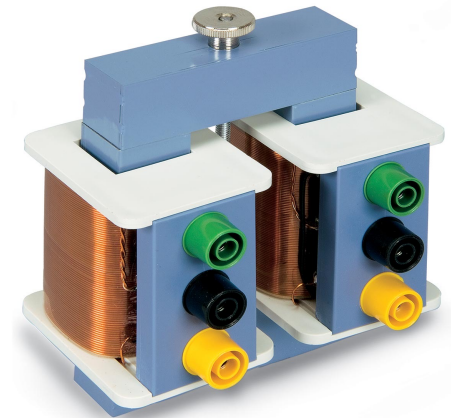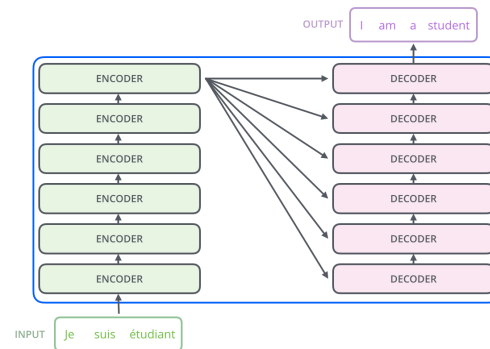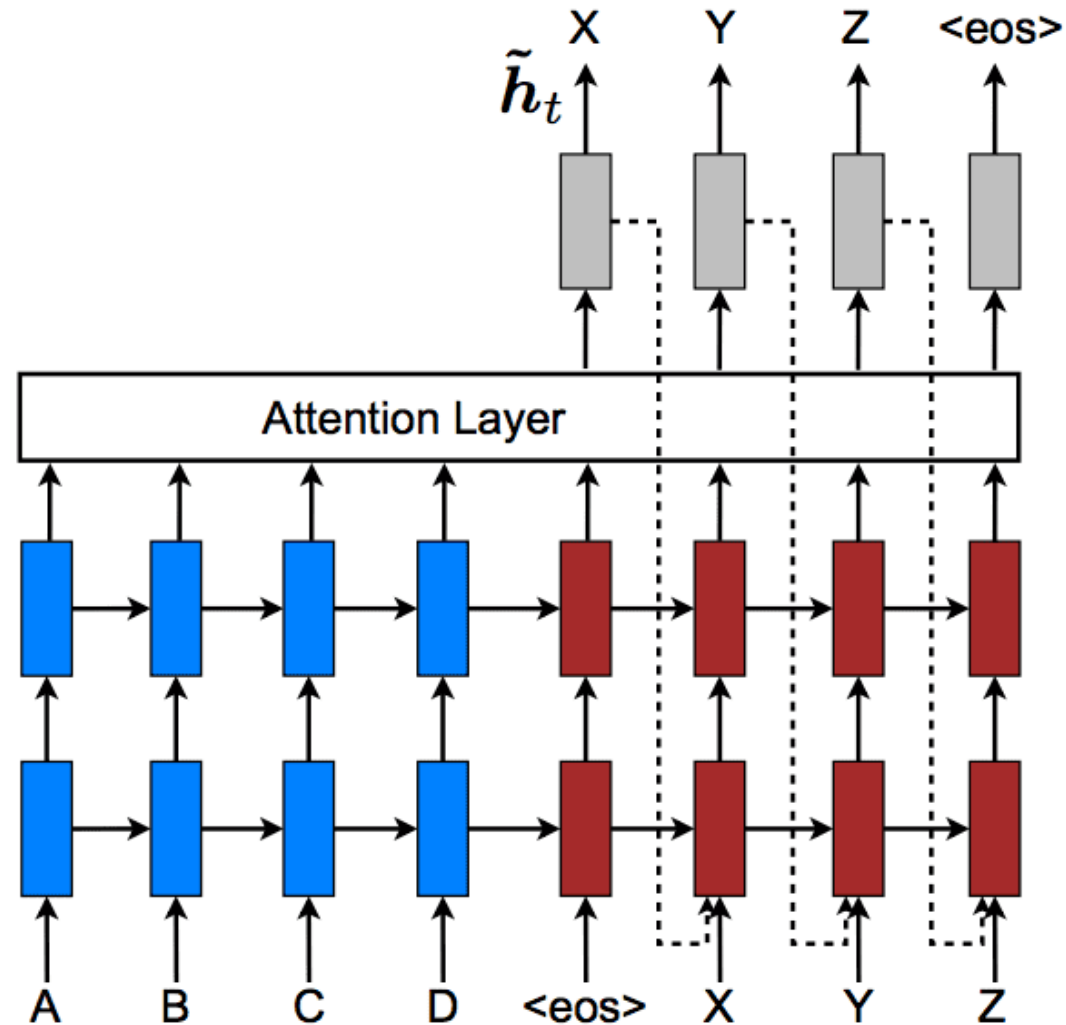# Transformers

