

Seq2Seq & Transformer Model

Indonesia AI

Proprietary document of Indonesia AI 2023



OBJECTIVE & OUTLINE

Proprietary document of Indonesia AI 2023

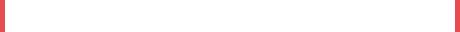


Seq2Seq & Transformer Model

Objektif: Memahami konsep dari Seq2Seq & Transformer Model dalam NLP

Outline:

1. Seq2Seq Model
2. Attention Model
3. Inside Transformer
4. Transformers Development

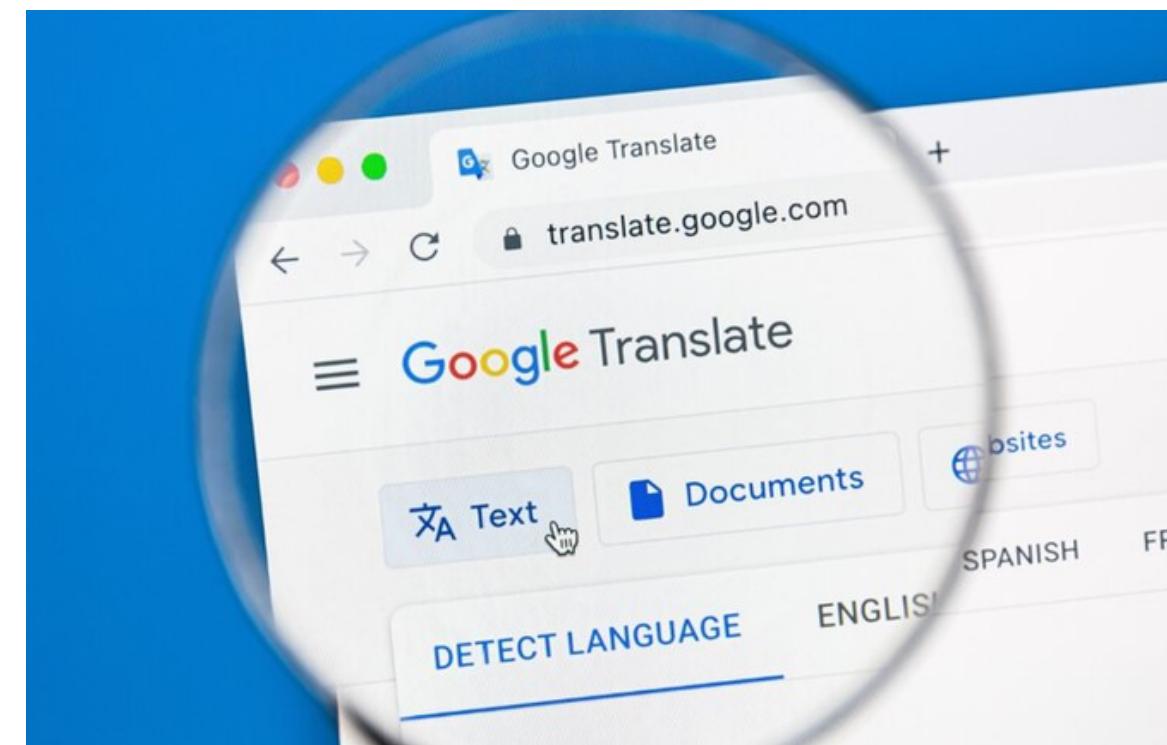


— Seq2Seq Model

EXAMPLE

Proprietary document of Indonesia AI 2023

das Haus ist gross



→ the house is big

EXAMPLE

Proprietary document of Indonesia AI 2023

Albert lives in Barlimore



→ PER NONE NONE LOC

NER

HANDLING SEQUENCES

Proprietary document of Indonesia AI 2023



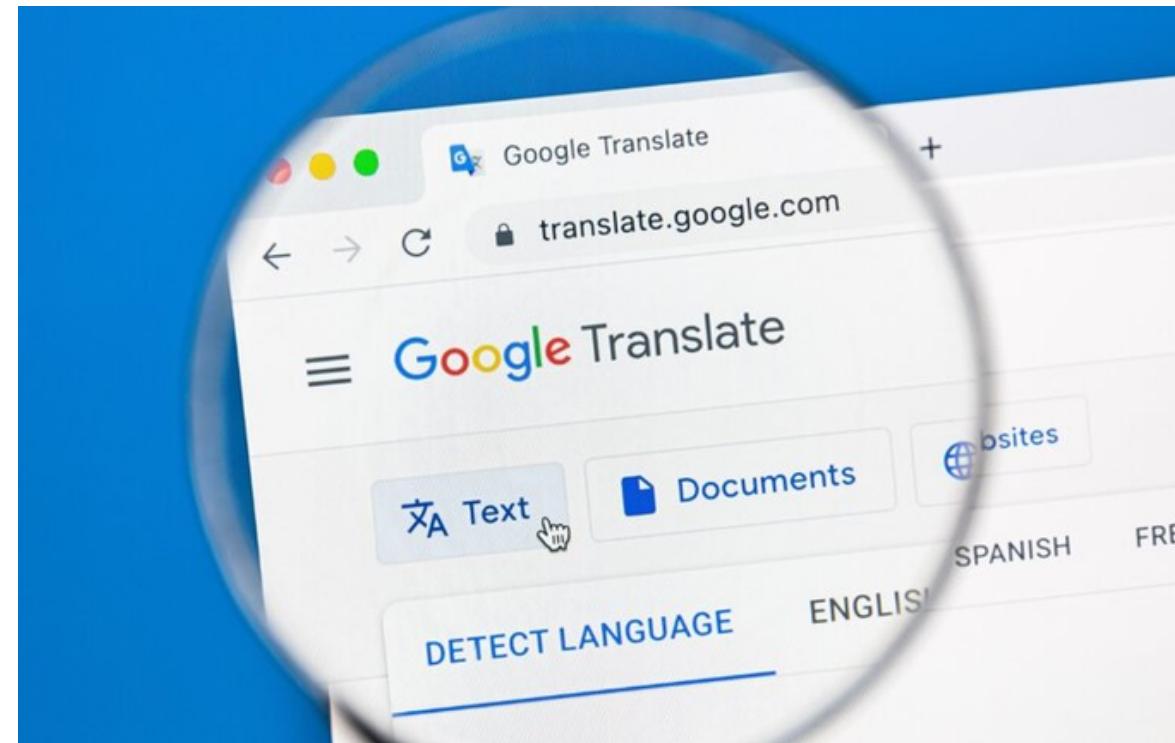
Encoder: Convert arbitrary length input to some fixed-length hidden representation

Decoder: generate arbitrary length output

EXAMPLE

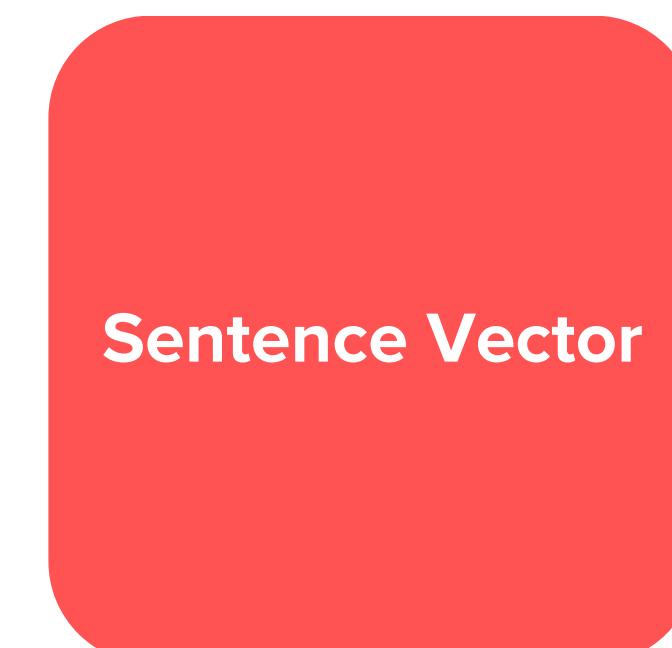
Proprietary document of Indonesia AI 2023

das Haus ist gross →



→ the house is big

Encoder
das Haus ist gross →



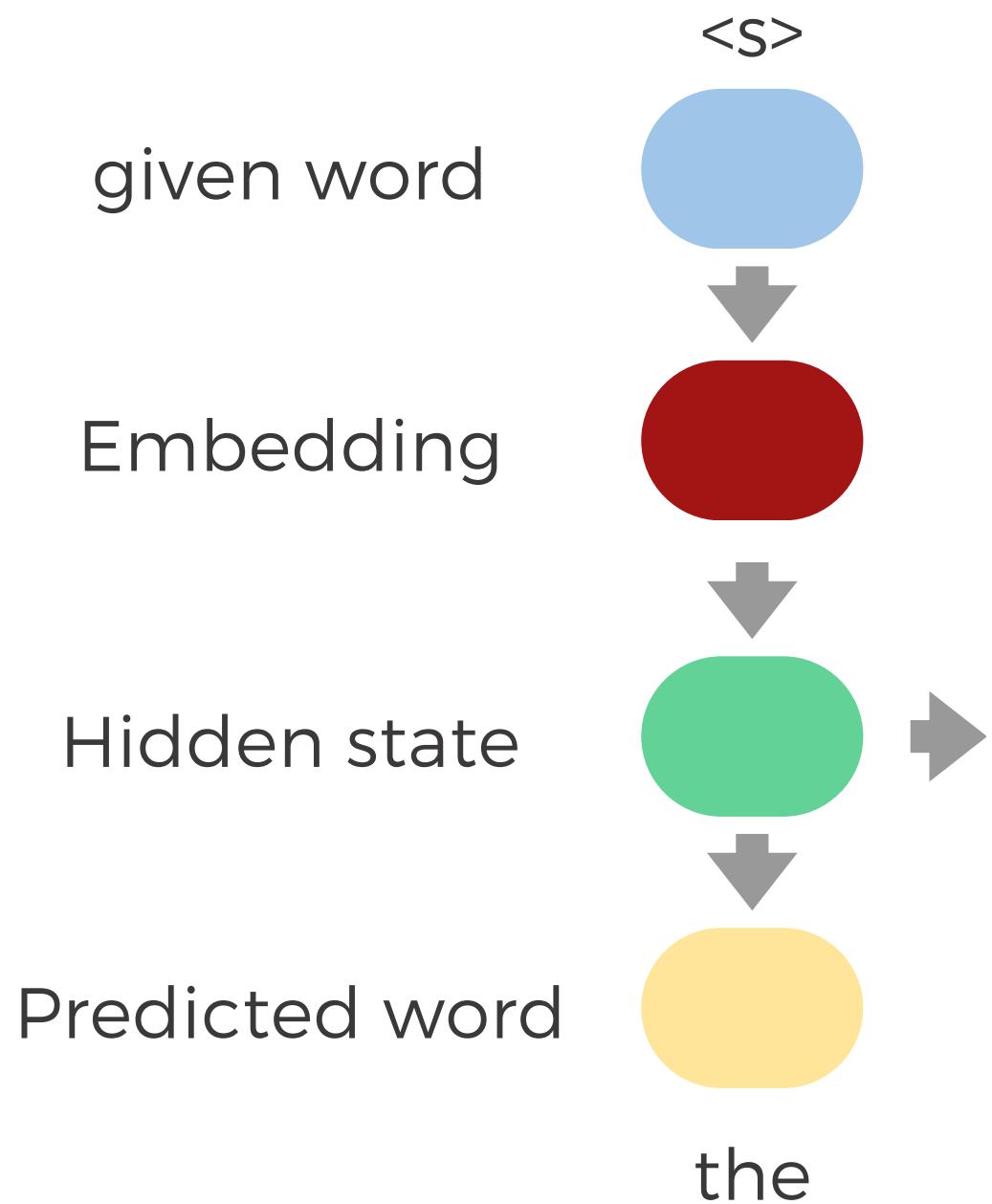
Decoder →
step 1: the
step 2: house
step 3: is
step 4: big
step 5: <stop>

