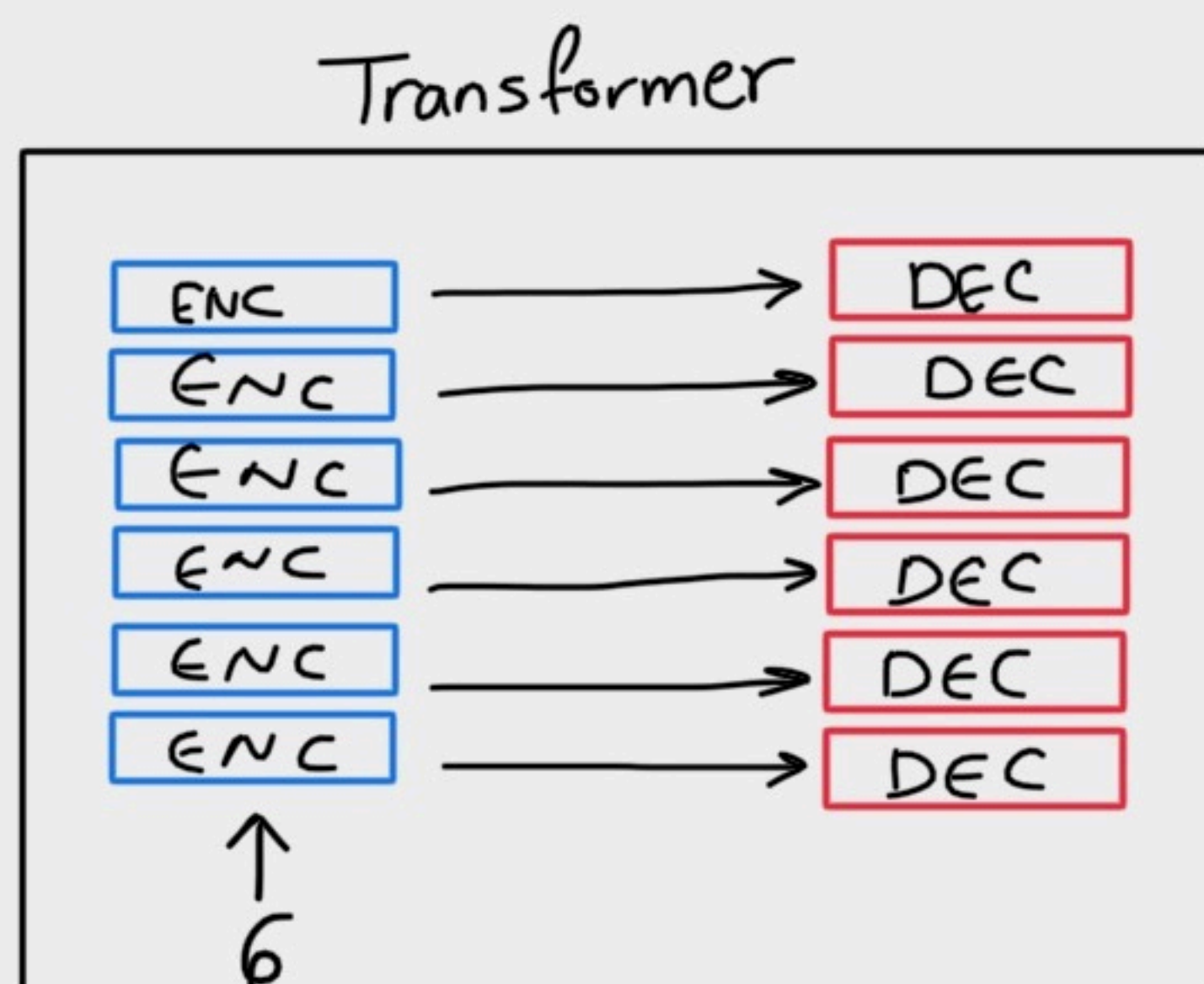# Transformer
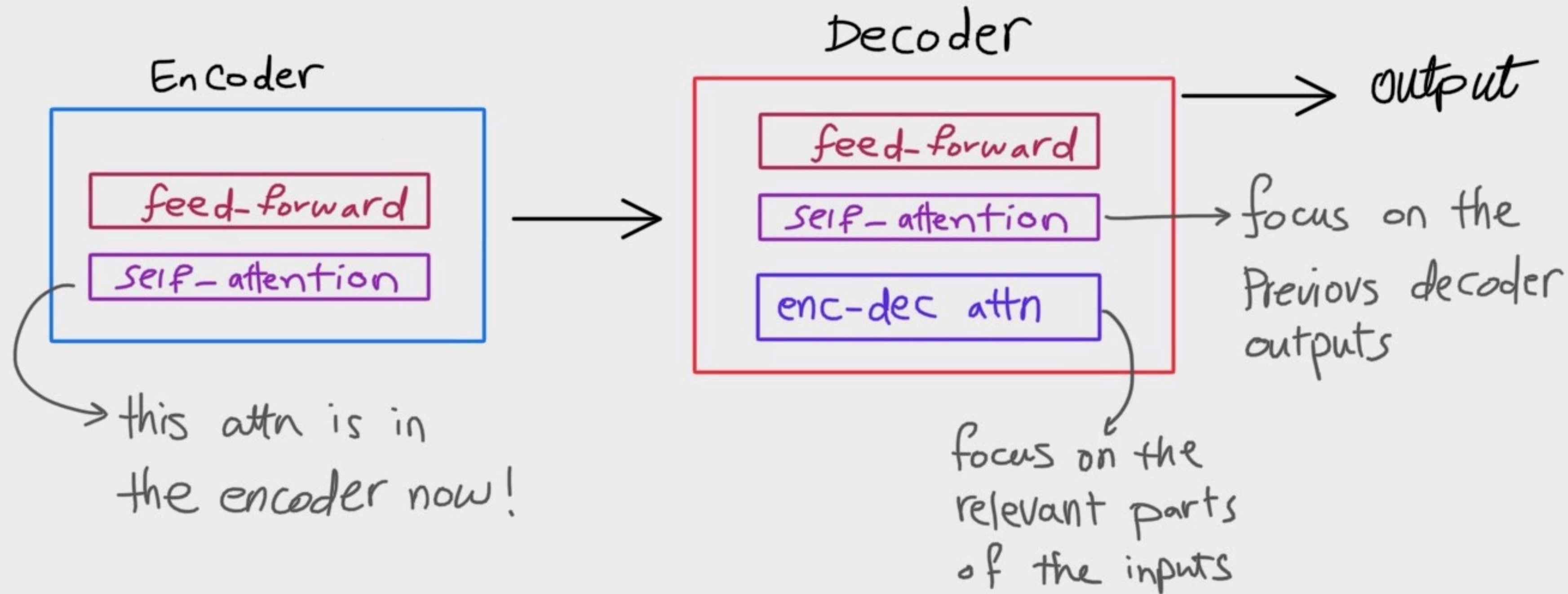
As opposed to RNN models that take the input sequence token-by-token, the transformer takes the input sequence as a whole (parallel), then generates the output one-by-one as before.