Julio Waissman

# Antes de los transformers

# El artículo que cambió todo

## Attention Is All You Need

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
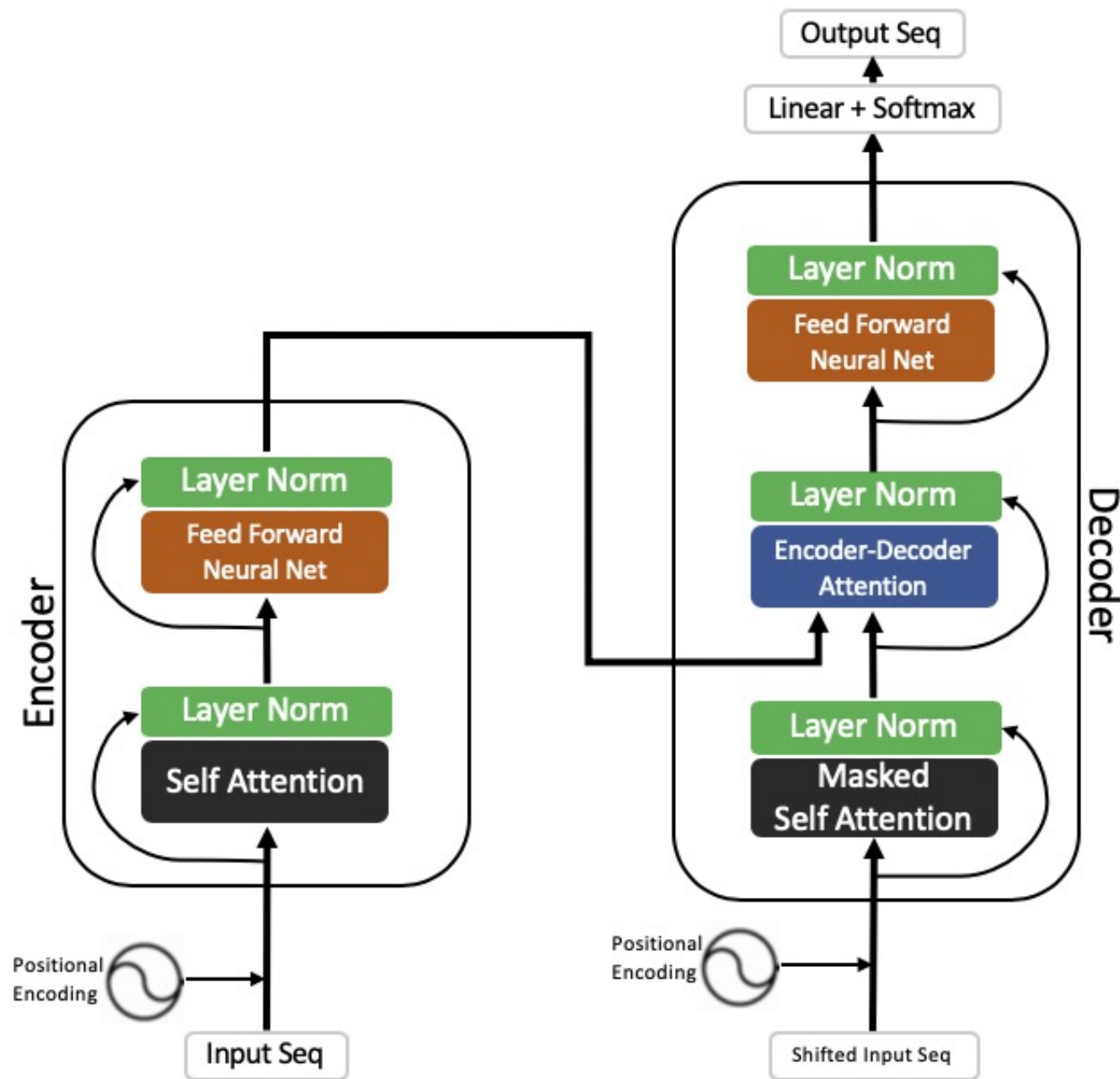Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

