

Attention

- a visual tour

Martin Stoffel

Transformer

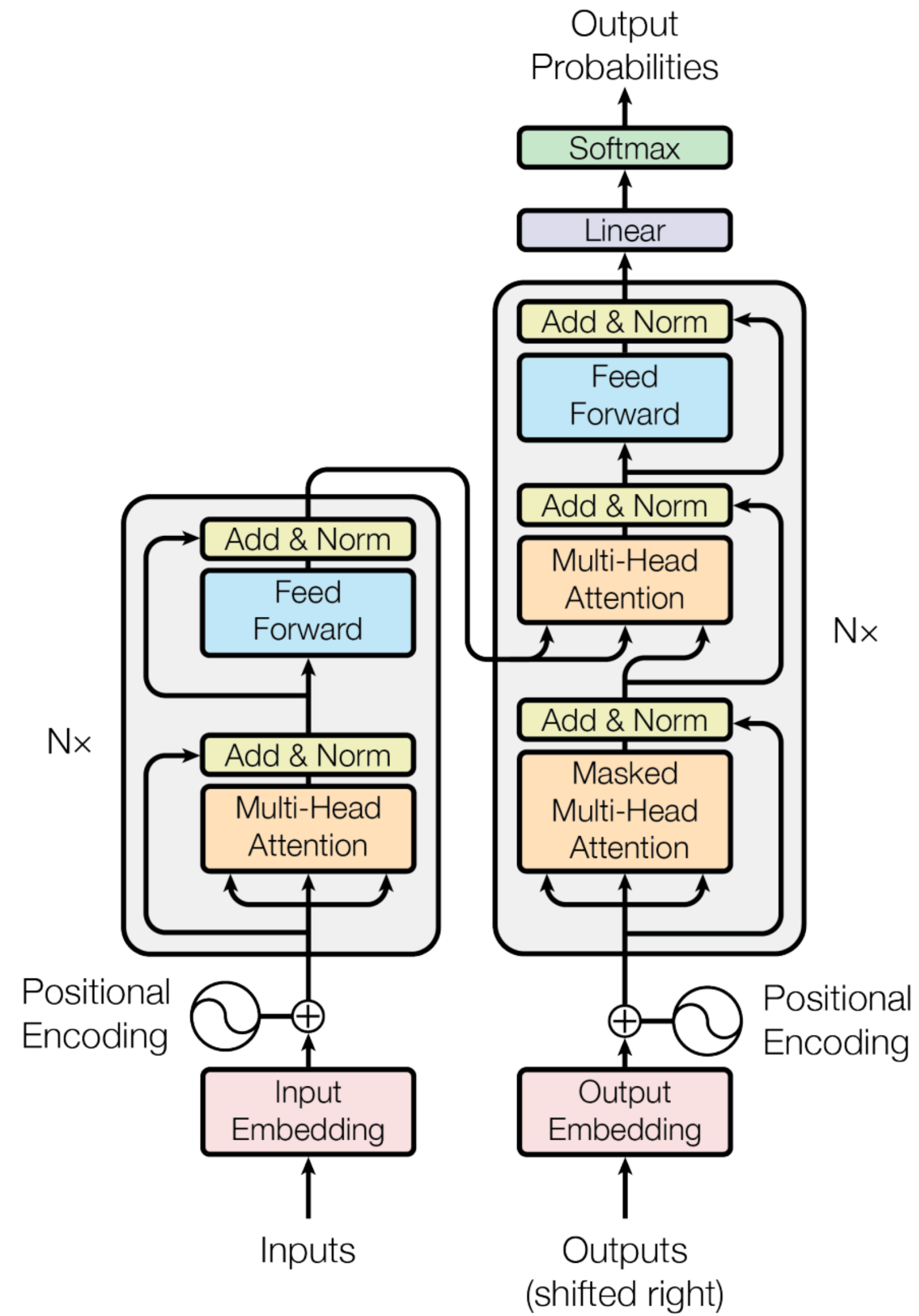
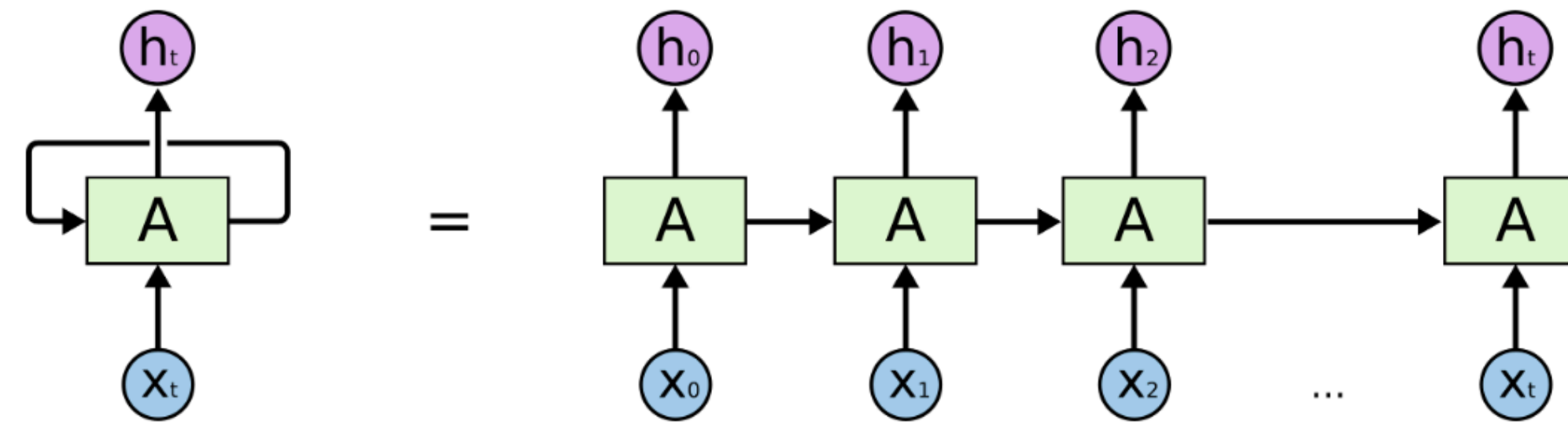


Figure 1: The Transformer - model architecture.

Why attention?



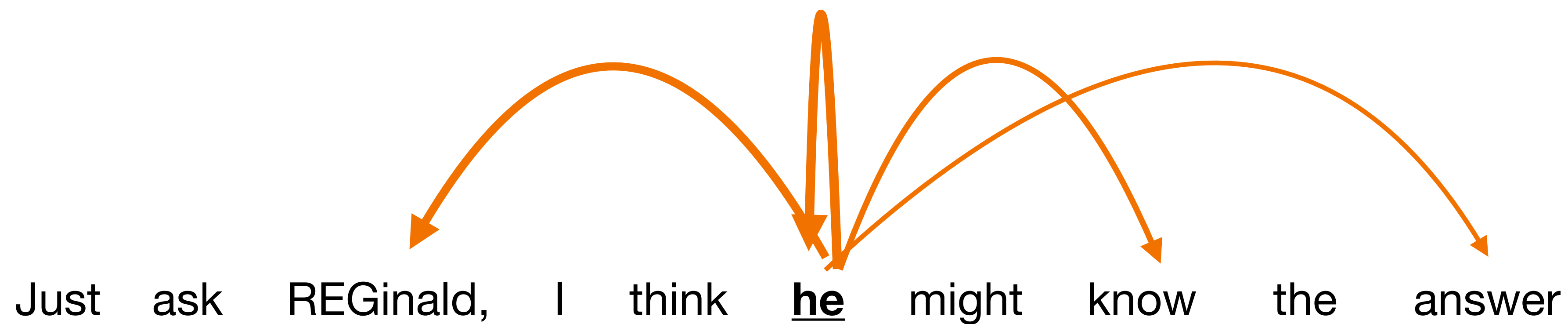
An unrolled recurrent neural network.

- Recurrence in RNNs and LSTMs has two fundamental weaknesses:
 - **Capturing long-range dependencies:** exploding/vanishing gradients
 - **Sequential computation:** makes parallelisation hard
- Also: attention seems to be a bit more *interpretable*

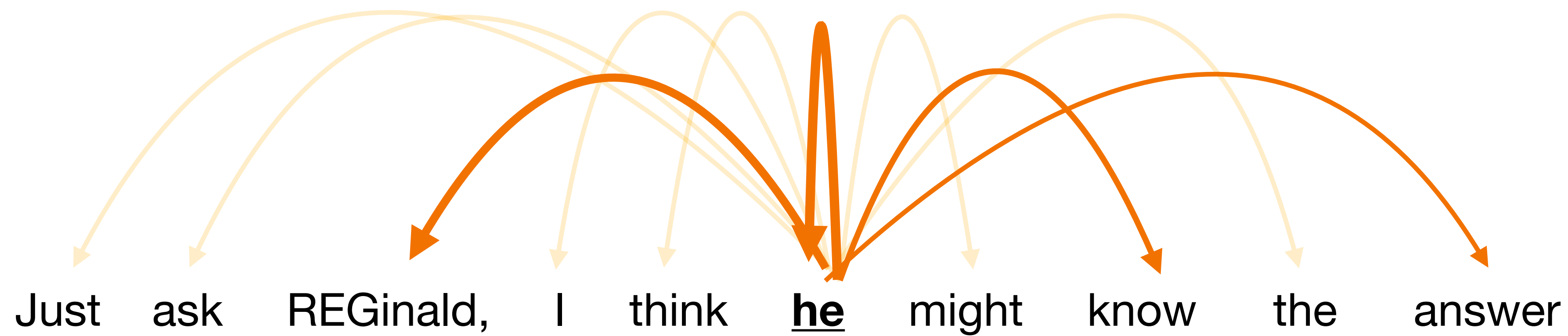
Attention: Goal

Just ask REGinald, I think he might know the answer

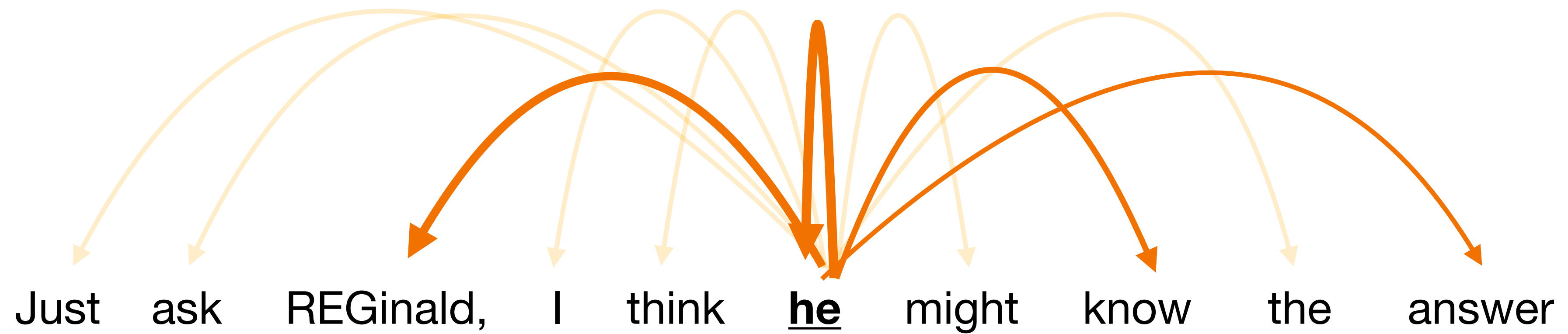
Attention: Goal



Attention: Goal

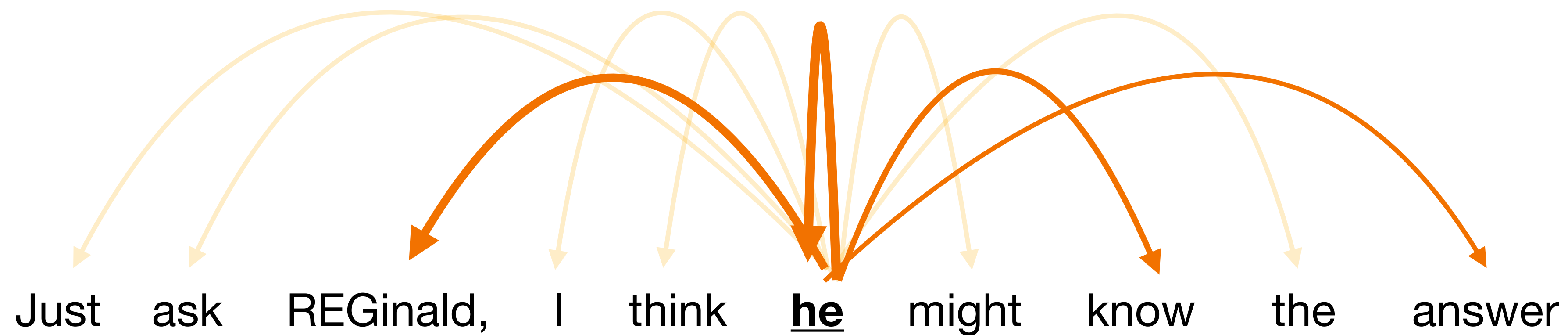


Attention: Goal



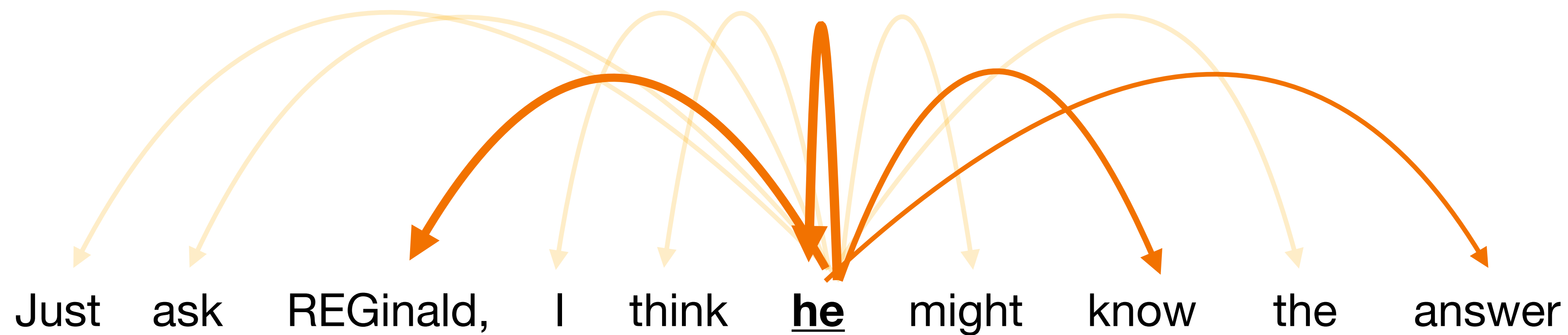
Me: What's the goal of the attention mechanism in transformers, in one illuminating sentence?

Attention: Goal



Me: What's the goal of the attention mechanism in transformers, in one illuminating sentence? Imagine you're not an AI, but a technologically advanced alien species, which doesn't use language to communicate.

