is *all you need*

Positional encoding



Attention is *all you need*

Positional encoding



Attention is *all you need*

Positional encoding

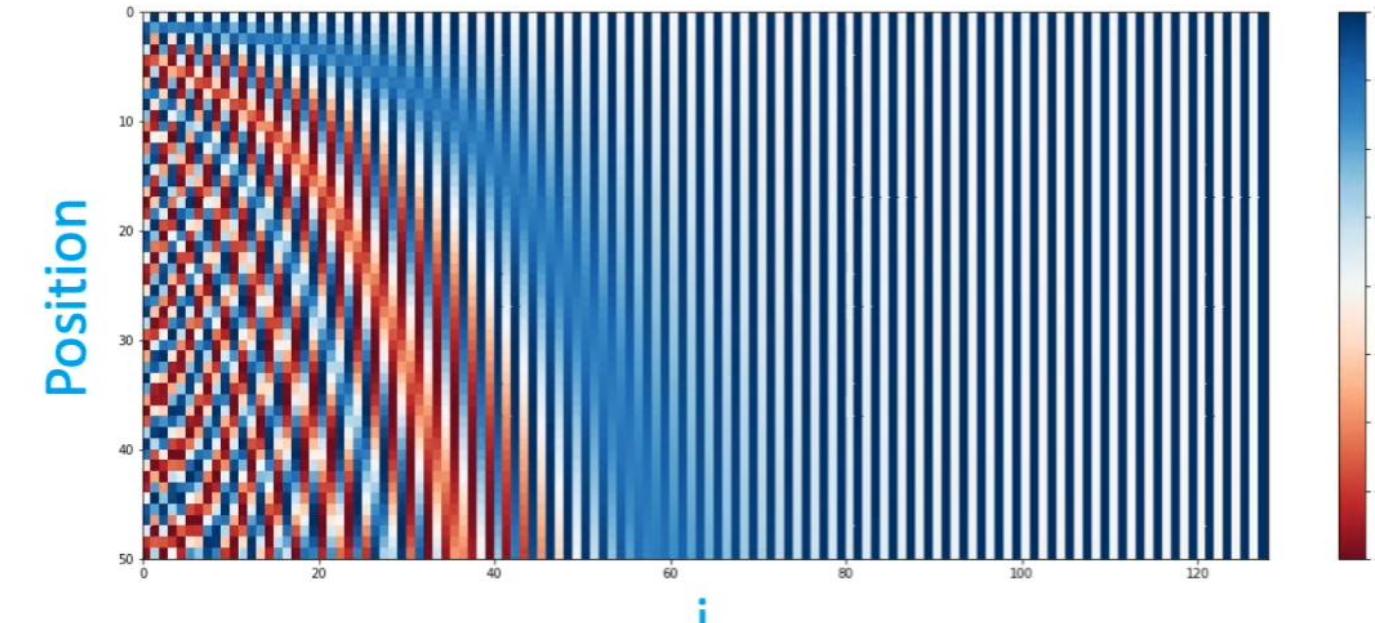
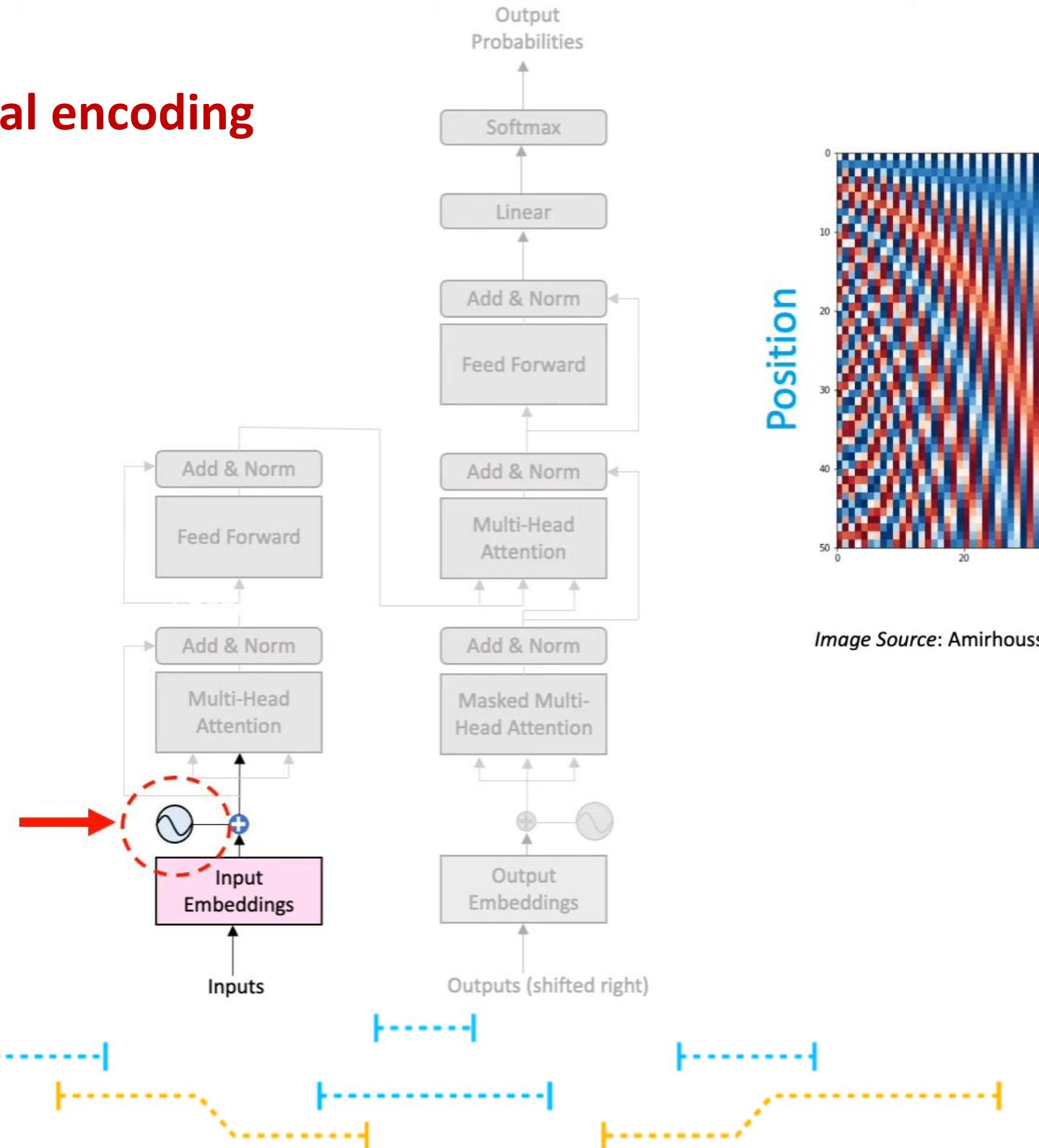


Image Source: Amiroussein Kazemnejad blog-post *Transformer Architecture: The Positional Encoding*

$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

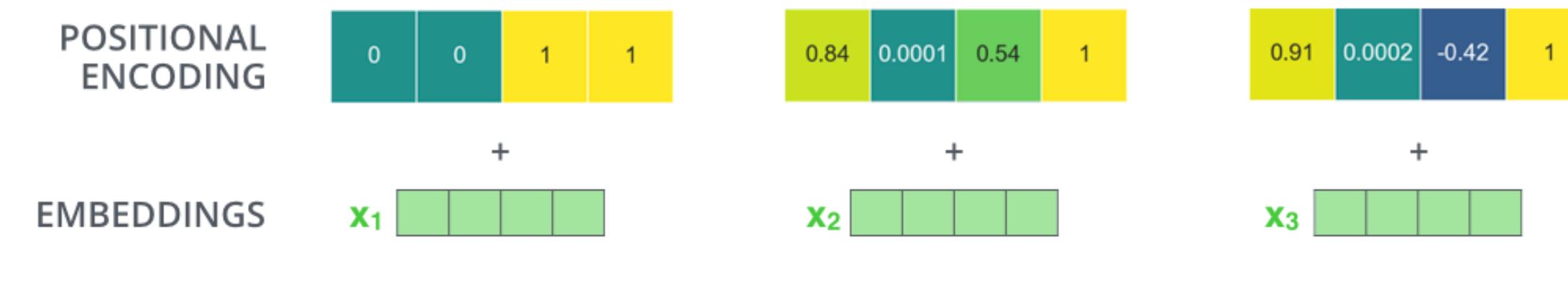
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000} \frac{2i}{d}\right)$$

TRANSFORMATEC

Ashish Vaswani et al. (2017) "Attention is all you need".
31st Conference on Neural Information Processing Systems (NIPS 2017)

