

Transformers - 1

Introduction to NLP

Rahul Mishra

IIIT-Hyderabad

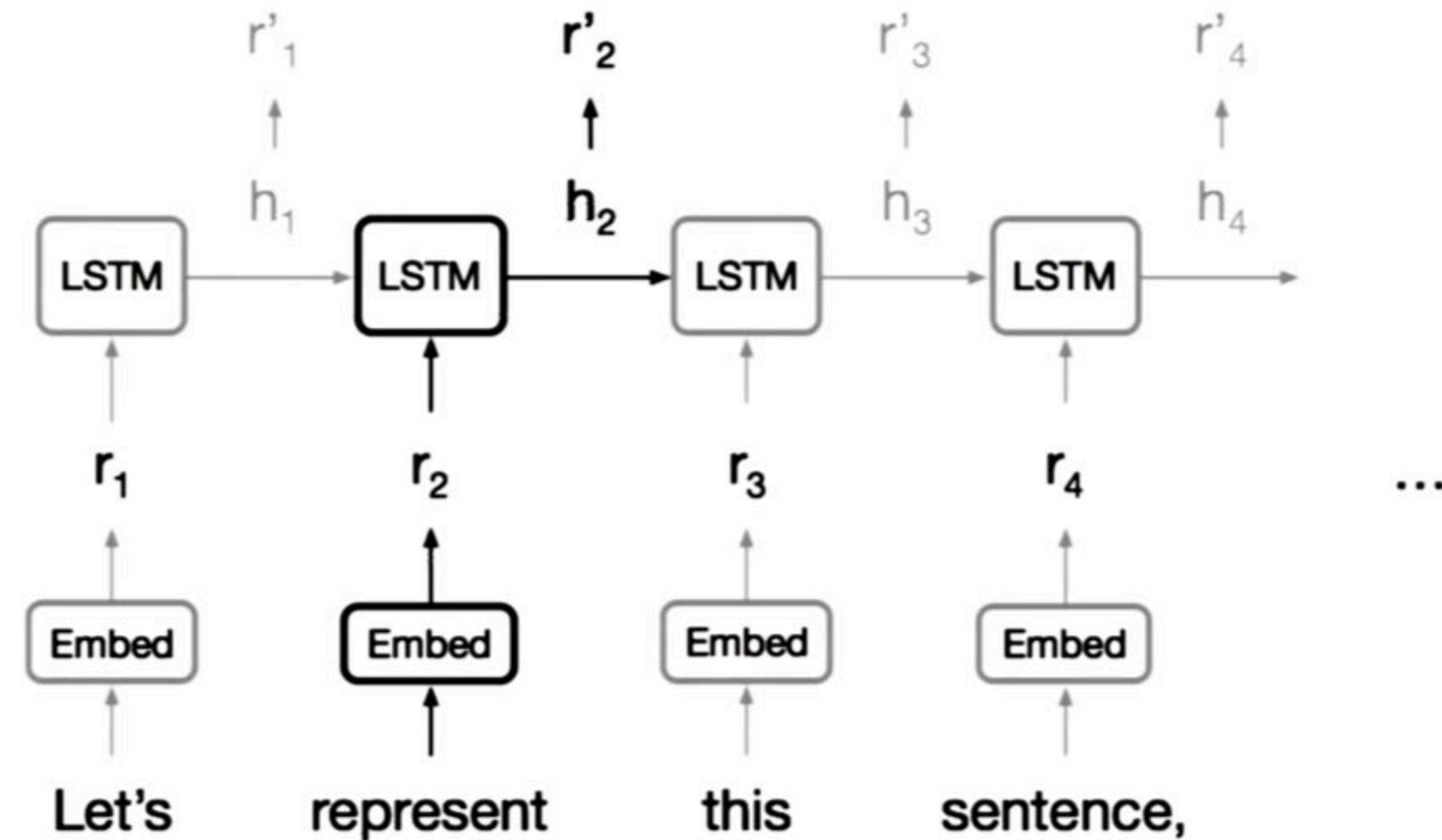
March 15, 2024



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

HYDERABAD

Motivation for a New Architecture



Motivation for a New Architecture

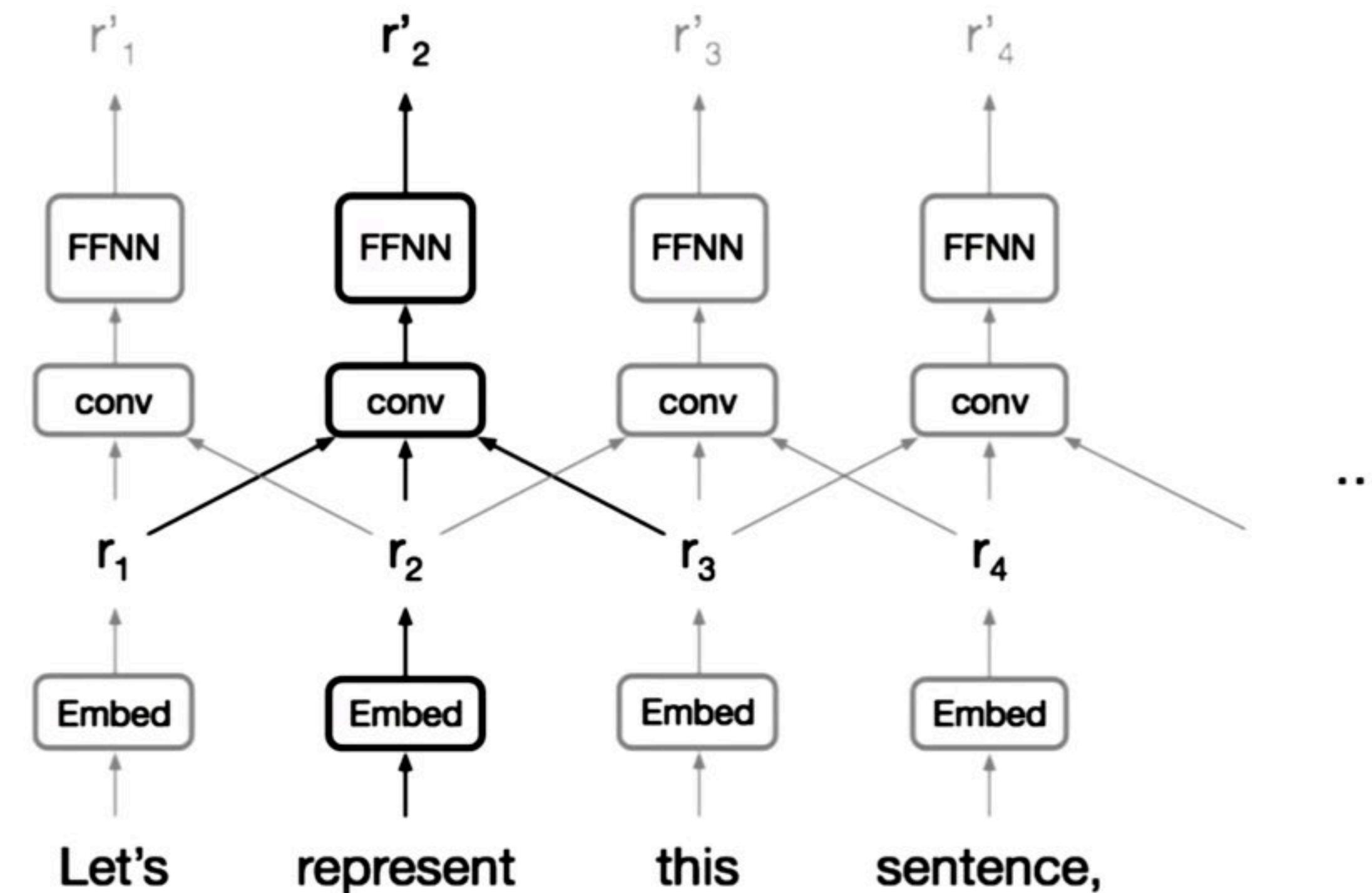
Sequential computation inhibits parallelization.

No explicit modeling of long and short range dependencies.

We want to model hierarchy.

RNNs (w/ sequence-aligned states) seem wasteful!

Motivation for a New Architecture



Motivation for a New Architecture

Trivial to parallelize (per layer).

Exploits local dependencies

‘Interaction distance’ between positions linear or logarithmic.

Long-distance dependencies require many layers.

Motivation for a New Architecture

- Vanishing gradient problem
- Parallel Processing
- No Recurrent Connections
- Better Capturing Long-Term Dependencies



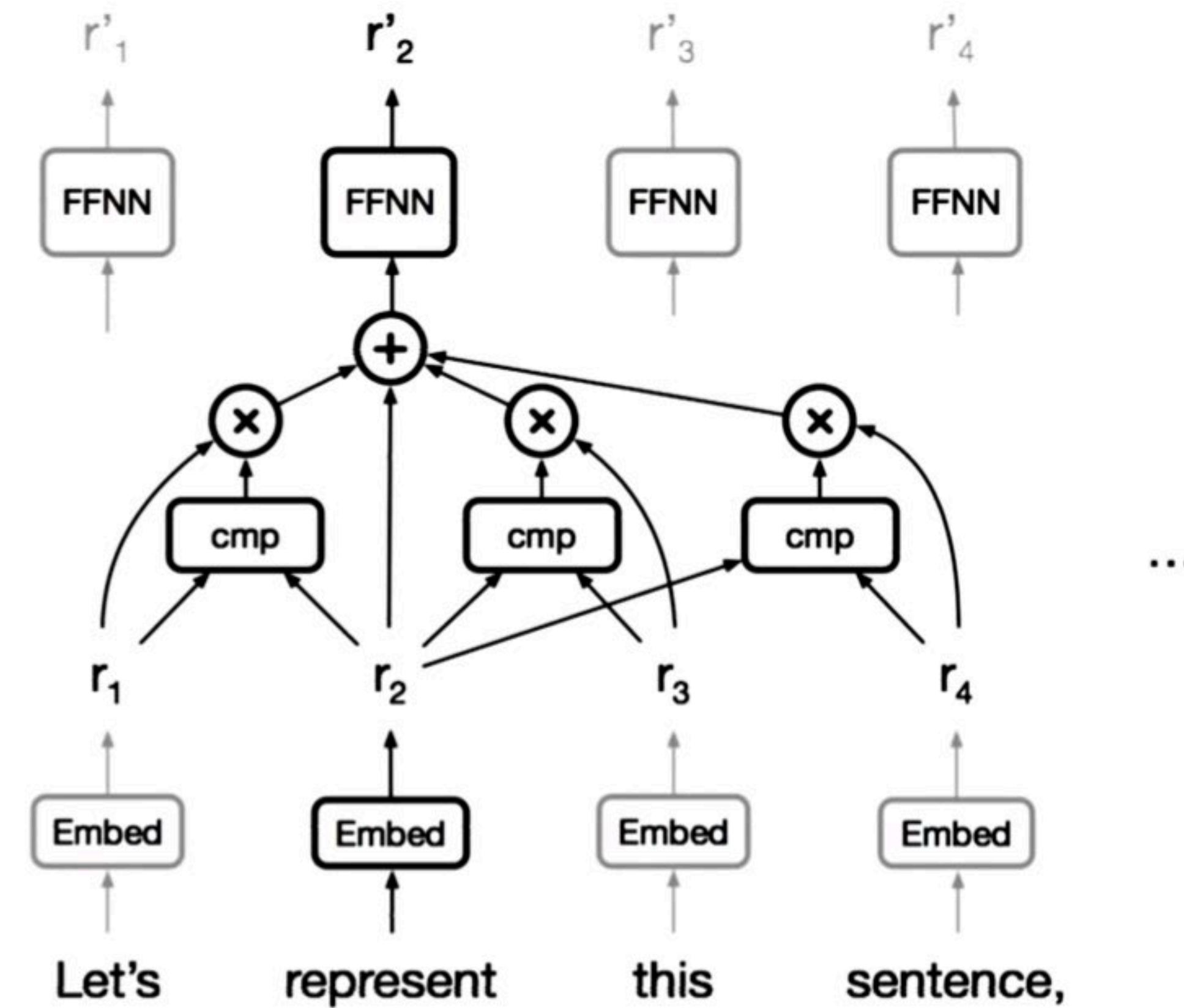
Attention is All You Need:Transformers



@grownupgaming 1 year ago

I would be scared for my life if I was Ashish. Someone from the future might travel back in time and kill me.

Self-Attention



Self-Attention

Classification & regression with self-attention:

Parikh et al. (2016), Lin et al. (2016)

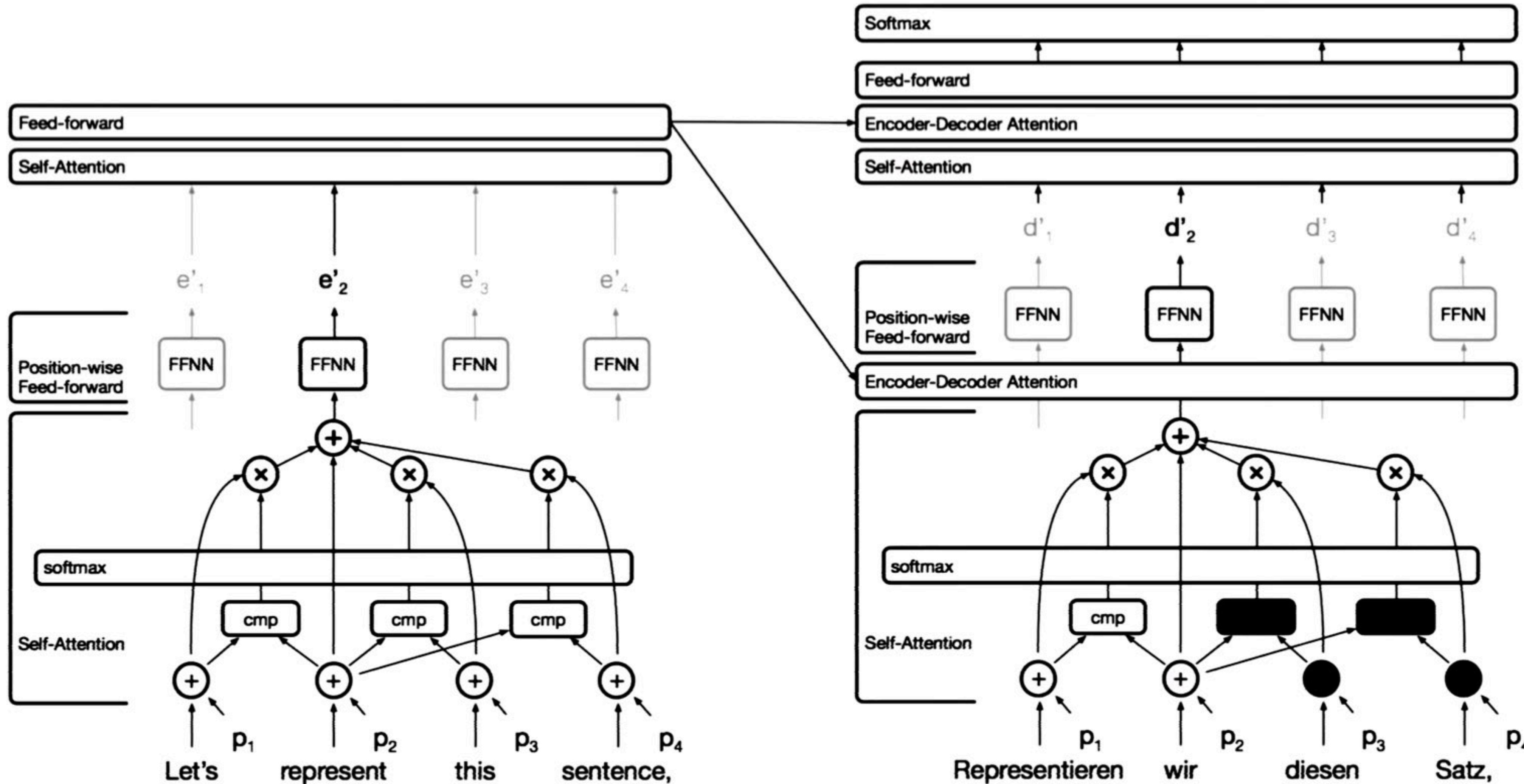
Self-attention with RNNs:

Long et al. (2016), Shao, Gows et al. (2017)

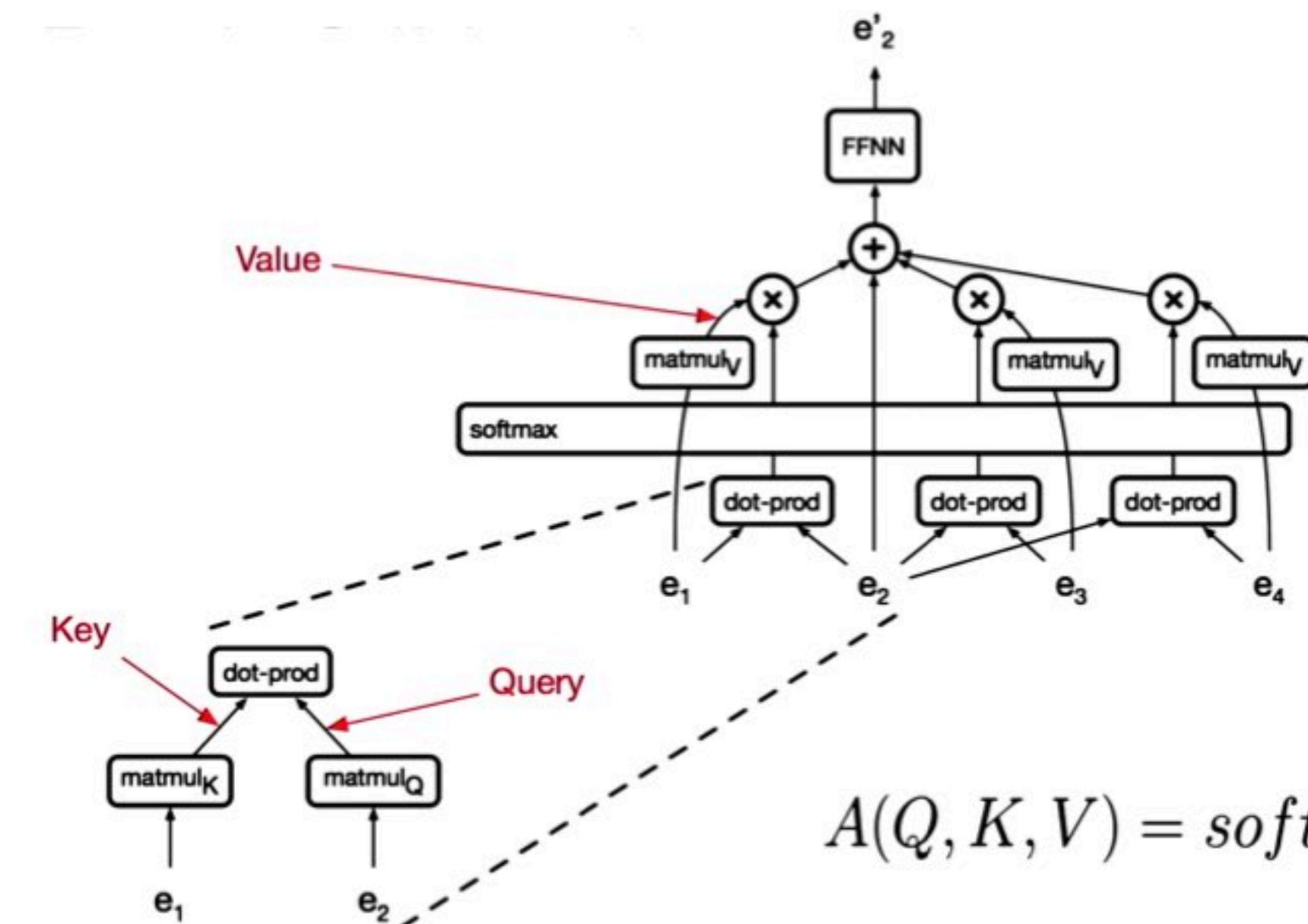
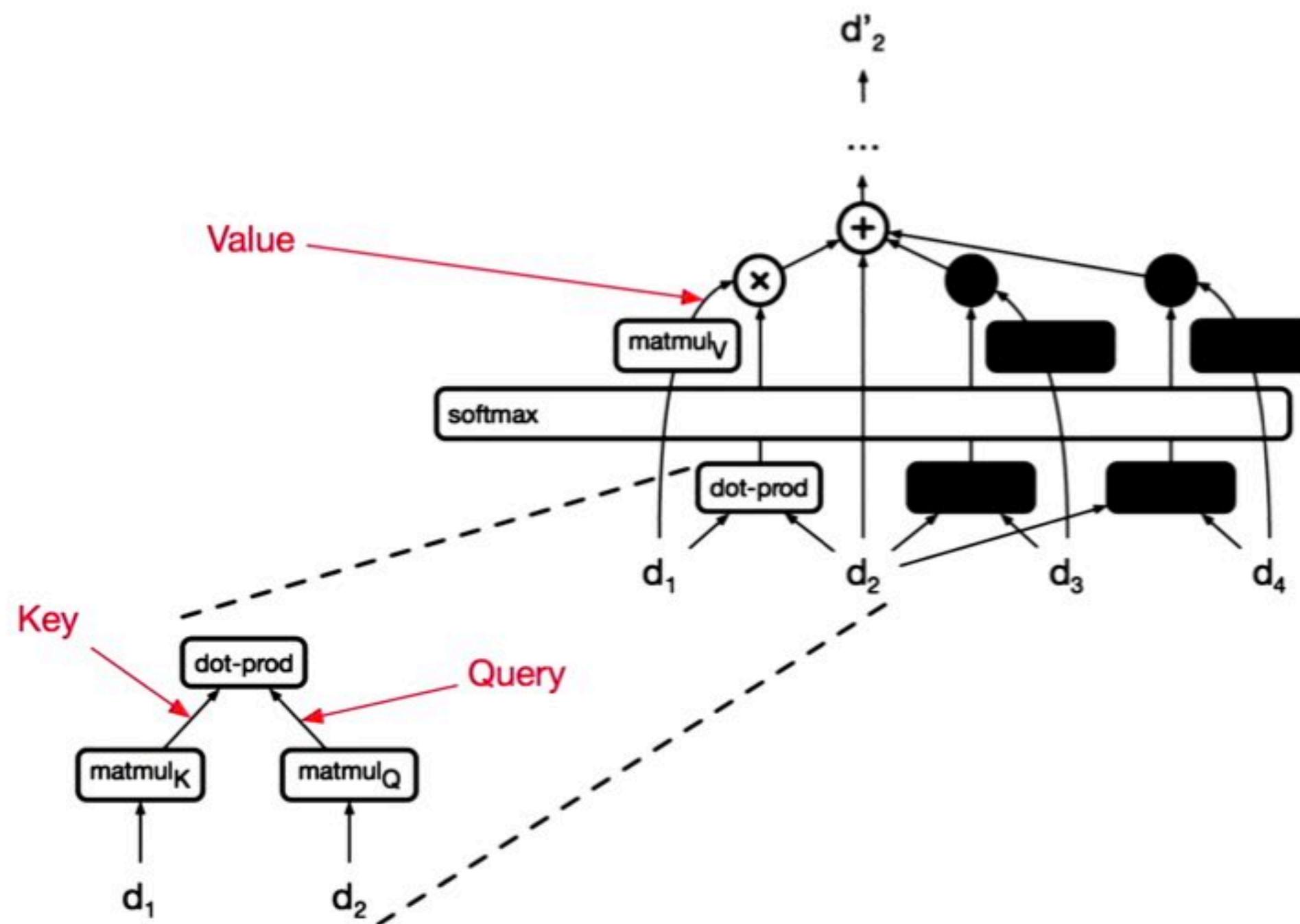
Recurrent attention:

Sukhbaatar et al. (2015)

Rough Idea of Transformers



Rough Idea of Transformers



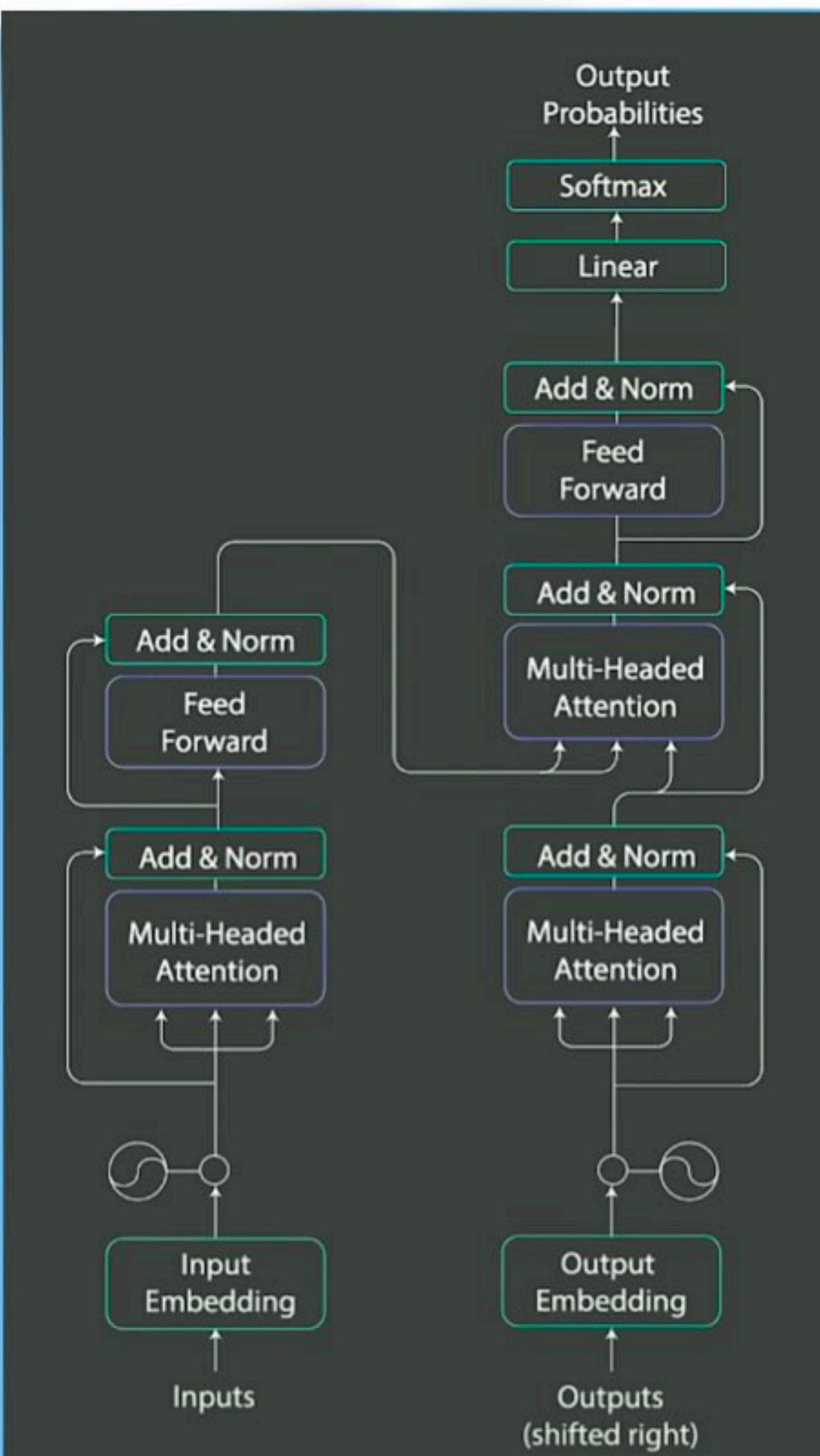
$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Is it complex? Yes and No

