# IMPORTANT

Spanish → Transformer → output English

Je suis étudiant → Encode → Decode → stack of decoders.

Decode → I am a student

Encoder
  └→ is actually a stack of encoders, no. 6
        └→ hyperparameter
              ├ encoders: 6
              └→ decoders: 6

output of final encoder will be input to all decoders.

Encode → Decoder
Encoder → Decode
Encode → Decoder
Encode → Decode
Encode → Decode
Encode → Decode
Encoder

↑
Input L...

# Inside Each encoder

### Encoder



Encoder's inputs first flow through a self attention layer — a layer that helps encoder look at other words in the input sentence as it encodes a specific word.

① Output of the self-attention layer are fed to feed forward NN. The exact same feed forward Net NN is applied to each position

### Decoder



↳ this second attention layer helps the decoder focus on relevant parts of the input sentence.
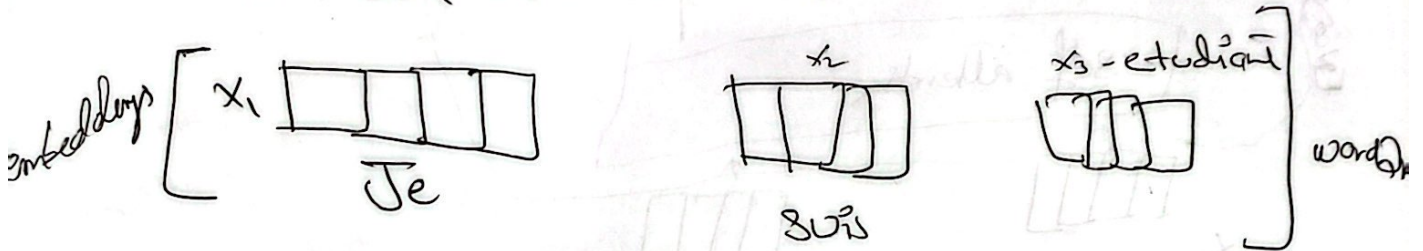
# Attend