Attention is *all you need*

Positional encoding



INPUT

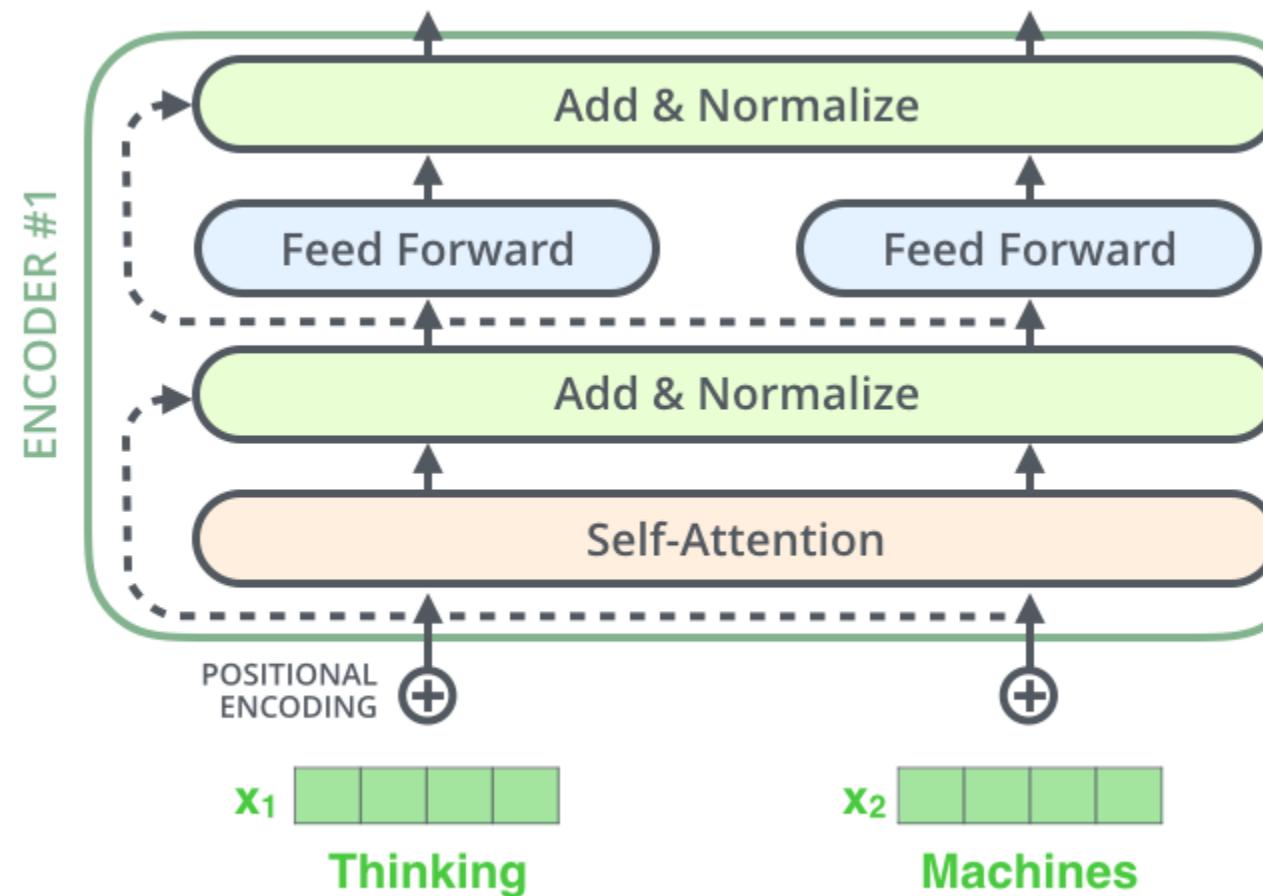
Je

suis

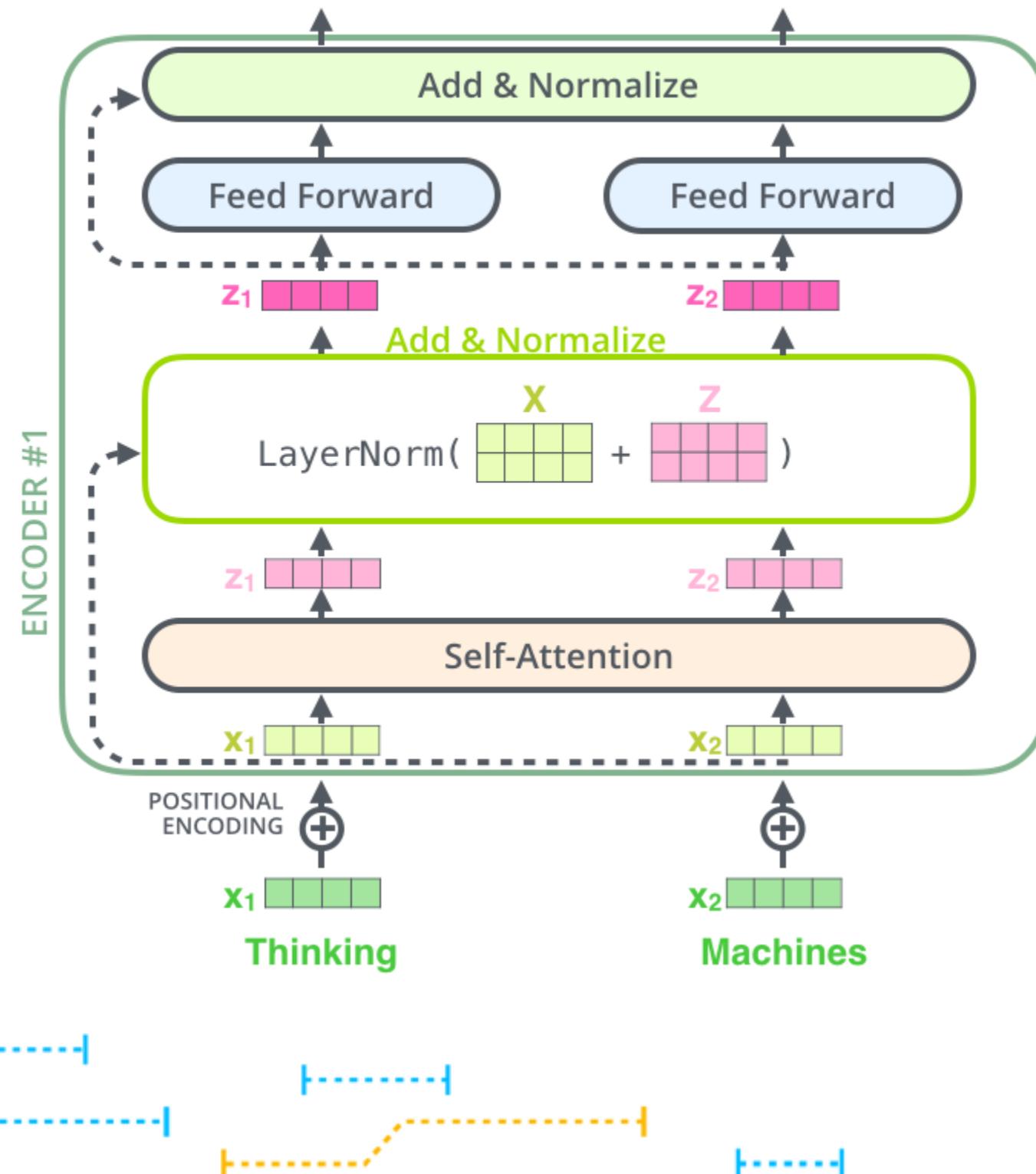
étudiant



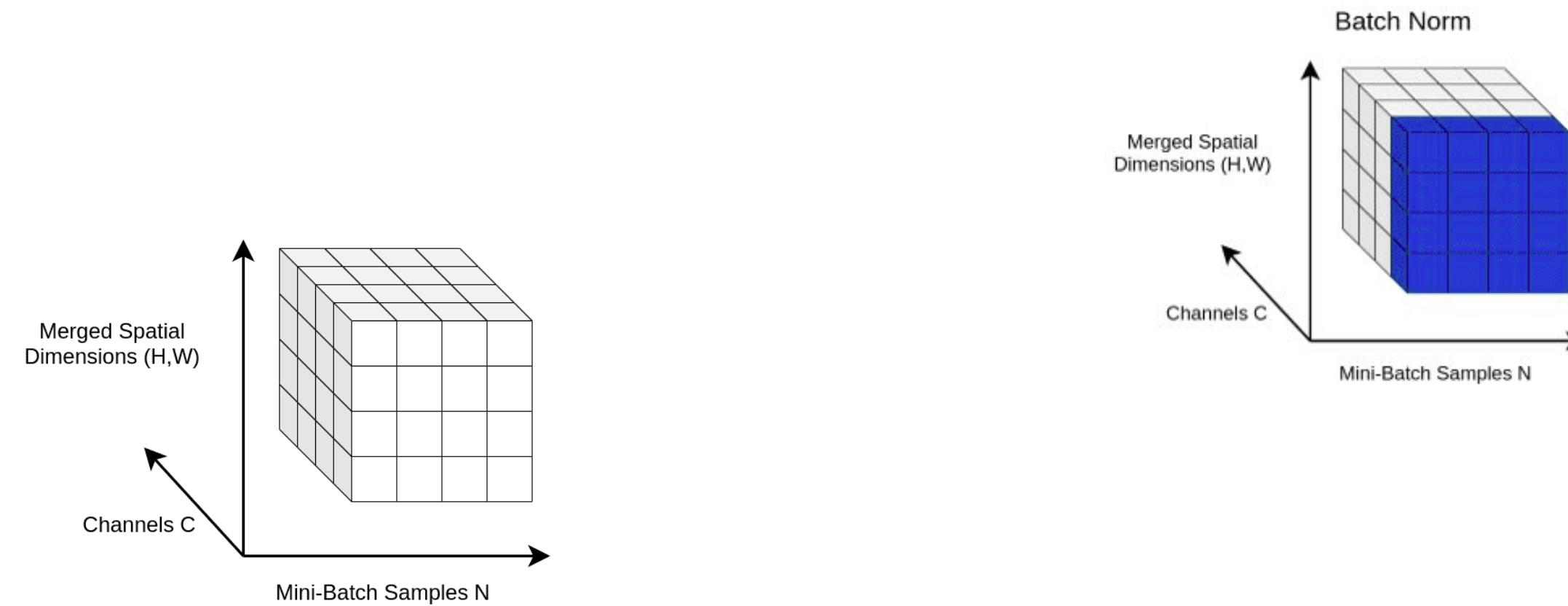
Attention is *all you need*



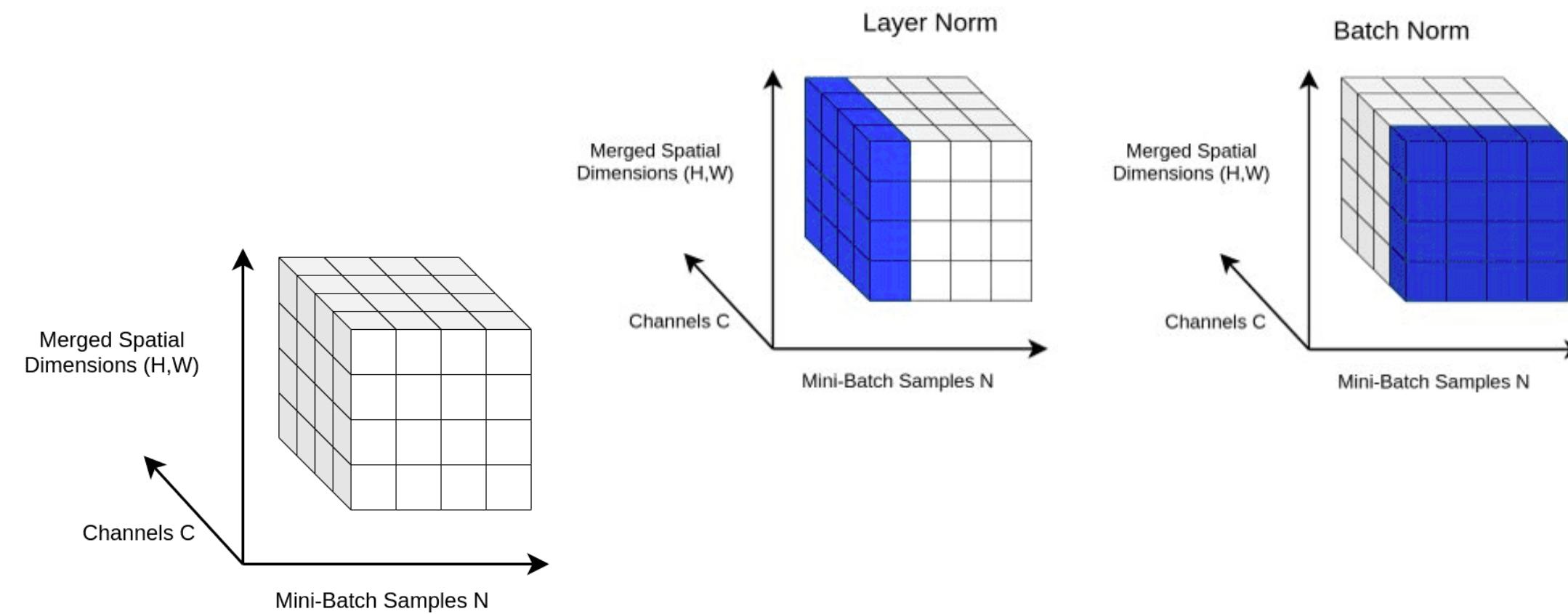
Attention is *all you need*



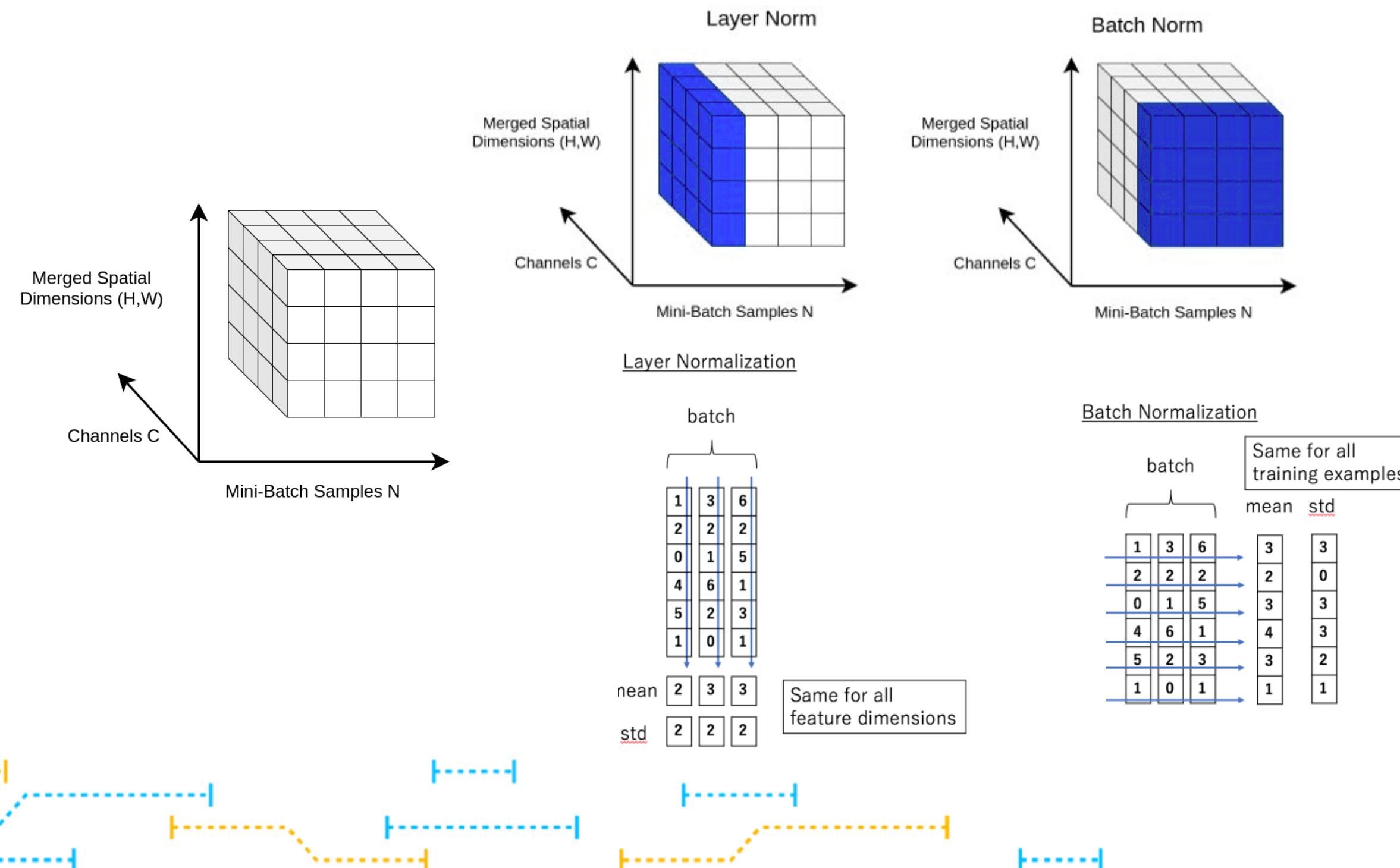
In-layer normalization techniques



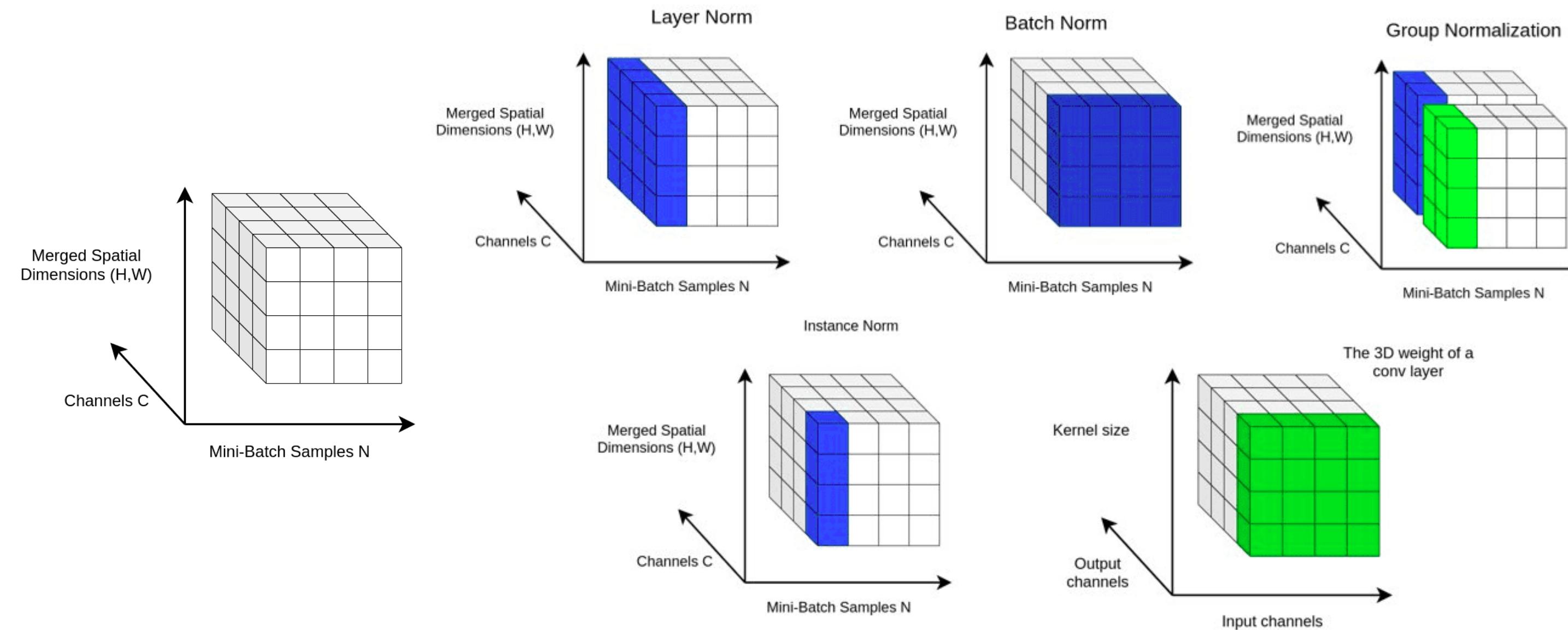
In-layer normalization techniques



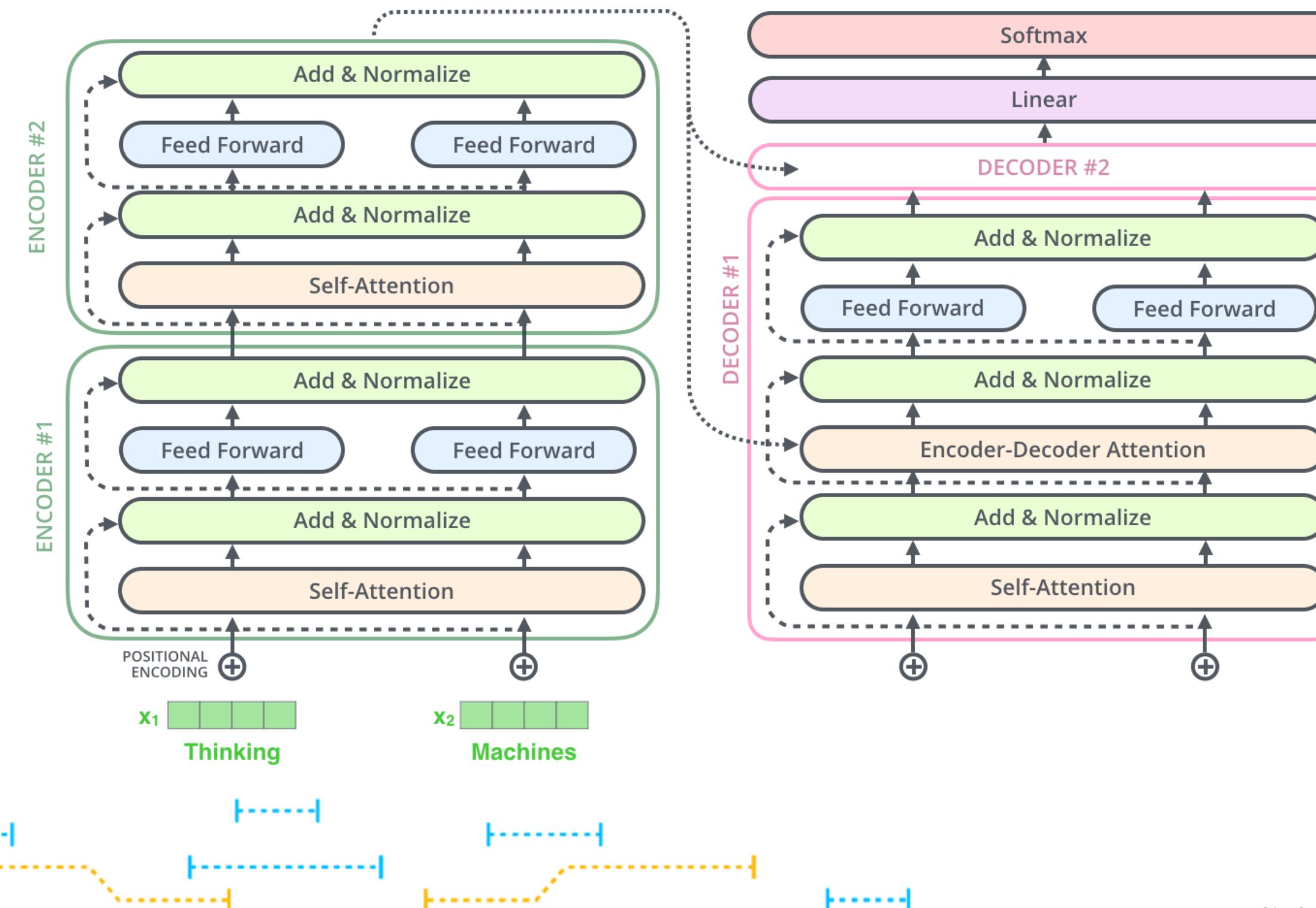
In-layer normalization techniques



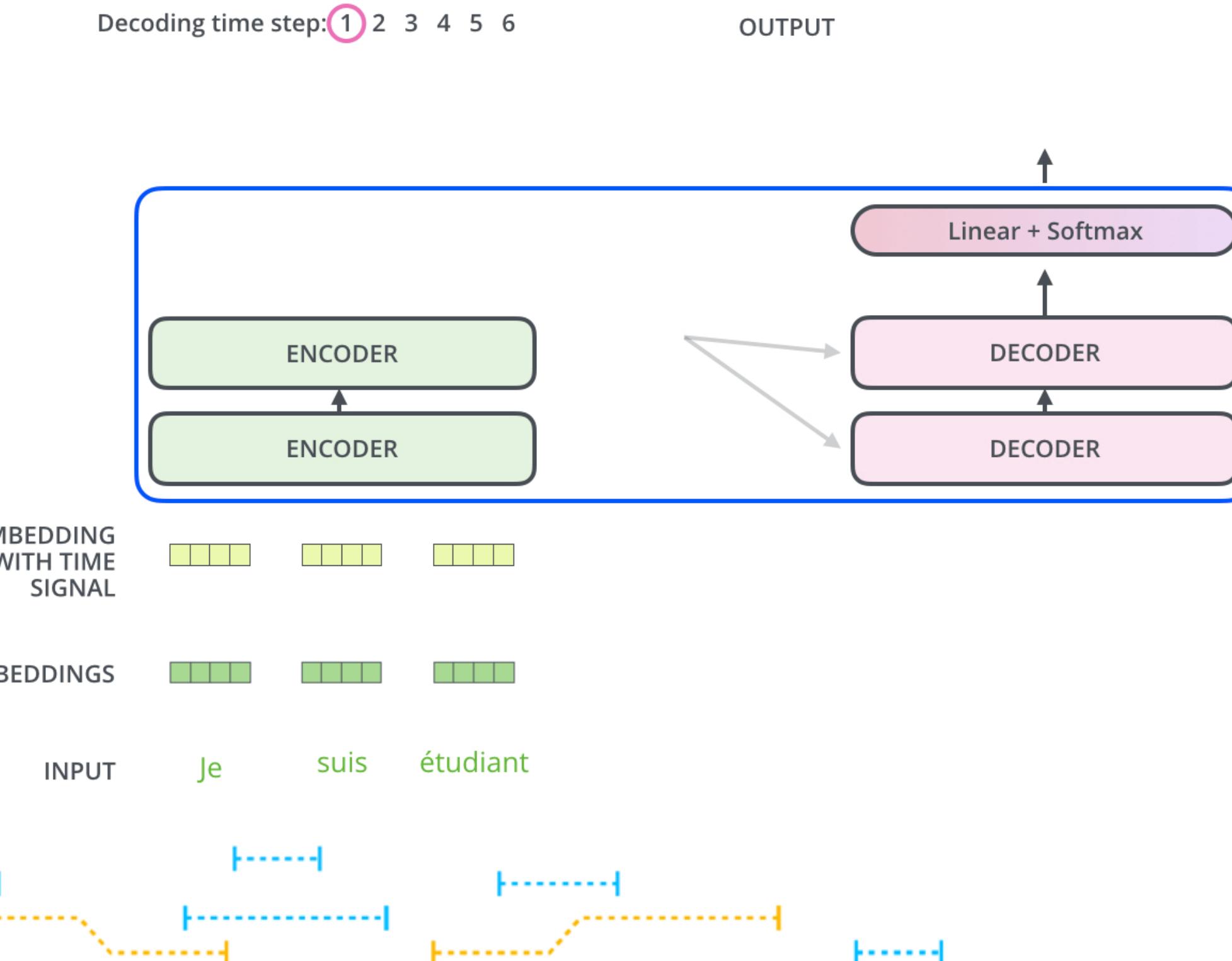
