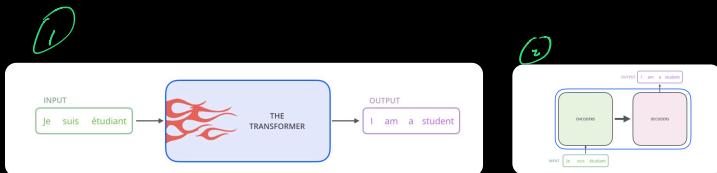


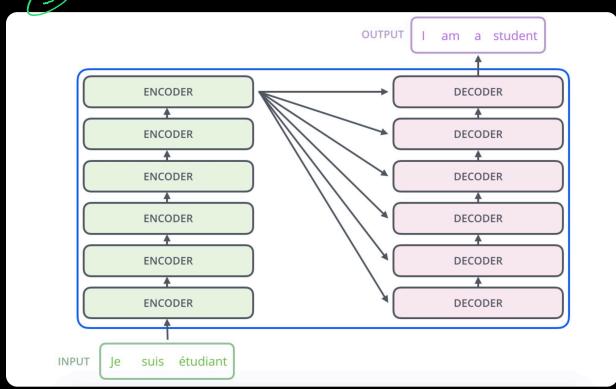
Assumption
You know Encoder People
& Attention



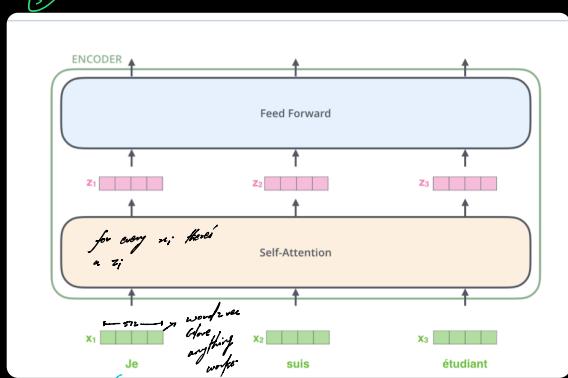
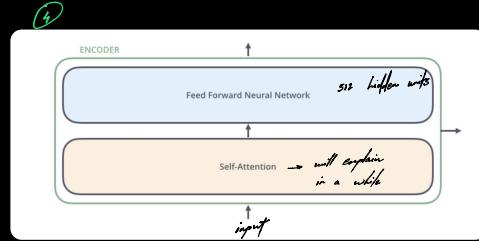
Jafarman.github.io
↳ Illustrated Transformer → better explanation
than anything out there

Transformer
↓
advanced seq2seq model





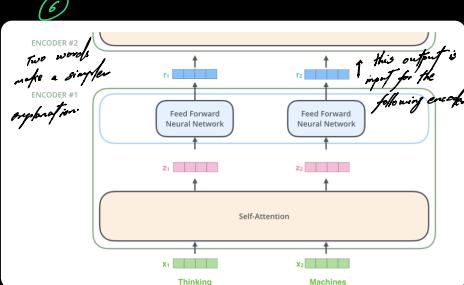
they are identical to each other.



→ all inputs are taken at once
we try to work around
the limitations of RNN

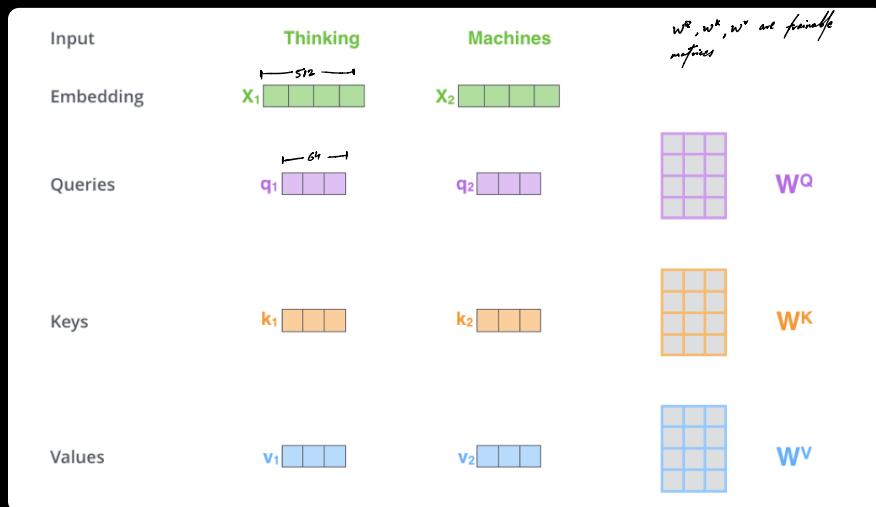
Self Attention

"The animal didn't cross the street because it was too tired."