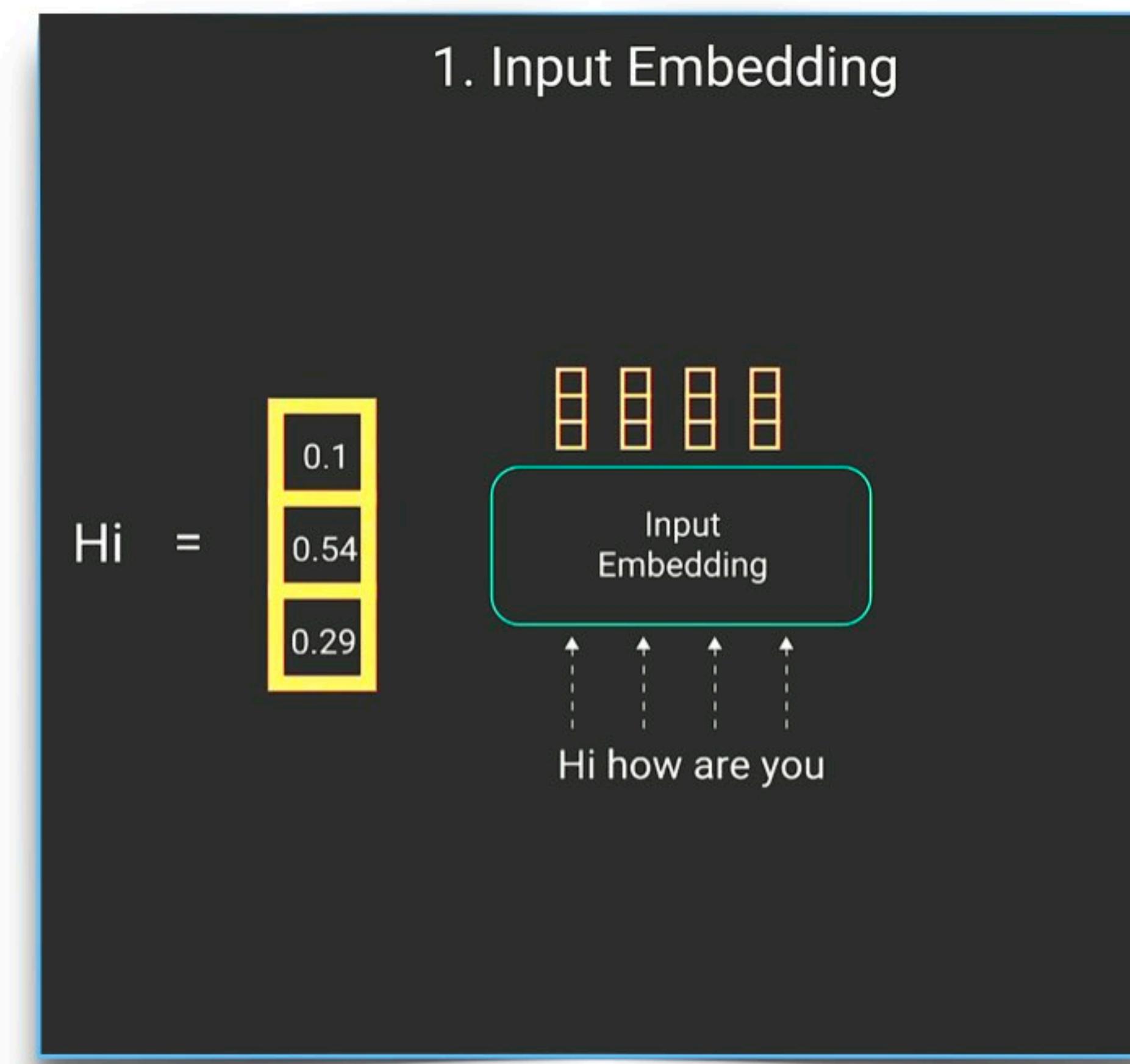
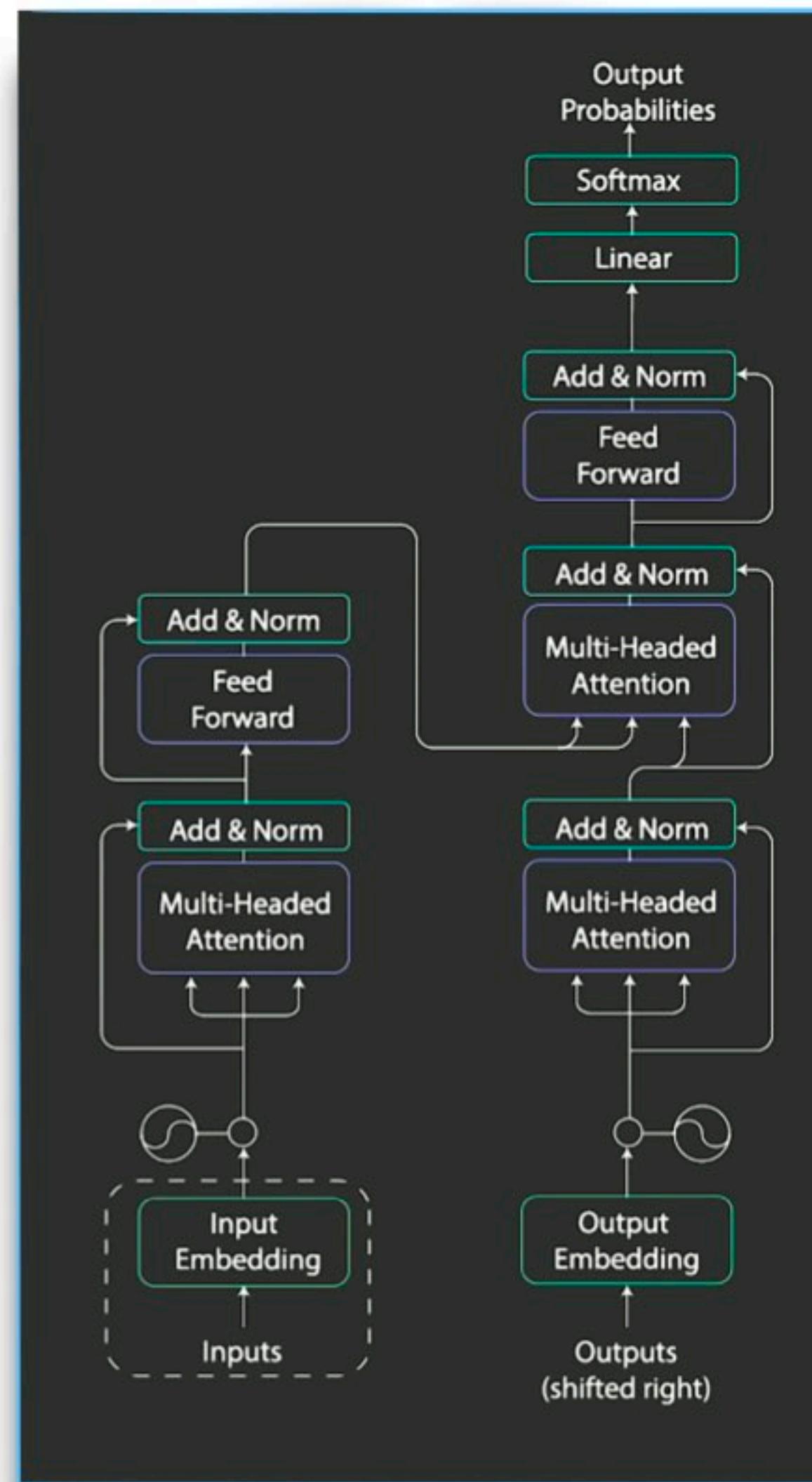
FLOPs

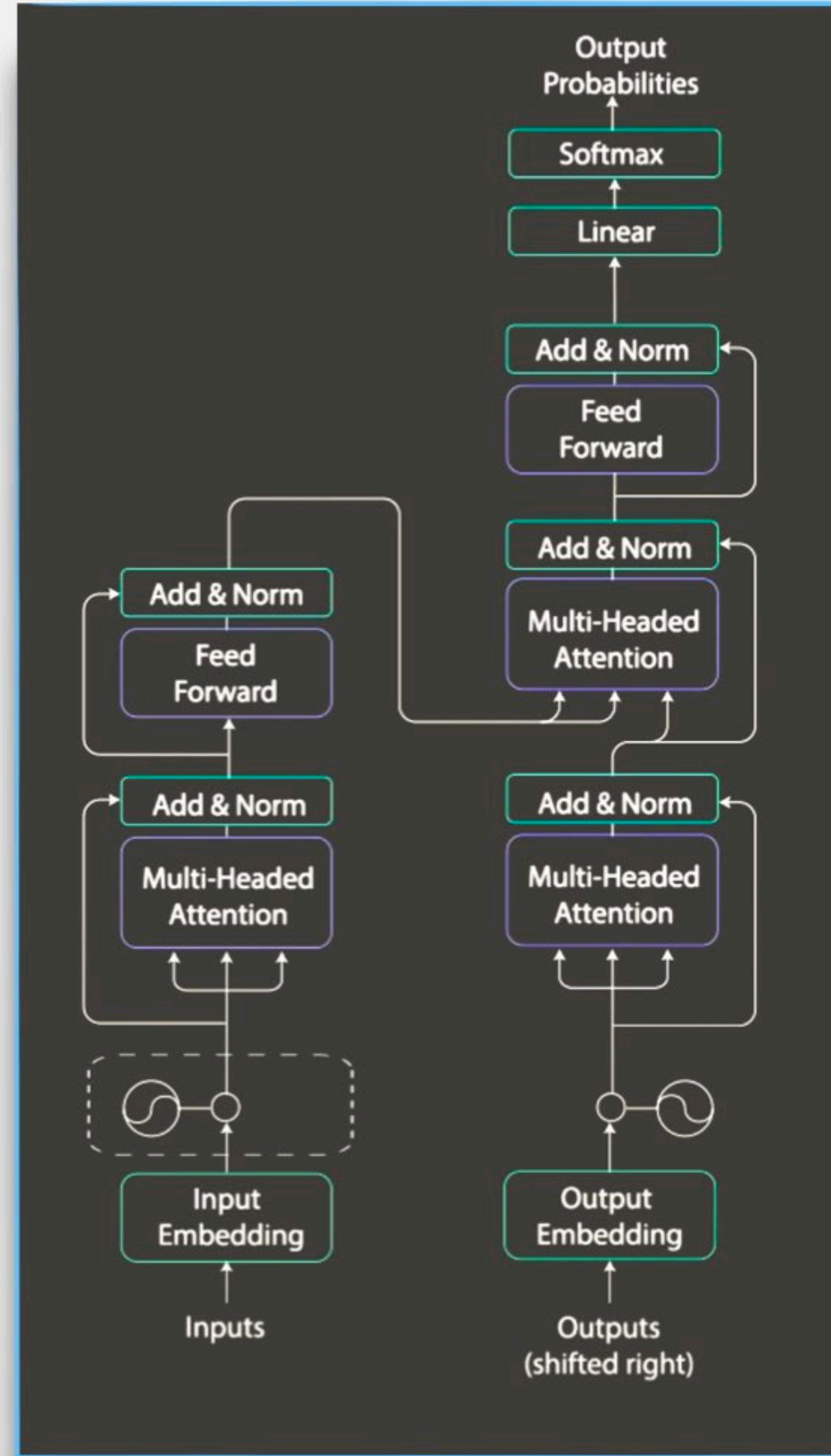
| | |
|----------------|------------------------------------------------------------------|
| Self-Attention | $O(\text{length}^2 \cdot \text{dim})$ |
| RNN (LSTM) | $O(\text{length} \cdot \text{dim}^2)$ |
| Convolution | $O(\text{length} \cdot \text{dim}^2 \cdot \text{kernel_width})$ |

Transformers: Attention Is All You Need



Transformers: Attention Is All You Need





2. Positional Encoding

Positional Input Embeddings



Positional Encoding

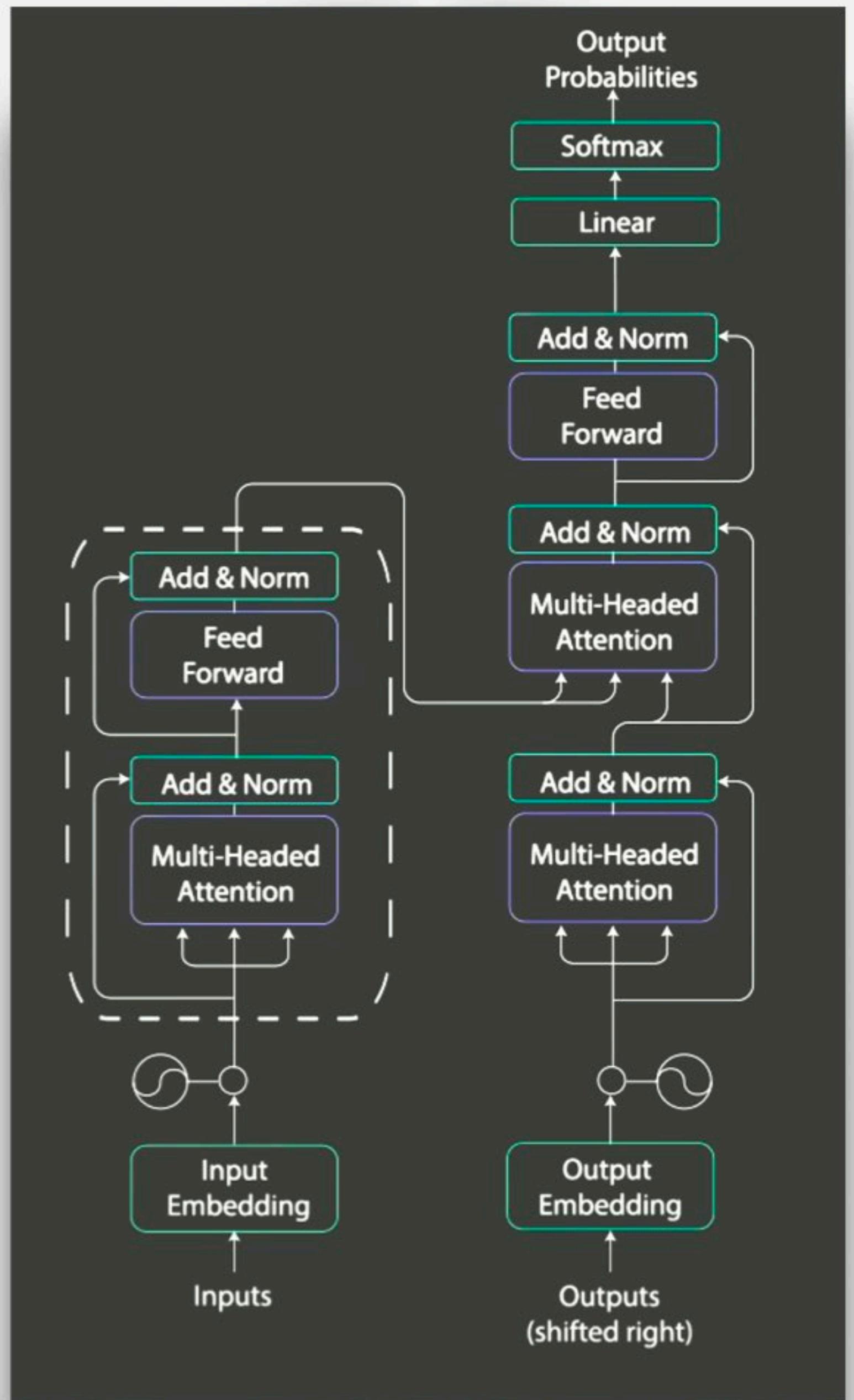


Time Step

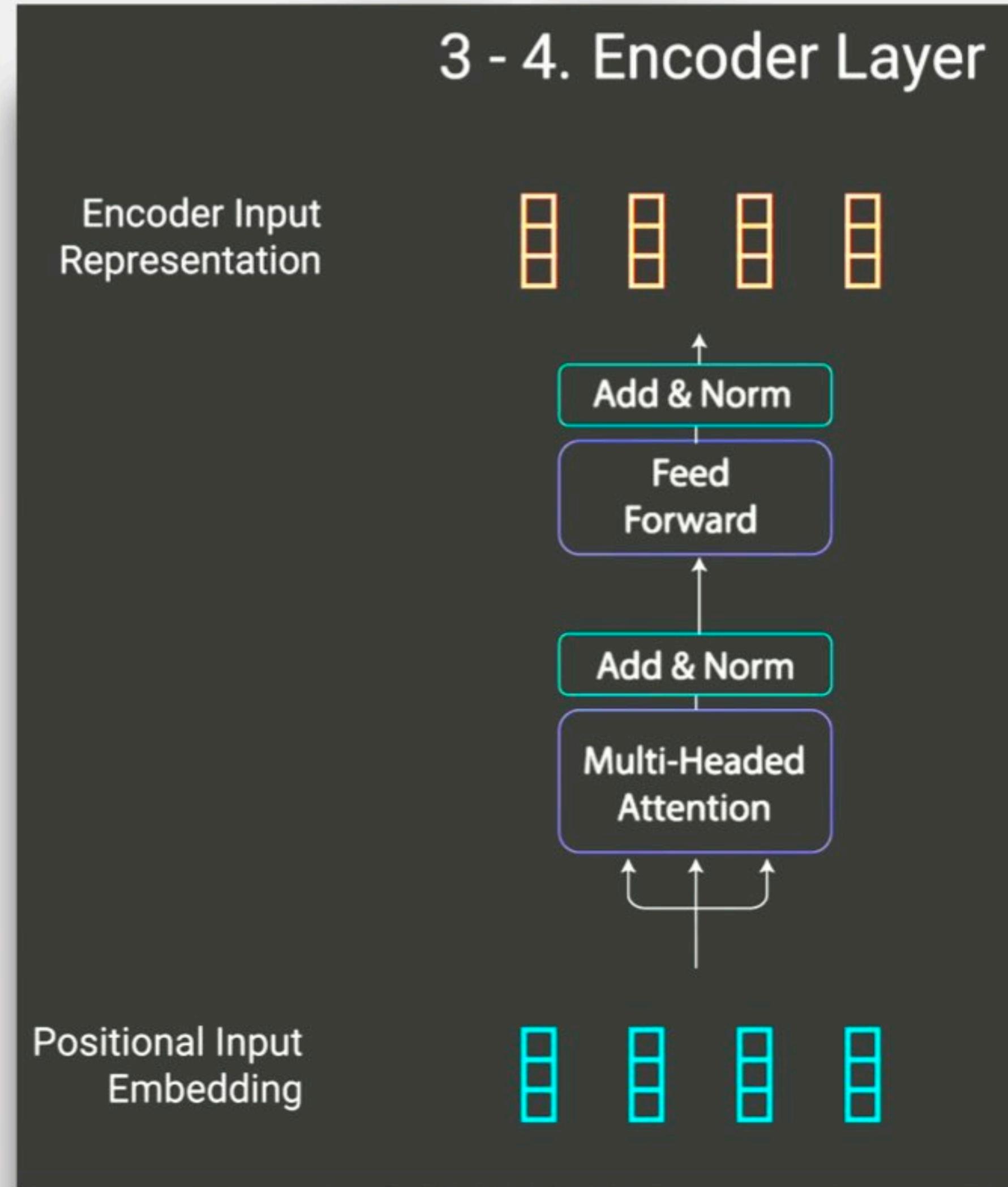
1 2 3 4

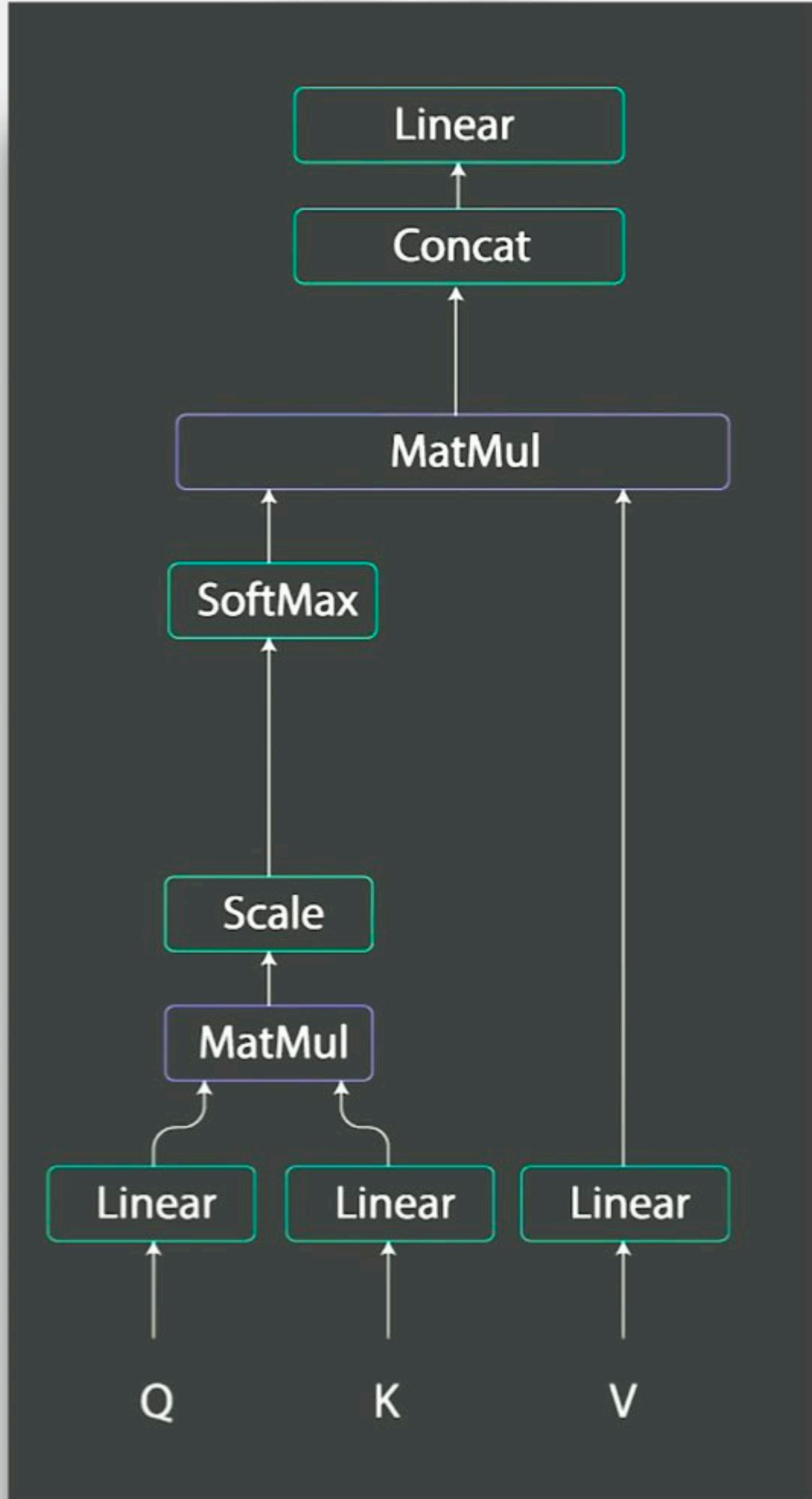
$$PE(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right)$$