**Aidan N. Gomez**[*][†]
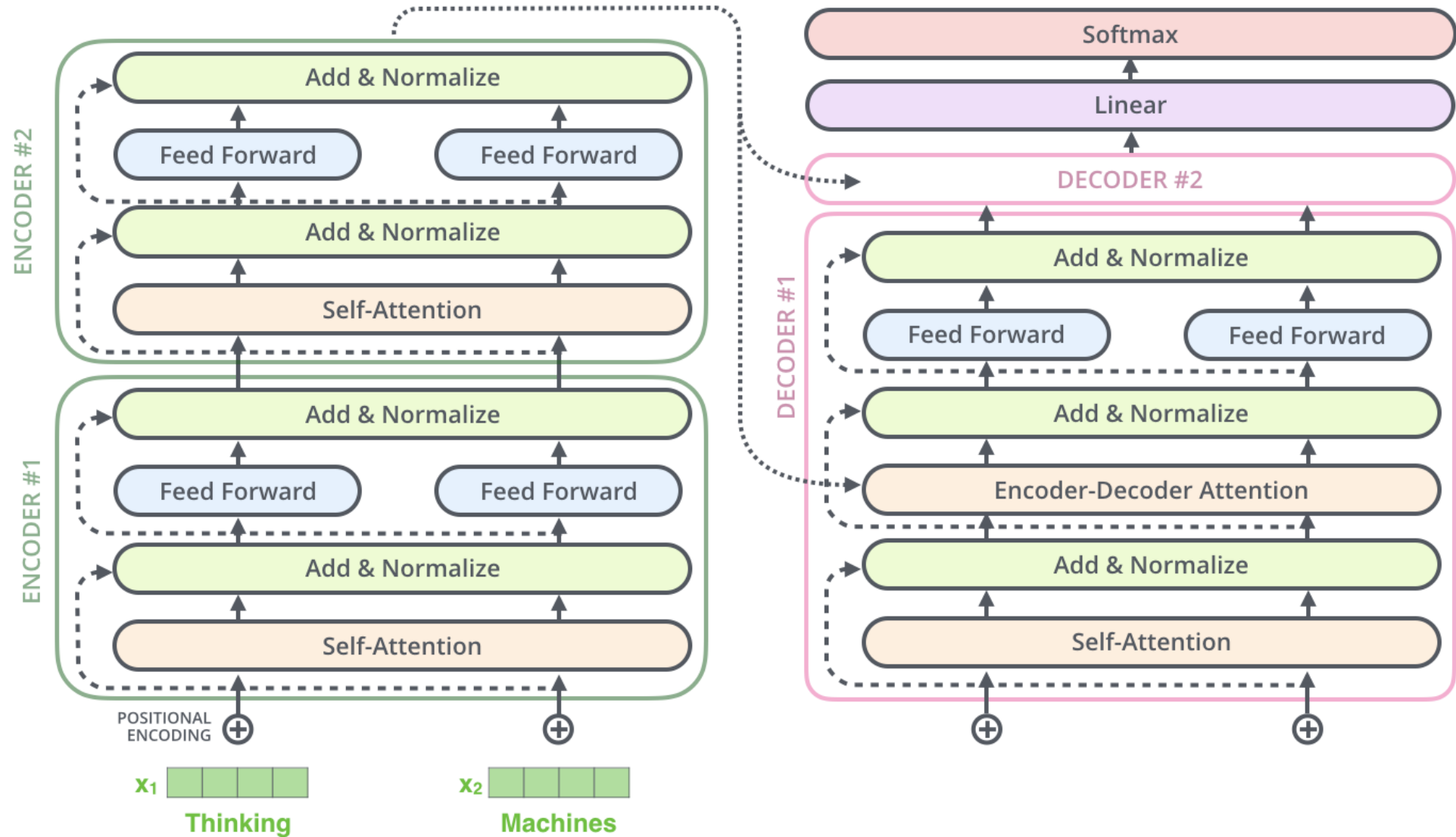University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[*][‡]