Indonesia AI

each step applies a softmax over all vocab

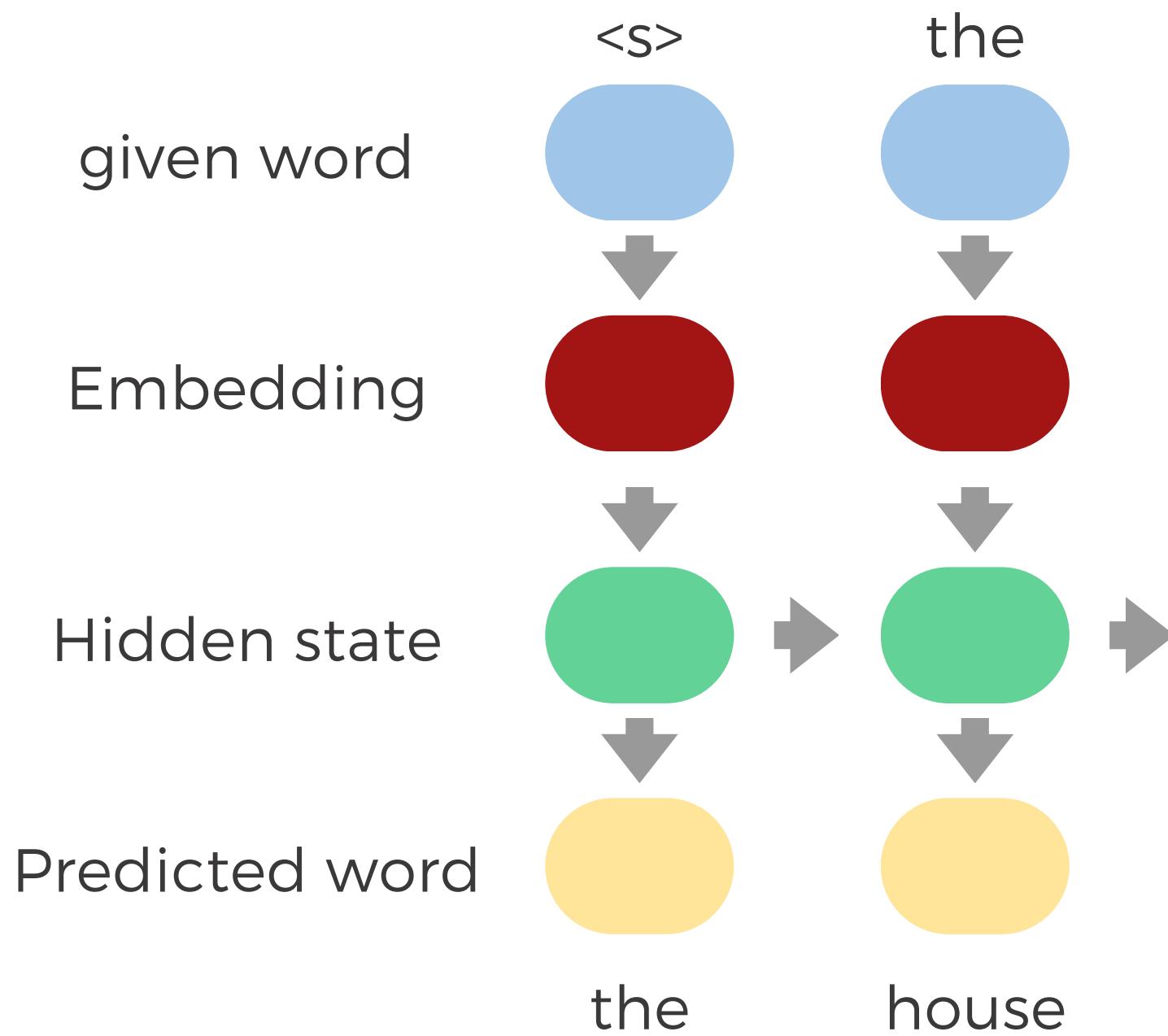
EXAMPLE

Proprietary document of Indonesia AI 2023



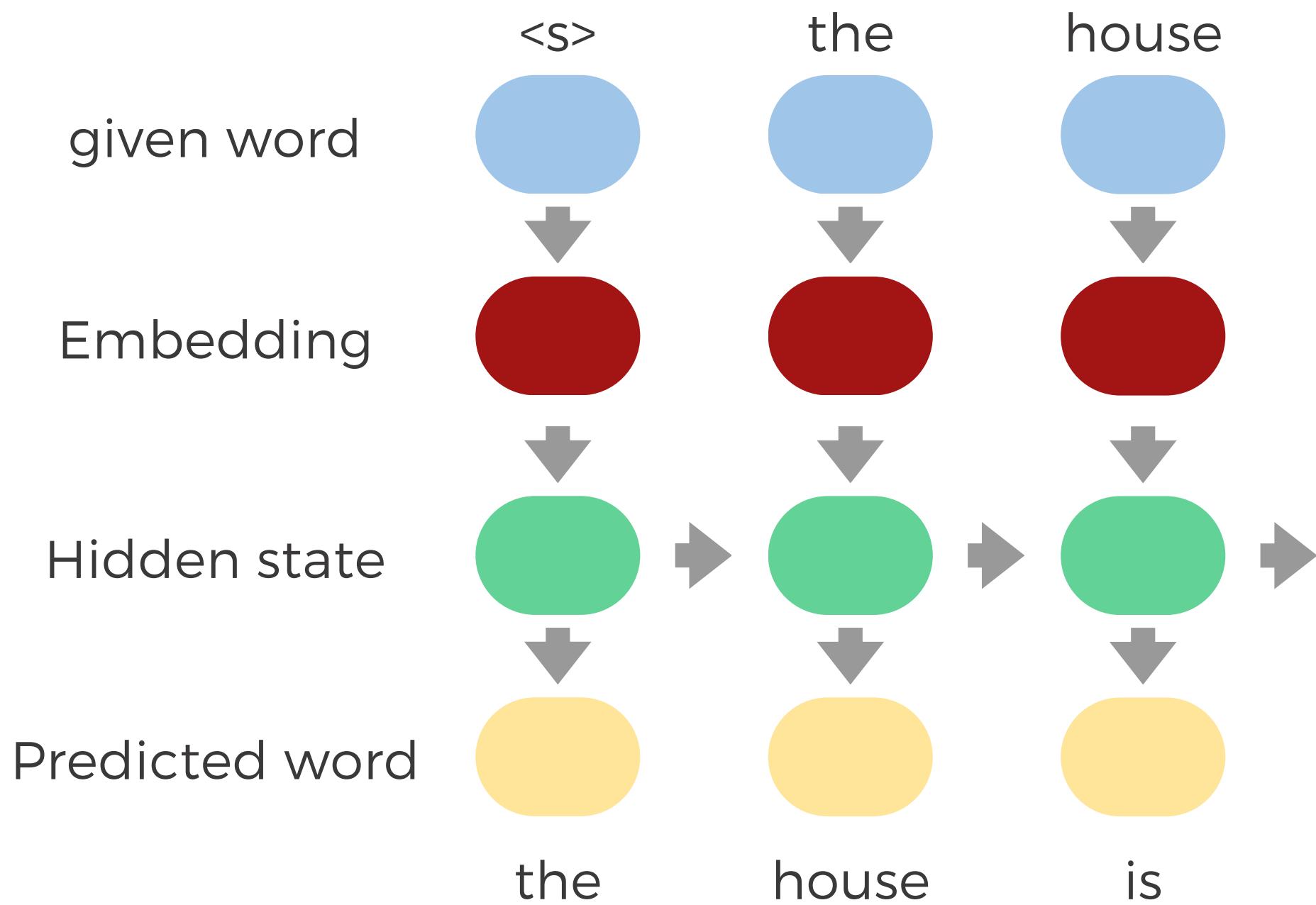
EXAMPLE

Proprietary document of Indonesia AI 2023



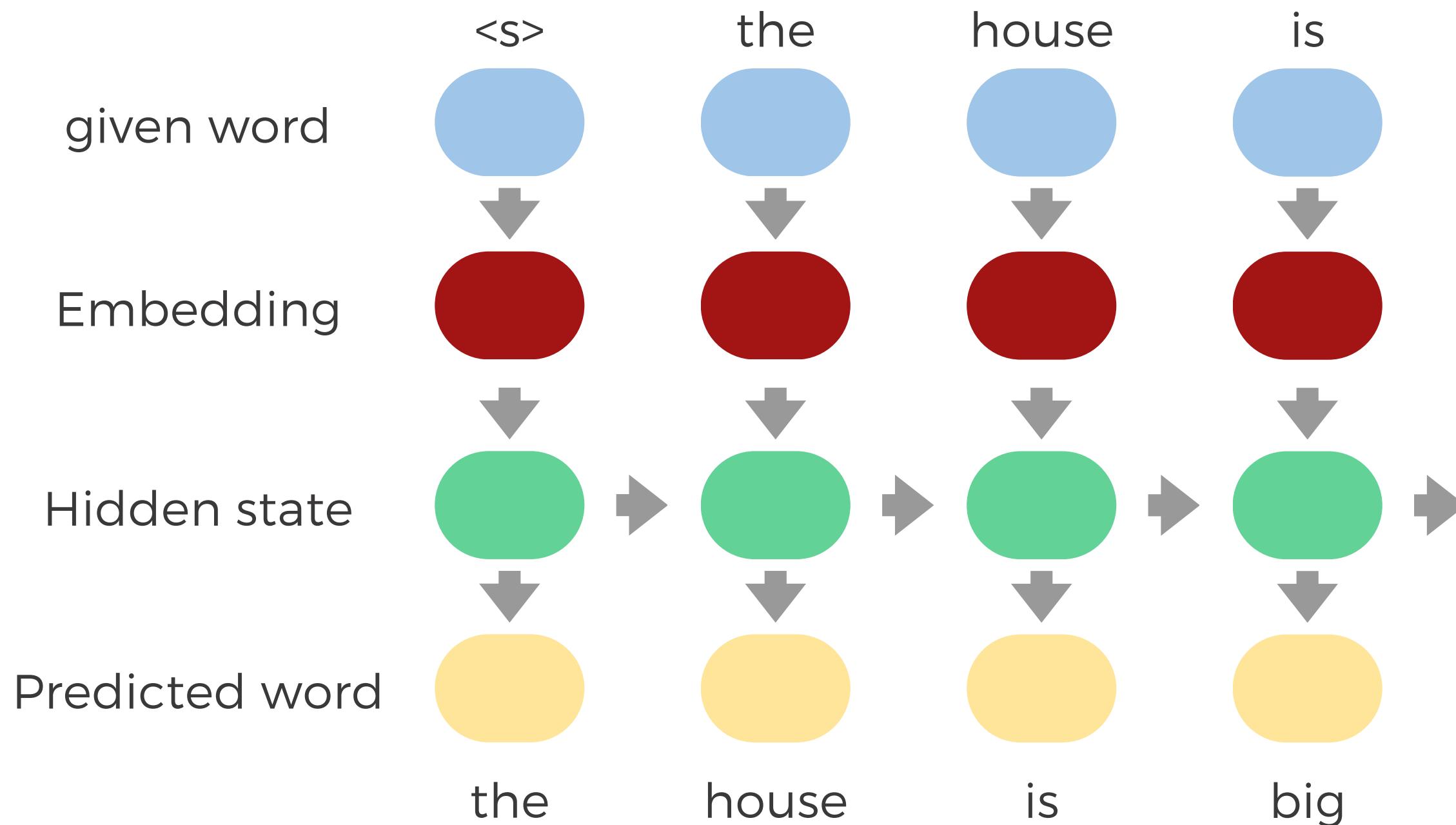
EXAMPLE

Proprietary document of Indonesia AI 2023



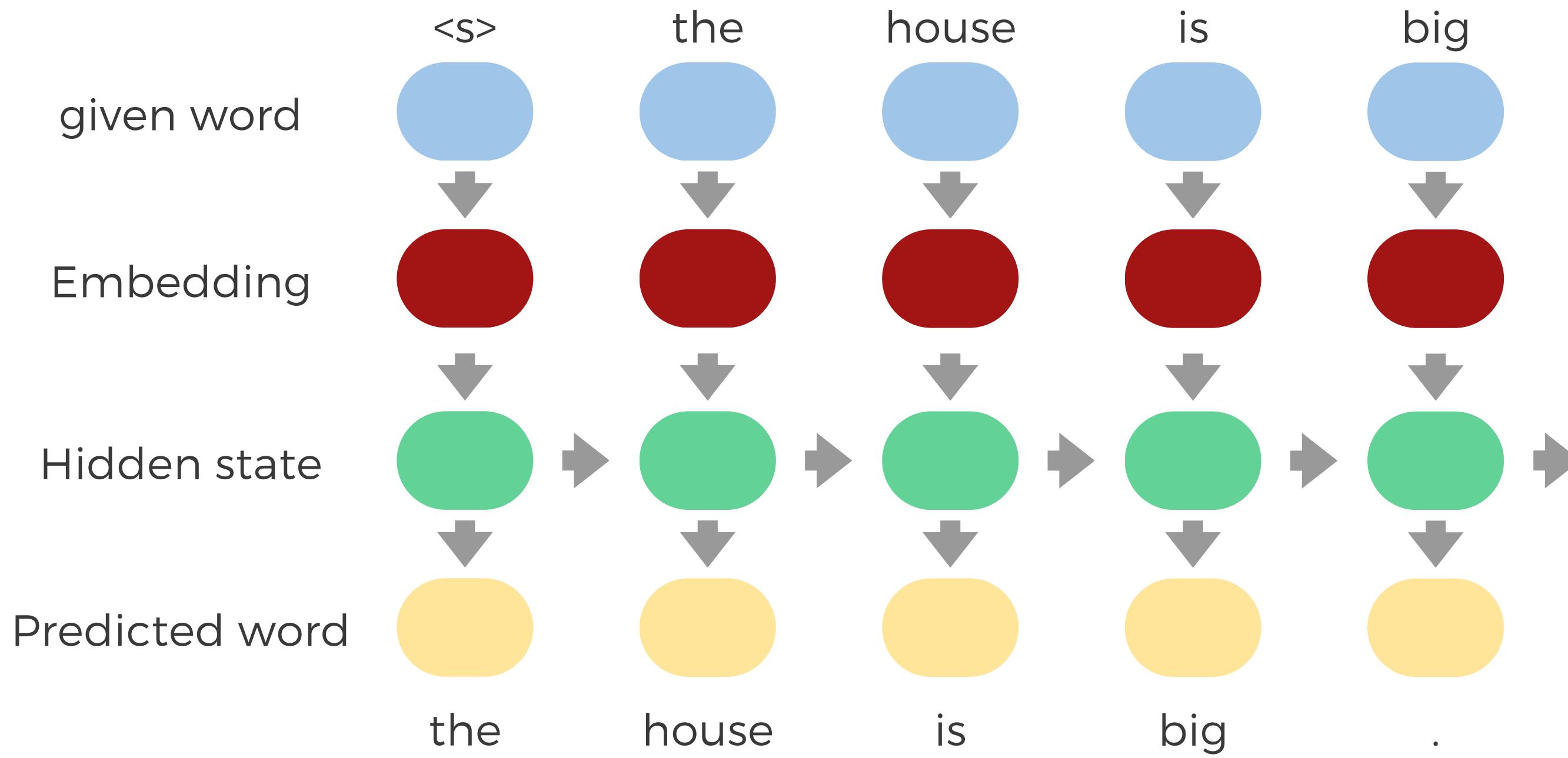
EXAMPLE

Proprietary document of Indonesia AI 2023