In-layer normalization techniques



Attention is *all you need*

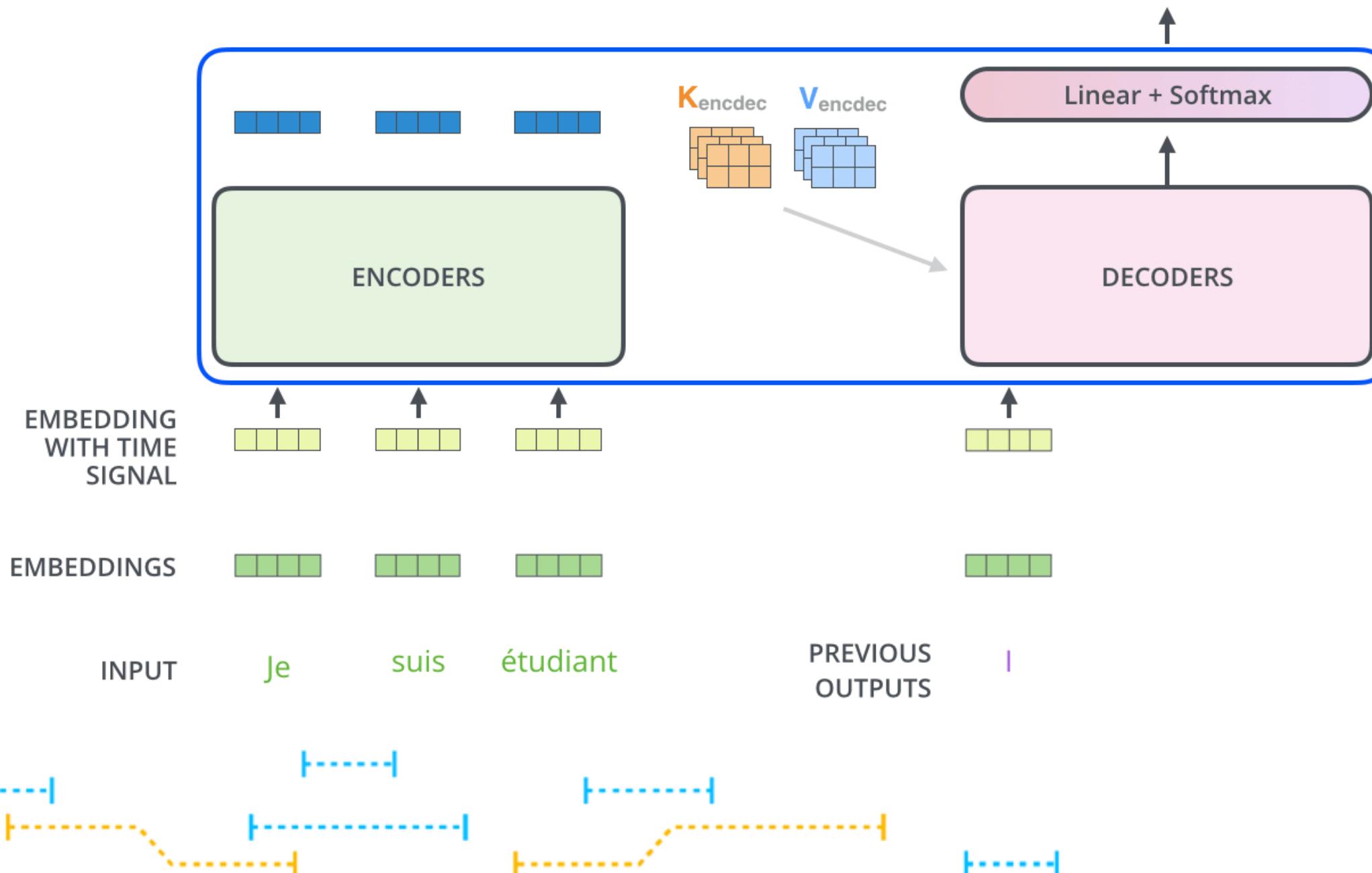


Attention is *all you need*



Attention is *all you need*

Decoding time step: 1 2 3 4 5 6 OUTPUT |



Attention is *all you need*

Which word in our vocabulary
is associated with this index?

am

Get the index of the cell
with the highest value
(argmax)

5

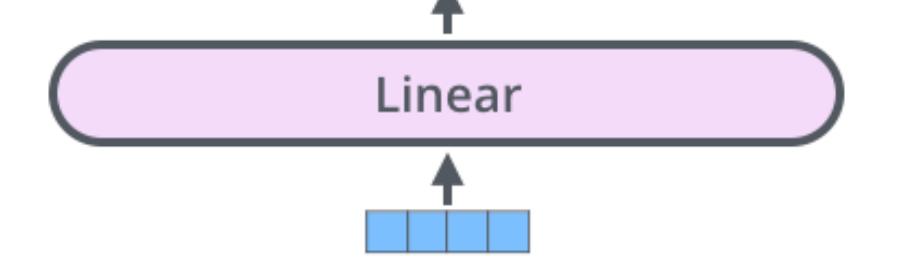
log_probs



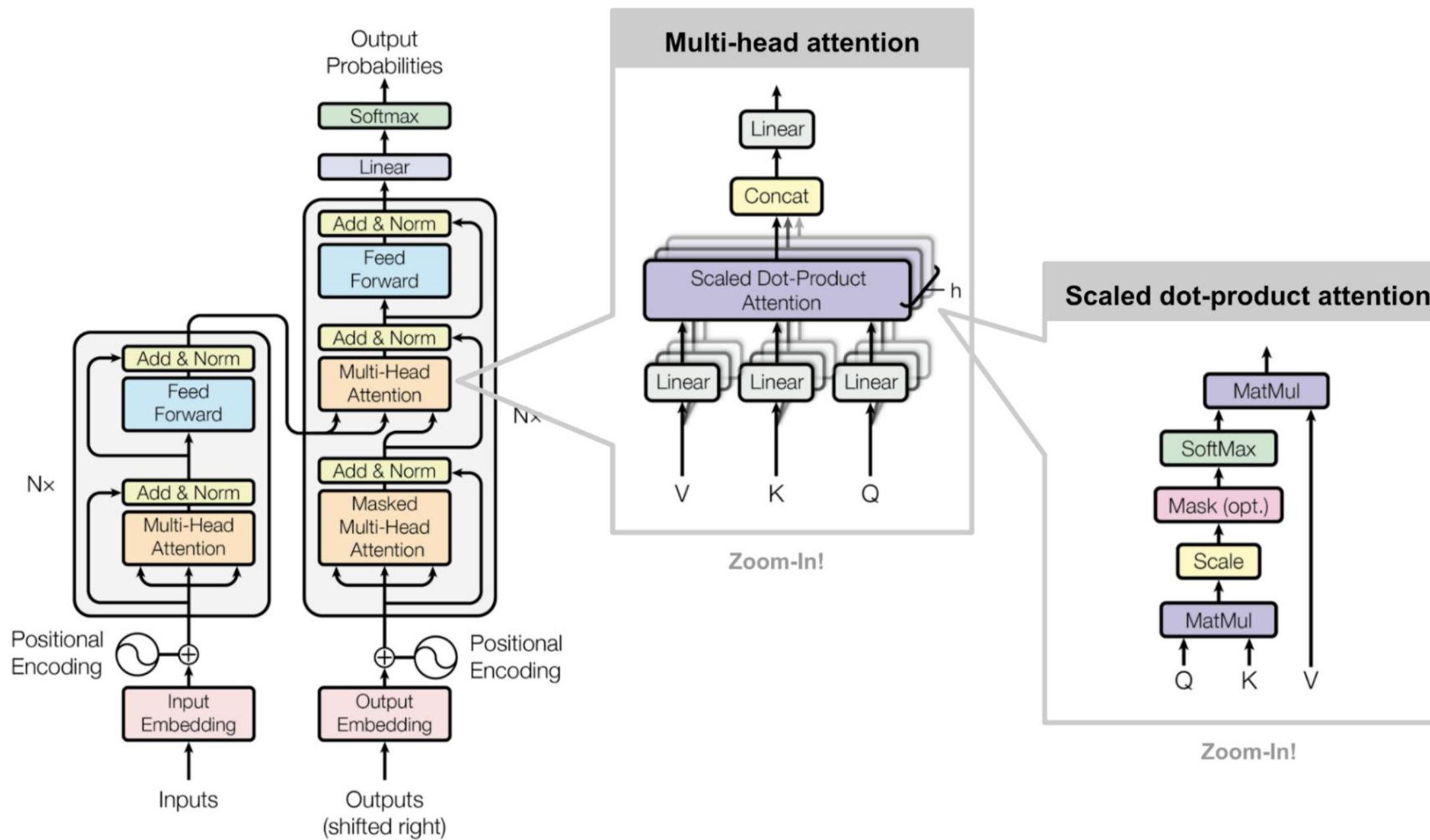
logits



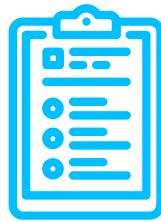
Decoder stack output



Attention is *all you need*



3.



