# Transformers

**1** **Why should you care about Transformers?**

ChatGPT

**2** RNNs: Problems and progress

Je      suis      content

<- Reverse of word, token embedding process

Encoder

| h1 | h2 | h3 |

Decoder

c

h1 → h2 → h3

Embedded vector  [0.3 ...]   [0.4 ...]   [0.1 ...]

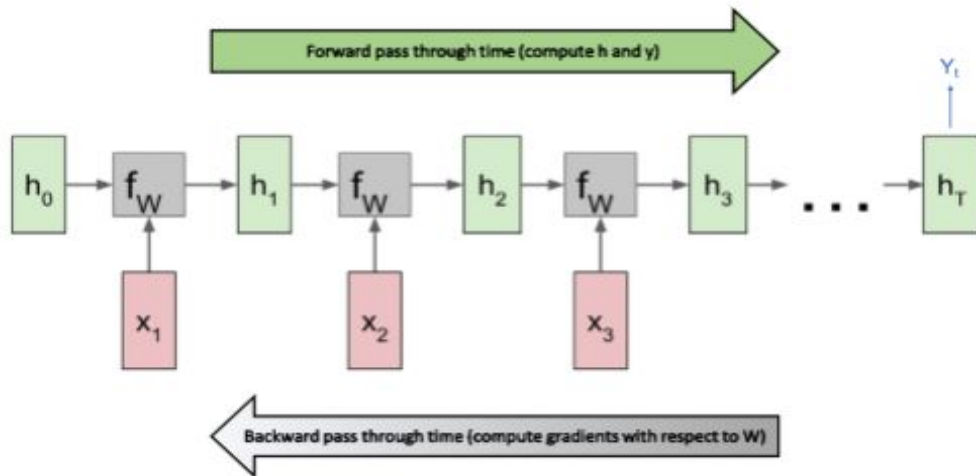Token    55       4       23

         I        am      happy

# What are the key issues we face here?

1. Information bottleneck at interface

2. Vanishing gradient problem

3. We have to compute the entire sequence recursively (makes scaling very hard!)

# RNNs suffer from vanishing gradient through time

# What does this mean for our performance?

A simplification of problems with RNNs:

> *Sally adored reading; when she received a book on her birthday she was <span style="color:red">older</span>*
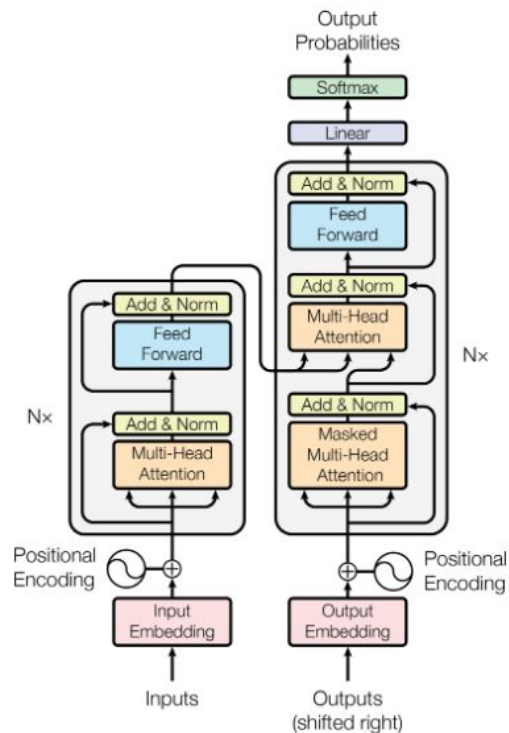
What we'd like:

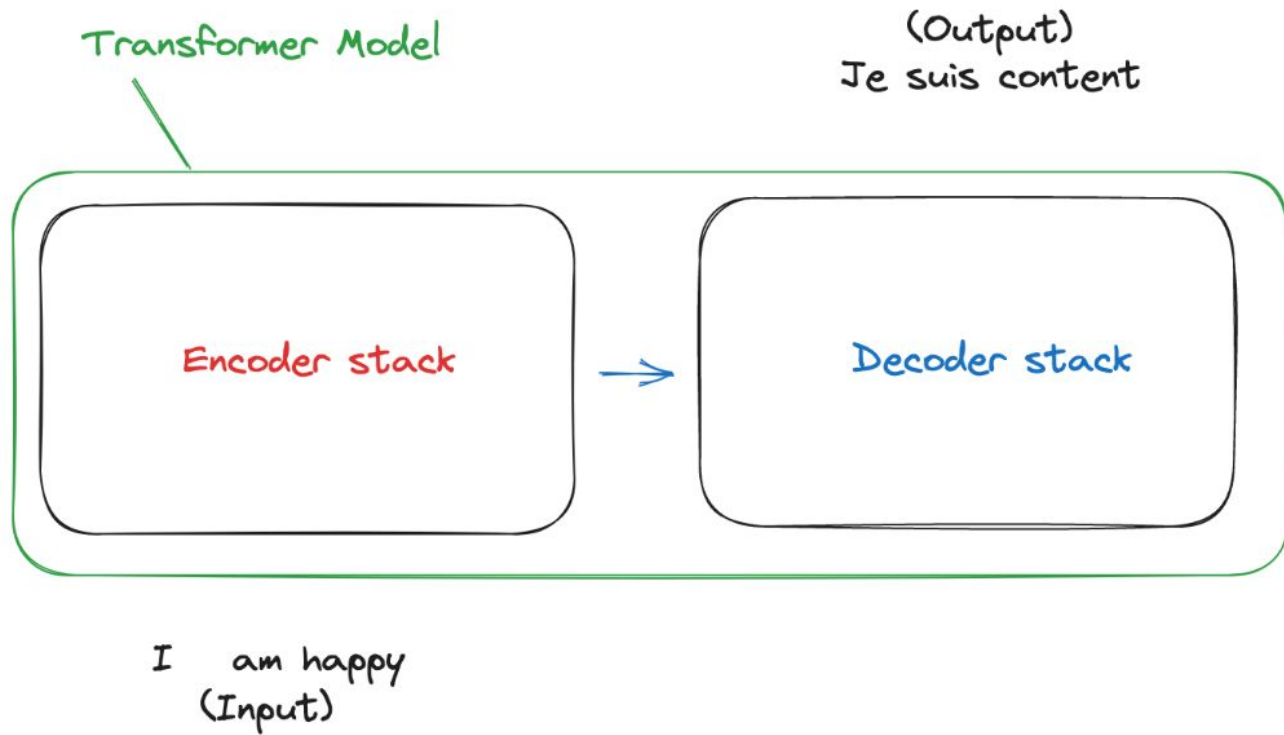> *Sally adored reading; when she received a book on her birthday she was <span style="color:green">happy!</span>*

RNNs are likely to miss out on **important context** from earlier in the sentence because of their recency bias 🫠
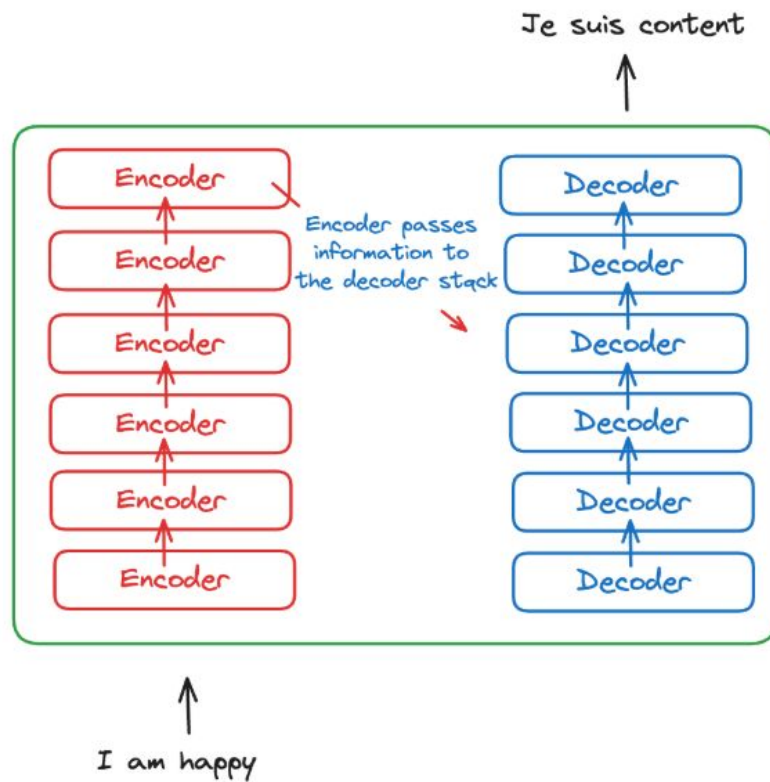
# The paper that started it all: Attention is All You Need

