

Transformer

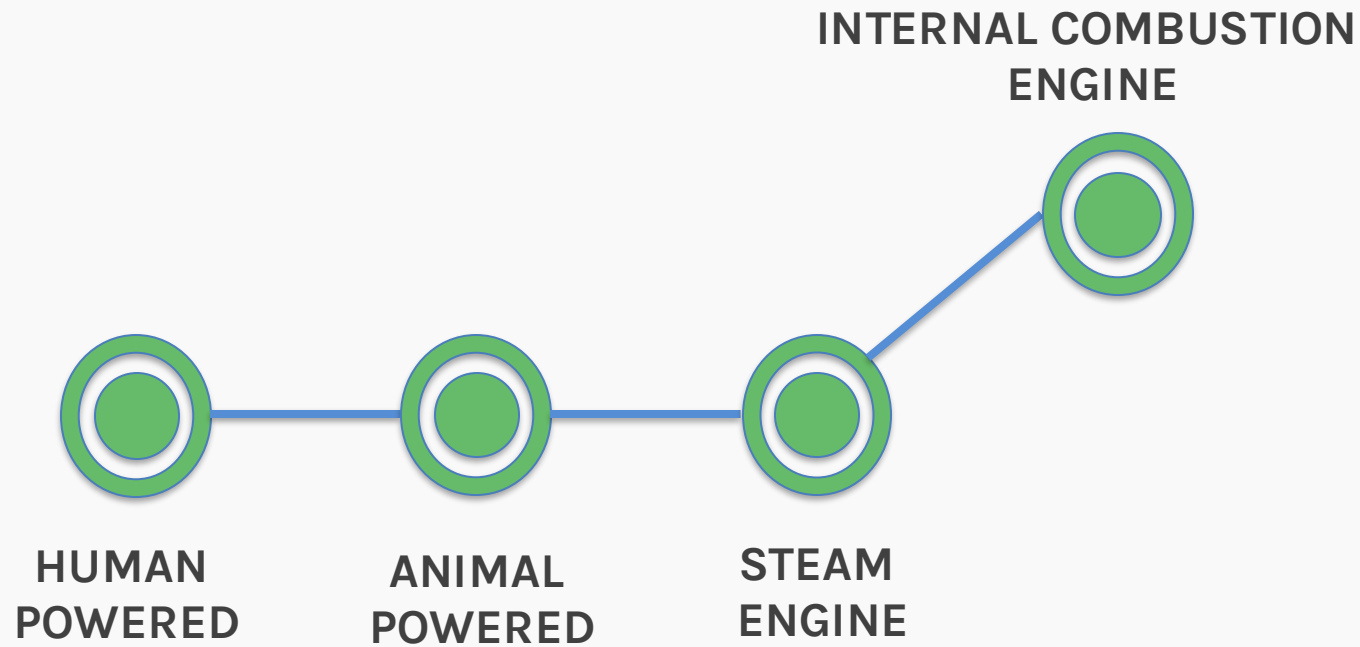
Pavlos Protopapas



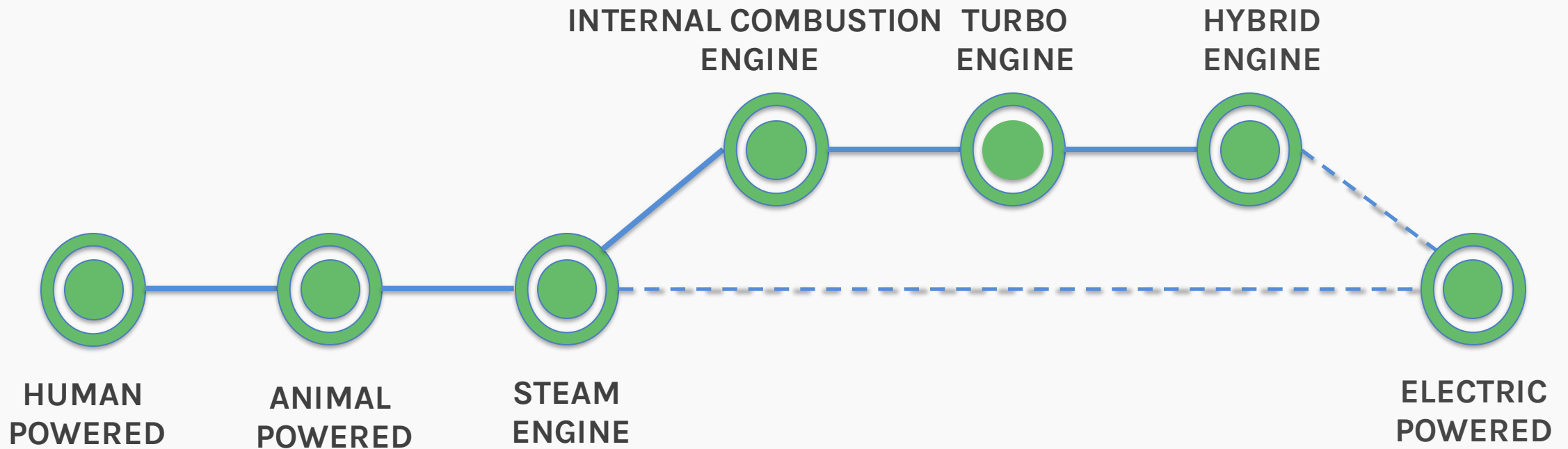
Outline

- **Motivation for Attention**
- Attention Basics
 - Using cosine similarity as a tool for contextual relations
 - Self-Attention
- Building blocks of Transformers and BERT
 - Multi-head attention block
 - Positional Encoding
 - Bringing it all together

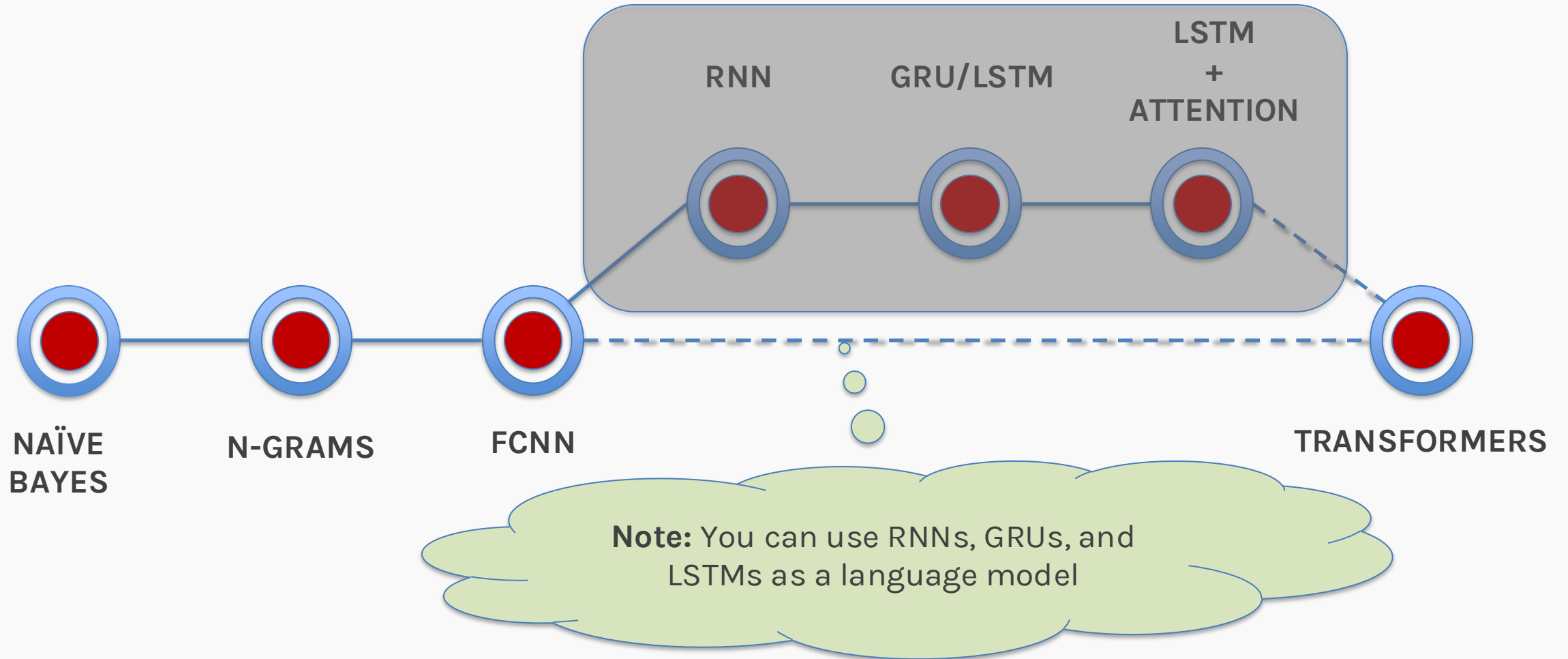
Electricity is all you need



Electricity is all you need



Language modeling



Limitations

Language Model Wishlist?

- Need long context
- We want to have strong **contextual** relations between words

FCNN have fixed context.

Word2Vec gives us one embedding per word

Ambiguity in static embeddings

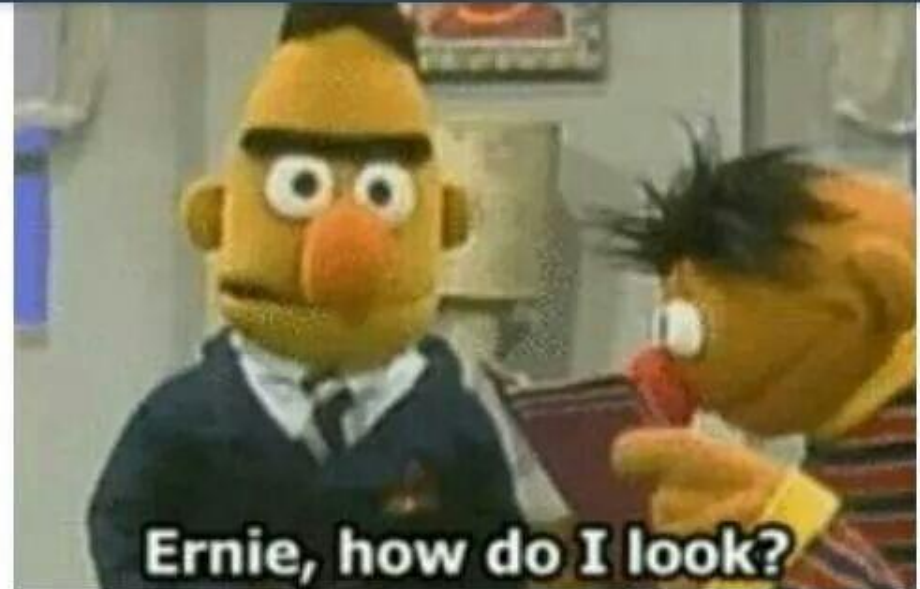
The **bank** is open on Fridays.

I went to the **bank** to take a walk by the river.

Limitations

Language Model Wishlist?

- Need long context
- We want to have strong **contextual** relations between words
- We want words to have **sequential** information
- We need an architecture that can be trained in **parallel** (non-Markovian property)



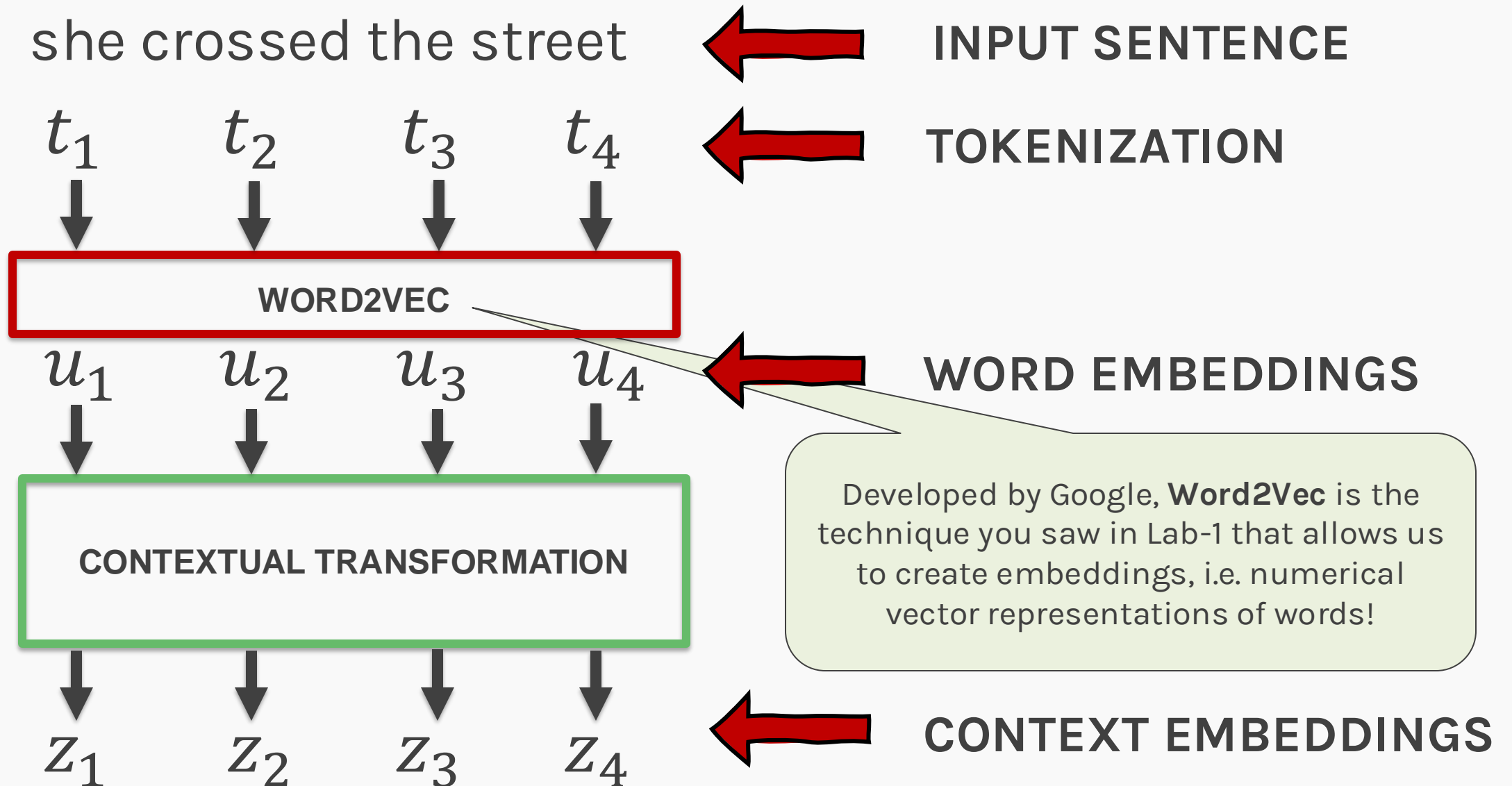
Attention: Example

Amelia was hit by a bus because she crossed the street

How do we find the context of the word ‘she’ in the sentence?

No teaching staff were harmed during the production of this lecture. This is a work of fiction. Names, characters, places and incidents either are products of the professor's imagination or are used fictitiously. Any resemblance to actual events or locales or persons, living or dead, is entirely coincidental.

Contextual Embeddings



Attention: Example

💡 IDEA #1: Distance relationship

Amelia was hit by a bus because she crossed the street

The diagram shows the sentence "Amelia was hit by a bus because she crossed the street". The word "she" is highlighted with a blue box. Four blue curved arrows originate from the box and point to the words "hit", "by", "a", and "bus". A large red curved arrow points from the word "she" back to the beginning of the sentence, labeled "True context".