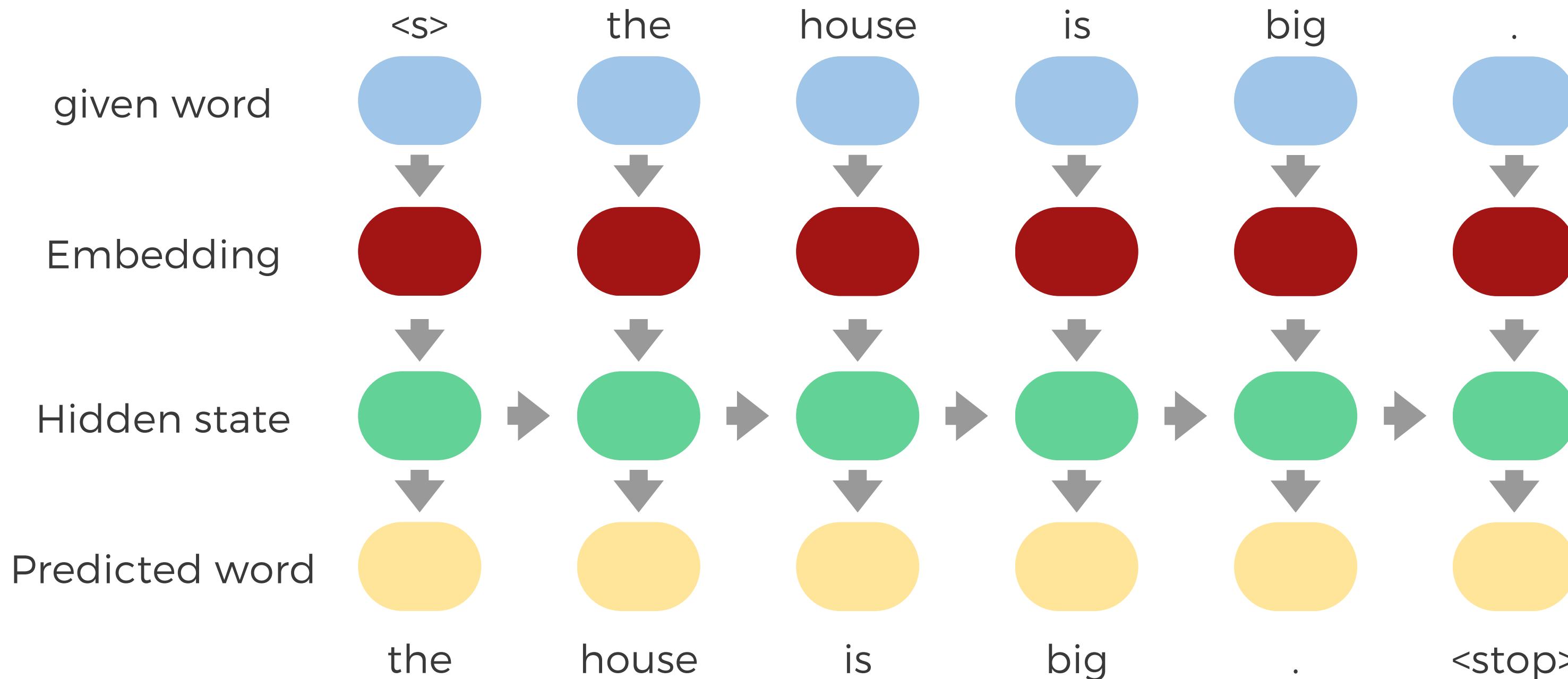
EXAMPLE

Proprietary document of Indonesia AI 2023



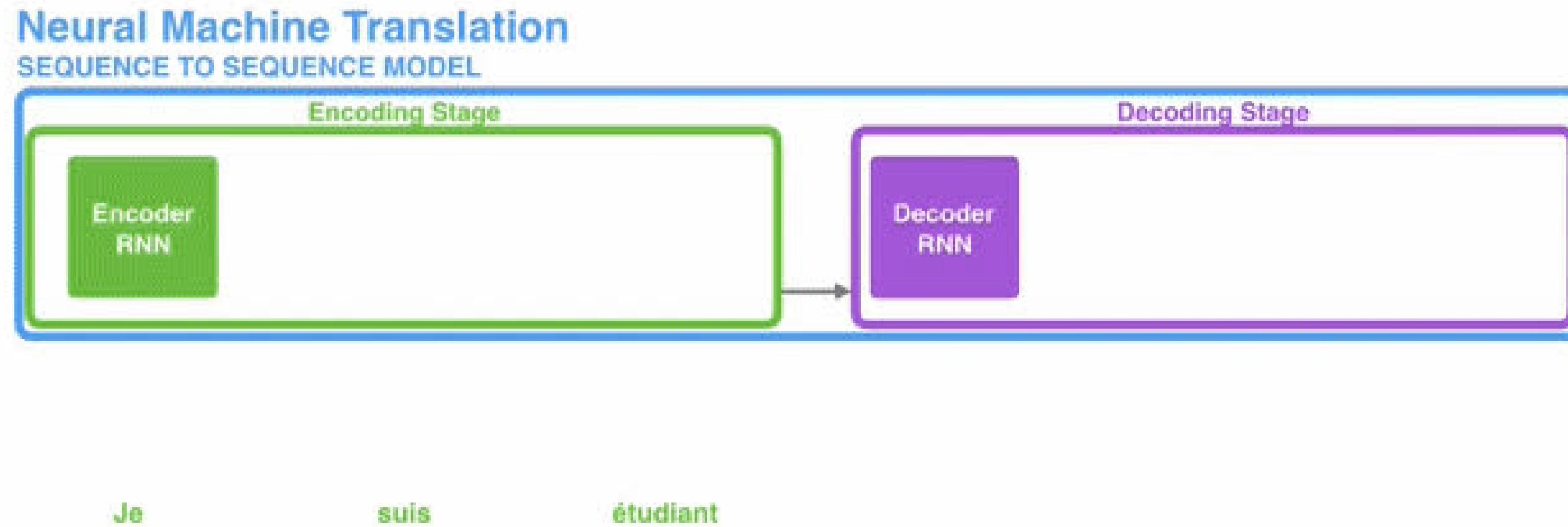
EXAMPLE

Proprietary document of Indonesia AI 2023



SEQ2SEQ ANIMATION

Proprietary document of Indonesia AI 2023





— Attention Model

Transformers

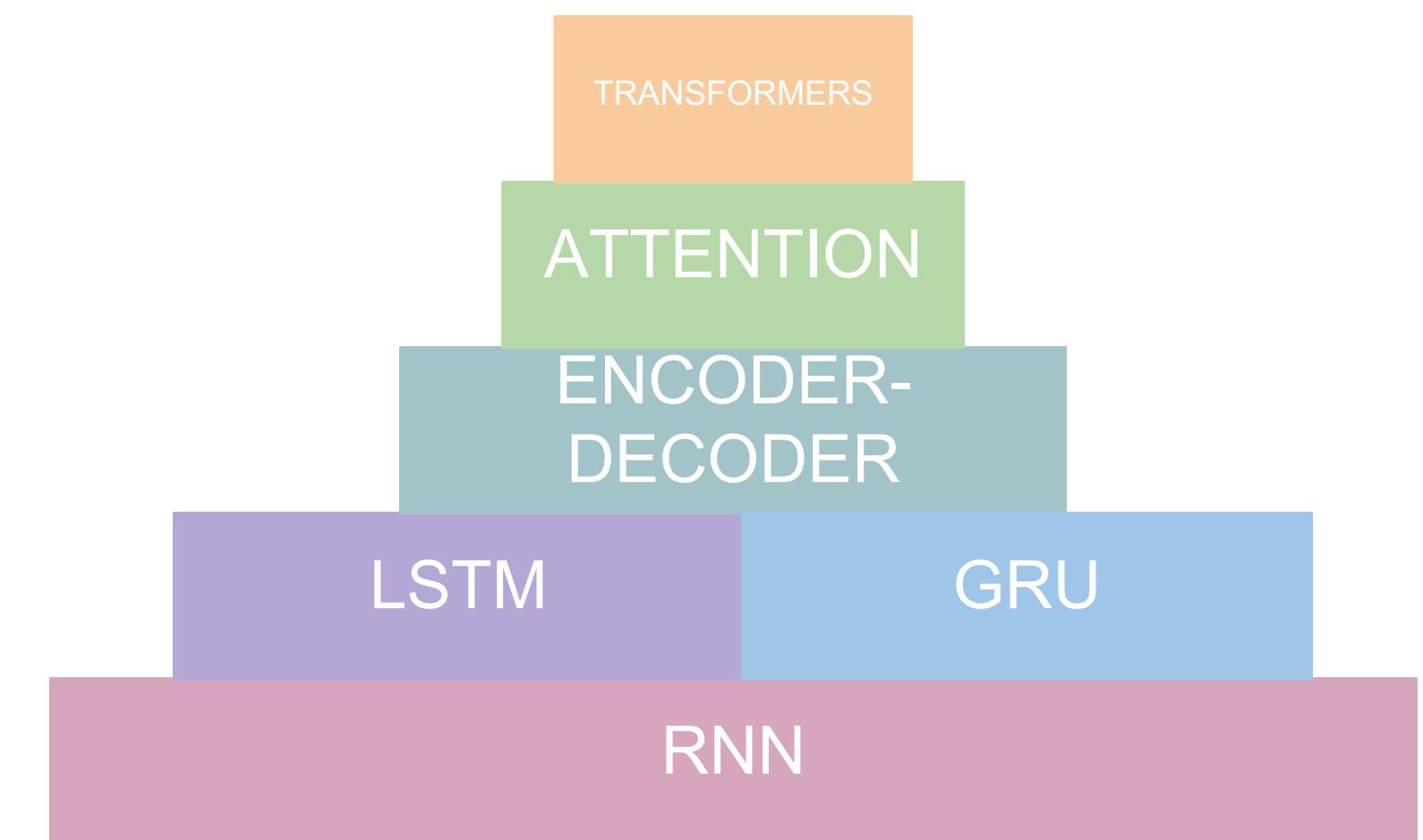
Paper

Attention Is All You
Need
(2017)

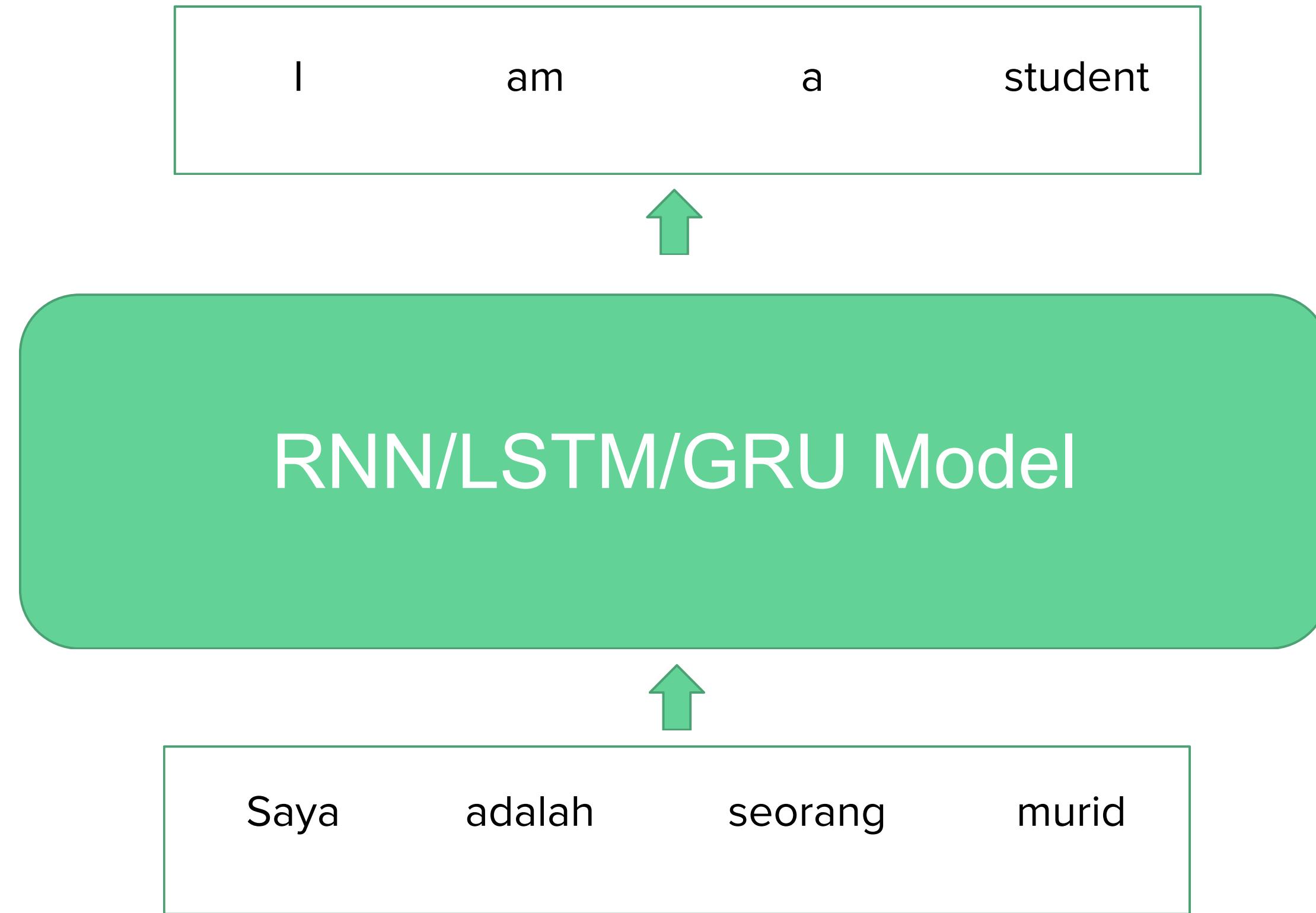
New family of neural
network architecture