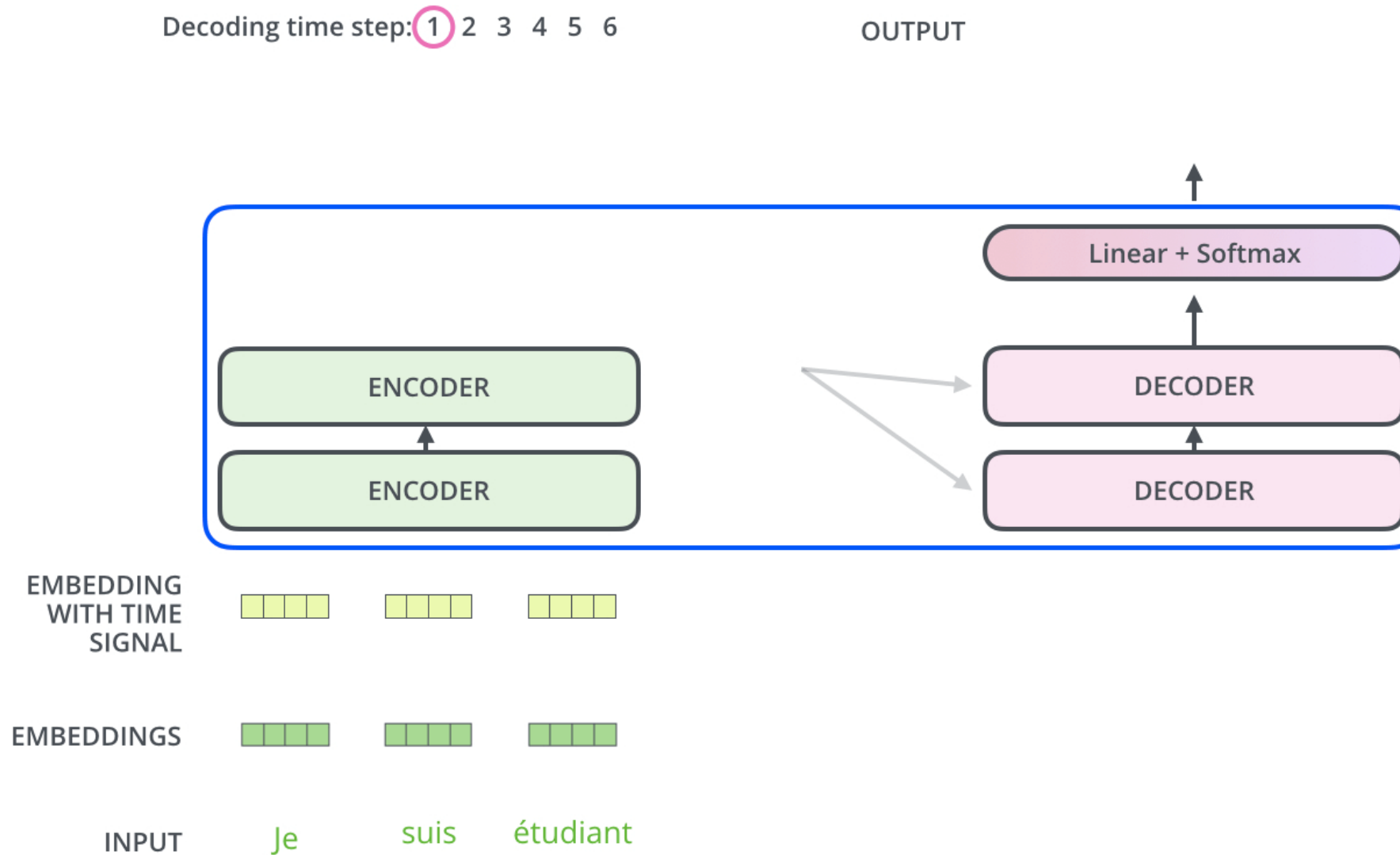illia.polosukhin@gmail.com
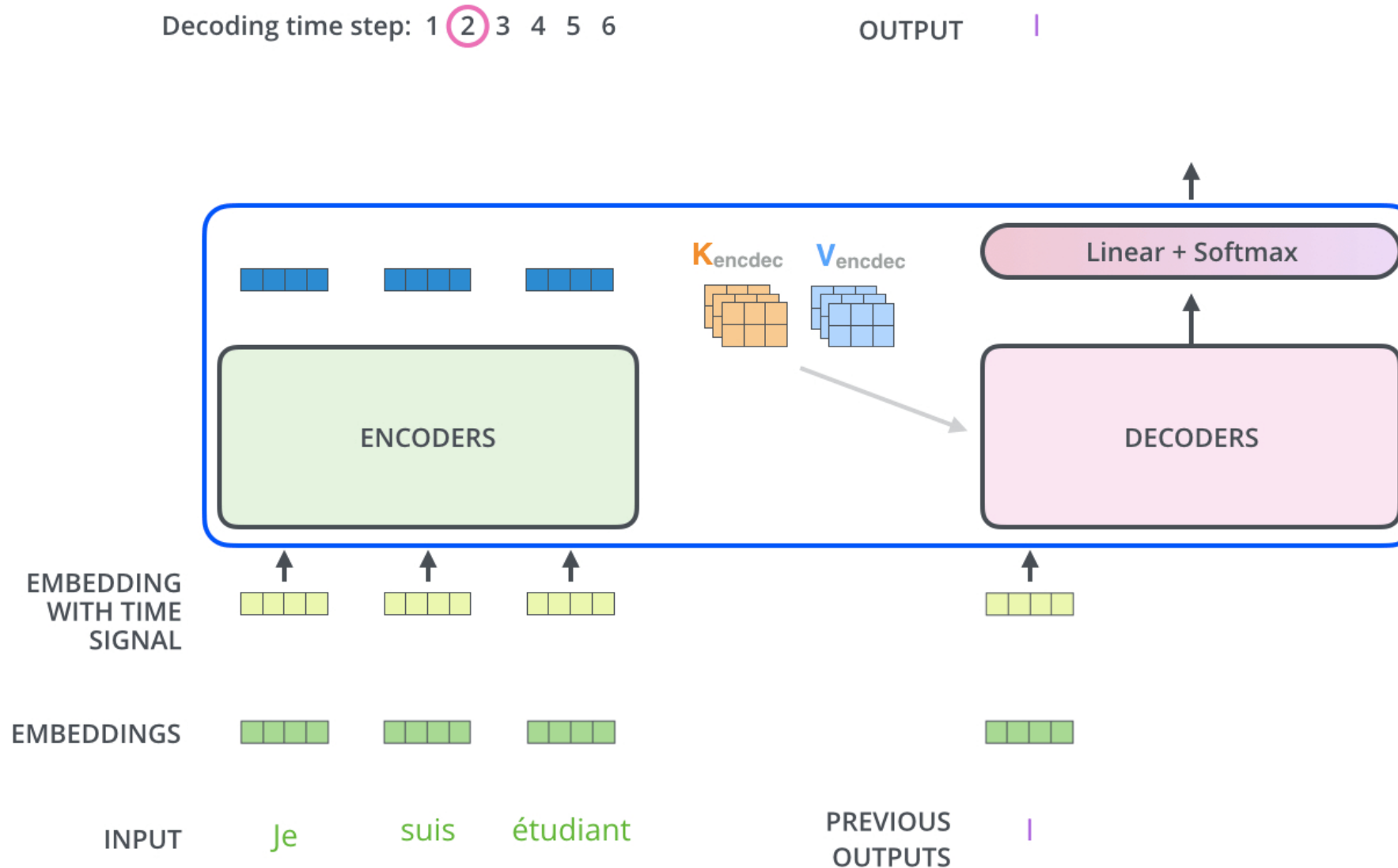
# Modelo Propuesto (Transformer)