❌ This idea does not work because context can be **unevenly** spread out in a sentence

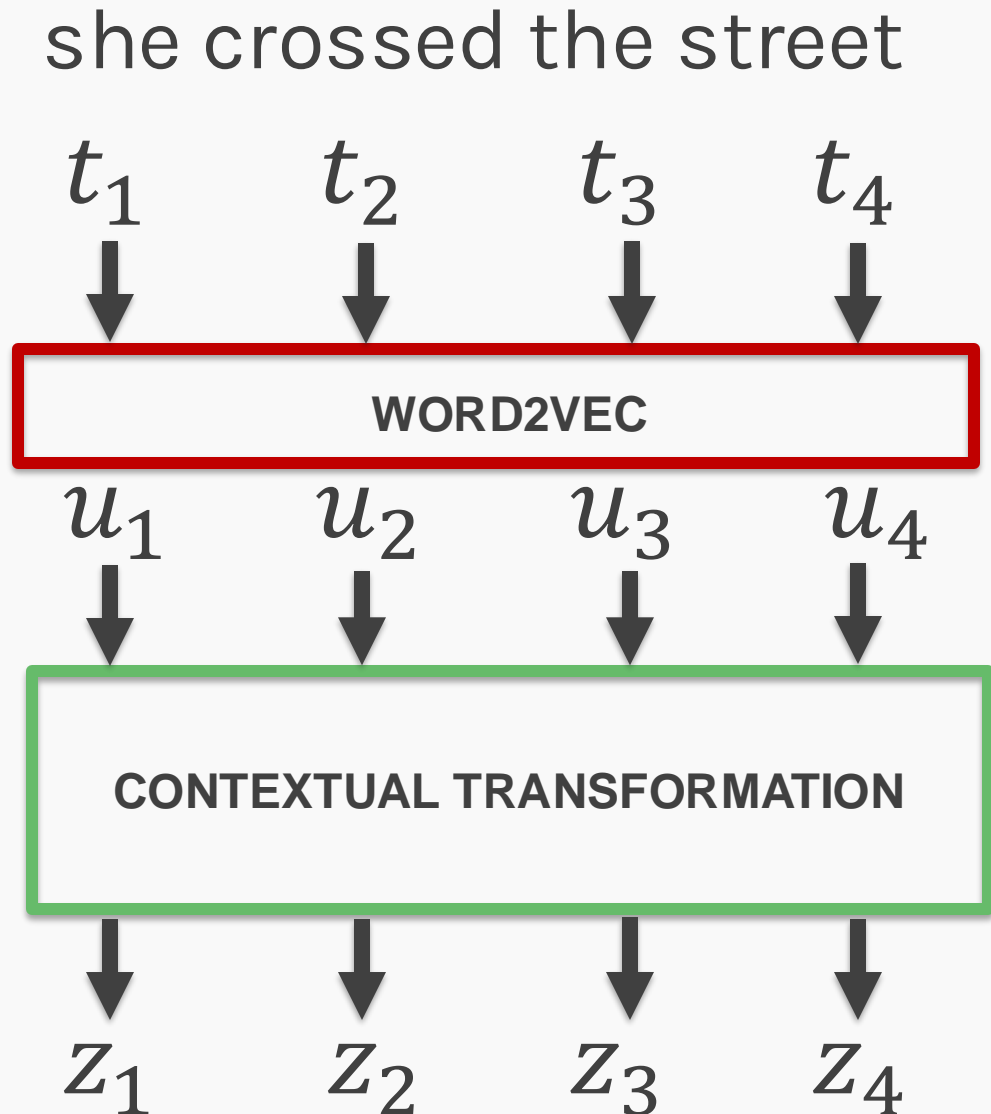
Attention: Example

Amelia was hit by a bus because she crossed the street

We would want the word ‘she’ to be related to every other word in the sentence, irrespective of the distance.

One way to do that would be to make the contextual embedding of the word ‘she’ to be a linear combination of the other word embeddings.

Contextual Embeddings



Contextual embeddings, z , can be obtained by applying a linear transformation to the word embeddings, u .

So, we want

$$z_1 = \alpha_{11}u_1 + \alpha_{12}u_2 + \alpha_{13}u_3 + \alpha_{14}u_4$$

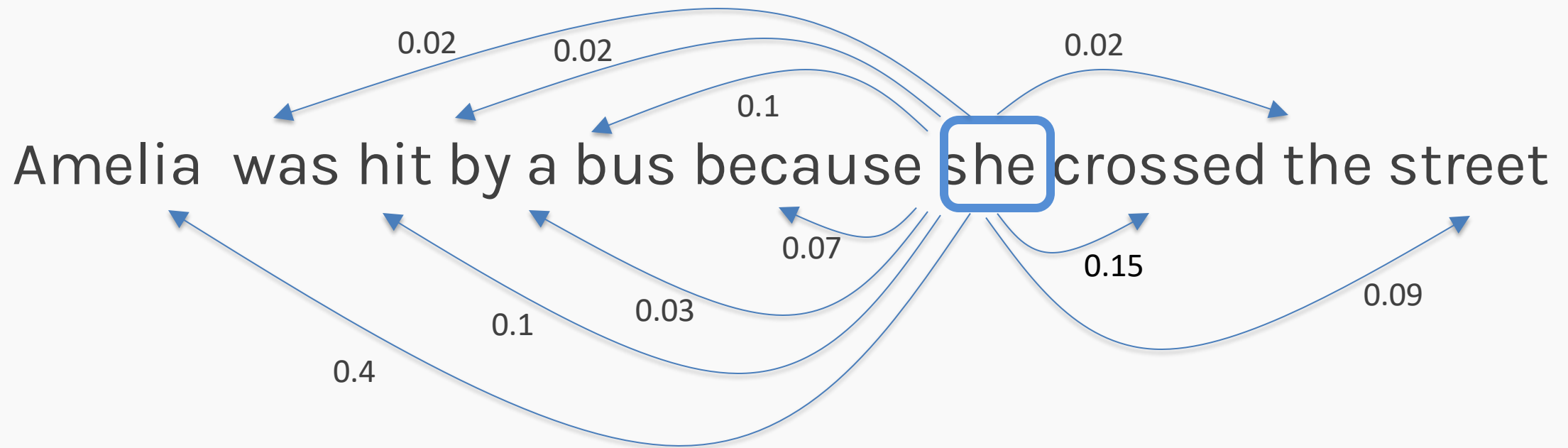
$$z_2 = \alpha_{21}u_1 + \alpha_{22}u_2 + \alpha_{23}u_3 + \alpha_{24}u_4$$

$$z_3 = \alpha_{31}u_1 + \alpha_{32}u_2 + \alpha_{33}u_3 + \alpha_{34}u_4$$

$$z_4 = \alpha_{41}u_1 + \alpha_{42}u_2 + \alpha_{43}u_3 + \alpha_{44}u_4$$

Contextual Embeddings: Coefficients

The coefficients are weights that encapsulates the degree to which new contextual embeddings should incorporate the **influence** of the other tokens



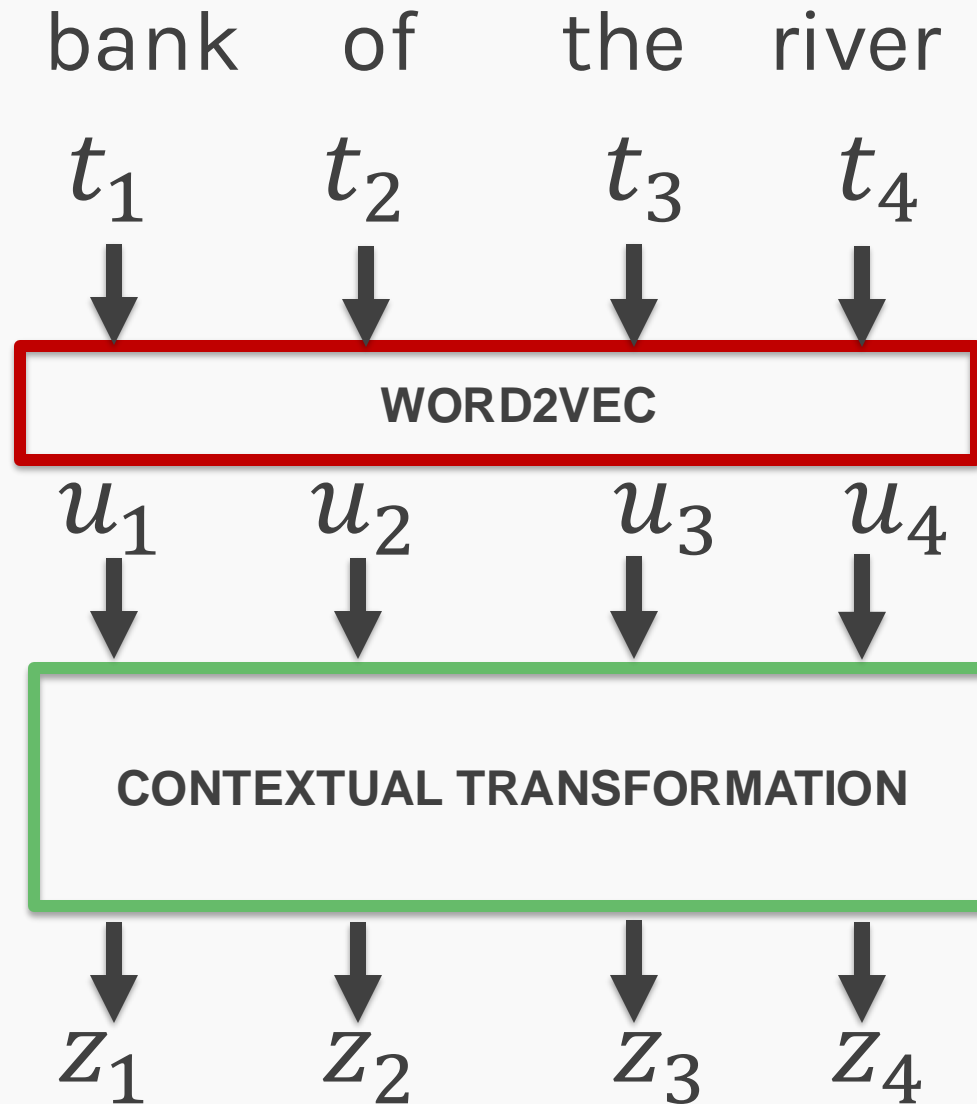
Contextual Embeddings: Coefficients

How do we calculate these coefficients?

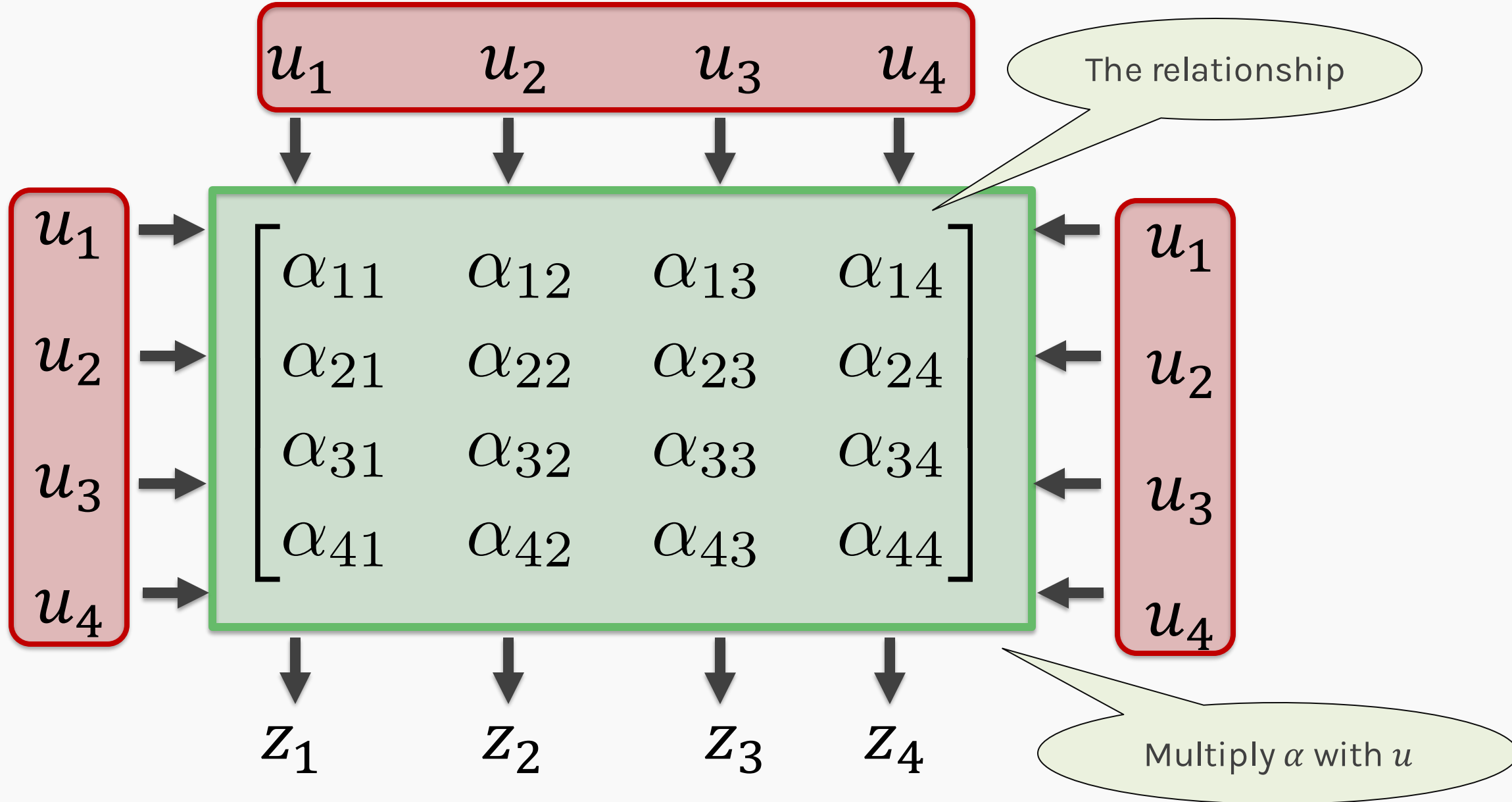
Attention: Example #2

Amelia was standing next to the bank of the river

Attention: Example #2

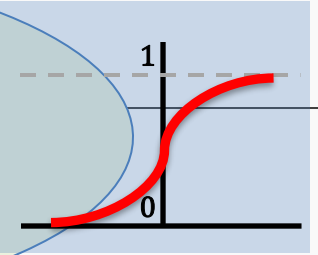


Attention: Behind the scenes



Attention: Behind the scenes

Softmax is a function that scales numbers (logits) into probabilities



$$u_1 \circ u_1 = a_{11}$$

$$u_1 \circ u_2 = a_{12}$$

$$u_1 \circ u_3 = a_{13}$$

$$u_1 \circ u_4 = a_{14}$$

softmax(.)