# The highest level view:
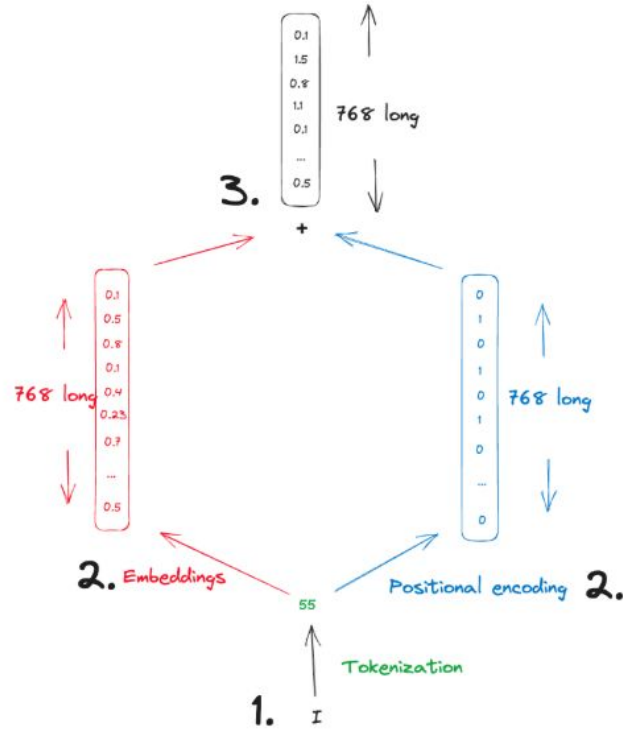
Transformer Model

(Output)
Je suis content

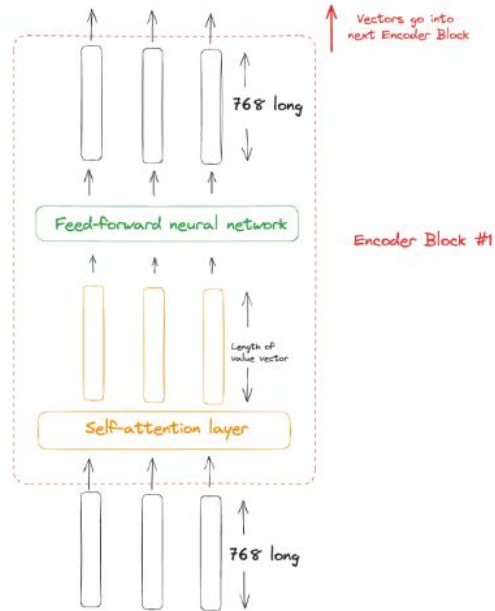Encoder stack → Decoder stack

I  am happy
(Input)

Broken down a bit more:

**Before we talk about what's going on inside the encoder layers, let's talk about what's going into it!**
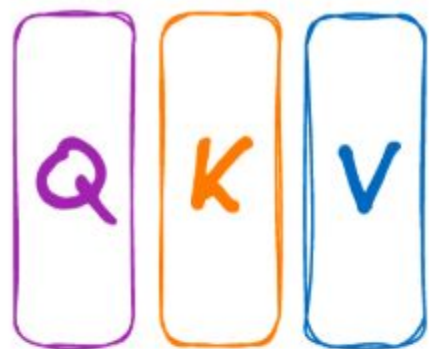
Zooming in on one of the encoders:



- Usually the **embedding size** is 768
- The **hidden dimension** (the length of the projected Q, K and V vectors) is also 768
- In the example we'll use size 2 for simplicity

❓ What's going on in this strange self-attention layer?

1) Each token (word) embedding gets **projected** ➡️ into 3 further vectors: the **query, key and value vectors**

2) We compute a **scaled dot-product** 🔴 on the query and key vectors to work out how much each word relates to those around it

3) Take these scores and **normalize with softmax** 🔃

4) **Multiply by our value vectors** ❎, sum and pass to our dense neural network.

# An example sentence

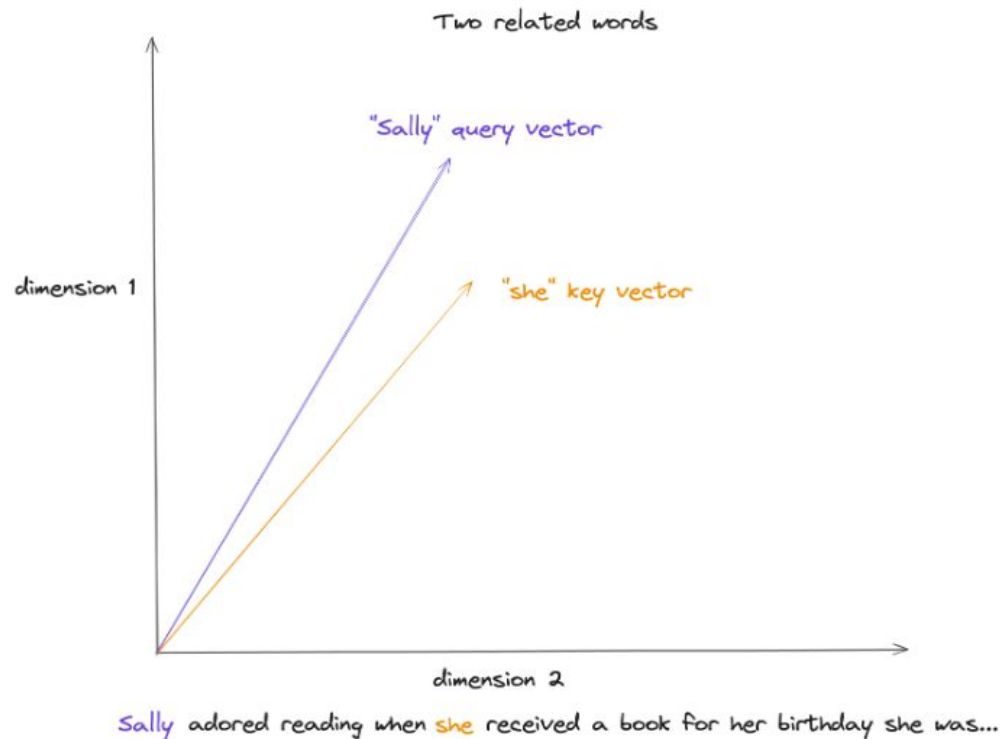Sally adored reading; when she received a book for her birthday she was...