SoTA for Machine
Translation, beats
Google MT

Designed for
parallelization



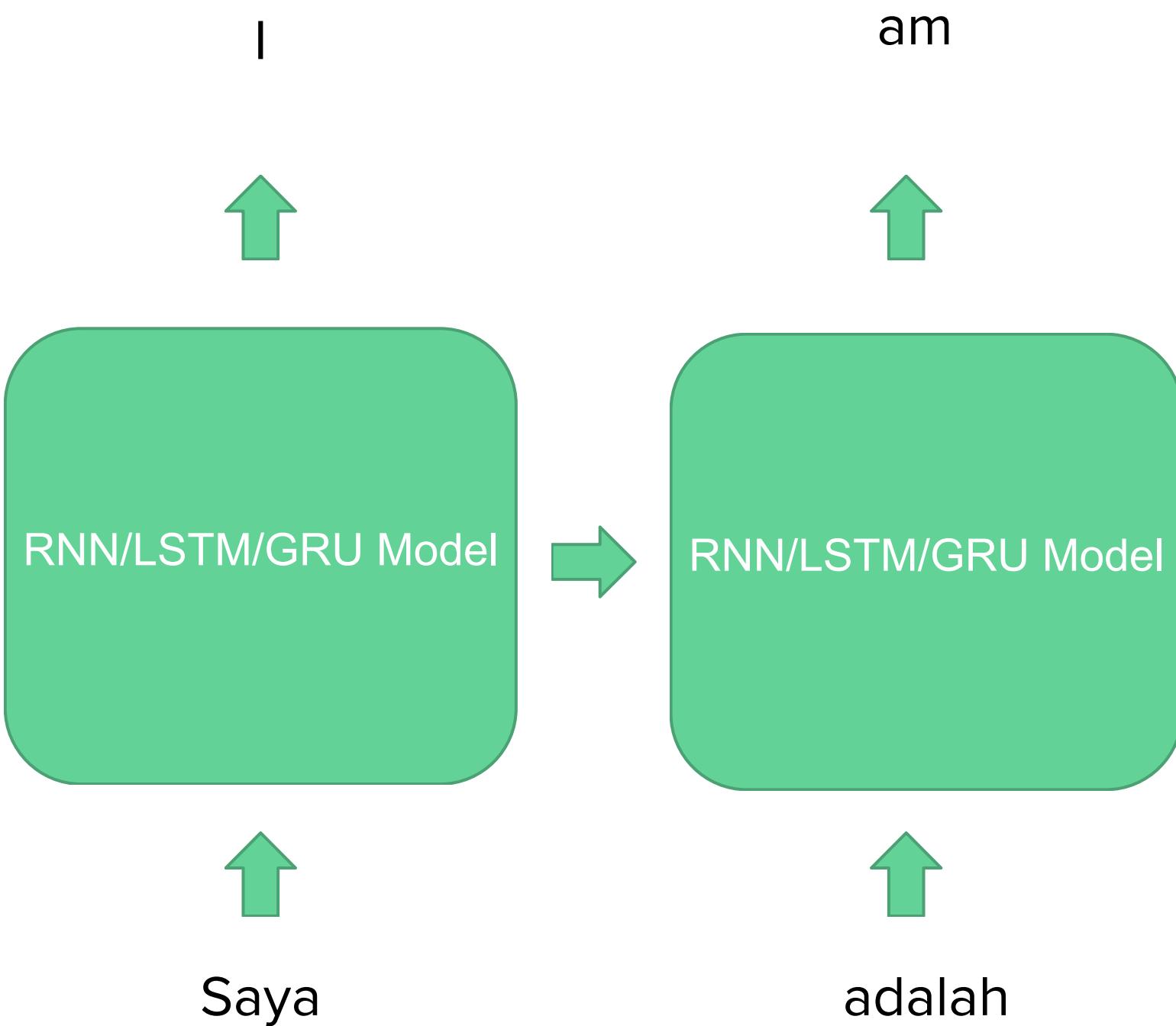
Old VS New Approach



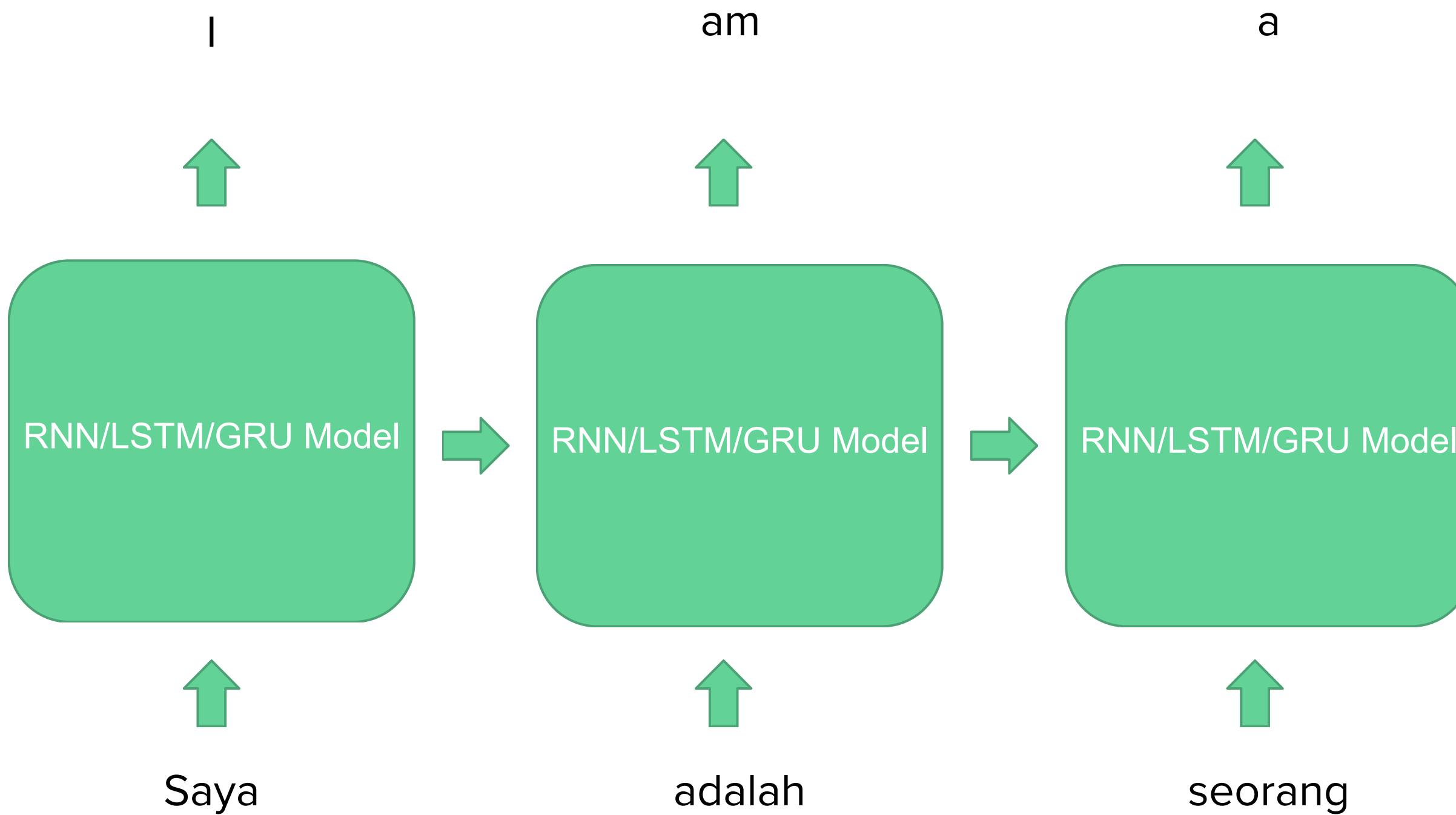
Old VS New Approach



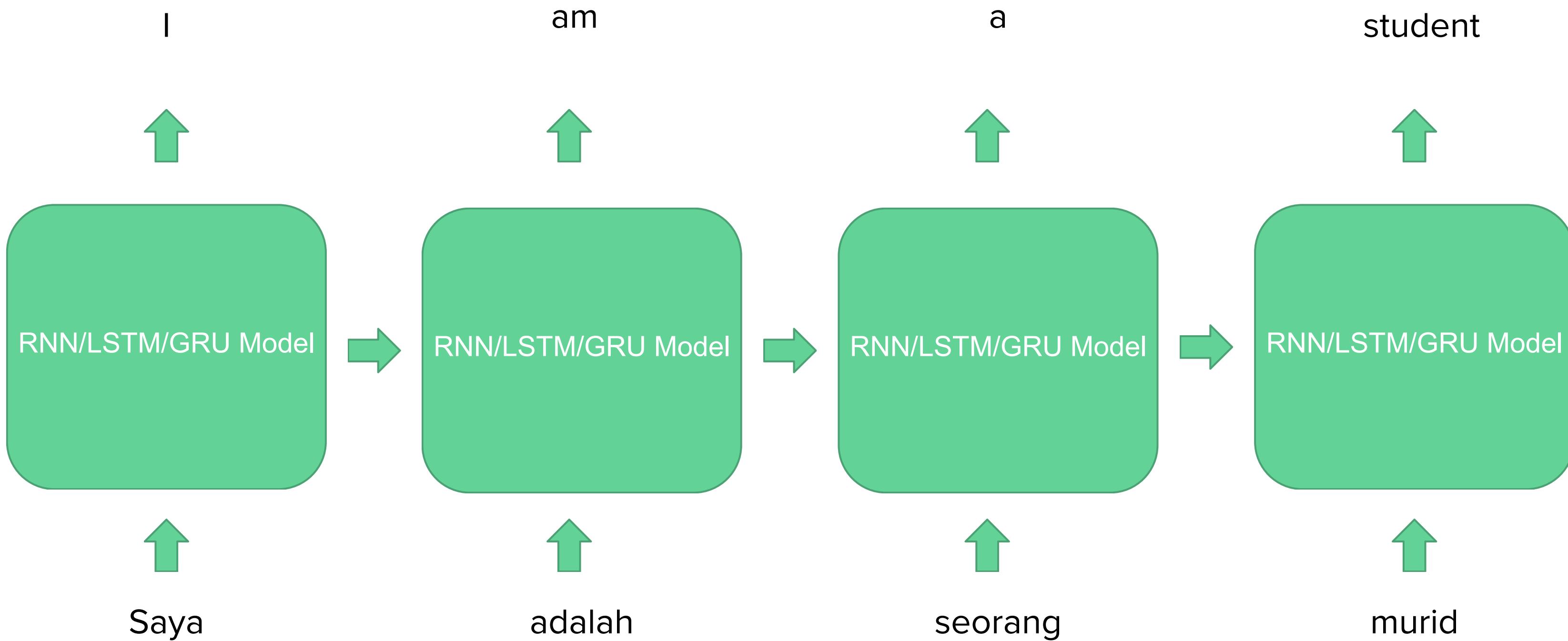
Old VS New Approach



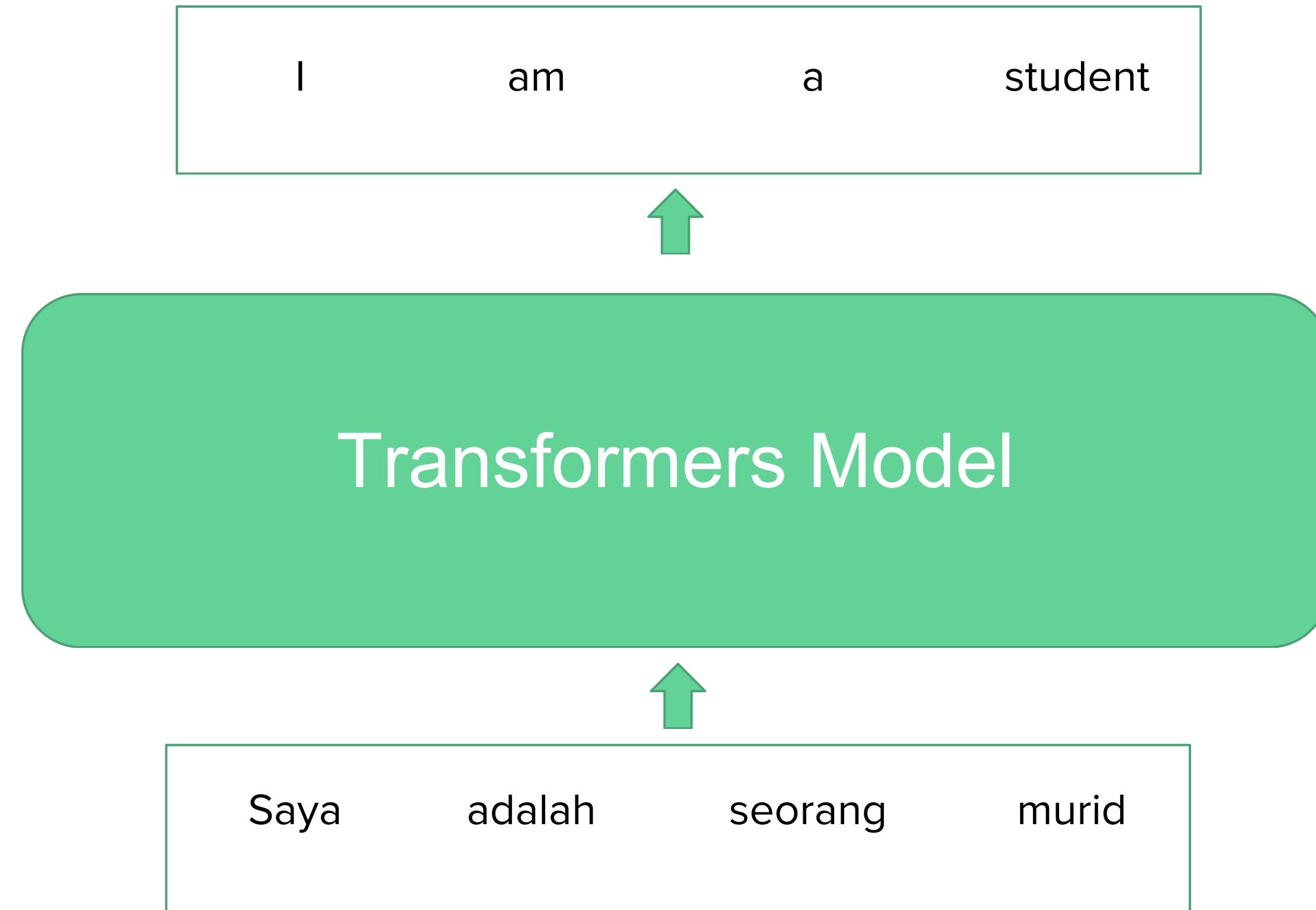
Old VS New Approach



Old VS New Approach



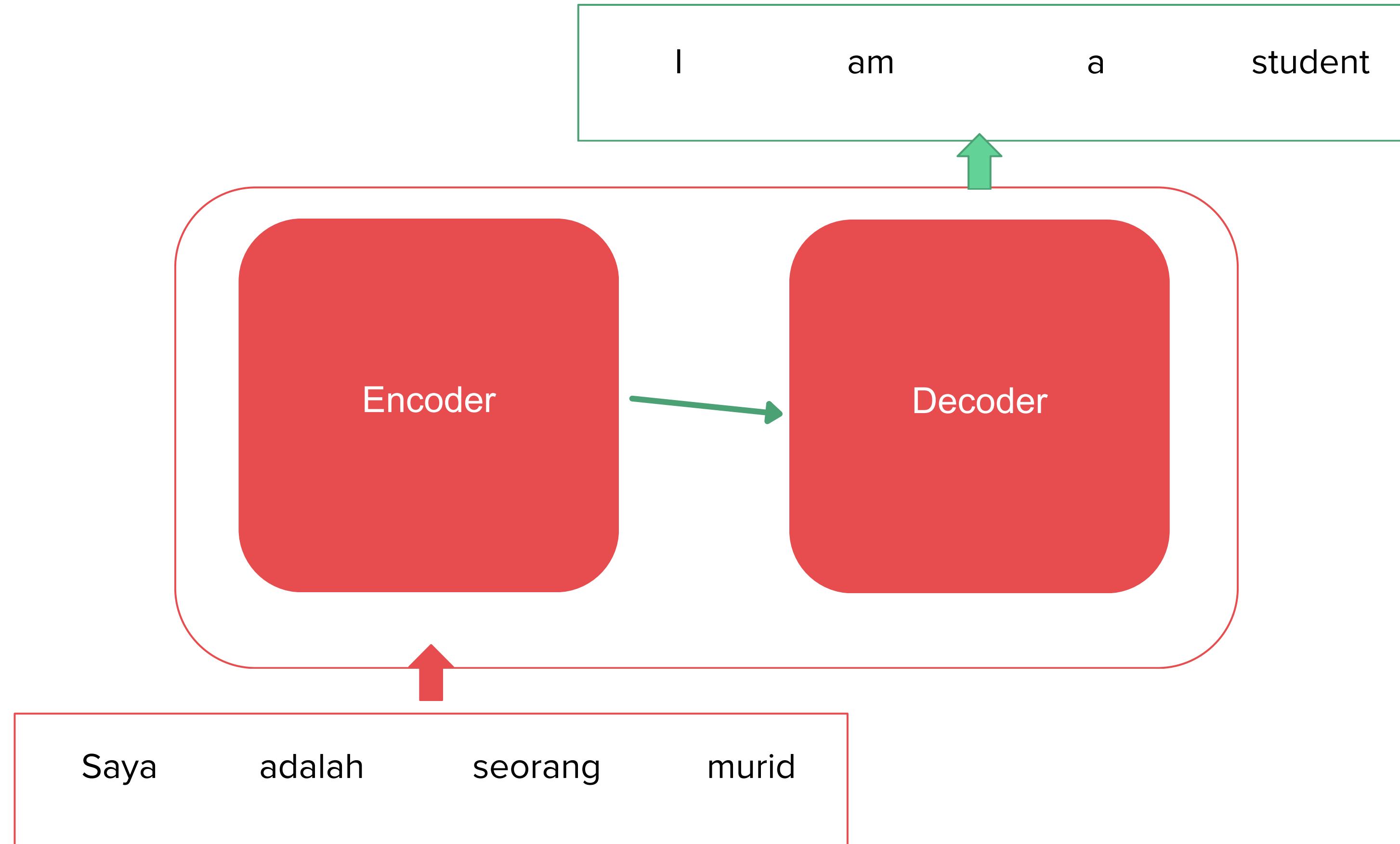
Old VS New Approach