It's called self affection because we're going to which part of ourself to focus on.

(x_i) → z_i
→ some words
need more affection.

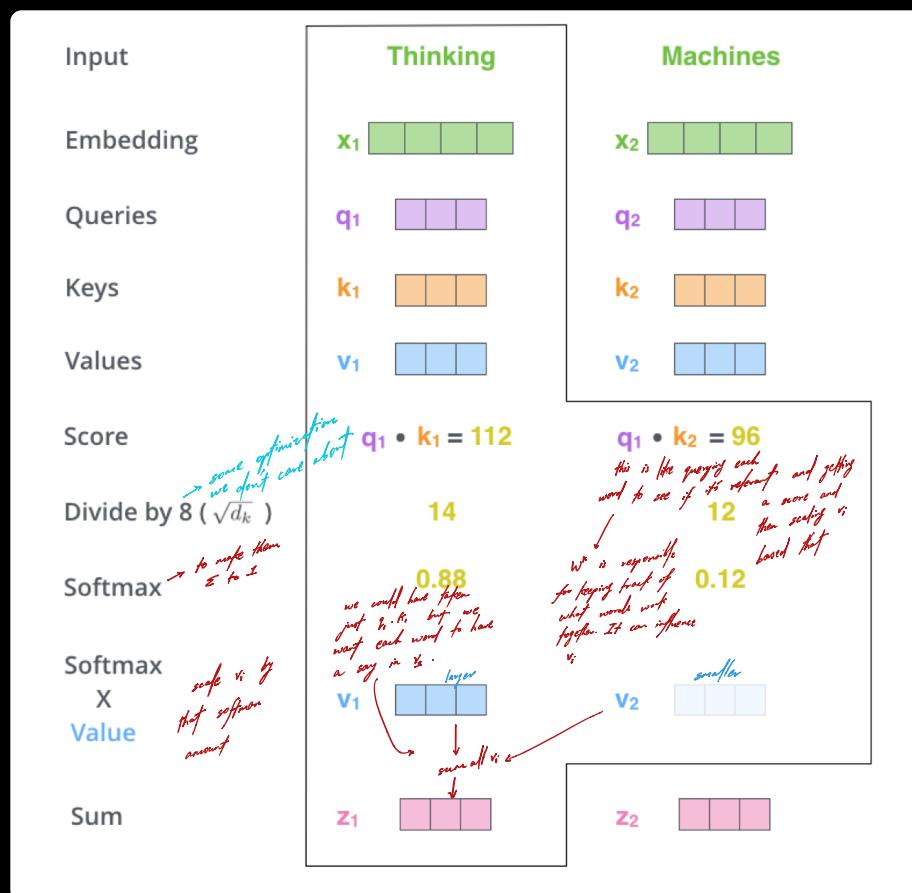


$$b_i = x_i \cdot w^Q \rightarrow$$
 this is the single row without activation.

$$k_i = x_i \cdot w^K$$

$$v_i = x_i \cdot w^V$$

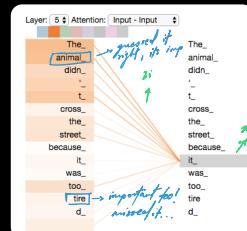
just take x_i & get
products of it $q_i, k_i \& v_i$



we going to use z_i values
to make v_i somehow.

$$\begin{array}{l} X \times W^Q = Q \\ X \times W^K = K \\ X \times W^V = V \end{array}$$