# Three zoomed-in examples:
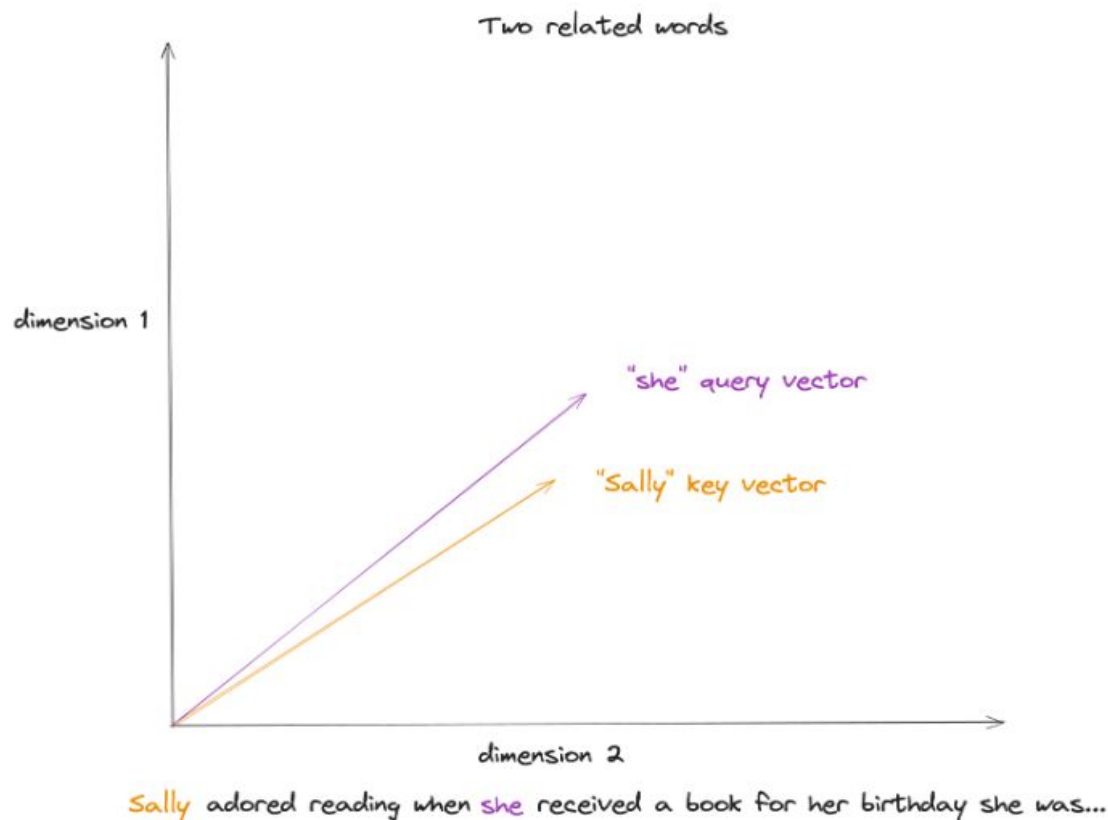
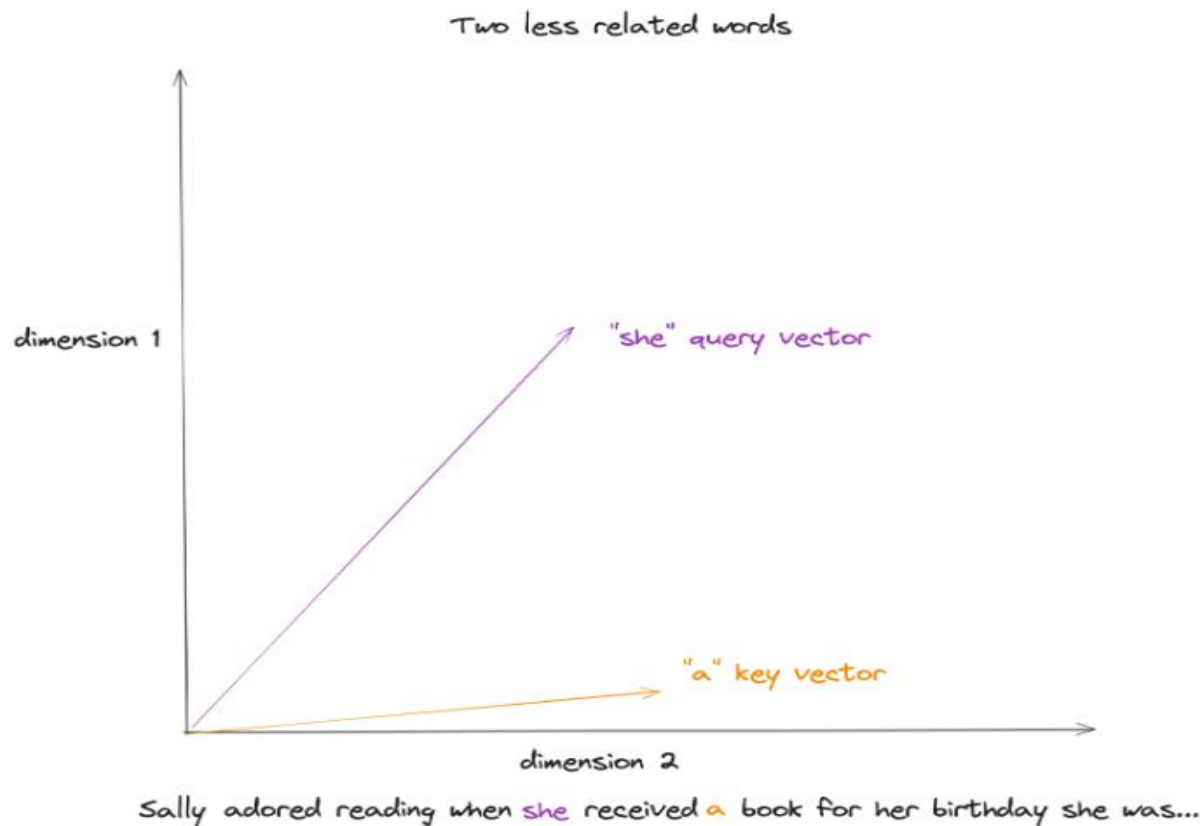Two words in a sentence that are closely related