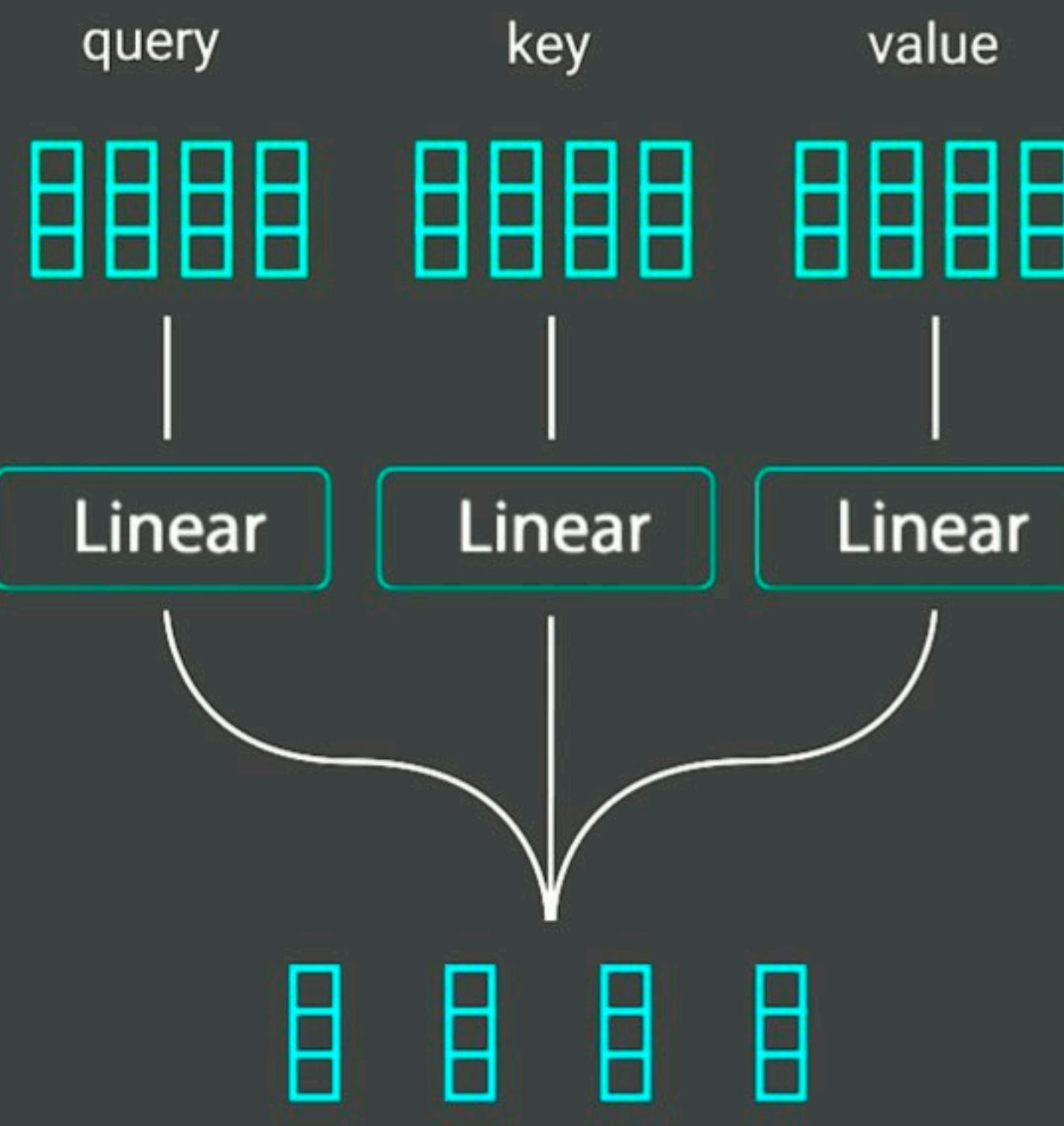
3 - 4. Encoder Layer

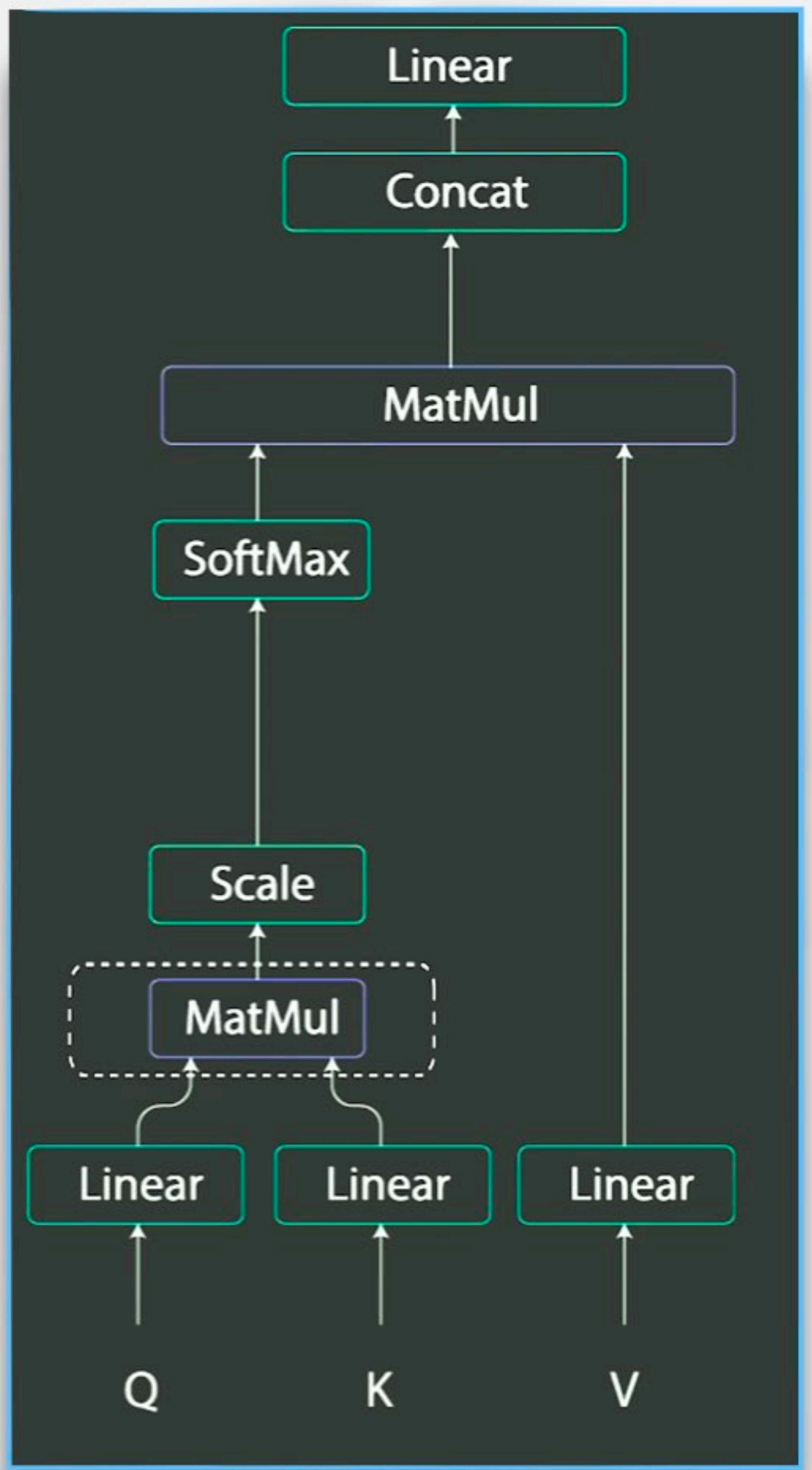




3. Multi-headed Attention

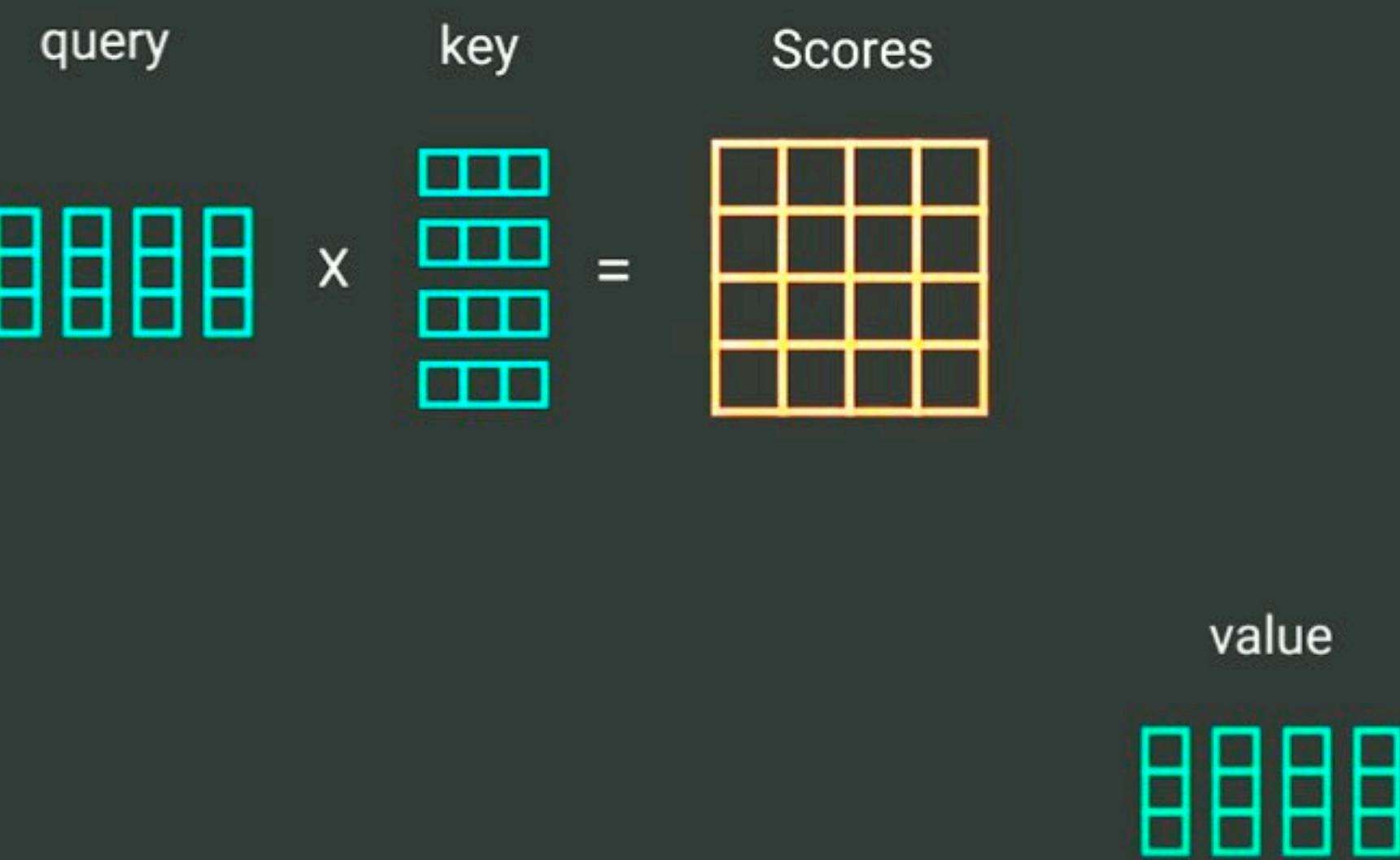
3.1. Self-Attention

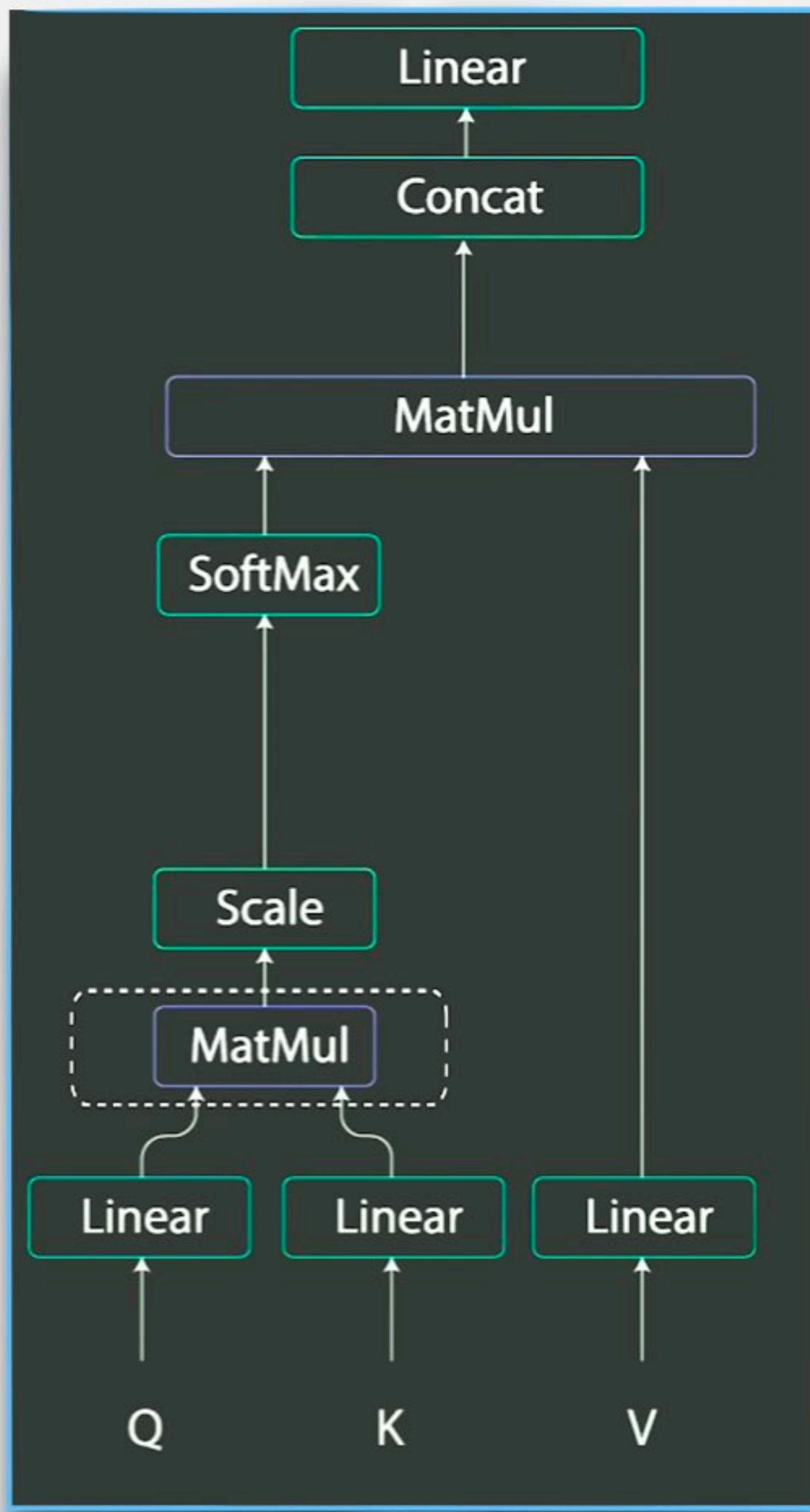




3. Multi-headed Attention

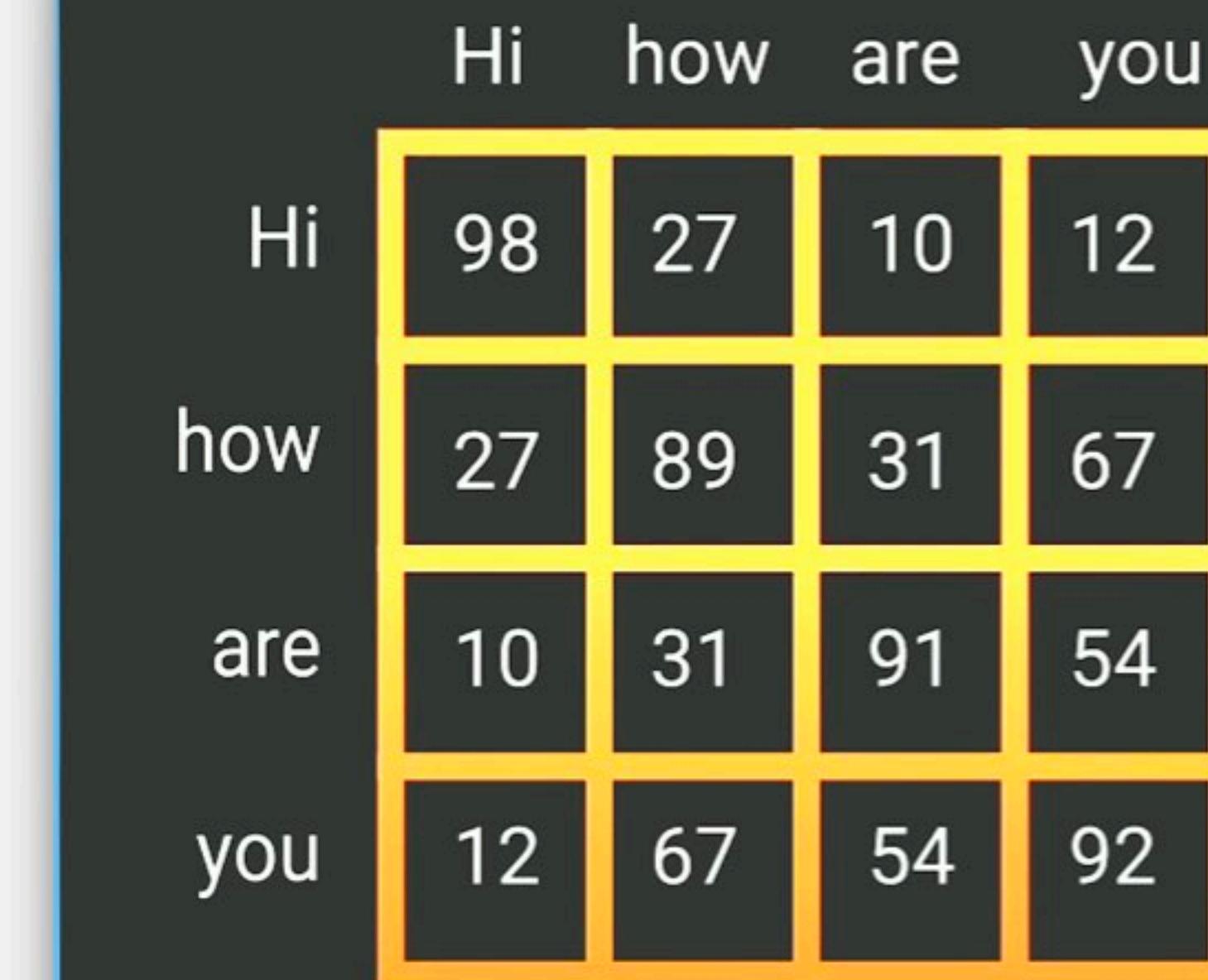
3.1. Self-Attention

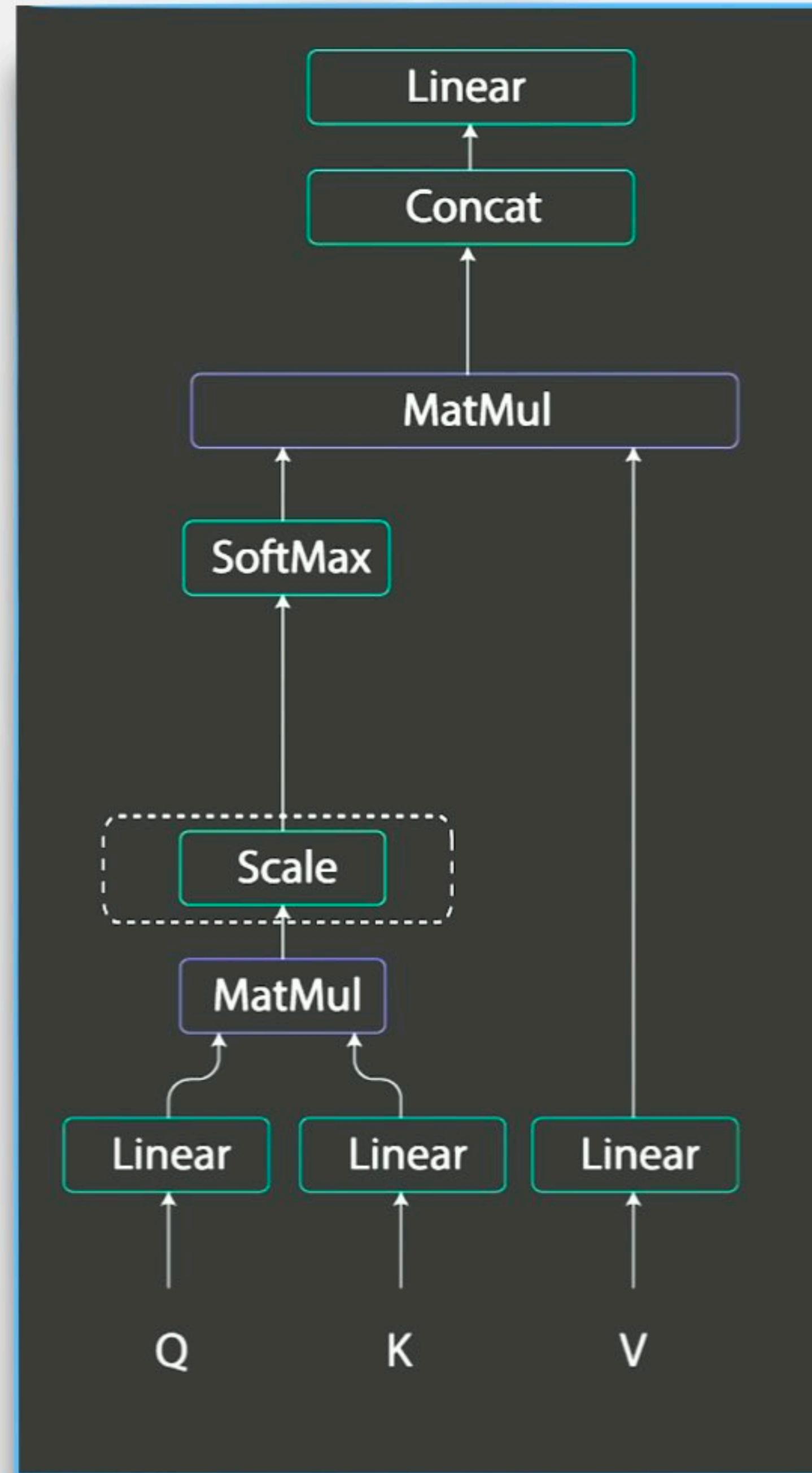




3. Multi-headed Attention

3.1. Self-Attention



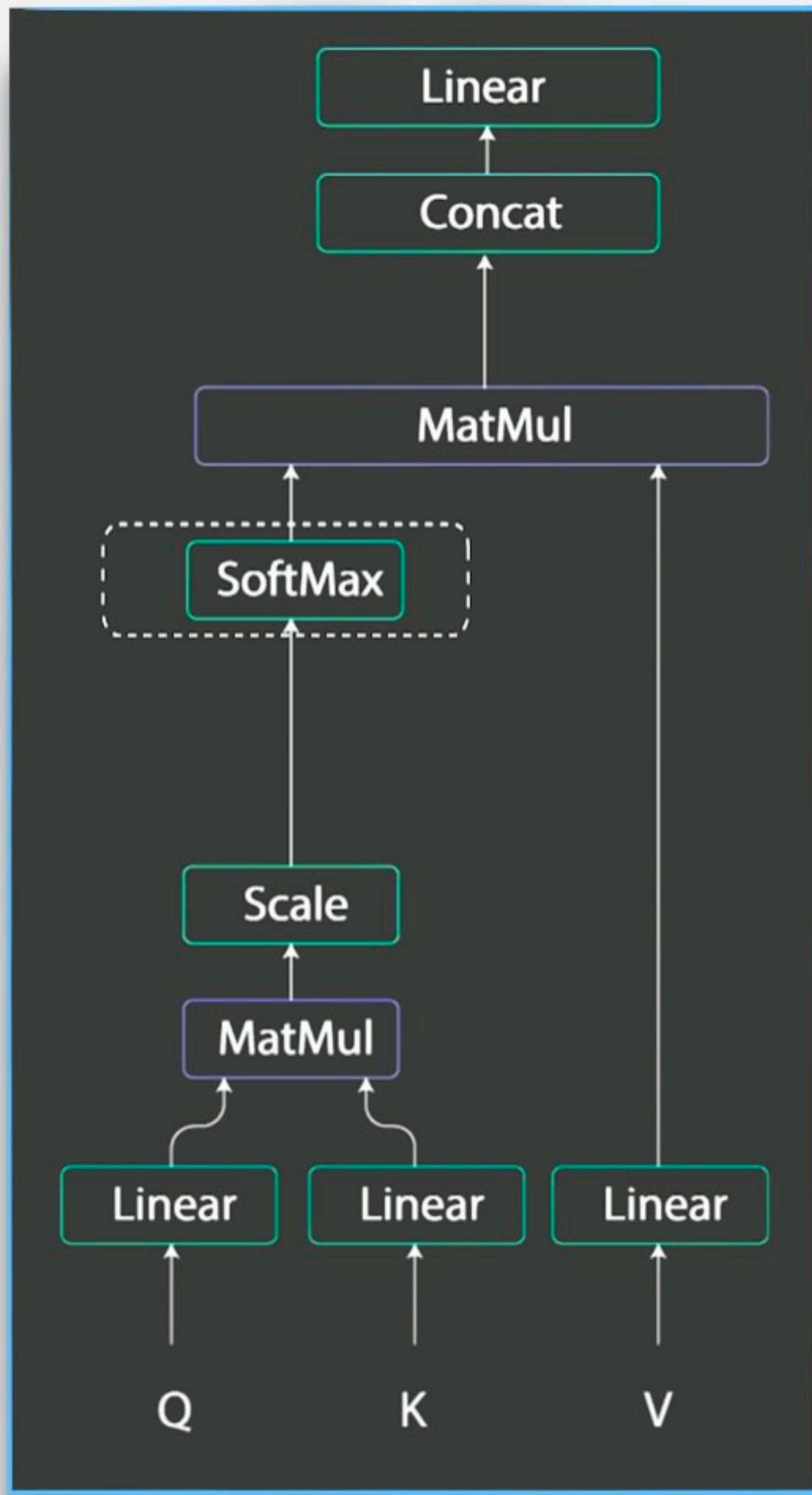


3. Multi-headed Attention

3.1. Self-Attention

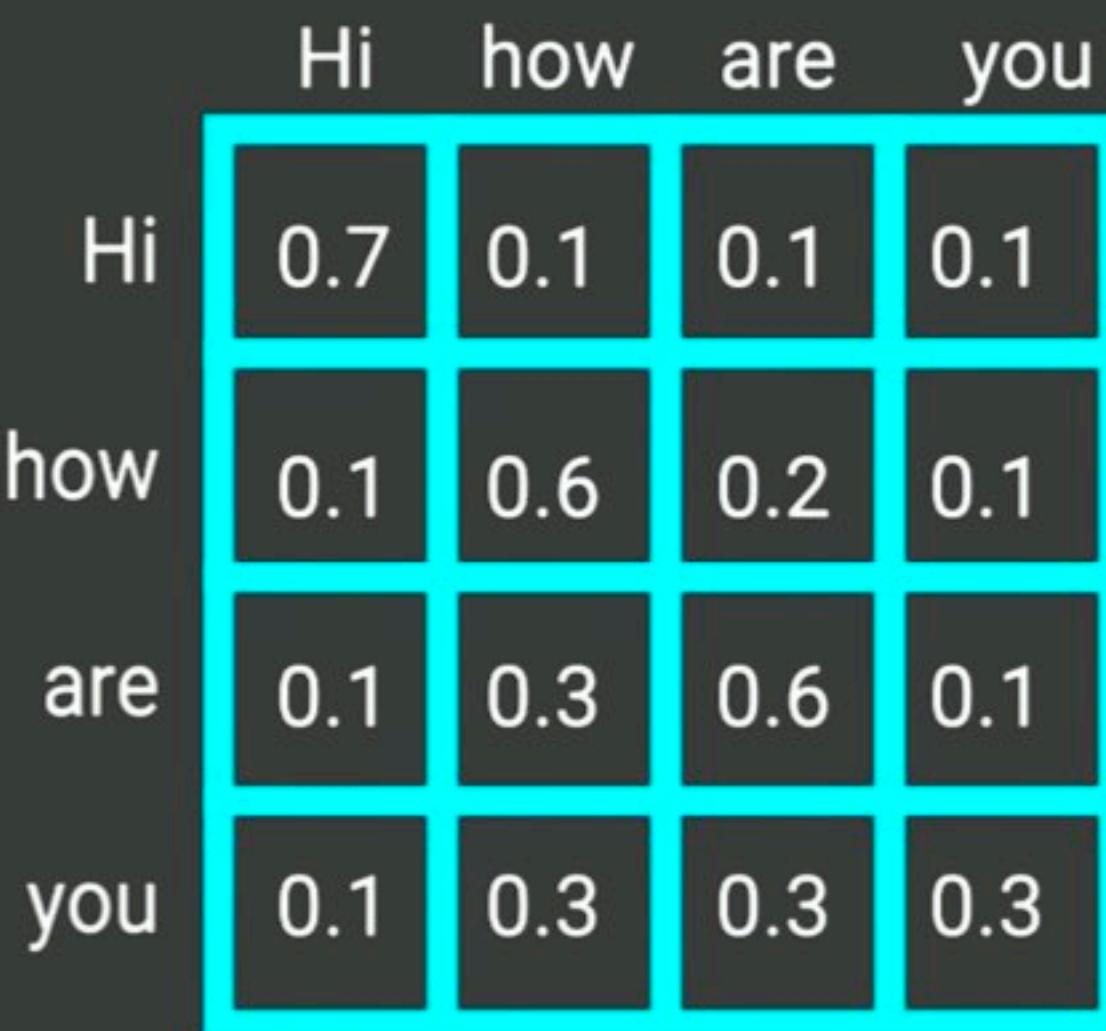
Scaled Scores

$$\frac{\text{Scaled Scores}}{\sqrt{d_k}} = \text{Output}$$

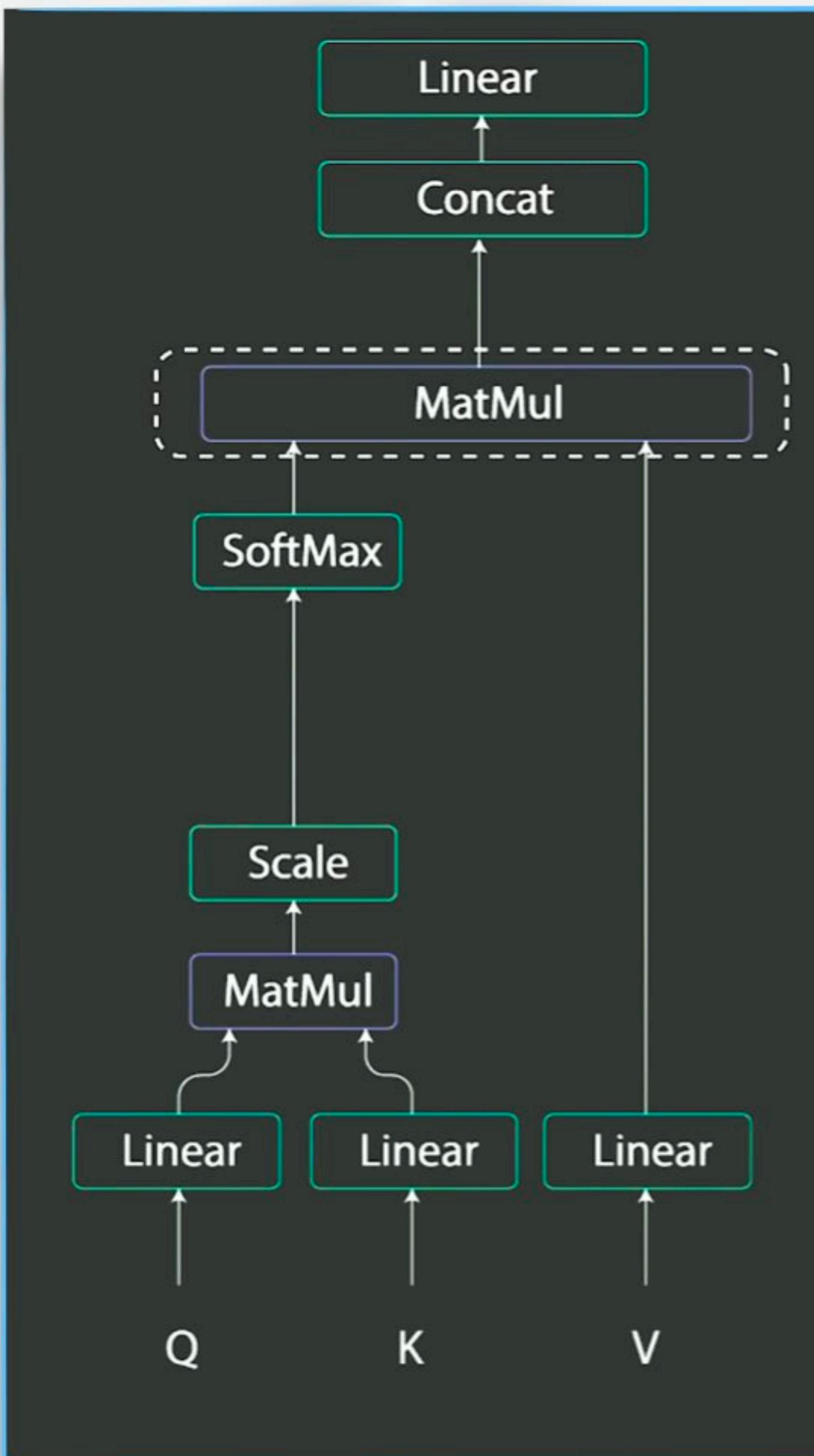


3. Multi-headed Attention

3.1. Self-Attention