**WEIGHTS
NORMALIZATION**

$$\alpha_{11}$$

$$\alpha_{12}$$

$$\alpha_{13}$$

$$\alpha_{14}$$

$$\Sigma \alpha_{1i} = 1$$

Cosine similarity

Attention: Behind the scenes

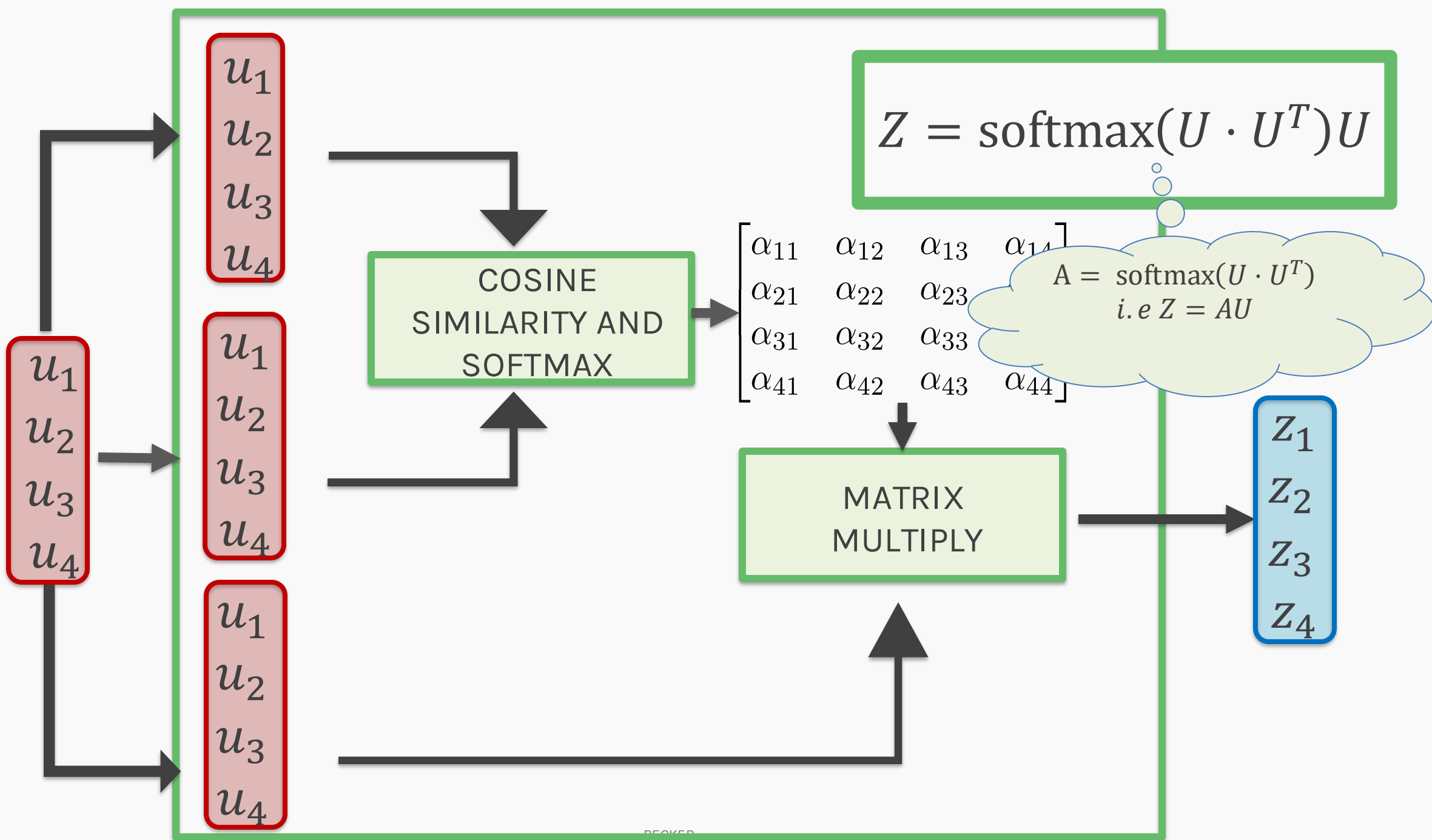
Similarly...

$$z_1 = \alpha_{11}u_1 + \alpha_{12}u_2 + \alpha_{13}u_3 + \alpha_{14}u_4$$

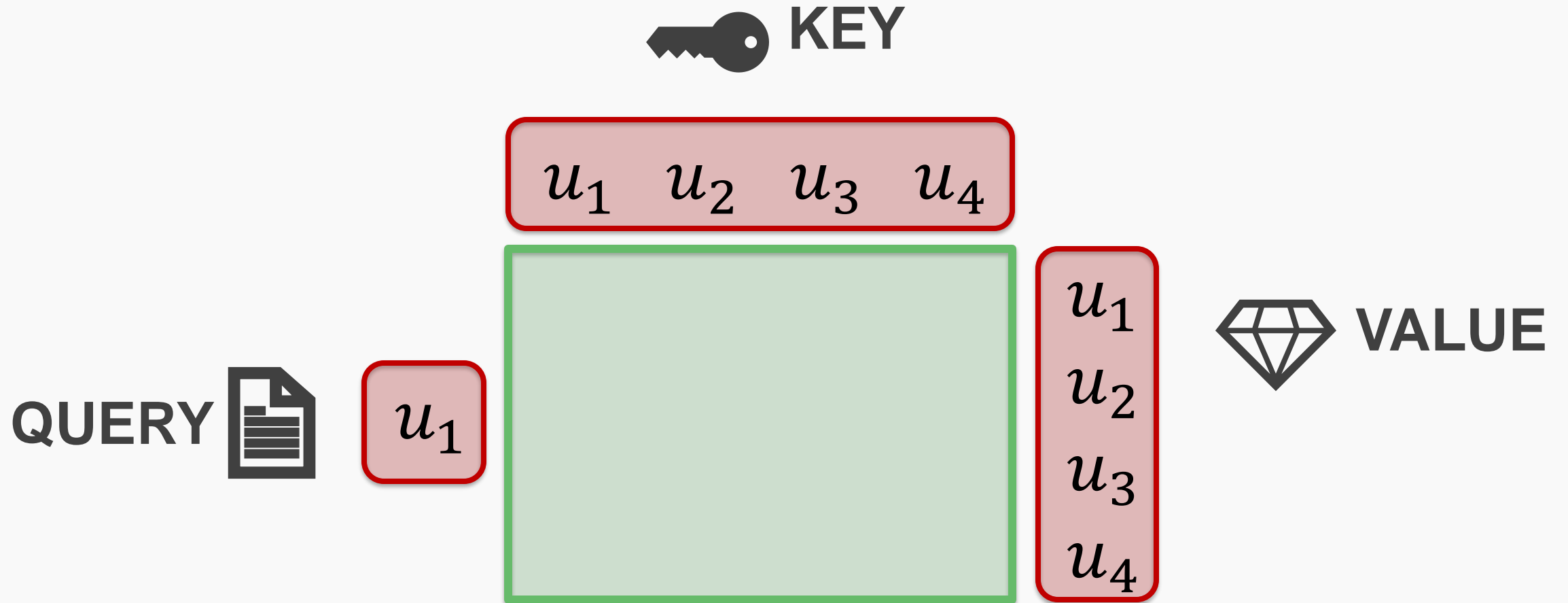
$$z_2 = \alpha_{21}u_1 + \alpha_{22}u_2 + \alpha_{23}u_3 + \alpha_{24}u_4$$

$$z_3 = \alpha_{31}u_1 + \alpha_{32}u_2 + \alpha_{33}u_3 + \alpha_{34}u_4$$

$$z_4 = \alpha_{41}u_1 + \alpha_{42}u_2 + \alpha_{43}u_3 + \alpha_{44}u_4$$

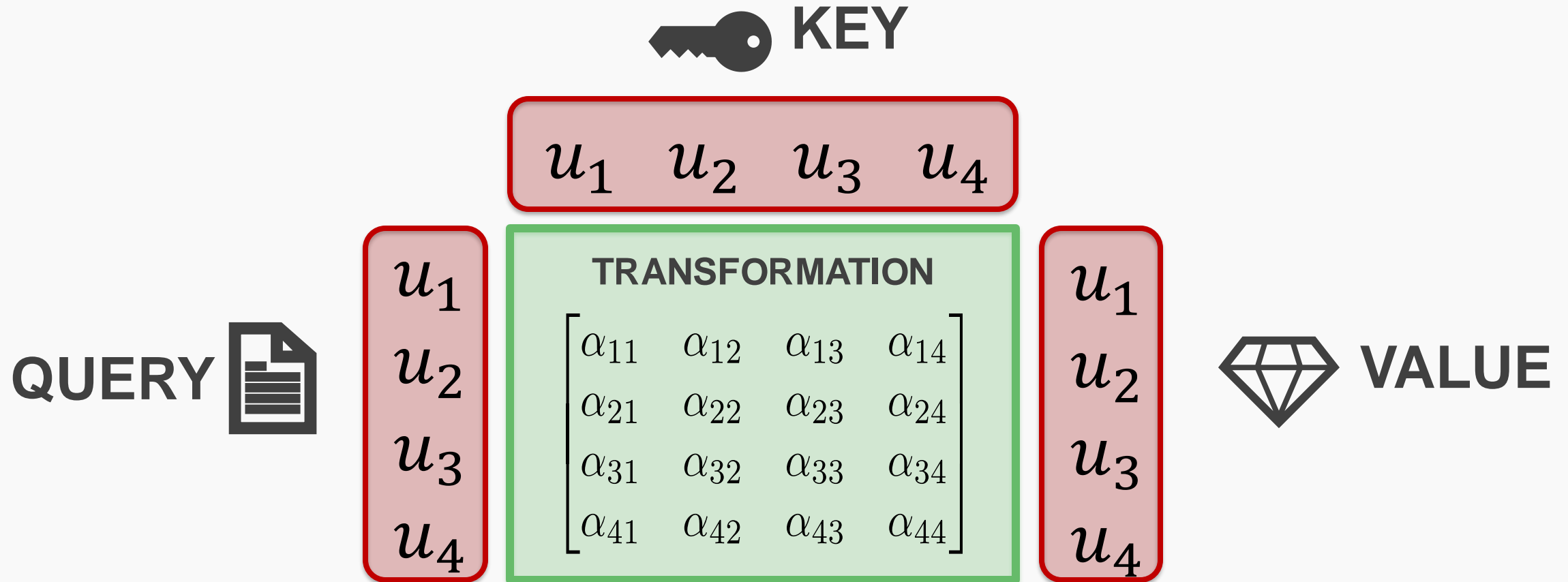


Attention: Database Analogy

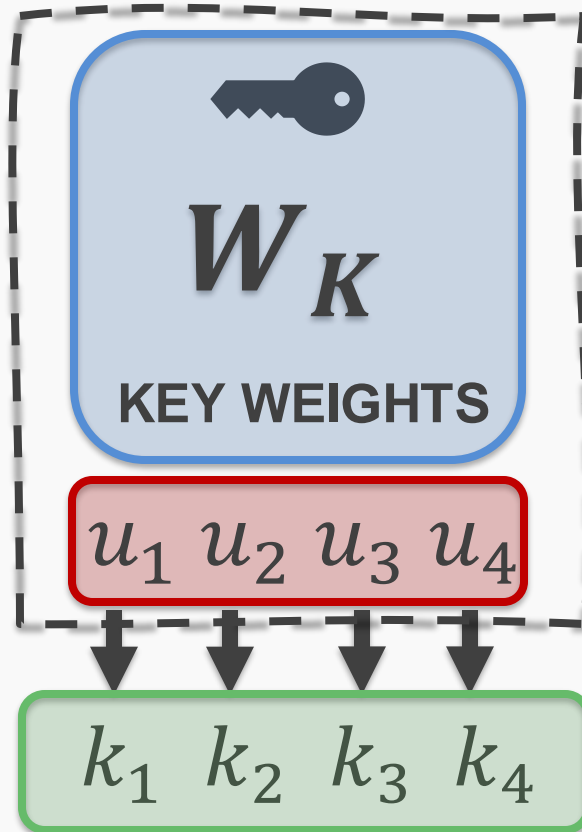
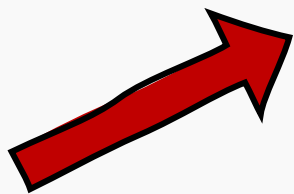


Attention: Database Analogy

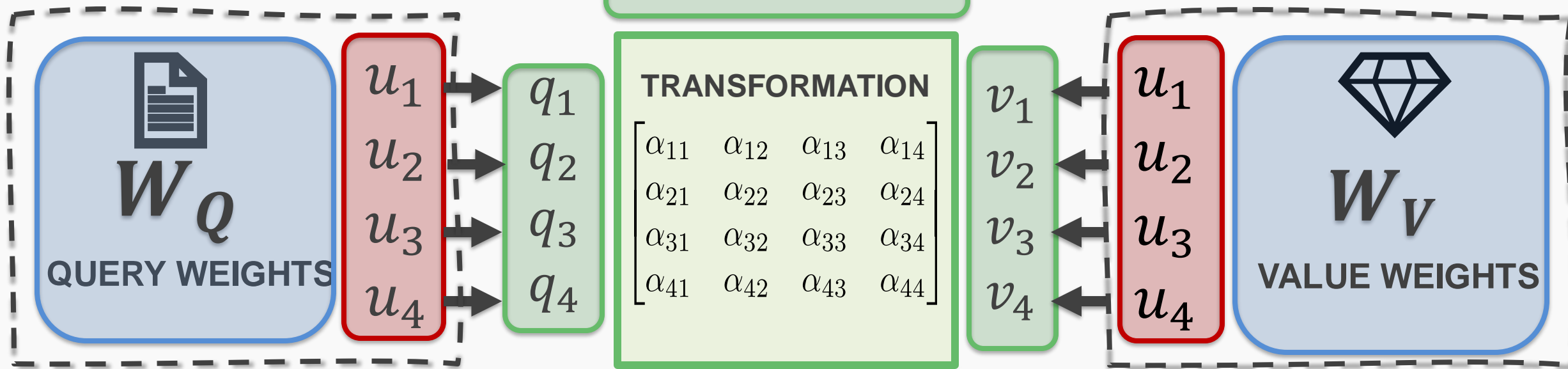
For simplicity, we stick to our database analogy of **QUERY, KEY & VALUE**



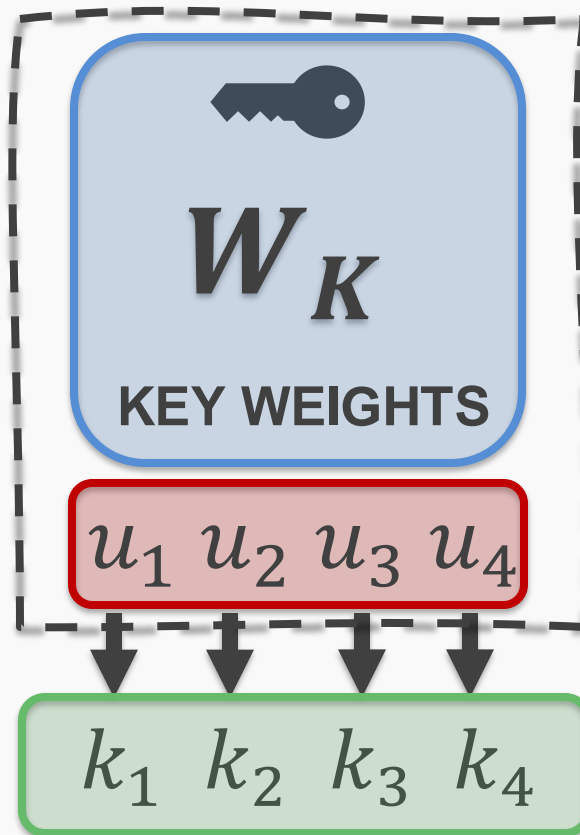
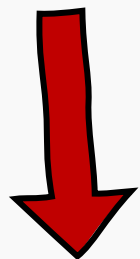
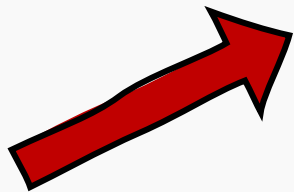
TRAINABLE WEIGHTS



TRAINABLE WEIGHTS

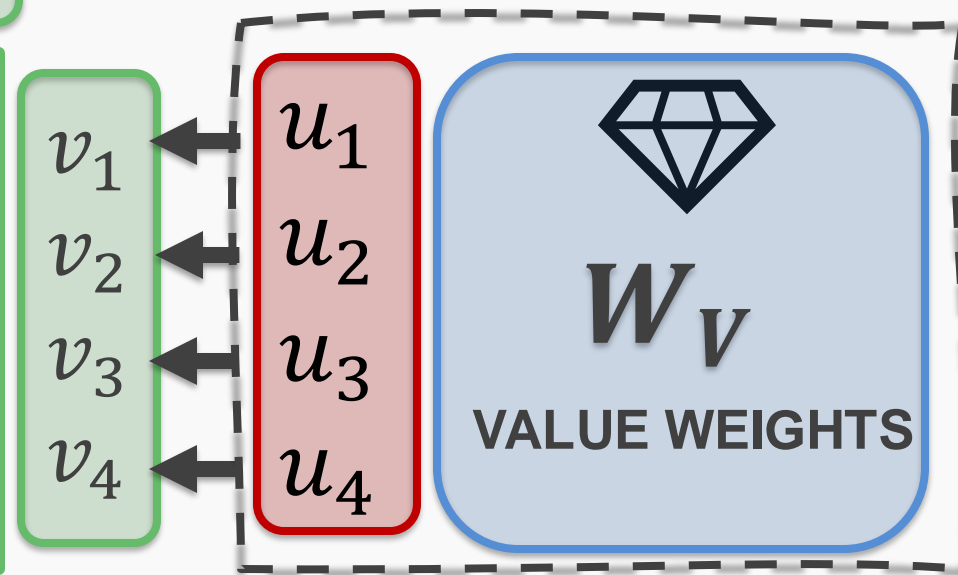
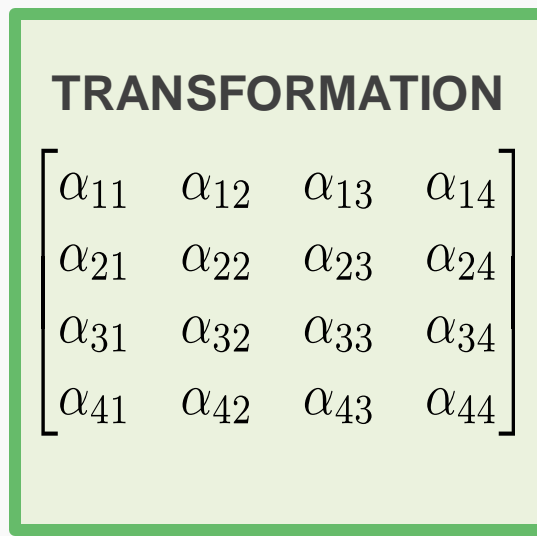
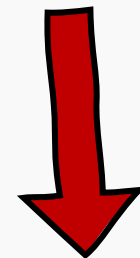