# ¿Cómo funciona todo?

# ¿Cómo funciona todo?



Decoding time step: (1) 2  3  4  5  6

OUTPUT

Linear + Softmax

ENCODER

ENCODER

DECODER

DECODER

EMBEDDING
WITH TIME
SIGNAL

EMBEDDINGS

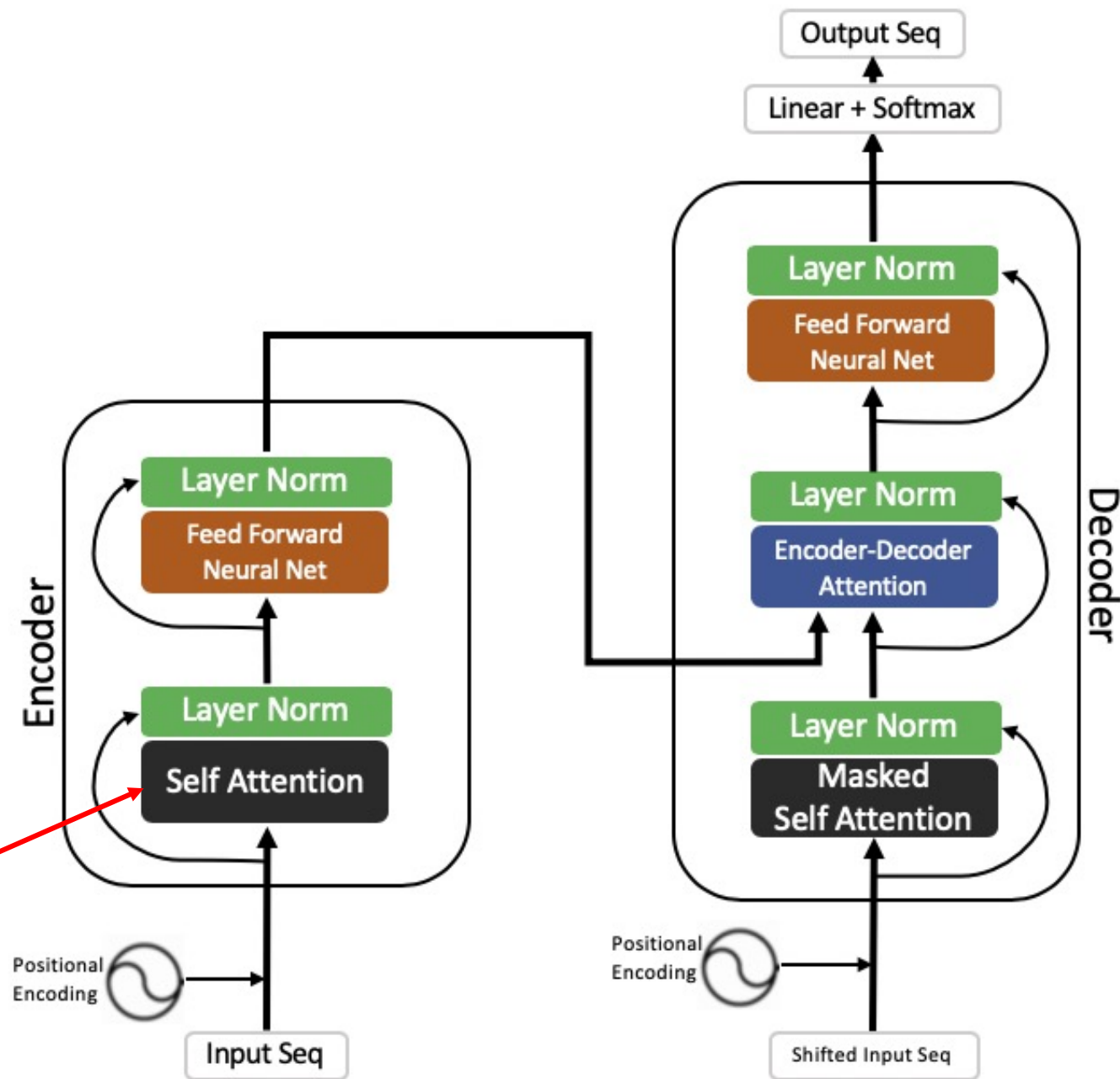INPUT     Je     suis    étudiant

# ¿Cómo funciona todo?

# Modelo Propuesto (Transformer)

La unidad básica

# Autoatención

Queries, keys and values come from the **same sentence**



Weight matrix

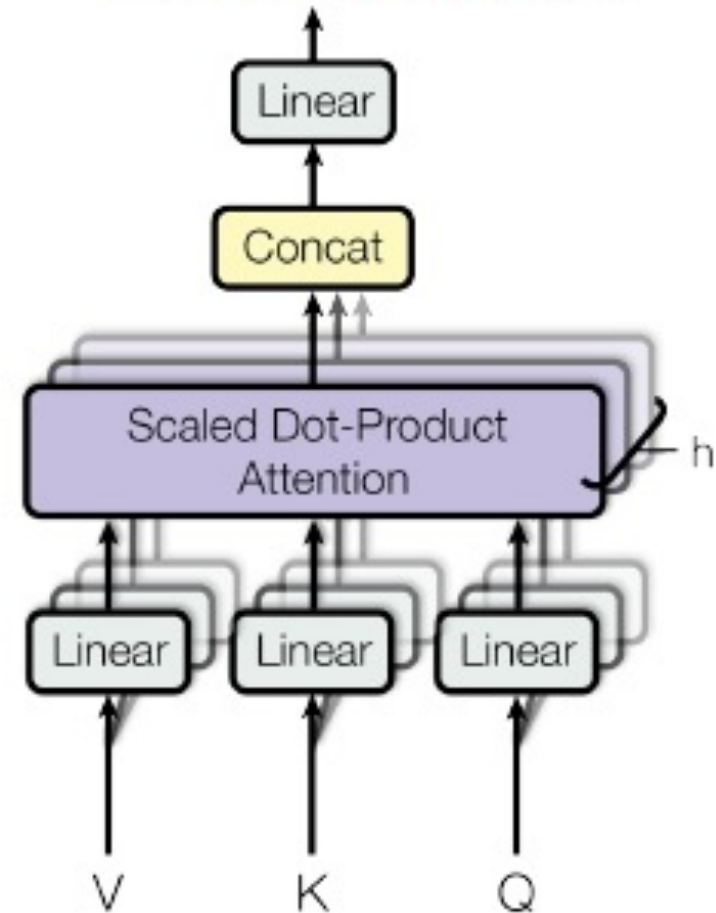Meaning of each word **within** the sentence

# Attention vs MultiHead Attention



Scaled Dot-Product Attention

Multi-Head Attention

# Self Attention

# Multi-head Self Attention

1) This is our input sentence*

2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

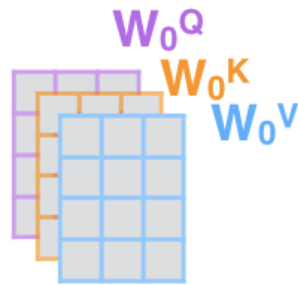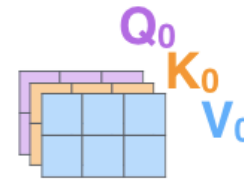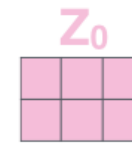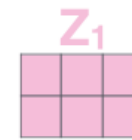4) Calculate attention using the resulting Q/K/V matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix $W^O$ to produce the output of the layer

Thinking Machines

$X$

$W_0^Q$
$W_0^K$
$W_0^V$

$Q_0$
$K_0$
$V_0$

$Z_0$

$W^O$

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

$W_1^Q$
$W_1^K$
$W_1^V$

$Q_1$
$K_1$
$V_1$

$Z_1$

$Z$

...

$R$

$W_7^Q$
$W_7^K$
$W_7^V$

$Q_7$
$K_7$
$V_7$

$Z_7$

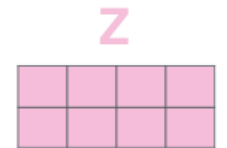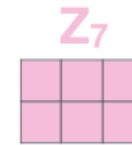# Modelo Propuesto (Transformer)

# Positional encoding

POSITIONAL ENCODING

| 0 | 0 | 1 | 1 |

| 0.84 | 0.0001 | 0.52 | 1 |

| 0.91 | 0.0002 | -0.42 | 1 |

**+**

EMBEDDINGS

INPUT

Je

suis

content

RNNs vs Transformer:  Positional Encoding

Modelo Propuesto (Transformer)