$$\text{softmax}(x)_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

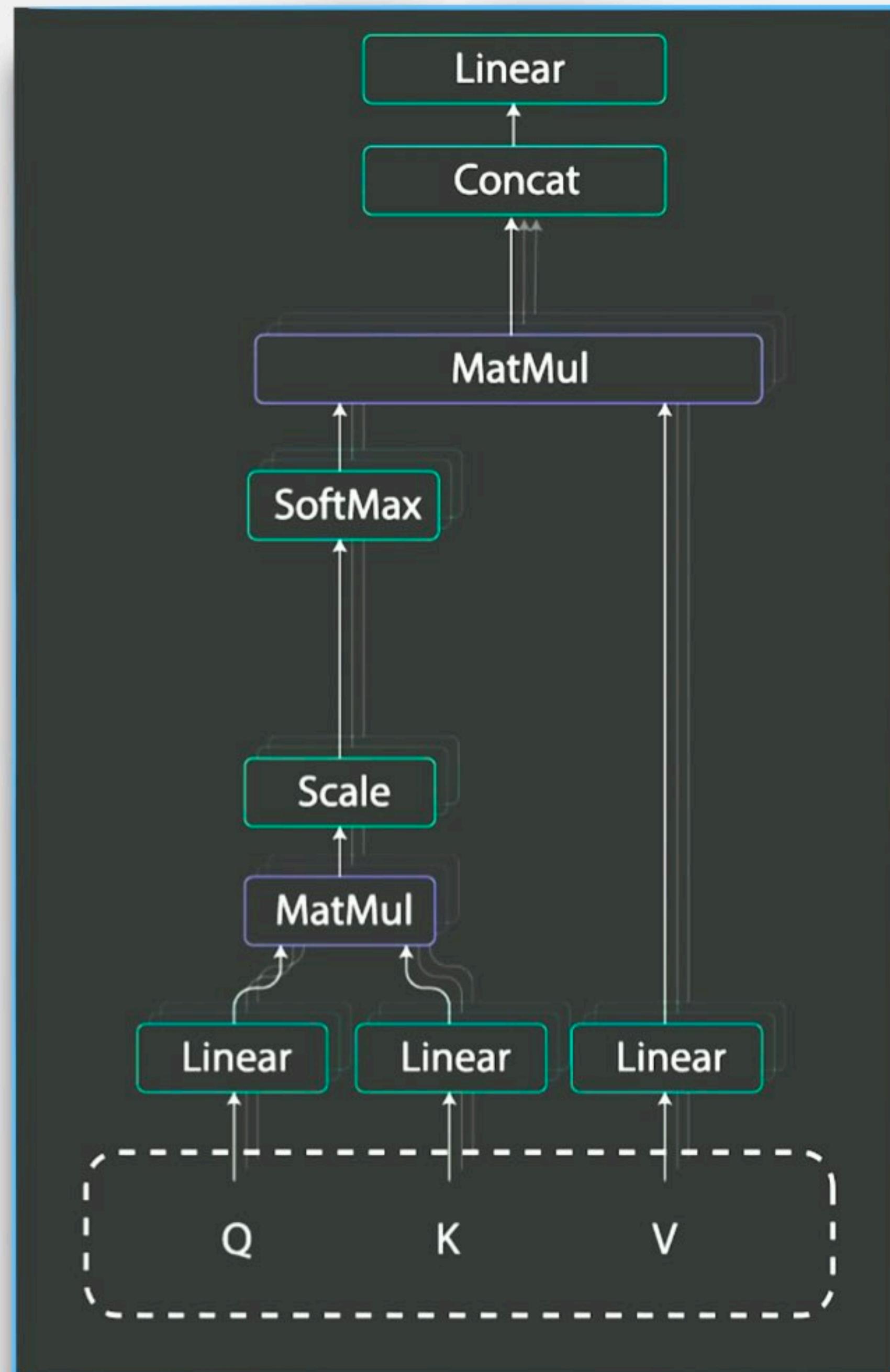


3. Multi-headed Attention

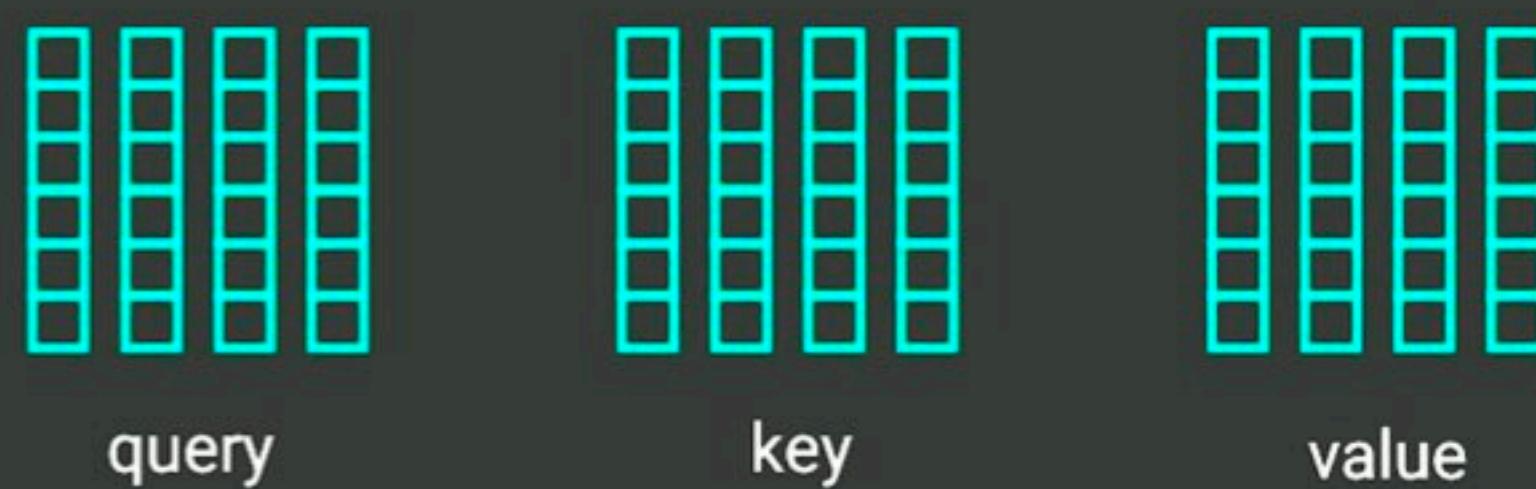
3.1. Self-Attention

attention weights value output

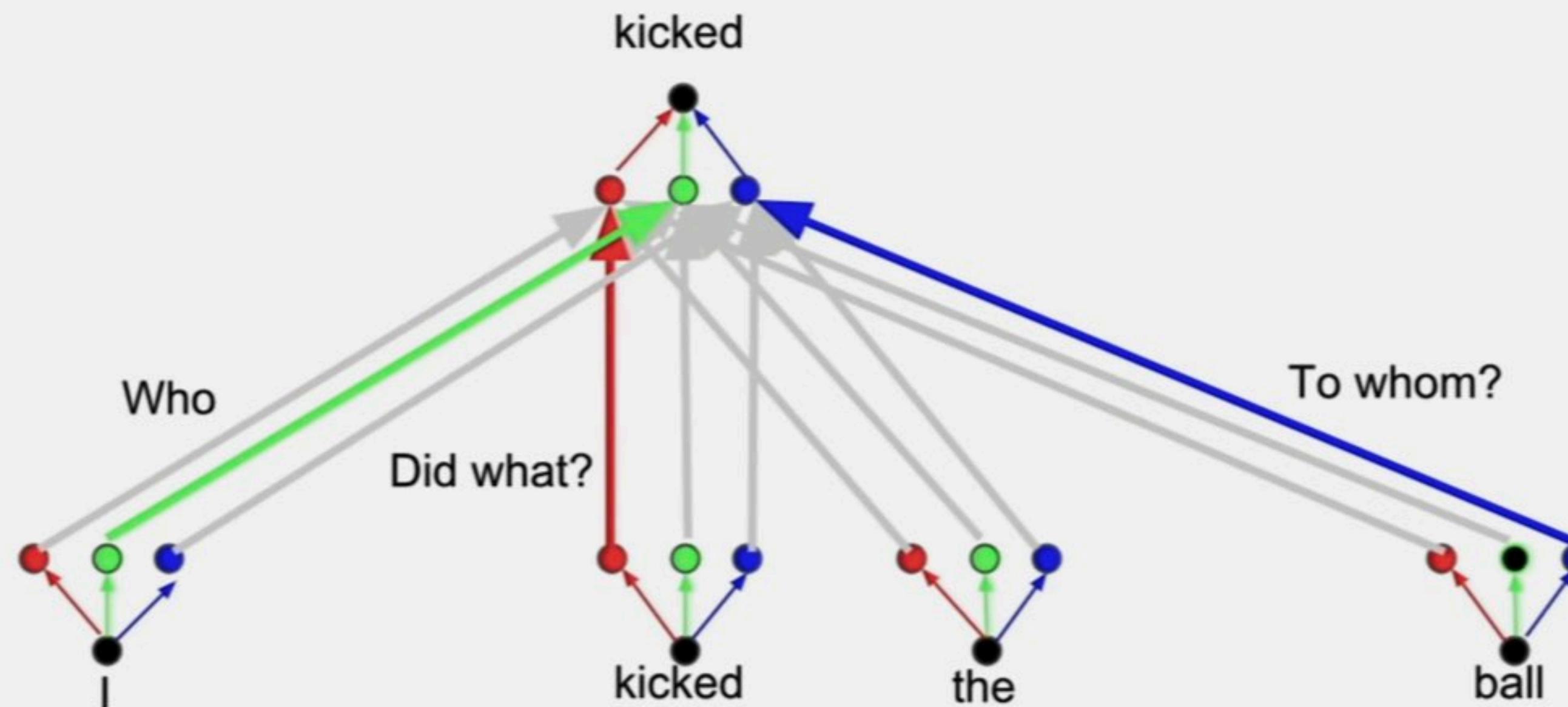
$$\begin{matrix} \text{attention weights} \\ \boxed{\text{grid}} \end{matrix} \times \begin{matrix} \text{value} \\ \boxed{\text{grid}} \end{matrix} = \begin{matrix} \text{output} \\ \boxed{\text{grid}} \end{matrix}$$

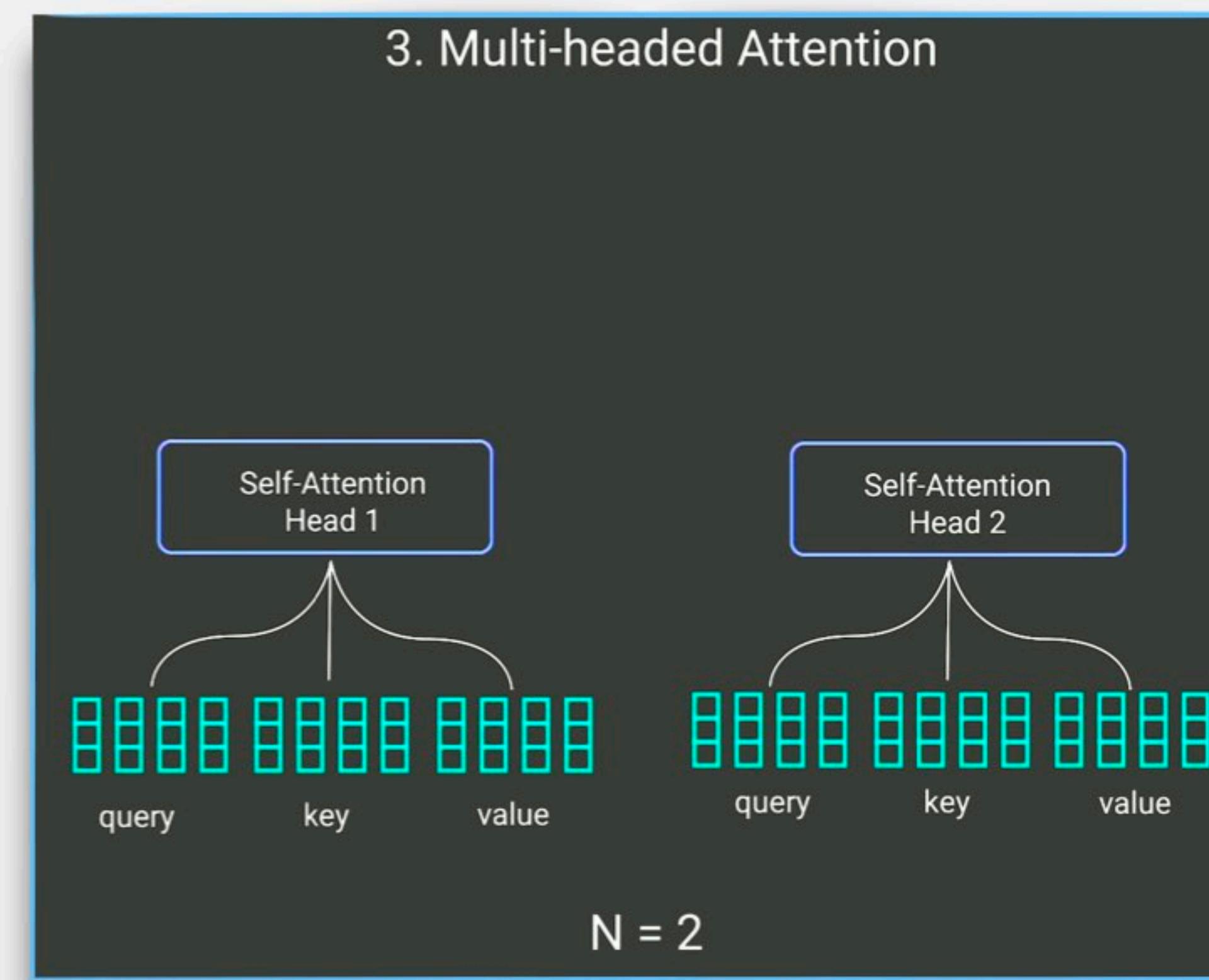
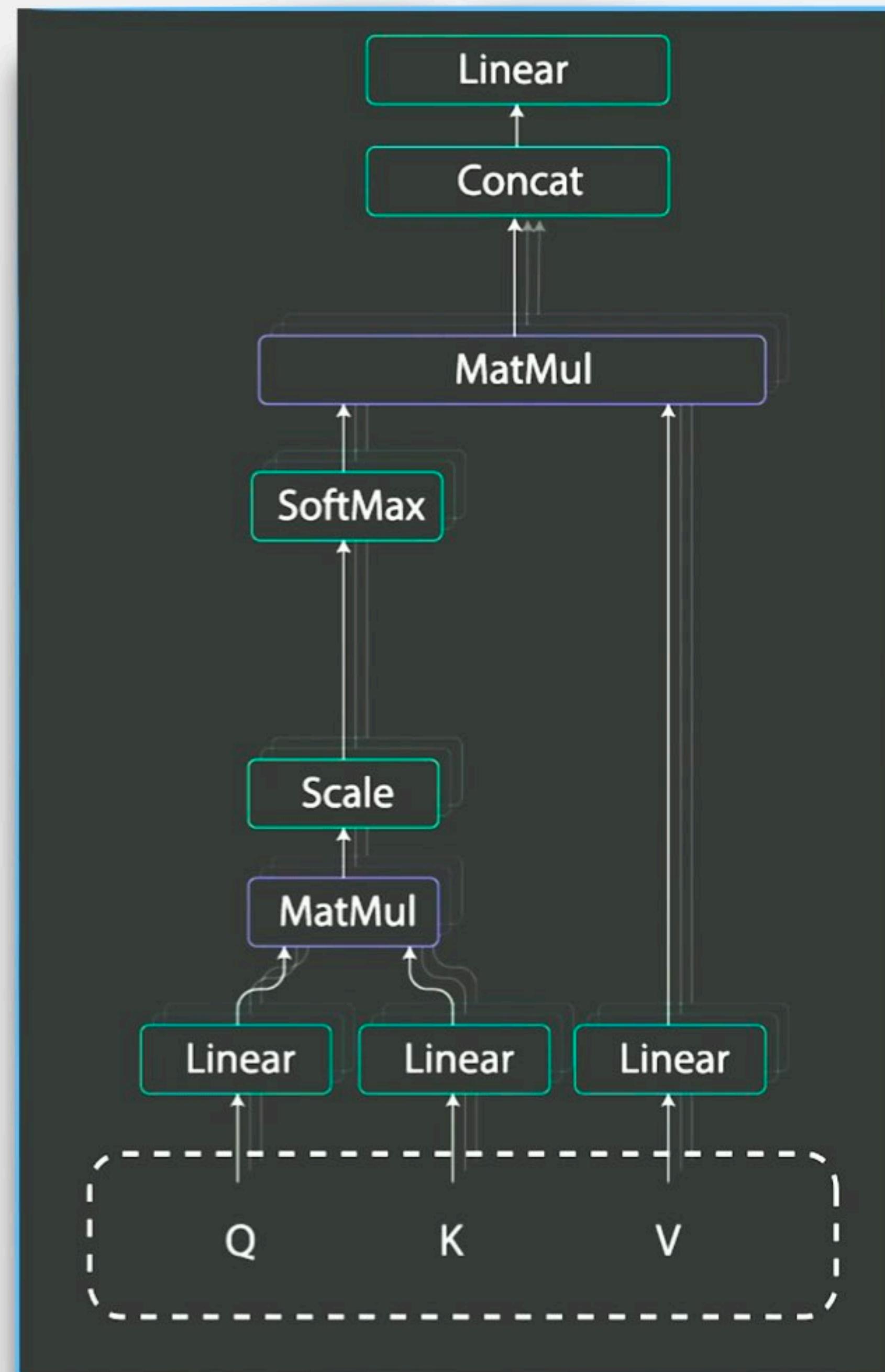


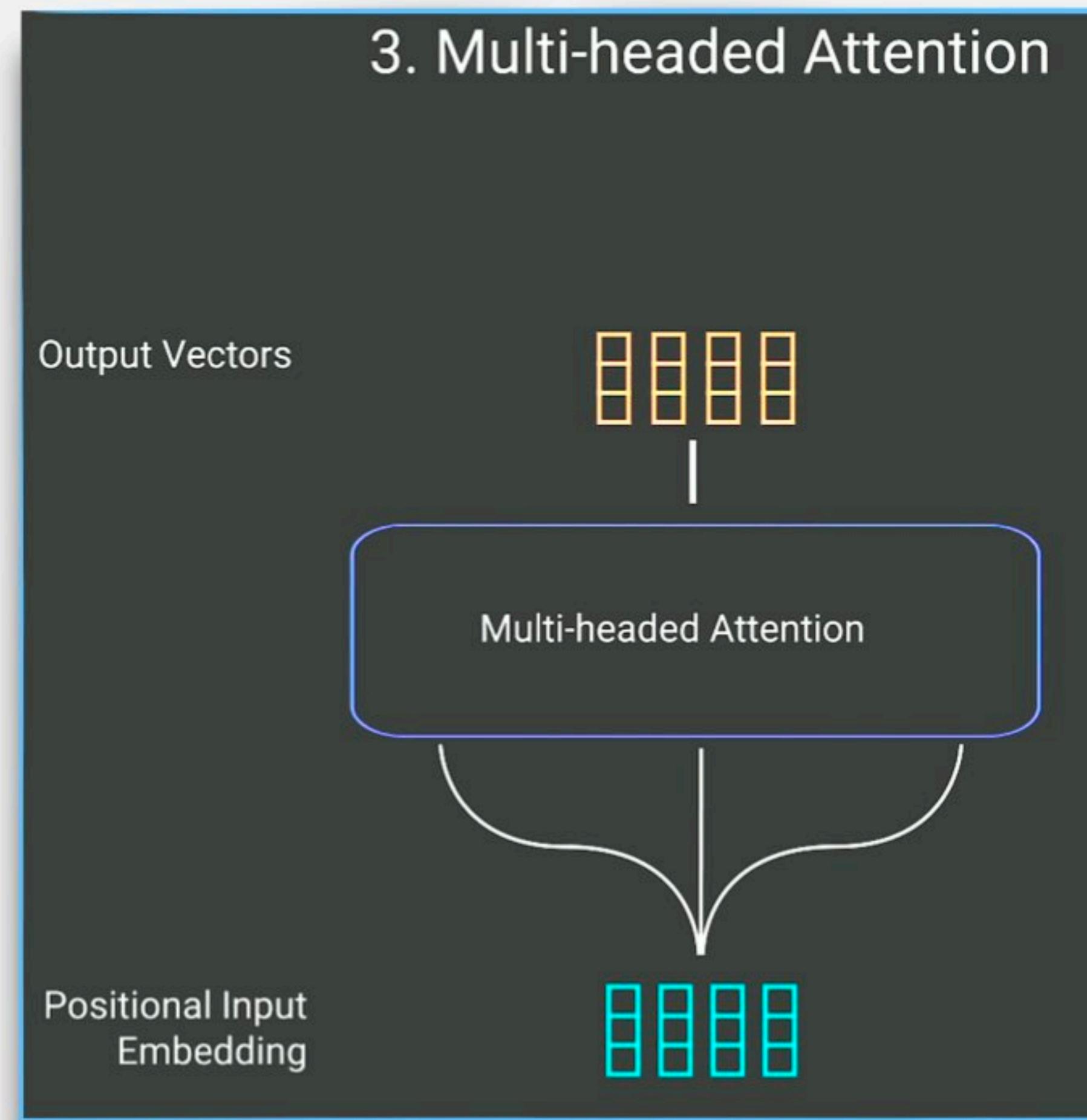
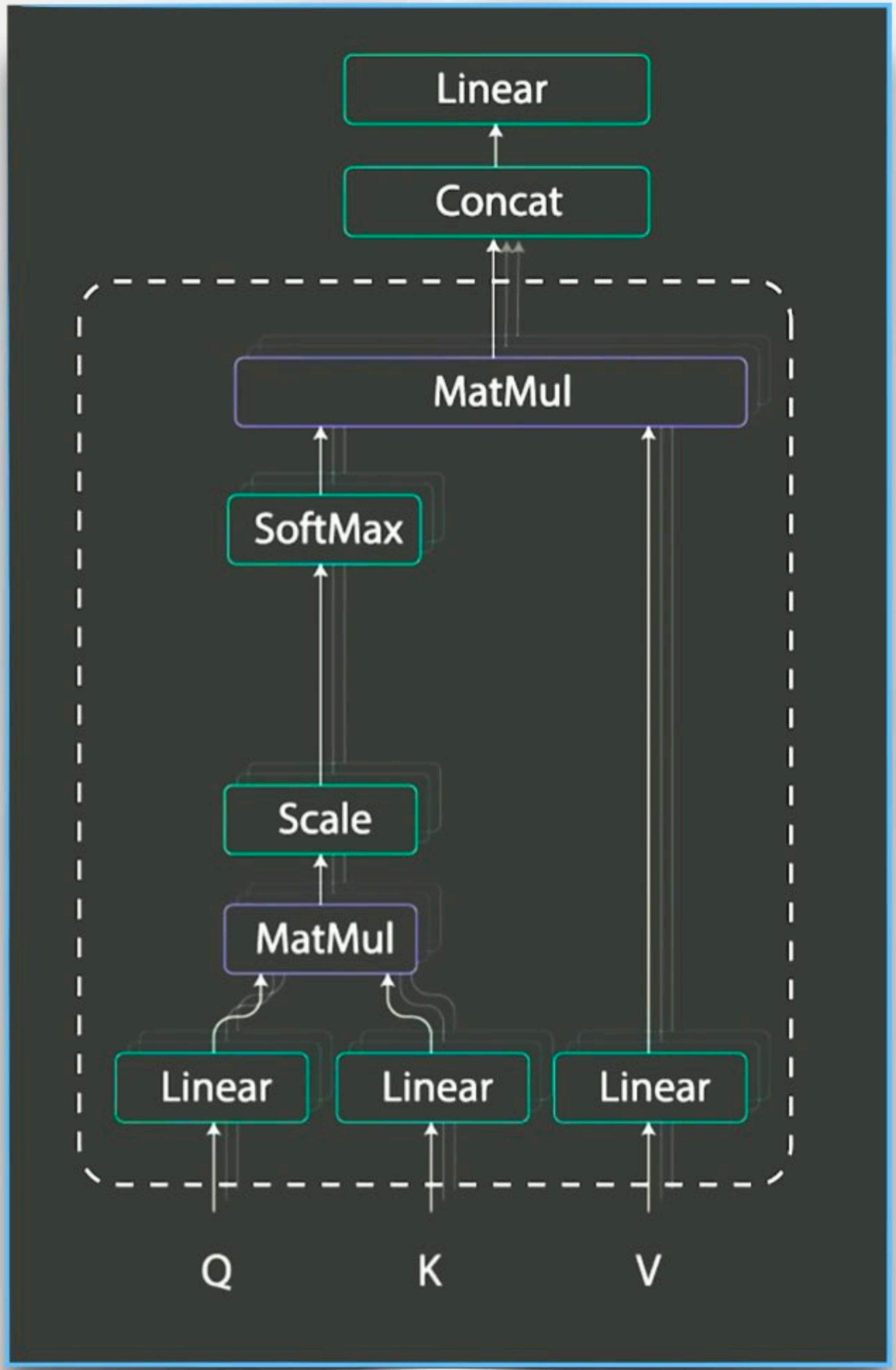
3. Multi-headed Attention

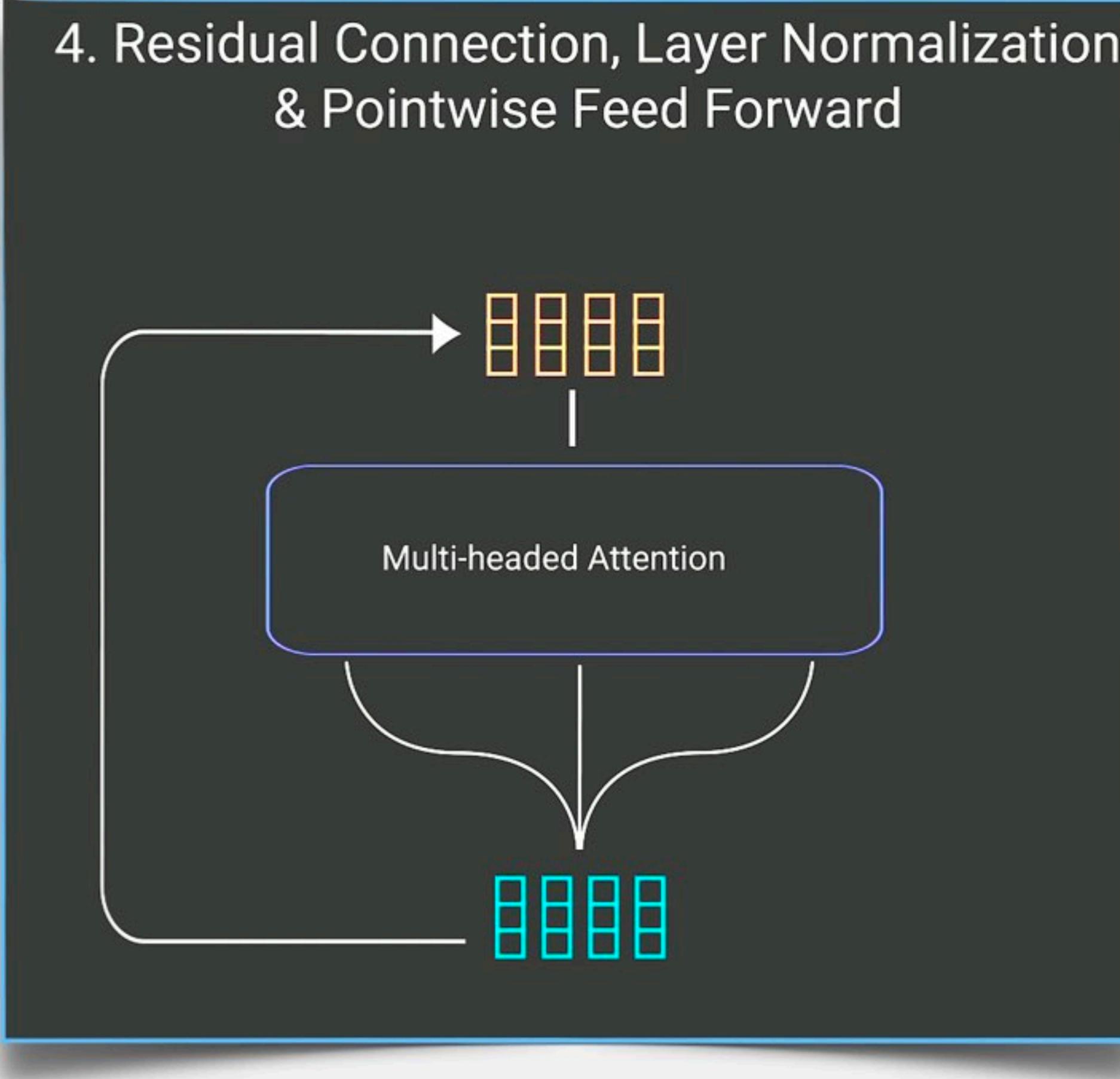
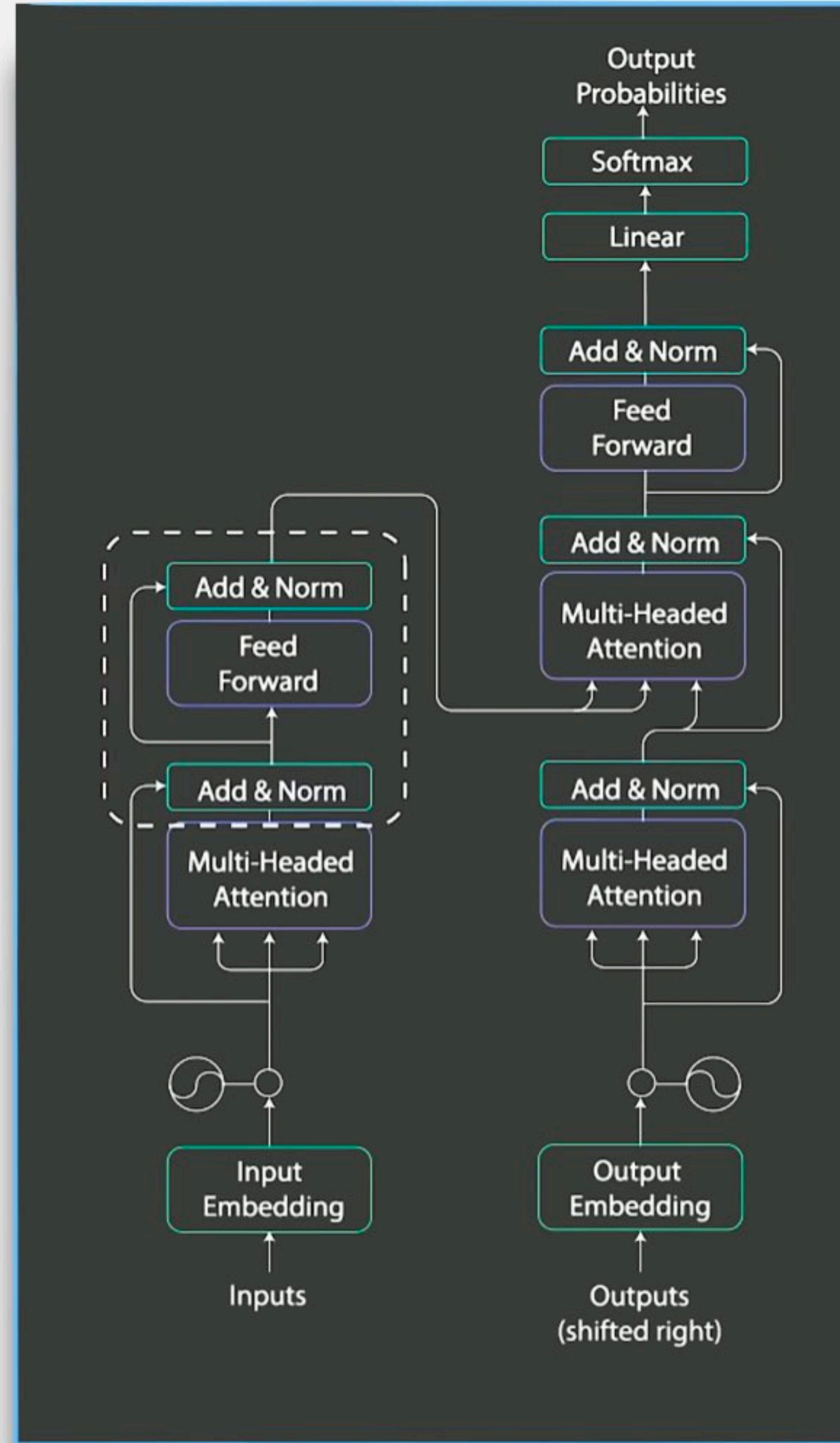