Another difference between transformer seq2seq models and RNN seq2seq is that they use feed-forward neural networks instead of RNNs.

The third major novelty in transformers is a concept called self-attention.

It has a stack of encoders and decoders.
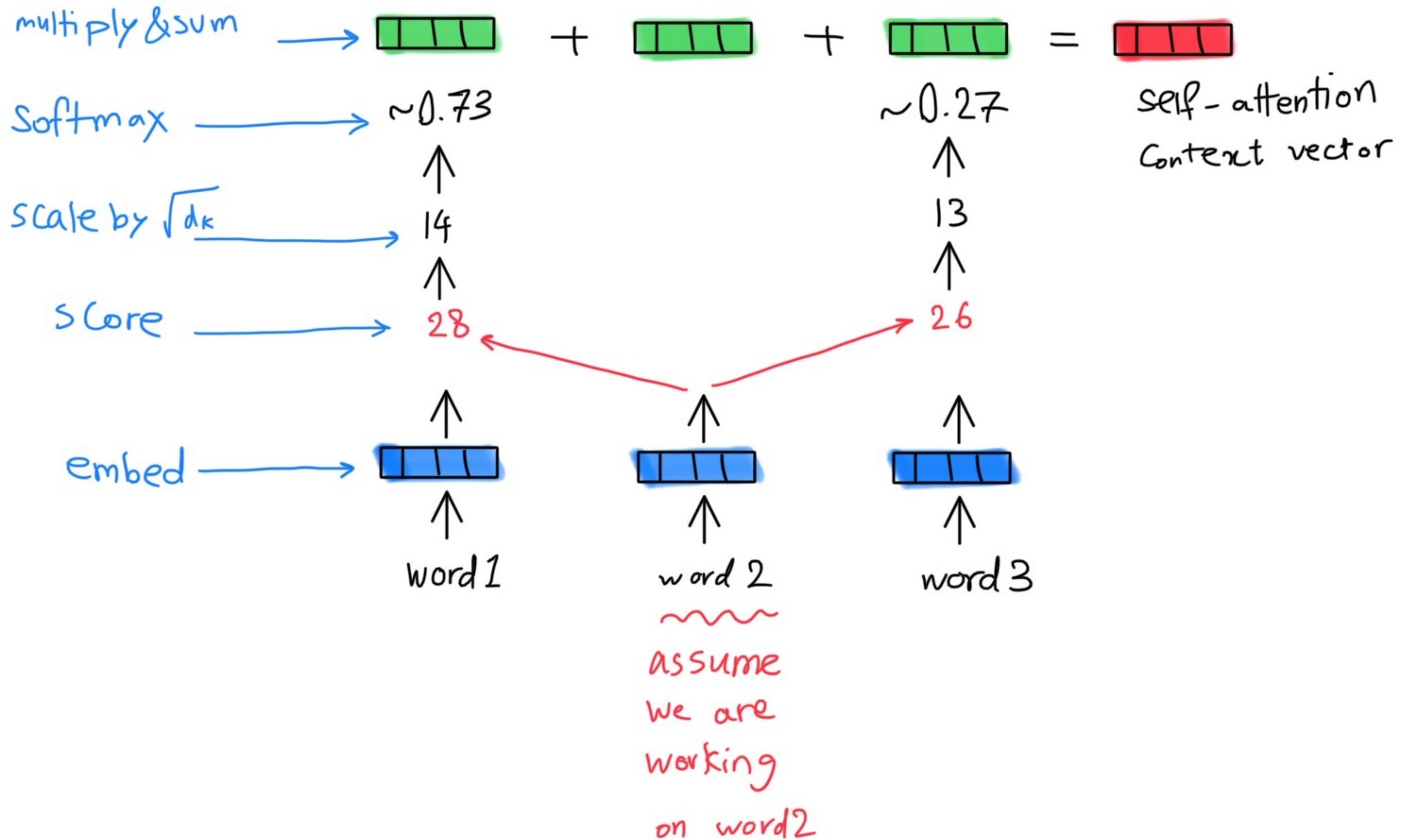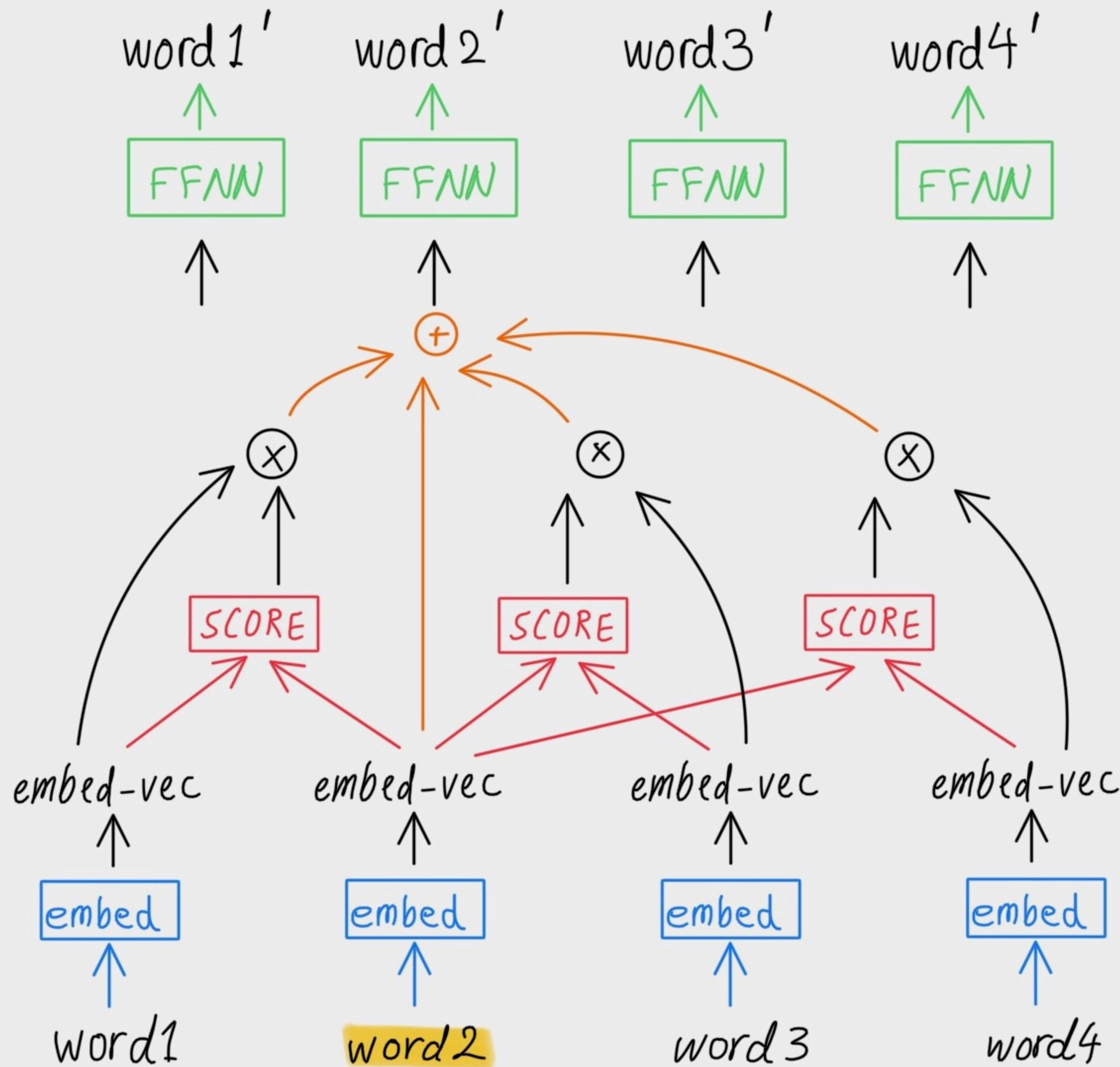
Transformer