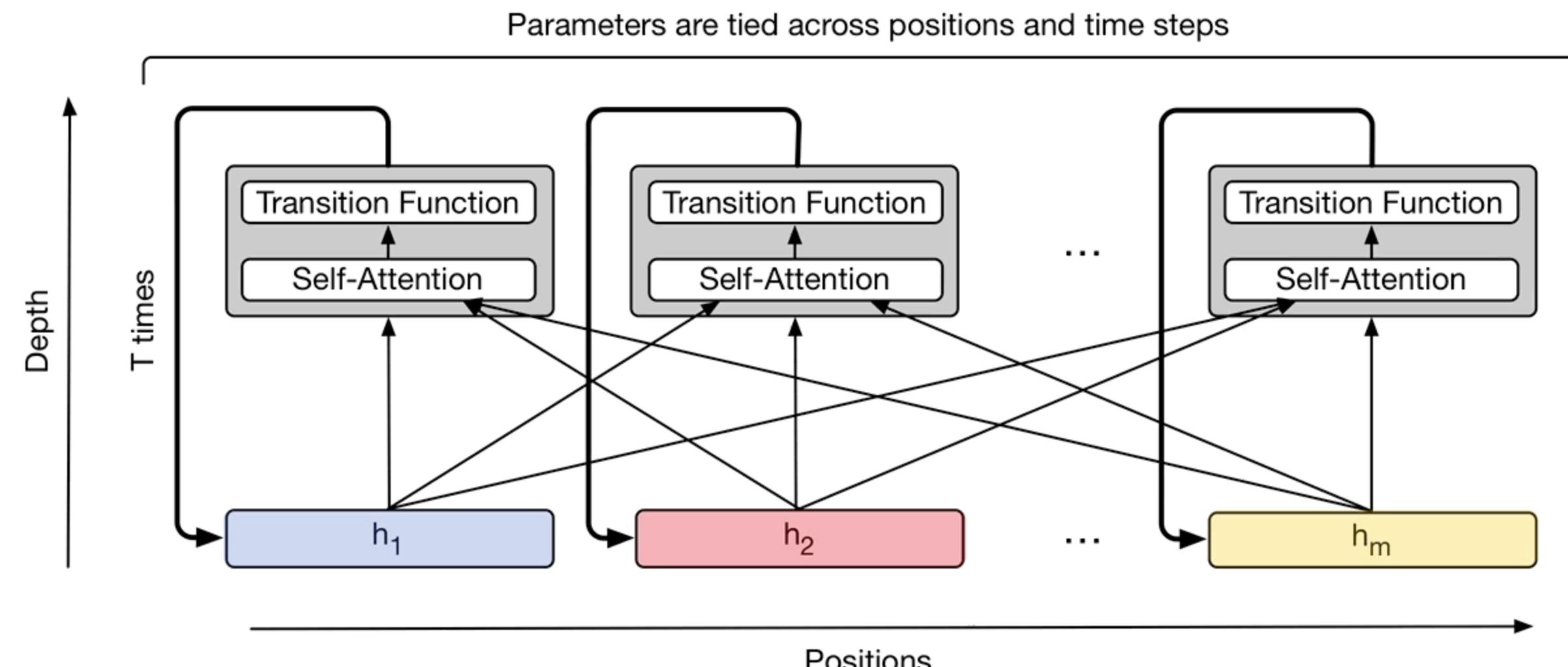
Universal *Transformers*

TRANSFORMATEC

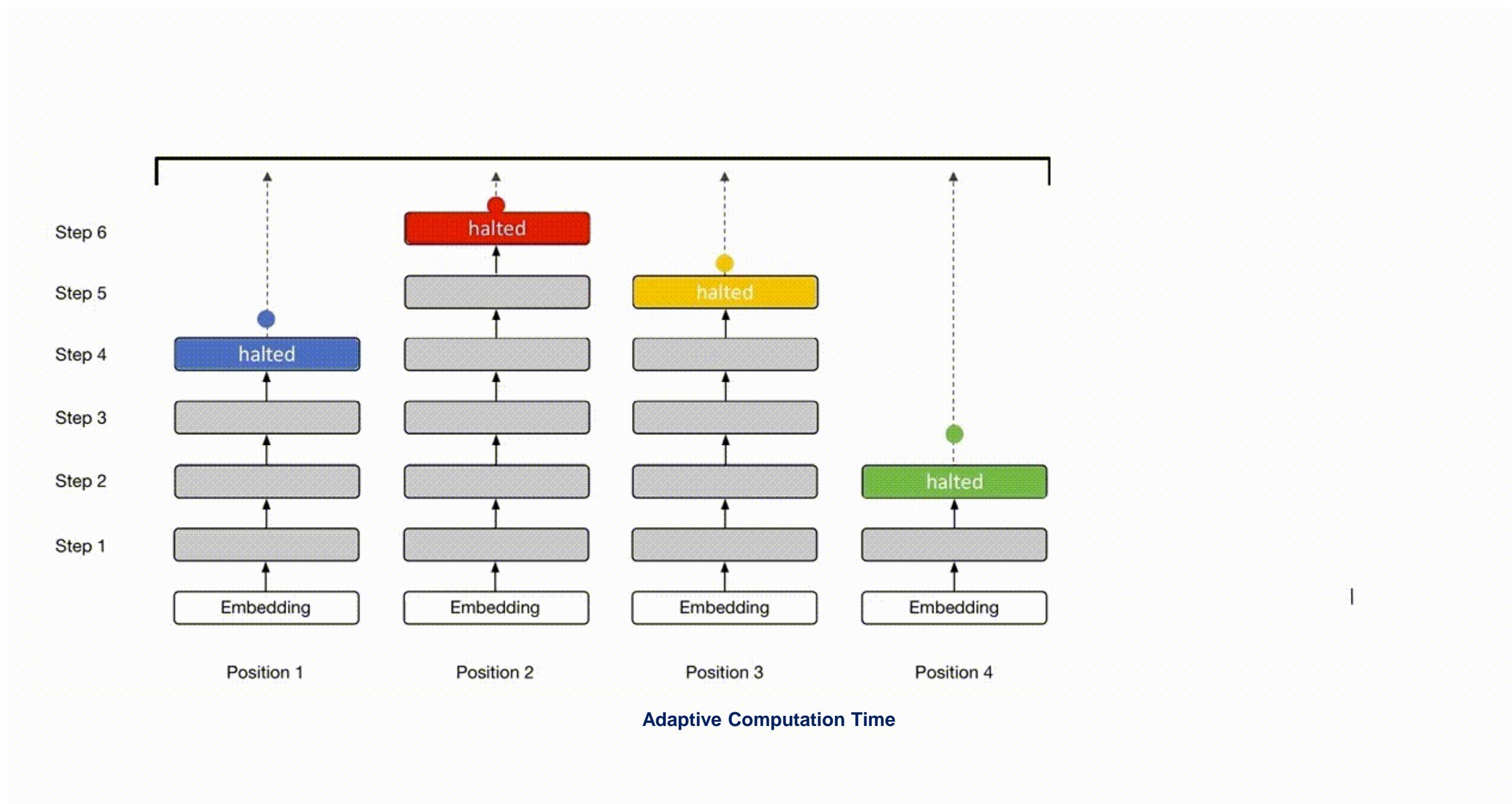
> Reinventa el mundo <



Universal *Transformers*



Universal *Transformers*



TRANSFORMATEC

Mostafa Dehghani et al. (2018) "Universal Transformers".
International Conference on Learning Representations (ICLR).

Universal Transformers

Adaptive Computation Time

Permite que el modelo determine de manera dinámica cuántas iteraciones son necesarias para procesar cada parte de un conjunto de datos.

En cada iteración t , el modelo calcula un valor de halting $h_t \in \mathbb{R}^L$ score para cada entrada:

$$h_t = \text{Sigmoid}(Wx_t)$$

El progreso acumulado P_t

$$P_t = P_{t-1} + h_t$$

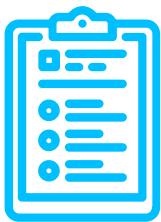
El modelo detiene el procesamiento de la entrada i cuando $P_t[i] \geq 1$.



TRANSFORMATEC

Alex Graves (2016) "Adaptive Computation Time for Recurrent Neural Networks".
arXiv preprint arXiv:1603.08983.

4.