Two related words

"Sally" query vector

"she" key vector

dimension 1

dimension 2

Sally adored reading when she received a book for her birthday she was...

# The same two words but seen from the other perspective:



Two related words

dimension 1

"she" query vector

"Sally" key vector

dimension 2

Sally adored reading when she received a book for her birthday she was...

# Finally, two words with a weak connection:



Two less related words

"she" query vector

"a" key vector

dimension 1

dimension 2

Sally adored reading when she received a book for her birthday she was...

# Let's look at one dot product

To keep it really simple, we're going to imagine our Q, K and V have only been projected into two dimensions

# What happens once we have our dot products?

Dot-product between vectors

| | Sally$_K$ | adored$_K$ | reading$_K$ | when$_K$ | she$_K$ |
|---|---|---|---|---|---|
| Sally$_Q$ | 112 | 75 | 60 | 12 | 105 |

# Then we scale:

Scaled dot-product between vectors

We divide by $\sqrt{hidden\_dimension}$ of embedding

In our case the root of $2 \cong 1.4$

|  | Sally$_K$ | adored$_K$ | reading$_K$ | when$_K$ | she$_K$ |
|---|---|---|---|---|---|
| Sally$_Q$ | 80 | 53 | 43 | 8.5 | 75 |

# Finally we apply softmax:

Apply softmax

$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}}$$

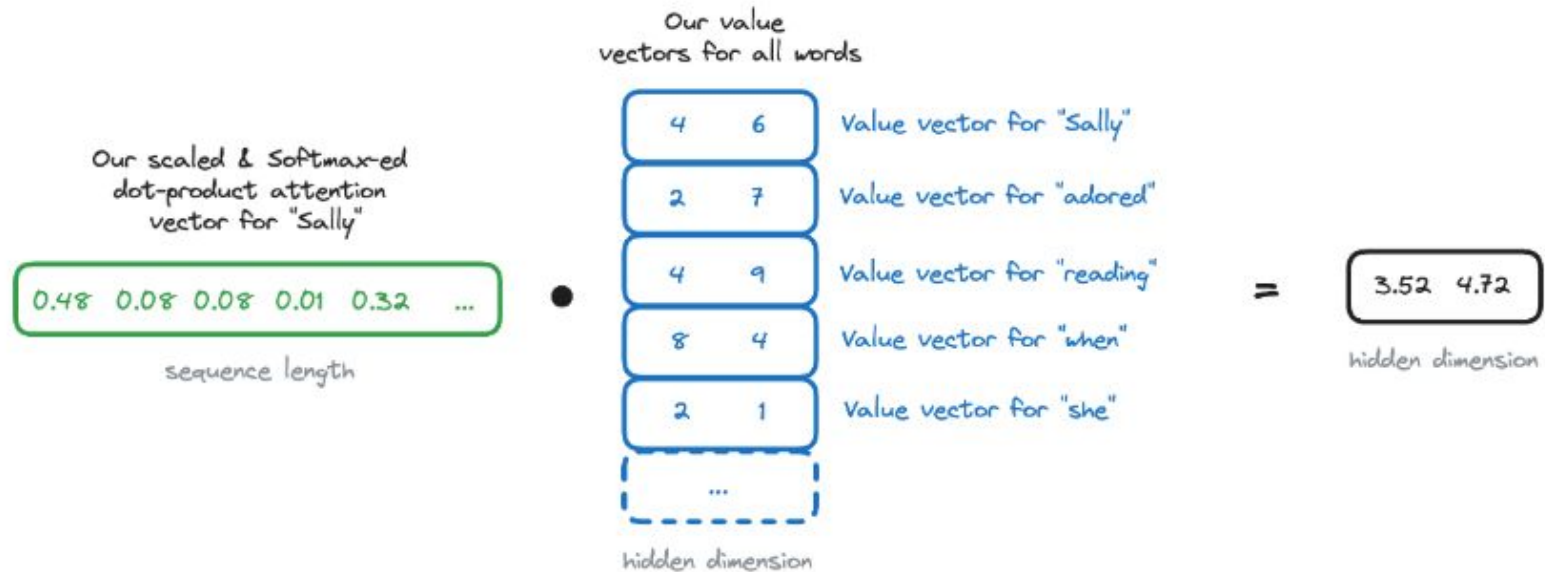|  | Sally$_K$ | adored$_K$ | reading$_K$ | when$_K$ | she$_K$ |
|---|---|---|---|---|---|
| Sally$_Q$ | 0.48 | 0.08 | 0.03 | 0.01 | 0.32 |