Attention: Goal

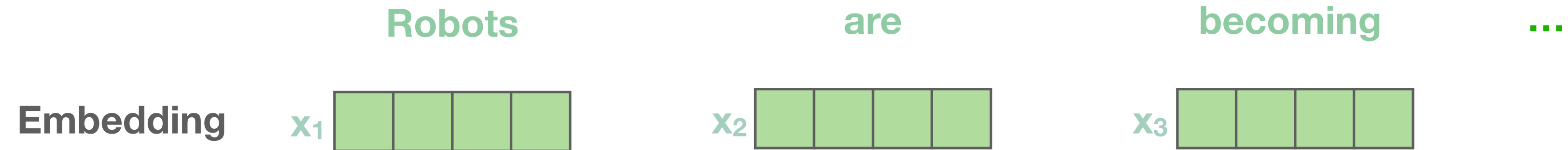


Me: What's the goal of the attention mechanism in transformers, in one illuminating sentence? Imagine you're not an AI, but a technologically advanced alien species, which doesn't use language to communicate.

ChatGP The attention mechanism in Transformers is akin to our advanced pattern recognition systems, **scanning and selectively amplifying** pertinent signals from a sea of information, regardless of their temporal sequence, to understand the **intricate dance of human symbolic sequences**, analogous to how we decipher cosmic vibrations.

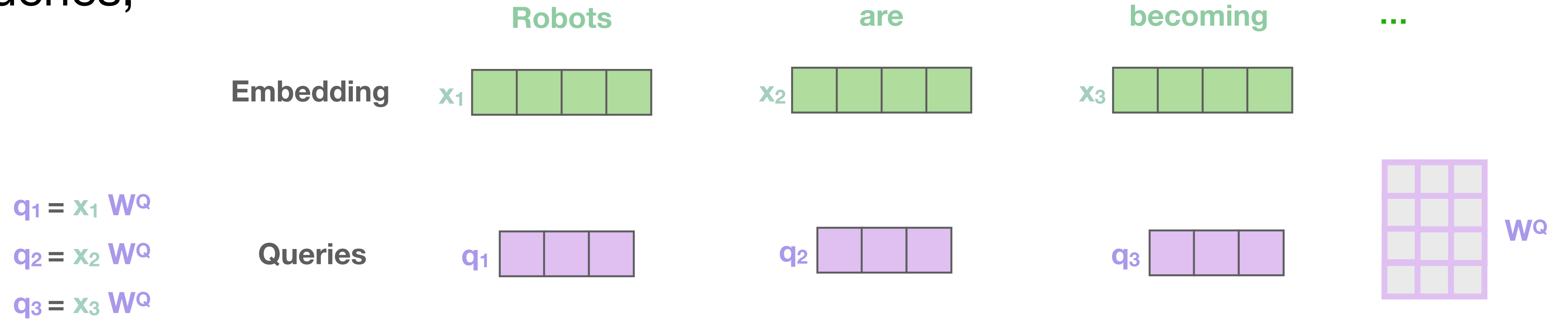
Attention 1)

Keys, Queries,
Values



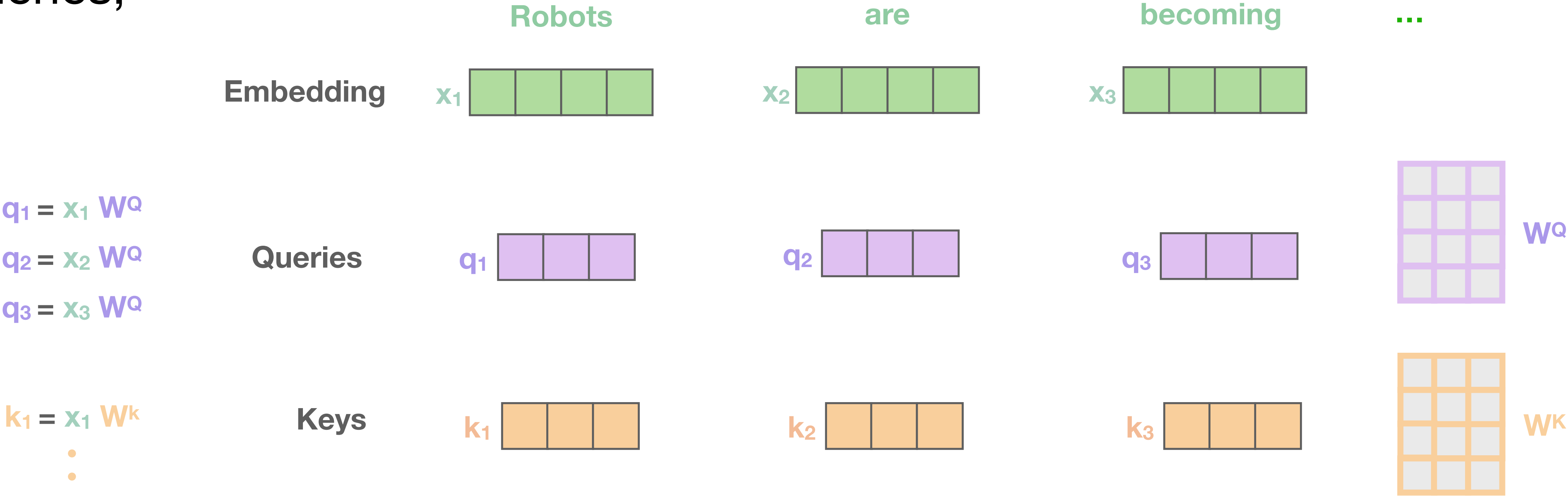
Attention 1)

Keys, Queries,
Values



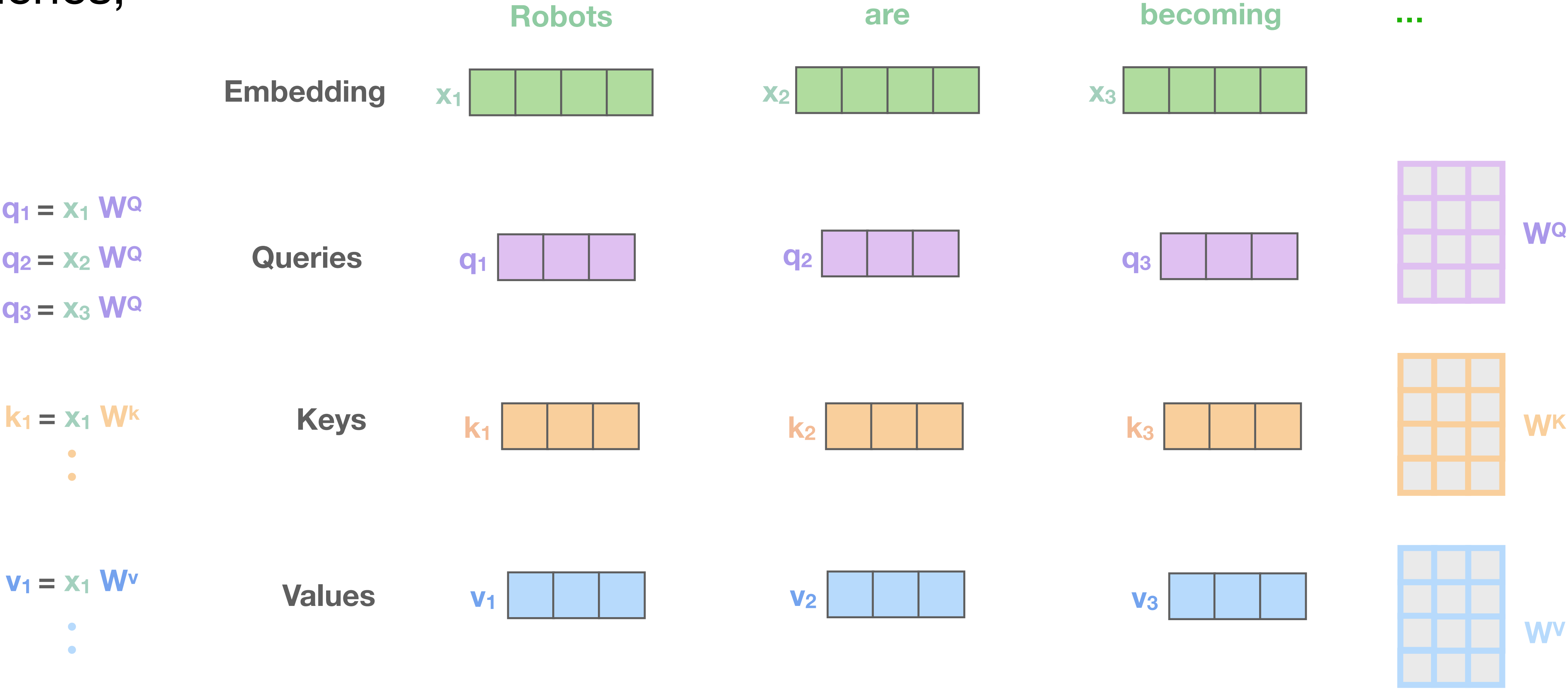
Attention 1)

Keys, Queries,
Values



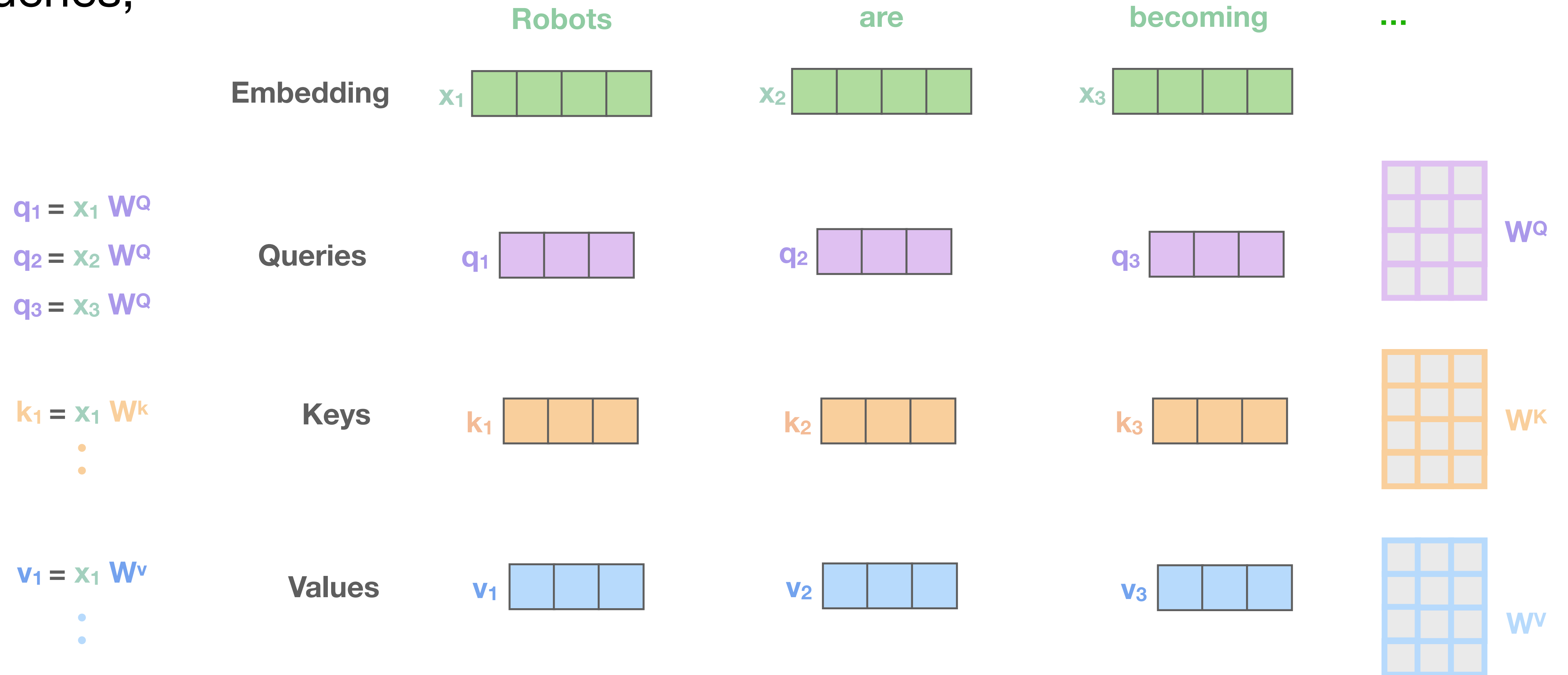
Attention 1)

Keys, Queries,
Values



Attention 1)

Keys, Queries,
Values



-> Linear projection of each input vector x into a query, key and value vector

Attention 2)

Calculations



Based on
<https://jalammar.github.io/illustrated-transformer/>

Attention 2)

Calculations



Attention 2)

Calculations

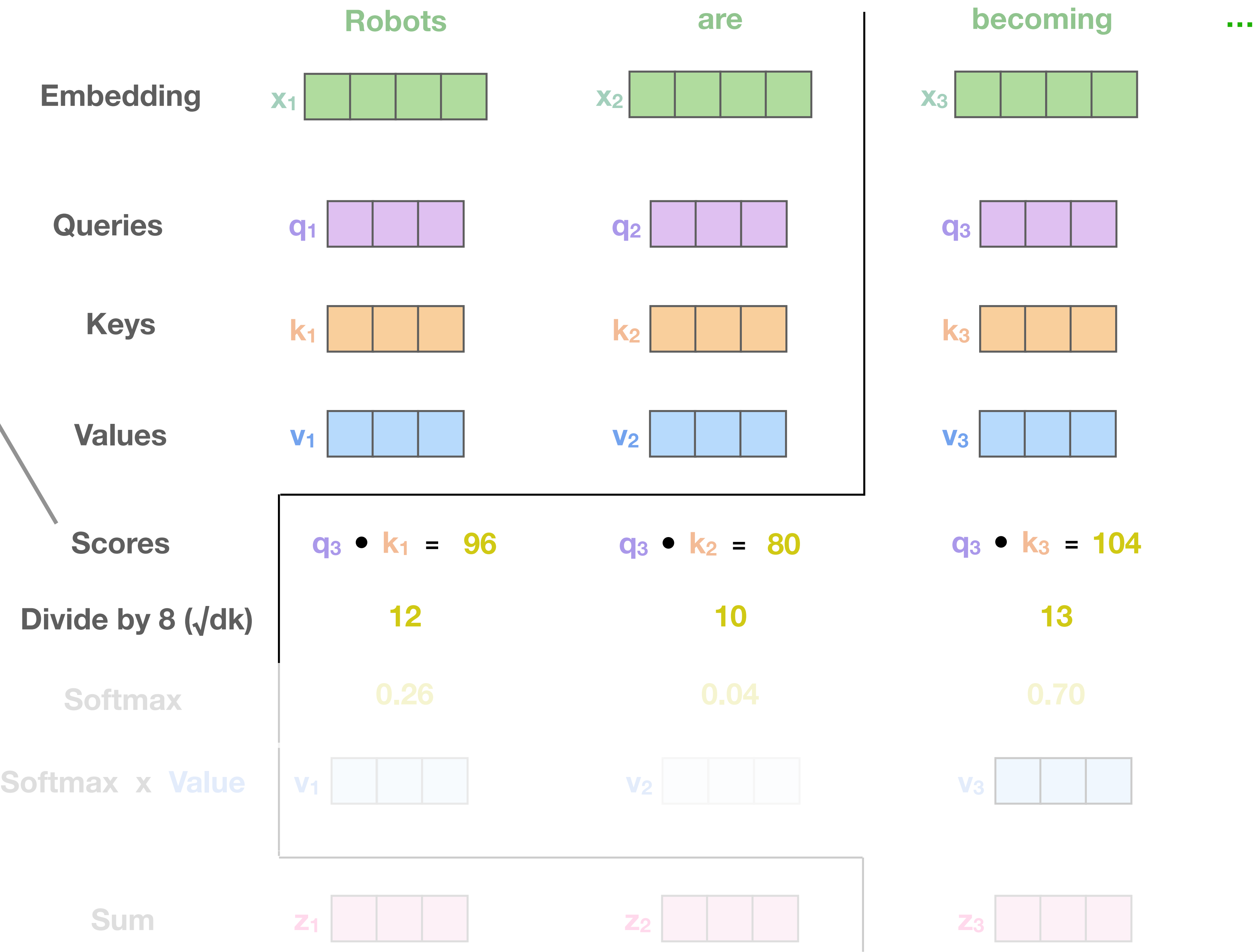
Scores:
How well does the current **query** align with all the **keys** in the sequence?
(*raw attention scores*)