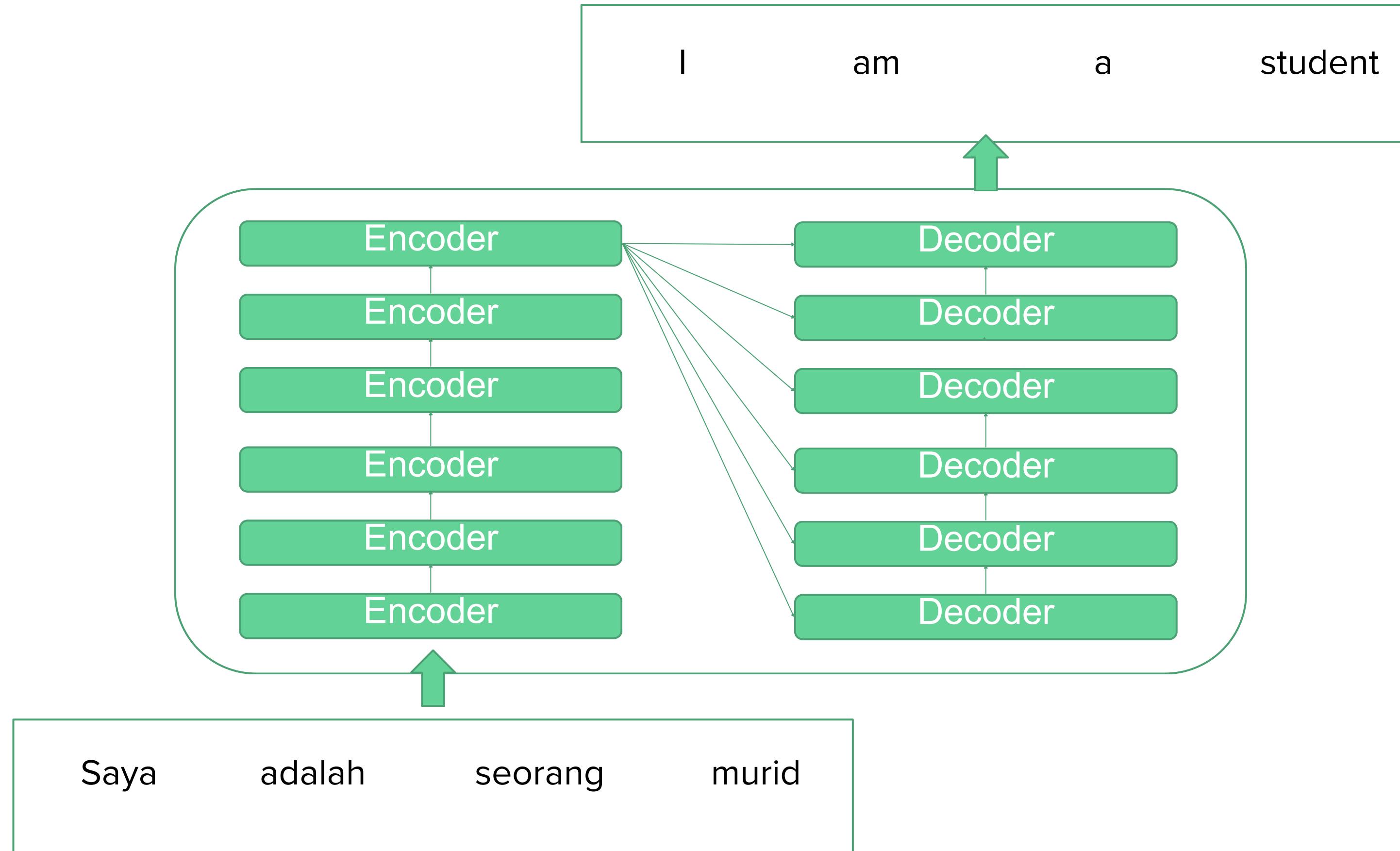
— Any Question Guys ~

— Inside The Transformers

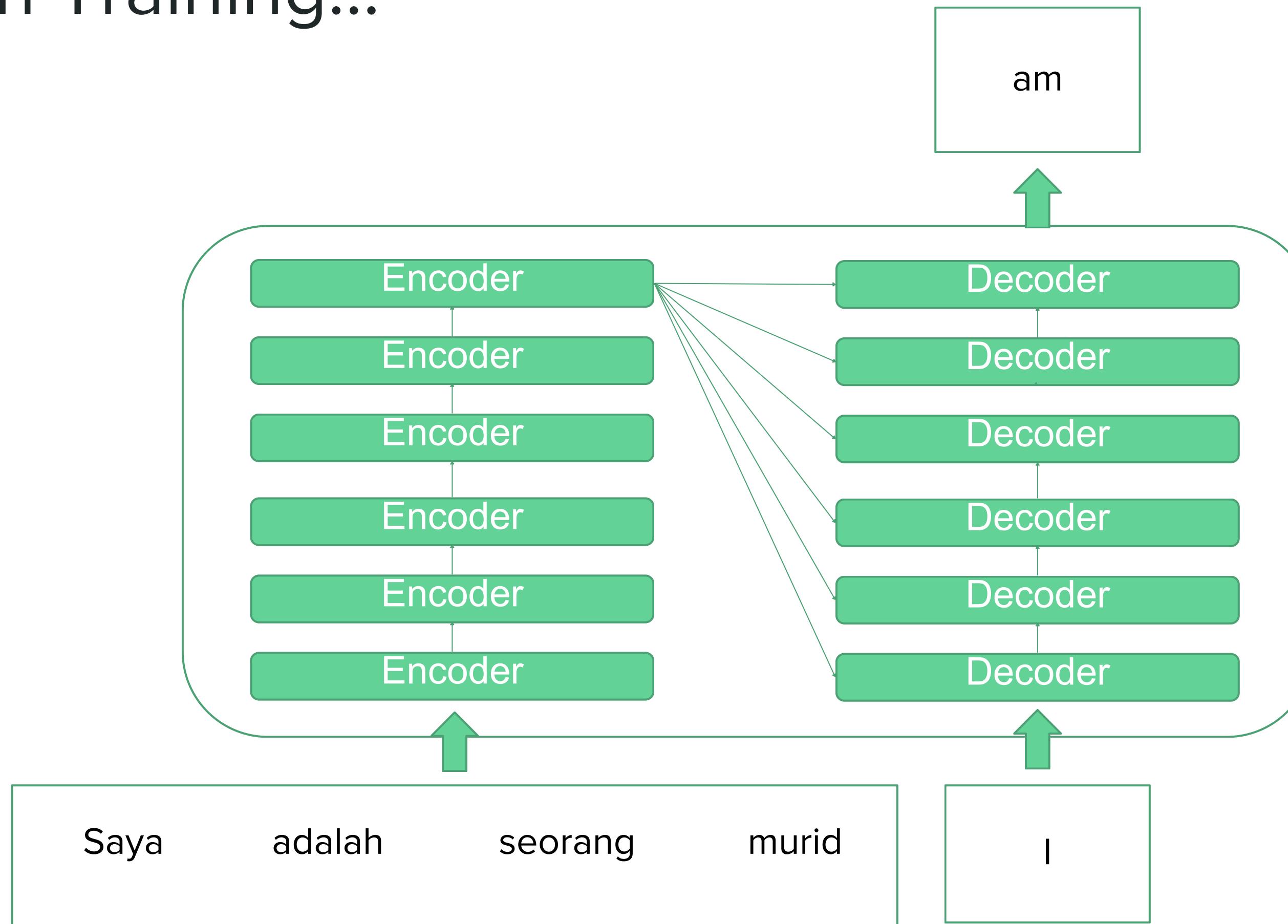
Inside Transformers



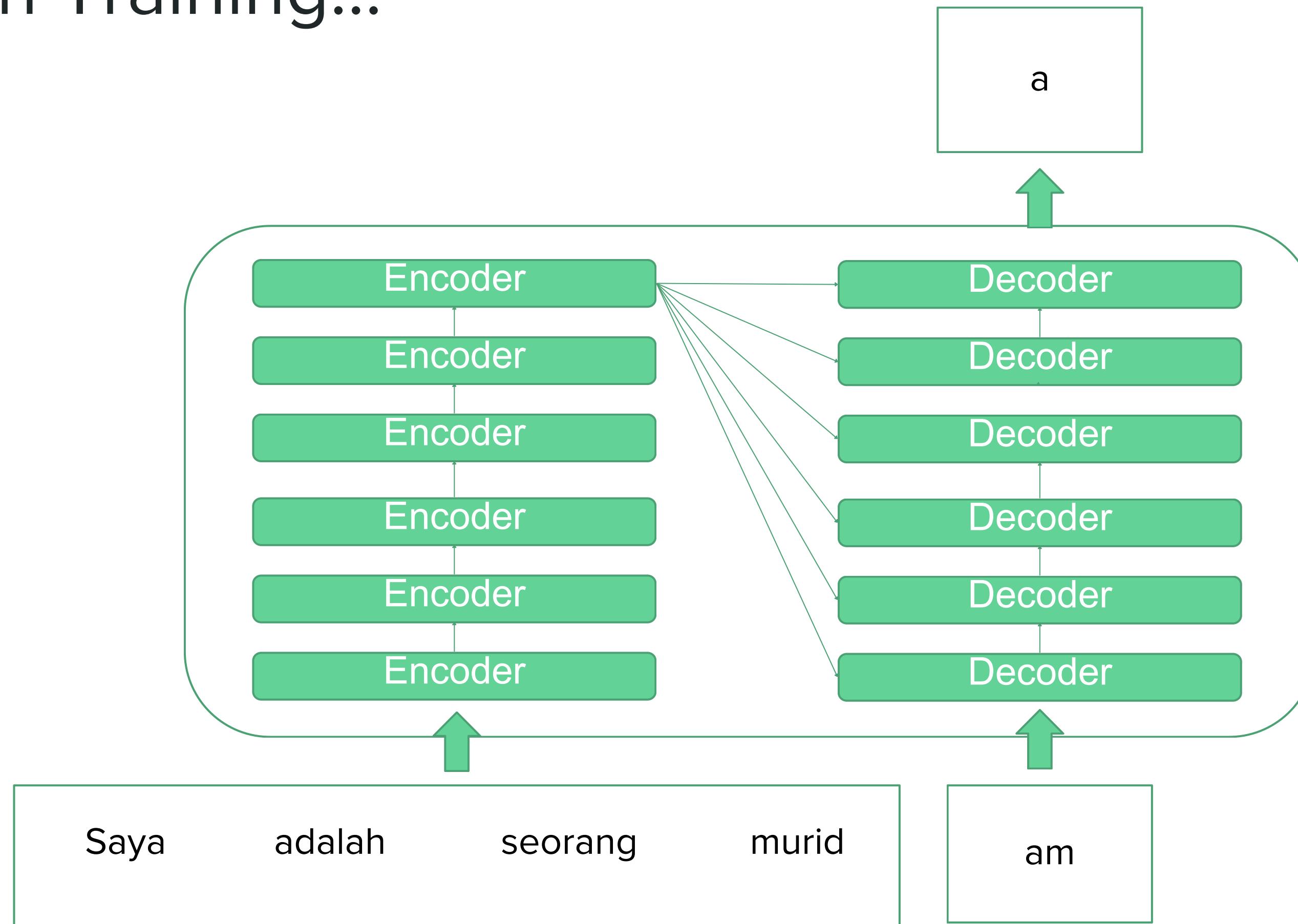
Inside Transformers



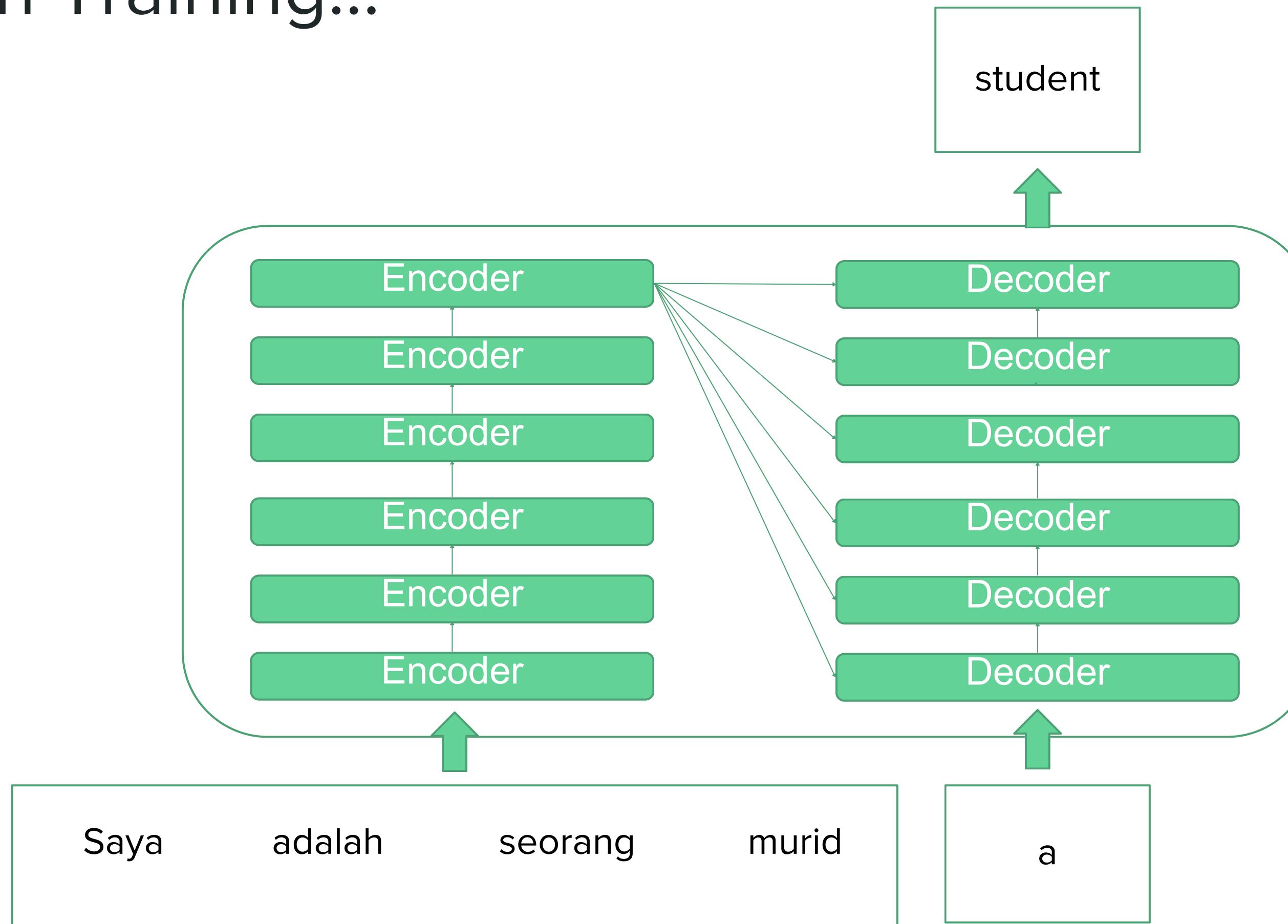
When Training...



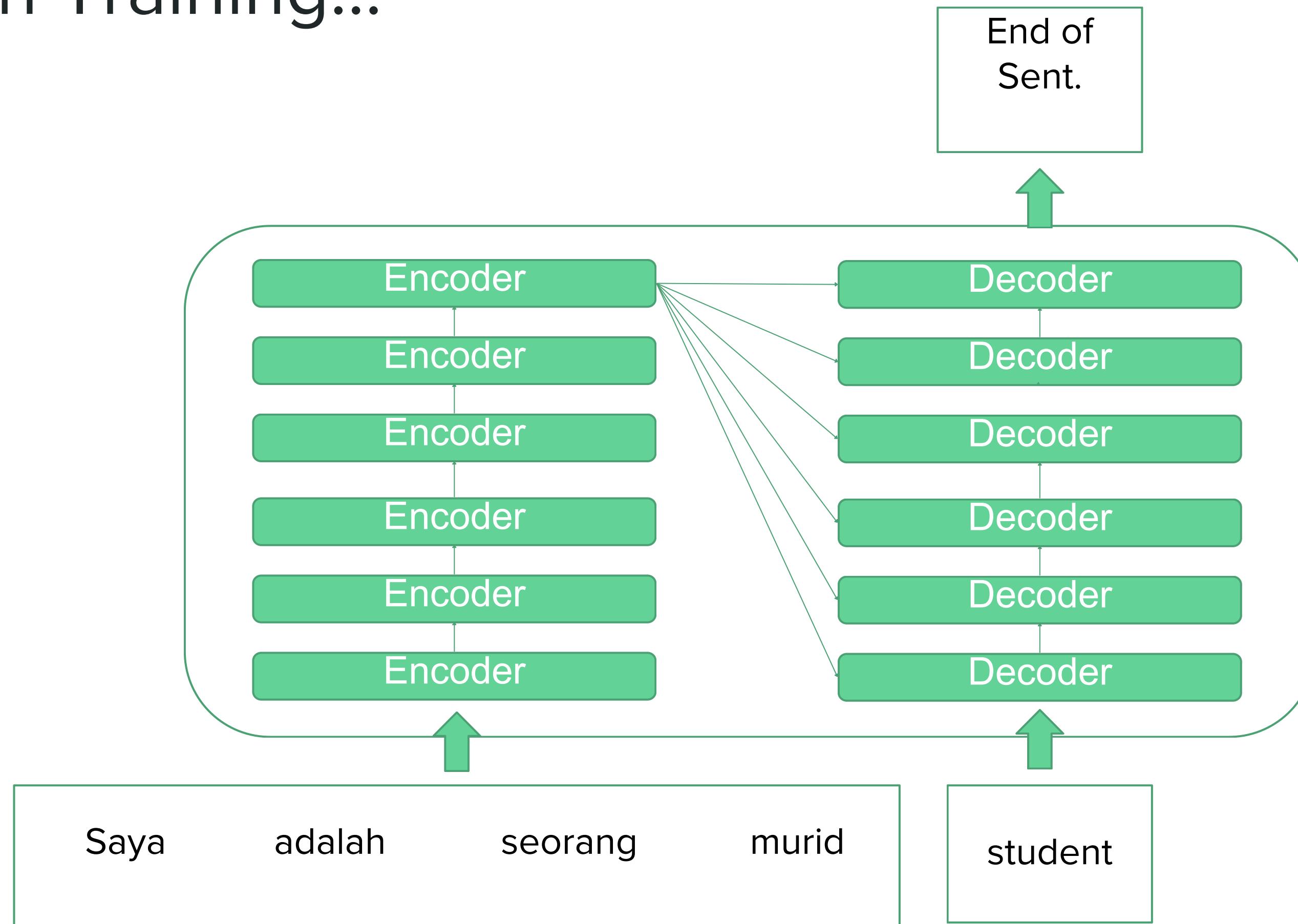
When Training...



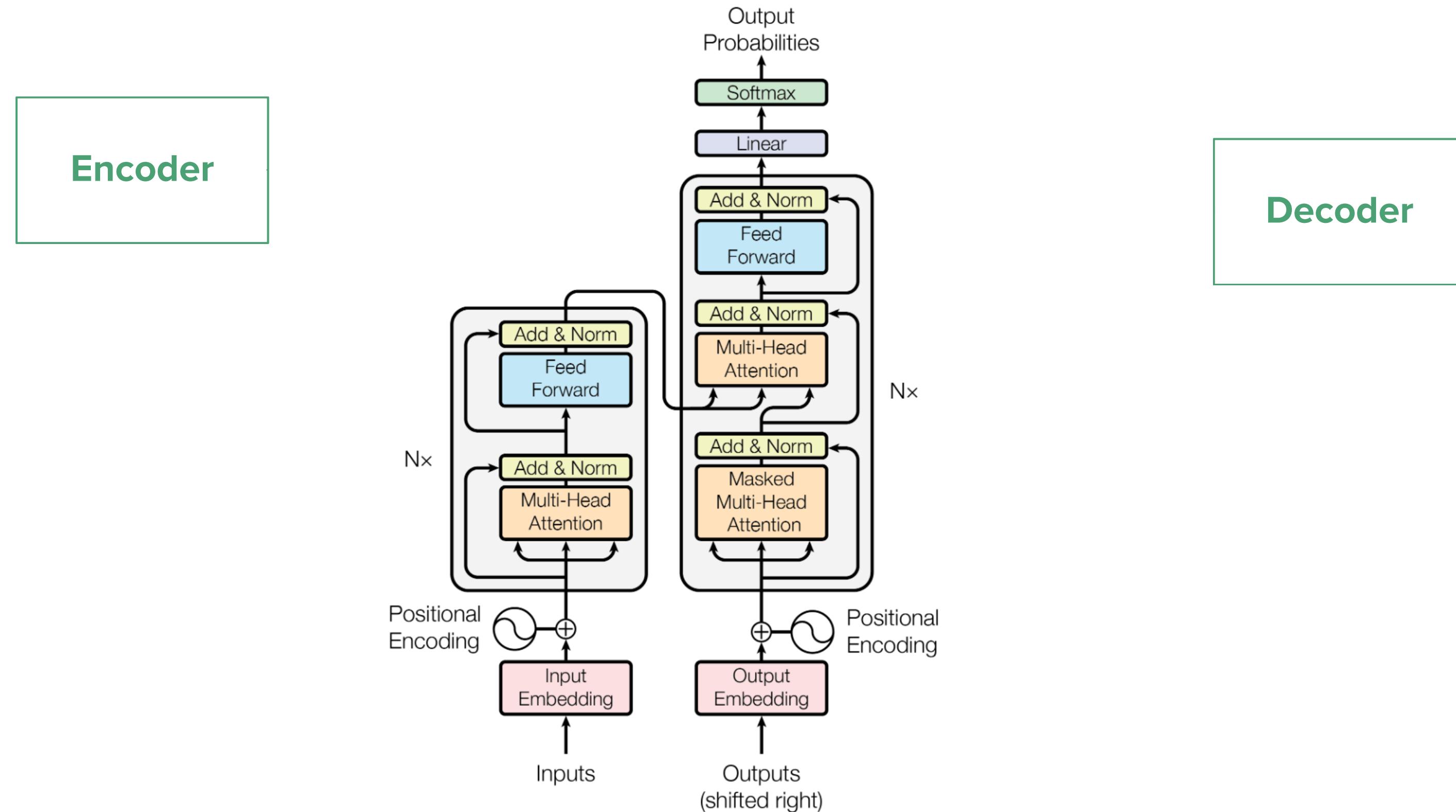
When Training...