TRAINABLE WEIGHTS

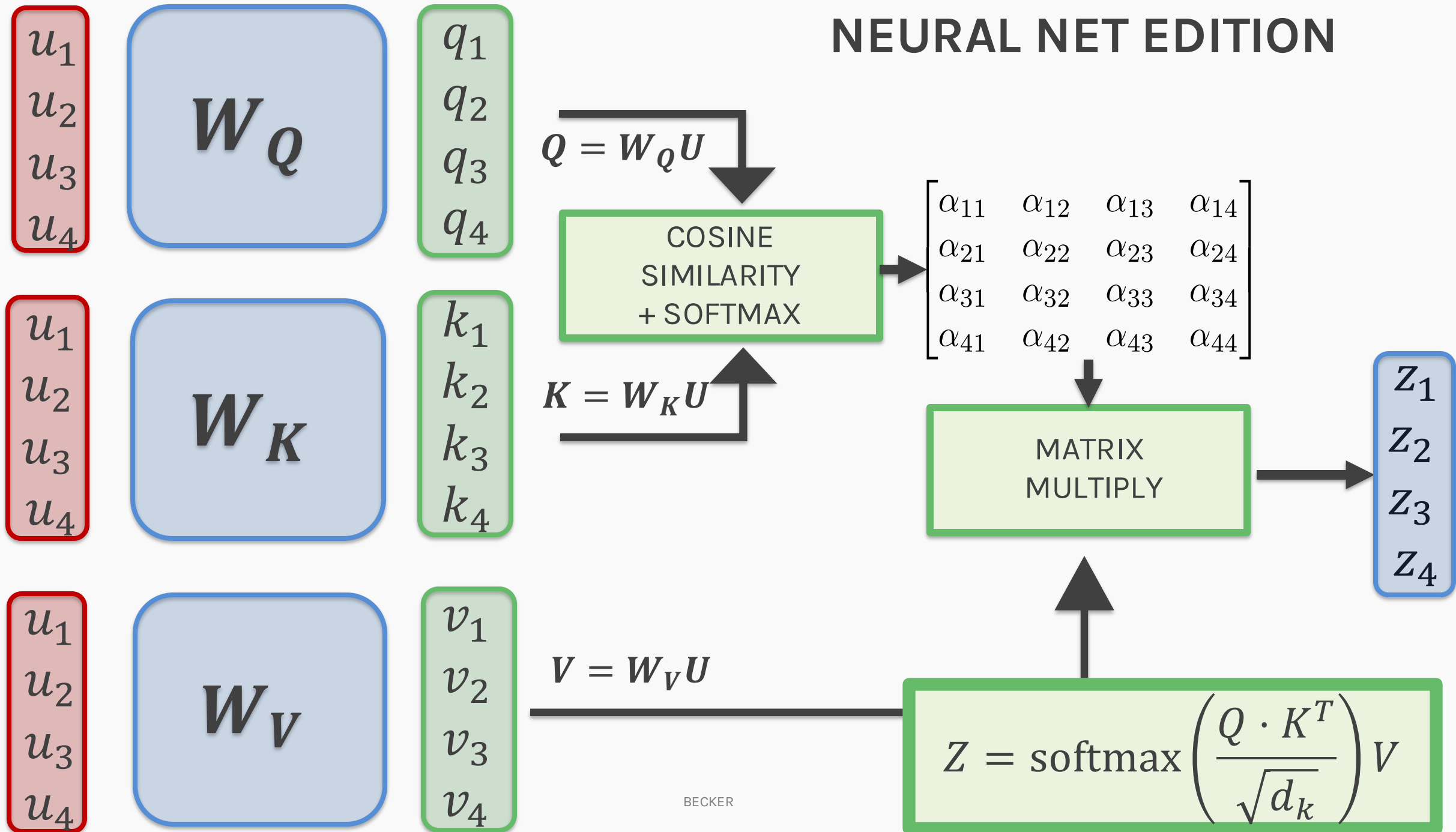


By adding trainable weights, we allow our transformations to be flexible to the task

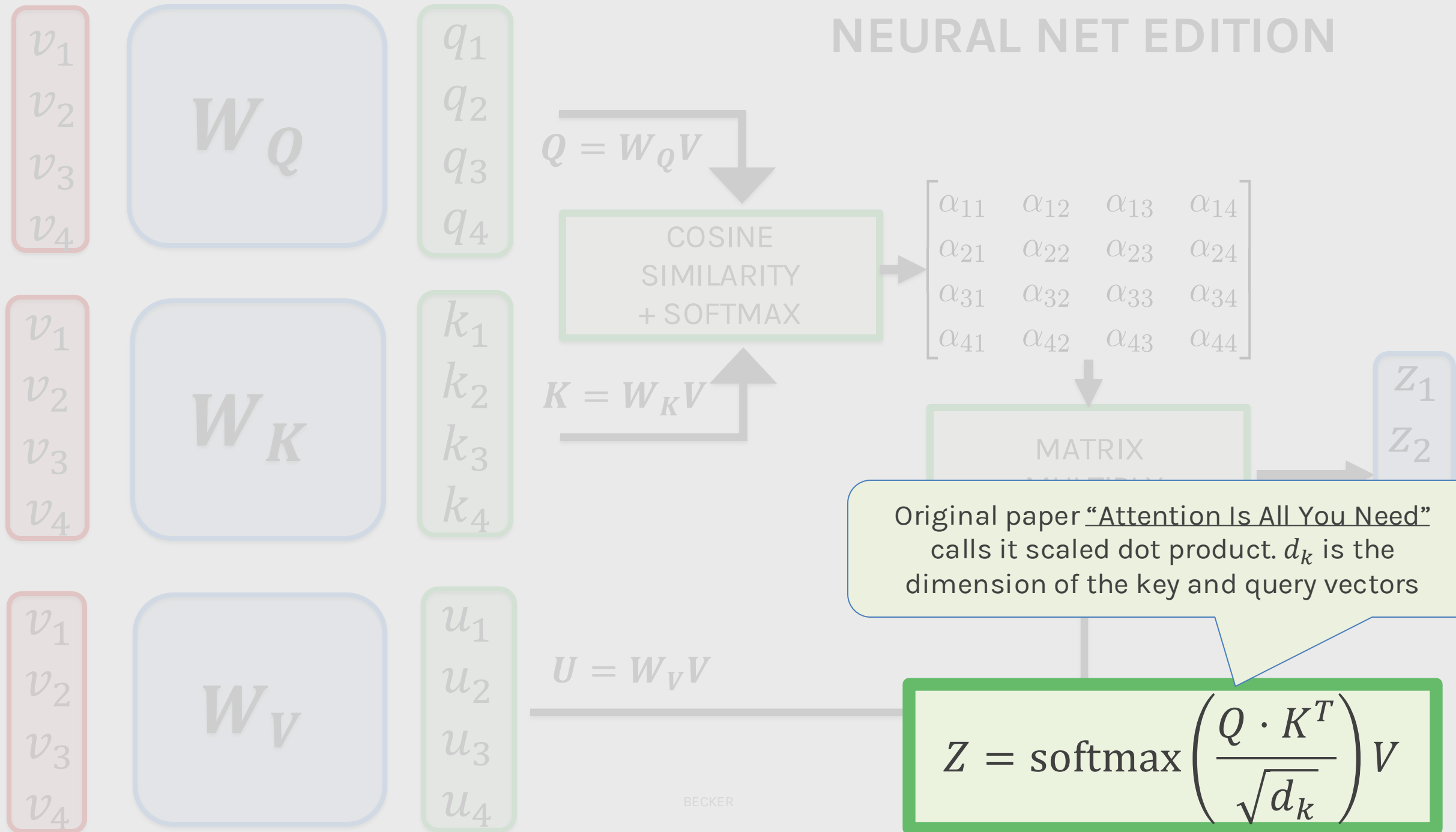
TRAINABLE WEIGHTS



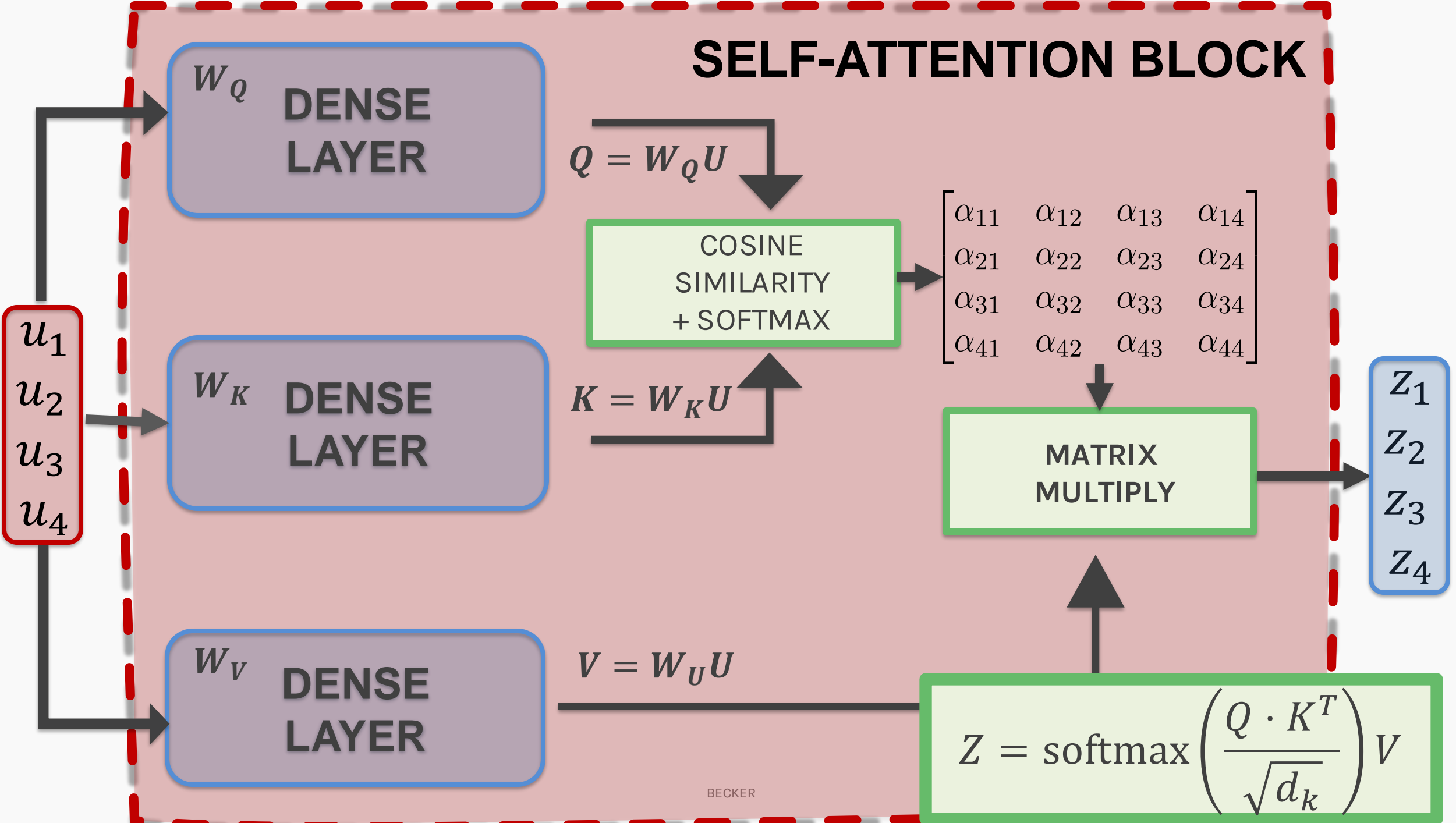
NEURAL NET EDITION



NEURAL NET EDITION



SELF-ATTENTION BLOCK



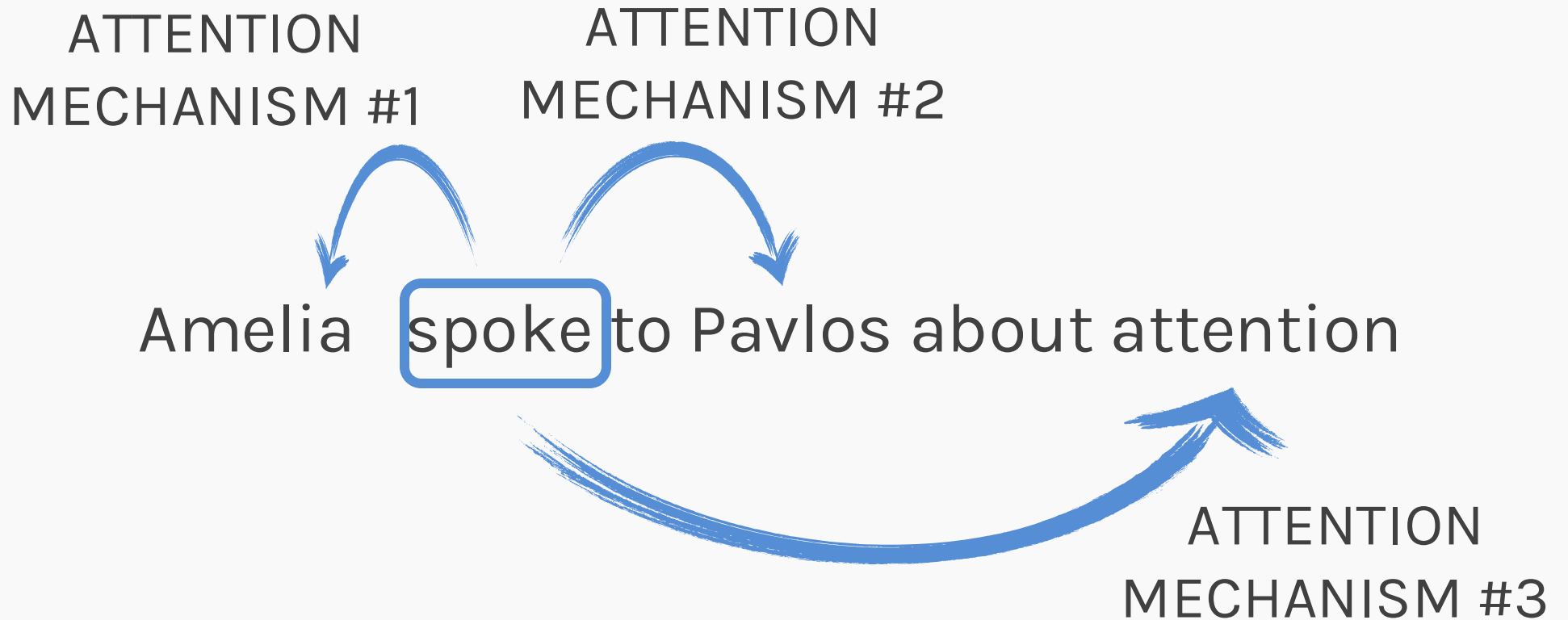
Attention

ATTENTION ISSUES?

- ~~No weights trained in the process~~
- **Attention leads to limited contextual mapping.**
- There is no positional information encoded.