Attention 2)

Calculations

Scores:
How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

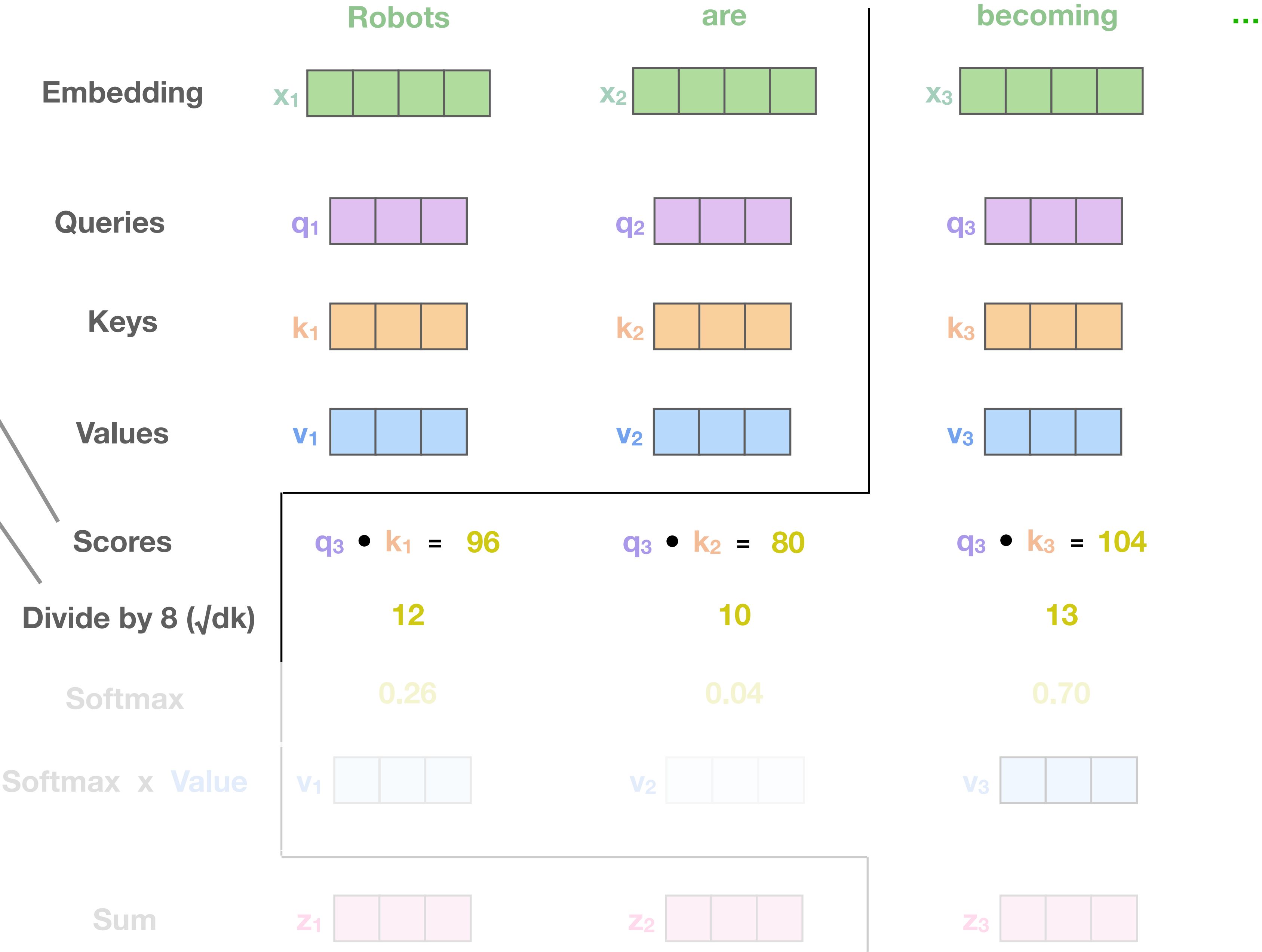


Attention 2)

Calculations

Scores:
How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:
Prevent dot-products from becoming too large and gradients too small.



Attention 2)

Calculations

Scores:
How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:
Prevent dot-products from becoming too large and gradients too small.



Attention 2)

Calculations

Scores:

How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:

Prevent dot-products from becoming too large and gradients too small.

Softmax:

Probability distribution (summing to 1)
How much *attention* to pay to each token?
(scaled attention scores)



Attention 2)

Calculations

Scores:

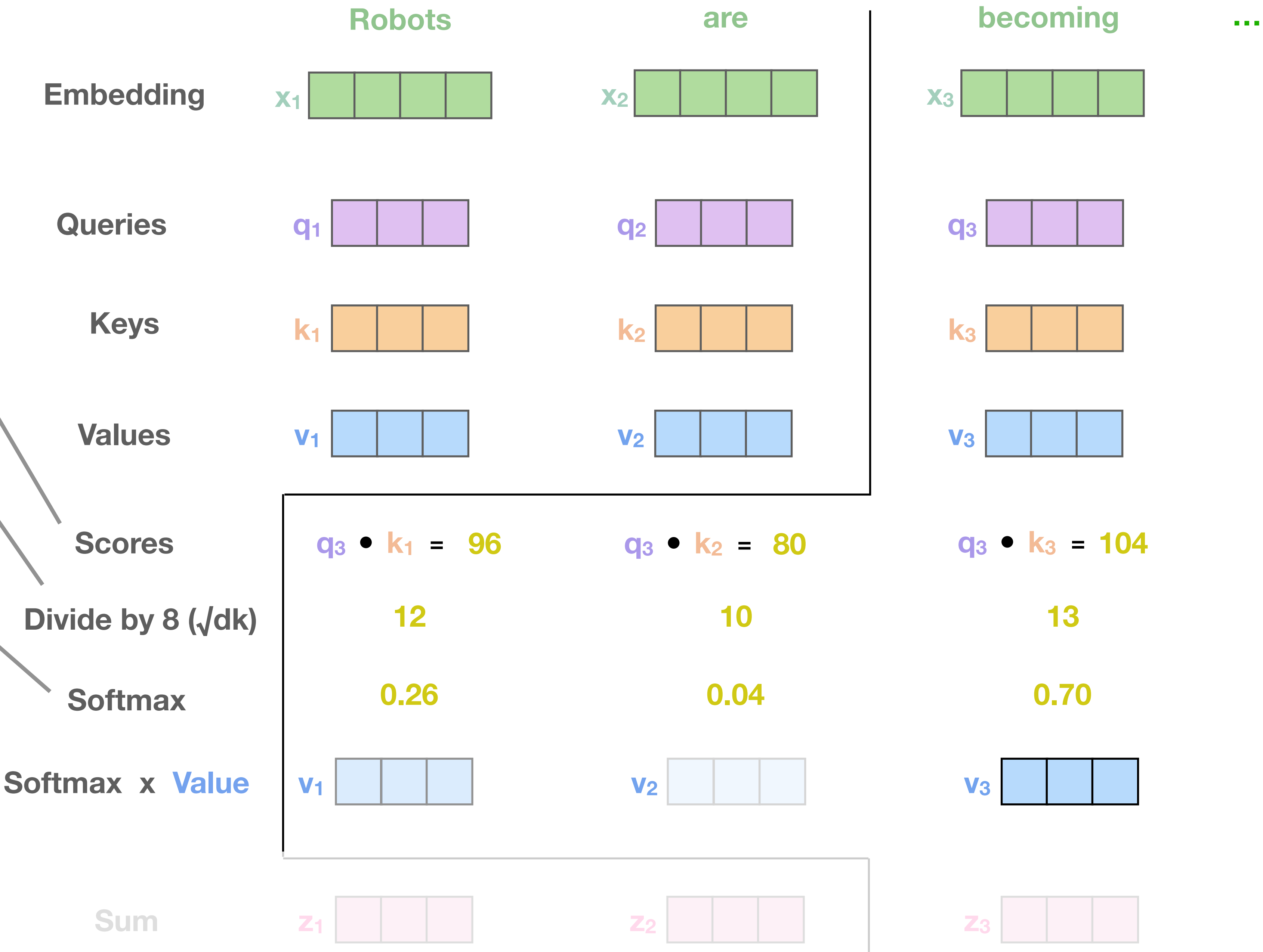
How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:

Prevent dot-products from becoming too large and gradients too small.

Softmax:

Probability distribution (summing to 1)
How much *attention* to pay to each token?
(scaled attention scores)



Attention 2)

Calculations

Scores:

How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:

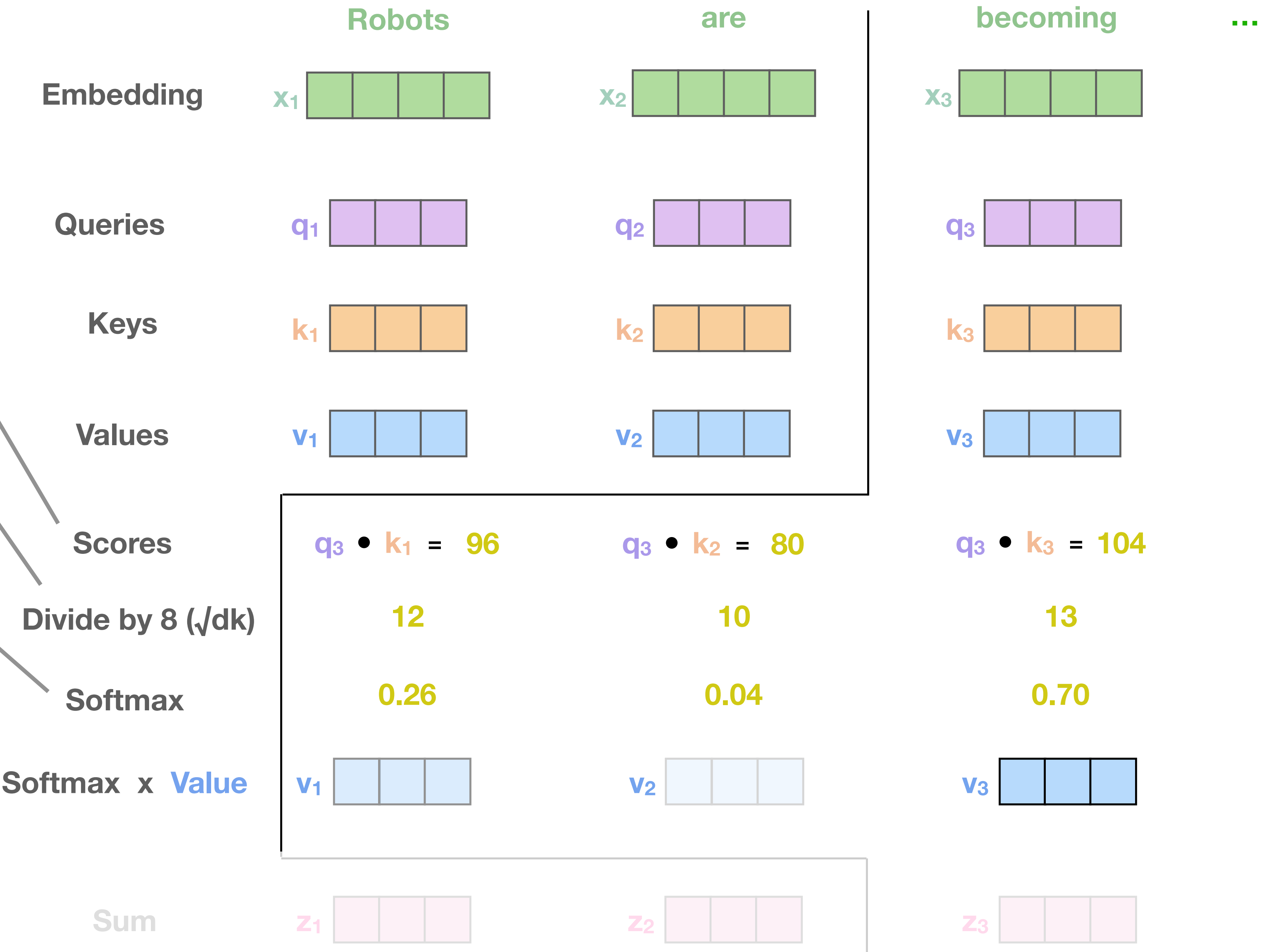
Prevent dot-products from becoming too large and gradients too small.

Softmax:

Probability distribution (summing to 1)
How much *attention* to pay to each token?
(scaled attention scores)

Attention-weighted values

Focus on relevant tokens, drown out irrelevant tokens.



Attention 2)

Calculations

Scores:

How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:

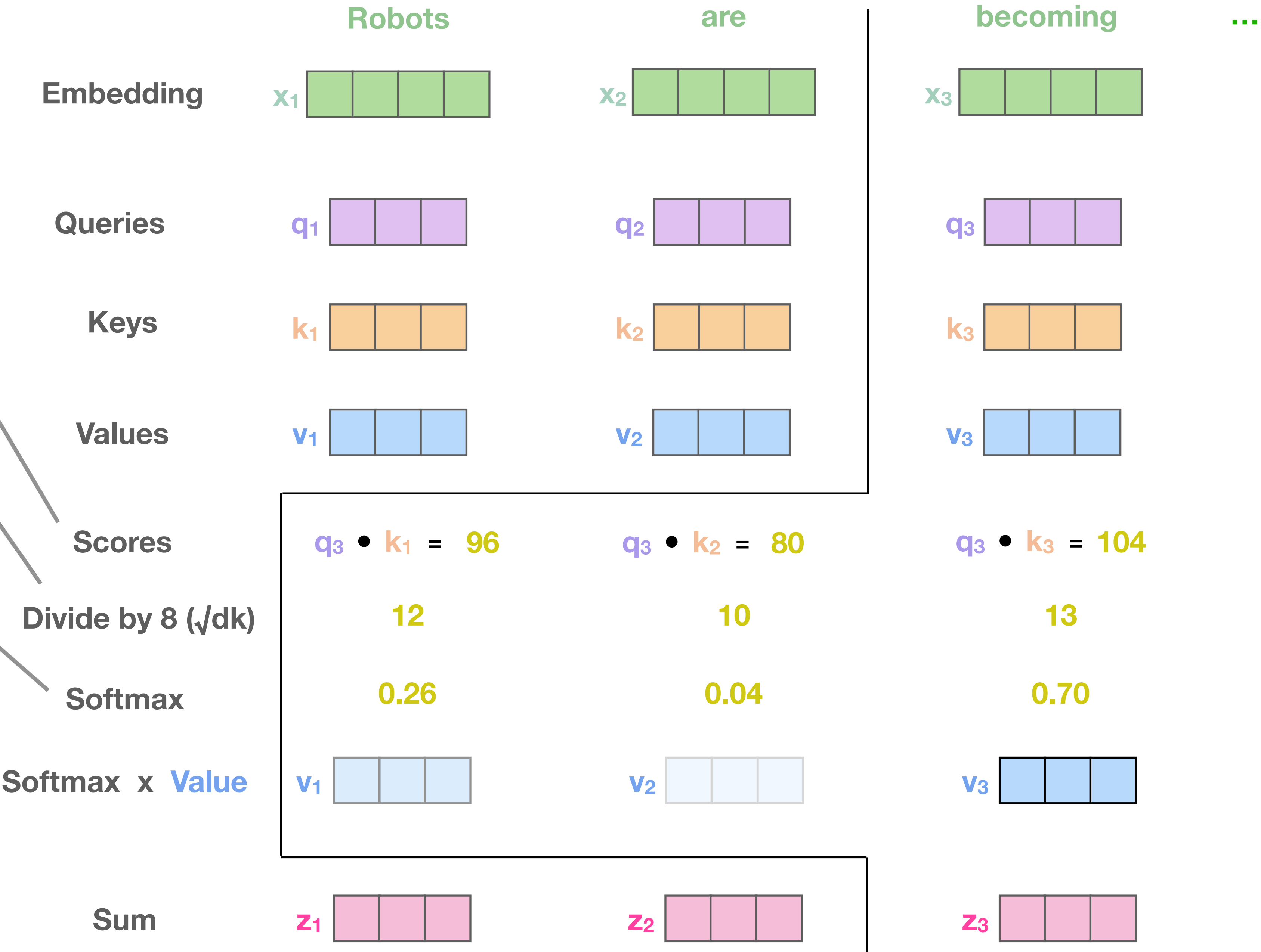
Prevent dot-products from becoming too large and gradients too small.