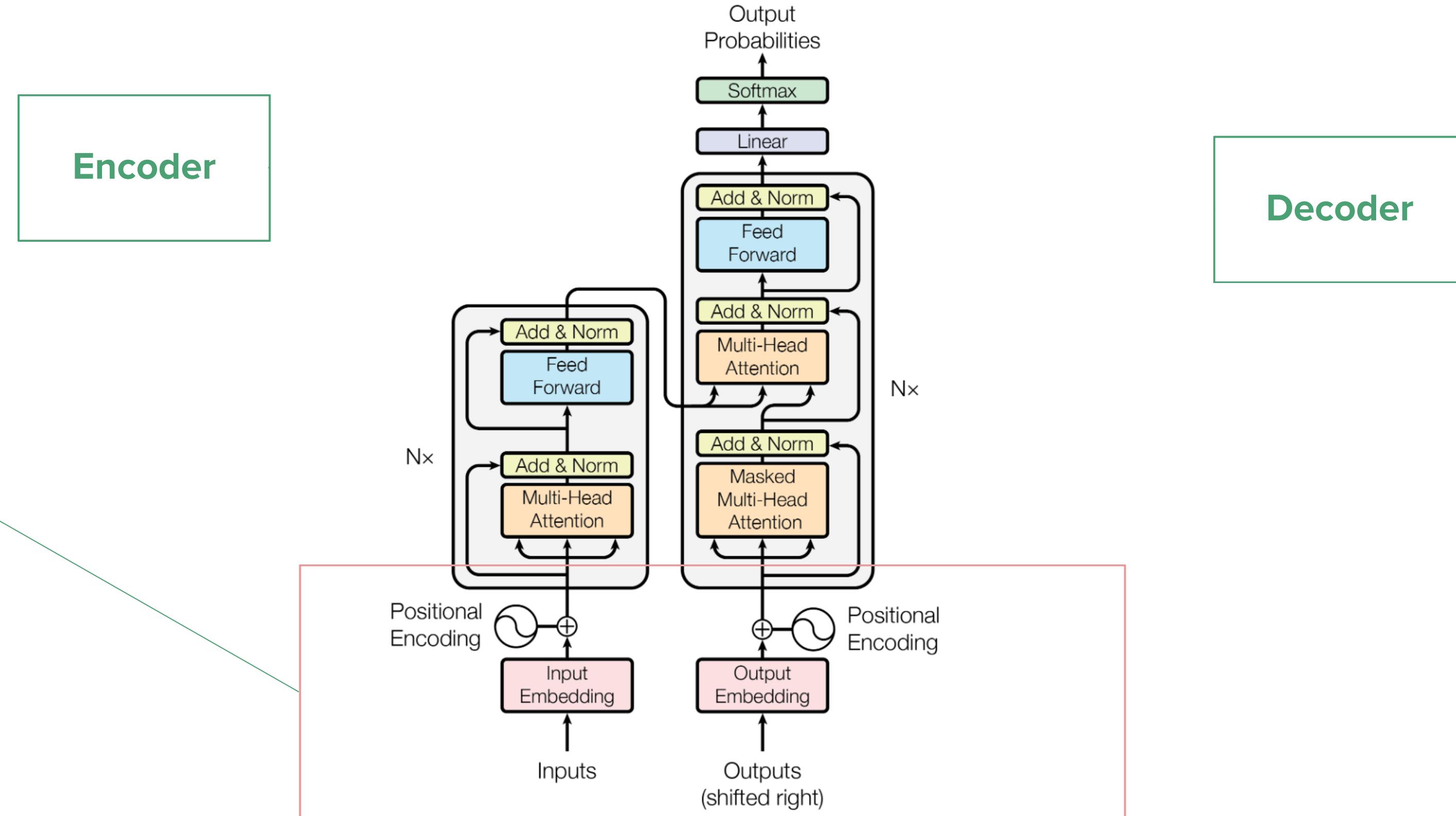
When Training...



Architecture



Architecture



— Any Question Guys ~

Transformer's Positional Embedding

Positional Embedding

Purpose:

Understanding **context** of words based on the **position** in a sentence.

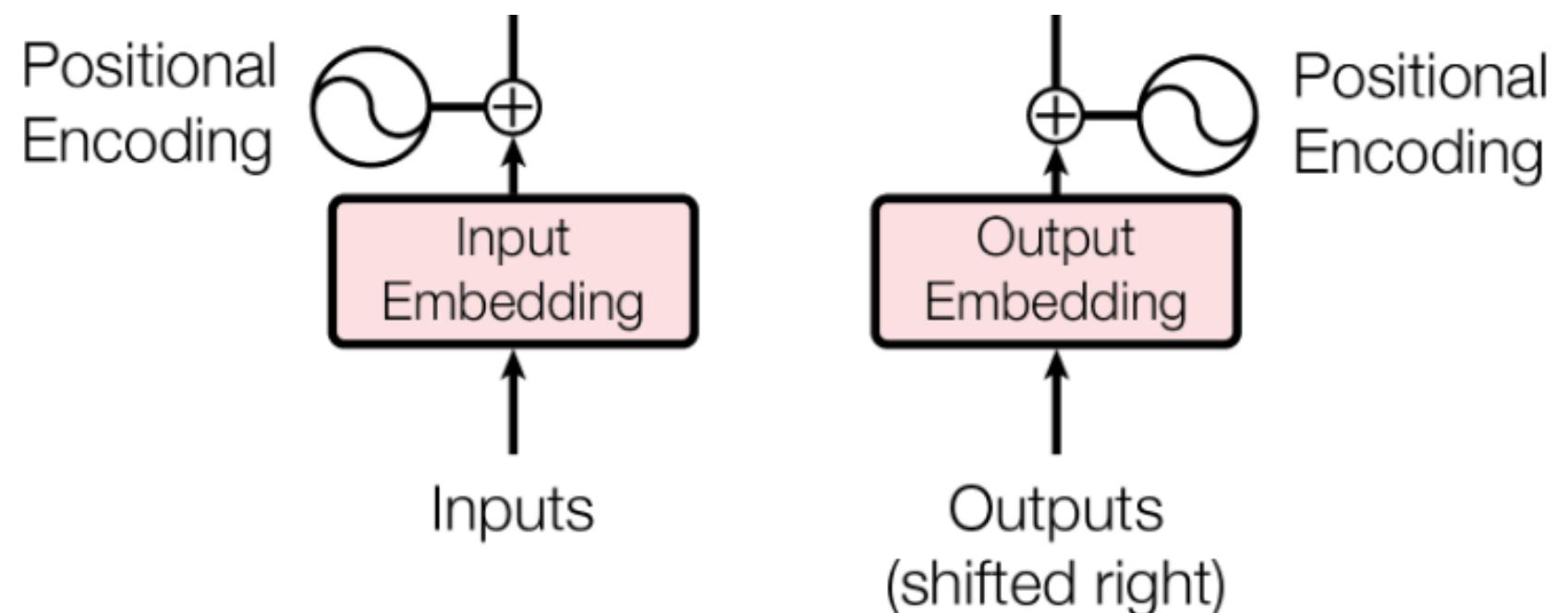
Input example:

- Pak **guru** sedang mengajar
- Fhadli adalah seorang **guru**

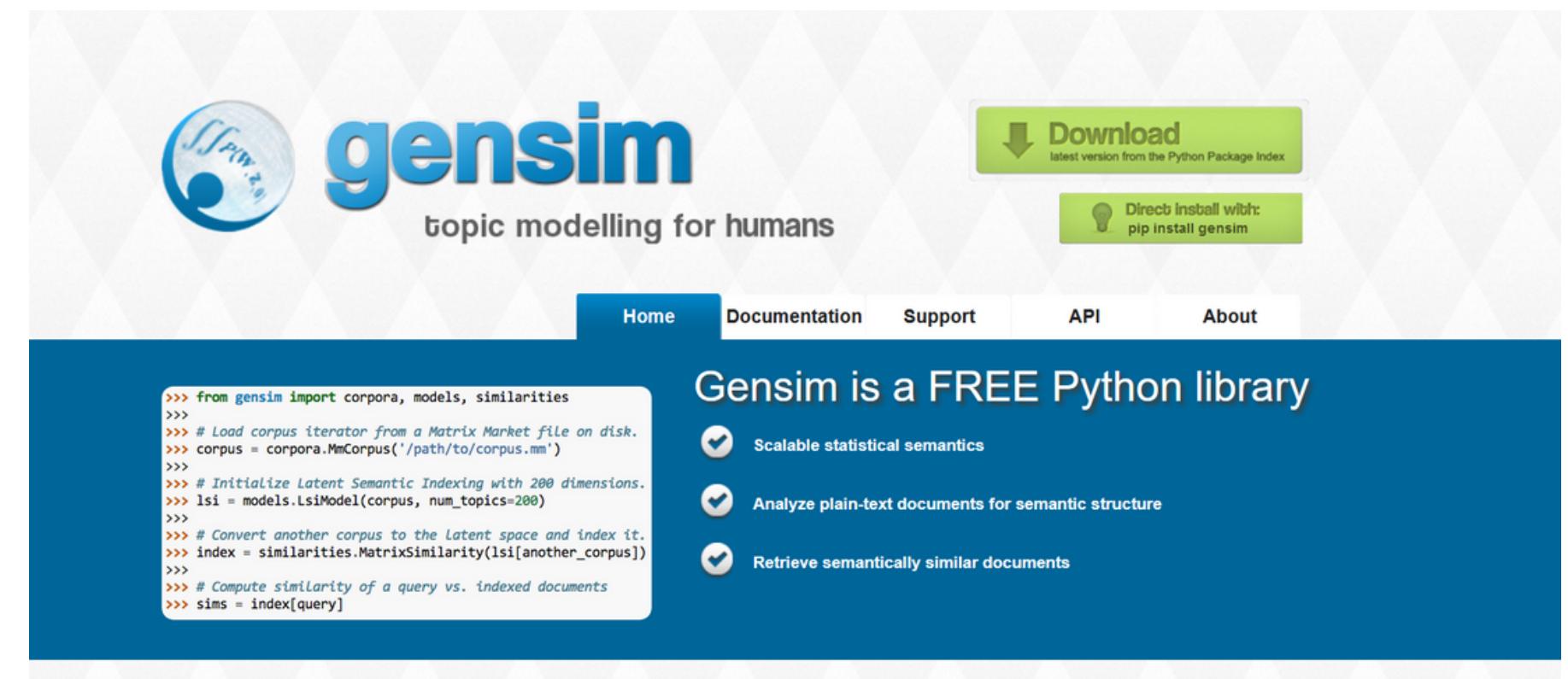
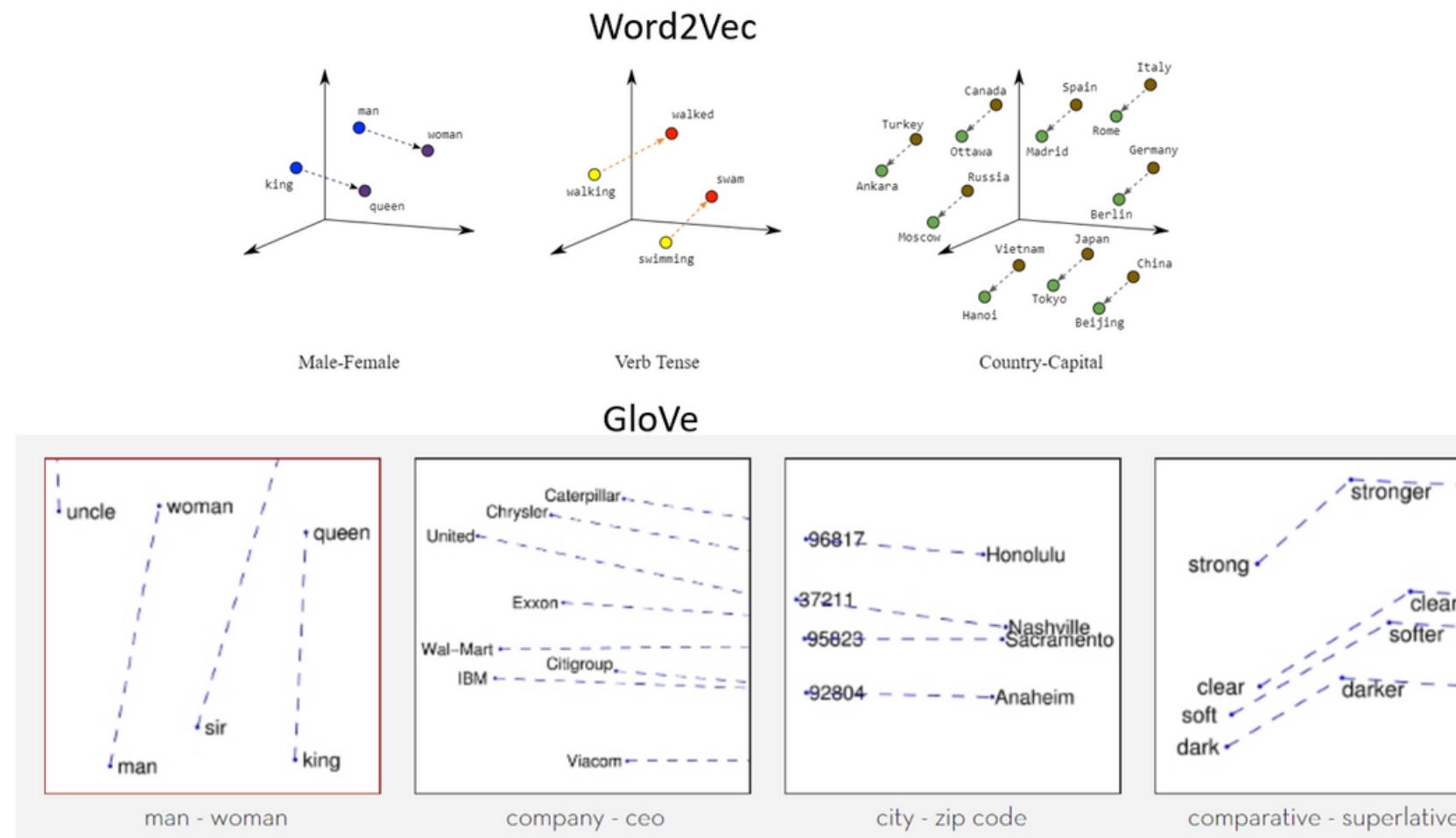
Formula:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

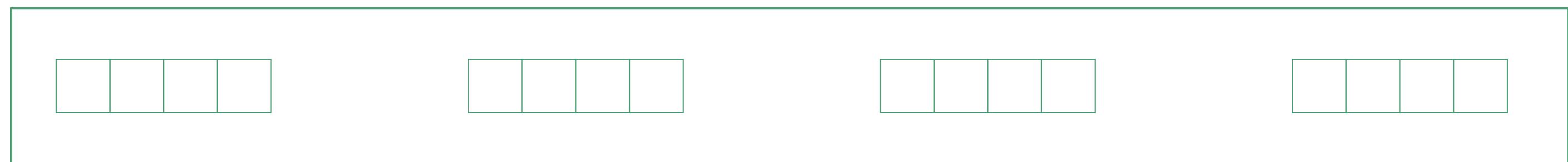
$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



Positional Embedding



Word Embedding



Saya

adalah

seorang

murid

Positional Embedding

Can be changed

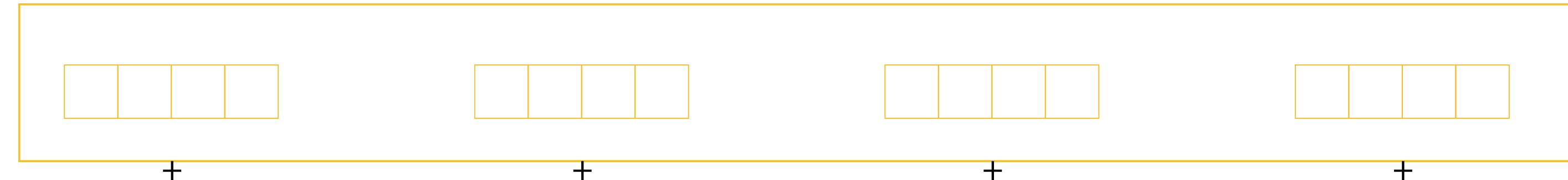
$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

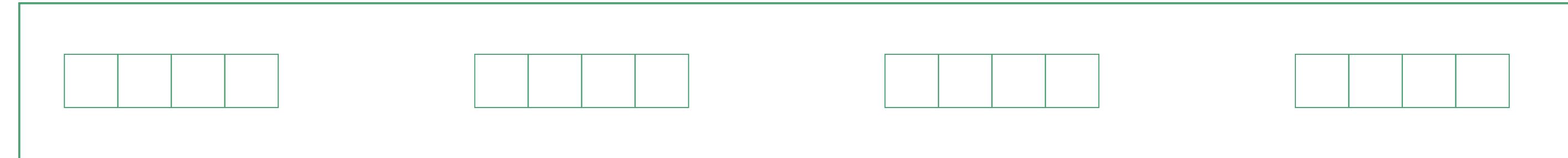
Original paper use this equation to generate positional embedding

model will easily learn to attend by relative positions rather than only 0 & 1

Positional Embedding



Word Embedding



Saya

adalah

seorang

murid

Positional Embedding

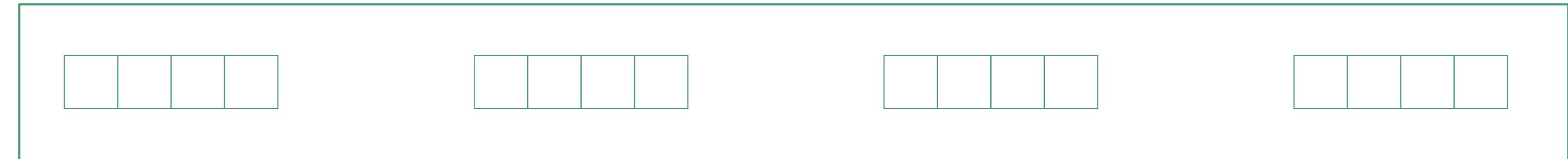
Real Embedding for
Input



Positional Embedding



Word Embedding



Saya

adalah

seorang

murid

Positional Embedding

Real Embedding for
Input

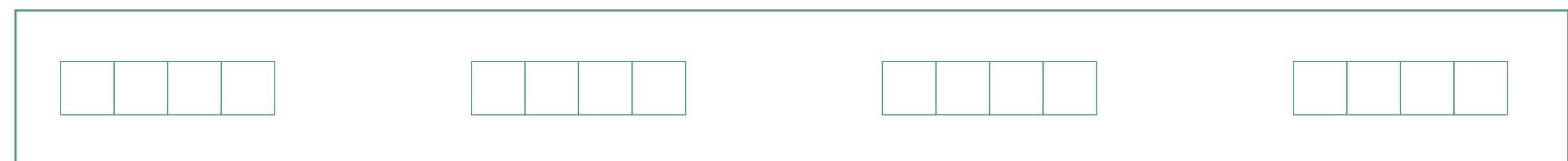


Positional Embedding



Why summation?