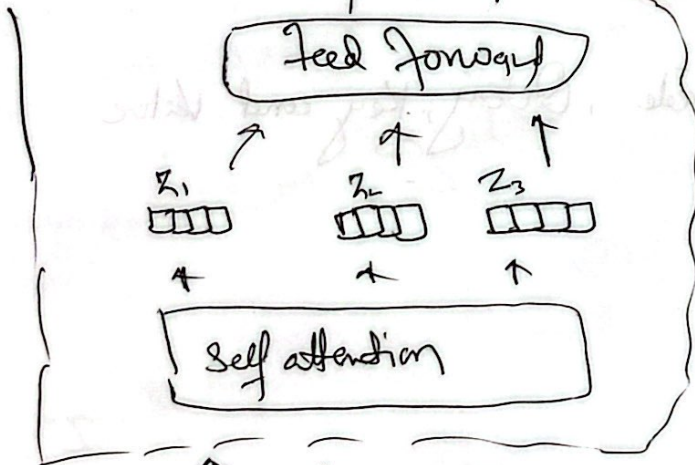## Attention v/s Self ATTENTION

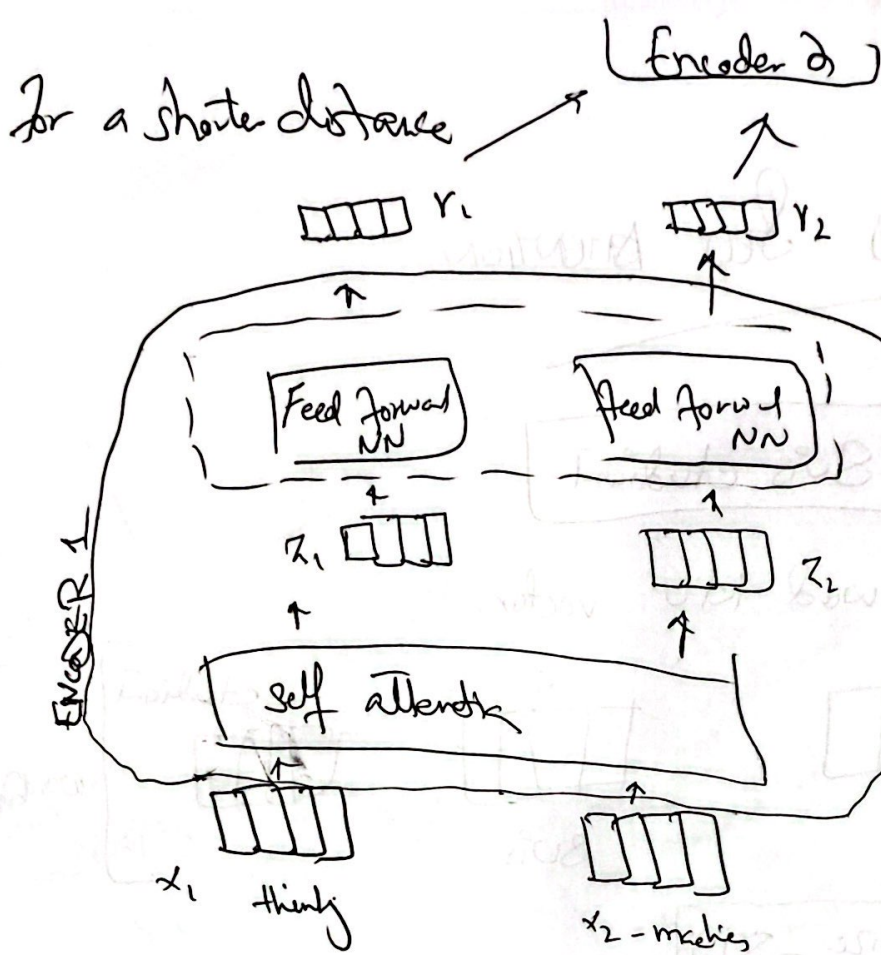input: | Je Sois etudiant |

turn each word into vector.

embeddings $\left[\begin{array}{c} x_1 \quad \square\square\square\square \\ Je \end{array} \quad \begin{array}{c} x_2 \\ \square\square\square \\ Suis \end{array} \quad \begin{array}{c} x_3 - etudiant \\ \square\square\square \end{array}\right]$ words

| vector size - 512 |

there vectors are passed to self attention layer

ENCODER

output will be gij to encoder 2.



Feed forward

$z_1$ $\square\square\square\square$   $z_2$ $\square\square\square$   $z_3$ $\square\square\square$

Self attention

$x_1$ $\square\square\square\square$   $x_2$ $\square\square\square\square$   $x_3$ $\square\square\square$

Je        Sois        etudiant

All embeddings are given concurrently

For a shorter distance → Encoder 2)