Softmax:

Probability distribution (summing to 1)
How much *attention* to pay to each token?
(scaled attention scores)

Attention-weighted values

Focus on relevant tokens, drown out irrelevant tokens.



Attention 2)

Calculations

Scores:

How well does the current **query** align with all the **keys** in the sequence?
(raw attention scores)

Scaling factor:

Prevent dot-products from becoming too large and gradients too small.

Softmax:

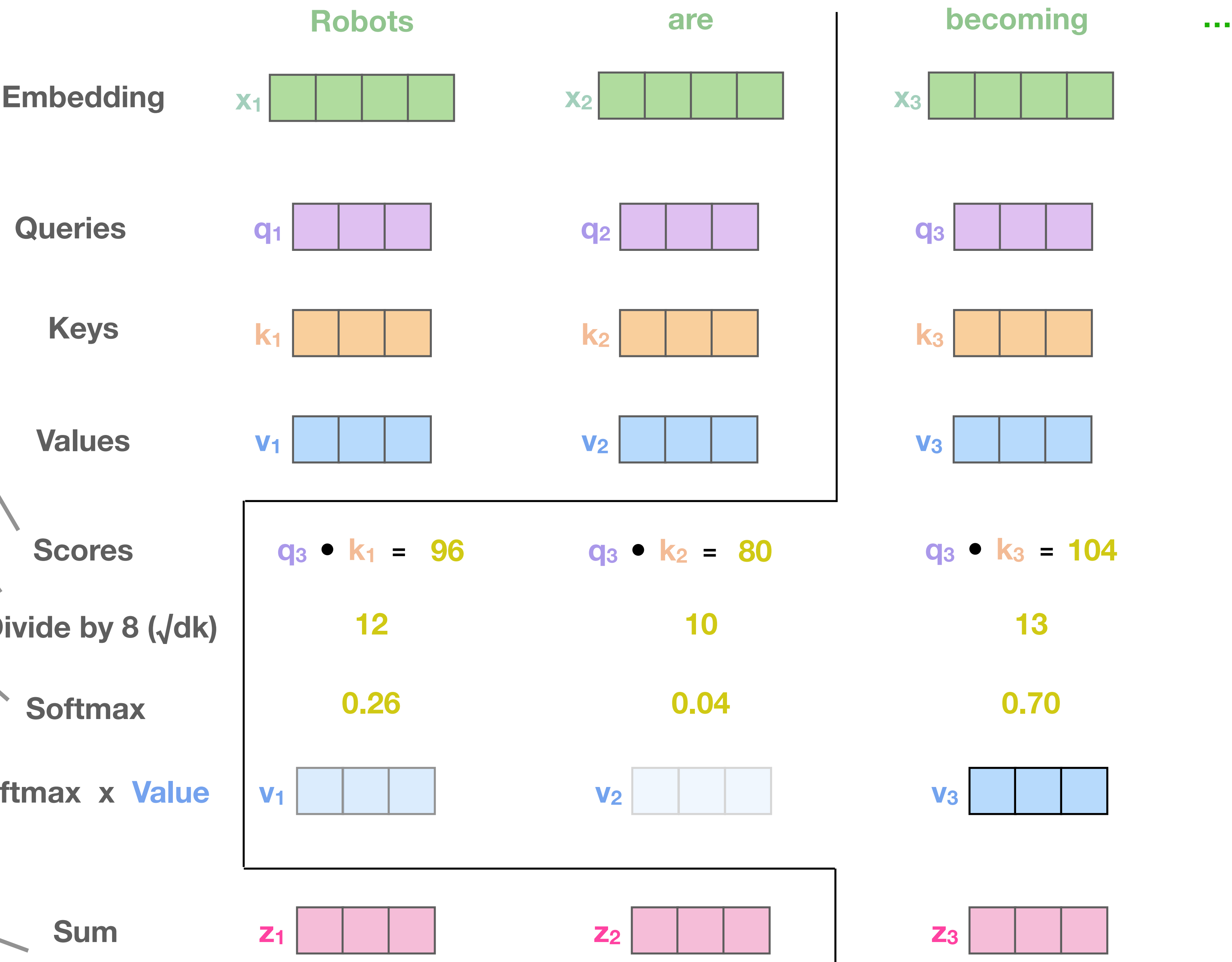
Probability distribution (summing to 1)
How much *attention* to pay to each token?
(scaled attention scores)

Attention-weighted values

Focus on relevant tokens, drown out irrelevant tokens.

Sum:

Aggregate attention-weighted **value** vectors into an output vector **z**

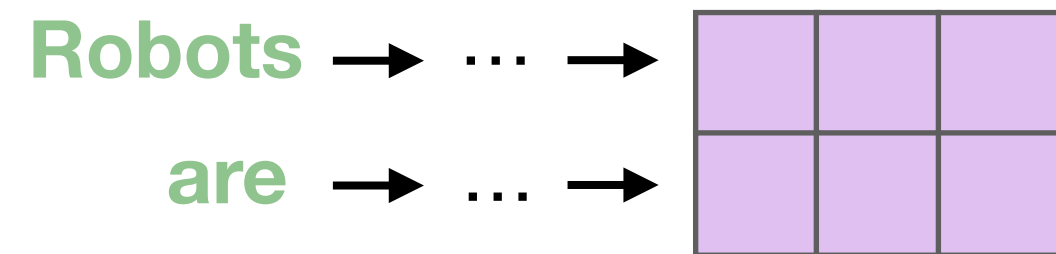


Attention 3)

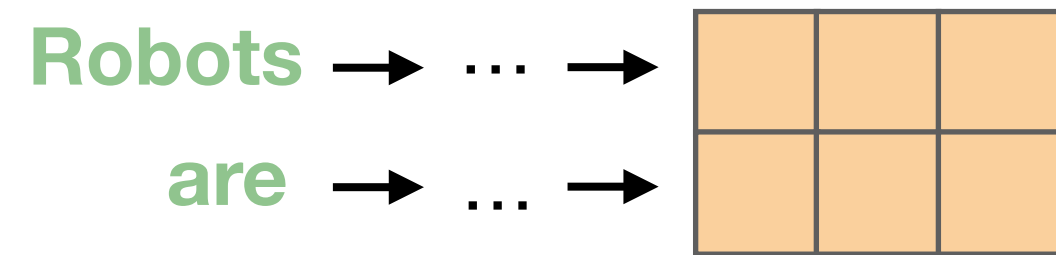
Matrix form

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

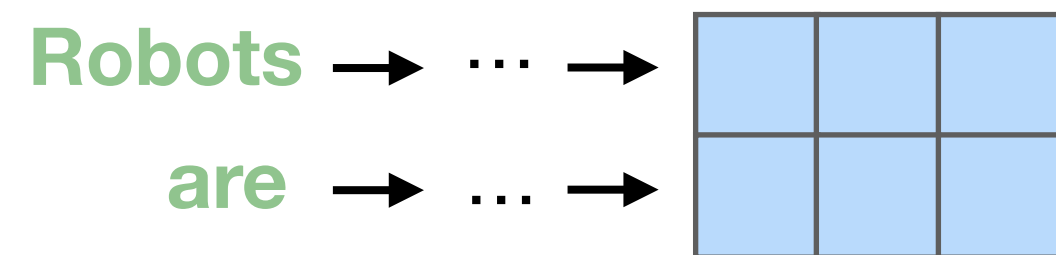
Q



K

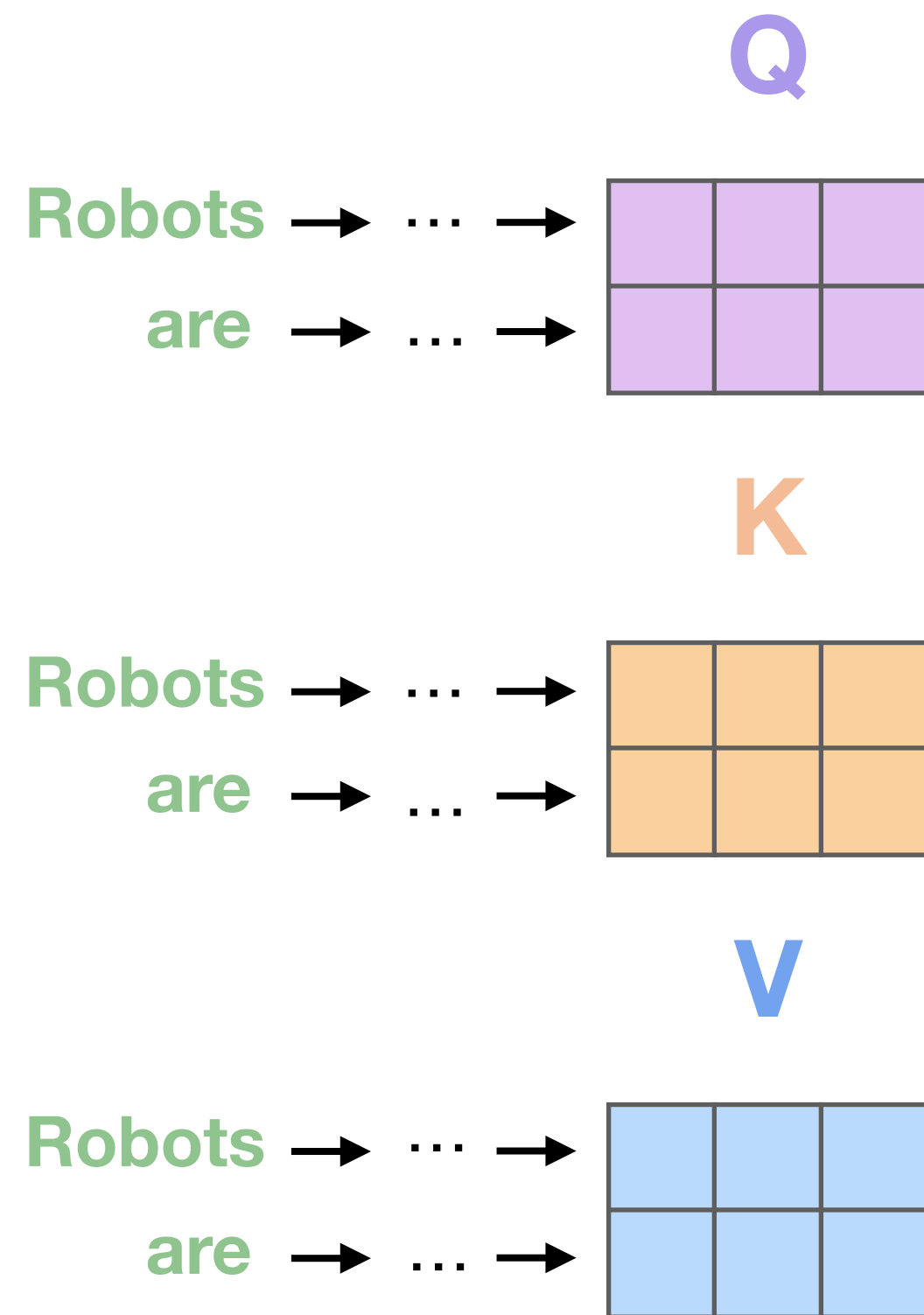


V



Attention 3)

Matrix form



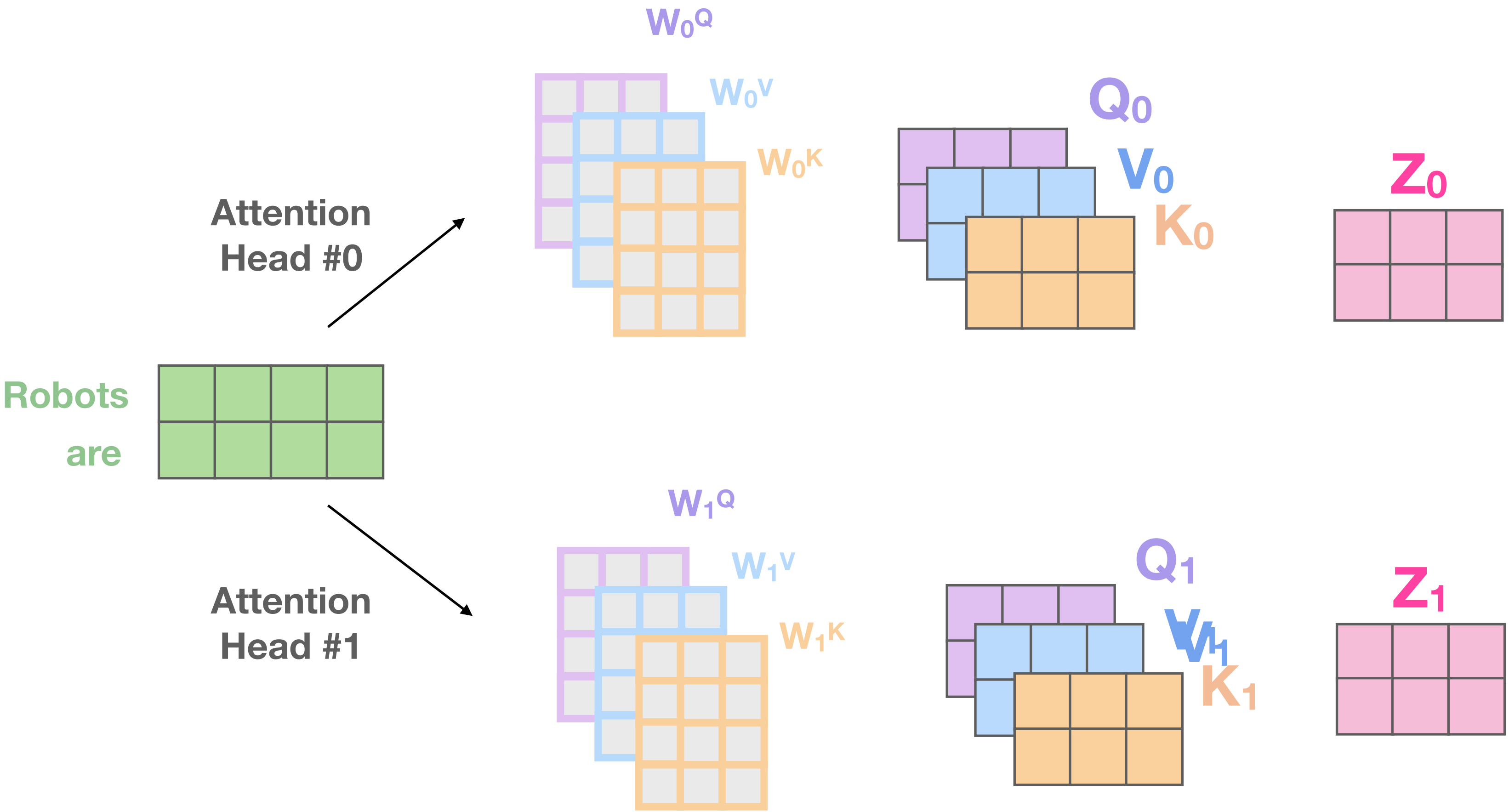
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Diagram illustrating the Attention mechanism:

- The Query matrix **Q** (purple) and the transposed Key matrix **K^T** (orange) are combined and divided by $\sqrt{d_k}$ to compute the attention weights.
- The result of the softmax operation is then multiplied by the Value matrix **V** (blue) to produce the final output matrix **Z** (pink).