We have to do this for each word in our sentence
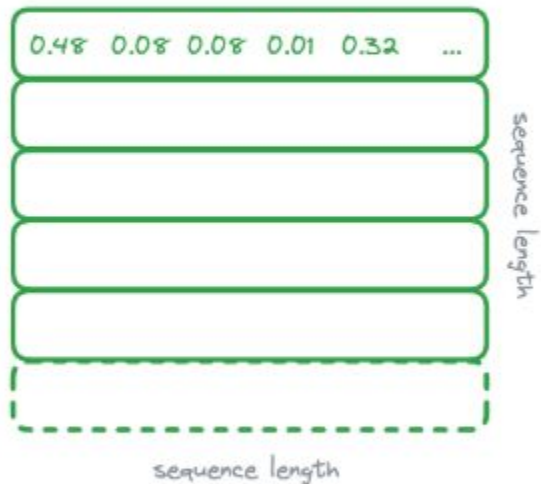You can see how this becomes a matrix operation 💪
We get the scaled dot-product attention for all of our Query and Key vectors:

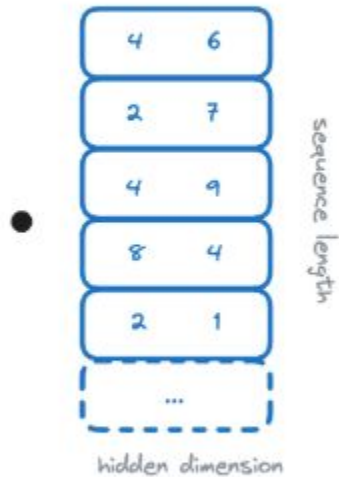| | $Sally_K$ | $adored_K$ | $reading_K$ | $when_K$ | $she_K$ | ... |
|---|---|---|---|---|---|---|
| $Sally_Q$ | 0.48 | 0.08 | 0.03 | 0.01 | 0.32 | ... |
| $adored_Q$ | 0.04 | 0.03 | 0.52 | 0.3 | 0.01 | ... |
| $reading_Q$ | 0.01 | 0.02 | 0.13 | 0.86 | 0.01 | ... |
| $when_Q$ | 0.01 | 0.03 | 0.02 | 0.58 | 0.01 | ... |
| $she_Q$ | 0.13 | 0.12 | 0.02 | 0.01 | 0.04 | ... |
| ... | ... | ... | ... | ... | ... | |

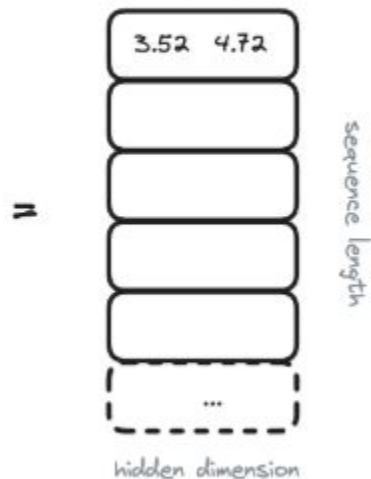# And then multiply our "similarity score" with all of our Value vectors



Our value vectors for all words

| 4 | 6 | Value vector for "Sally" |
| 2 | 7 | Value vector for "adored" |
| 4 | 9 | Value vector for "reading" |
| 8 | 4 | Value vector for "when" |
| 2 | 1 | Value vector for "she" |
| ... | | |

hidden dimension

Our scaled & Softmax-ed dot-product attention vector for "Sally"

0.48  0.08  0.08  0.01  0.32  ...

sequence length

=  3.52  4.72

hidden dimension

Our scaled & Softmax-ed
dot-product attention
vector for all words

| 0.48 | 0.08 | 0.08 | 0.01 | 0.32 | ... |

sequence length

sequence length

●

Our value
vectors for all words

| 4 | 6 |
| 2 | 7 |
| 4 | 9 |
| 8 | 4 |
| 2 | 1 |
| ... | |

hidden dimension

sequence length

=

Our final matrix
that then goes into our
feed-forward network

| 3.52 | 4.72 |
| ... | |

hidden dimension

sequence length

So really the entire thing can be written like this:

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

$$= Z$$

We are done with our multiplications
Now we just need to normalize and pass through to the feed-forward neural network
The neural network will output vectors of our original embedding dimension (e.g. 768)

🤯 Computing one set of all of these Q, K, V multiplications and processes is what we call "single-headed attention".