☐☐☐☐ $Y_1$          ☐☐☐☐ $Y_2$

ENCODER 1

| Feed Forward NN | Feed Forward NN |

$Z_1$ ☐☐☐☐          ☐☐☐☐ $Z_2$

Self attention

☐☐☐☐          ☐☐☐☐

$X_1$  thinly          $X_2$ - machine

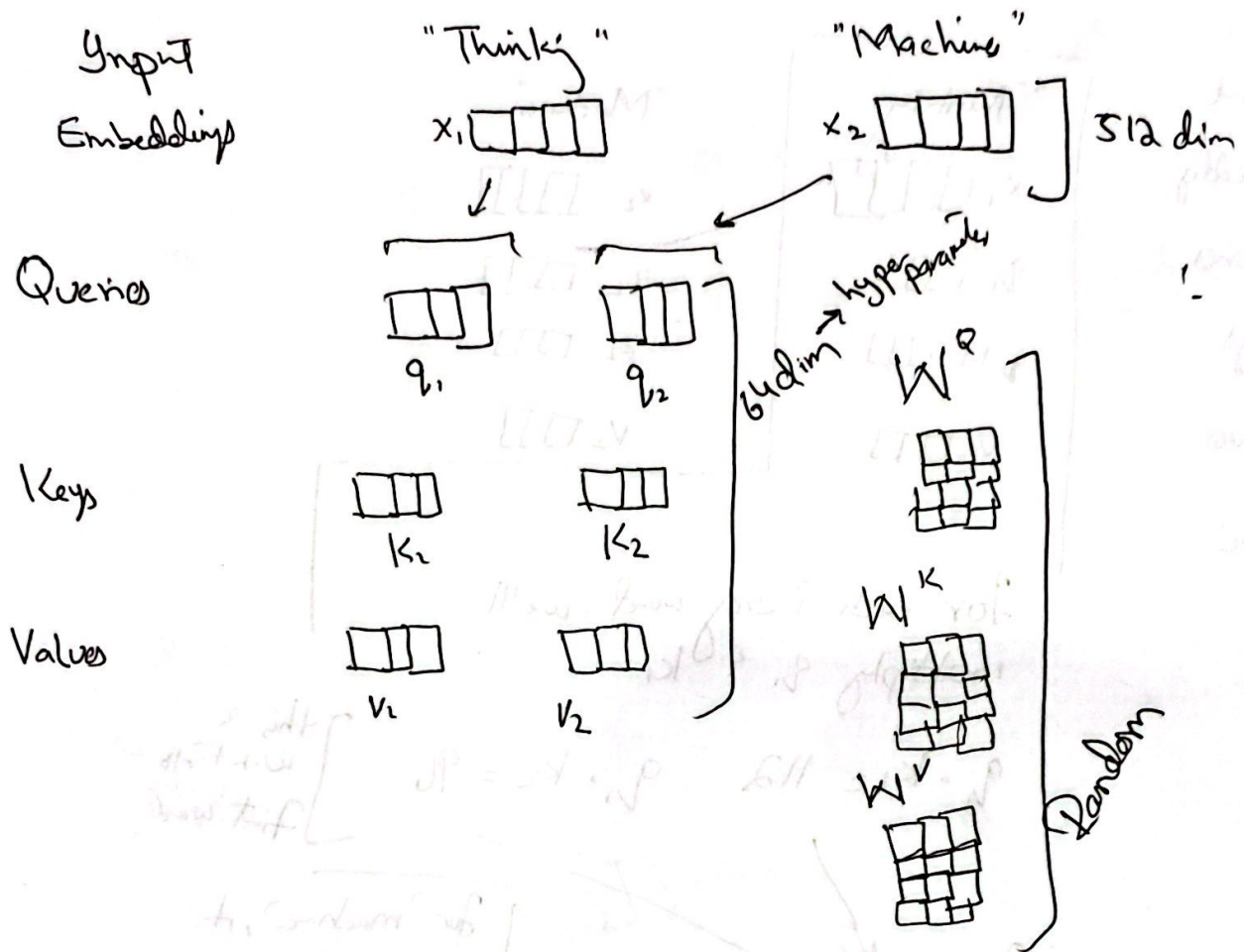## Self - attention,

but how will it give impol relate 'it' with animal?

For each word we create, Query, Key and Value vector

Attention 3

Input
Embeddings

"Thinky"  $x_1$ ☐☐☐    "Machine"  $x_2$ ☐☐☐ ] 512 dim

Queries

☐☐☐  ☐☐☐
$q_1$     $q_2$

64 dim → hyperparam

$W^Q$ ☐☐☐

Keys

☐☐ $K_1$   ☐☐ $K_2$

$W^K$ ☐☐☐

Values

☐☐☐ $V_1$   ☐☐ $V_2$

$W^V$ ☐☐☐

Random

$W^Q, W^K, W^V$ are 3 weights we randomly initialized. It will be updated based on back propagh.

So, first

$$q_1 = x_1 \cdot W^Q \qquad q_2 = x_2 \cdot W^Q \Big] \; 64\text{-}d$$

why multiplied? input multiplied by weight + bias
$Wx+b$, this is active
this how ANN operate.

Same for $W^K$ & $W^V$

$$K_1 = x_1 \cdot W^K$$
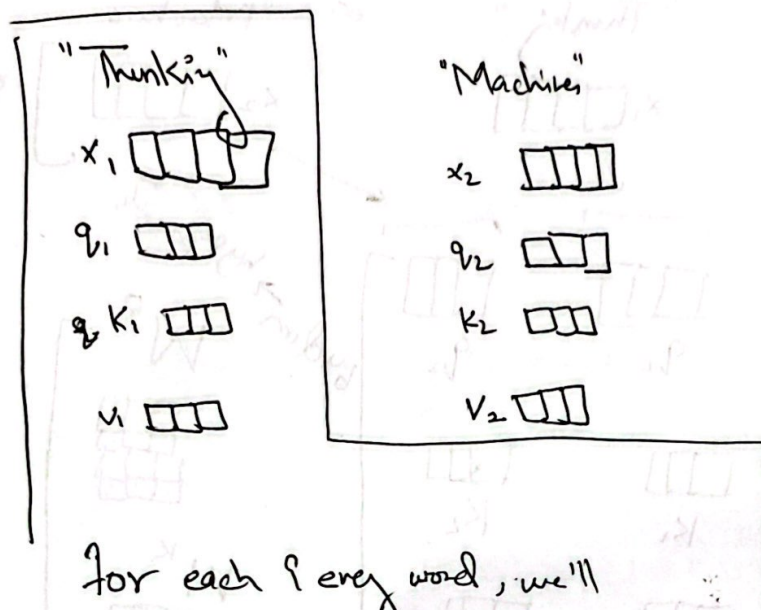$$V_1 = x_1 \cdot W^V$$

So,

Input

Embedding

Queries

Keys

Values

Score

"Thinking"                    "Machine"

$x_1$ ▭▭▭▭           $x_2$ ▭▭▭▭

$q_1$ ▭▭▭            $q_2$ ▭▭▭

& $k_1$ ▭▭▭          $k_2$ ▭▭▭

$v_1$ ▭▭▭            $v_2$ ▭▭▭

for each & every word, we'll
multiply $q_1$ & $k_1$

$q_1 \cdot k_1 = 112$     $q_1 \cdot k_2 = 96$    ⎤ this is
                                                 ⎥ w.r.t to
                                                 ⎦ First word.