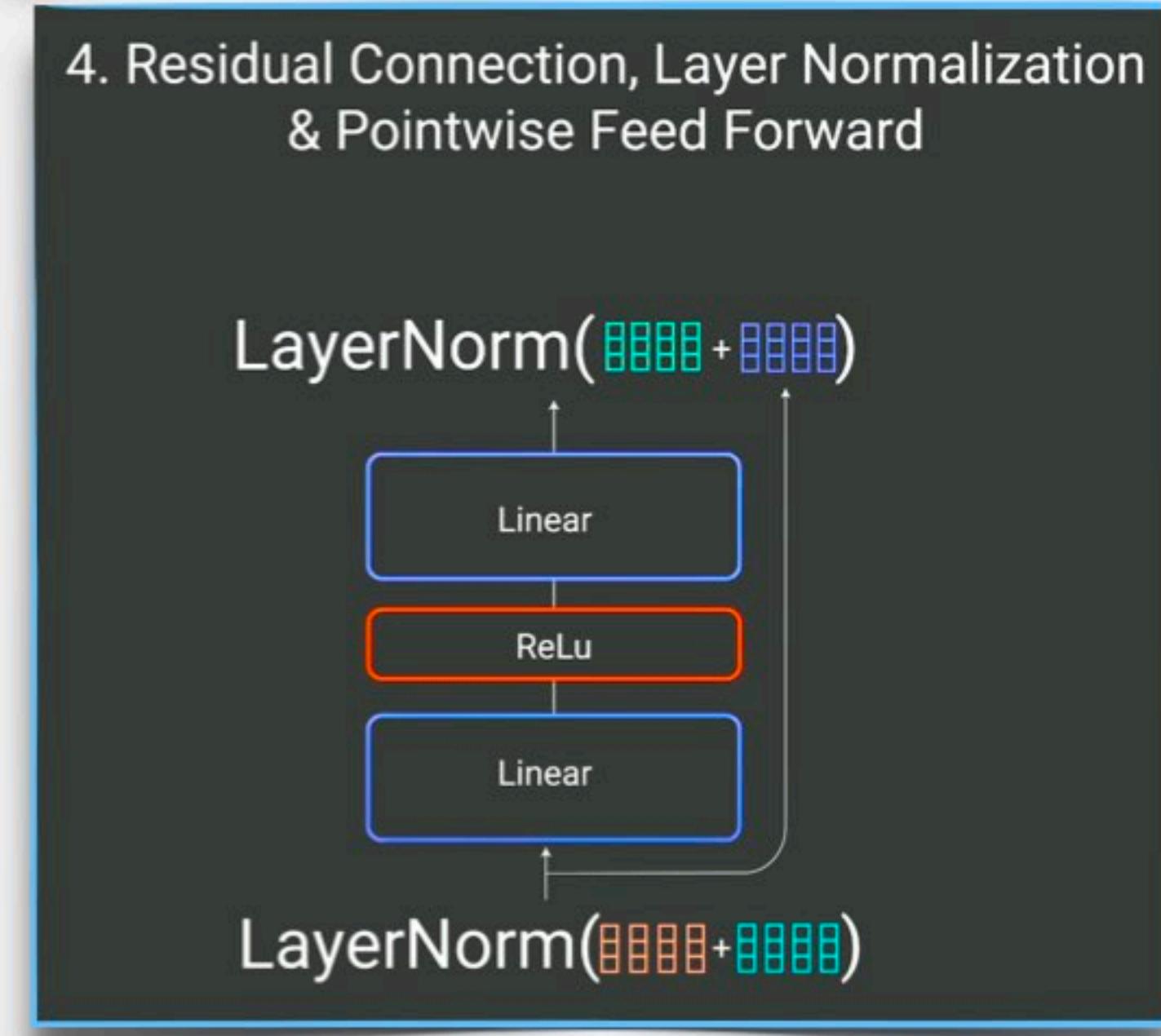
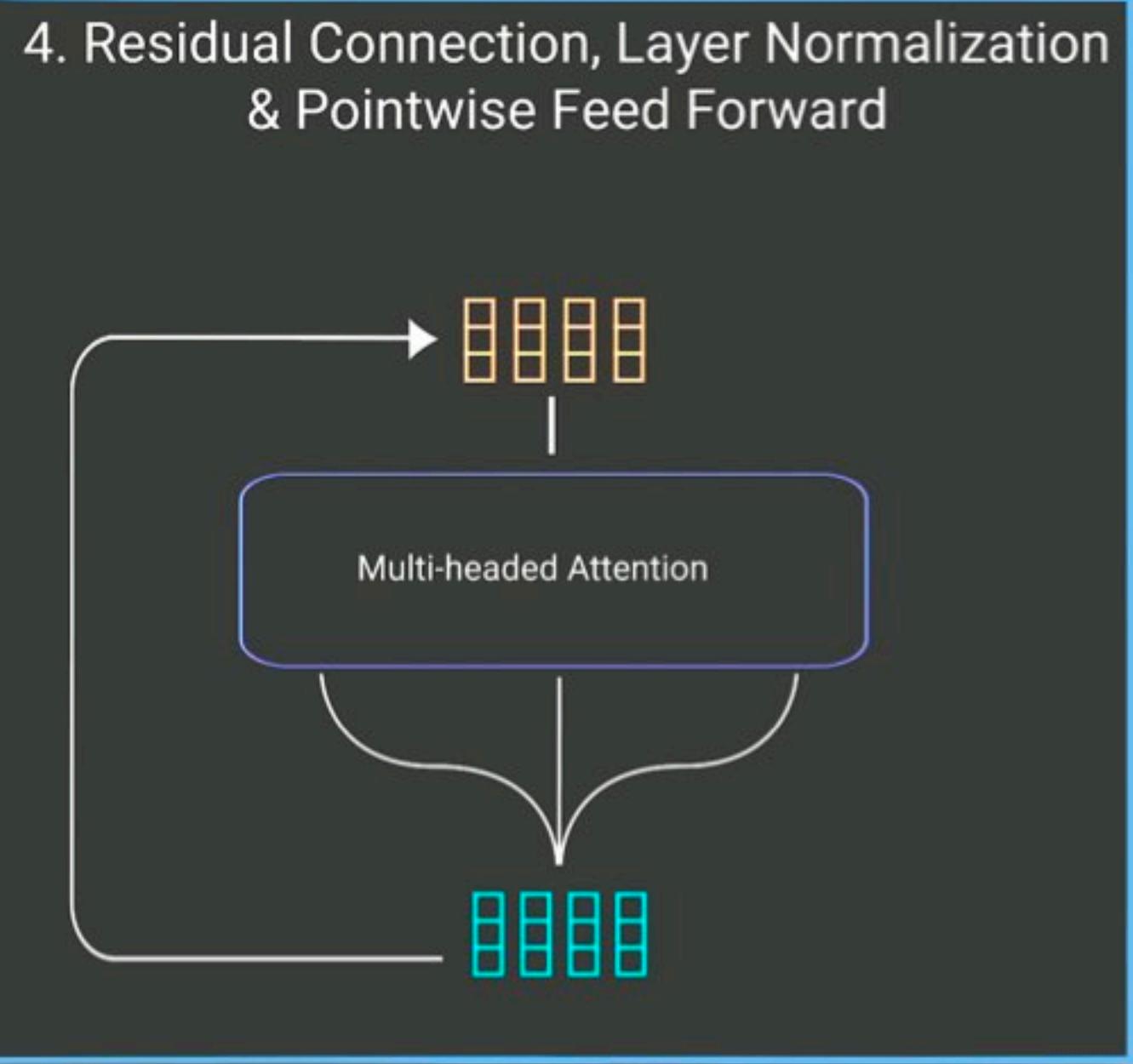
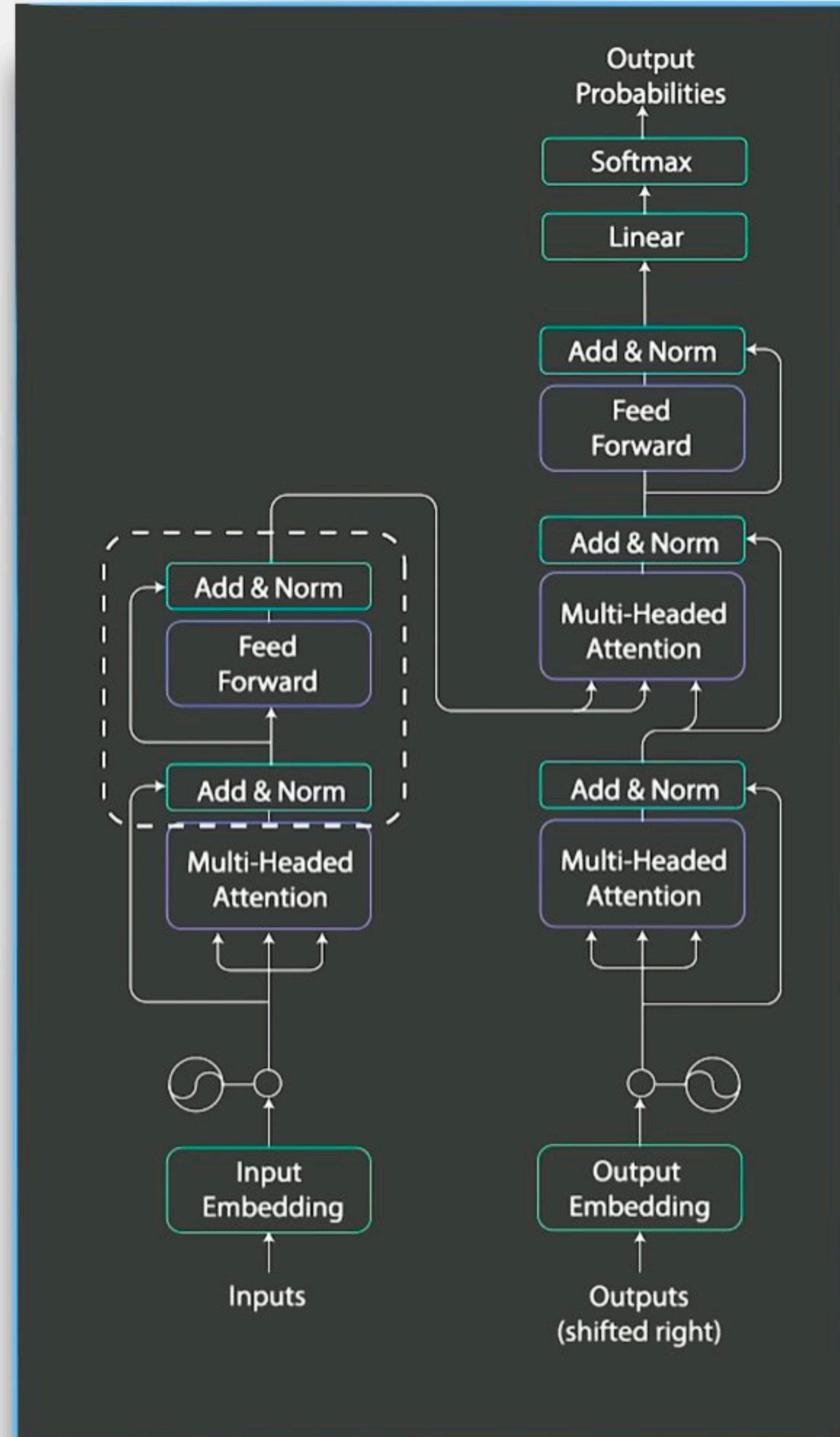
Attention Head-Who, What, and Whom

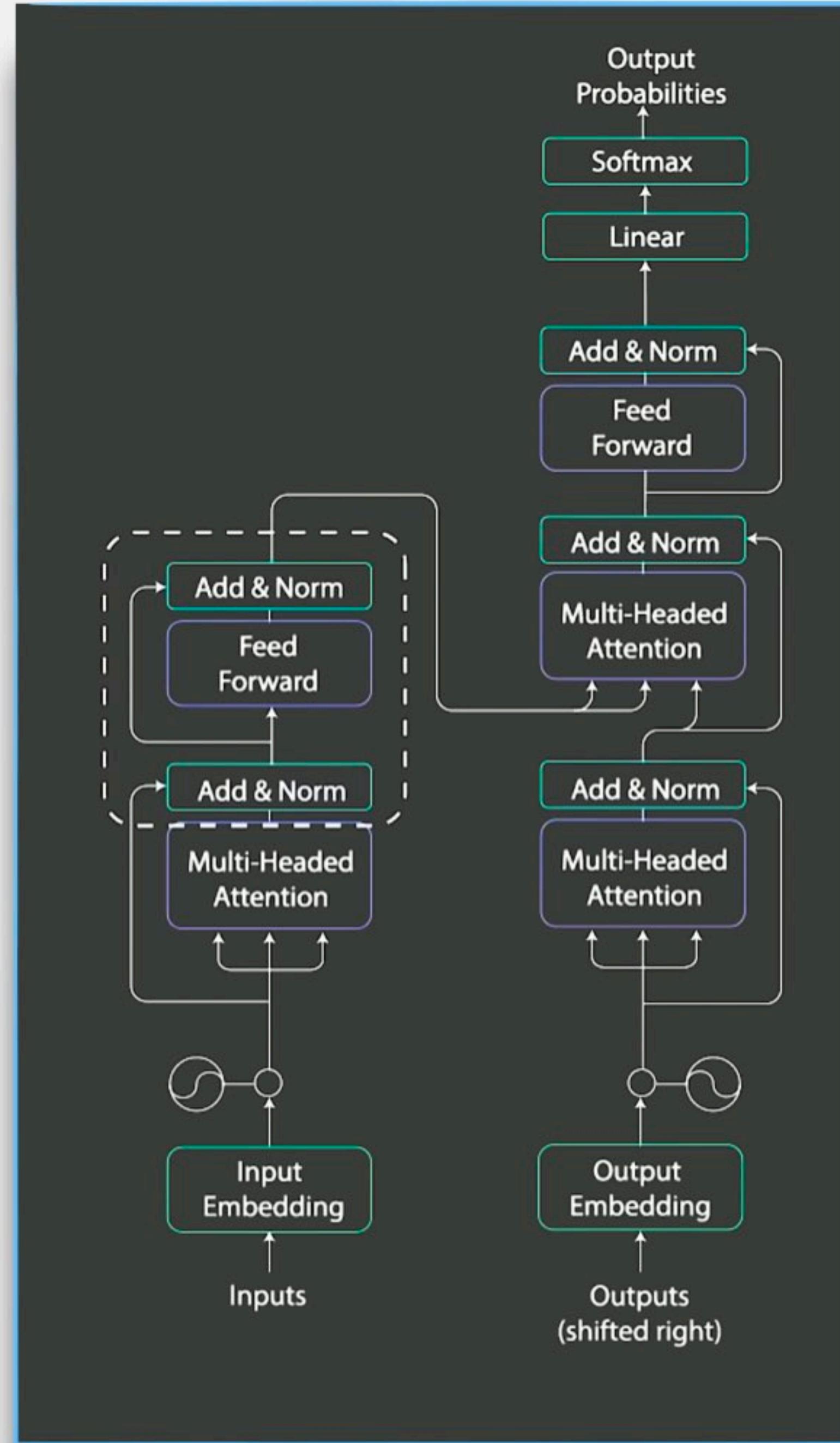








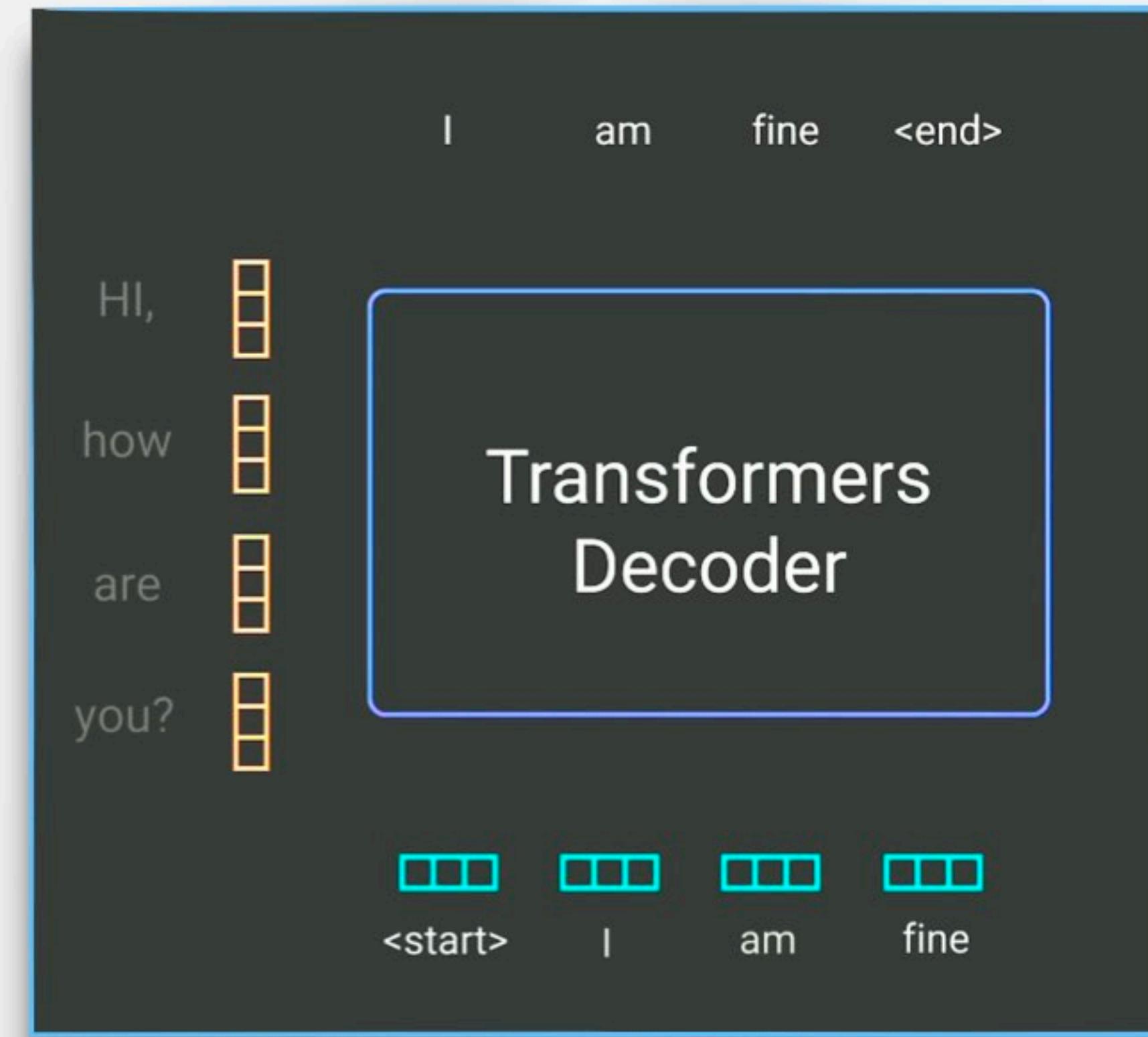
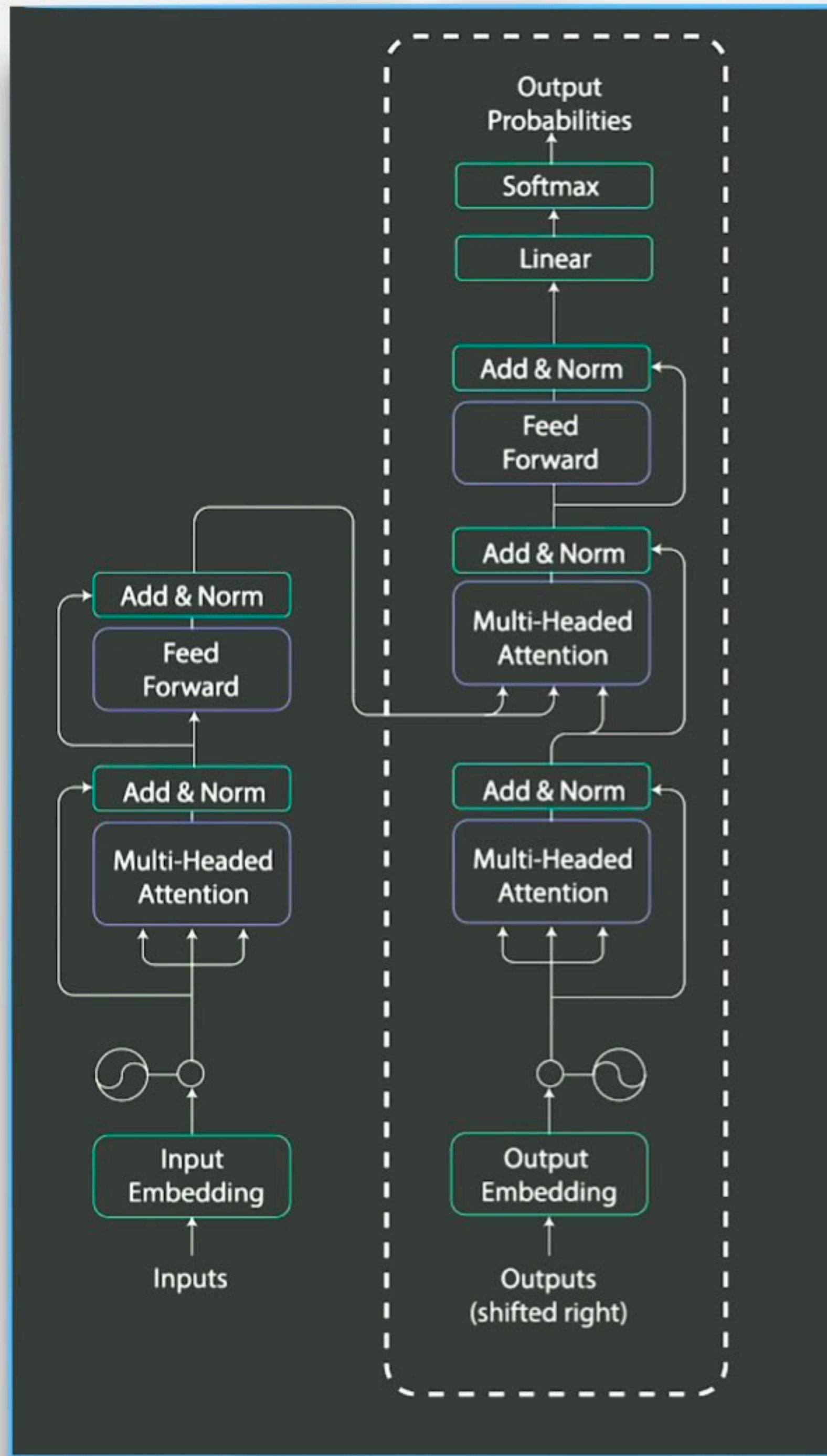




Transformer
Encoder

⋮
N times

Transformer
Encoder



6. Decoder Multi-Headed Attention 1



6. Decoder Multi-Headed Attention 1



6. Decoder Multi-Headed Attention 1

6.1. Look-Ahead Mask

Scaled Scores

| | | | |
|-----|-----|-----|-----|
| 0.7 | 0.1 | 0.1 | 0.1 |
| 0.1 | 0.6 | 0.2 | 0.1 |
| 0.1 | 0.3 | 0.6 | 0.1 |
| 0.1 | 0.3 | 0.3 | 0.3 |

Look-Ahead Mask

| | | | |
|---|------|------|------|
| 0 | -inf | -inf | -inf |
| 0 | 0 | -inf | -inf |
| 0 | 0 | 0 | -inf |
| 0 | 0 | 0 | 0 |