When we are doing our initial linear projections (used to create the Q, K and V vectors) we can express these operations as matrices of weights too!
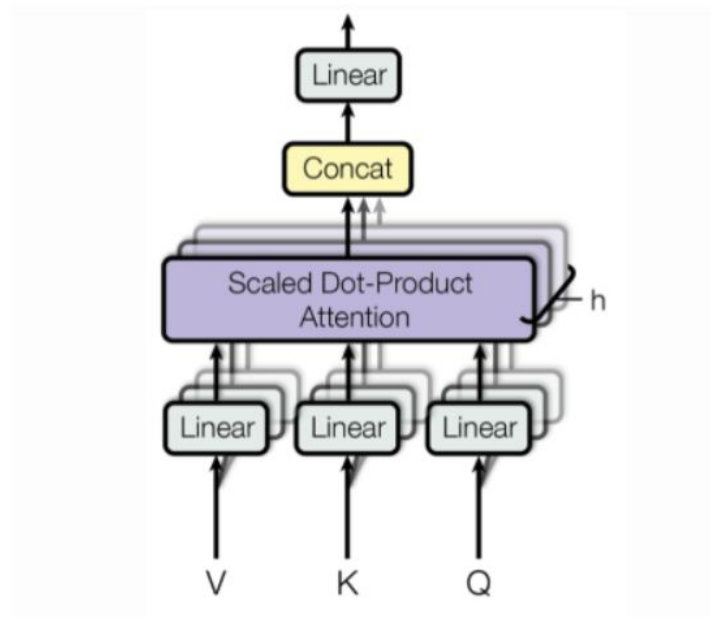
# Multi-headed?

- We can use multiple heads to split up and analyze different parts of our embedding
- Each can focus a different part of the embedding eg. working on semantic vs syntactic features of our sentences 🤯 🤯 🤯

# Let's check in with our original diagram:



One encoder block

# What about the decoder?

At the highest level view, its job is to choose the **most likely next token**:

Ex 1.

Je

↑

Decoder

↑

&lt;start&gt;

Ex 2.

suis

↑

Decoder

↑

&lt;start&gt; Je

Ex 3.

content

↑

Decoder

↑

&lt;start&gt; Je suis
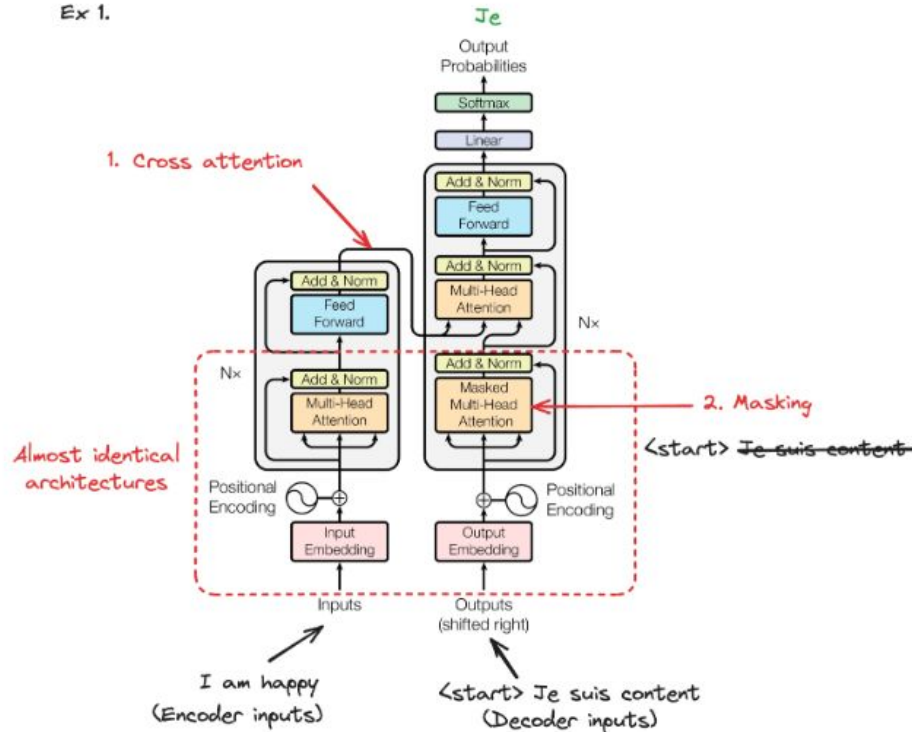
# How does it do this? And what happened to all the work the encoder did?



Ex 1.

1. Cross attention

2. Masking

Almost identical architectures

<start> ~~Je suis content~~

I am happy
(Encoder inputs)

<start> Je suis content
(Decoder inputs)

Je

Ex 2.