Deep Residual Learning for Image Recognition

El truco de las redes residuales para tener muchas capas de profundidad

Output Seq

Linear + Softmax

Layer Norm

Feed Forward Neural Net

Layer Norm

Encoder-Decoder Attention
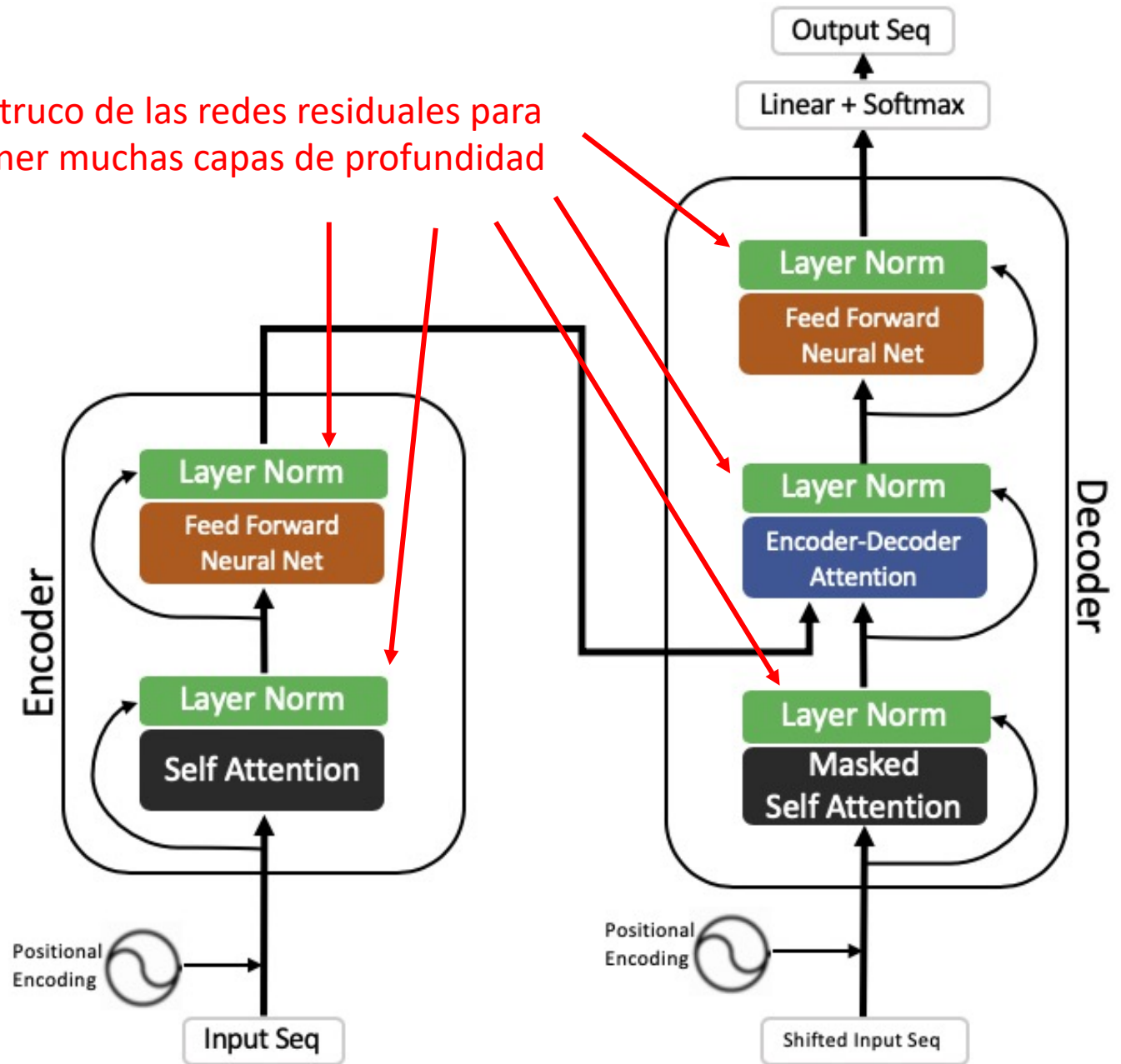
Layer Norm

Masked Self Attention

Decoder

Positional Encoding

Shifted Input Seq

Layer Norm

Feed Forward Neural Net

Layer Norm
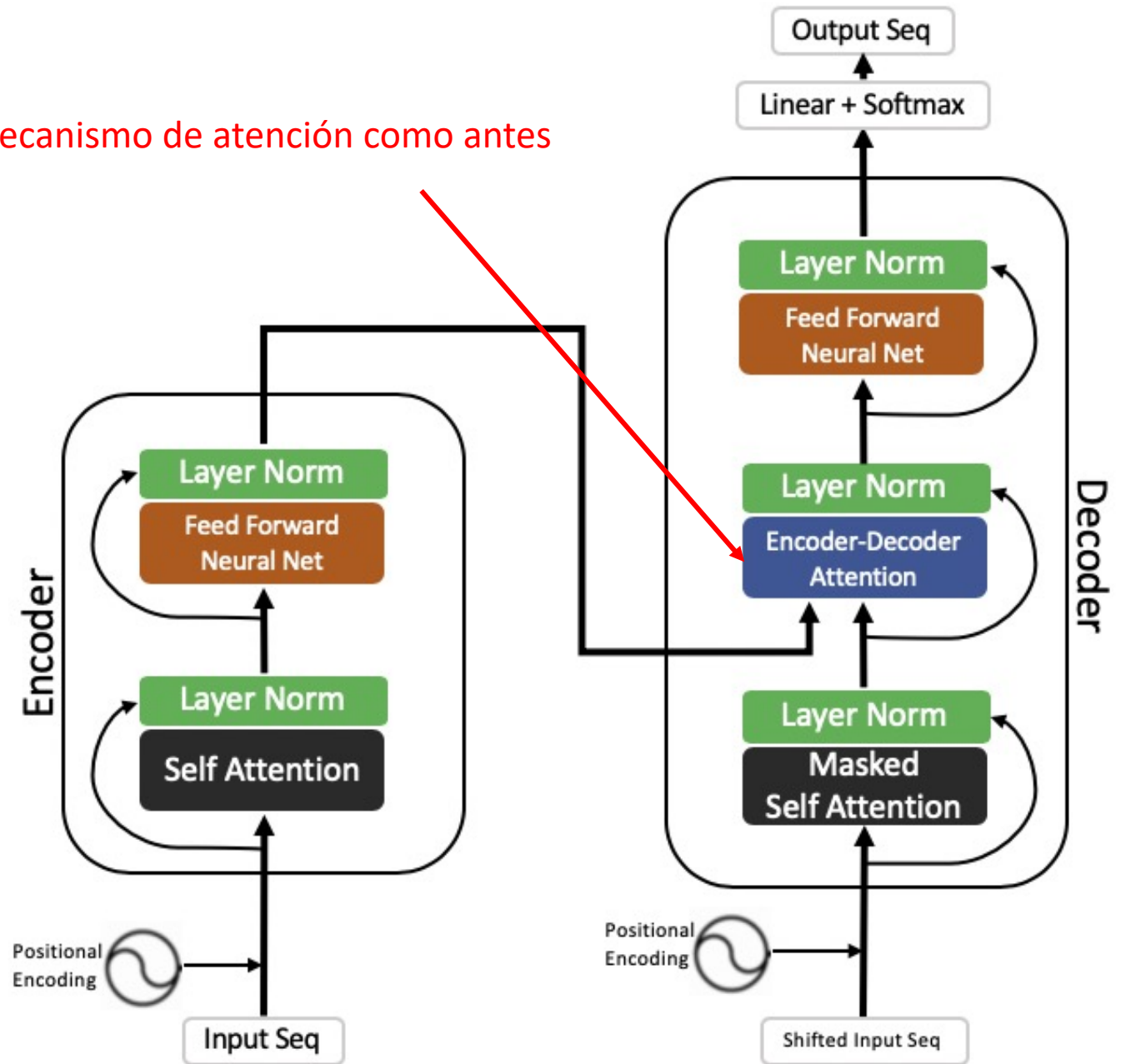
Self Attention

Encoder

Positional Encoding

Input Seq

# Add & Normalize

Modelo Propuesto (Transformer)