+

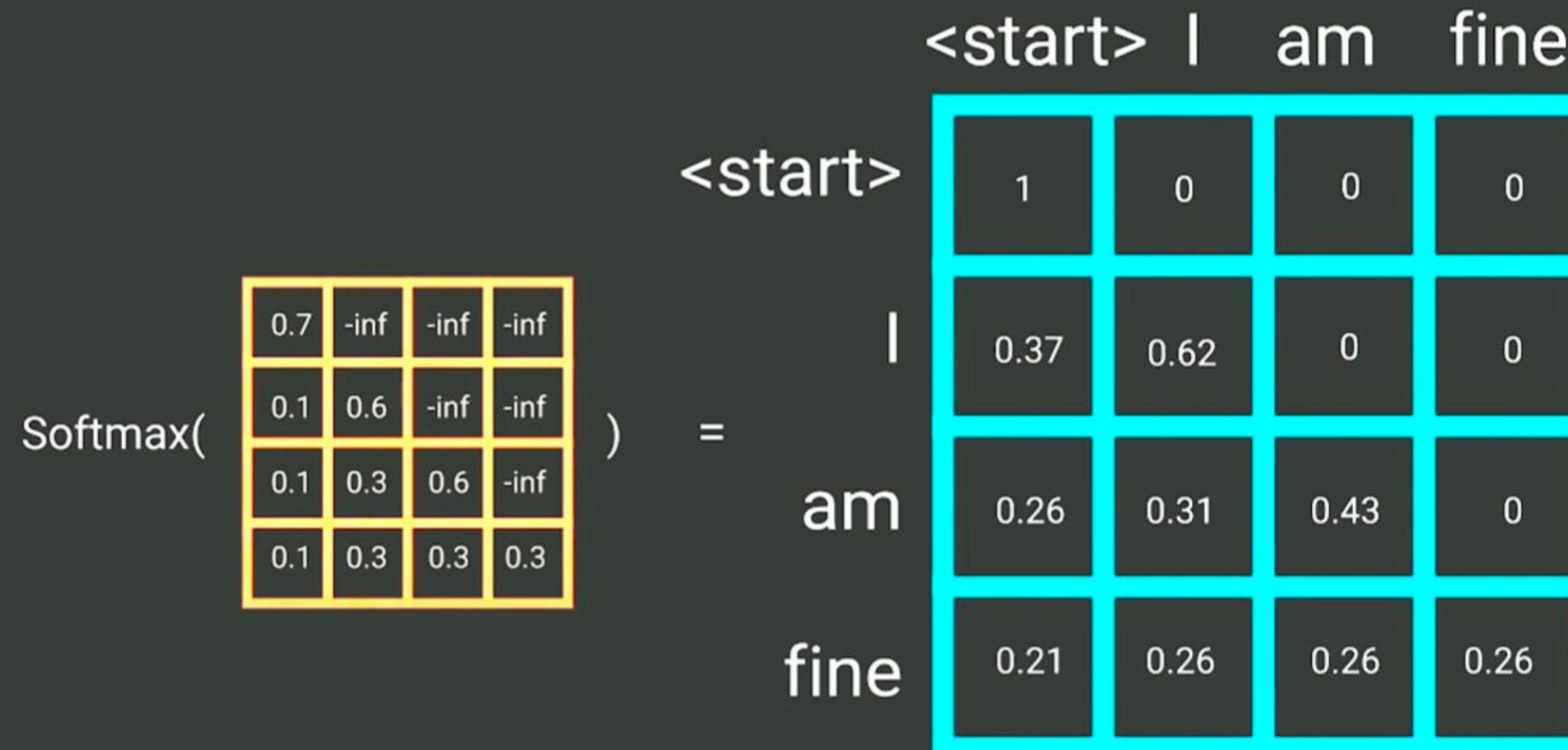
Masked Scores

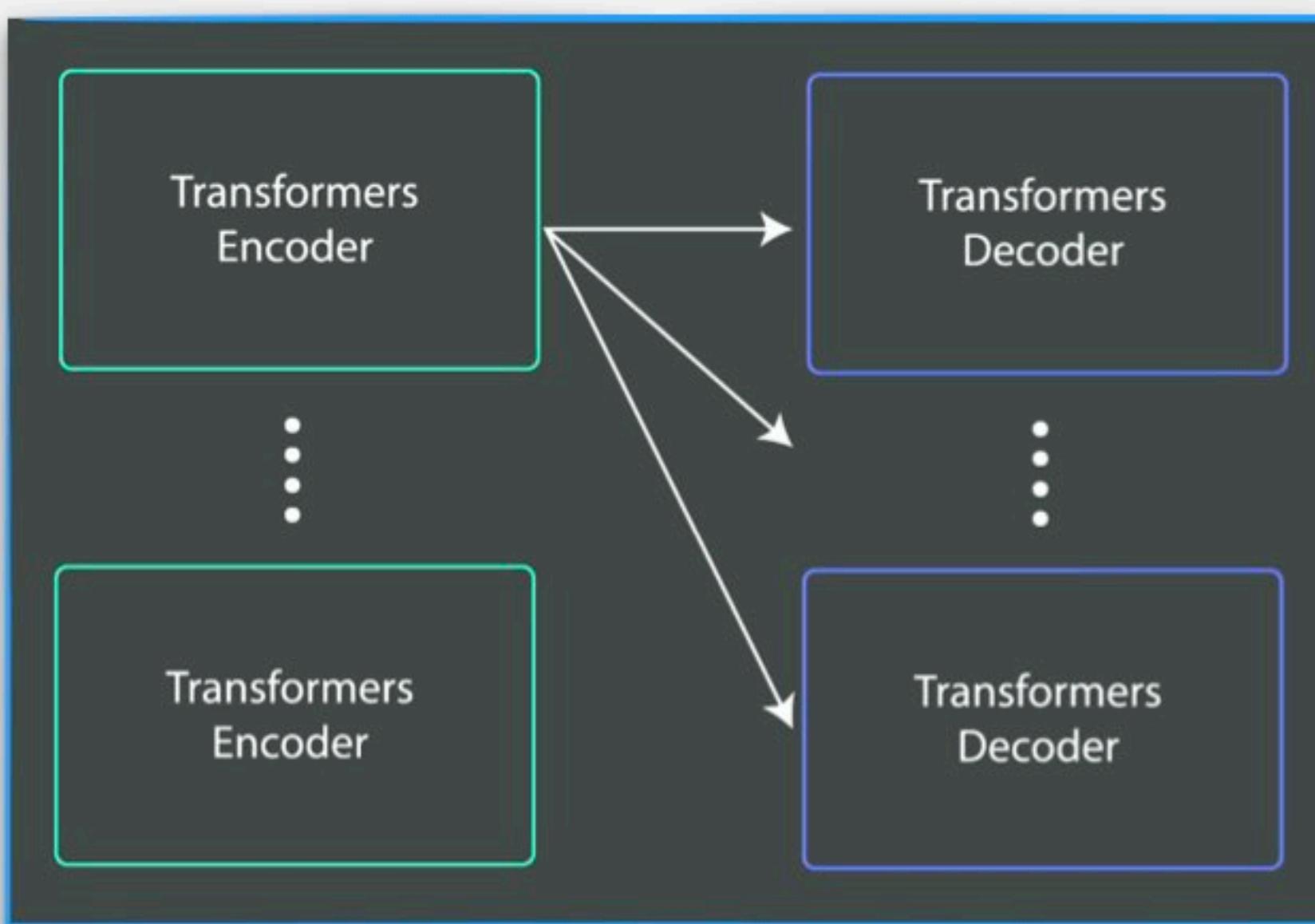
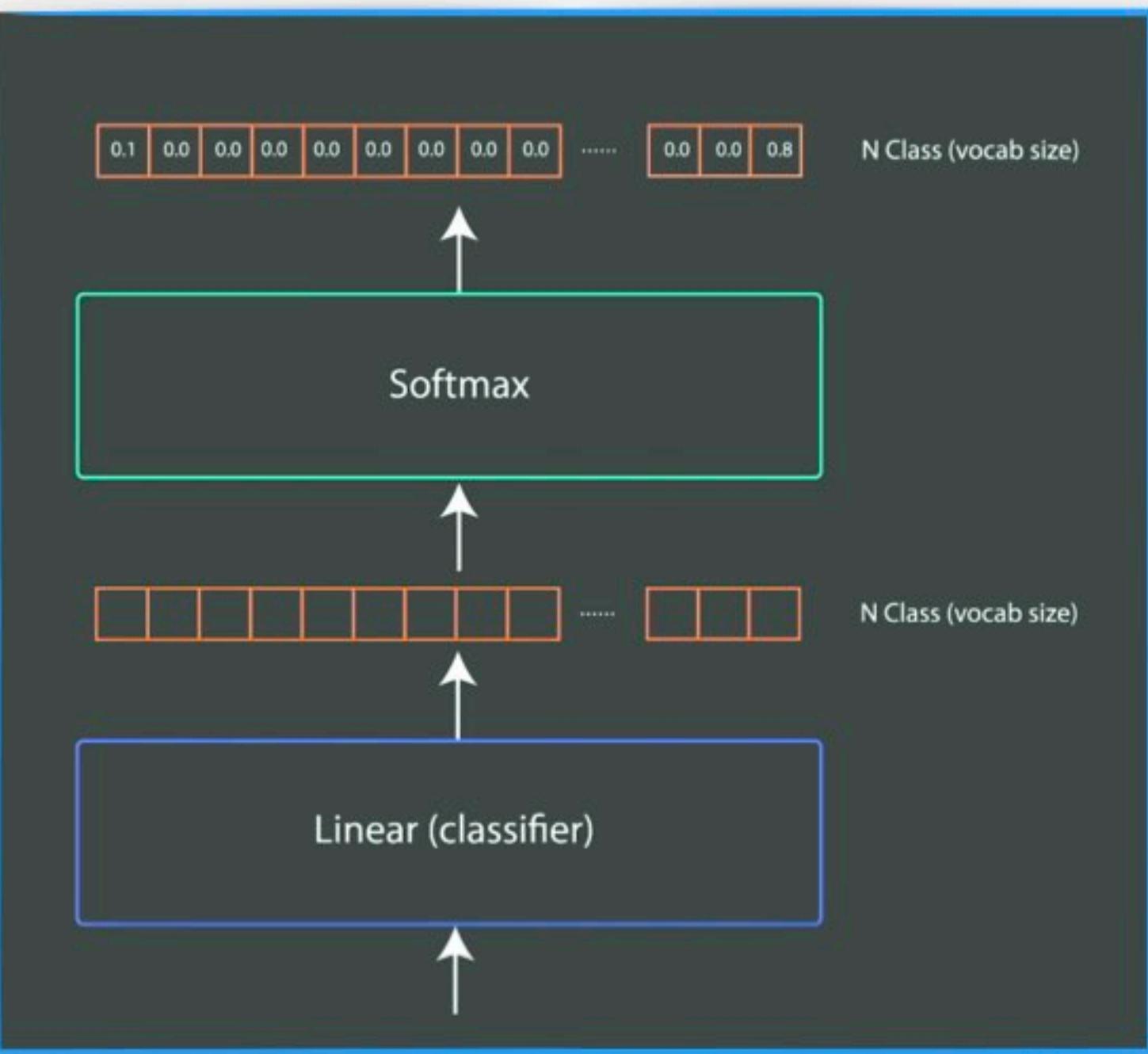
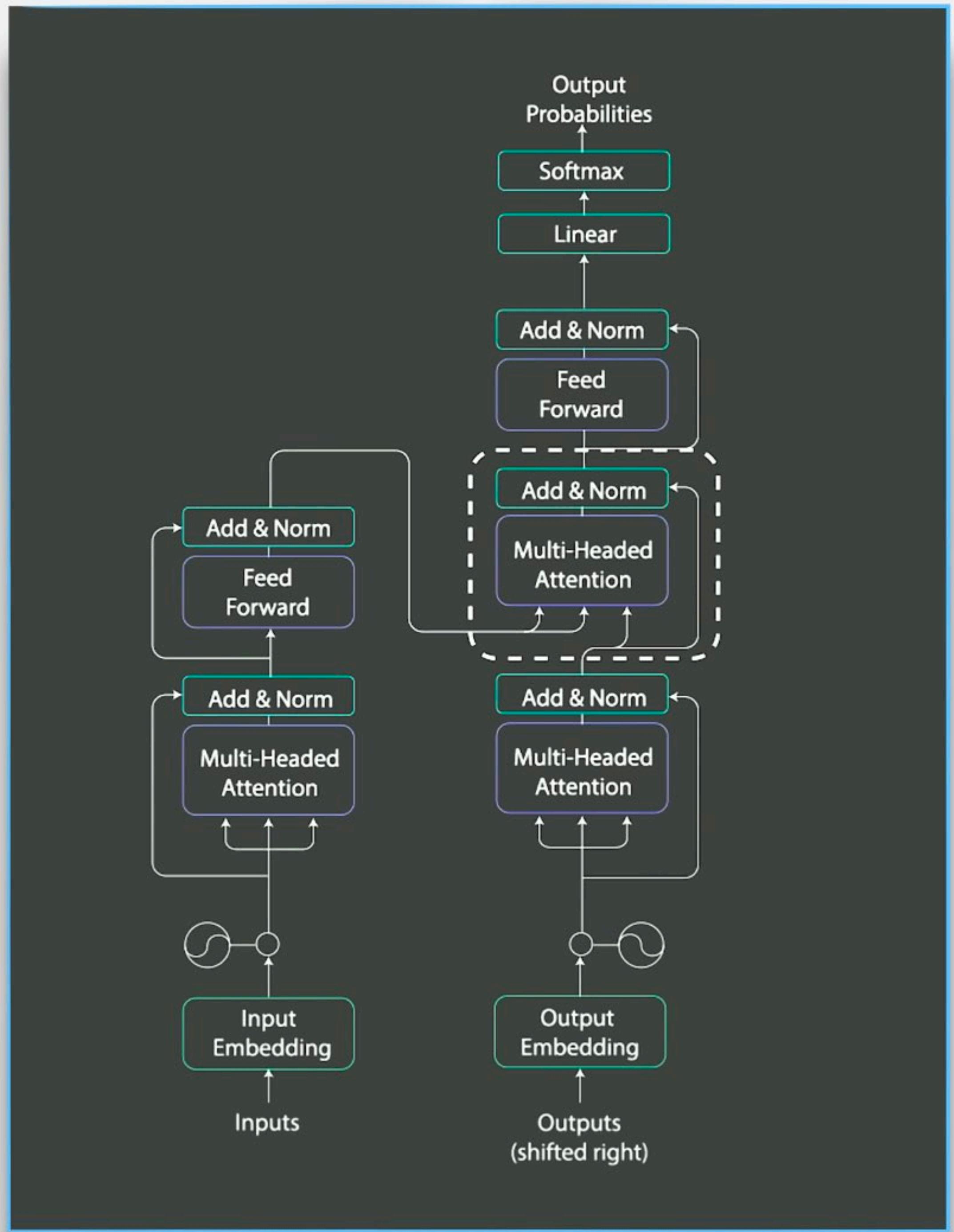
| | | | |
|-----|------|------|------|
| 0.7 | -inf | -inf | -inf |
| 0.1 | 0.6 | -inf | -inf |
| 0.1 | 0.3 | 0.6 | -inf |
| 0.1 | 0.3 | 0.3 | 0.3 |

=

6. Decoder Multi-Headed Attention 1

6.1. Look-Ahead Mask





Positional Encoding

what we see:

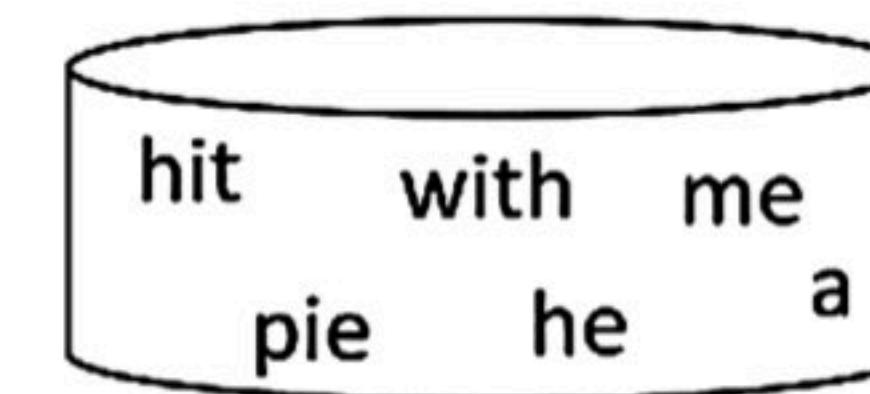
he hit me with a pie

what naïve self-attention sees:

a pie hit me with he

a hit with me he pie

he pie me with a hit

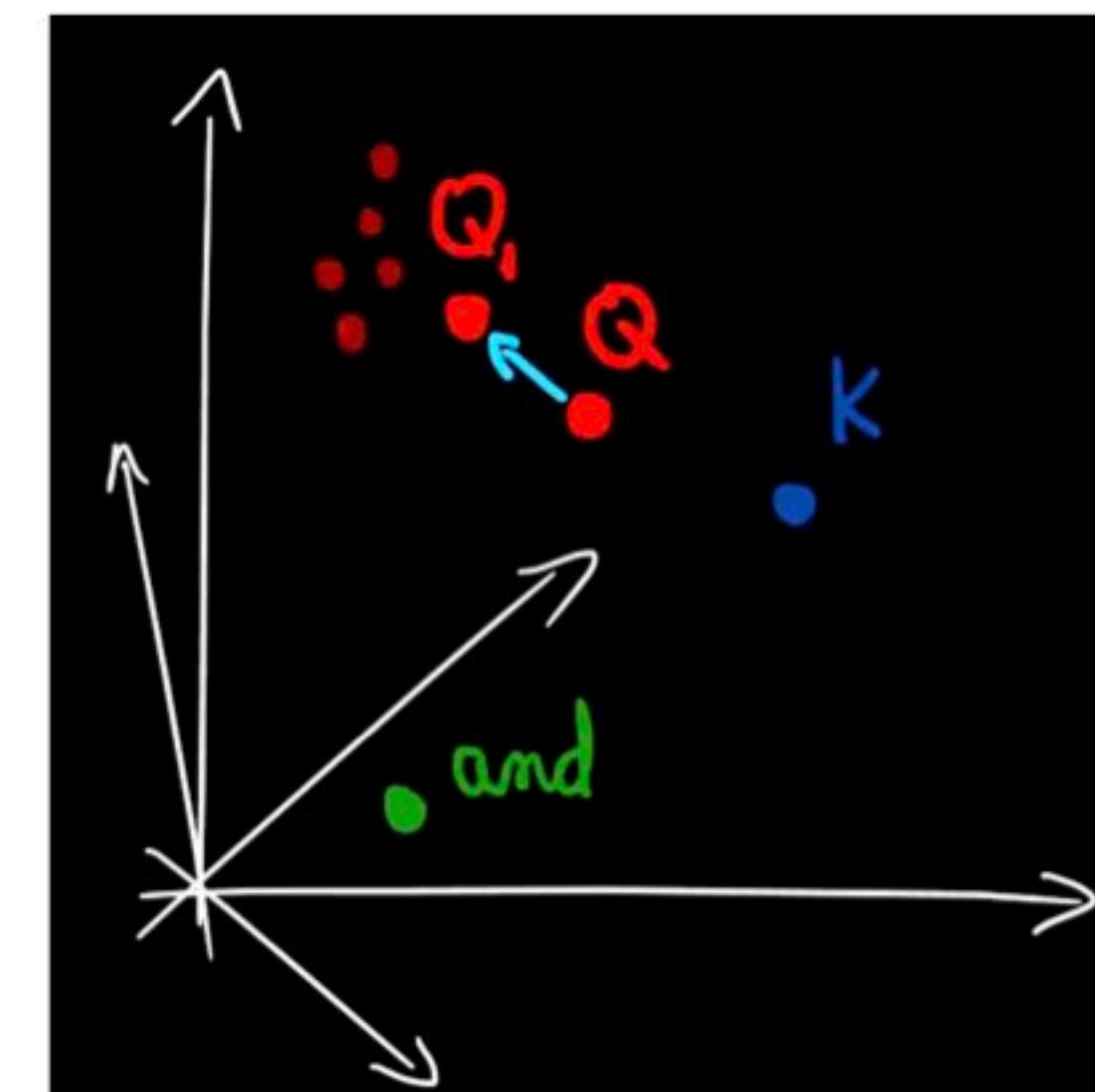
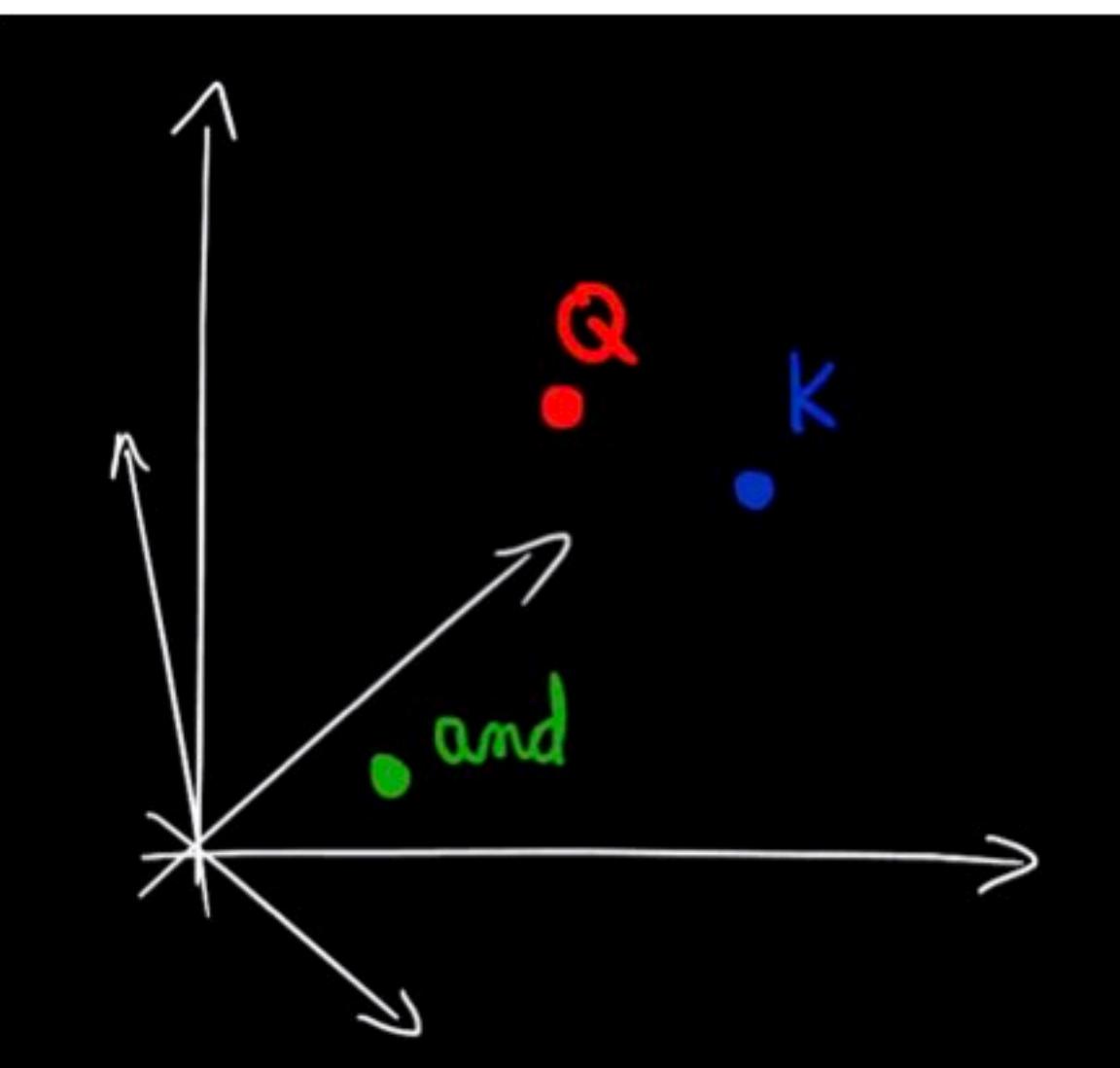


most alternative orderings are nonsense, but some change the meaning
in general the position of words in a sentence carries information!

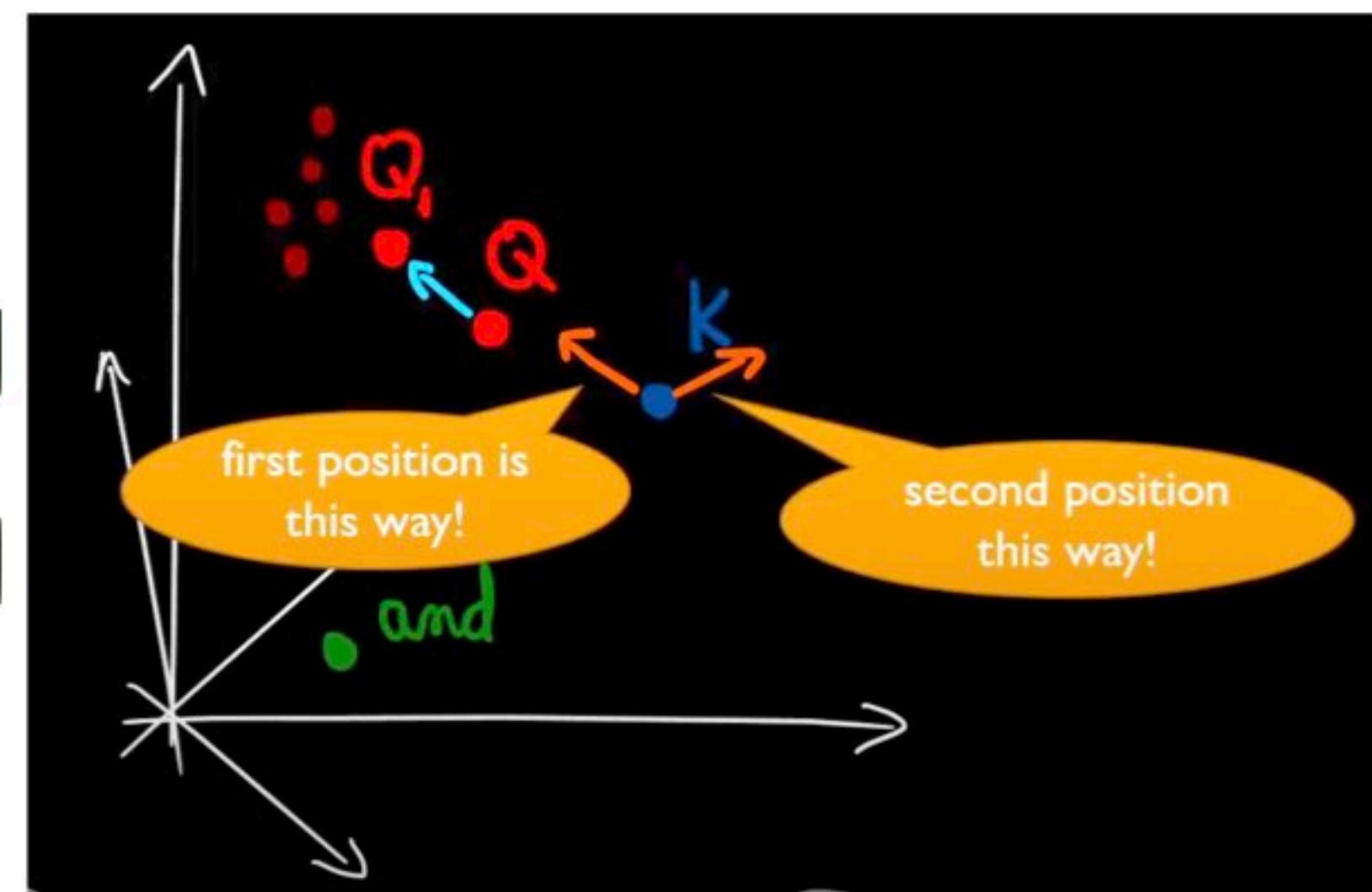
Idea: add some information to the representation at the beginning that indicates where it is in the sequence!