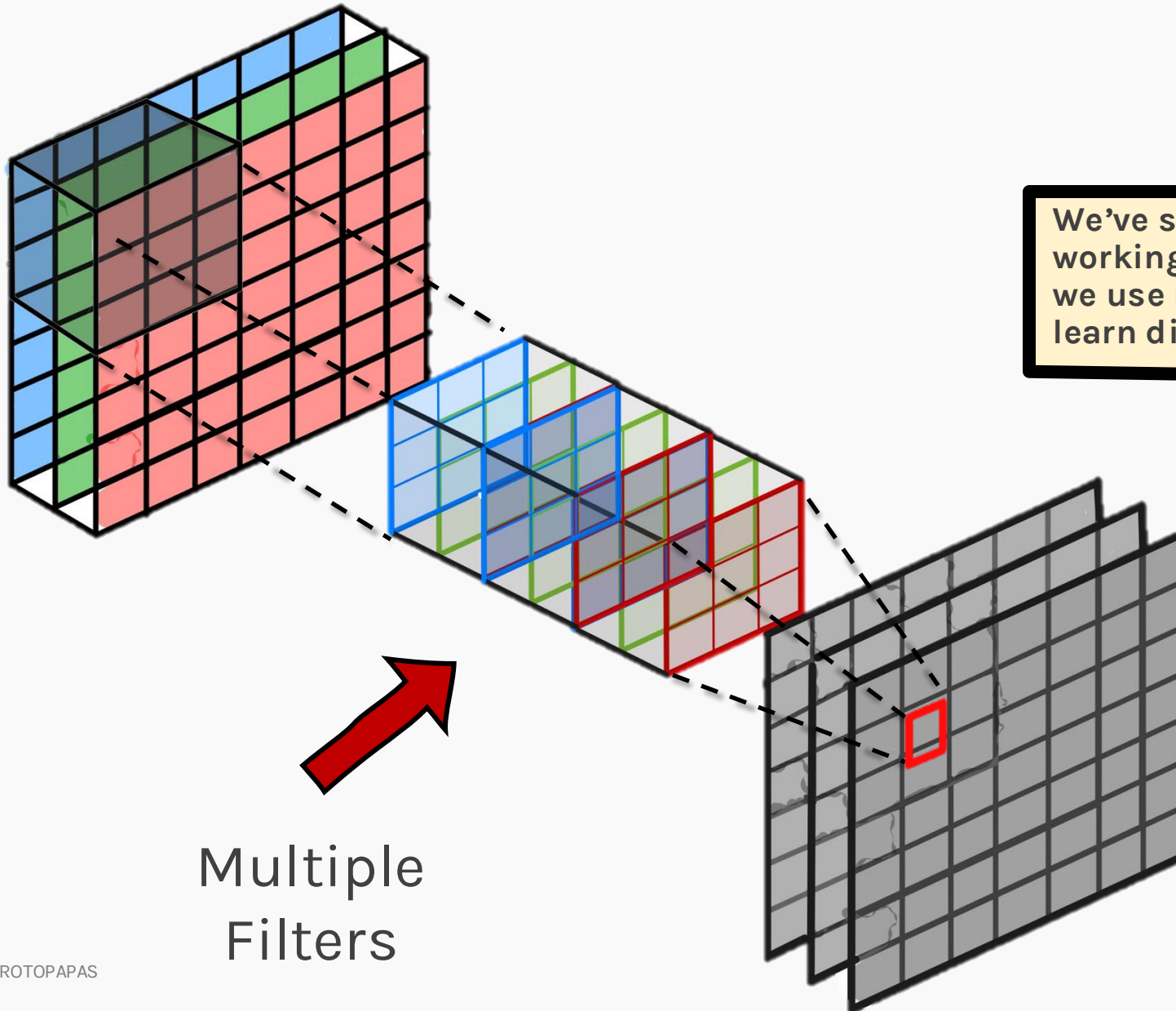


Attention: Do we have enough attention?



We need to have multiple attention mechanisms to look for different relations.

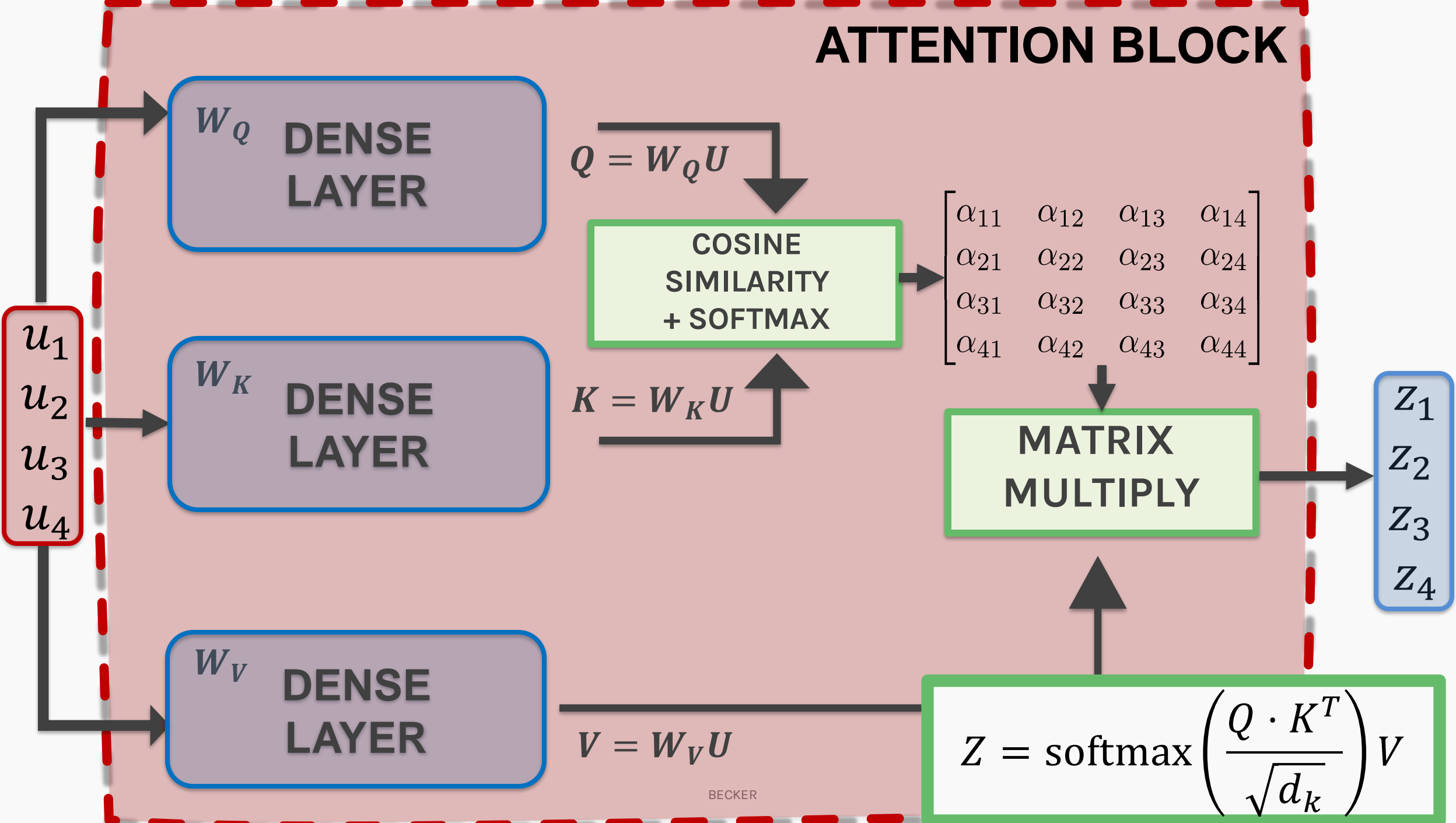
Analogy - CNN Filters



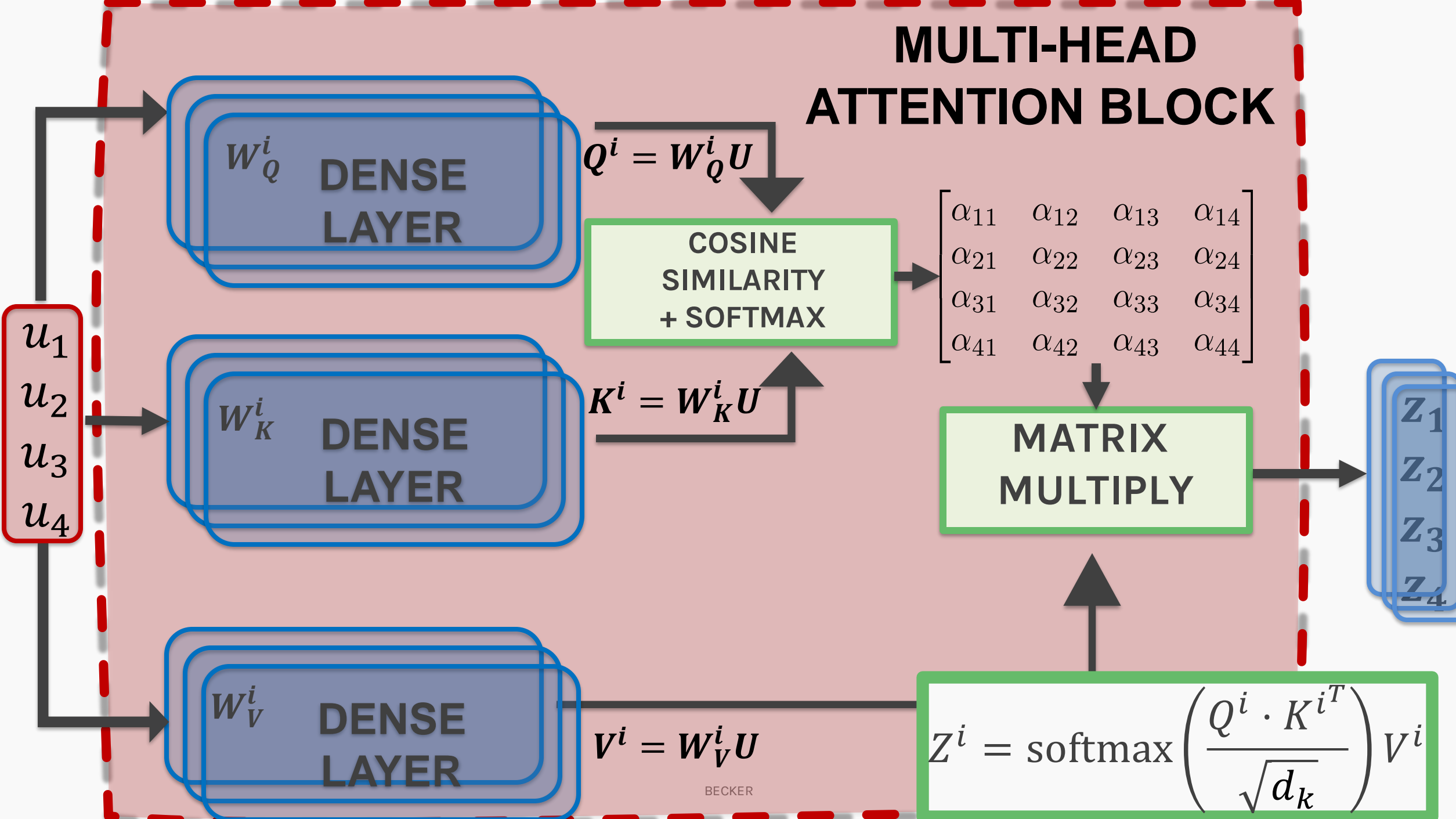
We've seen this before while working with CNNs, where we use multiple filters to learn different features



ATTENTION BLOCK



MULTI-HEAD ATTENTION BLOCK



CONTEXT EMBEDDINGS

z_1^* z_2^* z_3^* z_4^*

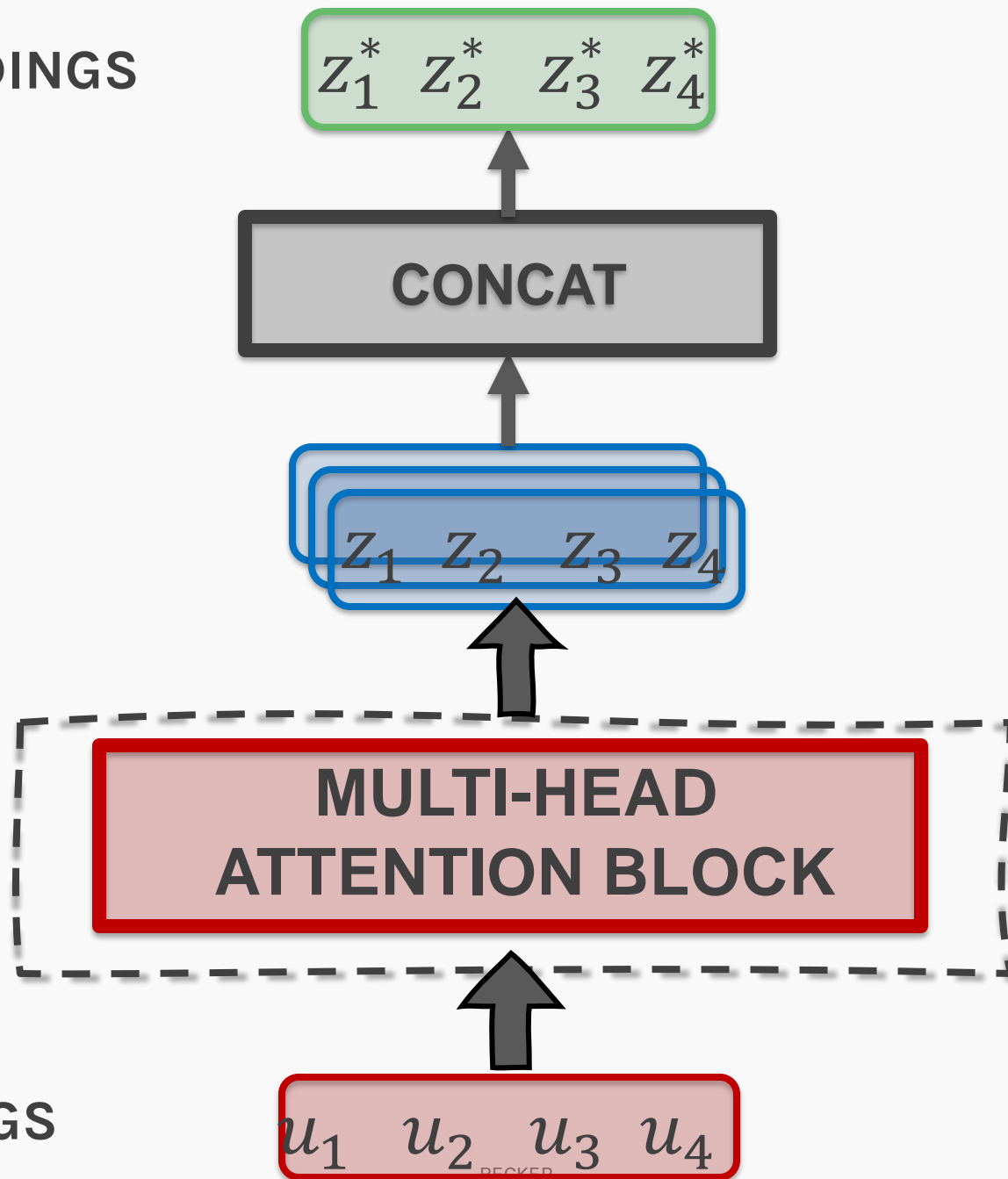
CONCAT

z_1 z_2 z_3 z_4

MULTI-HEAD
ATTENTION BLOCK

WORD EMBEDDINGS

u_1 u_2 u_3 u_4



We could even add skip connections to avoid the issue of vanishing gradients

Why not add Normalization layers to help gradients flow better during back propagation? Something like batch norm