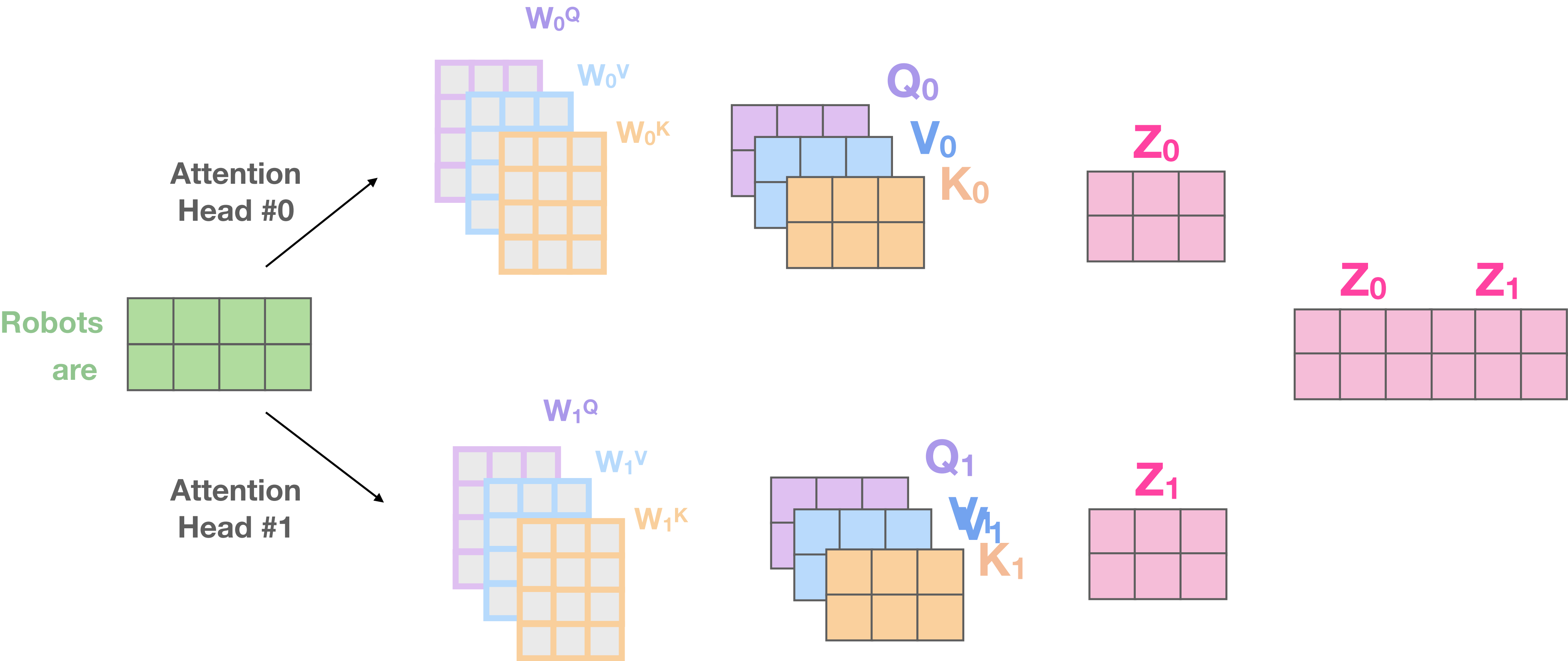
Attention 4)

Multiple heads



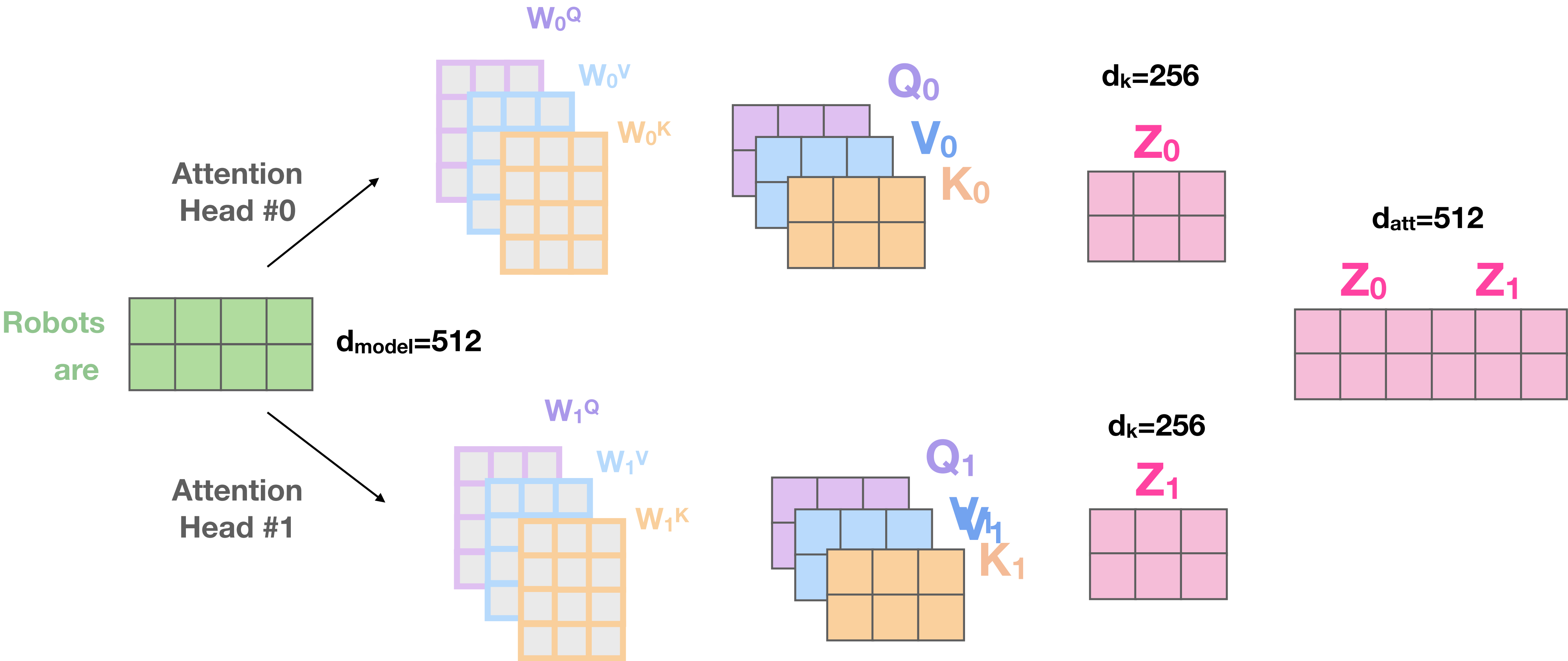
Attention 4)

Multiple heads



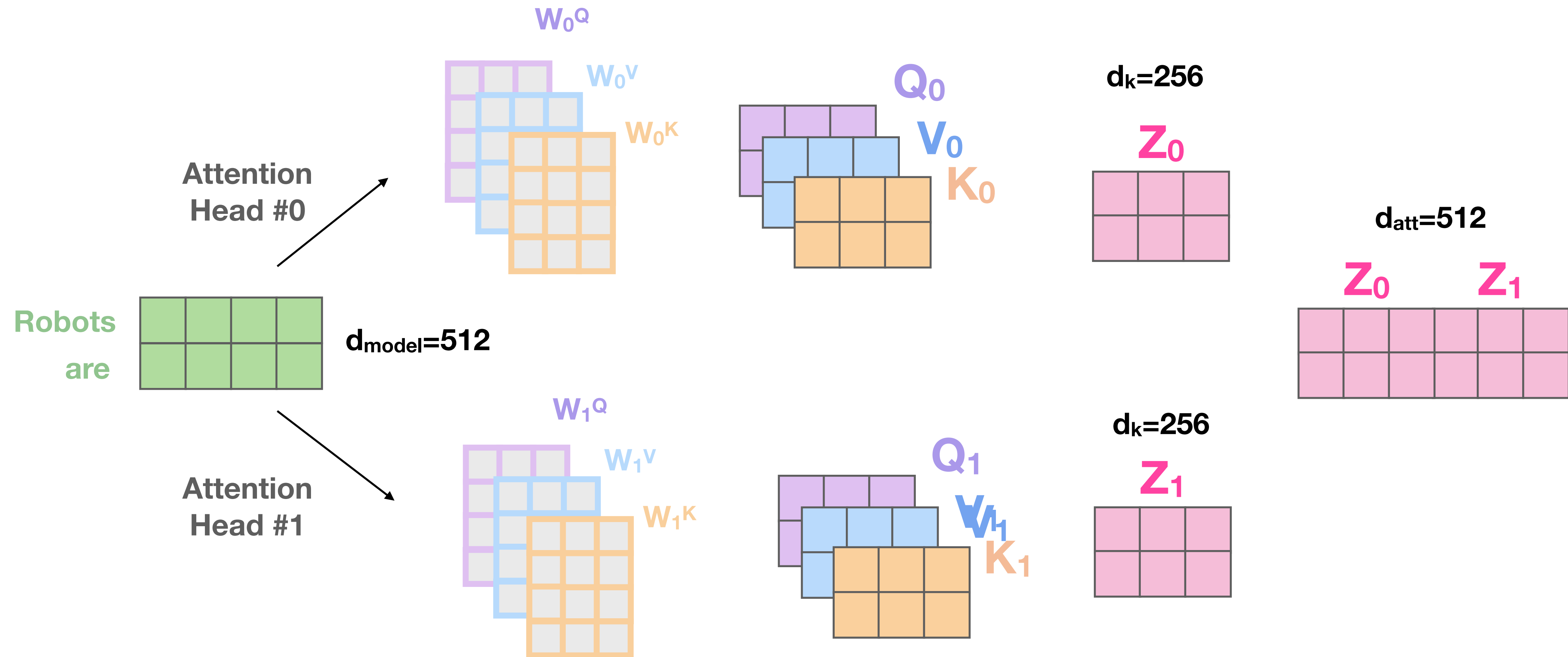
Attention 4)

Multiple heads



Attention 4)

Multiple heads



-> Allows the model to focus on different positions at each head

Three variants of Attention

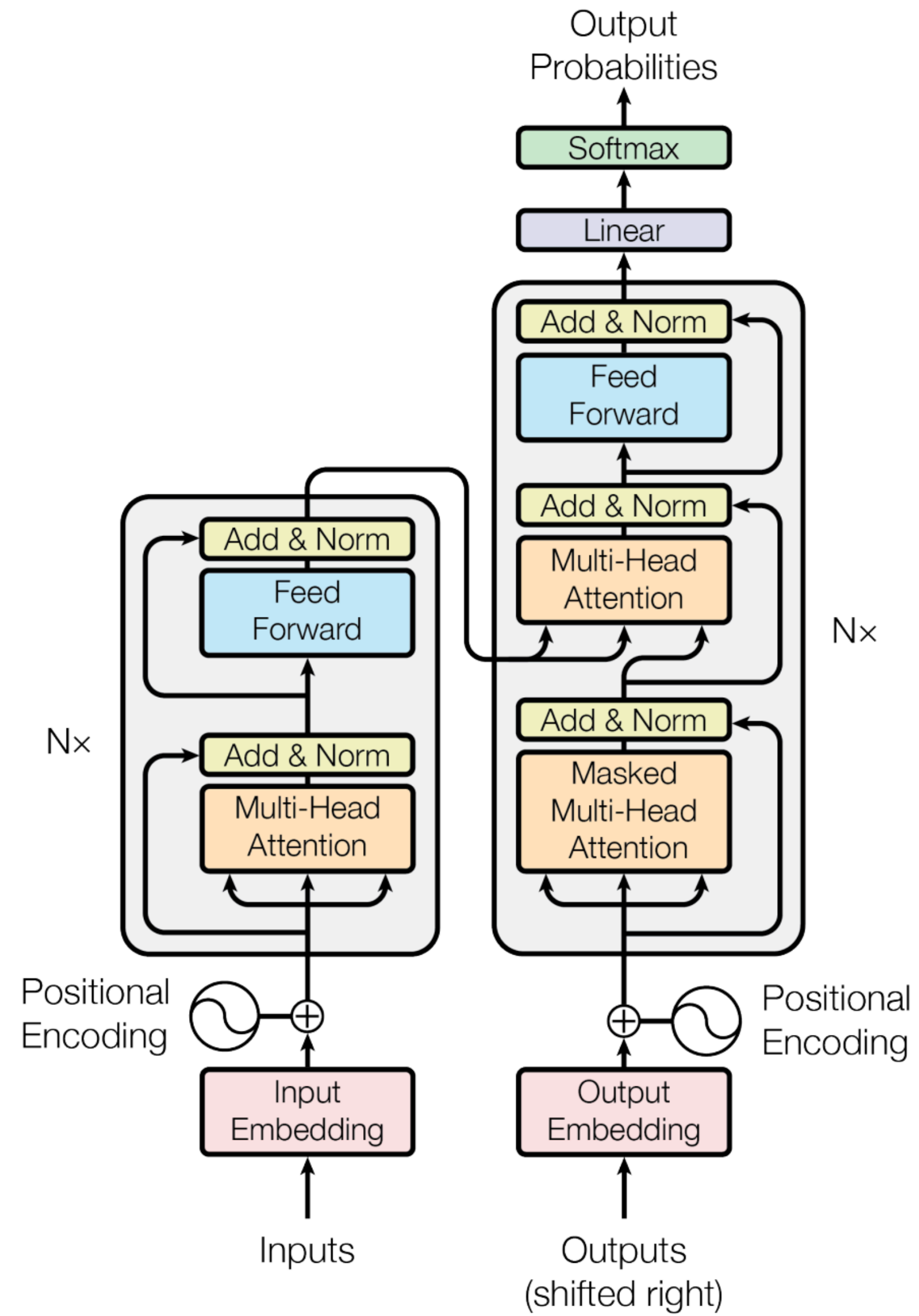


Figure 1: The Transformer - model architecture.