Positional Encoding

Naïve positional encoding



| | | | | |
|------|------|------|------|------|
| 0.33 | 0.71 | 0.91 | 0.23 | 0.15 |
| 0.01 | 0.01 | 0.01 | 0.01 | 0.01 |



I walk my dog every day



every single day I walk my dog



The fact that "my dog" is right after "I walk" is
the important part, not its absolute position

Positional Encoding

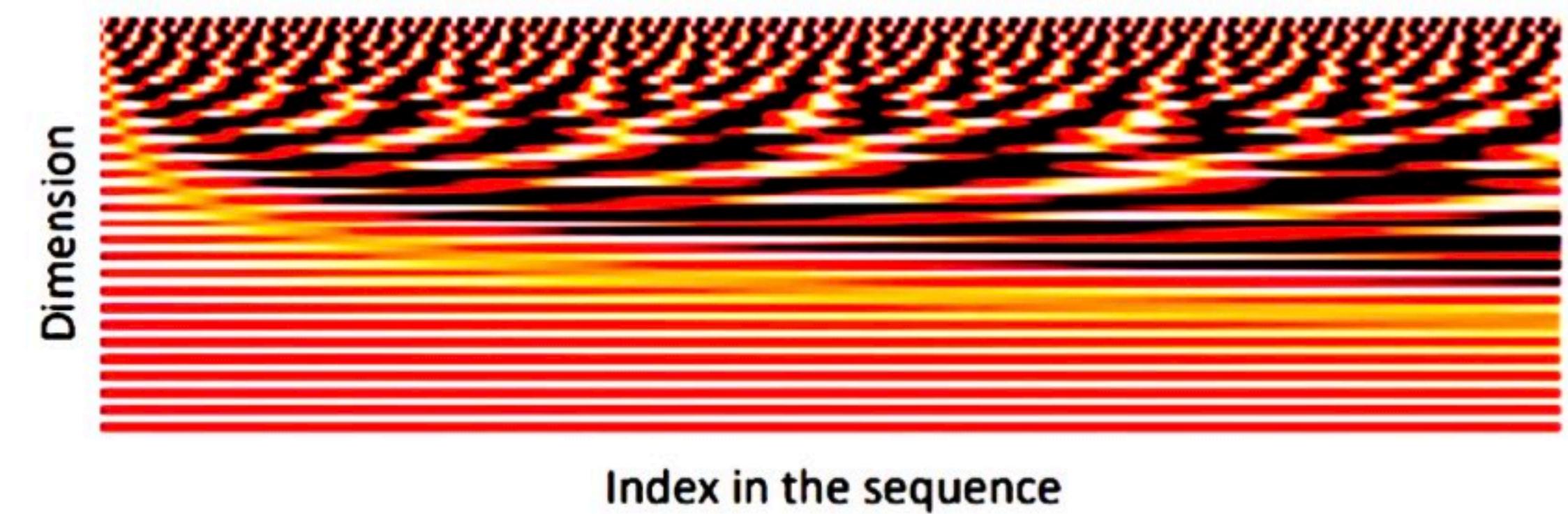
- Positional embeddings should be same irrespective of sequence length or content
- Shift in the embeddings space should not be large
- Ability to estimate distance between two tokens
- Deterministic based on some underlying rule

Positional Encoding

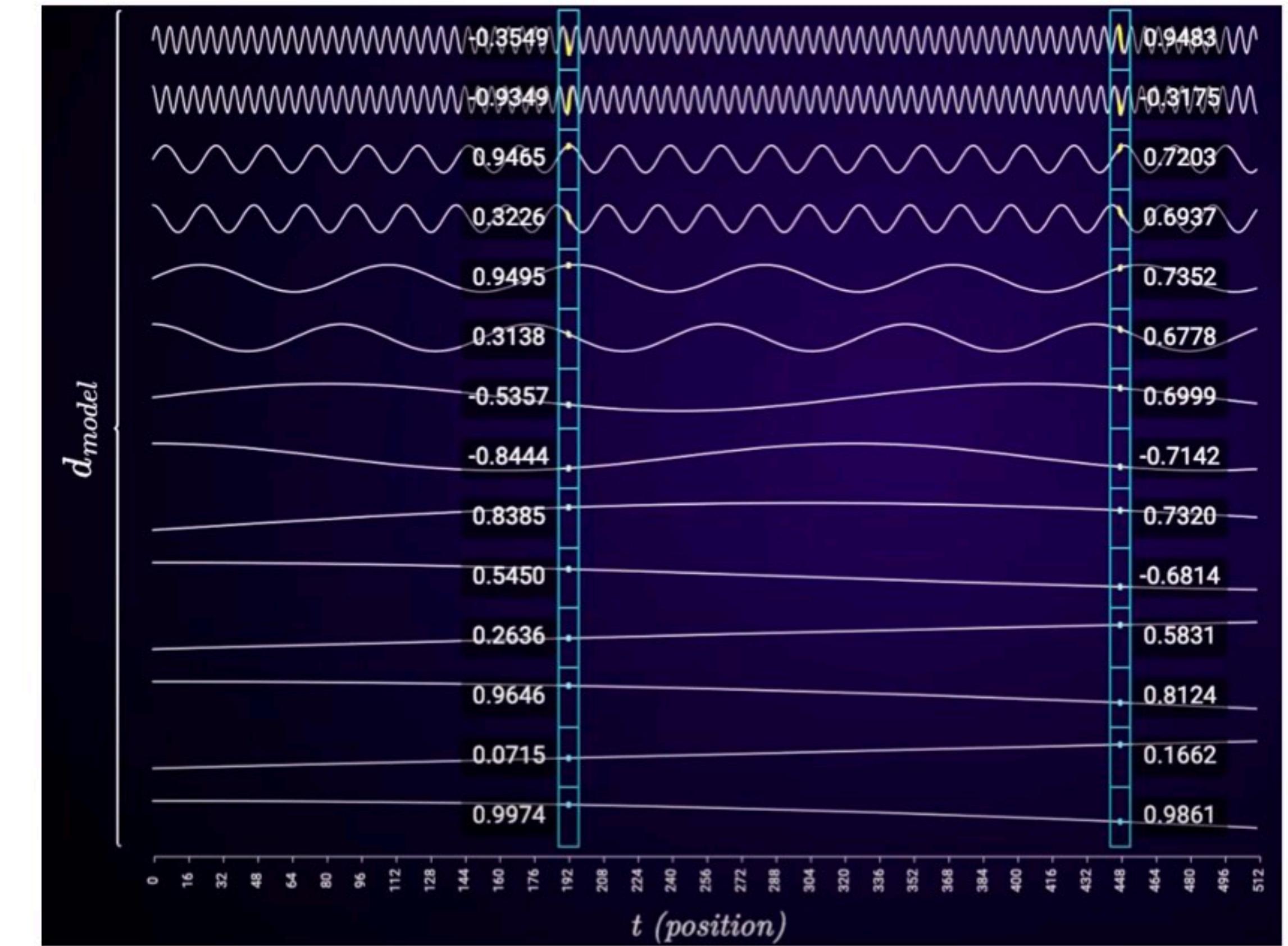
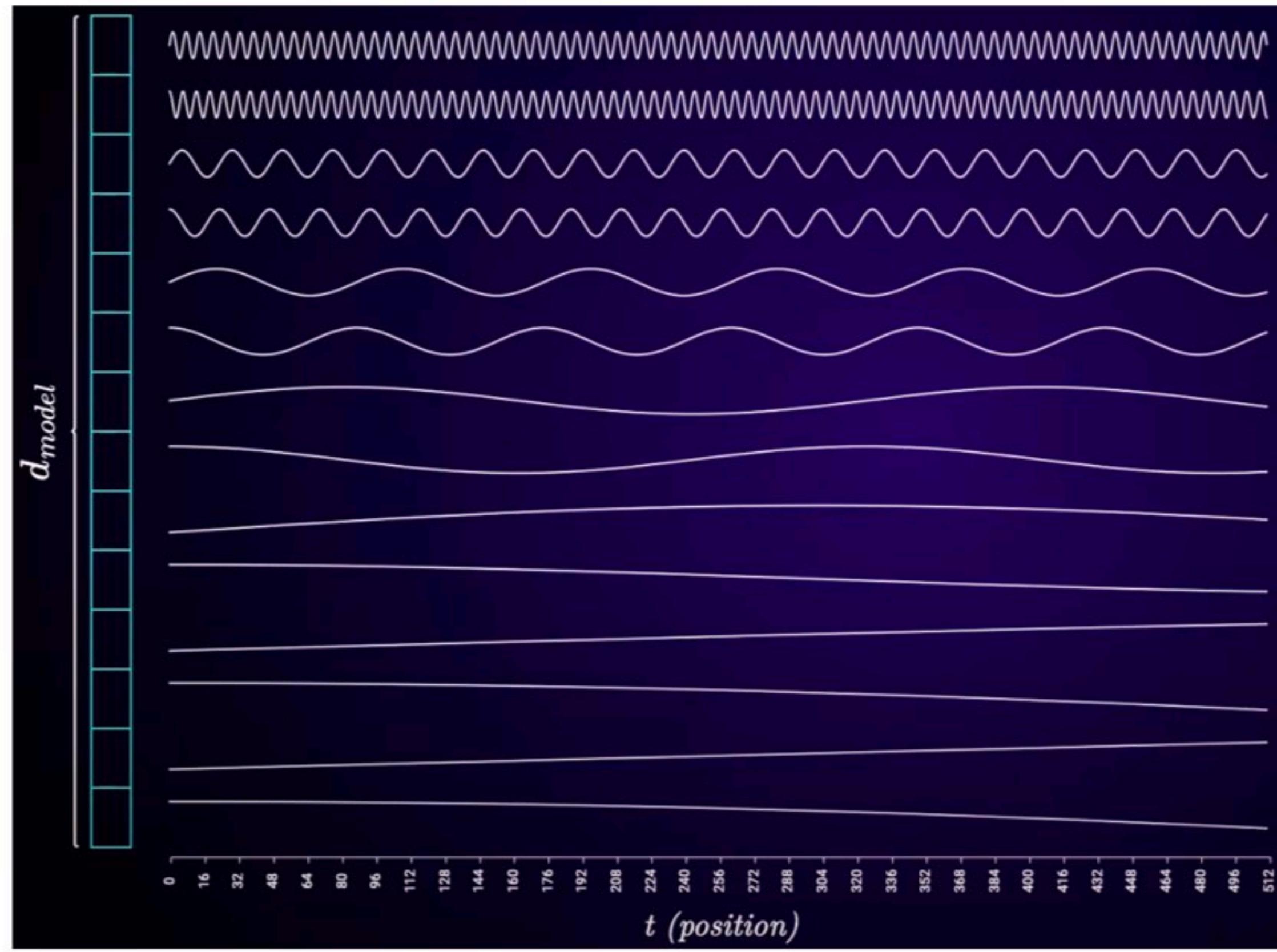
Idea: what if we use **frequency-based** representations?

$$p_t = \begin{bmatrix} \sin(t/10000^{2*1/d}) \\ \cos(t/10000^{2*1/d}) \\ \sin(t/10000^{2*2/d}) \\ \cos(t/10000^{2*2/d}) \\ \dots \\ \sin(t/10000^{2*\frac{d}{2}/d}) \\ \cos(t/10000^{2*\frac{d}{2}/d}) \end{bmatrix}$$

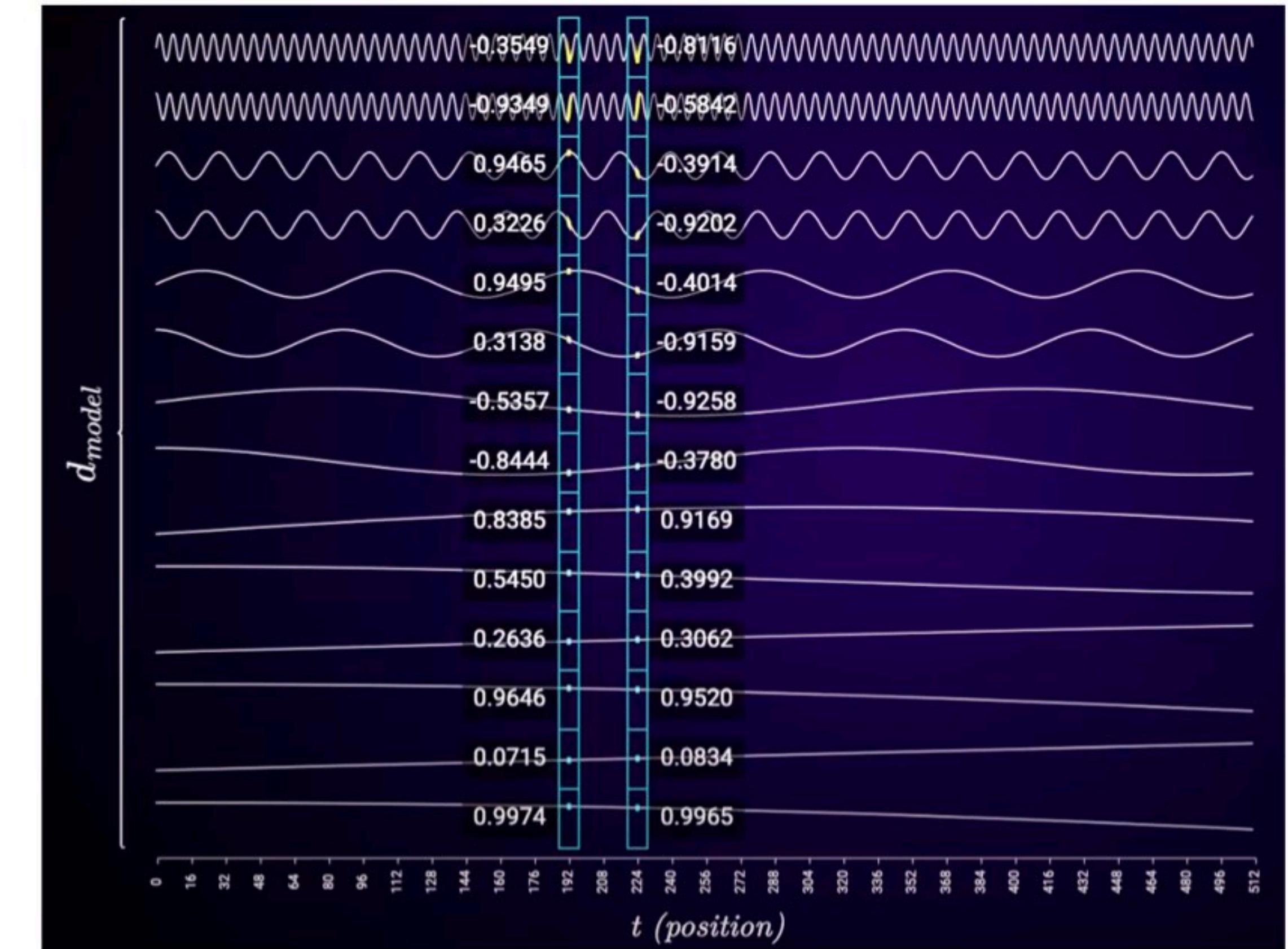
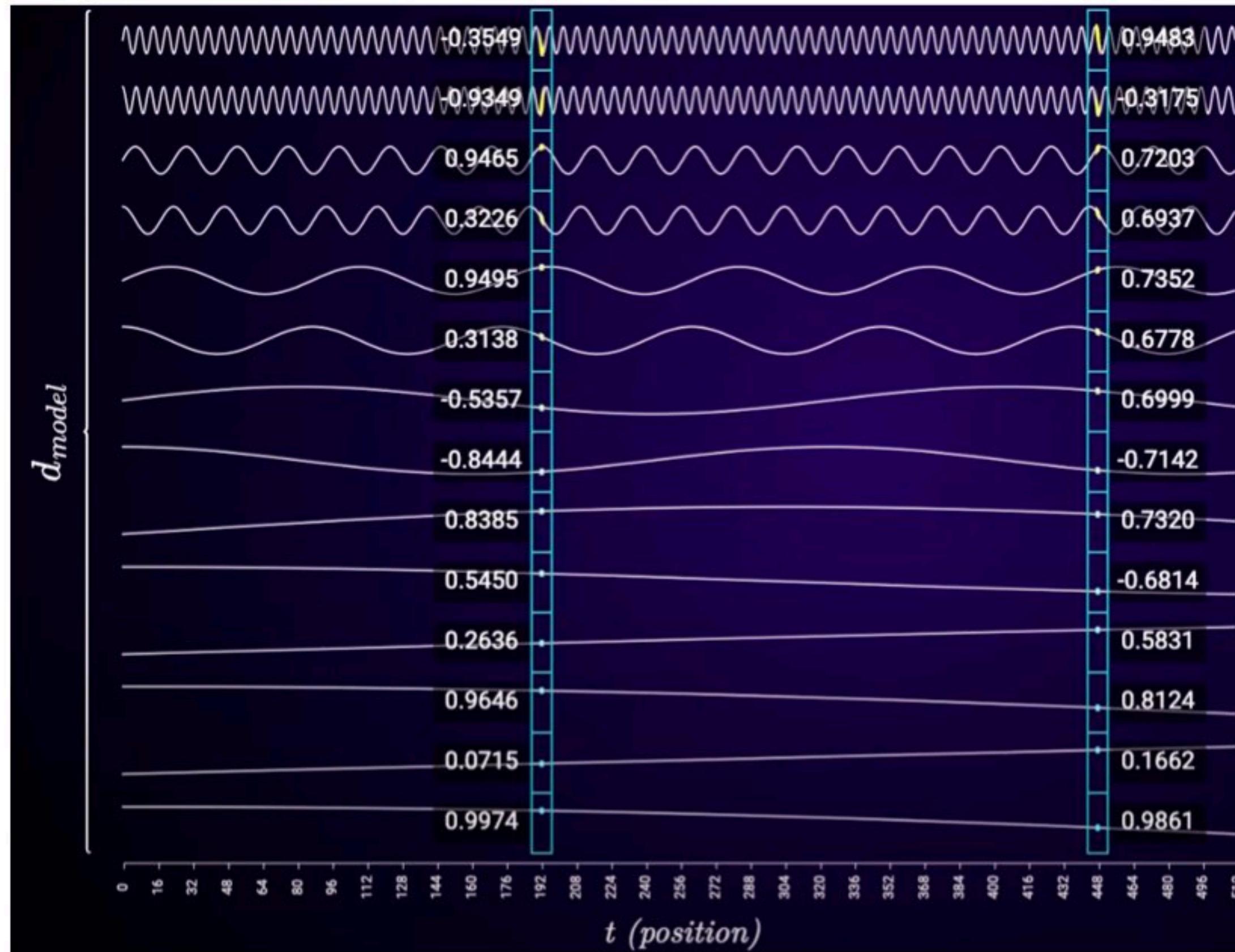
dimensionality
of positional
encoding



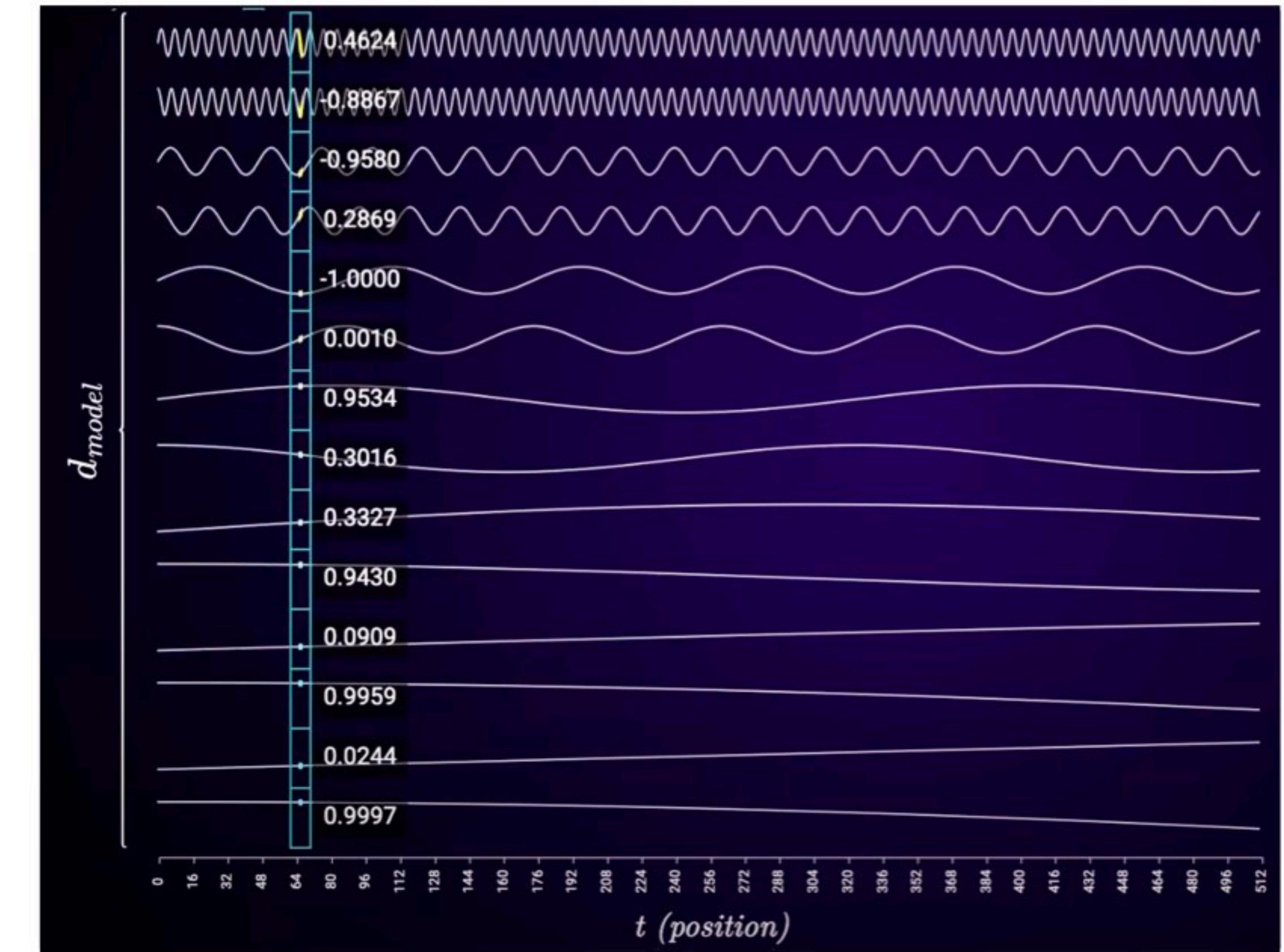
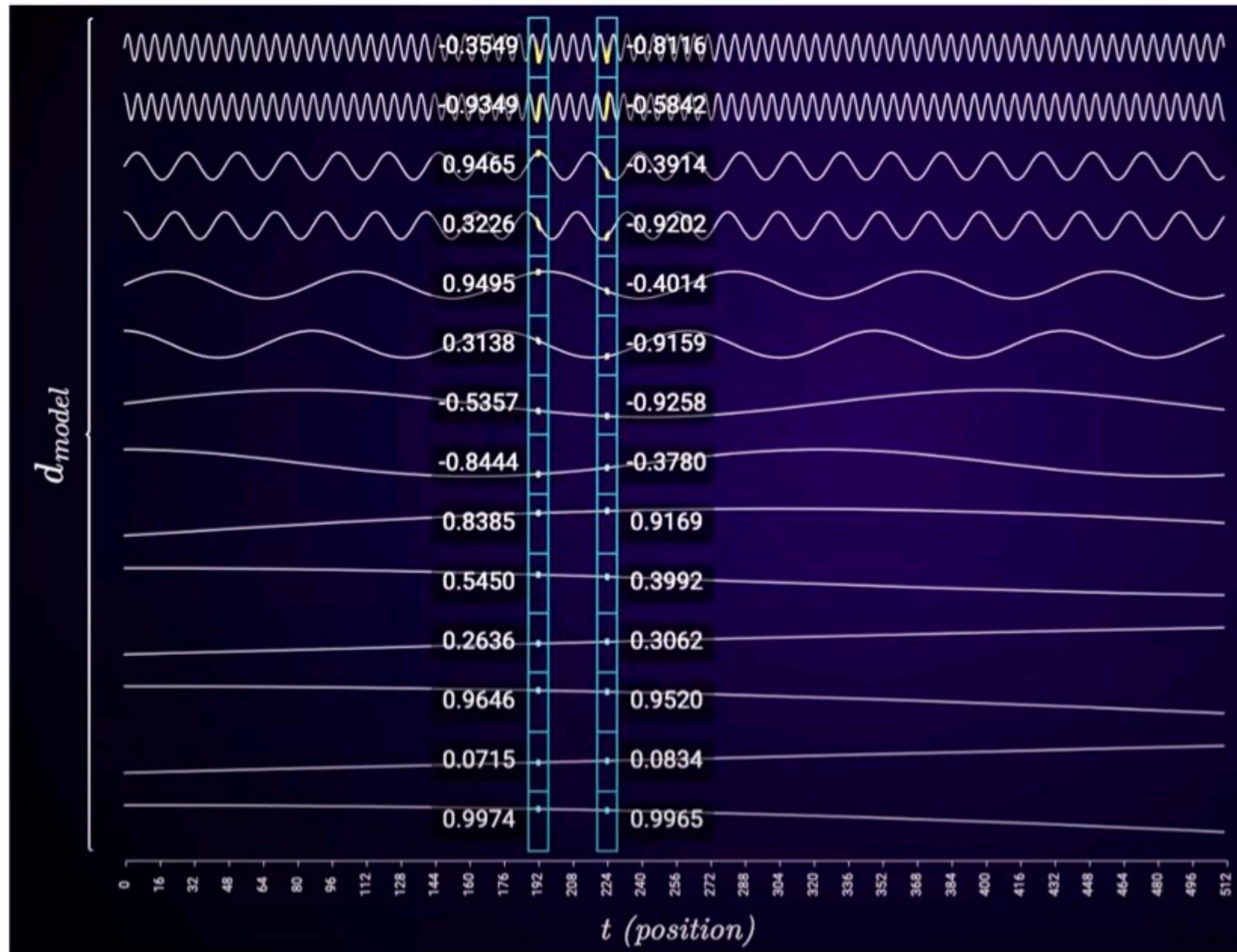
Positional Encoding