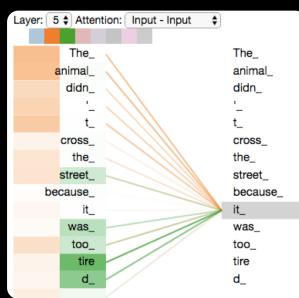
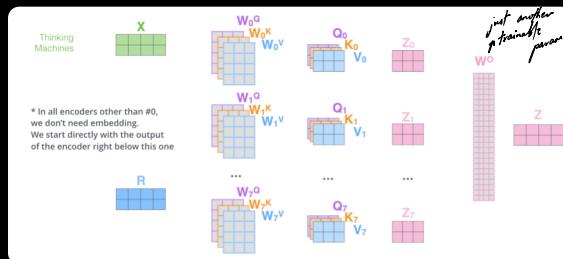
All this summed in a bunch of matrix multiplications.



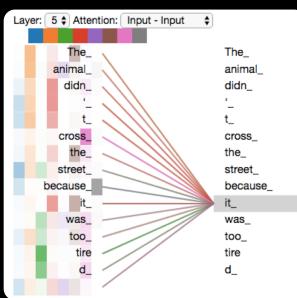
$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V = Z$$

The seq attention calculation in matrix form

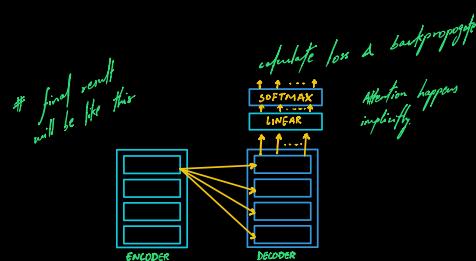
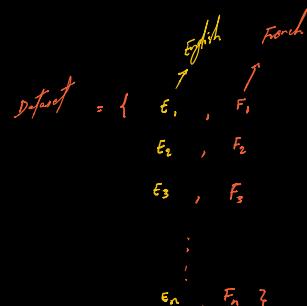
They say we should have multiple $[q, k, v]$ matrices → called **attention head**.
can focus on only one word
have 8 attention heads → can focus on multiple words effectively
so for each x_i we have 8 to vectors.



This is what happens for 2 attention heads



This is what happens with multiple attention heads





Loss function
conceptually correct

* big takeaway

- ① significantly faster than now → because all input are taken at once
- ② Just like now → it learns what part of input is important → on its own
- ③ state of the art → (BERT) is the simpler version of this architecture.
this is what google found out

* Sanketh Sir

* "If in an interview you say you know Transformer, It'll ask you all of this"

#BERT Bidirectional Encoder Representation from Transformers

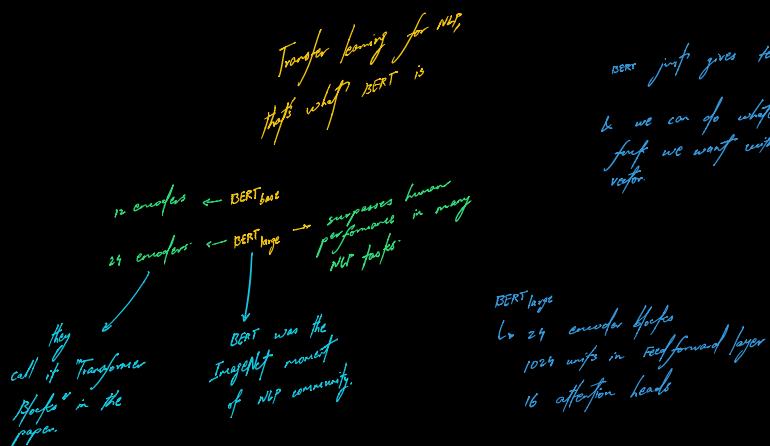
As long as you can frame a problem as text 2 text, BERT can be applied.

↳ text summarization
classification
question answer.

BERT is much simpler than transformer, but provides almost the same performance.

There's no decoder

now just gives text to vector
↳ we can do whatever we want with the vector



* It's bidirectional model

↓
so words before & after are taken into account

↓
BERT uses context & that helps with all those boundaries.

BERT_{base}

↳ 12 encoder blocks
768 units in feedforward layer → that'll be abuse of output
12 attention heads

Bert application
↳ transfer learning

* Example
given Example:

BERT can take upto 512 word surfaces