Three variants of Attention

Self-attention

- keys, queries, values come from the *same source* x
- each token attends to *all tokens* in the input sequence

-> e.g. sentiment analysis

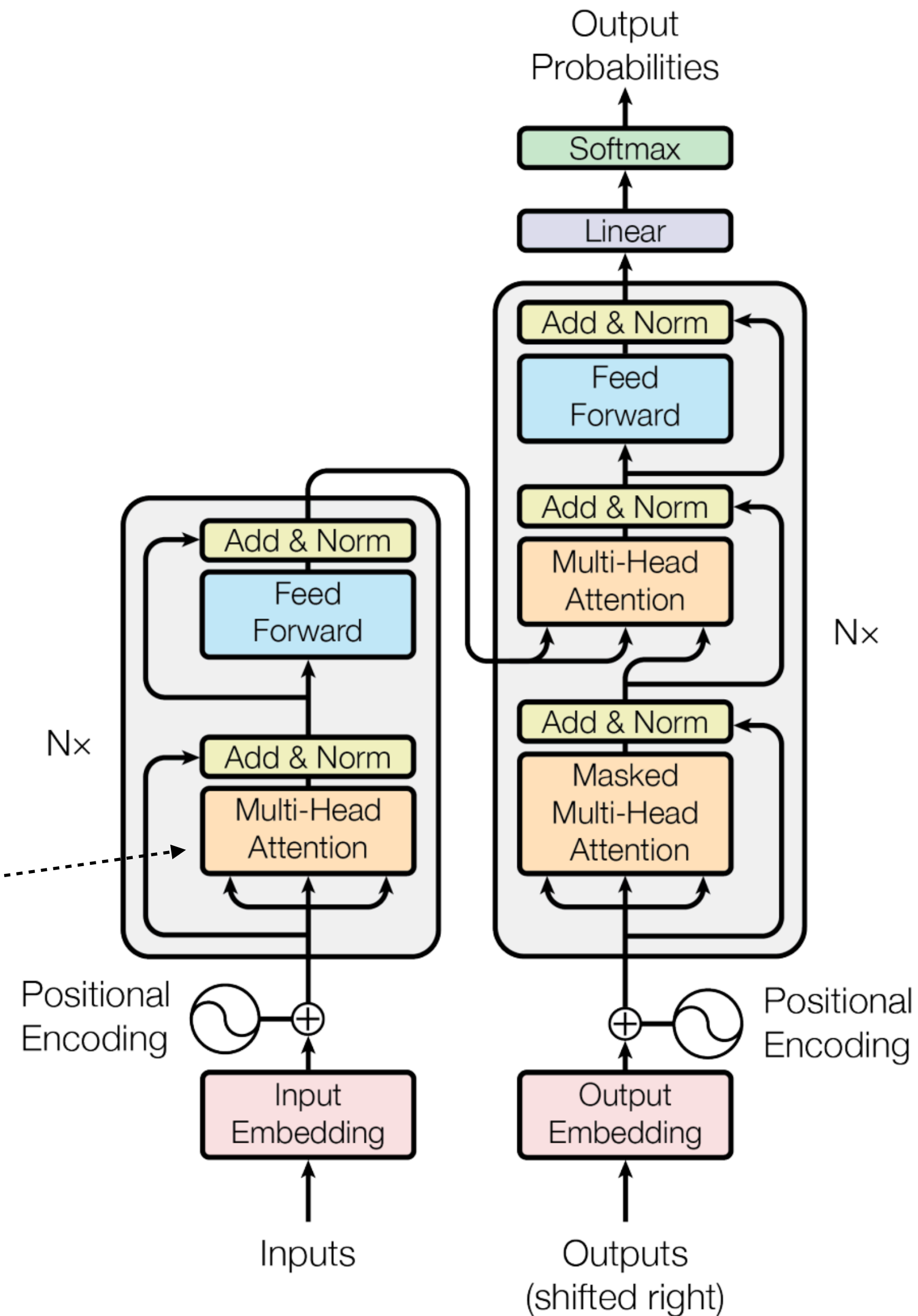


Figure 1: The Transformer - model architecture.

Three variants of Attention

Self-attention

- keys, queries, values come from the *same source* x
- each token attends to *all tokens* in the input sequence
- > e.g. sentiment analysis

Masked Self-attention

- keys, queries, values come from the *same source* x
- each token attends to all *previous tokens* in the input sequence
- > e.g. text generation

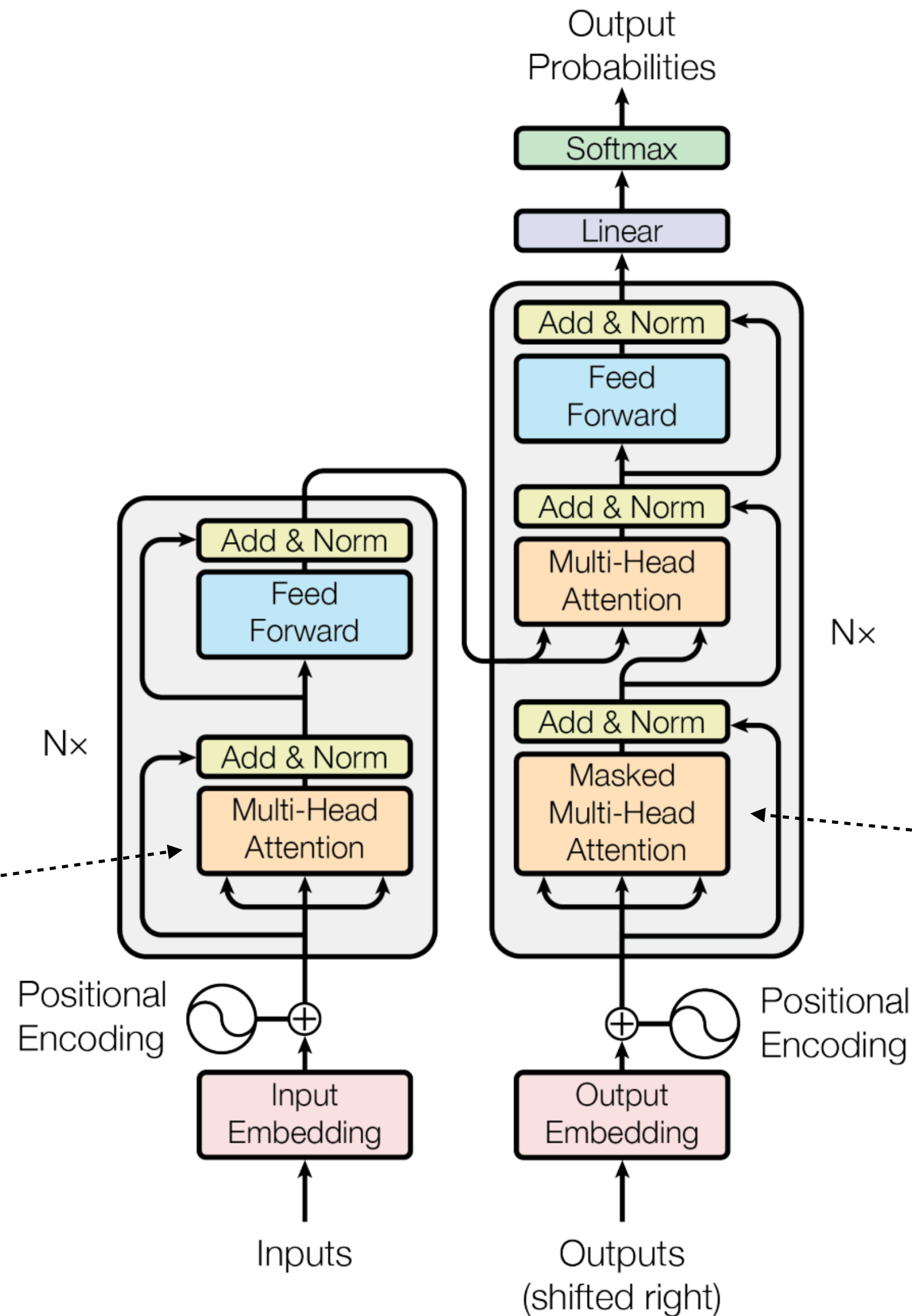


Figure 1: The Transformer - model architecture.

Three variants of Attention

Self-attention

- keys, queries, values come from the *same source* x
- each token attends to *all tokens* in the input sequence
- > e.g. sentiment analysis

Cross-attention

- query from decoder, keys/values from encoder
- Each token attends to *all tokens* in sequence

Masked Self-attention

- keys, queries, values come from the *same source* x
- each token attends to all *previous tokens* in the input sequence
- > e.g. text generation

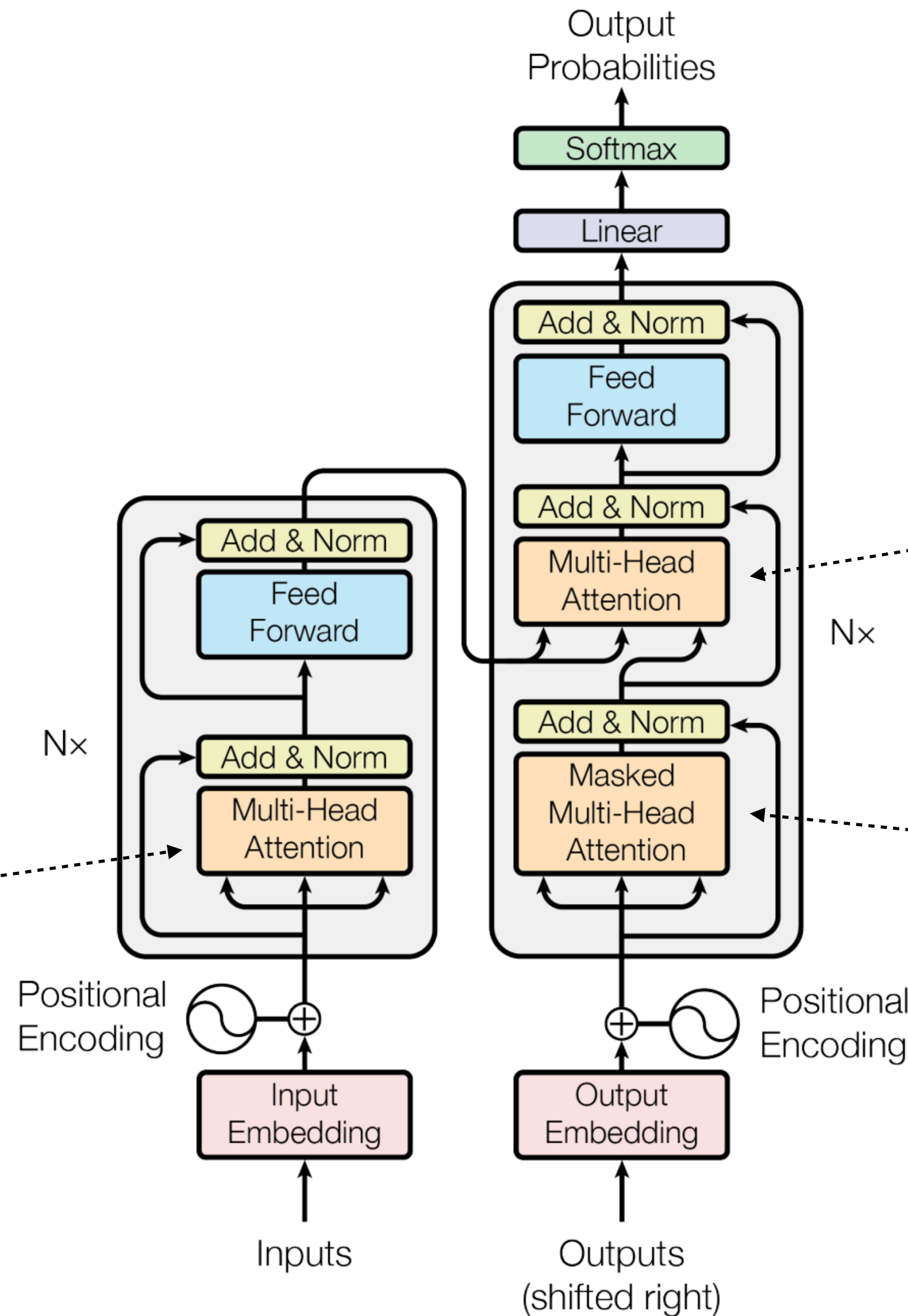


Figure 1: The Transformer - model architecture.

Summary

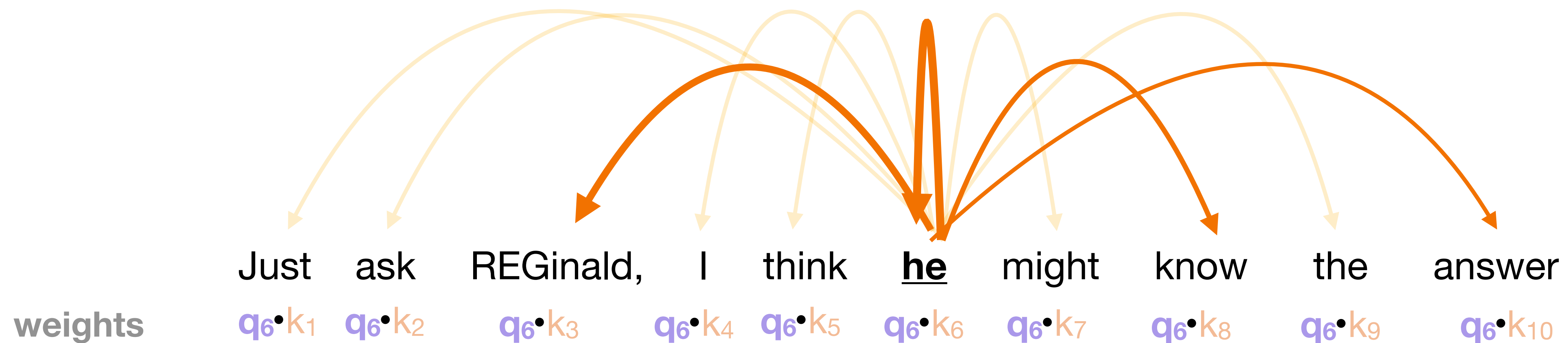
- Attention captures dependencies between words in a sequence, regardless of their distance

Summary

- Attention captures dependencies between words in a sequence, regardless of their distance
- **Query** **Key** dot-products produce attention scores / weights

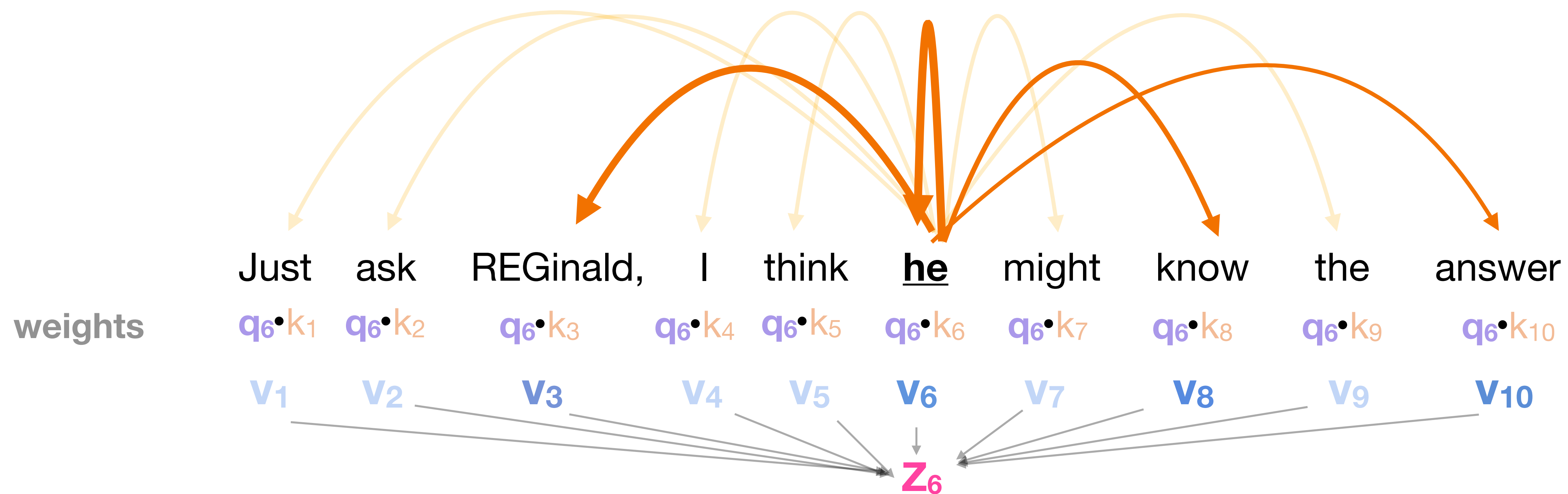
Summary

- Attention captures dependencies between words in a sequence, regardless of their distance
- **Query** **Key** dot-products produce attention scores / weights