Word Embedding



Saya

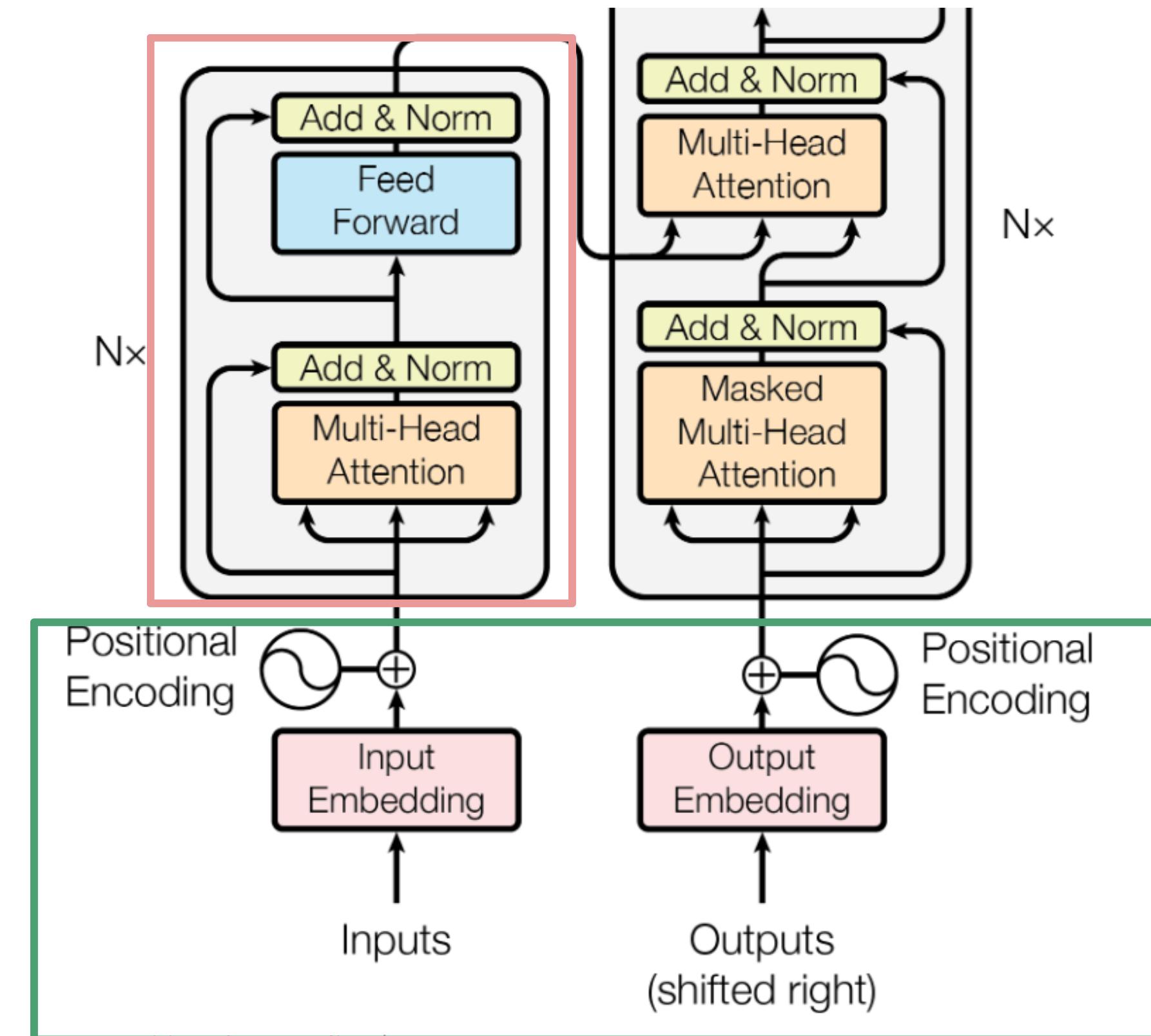
adalah

seorang

murid

Positional Embedding

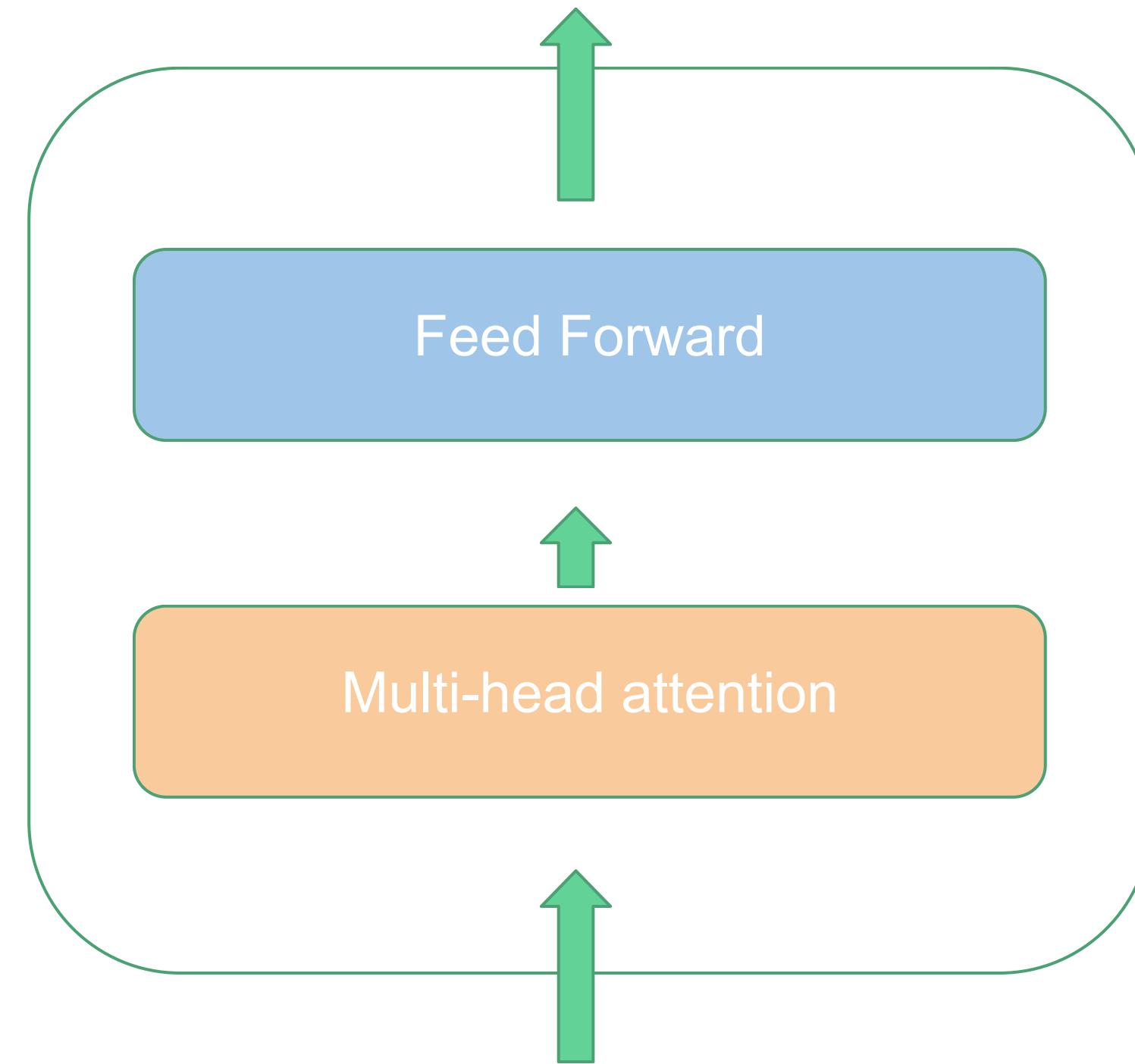
Let's talk about the encoder



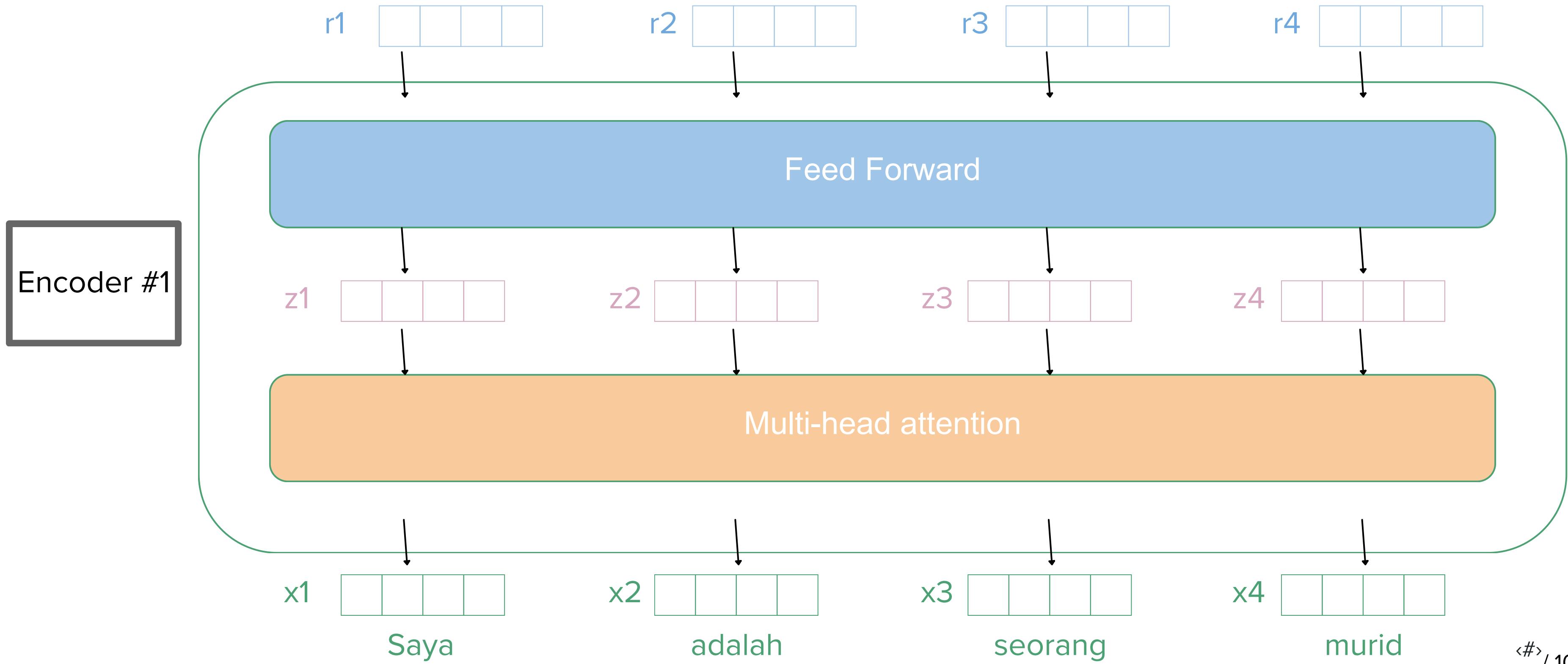
— Any Question Guys ~

— Transformer's Encoder

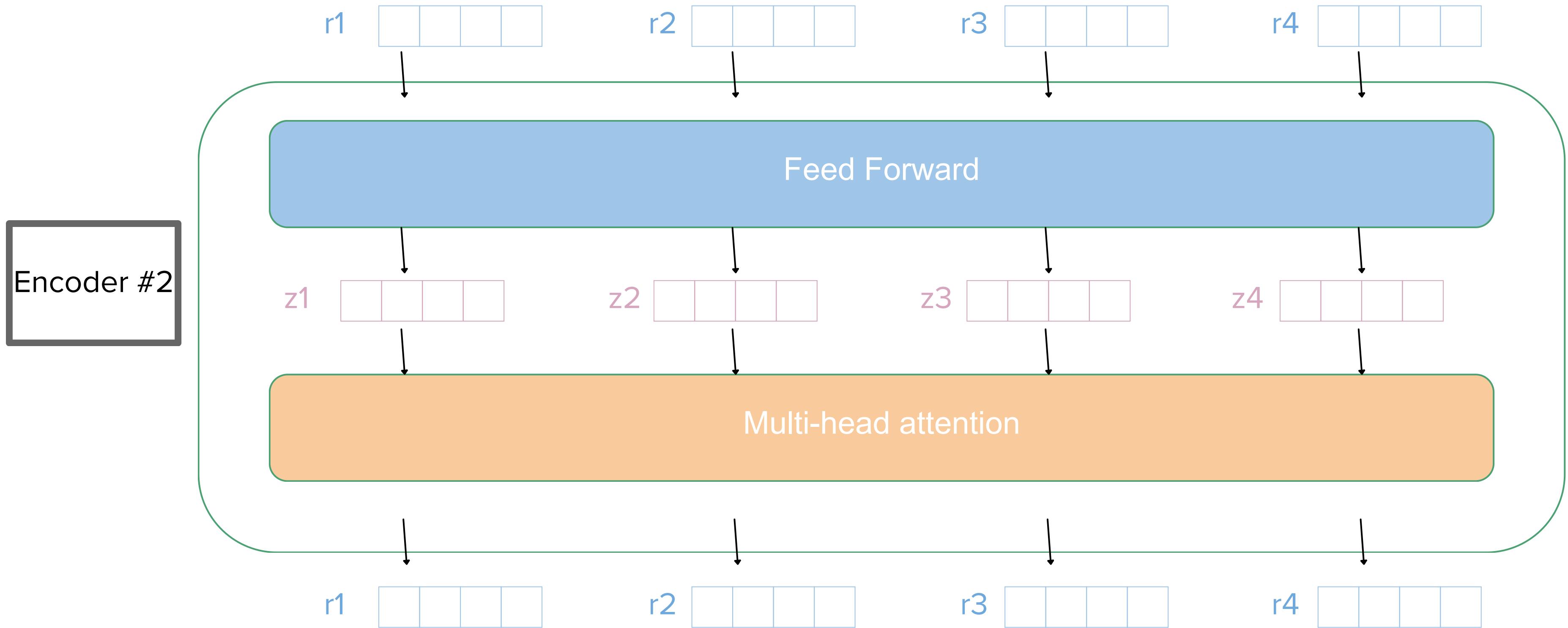
Encoder



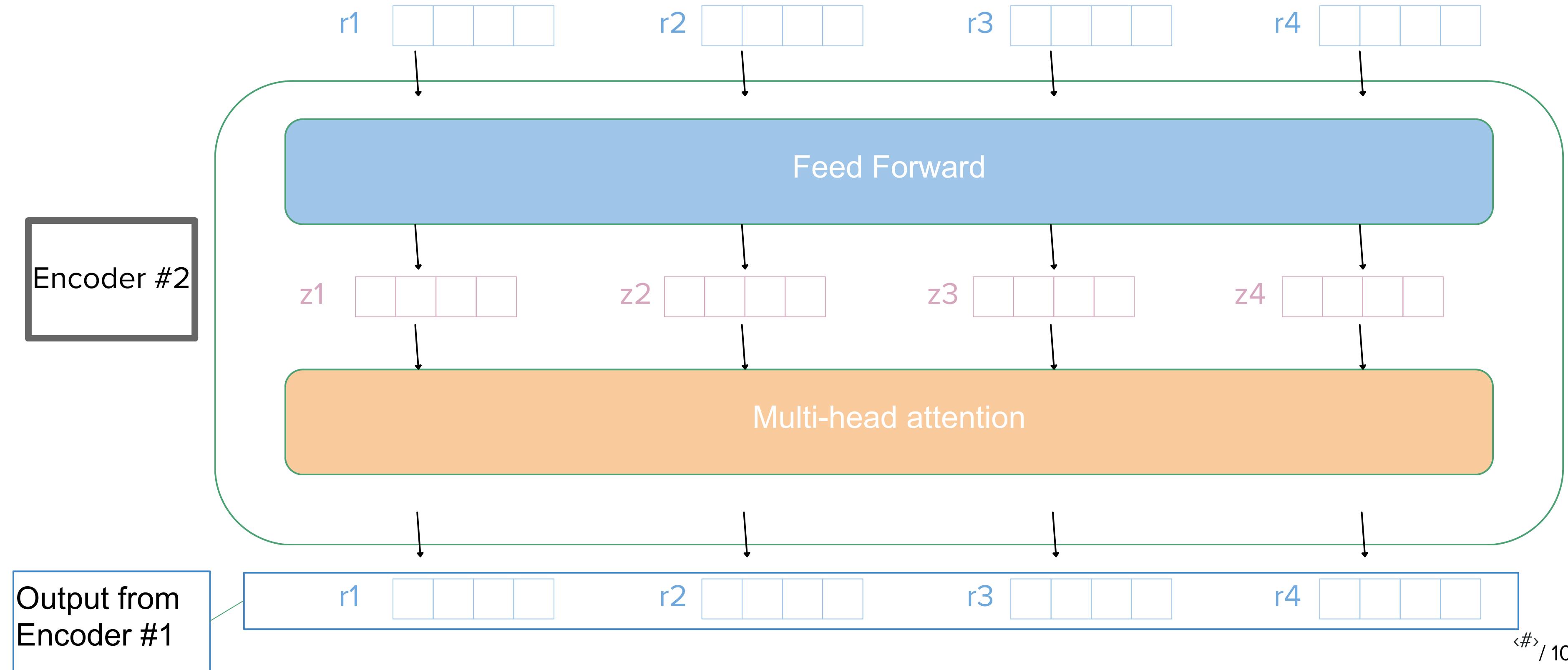
Encoder



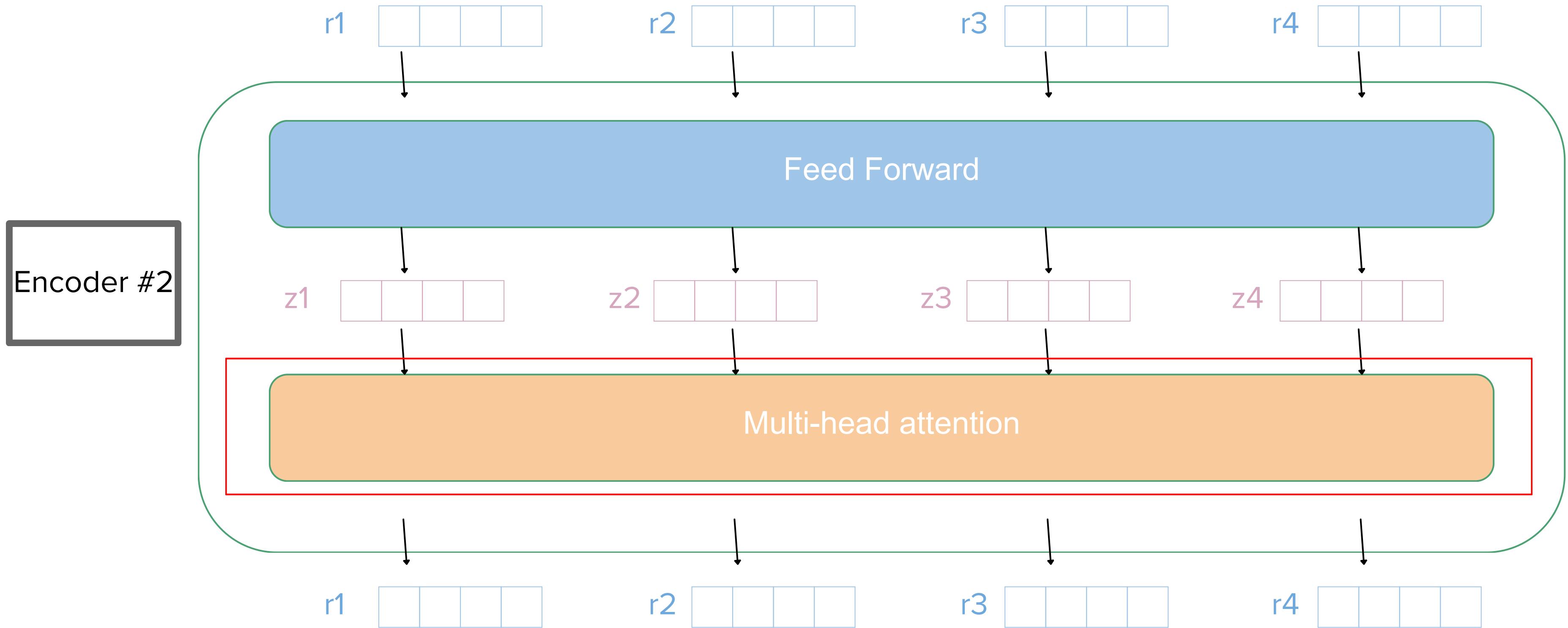
Encoder



Encoder



Encoder



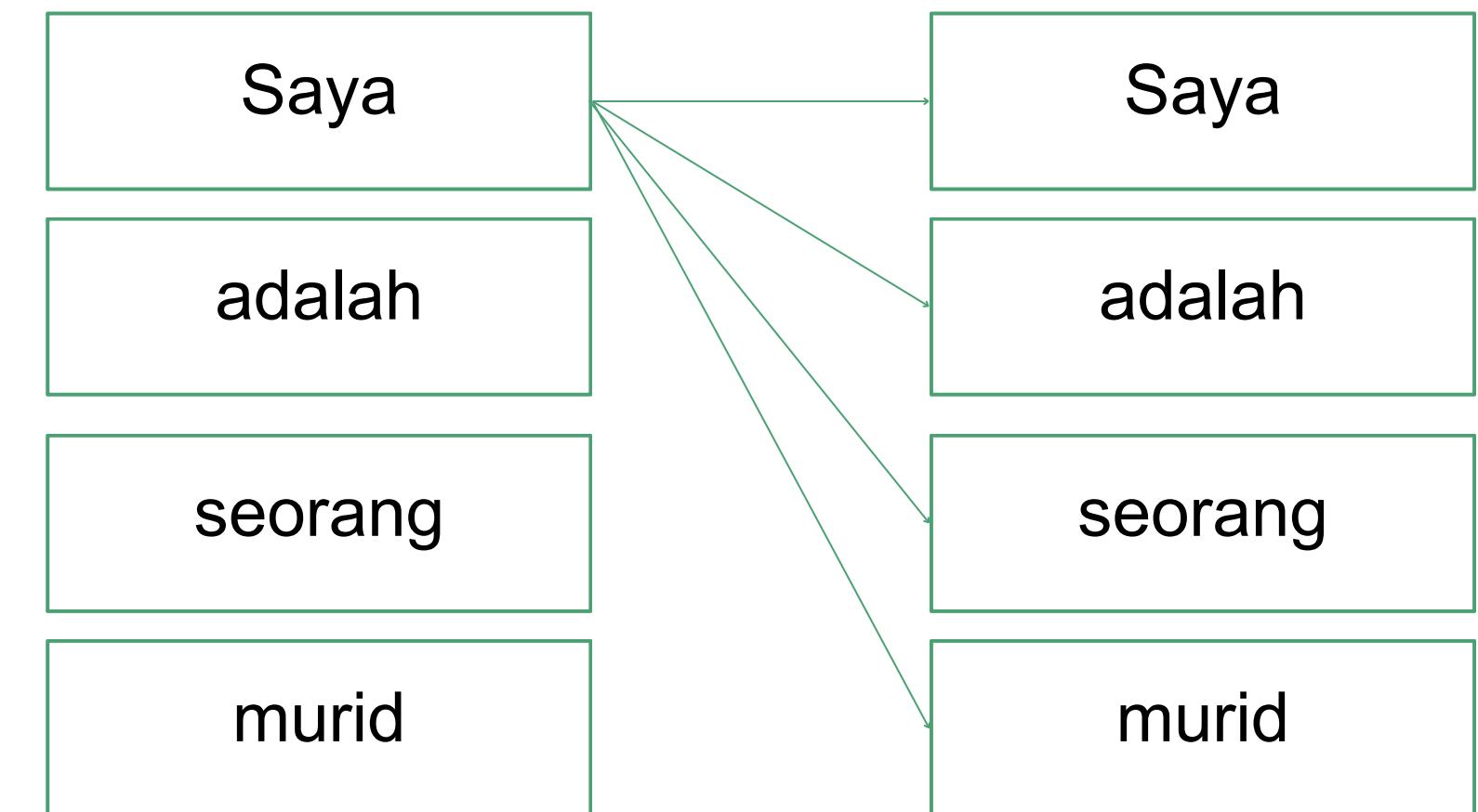
Encoder: Multi-head attention

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words in a sentence

Example:

Ari melihat pak **guru** sedang bersama anak **beliau**