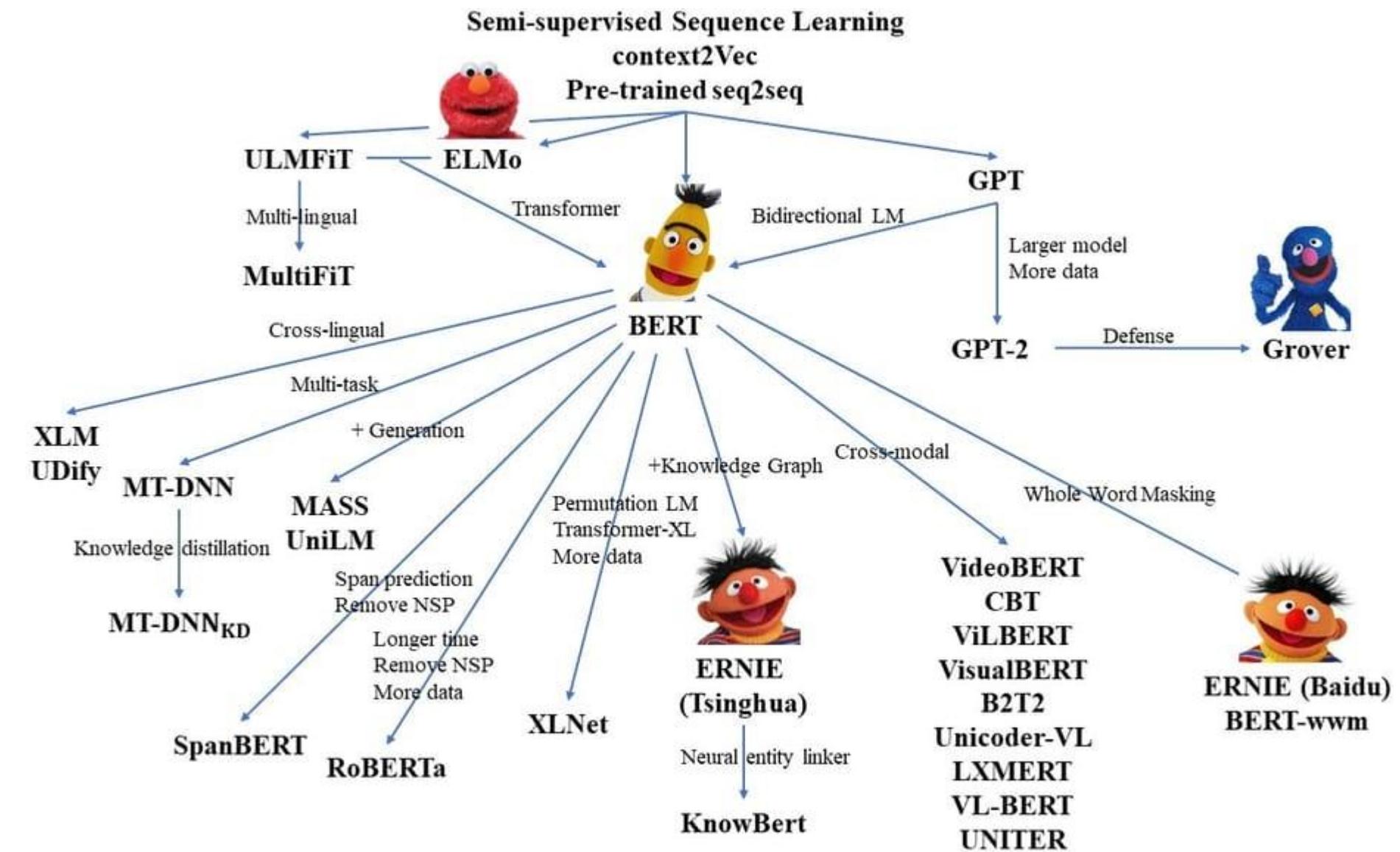


BERT

TRANSFORMATEC



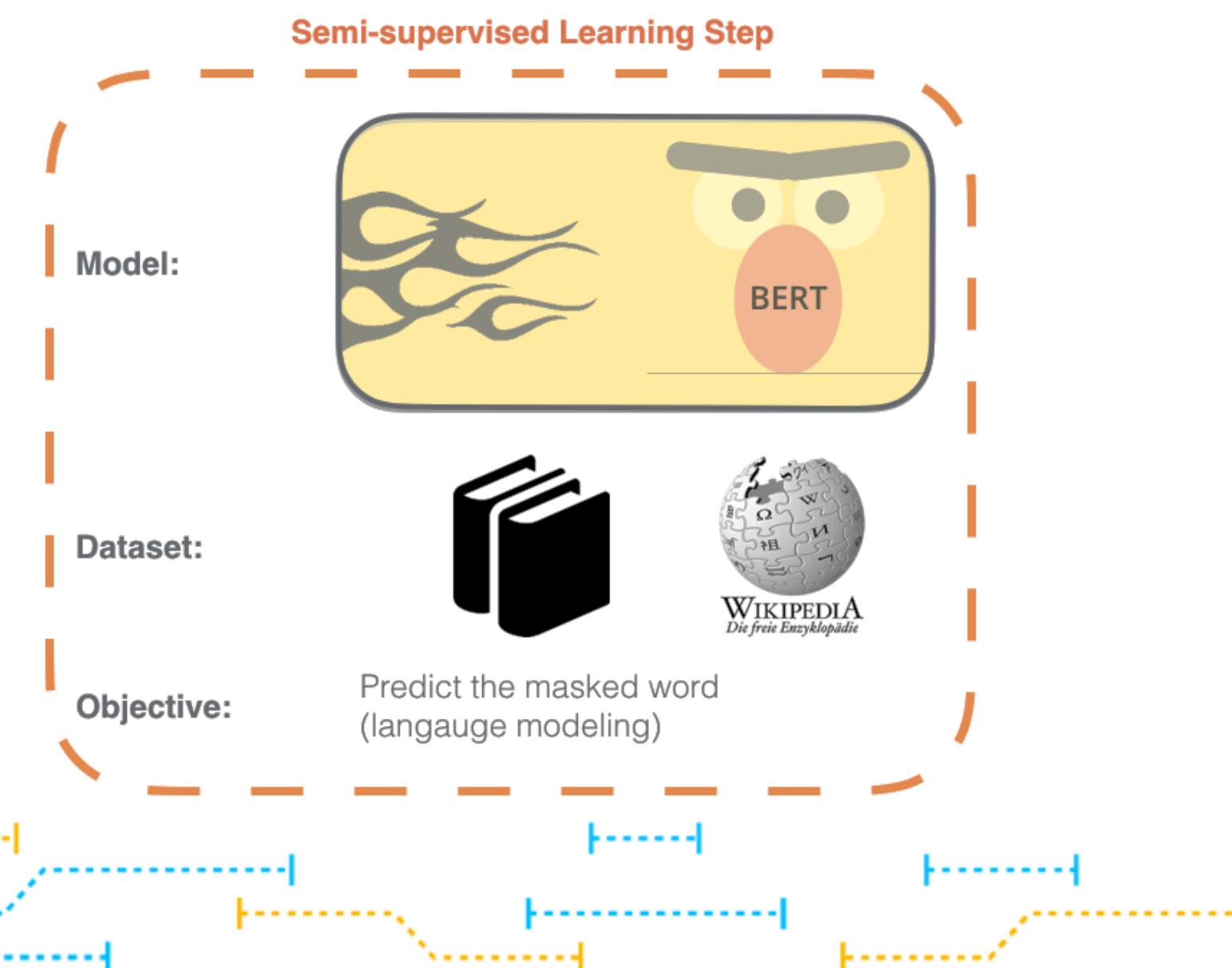
Sesame Street



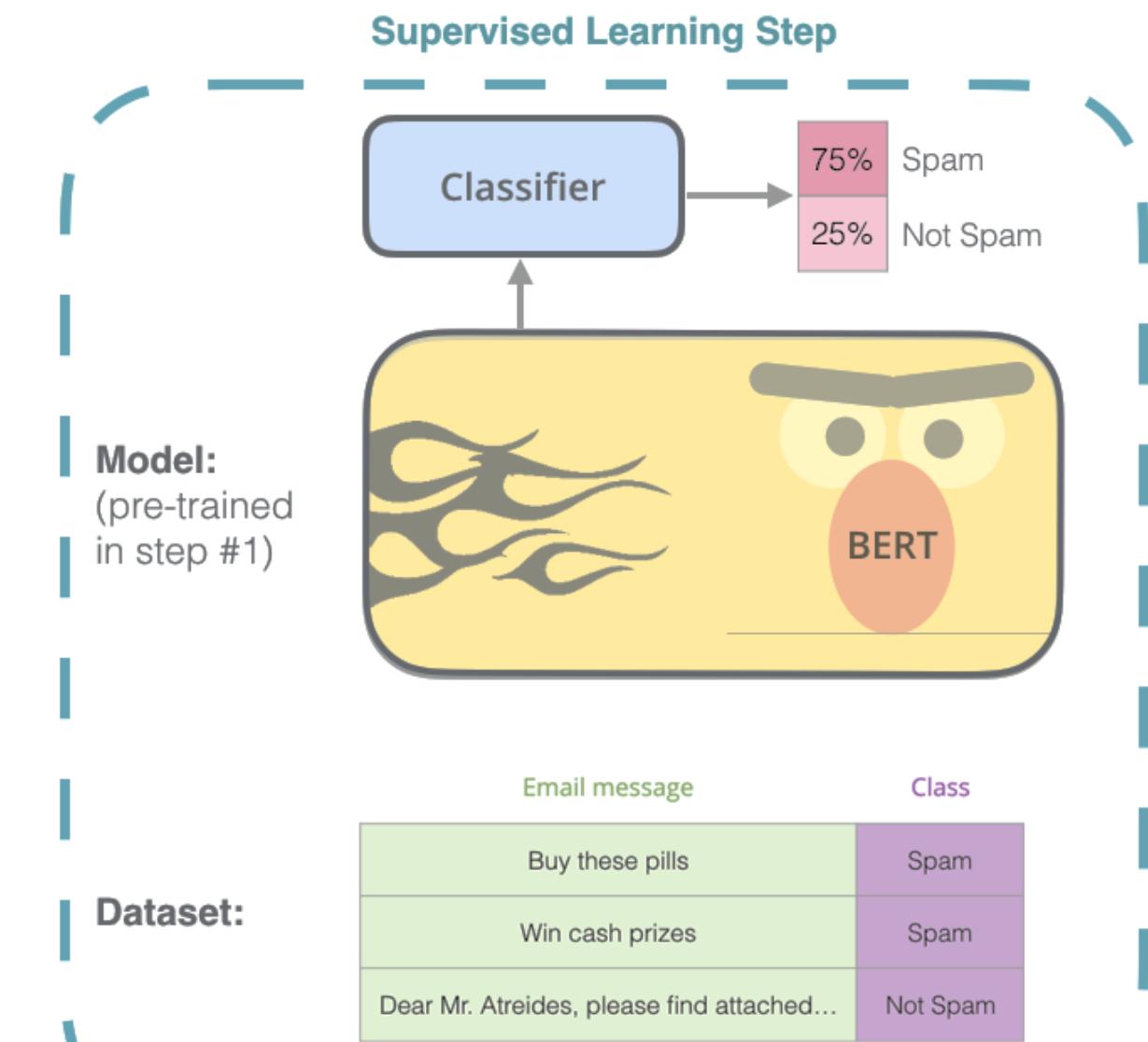
BERT

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

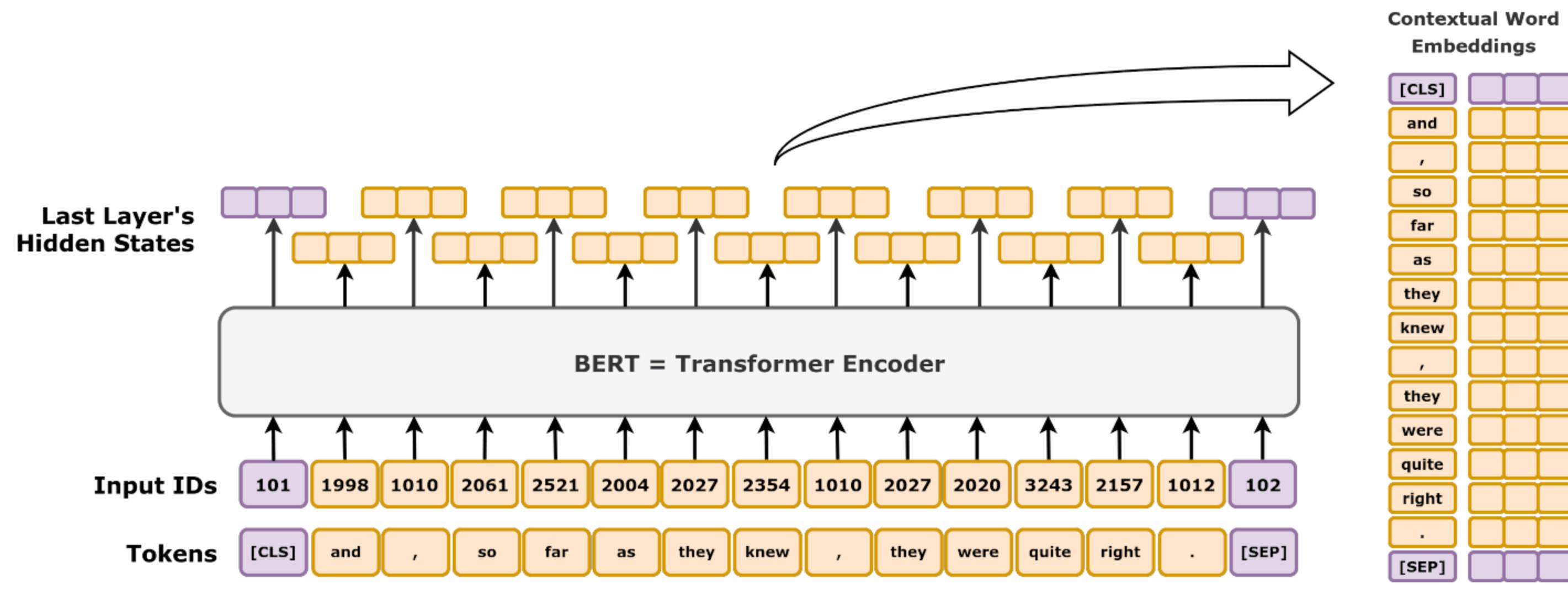
The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



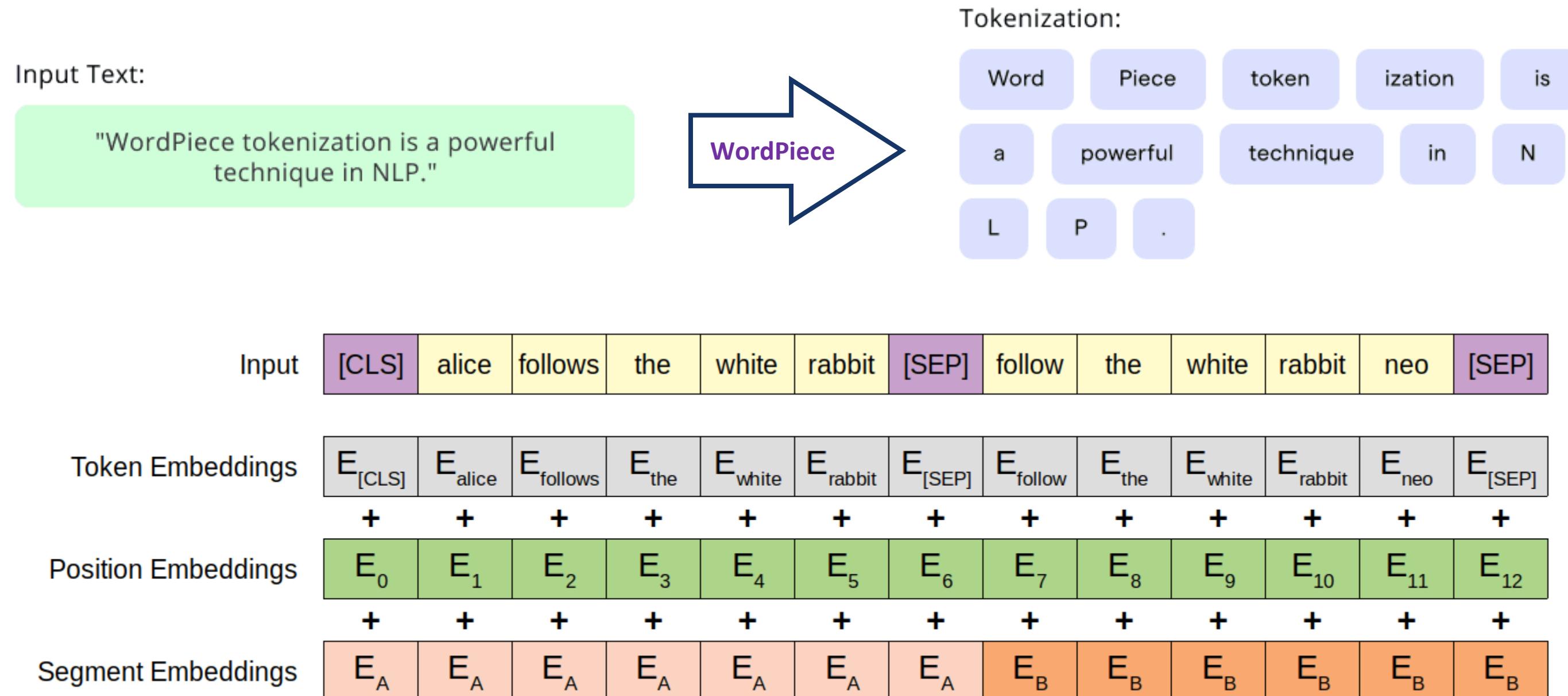
2 - Supervised training on a specific task with a labeled dataset.