Summary

- Attention captures dependencies between words in a sequence, regardless of their distance
- **Query** **Key** dot-products produce attention scores / weights



- Attention output **z** is a sum of weighted **values**

Notes

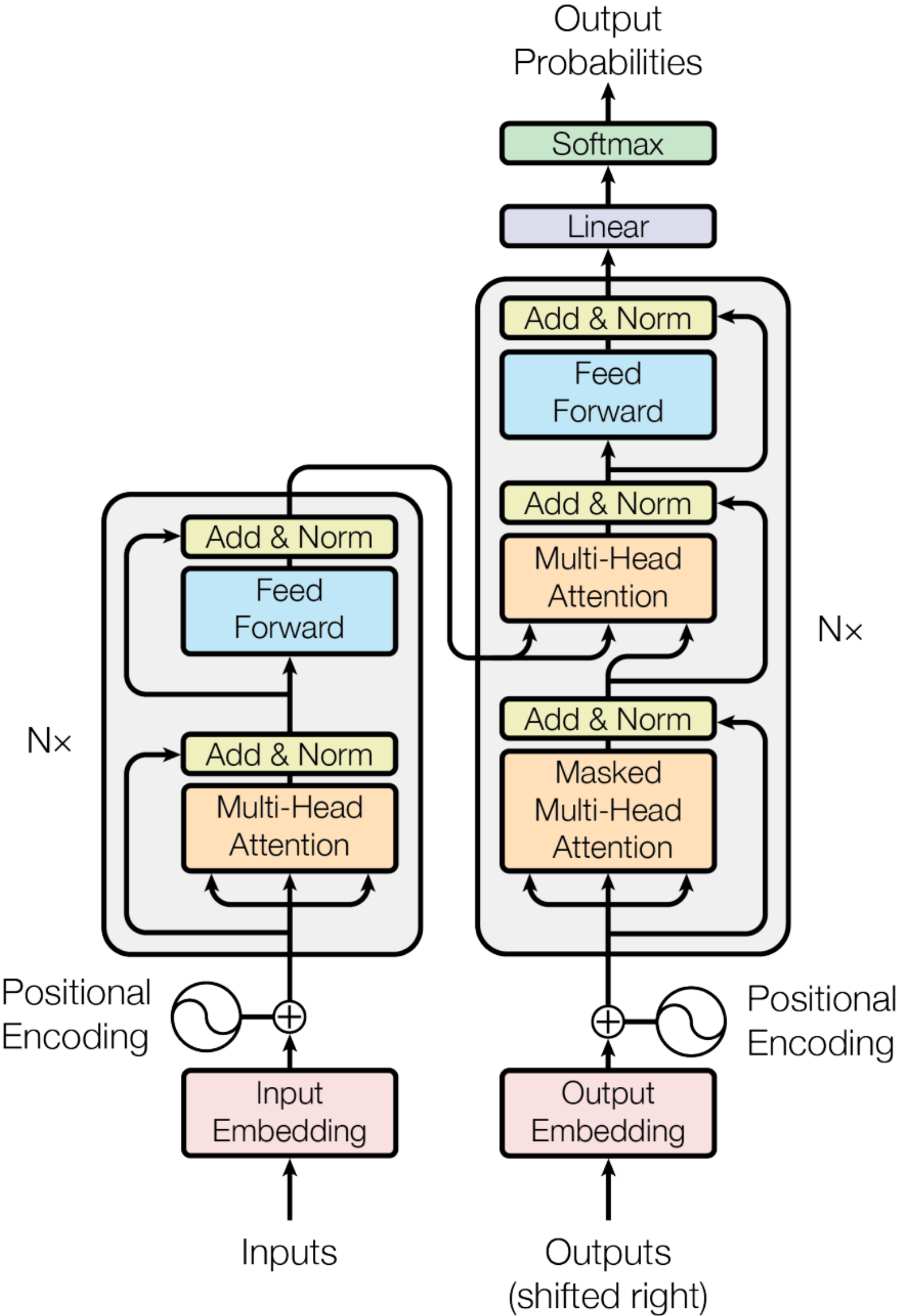


Figure 1: The Transformer - model architecture.

* https://colab.research.google.com/drive/1JMLa53HDuA-i7ZBmqV7ZnA3c_fvtXnx-?usp=sharing

Notes

- Attention is indifferent to the *position* of words, it simply acts over a set of vectors

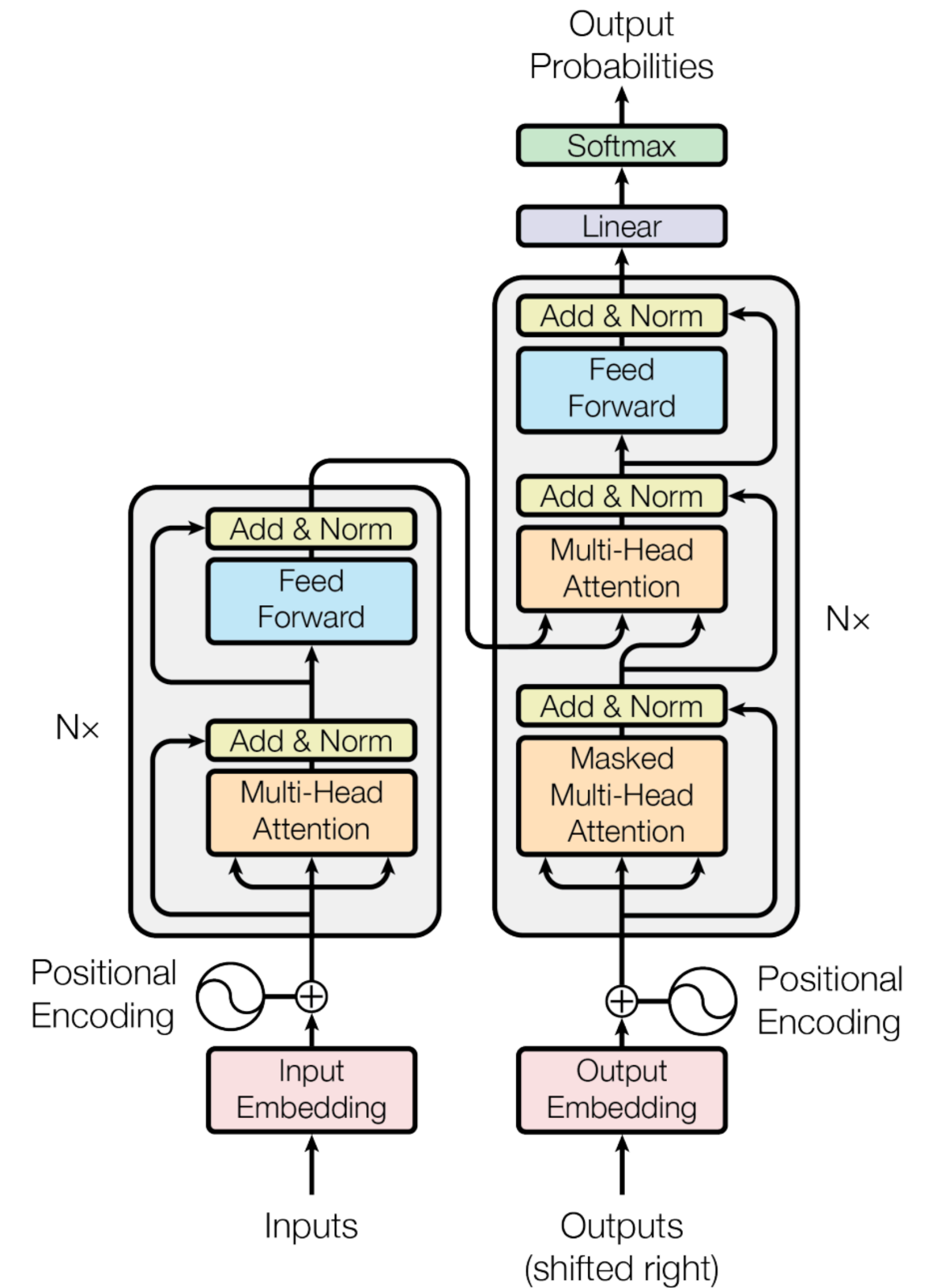


Figure 1: The Transformer - model architecture.

Notes

- Attention is indifferent to the *position* of words, it simply acts over a set of vectors
 - -> Needs positional encoders

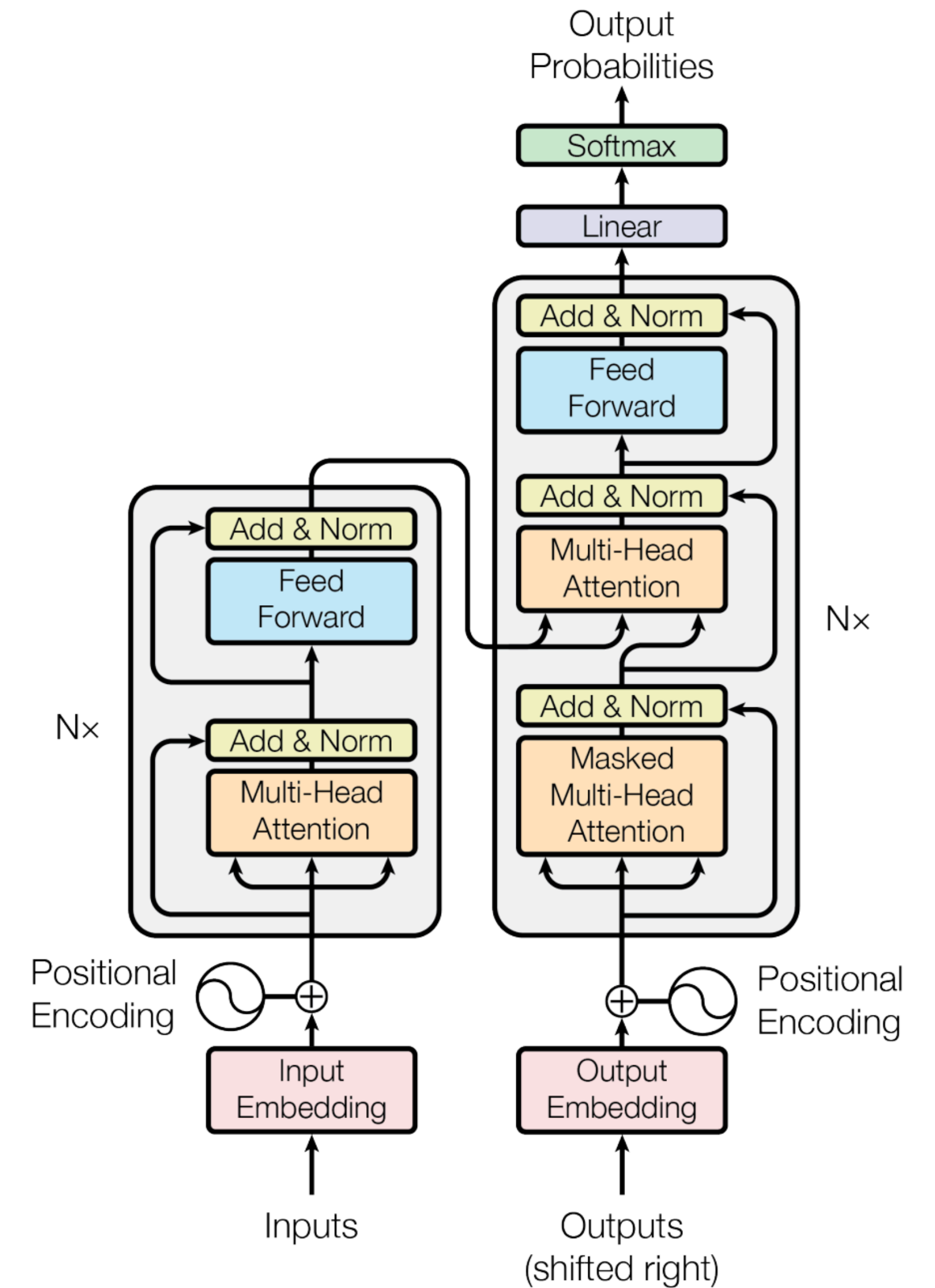


Figure 1: The Transformer - model architecture.

Notes

- Attention is indifferent to the *position* of words, it simply acts over a set of vectors
 - -> Needs positional encoders
- Attention weights are data-dependent and *change during runtime* (unlike Feed Forward NNs / MLPs)

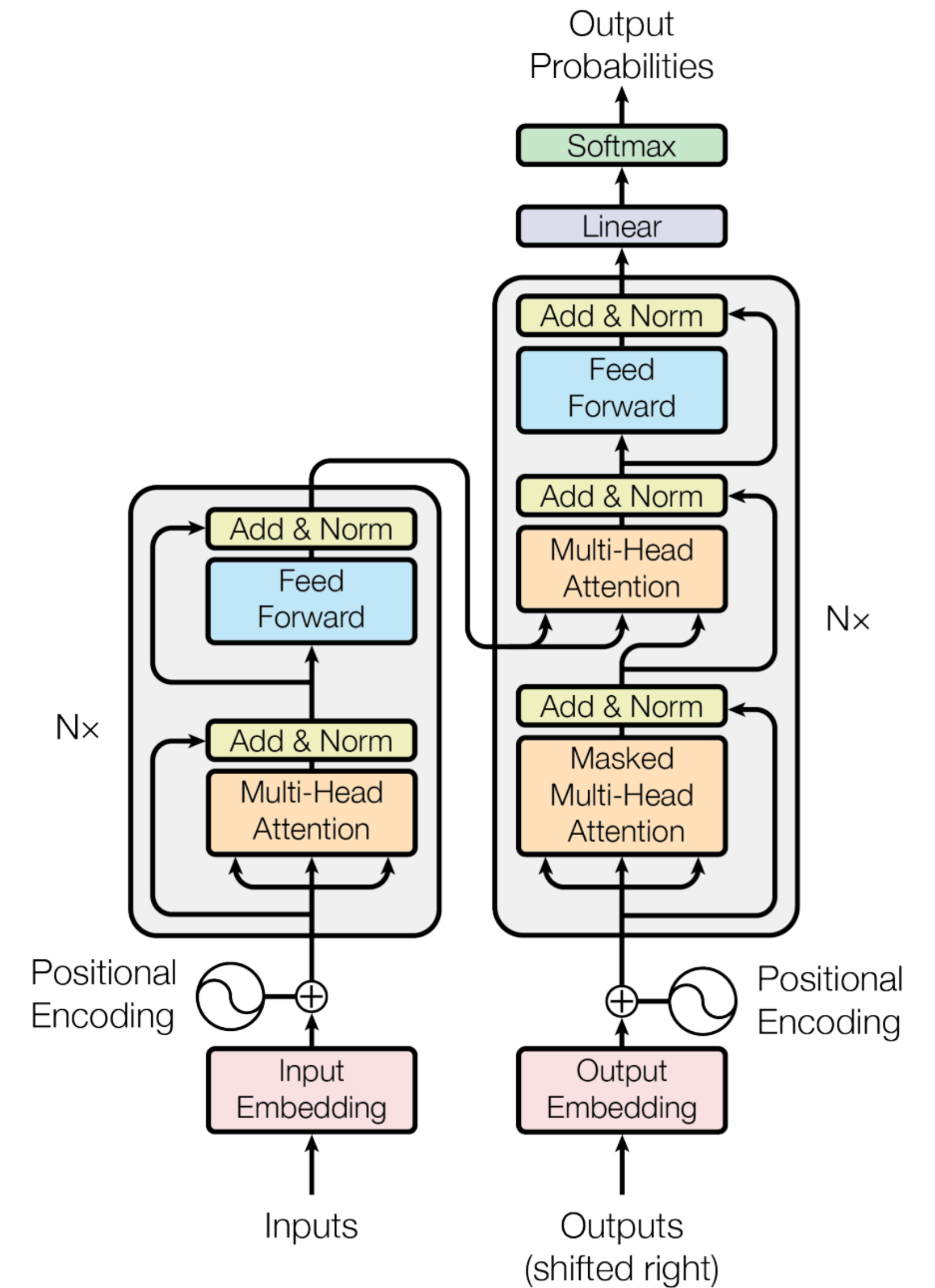


Figure 1: The Transformer - model architecture.