Mecanismo de atención como antes

# Encode-Decoder Self-Attention

Decoding time step: 1 (2) 3 4 5 6    OUTPUT    I

K<sub>encdec</sub>  V<sub>encdec</sub>

Linear + Softmax

ENCODERS

DECODERS

EMBEDDING WITH TIME SIGNAL

EMBEDDINGS

INPUT    Je    suis    étudiant
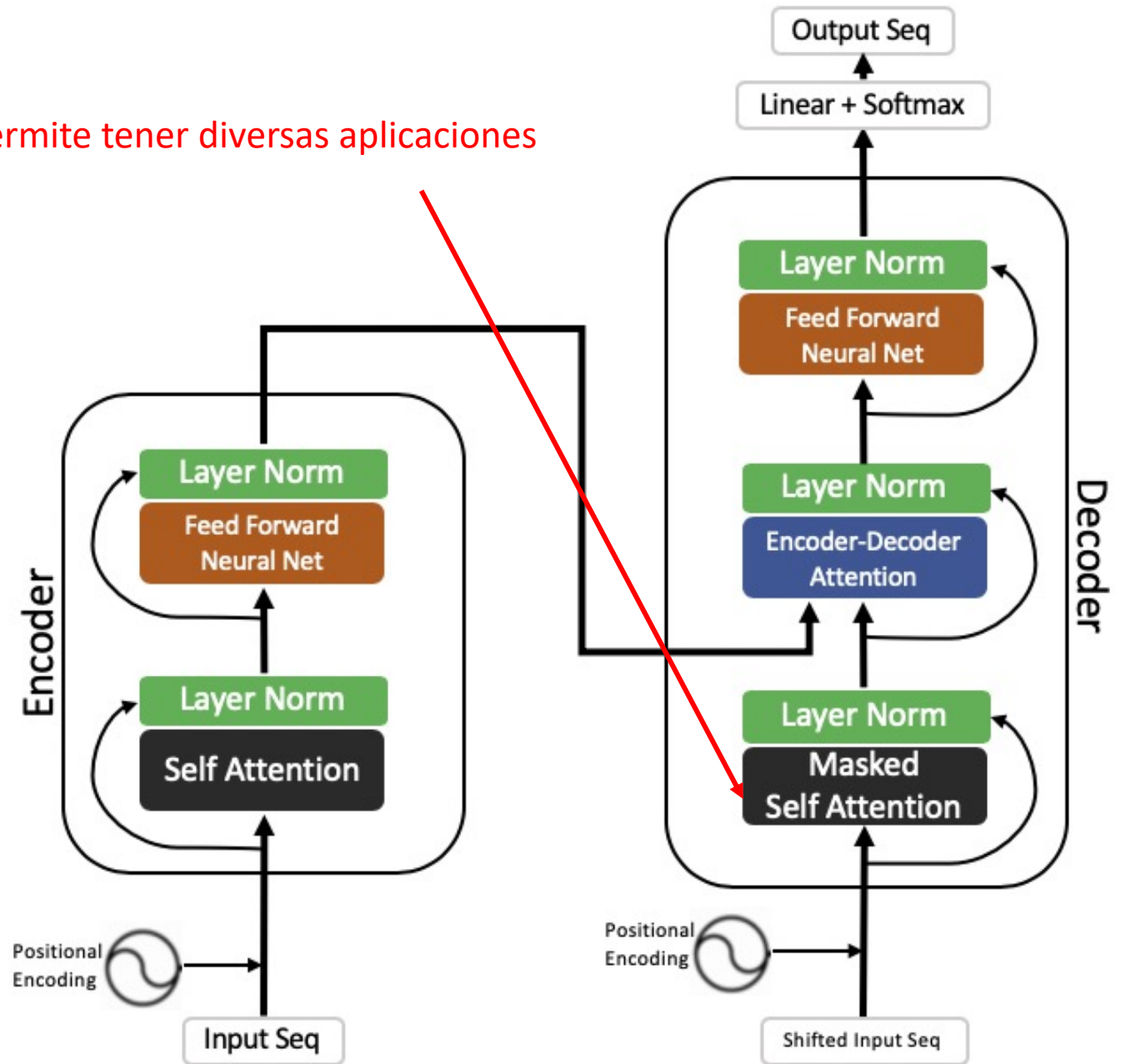
PREVIOUS OUTPUTS    I

# Encoder-Decoder Attention

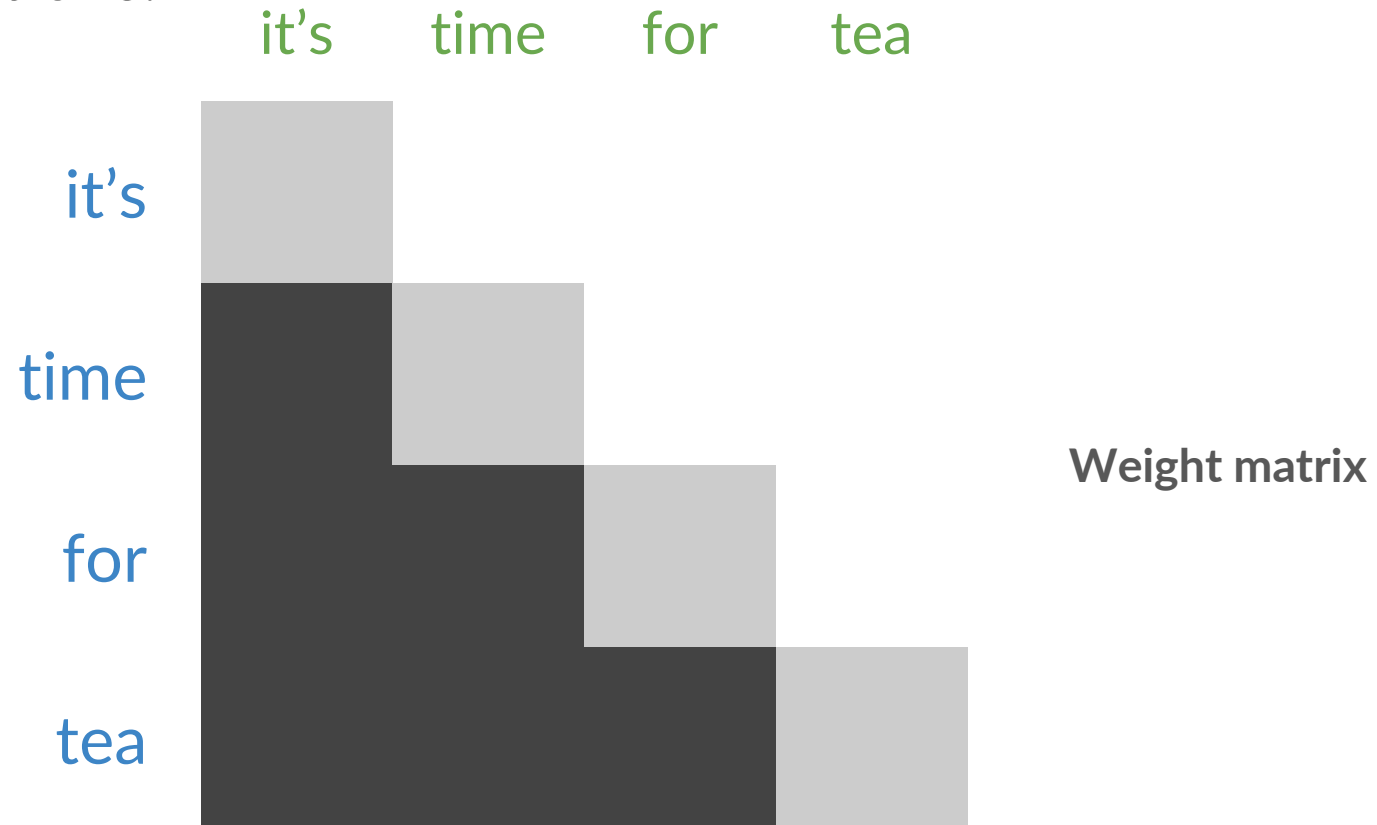Queries from one sentence, keys and values from another