BERT



BERT



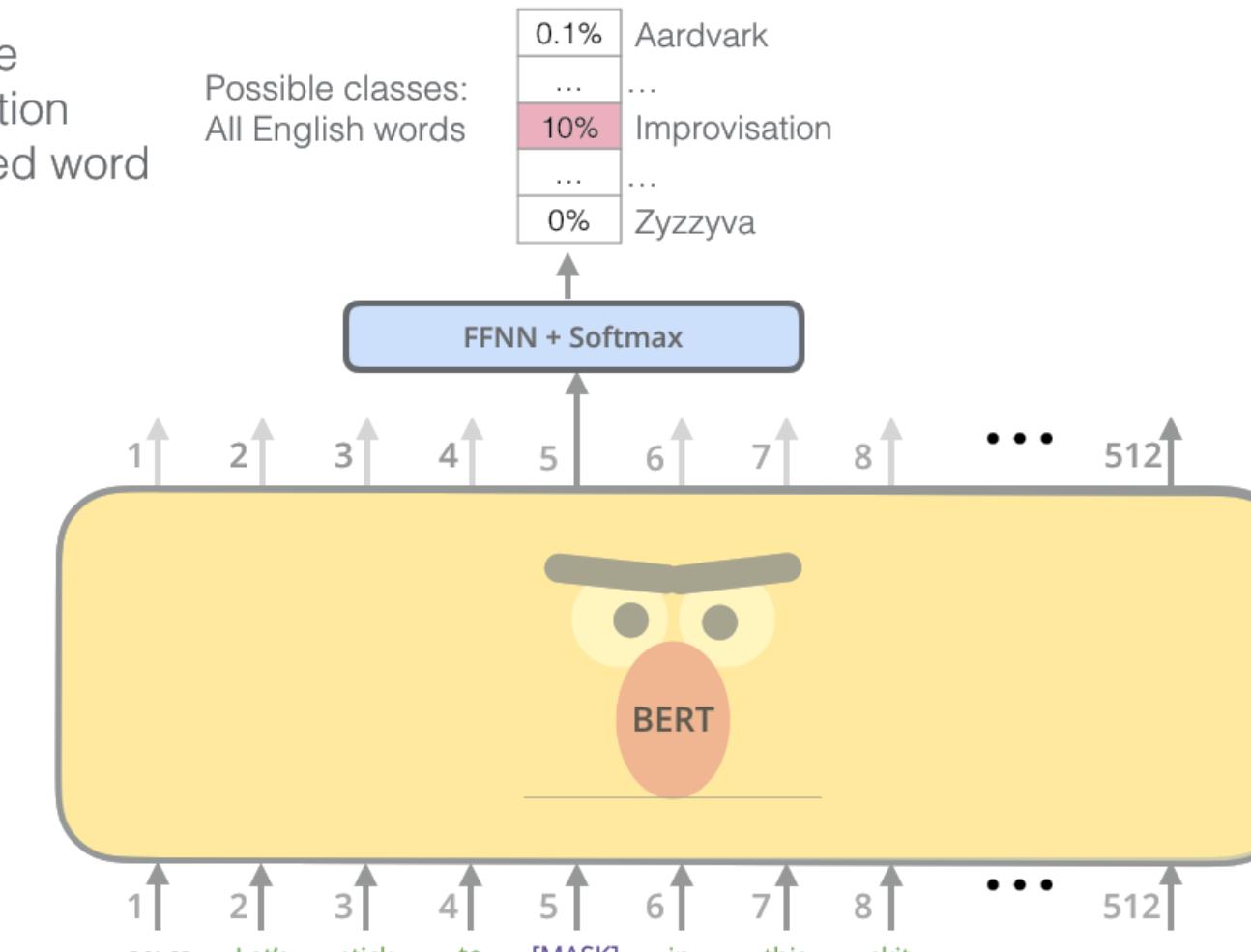
BERT: Pre-training Tasks

Masked Language Model

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

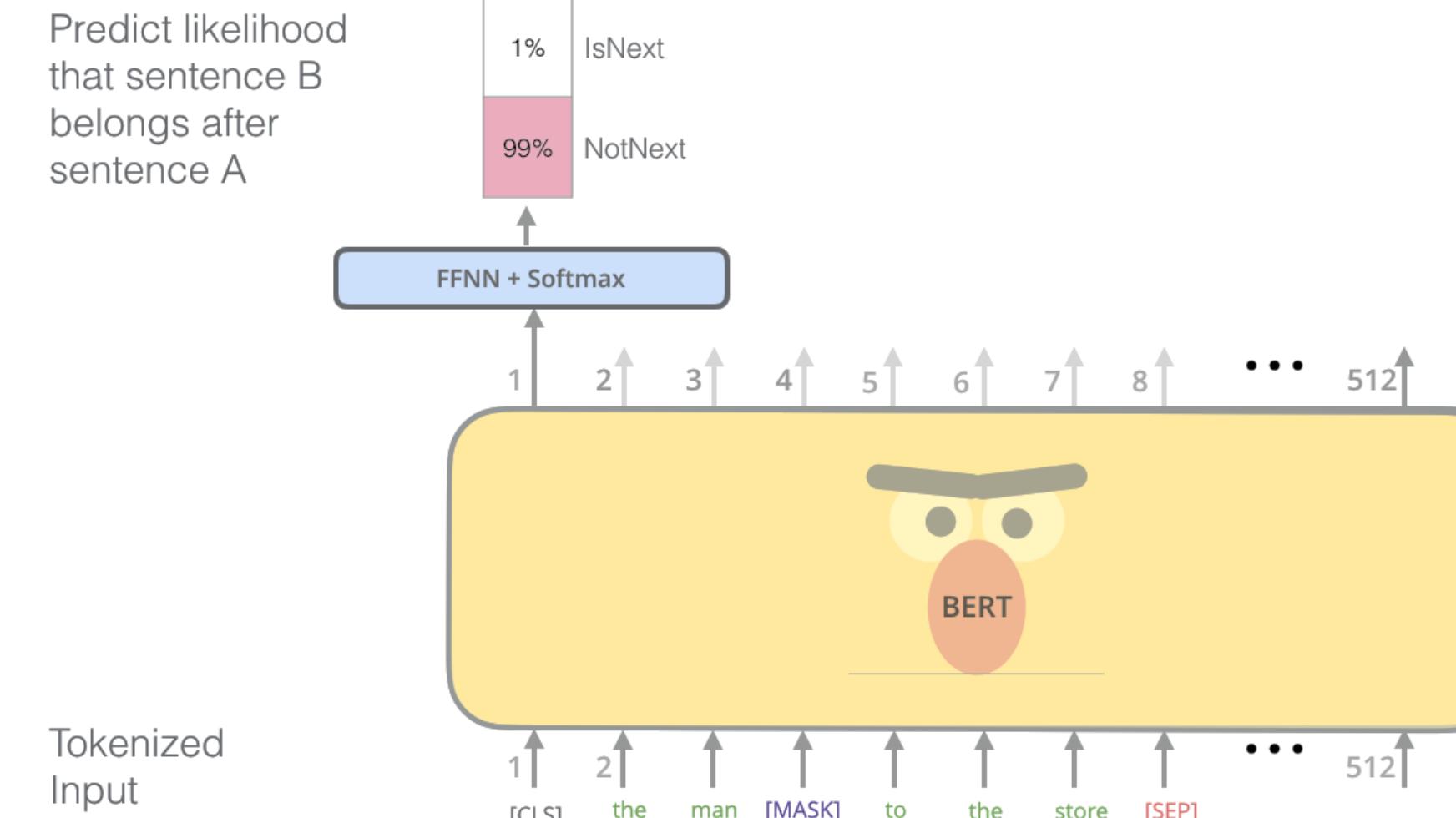
Randomly mask 15% of tokens

Input



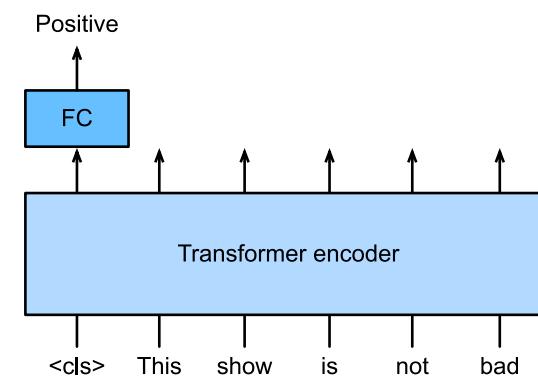
BERT: *Pre-training Tasks*

Two-sentence Tasks

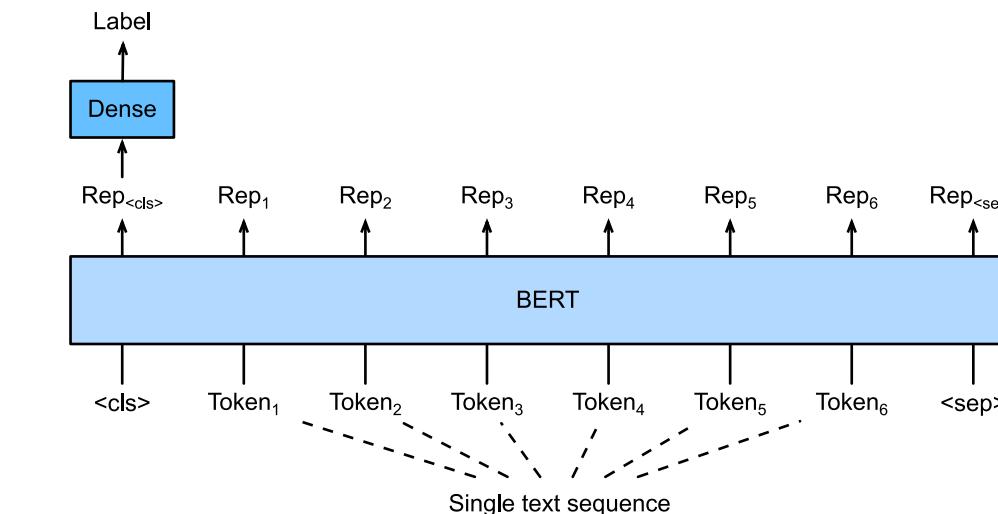


BERT: Pre-training Tasks

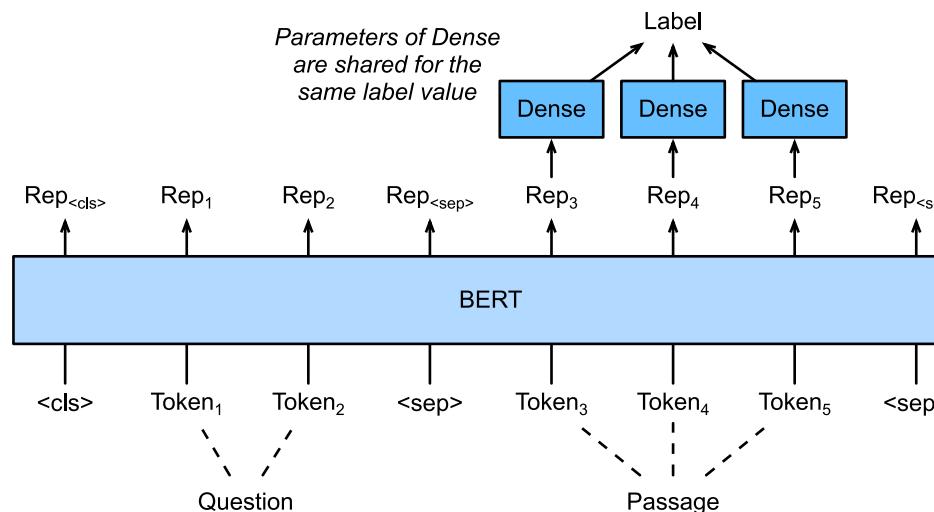
Sentiment classification



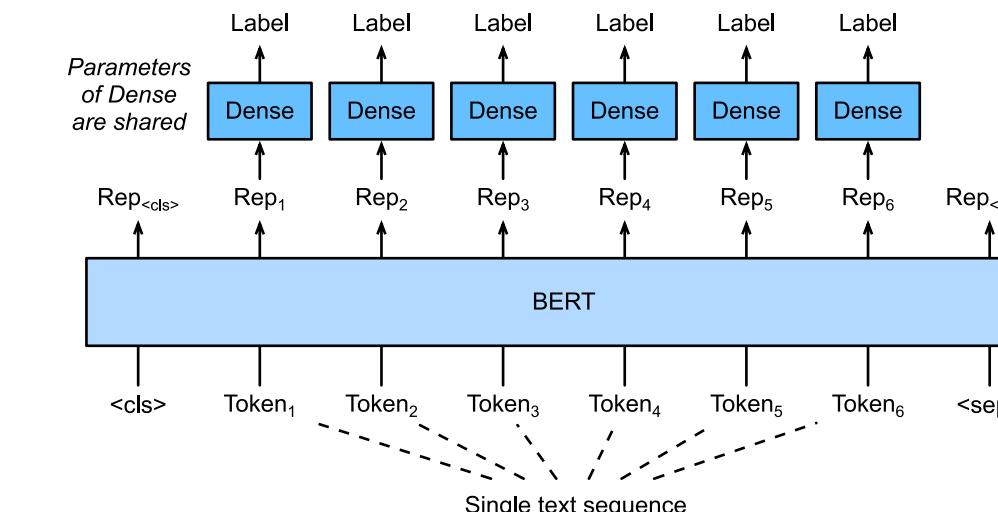
Sentence classification



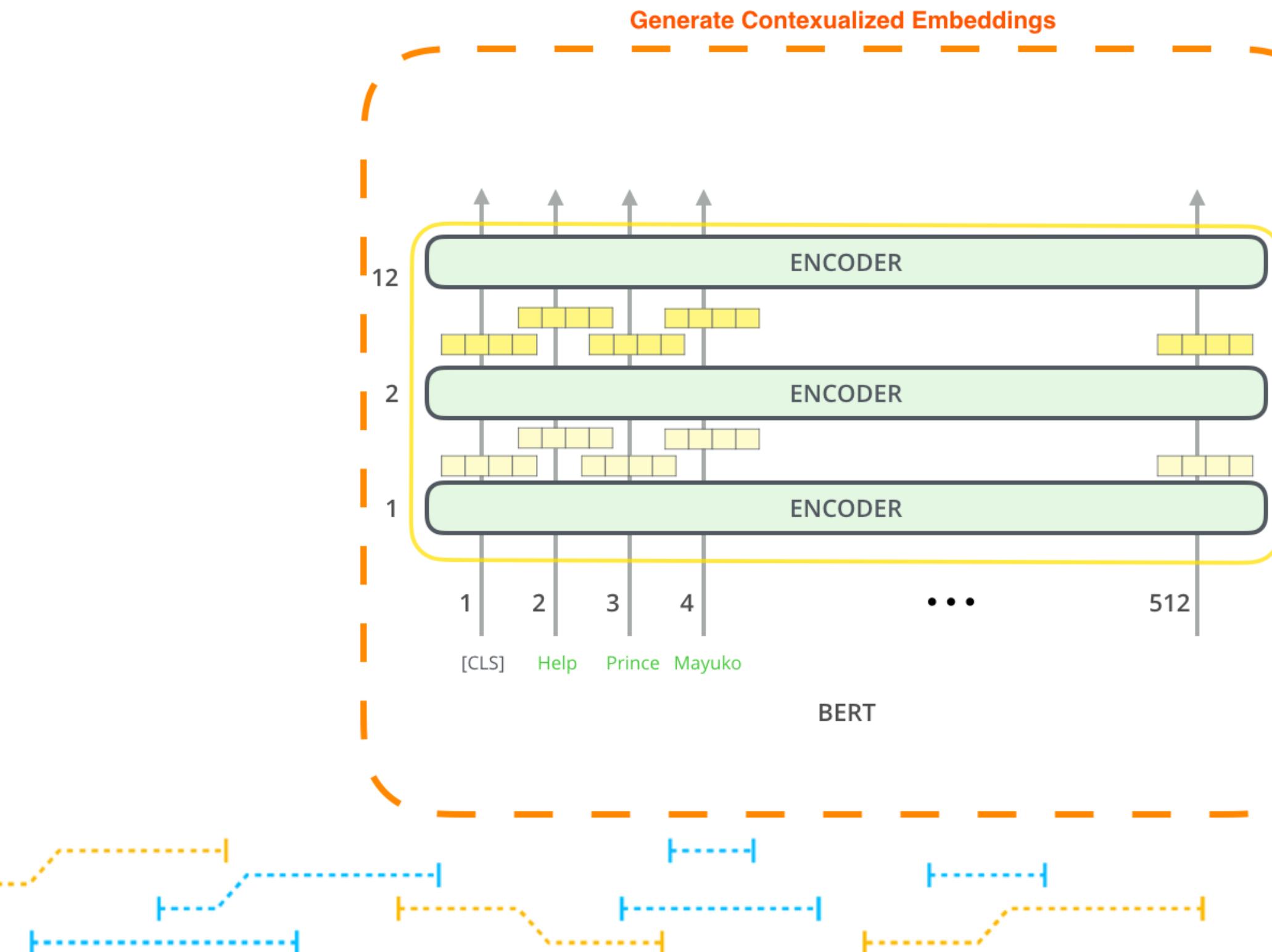
Answering multiple-choice questions



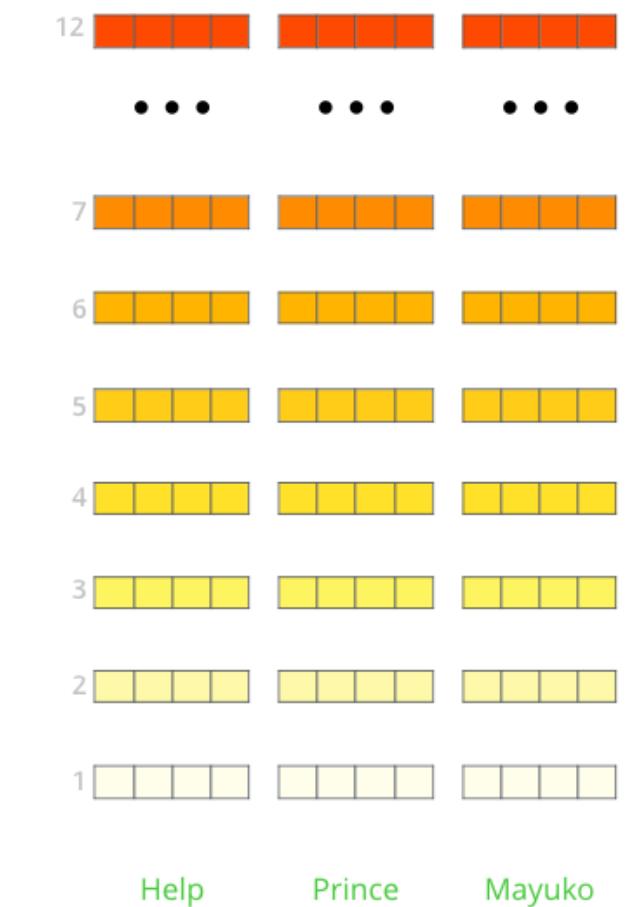
Part-of-speech tagging



BERT: Feature extraction

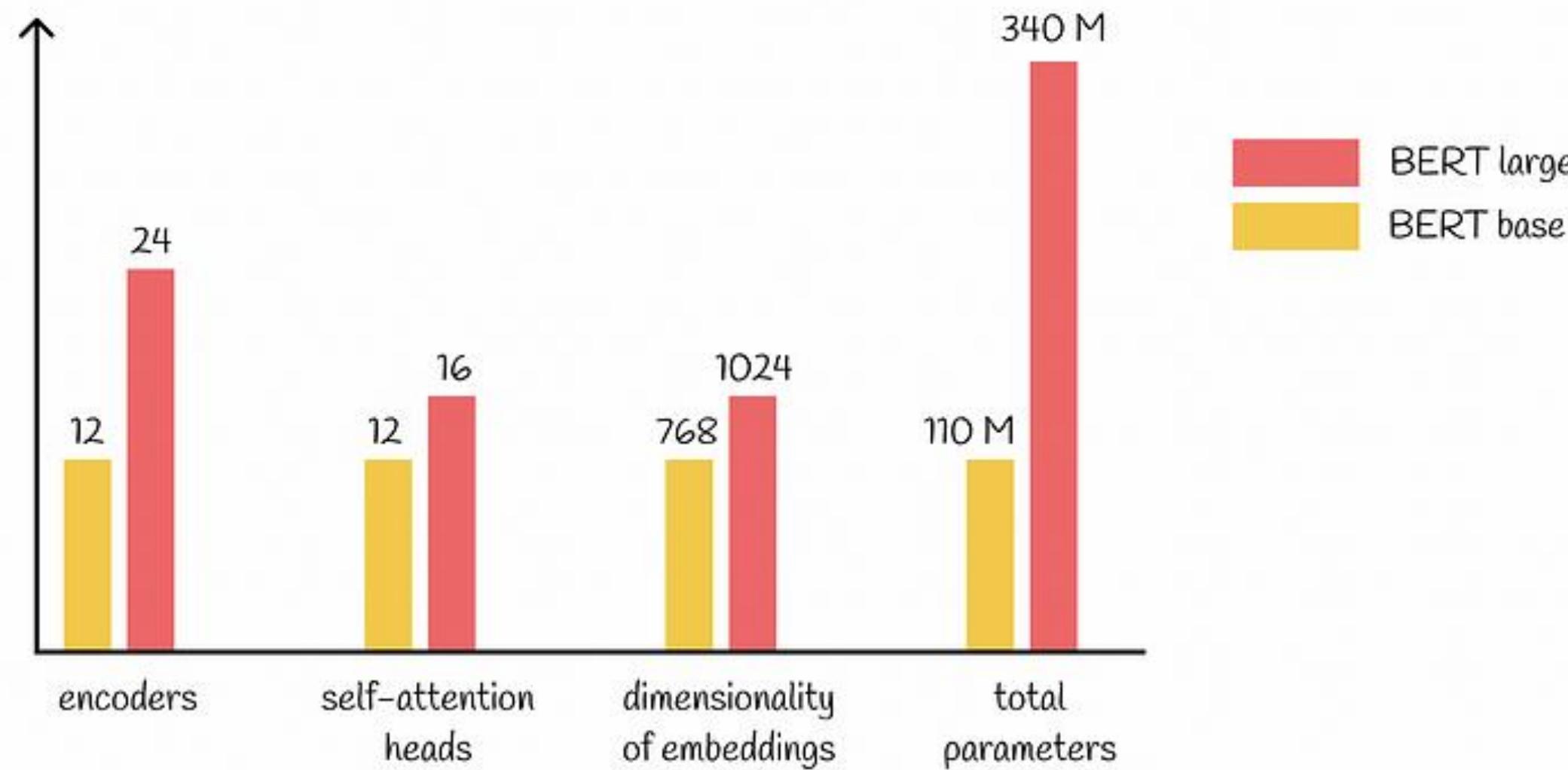


The output of each encoder layer along each token's path can be used as a feature representing that token.



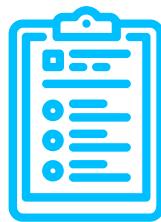
But which one should we use?

BERT



Jacob Devlin et al. (2018) "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding".
Proceedings of naacl-HLT (Vol. 1, p. 2).

5.



Transformer-XL

TRANSFORMATEC

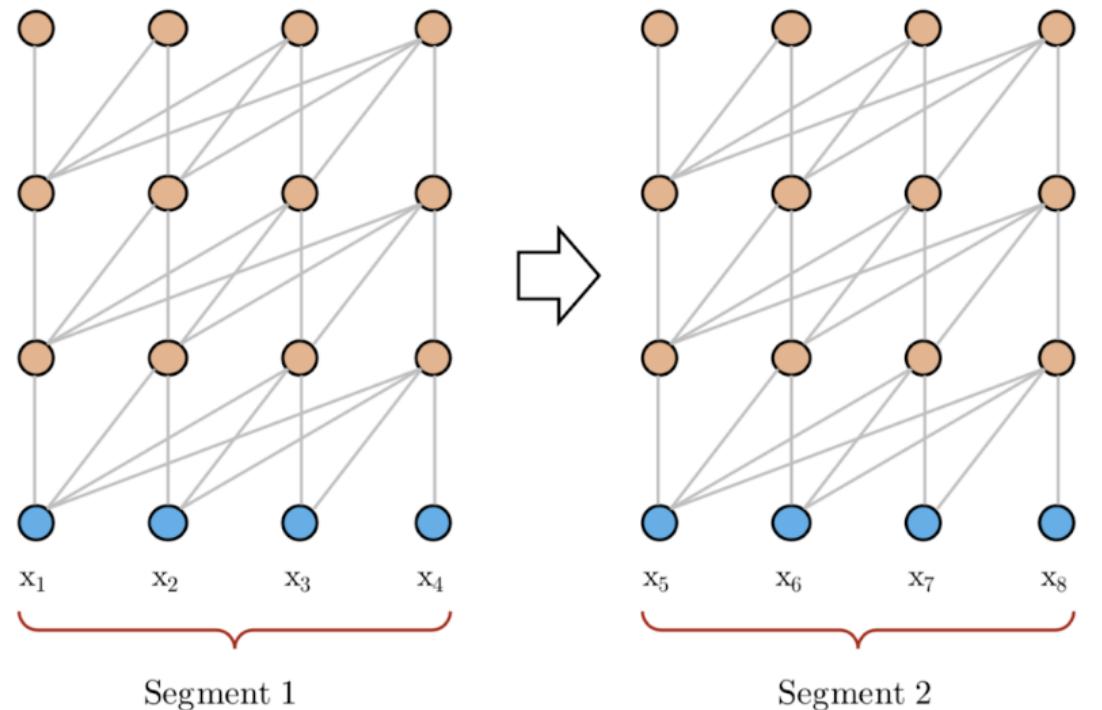
> Reinventa el mundo <



Transformer-XL

Ventana de contexto fija

Transformer (Training)