Notes

- Attention is indifferent to the *position* of words, it simply acts over a set of vectors
 - -> Needs positional encoders
- Attention weights are data-dependent and *change during runtime* (unlike Feed Forward NNs / MLPs)
- No communication *between batches* of data

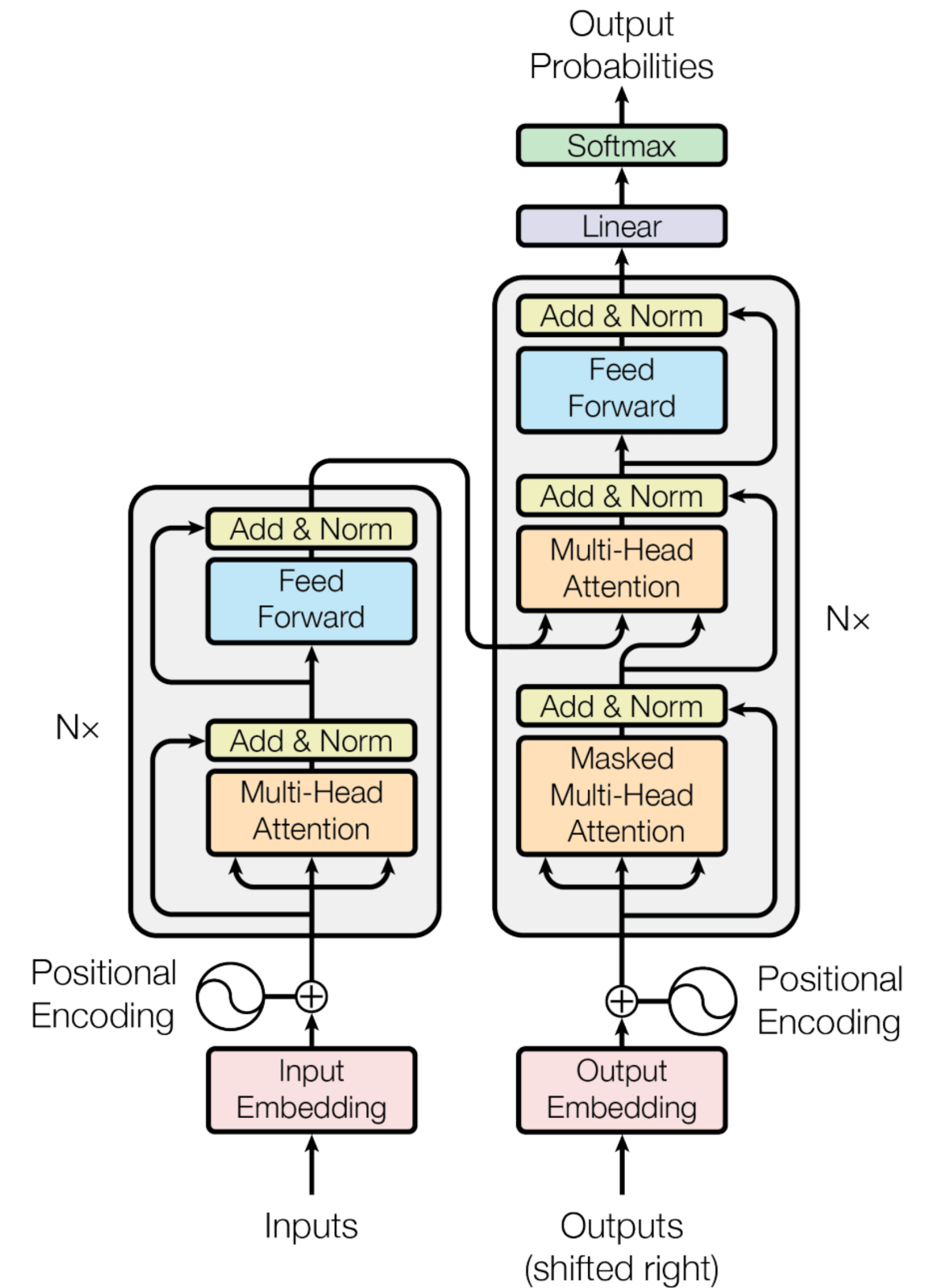


Figure 1: The Transformer - model architecture.

Notes

- Attention is indifferent to the *position* of words, it simply acts over a set of vectors
 - -> Needs positional encoders
- Attention weights are data-dependent and *change during runtime* (unlike Feed Forward NNs / MLPs)
- No communication *between batches* of data
- Attention is a general communication mechanism, tokens can be seen as nodes in a directed graph*

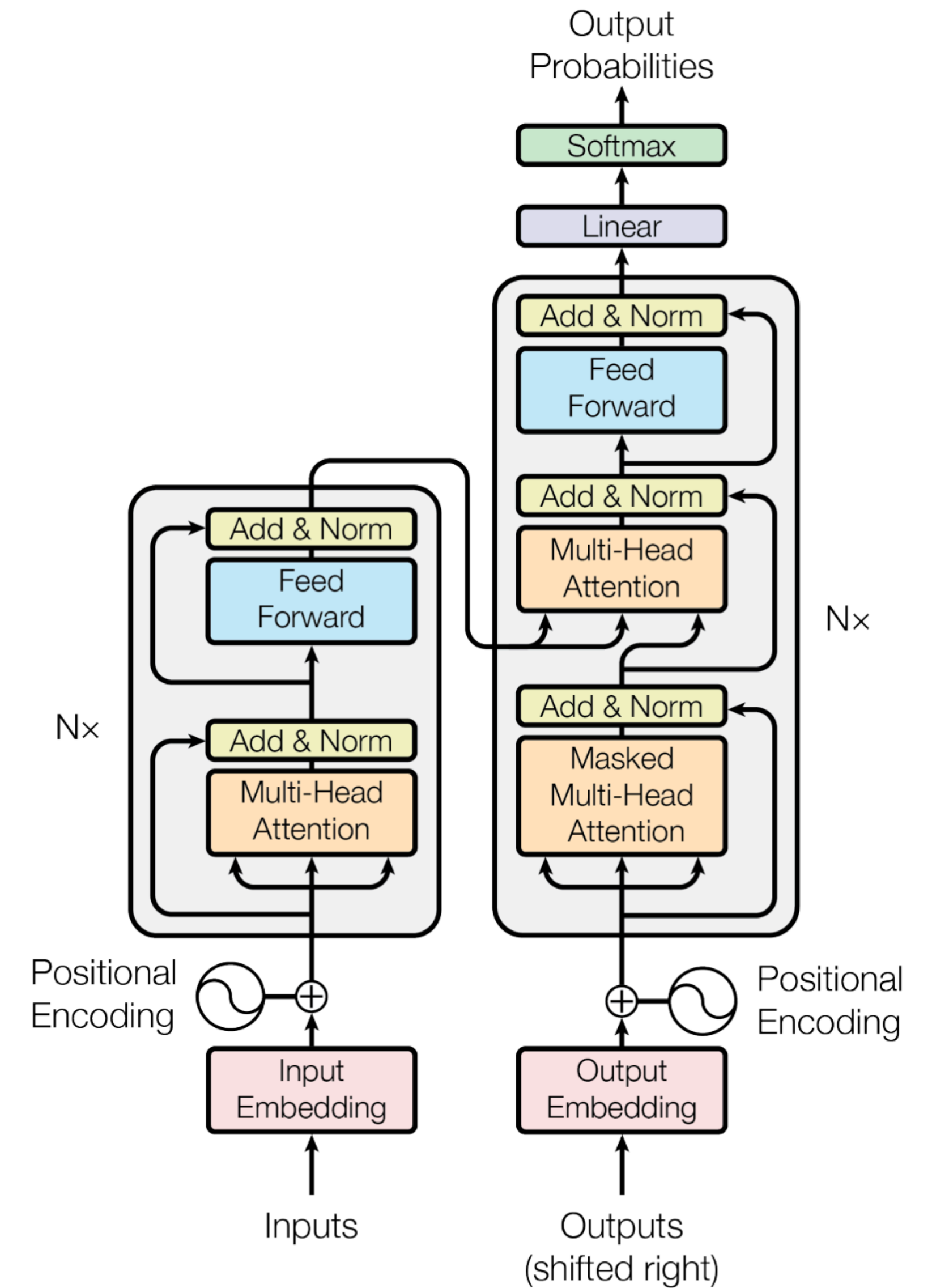


Figure 1: The Transformer - model architecture.

* https://colab.research.google.com/drive/1JMLa53HDuA-i7ZBmqV7ZnA3c_fvtXnx-?usp=sharing

Appendix

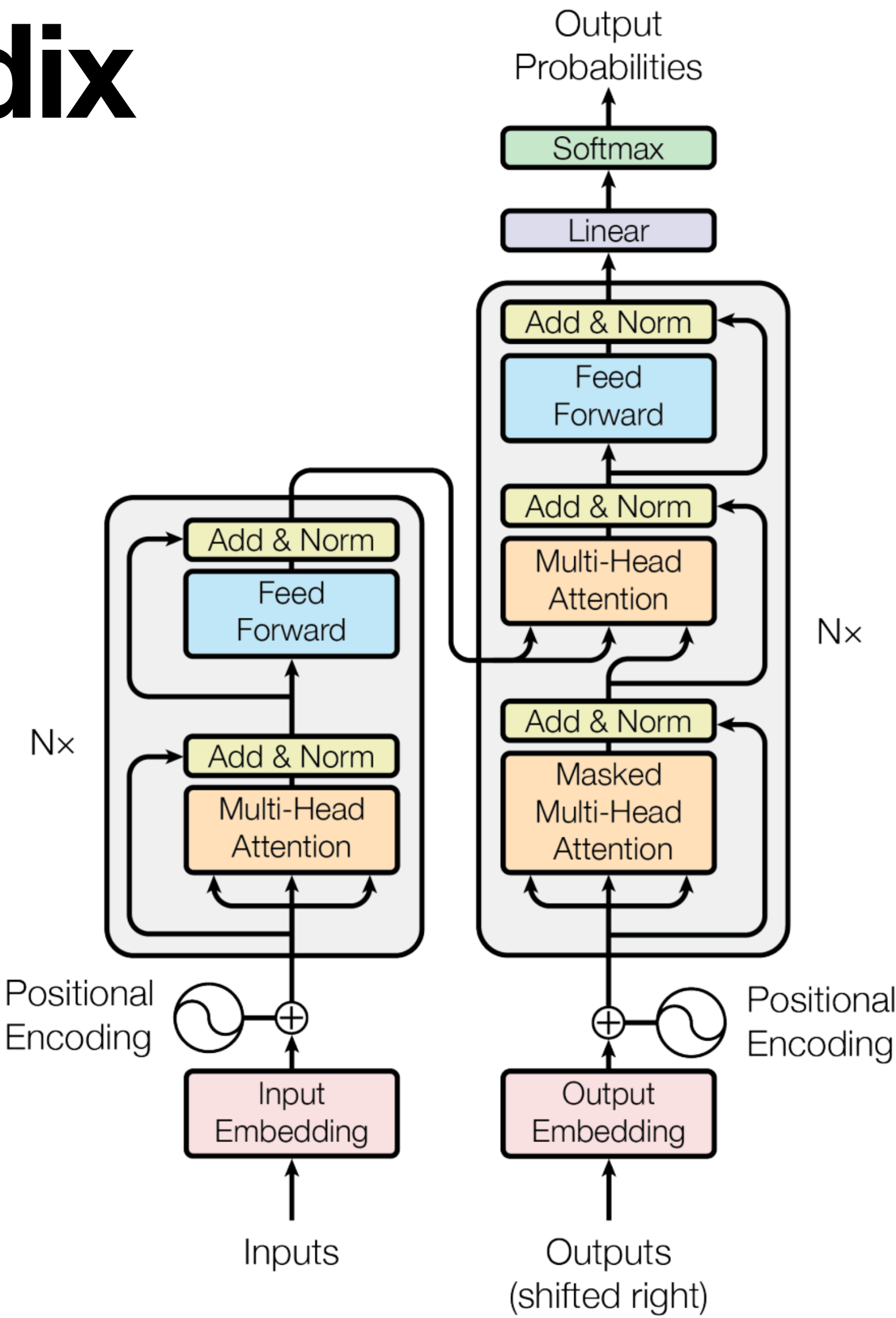
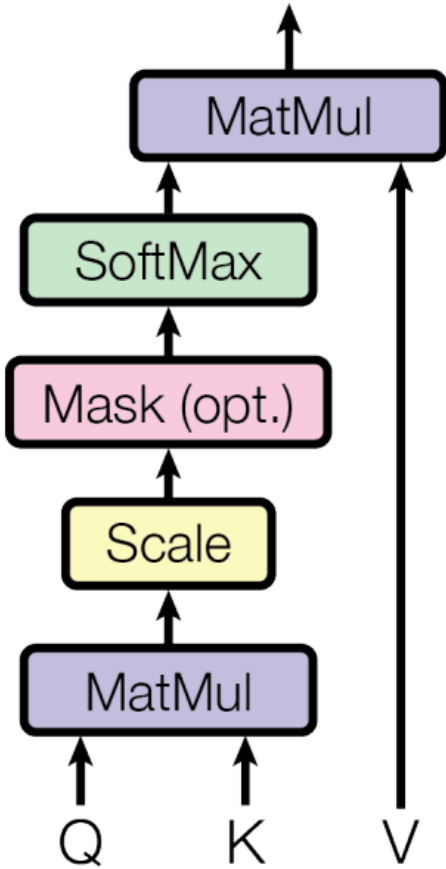


Figure 1: The Transformer - model architecture.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Scaled Dot-Product Attention



Multi-Head Attention