Weight matrix

Modelo Propuesto (Transformer)

Permite tener diversas aplicaciones
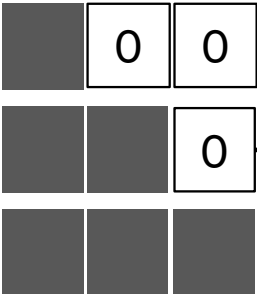
# Masked Self-Attention

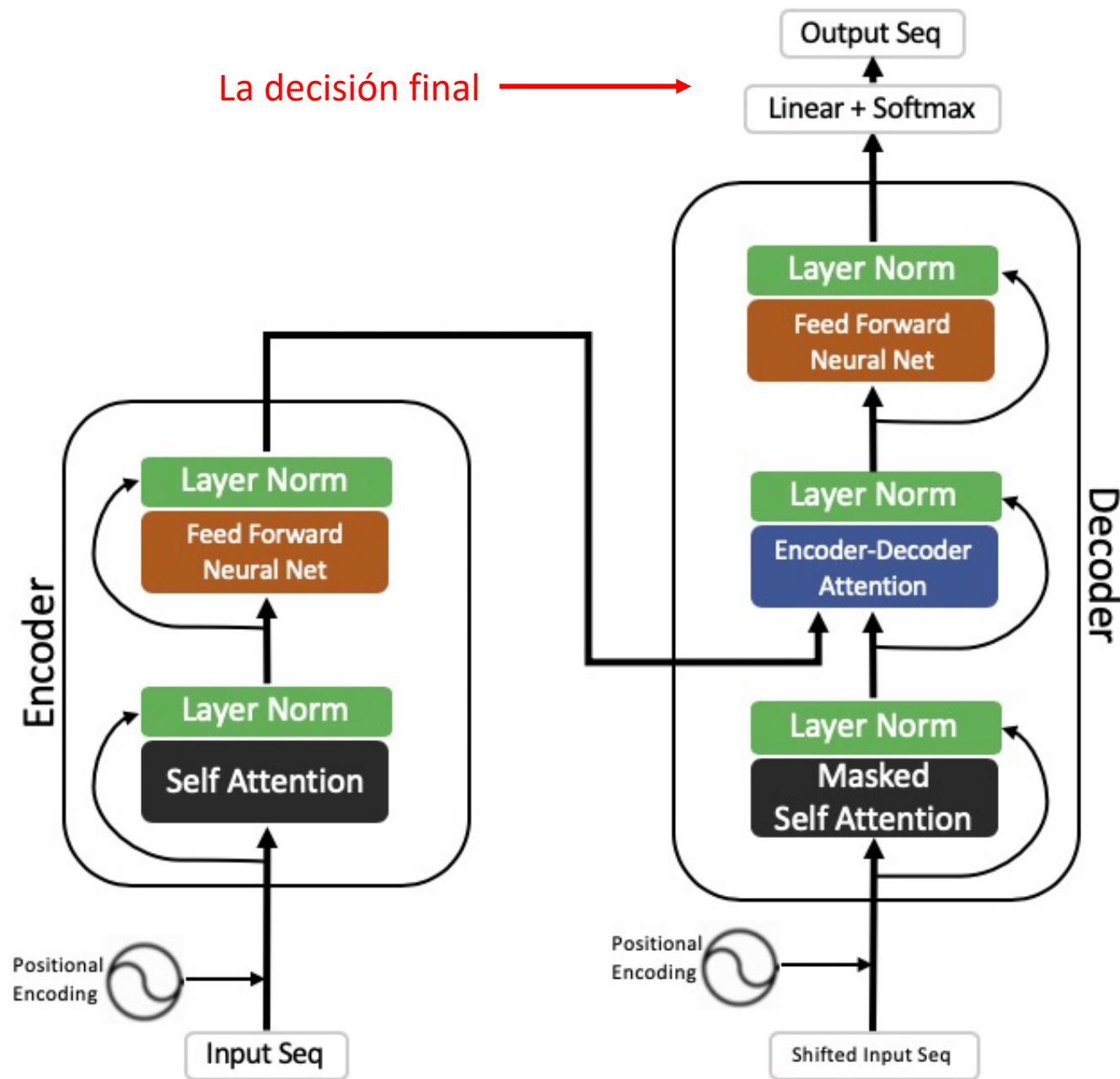Queries, keys and values come from the **same sentence**. Queries don't attend to future positions.



Weight matrix

# ¿Y cómo se hace el Masked Self-Attention?



Minus infinity

$$\mathrm{softmax}\left(\frac{Q\,K^{\top}}{\sqrt{d_k}} + \begin{matrix} 0 \\ 0 & 0 \\ 0 & 0 & 0 \end{matrix}\right) V$$

Weights assigned to future positions are equal to 0

Modelo Propuesto (Transformer)

# Capa final de los decoders

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

0 1 2 3 4 5 ... vocab_size

Softmax

logits

0 1 2 3 4 5 ... vocab_size

Linear

Decoder stack output