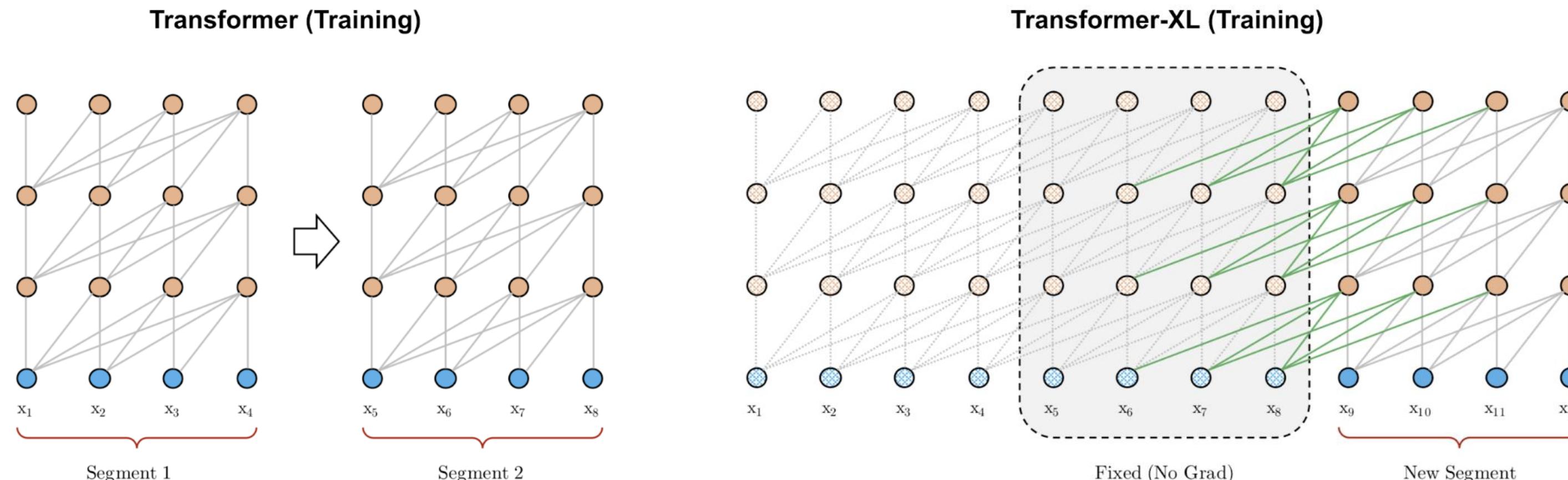


TRANSFORMATEC

Zihang Dai et al. (2018) "Transformer-XL: Language Modeling with Longer-Term Dependency".
International Conference on Learning Representations (ICRL).

Transformer-XL

Segment-Level Recurrence



Transformer-XL

Segment-Level Recurrence

$$h_t^l = [\text{detach}(M_{t-1}^l), h_t^l]$$

Atención extendida

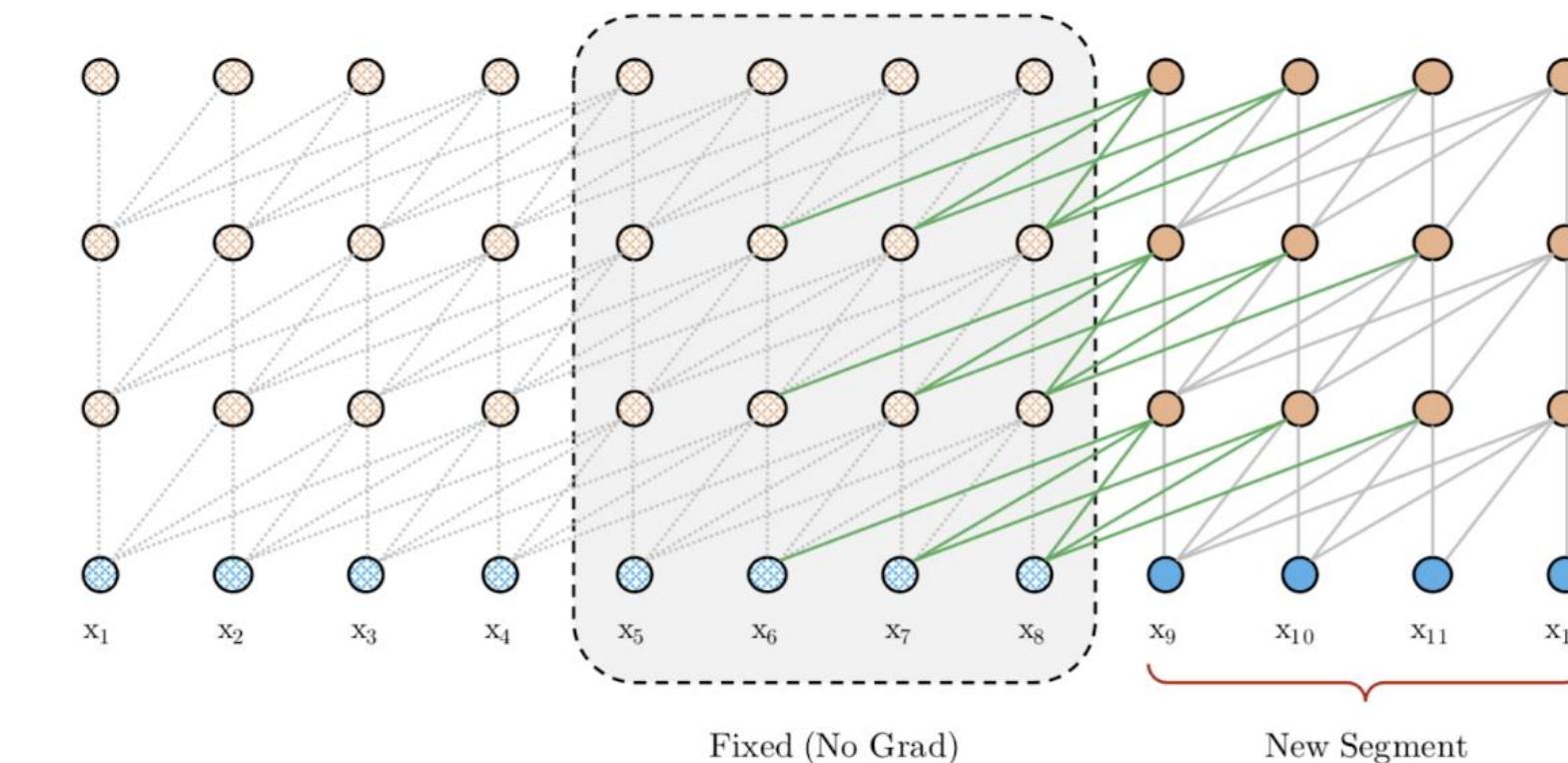
$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V$$

K y V Generados tanto por el contexto recurrente (memoria) como el segmento actual

Q Generados por la secuencia actual



Transformer-XL (Training)



Transformer-XL

Relative Positional Embedding (RPE)

$$A_{ij} = \frac{Q_i K_j^\top + Q_i r_{i-j}^\top + u K_j^\top + v r_{i-j}^\top}{\sqrt{d_k}}$$

- $Q_i r_{i-j}^\top$: Modela explícitamente la relación relativa de posición entre las entradas i (query) y j (key) mediante el embedding relativo r_{i-j} .
- $u K_j^\top$ y $v r_{i-j}^\top$: Es el **sesgo global** que capturan efectos posicionales independientes del query y key.
- r_{i-j} : Un embedding que representa la relación relativa de posición entre i y j .
- u, v : Vectores entrenables, que permiten al modelo ajustar los efectos globales de la posición relativa.