z_1^* z_2^* z_3^* z_4^*

CONCAT

MULTI-HEAD ATTENTION

u_1 u_2 u_3 u_4



CONTEXT EMBEDDINGS

z_1^* z_2^* z_3^* z_4^*

TRANSFORMER LAYER

ADD+NORMALIZE

DENSE x 2 + RELU

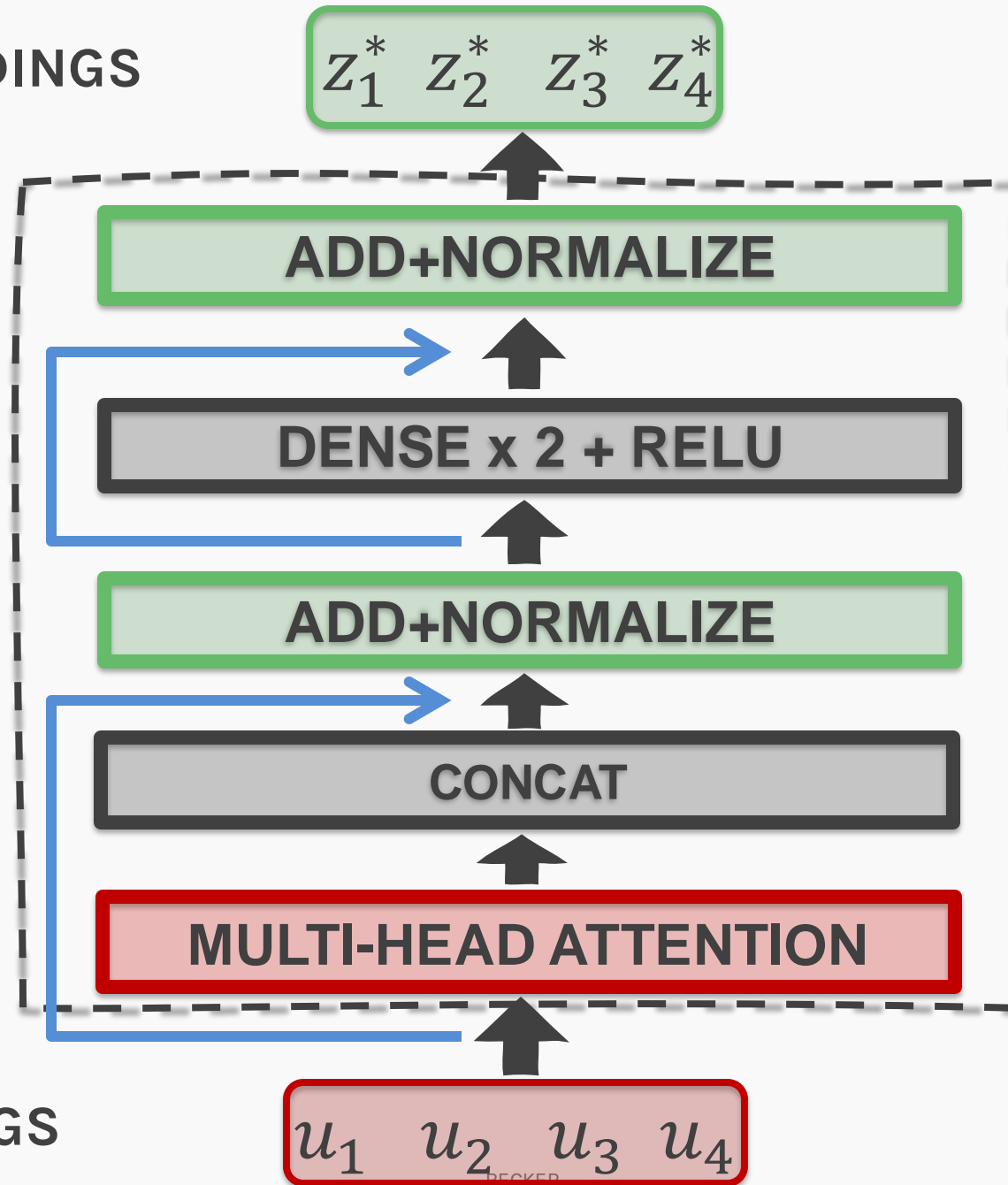
ADD+NORMALIZE

CONCAT

MULTI-HEAD ATTENTION

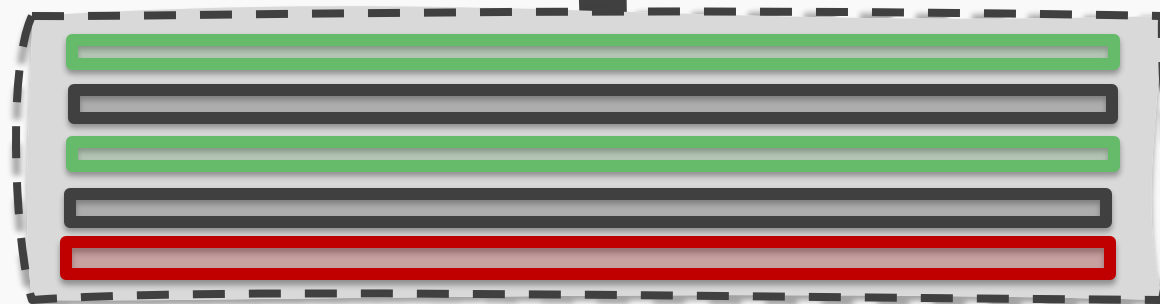
WORD EMBEDDINGS

u_1 u_2 u_3 u_4



Why not stack multiple transformer layers?

z_1^* z_2^* z_3^* z_4^*



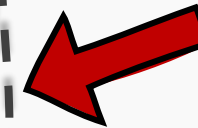
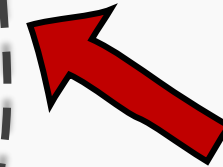
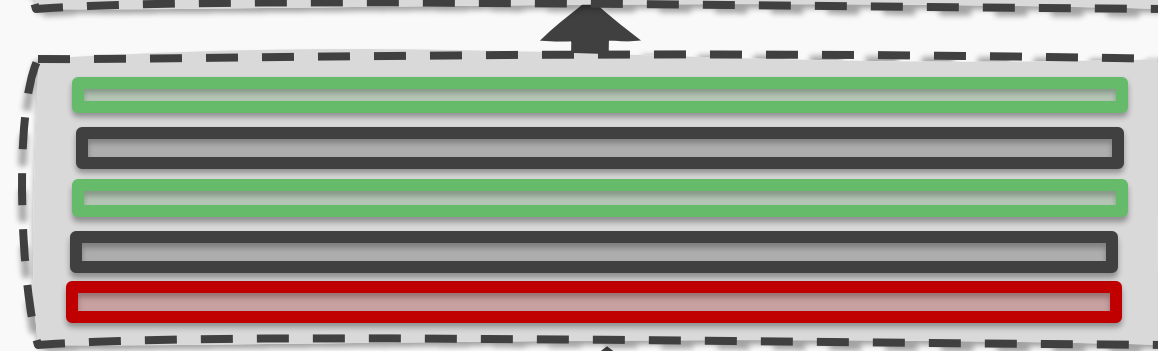
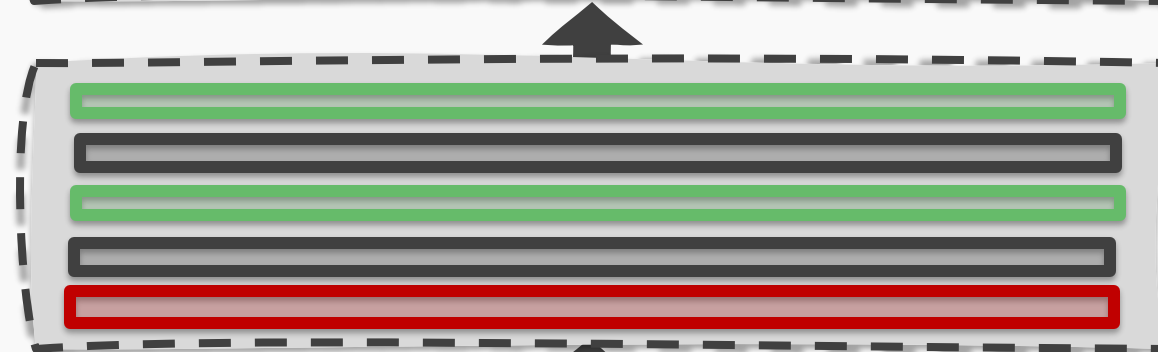
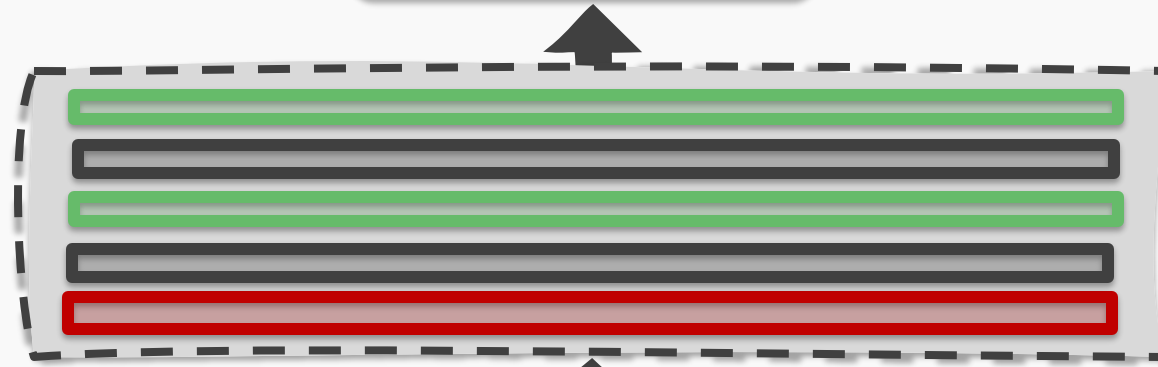
TRANSFORMER LAYER

u_1 u_2 u_3 u_4



CONTEXT EMBEDDINGS

z_1^* z_2^* z_3^* z_4^*

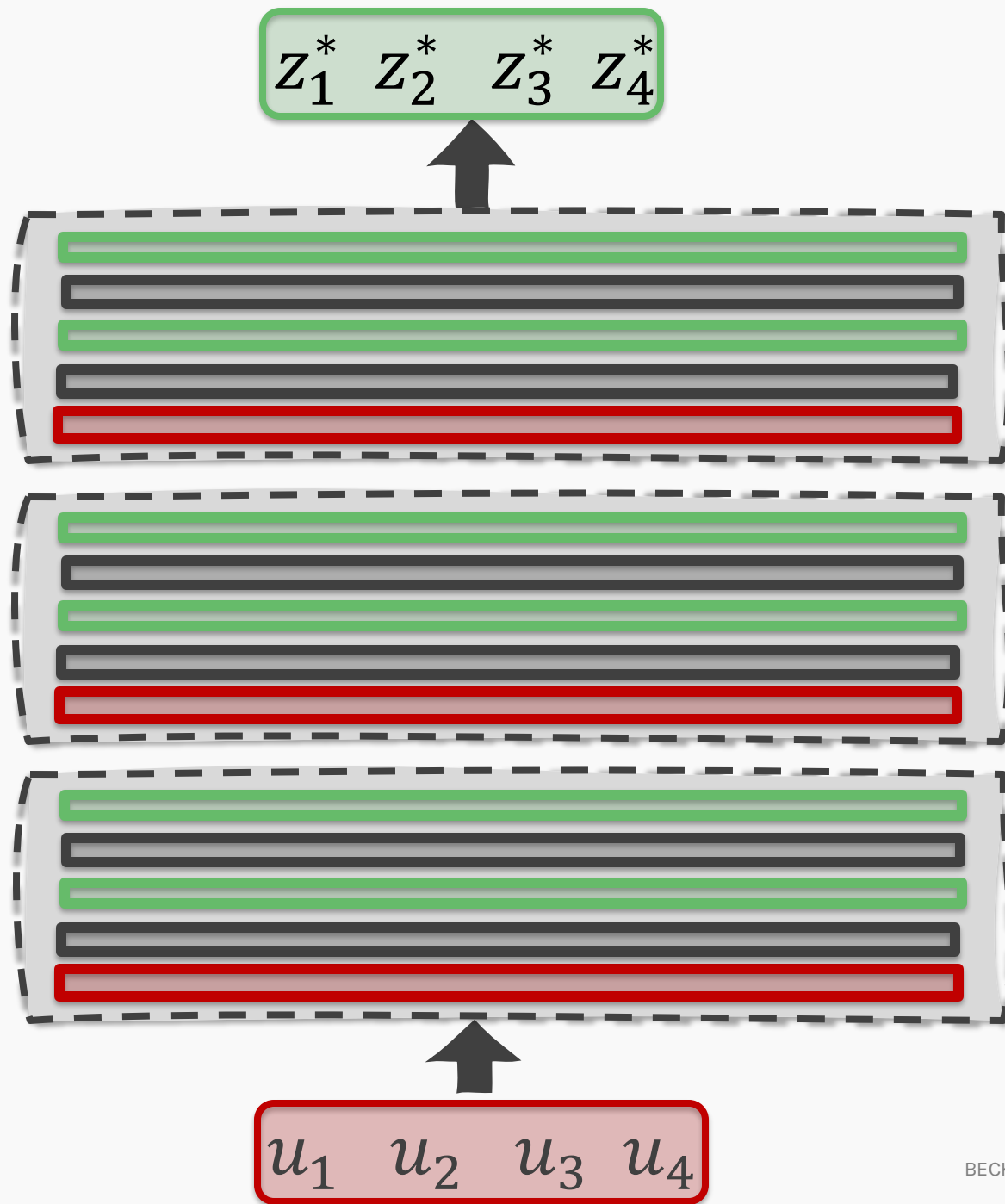


TRANSFORMER
LAYERS

WORD EMBEDDINGS

u_1 u_2 u_3 u_4

BOOKER



Now we can potentially look for richer contextual mappings in a sentence without worrying about vanishing gradients



Attention

MULTI-HEAD ATTENTION ISSUES?

- ~~No weights trained in the process~~
- ~~Attention leads to limited contextual mapping~~
- There is no positional information encoded



Amelia spoke to Pavlos about attention



Attention

MULTI-HEAD ATTENTION ISSUES?

- ~~No weights trained in the process~~
- ~~Attention leads to limited contextual mapping~~
- There is no positional information encoded



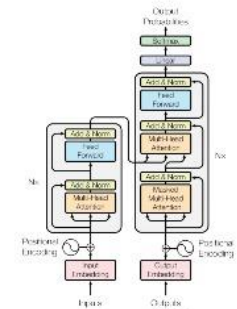
Pavlos spoke to Amelia about attention



What we want

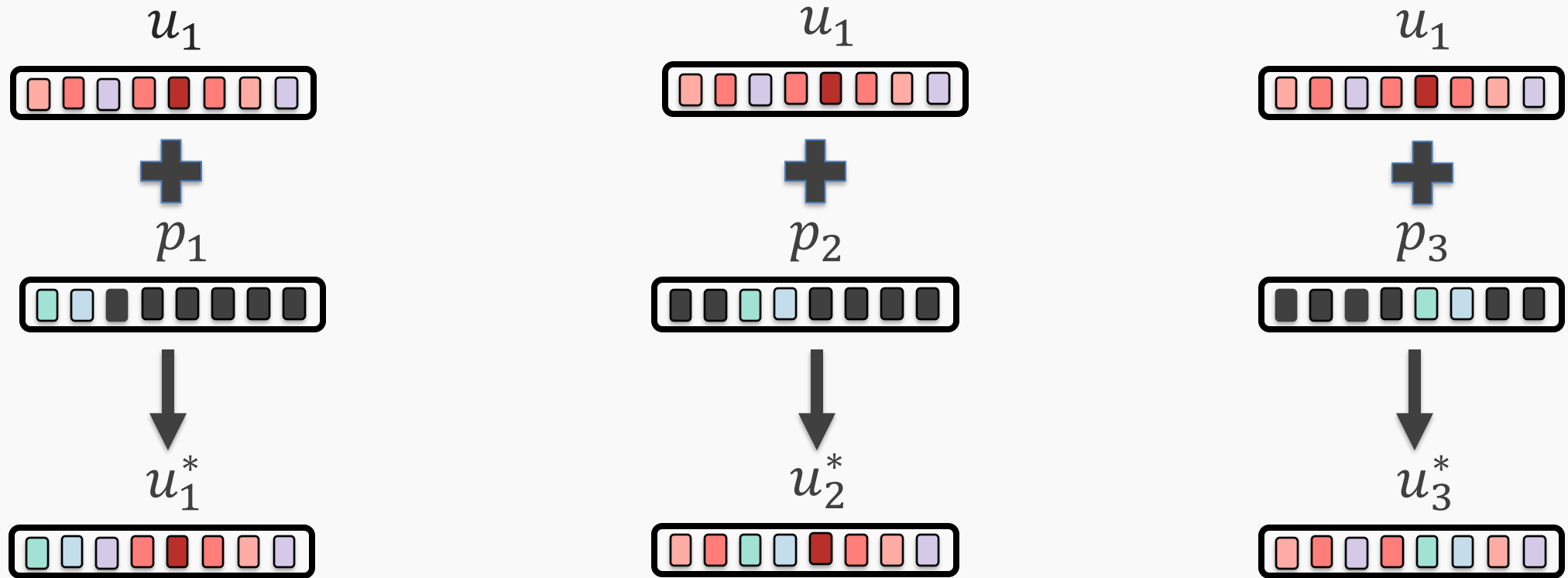
LANGUAGE MODEL WISHLIST

- **Position** and **order** of words are the essential parts of any language
- Recurrent Neural Networks (RNNs) inherently take the **order** of word into account
- Multi-head attention blocks do **not** take such an order by design, so there's the need to incorporate the order of the words separately



💡 **IDEA #255: Positional Encoding**

Positional Encoding



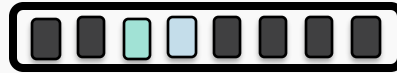
In the above case, $u_1^* \neq u_2^* \neq u_3^*$

Positional Encoding

p_1



p_2



p_3



But how do we find these positional embedding vectors?

Don't worry Ignacio, we have several options to choose from...