Encoder

Decoder

feed-forward

self-attention

this attn is in the encoder now!

feed-forward

self-attention → focus on the previous decoder outputs

enc-dec attn → focus on the relevant parts of the inputs

output

# Self-Attention

multiply & sum → ▭▭▭ + ▭▭▭ + ▭▭▭ = ▭▭▭

Softmax → ~0.73          ~0.27          self-attention context vector

scale by $\sqrt{d_k}$ → 14          13

score → 28 ⟷ 26

embed → ▭▭▭▭          ▭▭▭▭          ▭▭▭▭

word 1          word 2          word 3

assume we are working on word2

# De-crypt the transformer paper



**word1'**  **word2'**  **word3'**  **word4'**

FFNN  FFNN  FFNN  FFNN

embed-vec  embed-vec  embed-vec  embed-vec

embed  embed  embed  embed

word1  word2  word3  word4
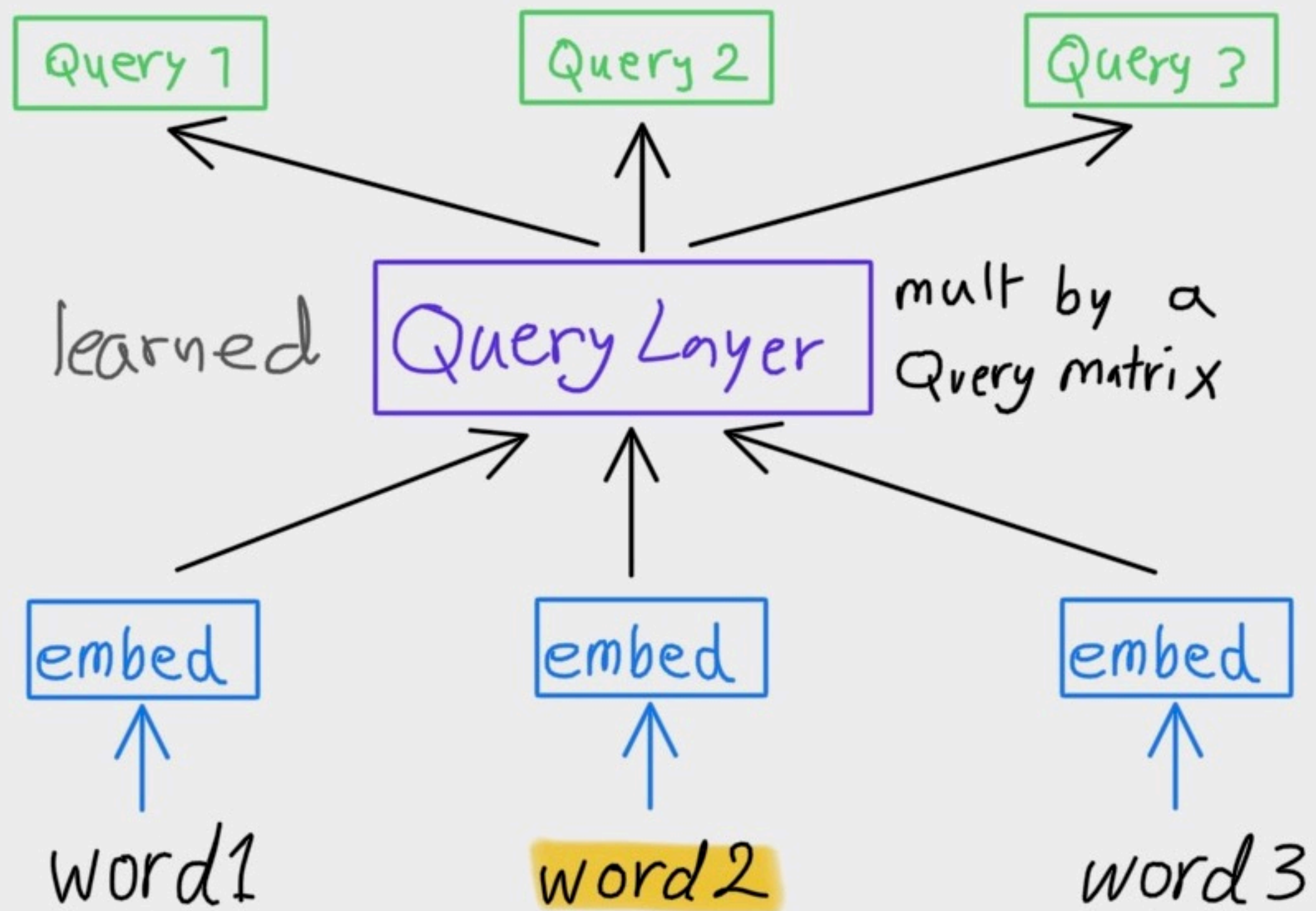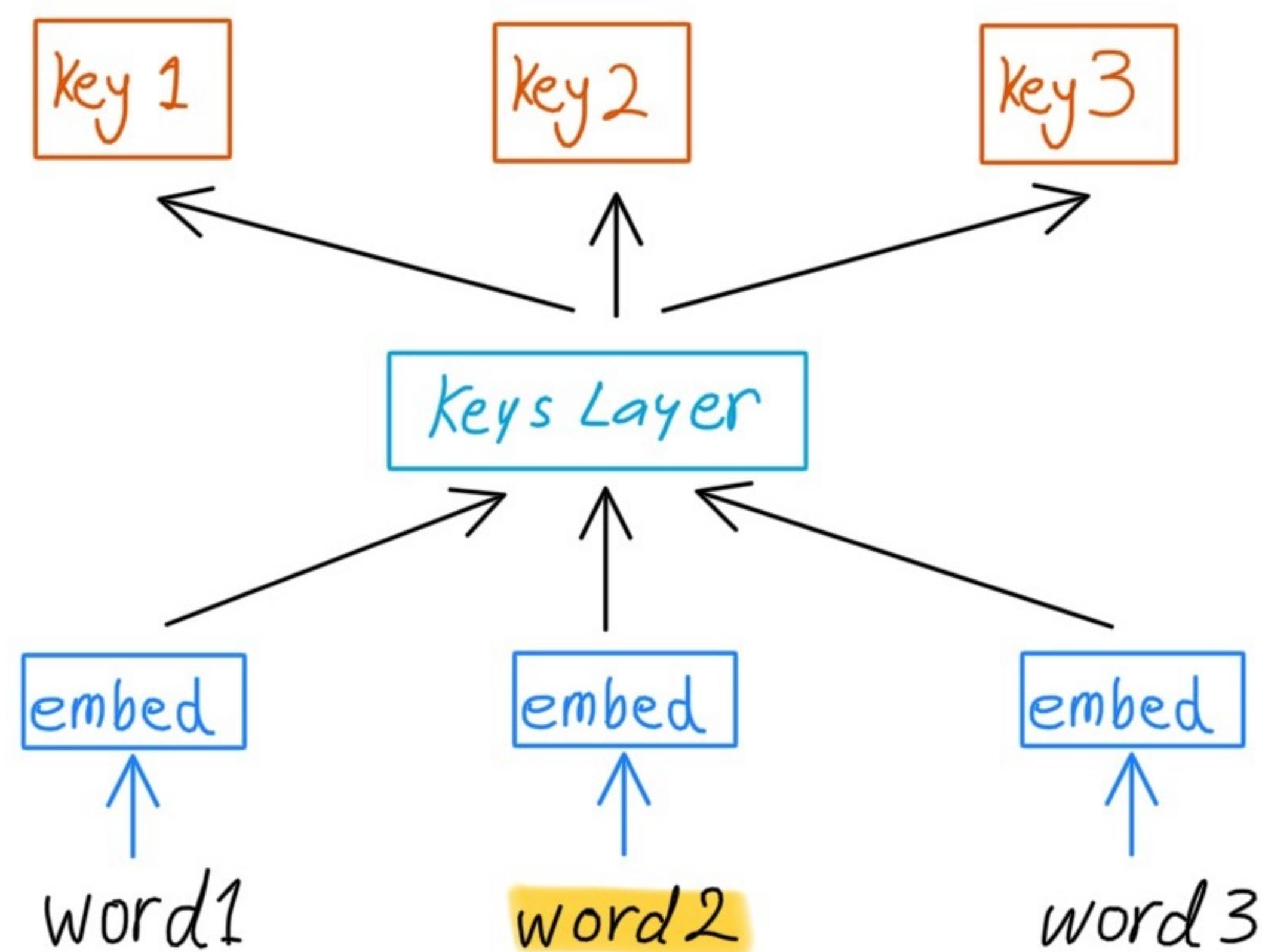
because we are directly going from embedding