Transformer-XL

Relative Positional Embedding (RPE)

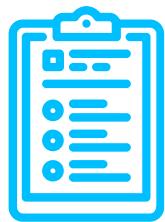
$$A_{ij} = \frac{Q_i K_j^\top + Q_i r_{i-j}^\top + u K_j^\top + v r_{i-j}^\top}{\sqrt{d_k}}$$

Algunos autores omiten los dos últimos términos, asumiendo que las dependencias globales no son relevantes.

Esta aproximación puede ser valida en secuencias cortas.



6.



PaLM

TRANSFORMATEC

> Reinventa el mundo <



PaLM

(Pathways Language Model)



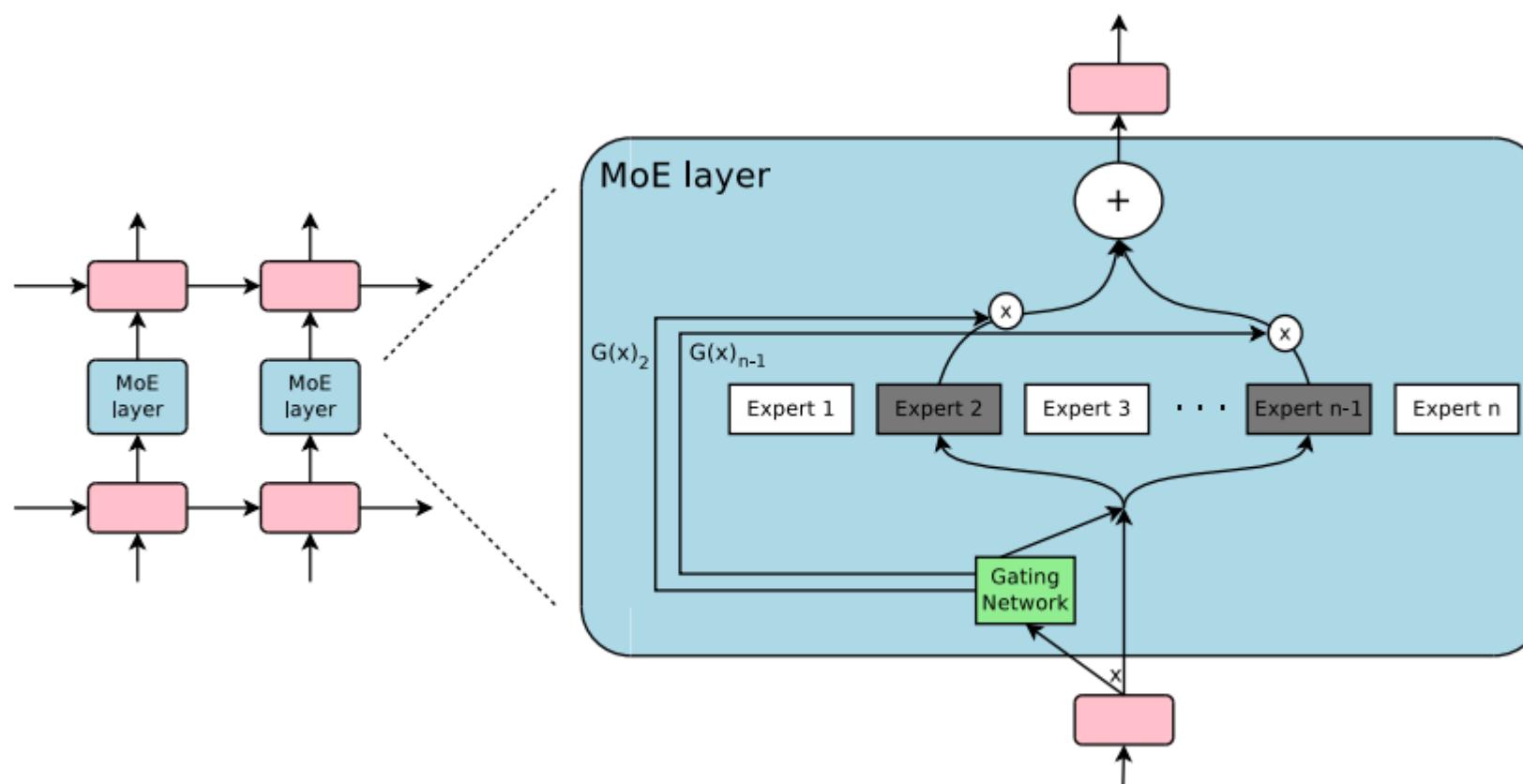
540 billion parameters



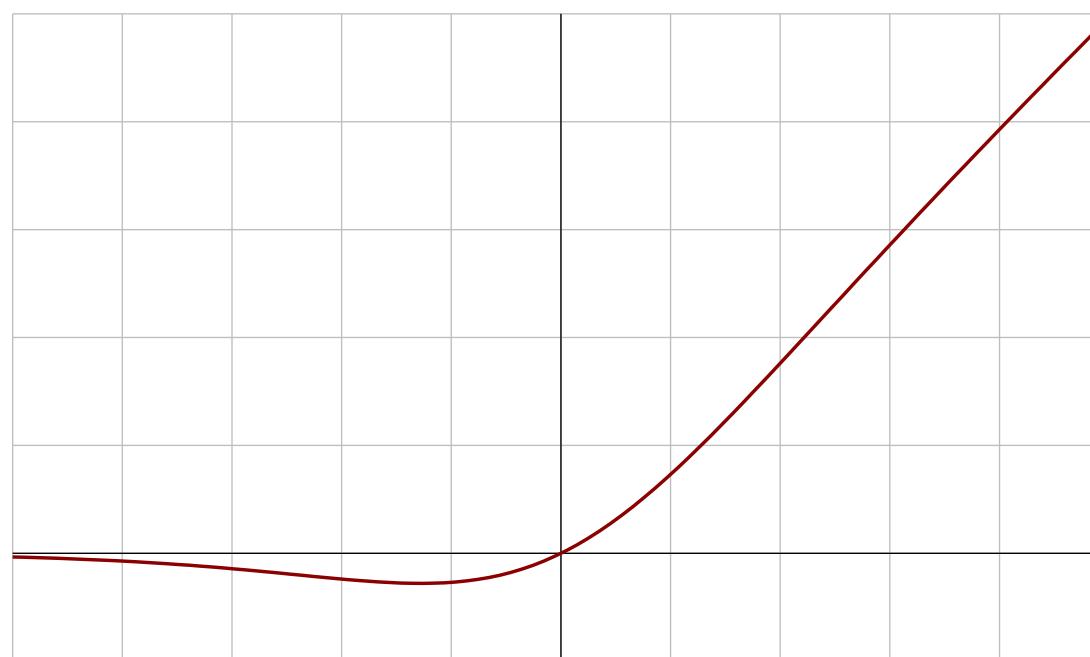
PaLM

(Pathways Language Model)

Mixture-of-Experts (MoE)



SwiGLU activation

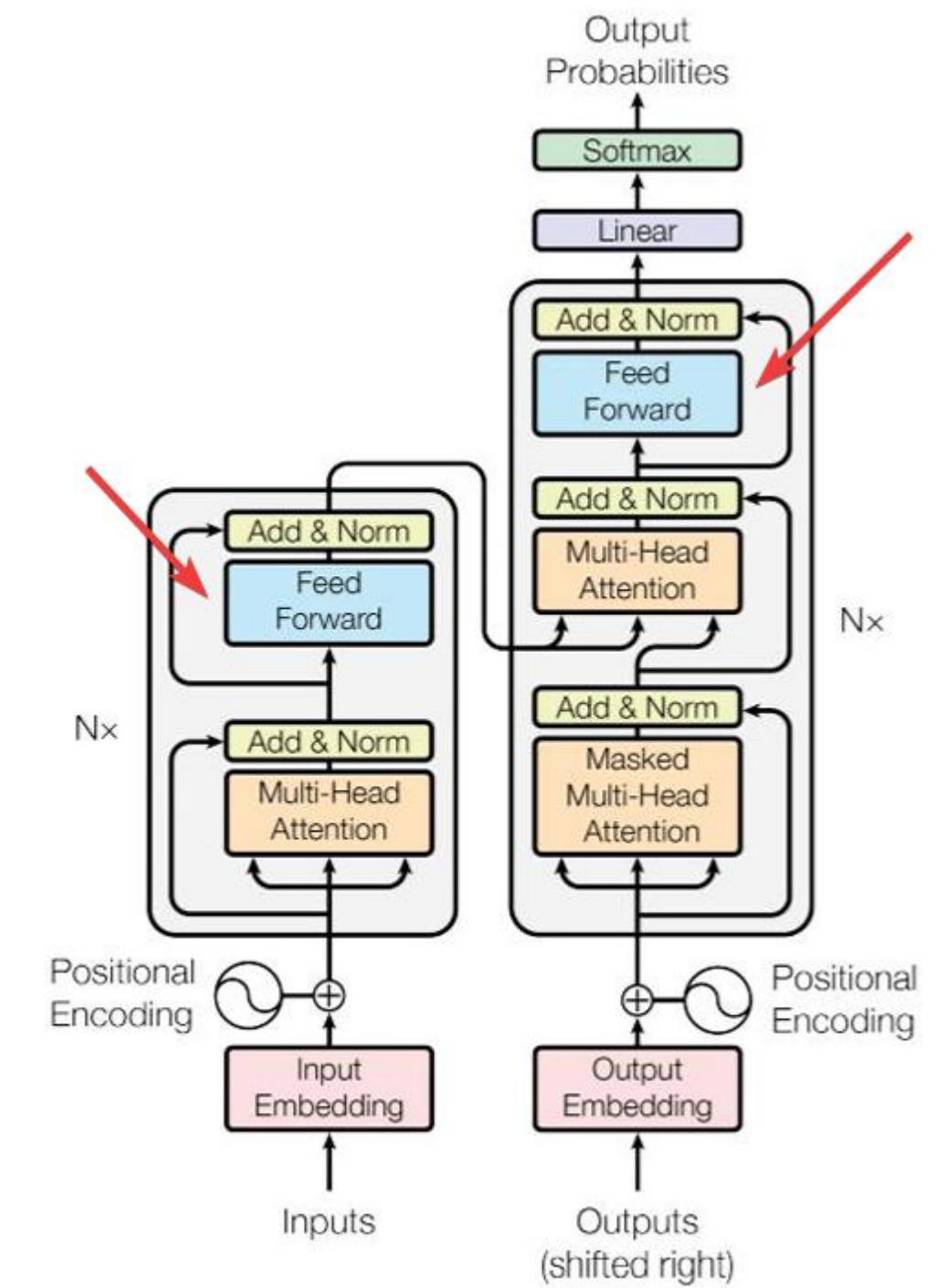


$$\text{Swish}(x) = x \text{ Sigmoid}(x)$$



PaLM

(Pathways Language Model)

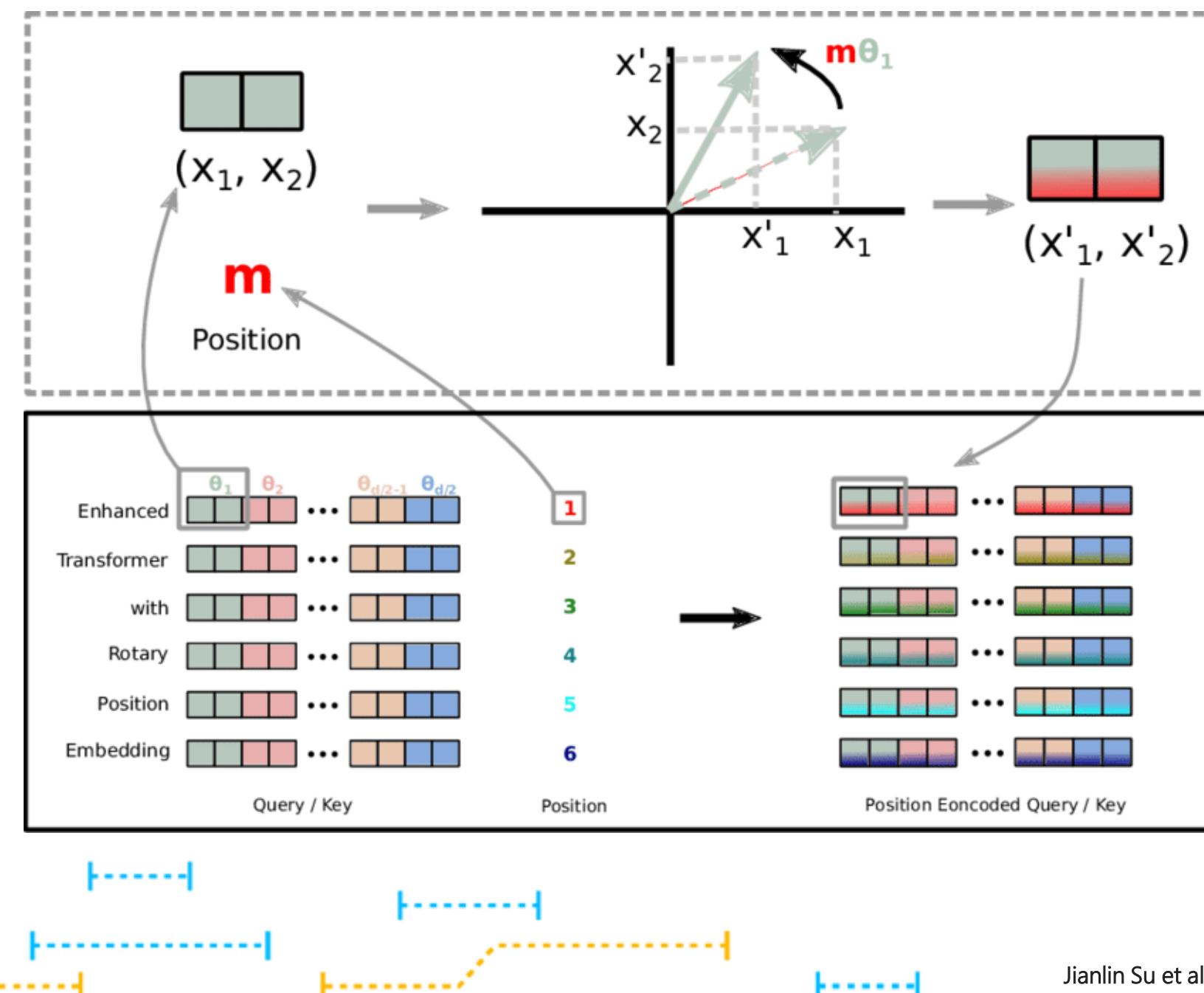


TRANSFORMATEC

PaLM

(Pathways Language Model)

RoPE (Rotary Positional Embeddings)



GRACIAS

Victor Flores Benites