Positional Encoding

POSITIONAL EMBEDDING WISHLIST

- We should have a **unique** embedding for each timestep
- Relative embedding must remain **consistent** across sentences of different lengths
- It should **generalize** to longer sentences
- It should be **bounded** & **deterministic**

Positional encoding utilizes a combination of sine and cosine functions to generate unique positional embeddings.

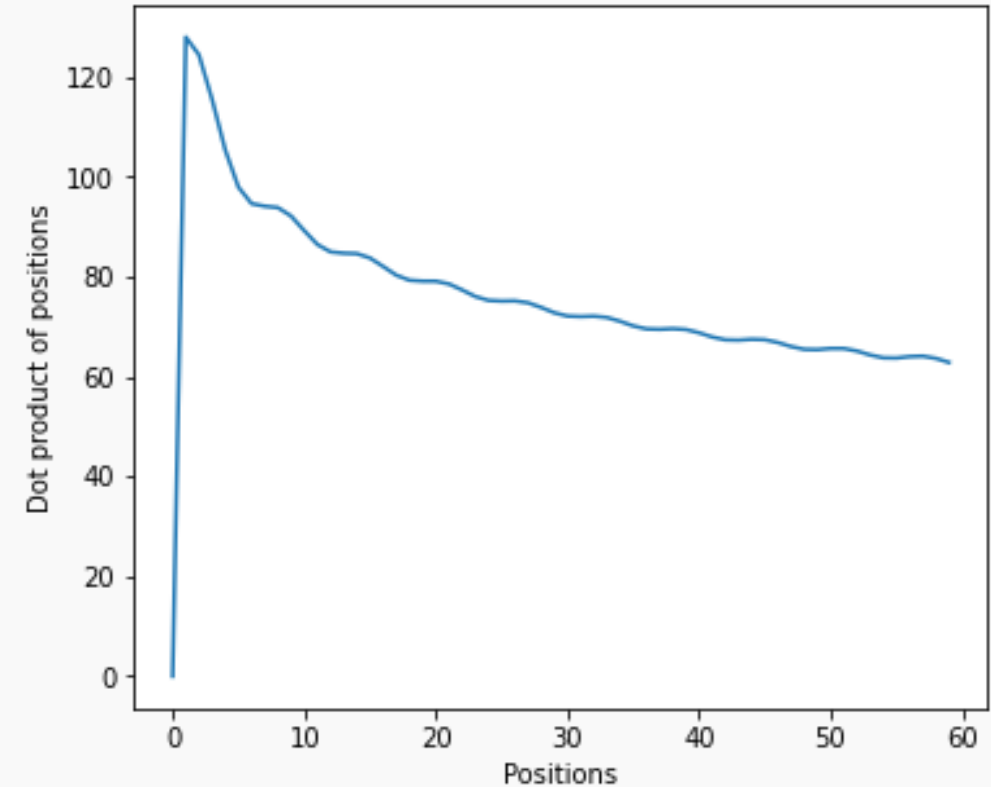


Positional Encoding

TECHNICAL DETAILS

Positional encoding is a method used to help models like transformers understand the order and relationship of words in a sentence.

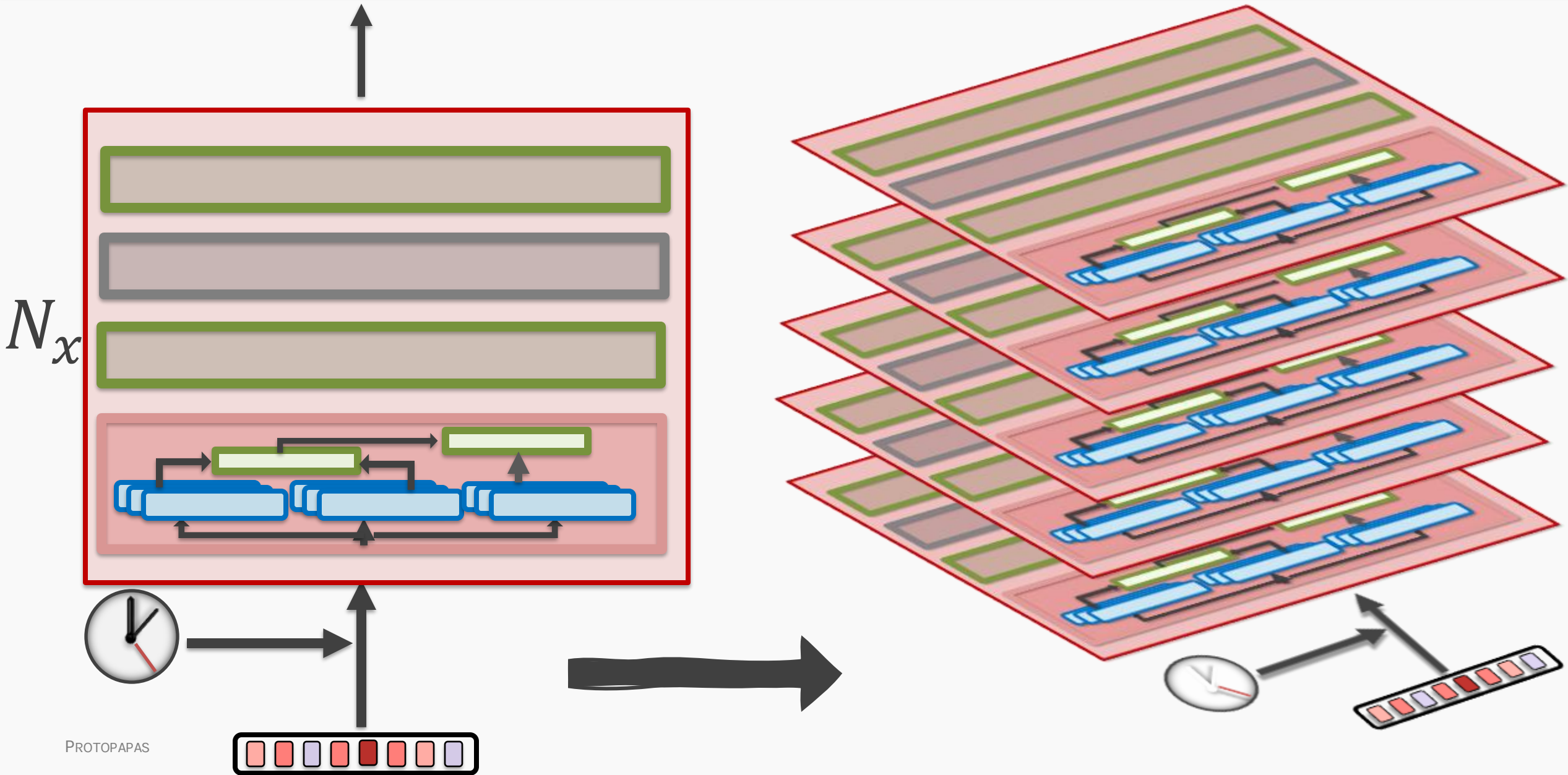
- Each position in a sentence is assigned a vector that encodes its position.
- The encoding uses a mix of sine and cosine functions to ensure each position has a unique, yet repeatable, pattern.



values of sine and cosine change with the position in the sequence

Bringing it all together

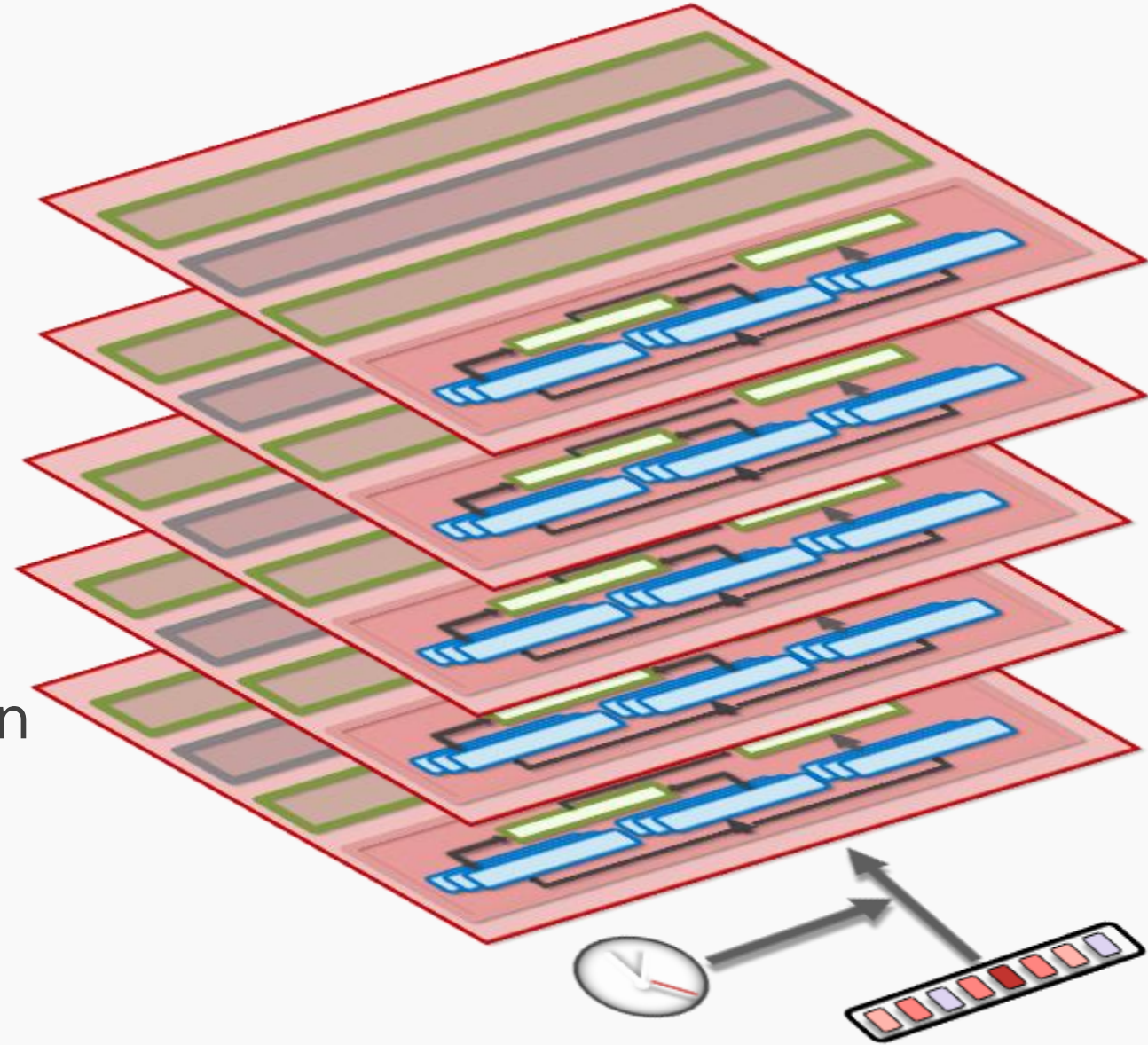
Transformer as a 3-D model



Transformers - Summary

Language Model Wishlist?

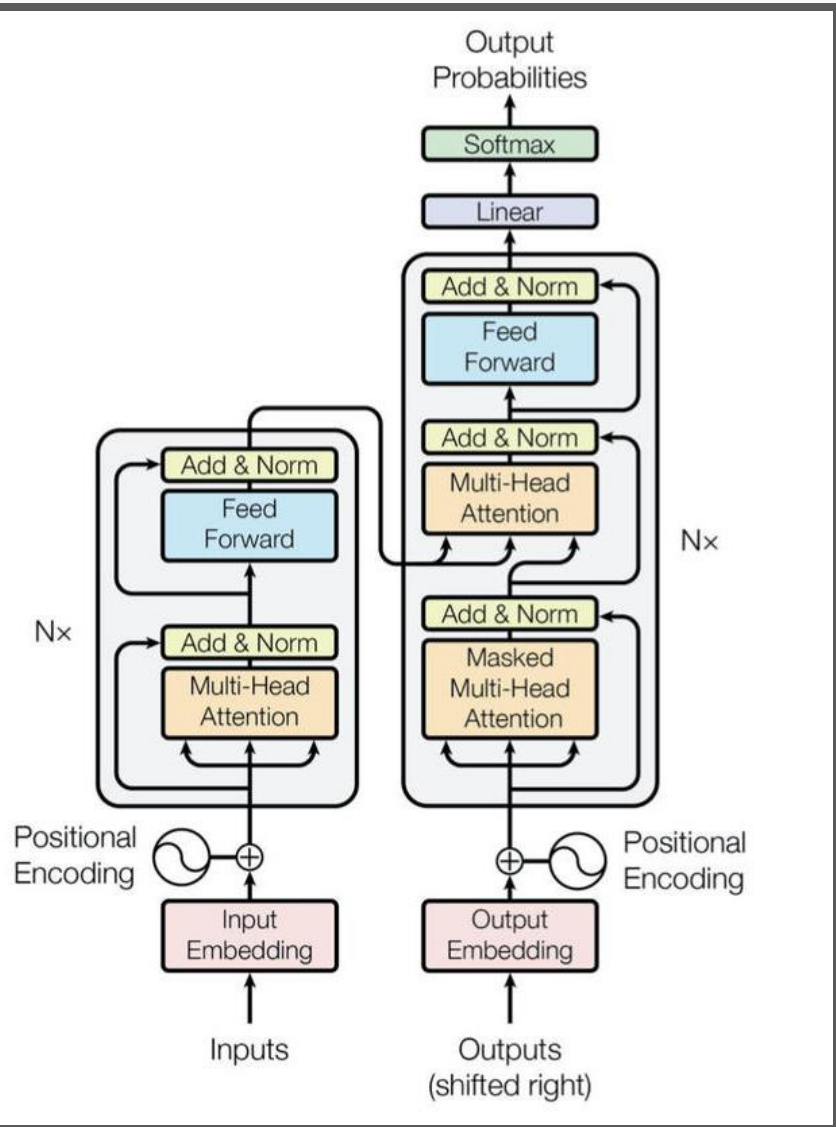
- ✓ We want to have strong contextual relations between words - **DONE**
- ✓ We want words to have sequential information - **DONE**
- ✓ We need an architecture that can be trained in parallel (non-Markovian property) - **DONE**



Transformers - Summary

Let's look at the diagram of the transformer architecture from the original paper