Score ←                          for "machine", it
                                 would be
                                 $q_2 \cdot k_1 =$
We need to find out most          $q_2 \cdot k_2 =$
related words.

Divide score by $\sqrt{64}$ — → dim of $q_1$, term : $d_k$

Divide by $\sqrt{64}$ $\left(\sqrt{d_k}\right)$ = 0.88     0.12
                                            112/8     96/8

softmax +                    =  0.88          0.12       ⎤ to get
                                "thinking"    "machine"  ⎦ prob

# Attention 4

Softmax      0.88         0.12    → to neglgie low prob to less rel word

Softmax
x
value     $V_1 \times$ Softmax      $V_2 \times$ softmax ] to keep relev words only

↳
multiply
by $V_1$ $V_2$

Sum :     $z_1$ ☐☐☐ $\times V_1 d V_2$

↳ output of "thing"

Matrix multip



$X \times W^Q = Q$

2×4
↳ emb
2 word

$2 \times 4$     $2 \times 6 \cdot 4$

$X \times W^K = K$

$X \times W^V = V$

No window Sie becaue all words are being passed parallely.

Condensing

softmax $\left( \dfrac{Q \times K^{T}}{\sqrt{d_k}} \right) V = Z$

( $K^T \rightarrow$ transposed )

## Problem?

↳ We are using one $W^Q$, one $W^K$, one $W^V$ for all words. This is single head attention.

   Sure $W^K, W^Q, W^V$ for all words.

   but problem lies when there might might be important connection.

   So, we might use multiple weights

## ⤬MULTI-HEAD Attention

   ↳ to find importance of other words too
   —mean not only finding single word but multiple relevant words

Attention head 0 $\left[ \begin{array}{cc} Q_0 & W_0 Q, W_0^Q \\ K_0 & W_0^K \\ V_0 & W_0^V \end{array} \right.$

head 1 $\left[ \begin{array}{cc} Q_1 & W_1^Q \\ K_1 & W_1^K \\ V_1 & W_1^V \end{array} \right.$

DIFFERENT HEADS FOR DIFFERENT WORDS.

Attention 5

eight heads
└ eight diff$^T$ weight matrices.
will give 8 z matrices.

So, now

X
"thinly machi"


→

Attn head 0
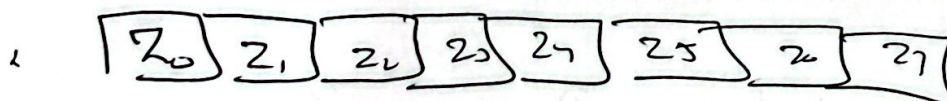head 1
head 2
3
4
5
6
7

⇒

$z_0$ (3×2
$z_1$
$z_2$
$z_3$
$z_4$
$z_5$
$z_6$
$z_7$

But we are supposed to get 1 z value not 8.

1 — we'll concatenate all $z$s

| $z_0$ | $z_1$ | $z_2$ | $z_3$ | $z_4$ | $z_5$ | $z_6$ | $z_7$ |

2 — initialize another $W^0$ matrix → would be updated in back propagation
$W^0$ will be multiplied with concatent z

result would be [ Z ⊞ ]

3 POSITIONAL ENCODING ~~after~~ ~~before~~ ~~convert~~ words
into embedding

4 to preserve oder.

| INPUT | Je | Suis | ETUDIANT |
|---|---|---|---|
| Embedding | $x_1$ | $x_2$ | $x_3$ |
| | $+$ | $+$ | $+$ |
| | Positional encoding $t_1$ | $t_2$ | $t_3$ |

| Embedding with the signal | $x_1$ | $x_2$ | $x_3$ |

$t_1$ would be nearer to $\underline{t_2}$ than $\underline{t_3}$.

This happen before pass to encoder

In encoder, before passing to feed forward
~~Each~~ encoder and after self atten, we'll
be ~~part~~ apply normalizati. and
a <u>residual</u> connecti.
 └ like in RNN, you are skipping one step.

  ~~If~~ i.e if self atten is not useful,
  you are skipping it thgh residual network

In add & Normalize
  └ layer norm $\left( X + z \right)$
             $\downarrow$      $\downarrow$
            input    encoder o/p
                       └ by

In Decoder

Encoder Decoder Attention

↳ Save only output of encoder
↳ passed into the layer of decoder

input - all words at once
output - one word at a time.
↳ after each word predicts that
whole output will be give input
to decoder. It won't stop
untill <EOS>.