suis

Output
Probabilities

Softmax

Linear

1. Cross attention

Add & Norm

Feed
Forward

Add & Norm

Add & Norm

Multi-Head
Attention

Feed
Forward

Nx

Add & Norm

Nx

Add & Norm

Multi-Head
Attention

Masked
Multi-Head
Attention

2. Masking

<start> Je suis content

Almost identical
architectures

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

I am happy   <start> Je suis content

# Cross attention



Cross-attention in Transformer

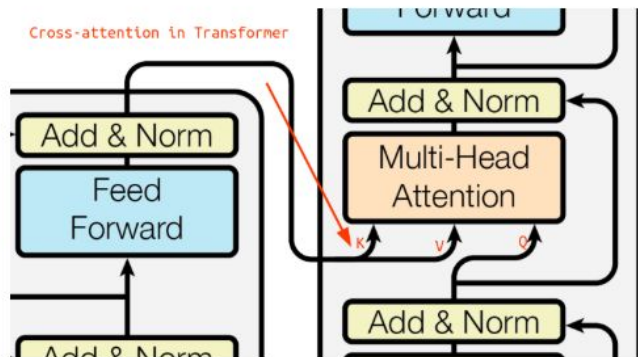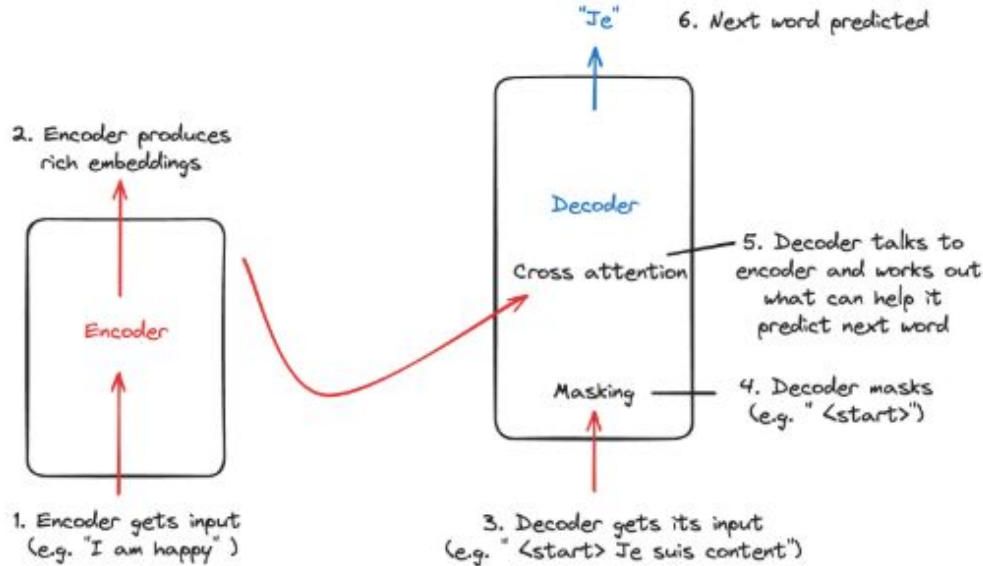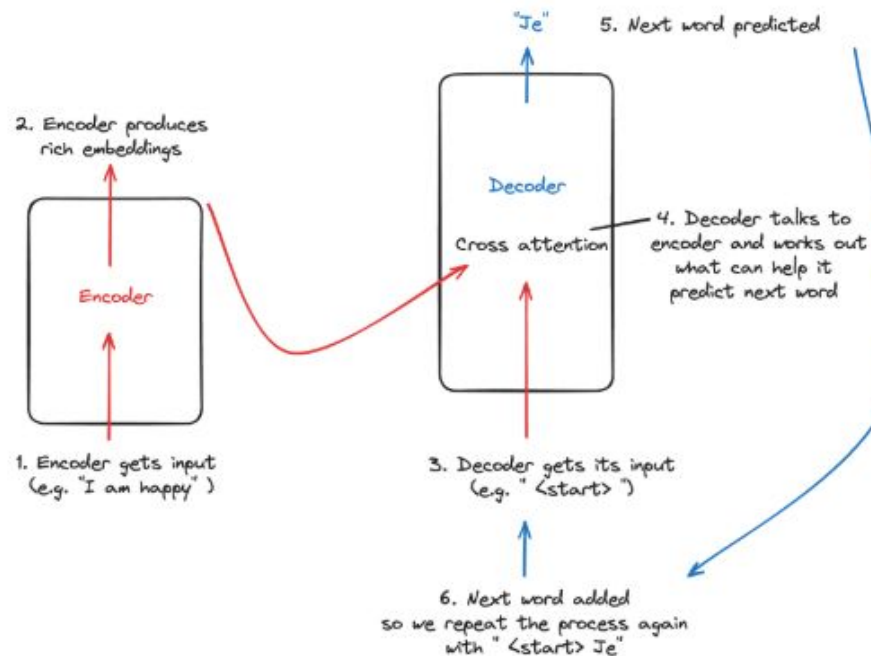Attention between encoder and decoder (a.k.a. cross-attention):

- **Self-attention** operates within a single sequence and captures the relationships between tokens within that sequence.
- **Cross-attention** operates between two different sequences and captures the relationships between tokens from the source sequence and tokens from the target sequence, allowing the model to generate relevant output based on the information in the source sequence.

# Let's recap the training process at a high level



"Je"    6. Next word predicted

2. Encoder produces
rich embeddings

Decoder

Encoder

Cross attention    5. Decoder talks to
encoder and works out
what can help it
predict next word

Masking    4. Decoder masks
(e.g. " <start>")

1. Encoder gets input
(e.g. "I am happy")

3. Decoder gets its input
(e.g. " <start> Je suis content")

# So what does inference look like?

# That's all there is to it!