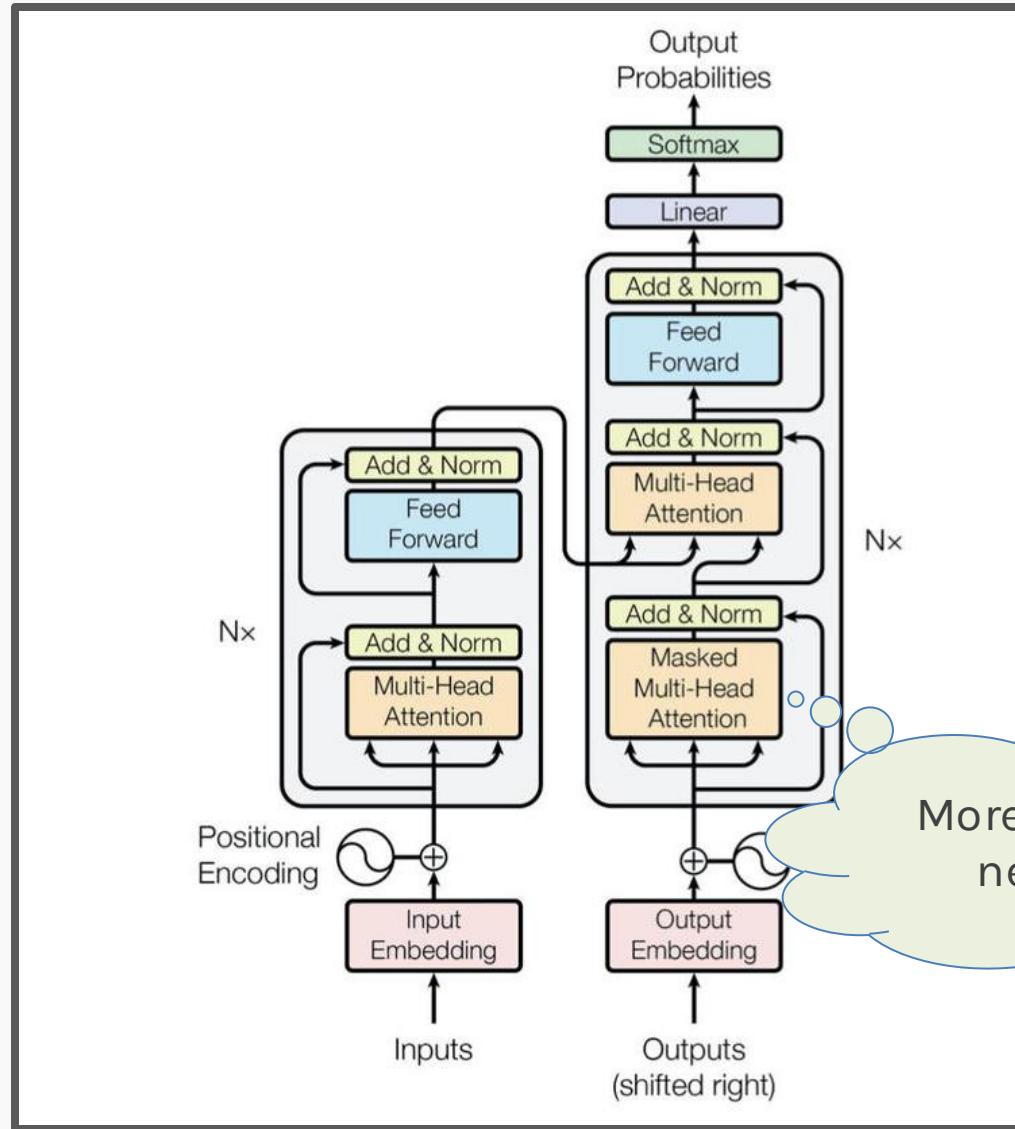
Transformers - Summary



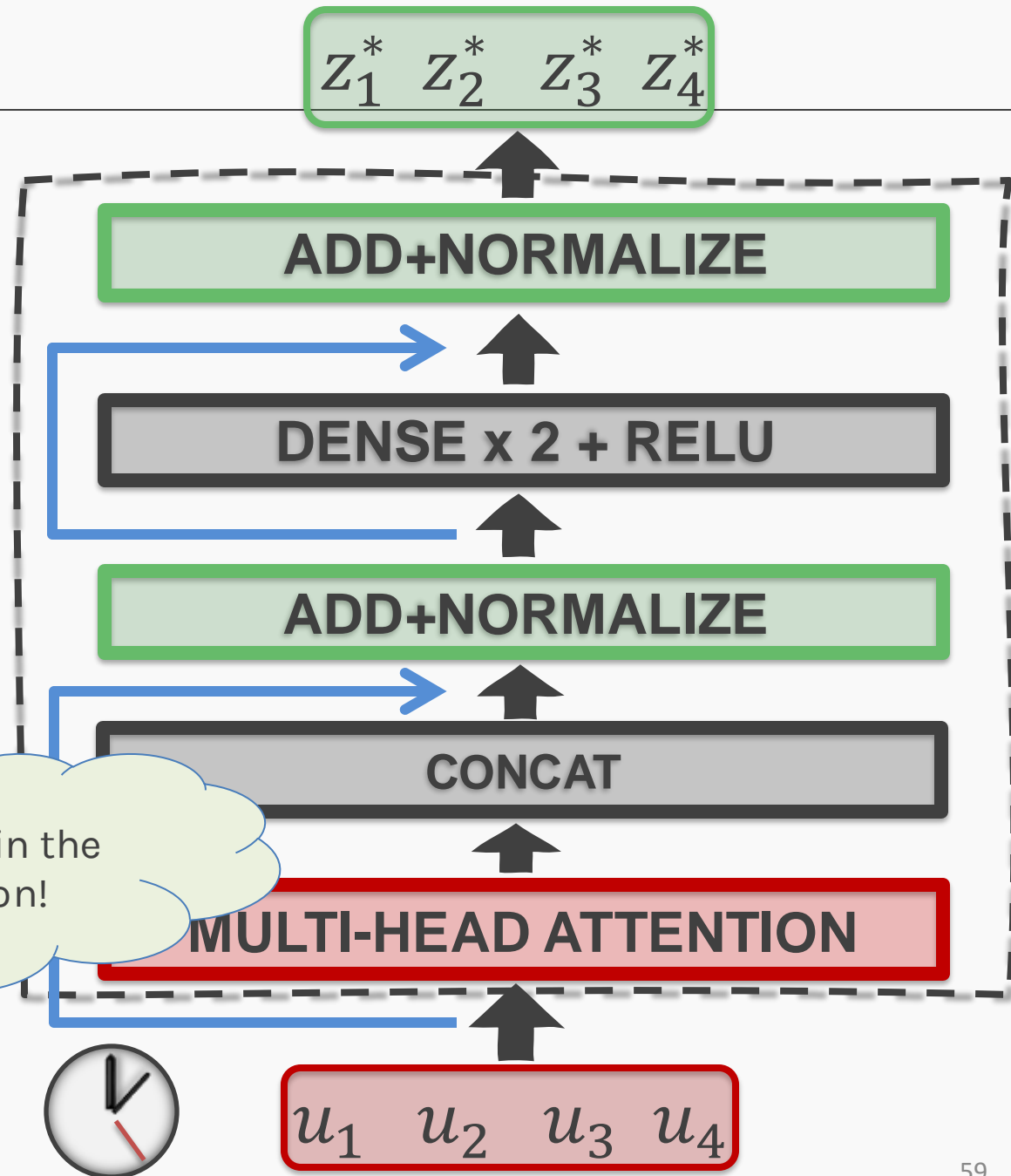
Let's look at the diagram of the transformer architecture from the original paper

And now, compare it to what we have

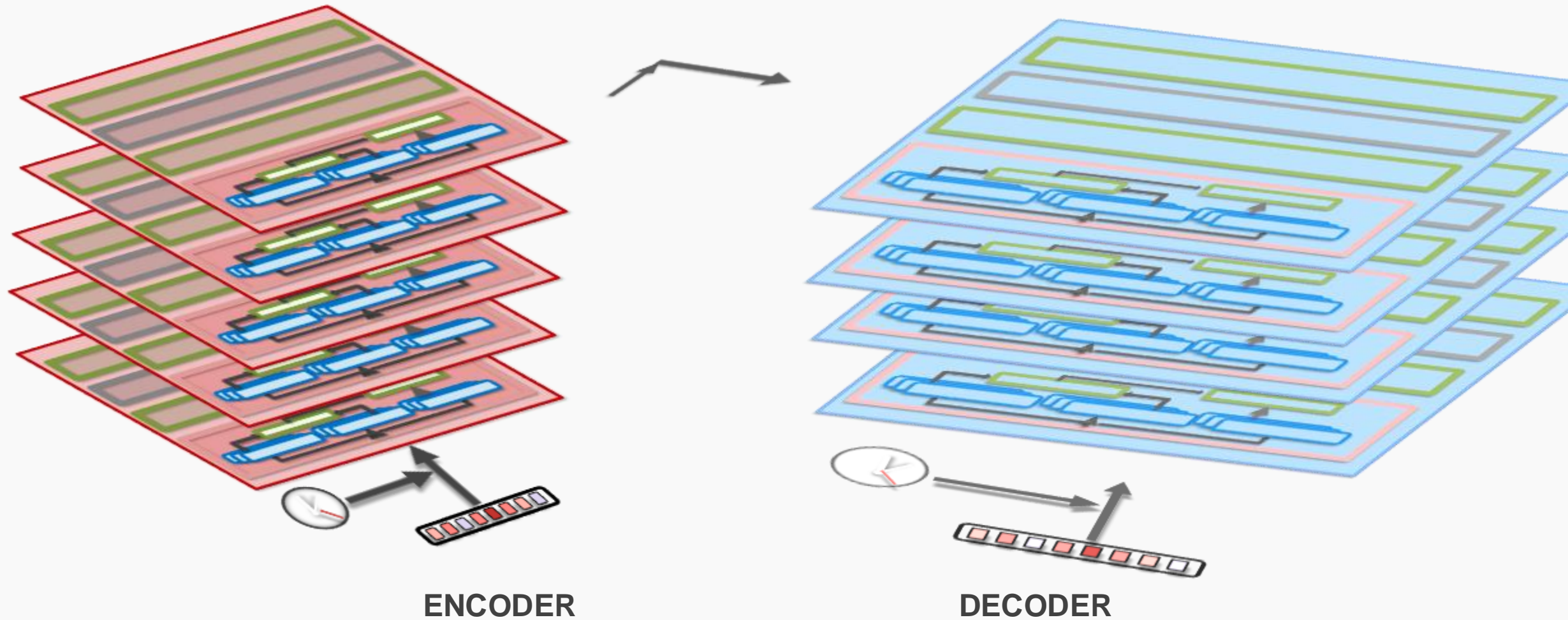
Transformers - Summary



More on this in the next session!

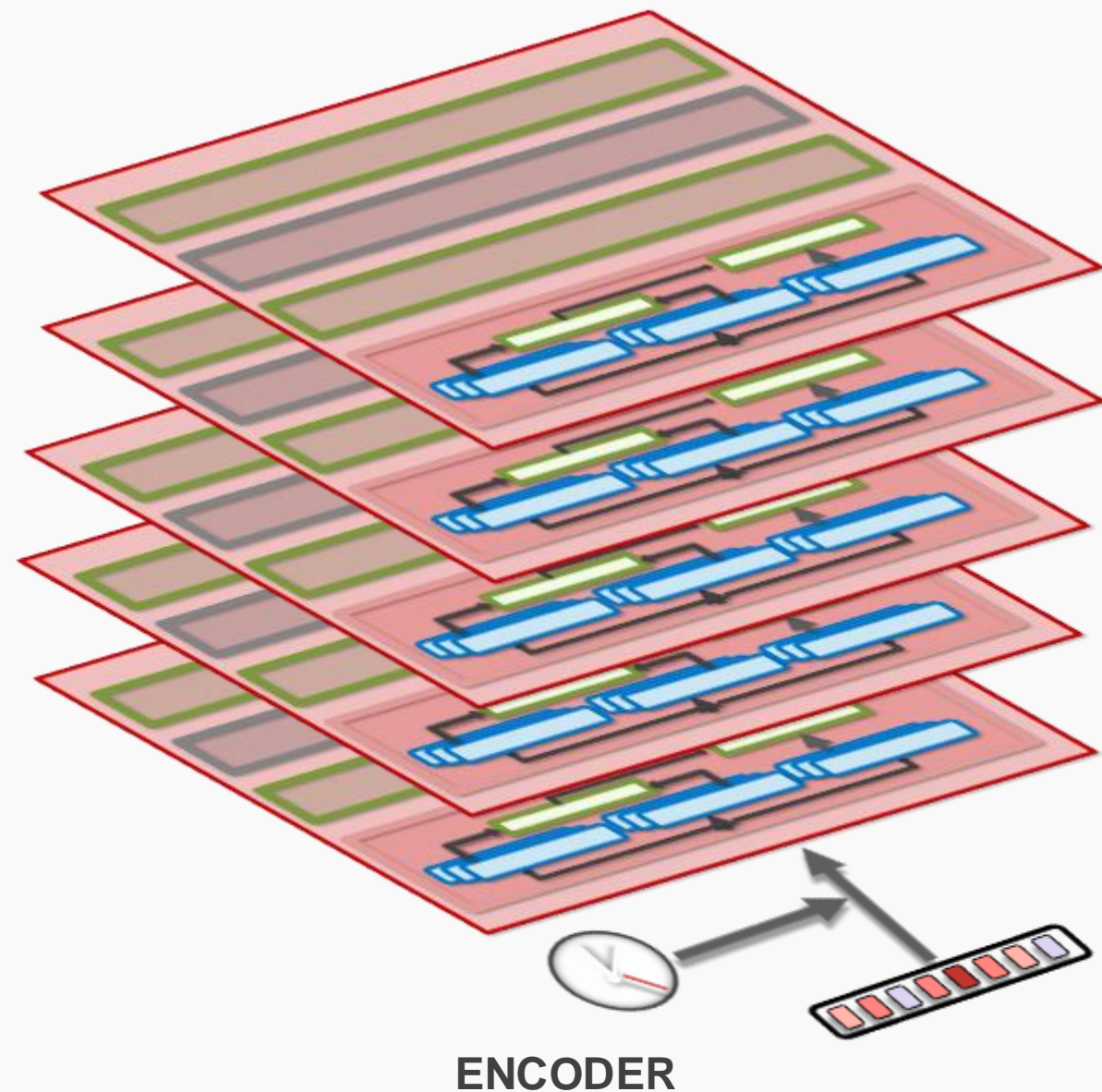


Transformers - Summary



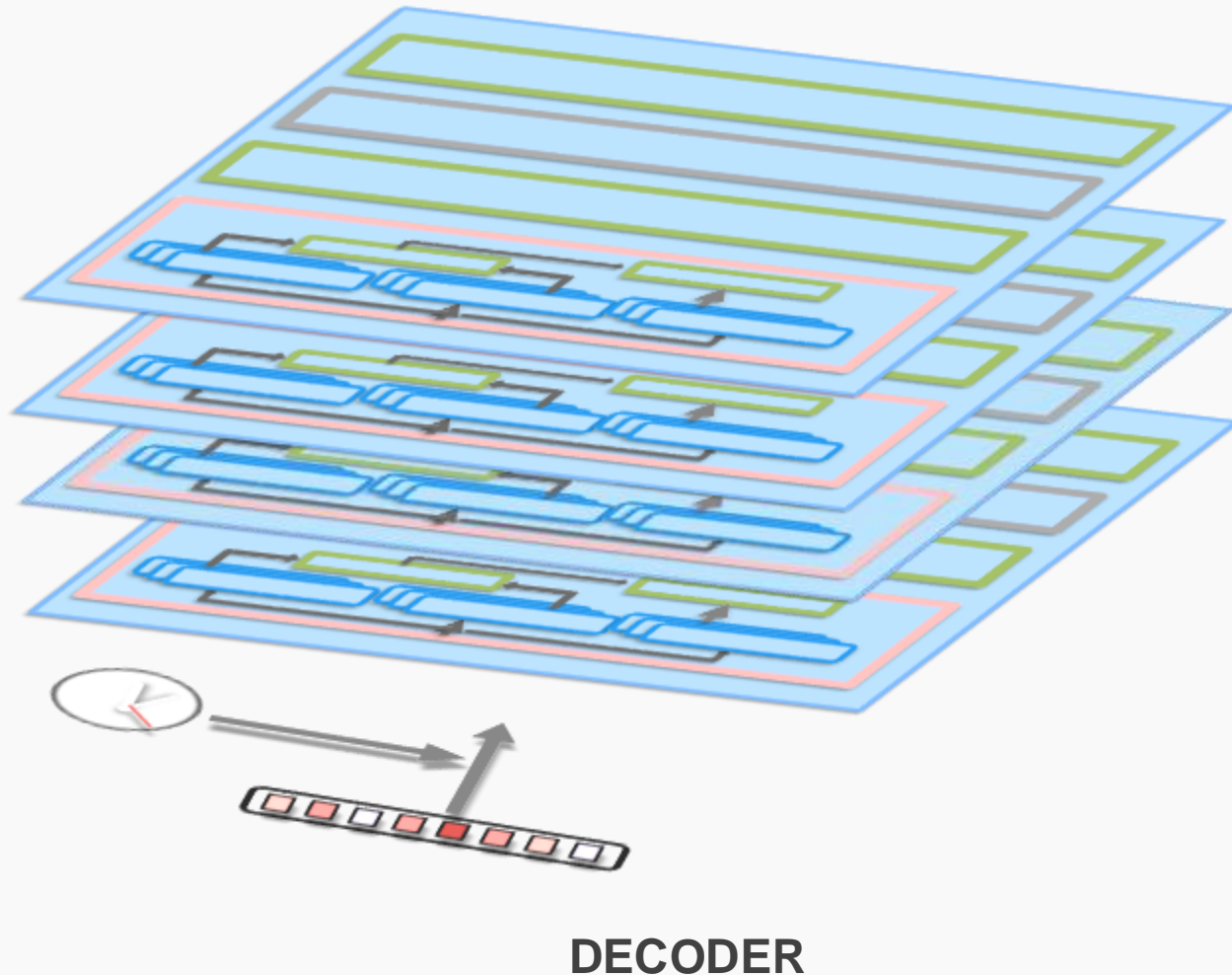
From Transformers to BERT

- Instead of an **Encoder-Decoder** architecture for machine translation, what if we just use the encoder for Language Model.
- This led to the new architecture called **Bidirectional Encoder Representations from Transformers**, or more commonly known as BERT.



From Transformers to GPT

- Now if we use just the decoder for Language Model
- We get a Causal Language Model (A word is predicted using words from its left context)
- Also known as autoregressive model
- Generative Pre-trained Transformer, or more commonly known as GPT



Thank you!