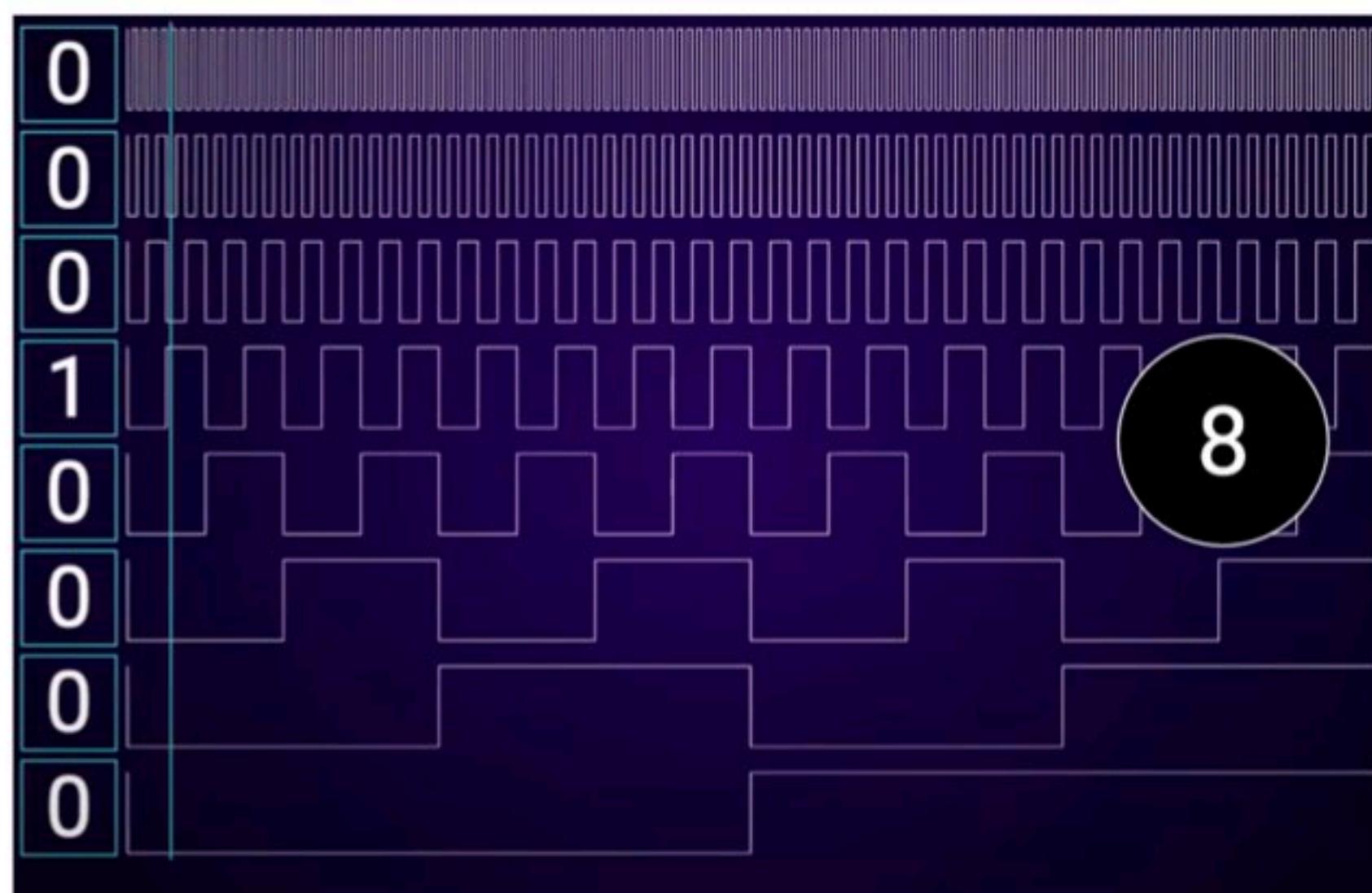
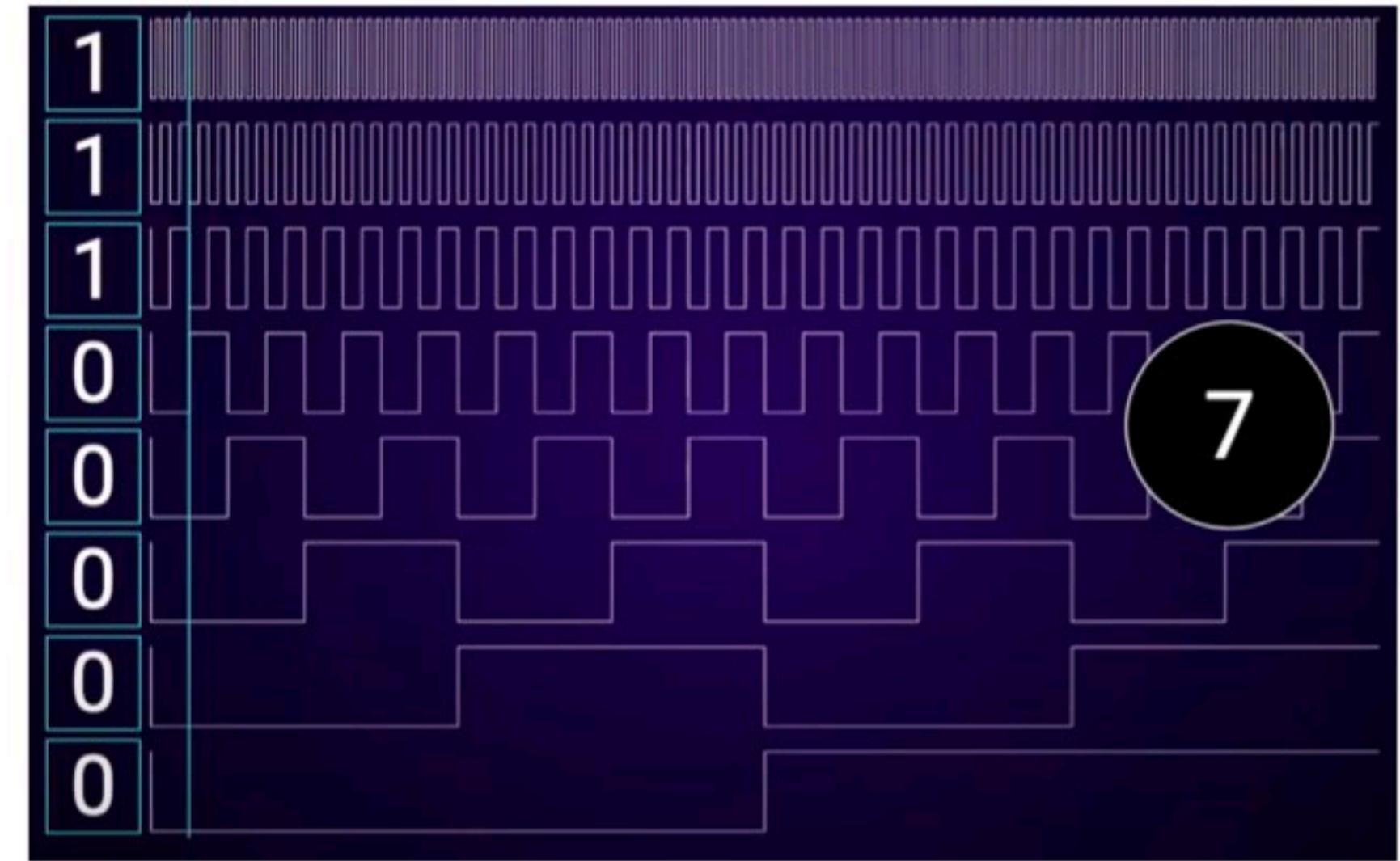
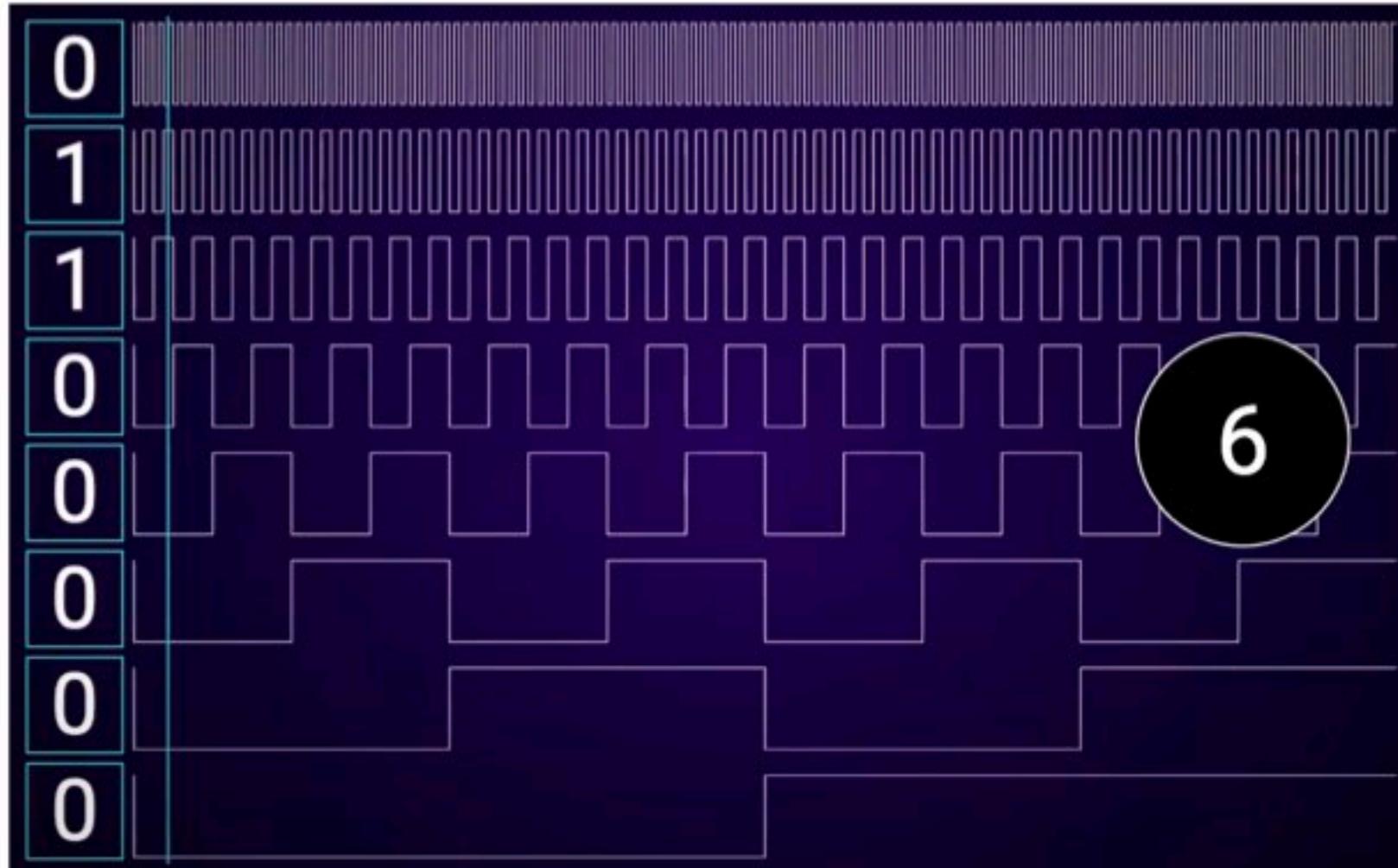
Positional Encoding



Positional Encoding

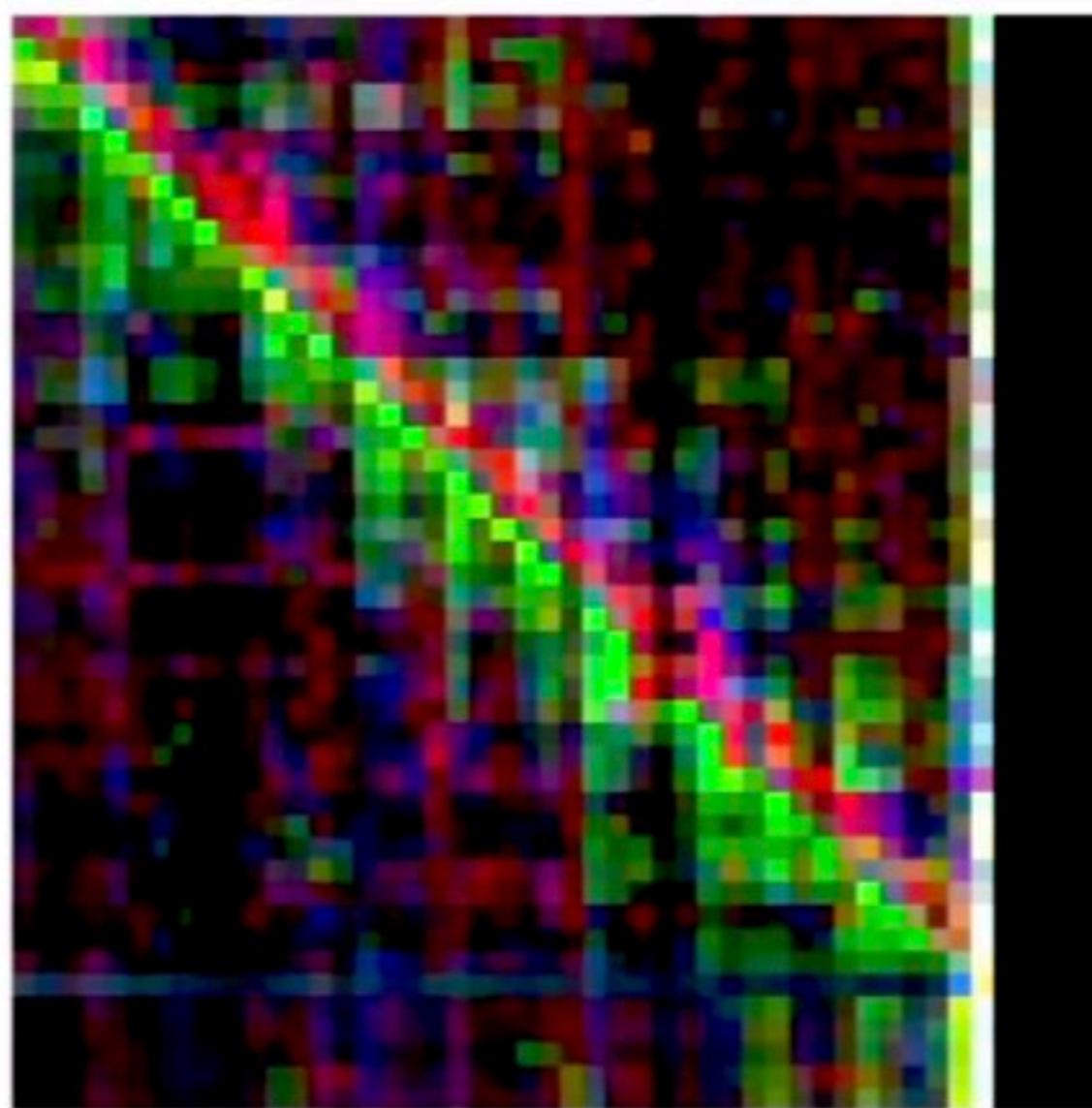


Positional Encoding

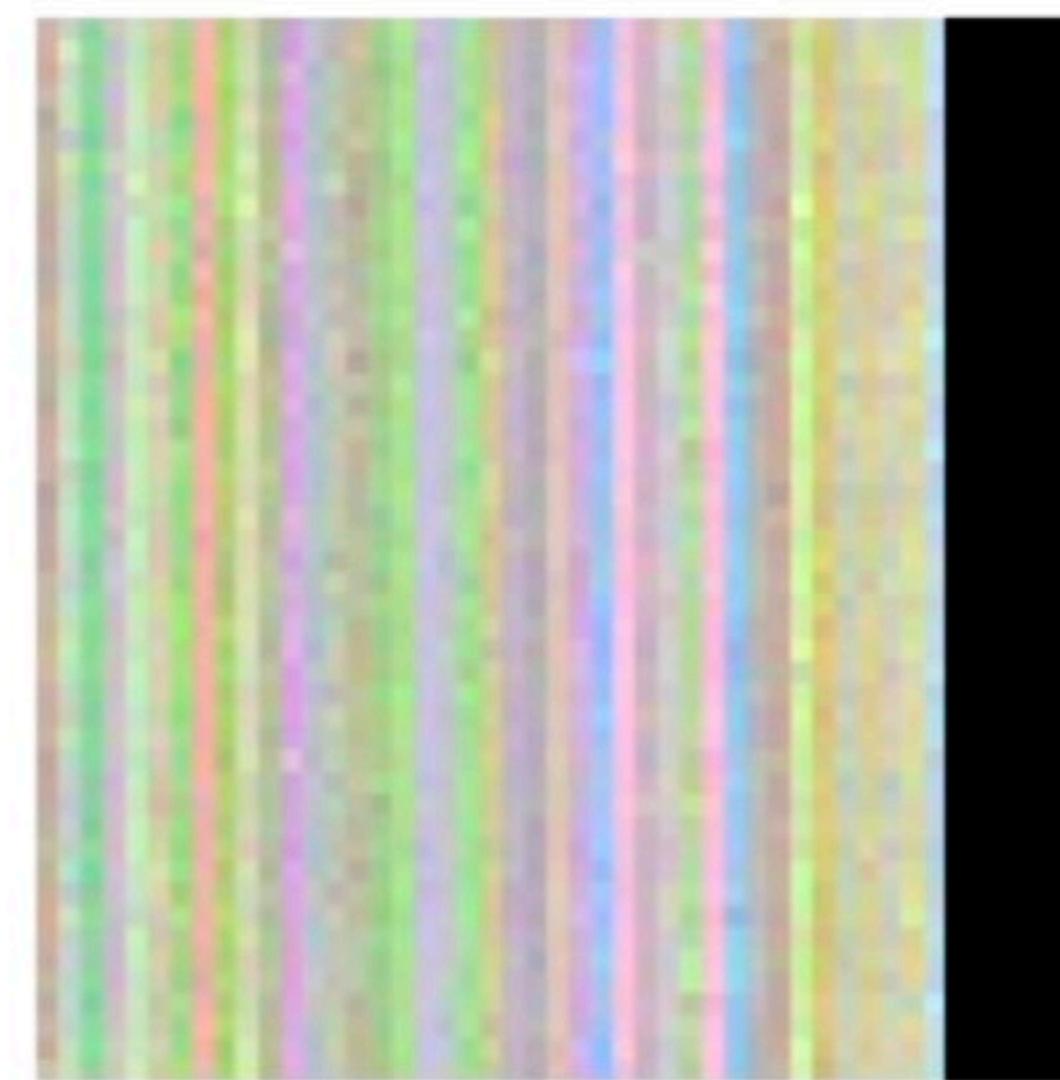


Importance of Residual Connection

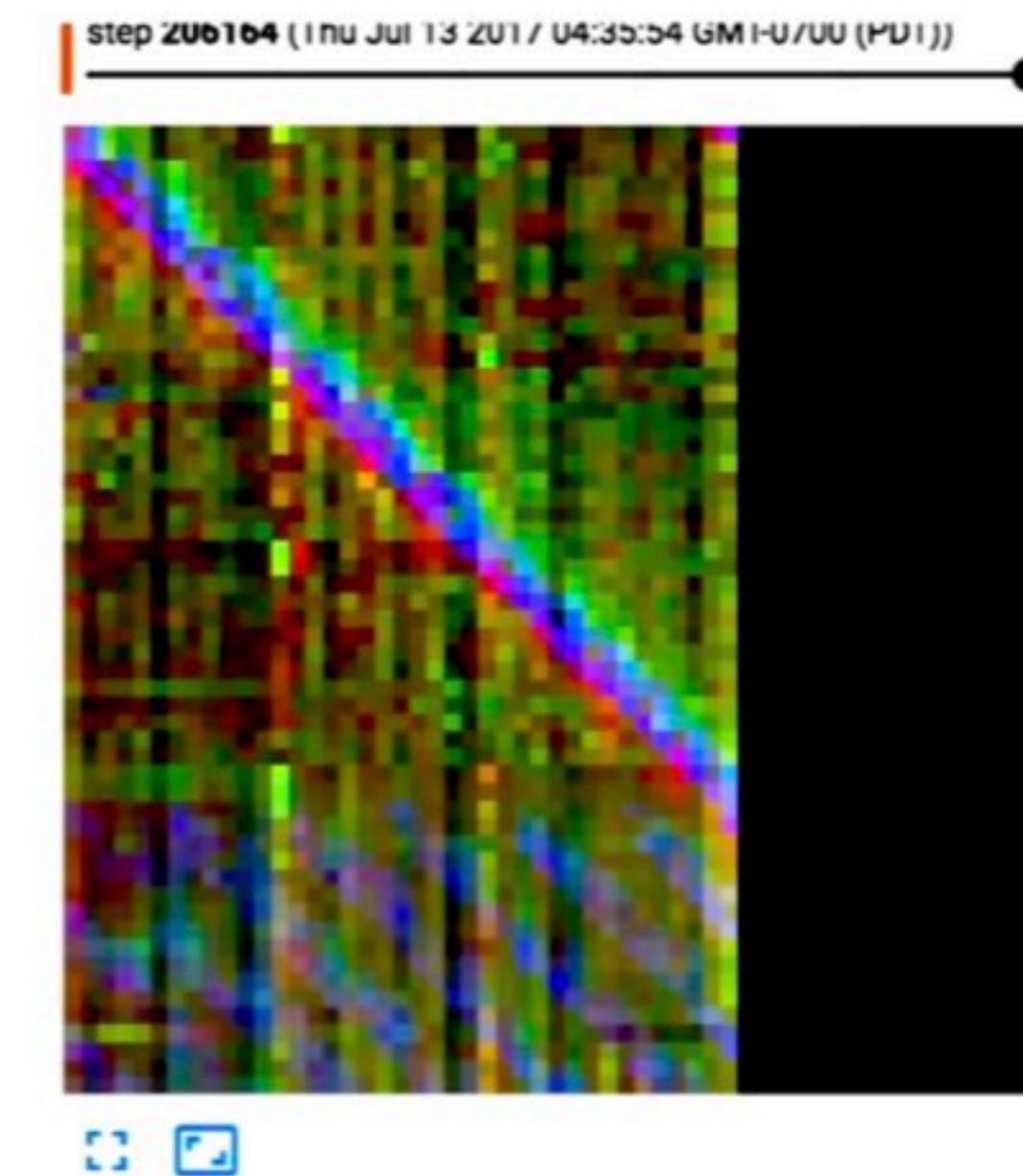
Residuals carry positional information to higher layers, among other information.



With residuals

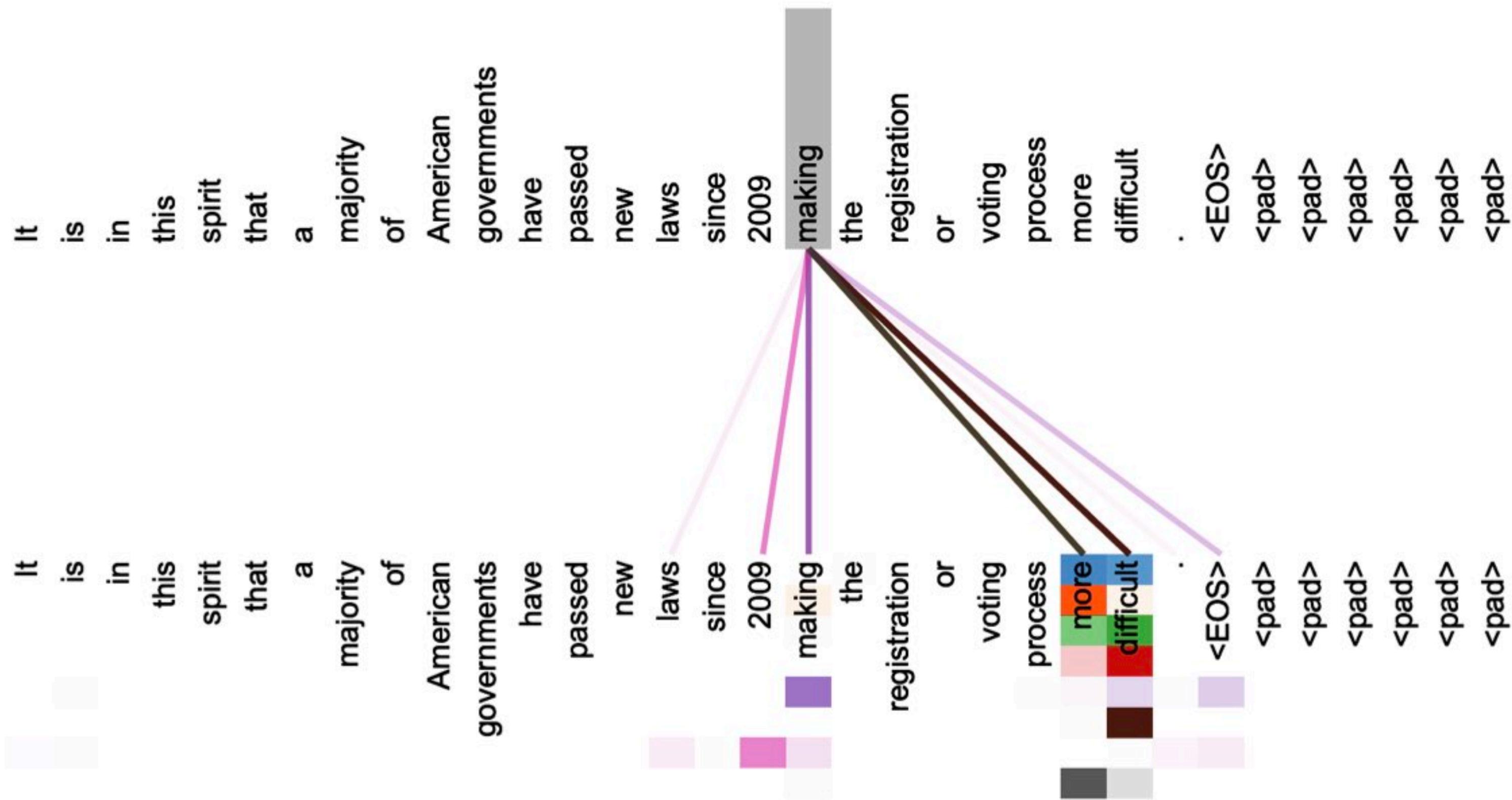


Without residuals

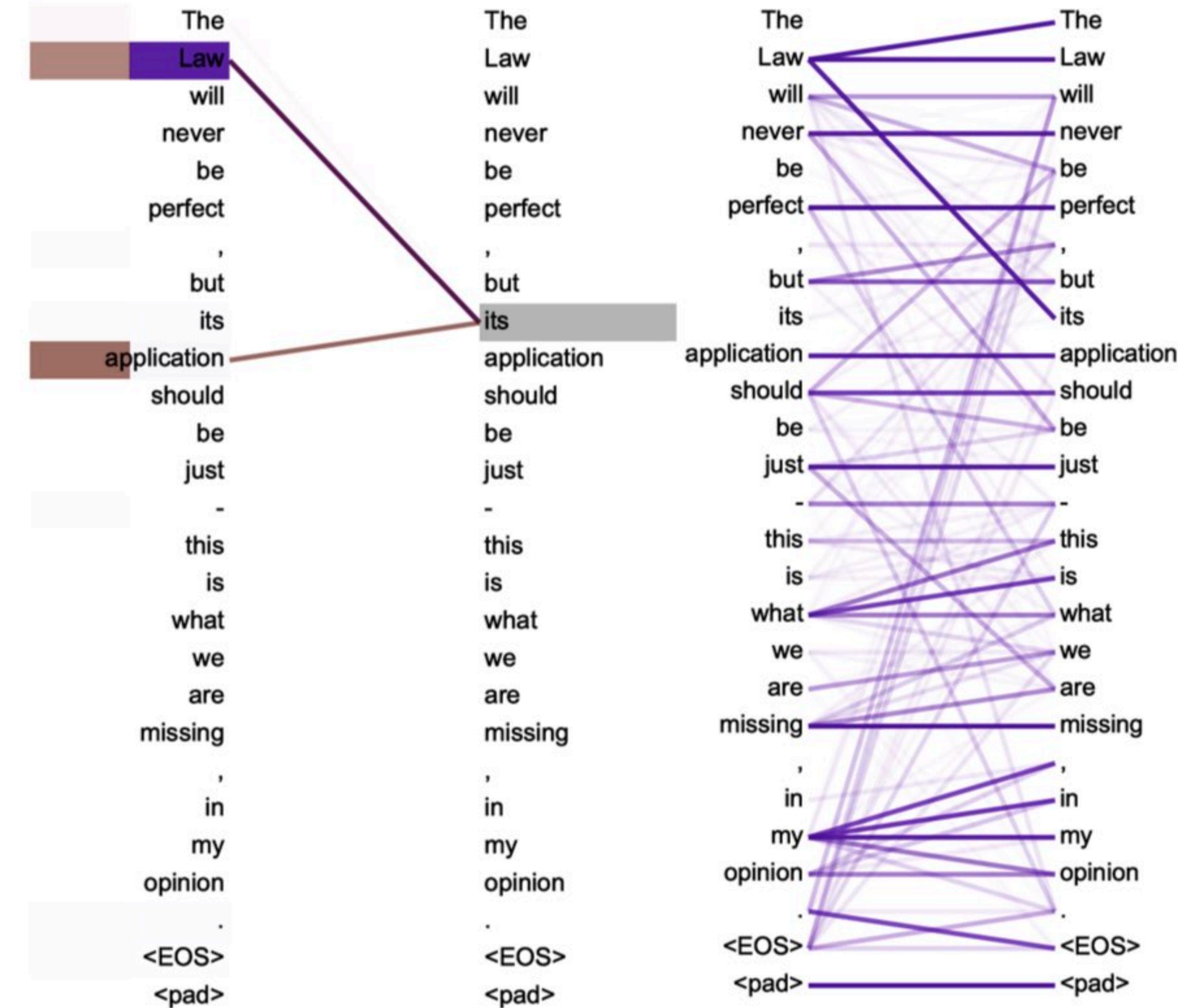


Without residuals,
with timing signals

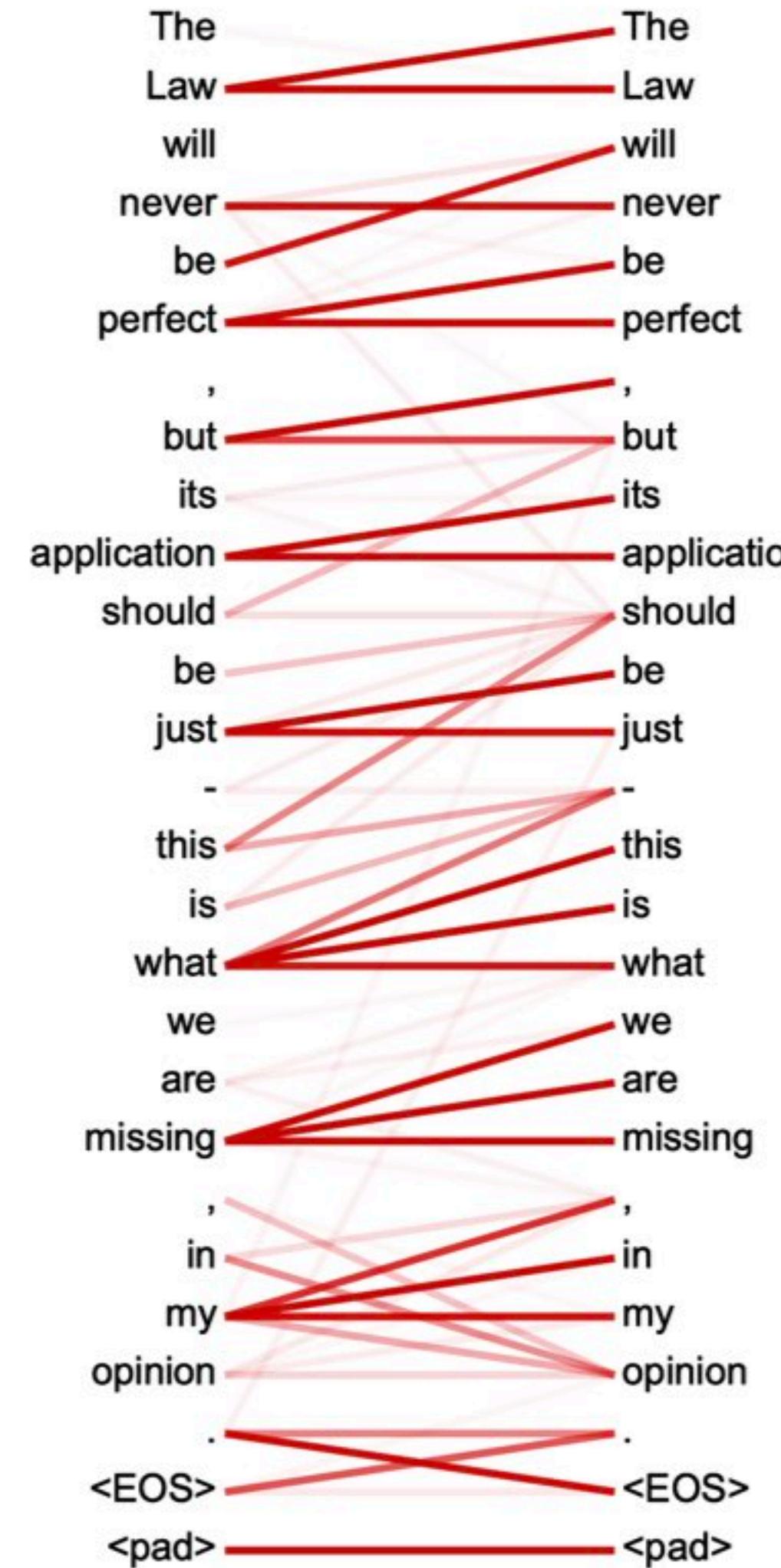
Multi-head Attention intuition



Multi-head Attention intuition



Multi-head Attention intuition



Limitations of Transformers

- Computationally Expensive
- Lack of interpretability
- Limited sequence length