● Problem in this figure: model is mainly focusing on other similar words. ⟶ Let's modify it!

Keys Layer

key 1    key 2    key 3

embed ← word1    embed ← word2    embed ← word3

● Using the keys and queries, we calculate like this:

|  | Q | K | score(Q,k) | *$\sqrt{d_k}$ | softmax | softmax *K |
|---|---|---|---|---|---|---|
| word1   embed¹ |  | key1 | 28 | 14 | 0.2 | vec1 |
| word2   embed² | Query | key2 |  |  |  | vec2 |
| word3   embed³ |  | key3 | 26 | 13 | 0.8 | vec3 |

vec1 + vec2 + vec3    self-attn Context vector → Passed to FFNN

● The authors also add a third thing: values. So:

| | Q | K | V | score($Q$, $K$) | * $\sqrt{d_k}$ | softmax | softmax * $V$ |
|---|---|---|---|---|---|---|---|
| word1 embed¹ | | Key1 | Val 1 | 28 | 14 | 0.2 | vec1 |
| word2 embed² | Query | Key2 | Val 2 | | | | vec2 |
| word3 embed³ | | Key3 | Val 3 | 26 | 13 | 0.8 | vec3 |

$Q$ matrix: learned during training process        learned

$K$ : by multiplying embeddings by $K$ matrix

$V$ : by multiplying embeddings by $V$ matrix        learned