It looks into the relationship of a word in a **sentence** and every words in that **sentence itself**. That's why it's also called **self-attention**

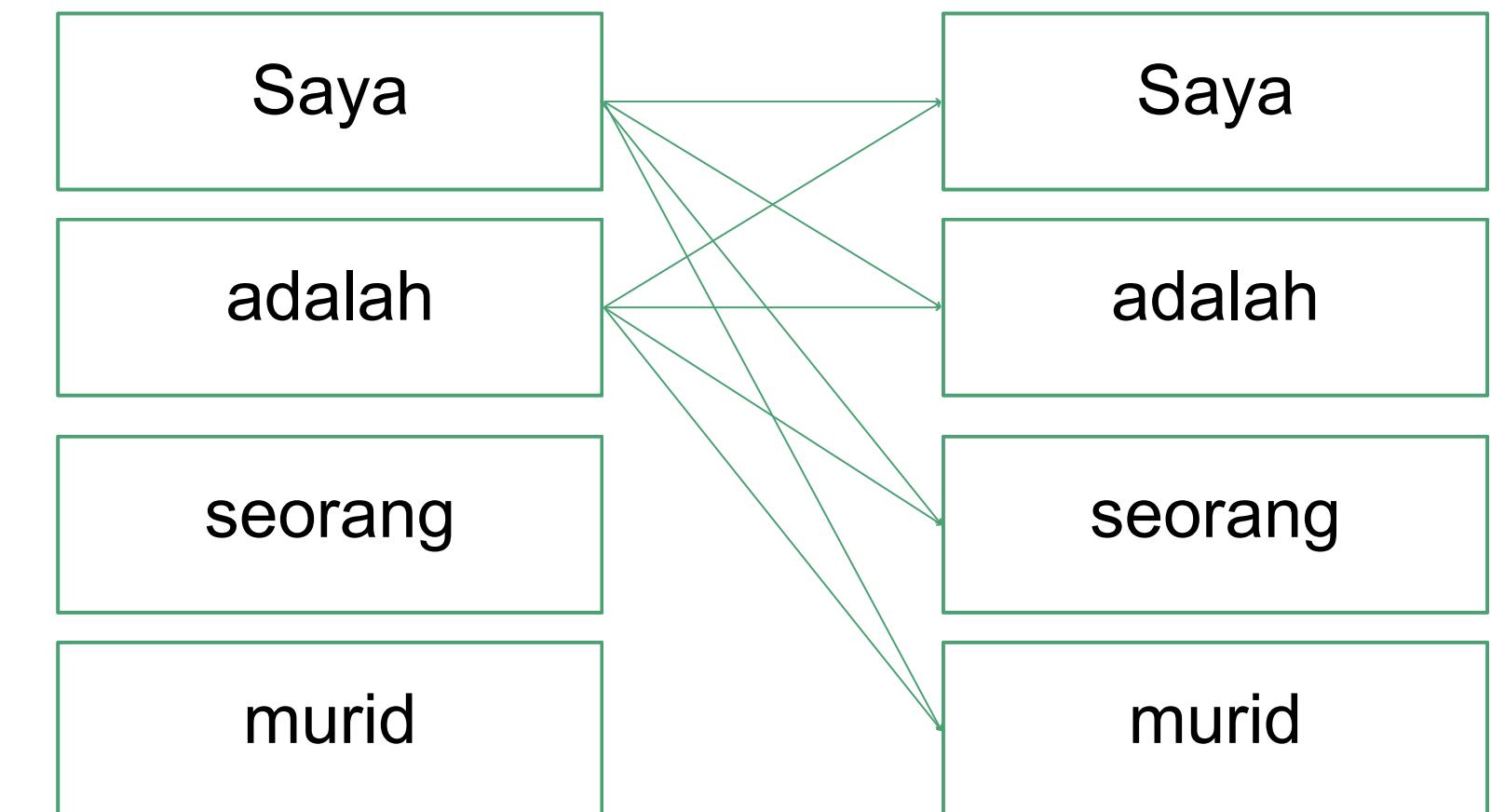
Encoder: Multi-head attention

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words in a sentence

Example

Ari melihat pak **guru** sedang bersama anak **beliau**



This is done between every words in the sentence. That's why it's also called **multi-headed self-attention**

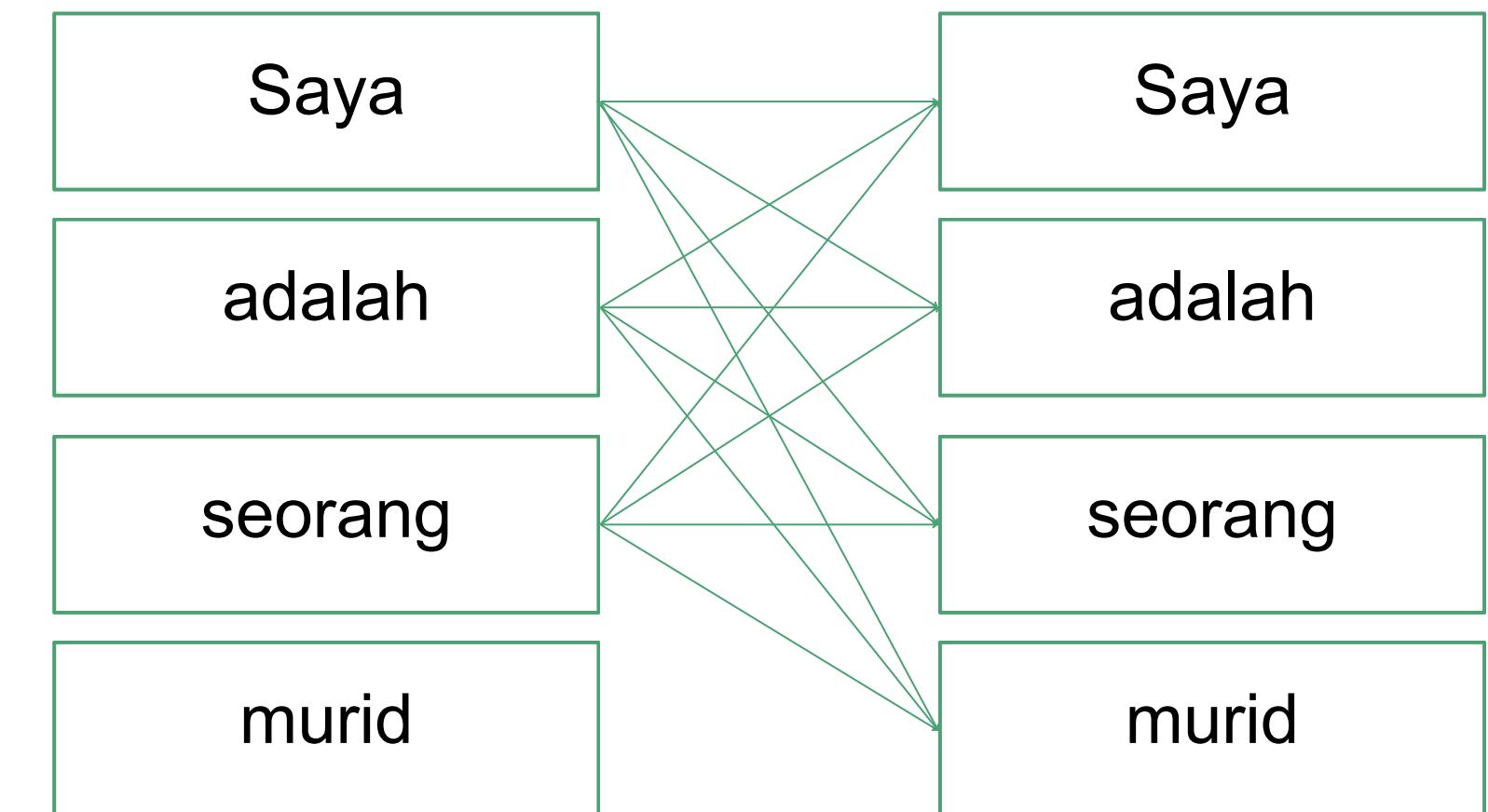
Encoder: Multi-head attention

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words in a sentence

Example

Ari melihat pak **guru** sedang bersama anak **beliau**



This is done between every words in the sentence. That's why it's also called **multi-headed self-attention**

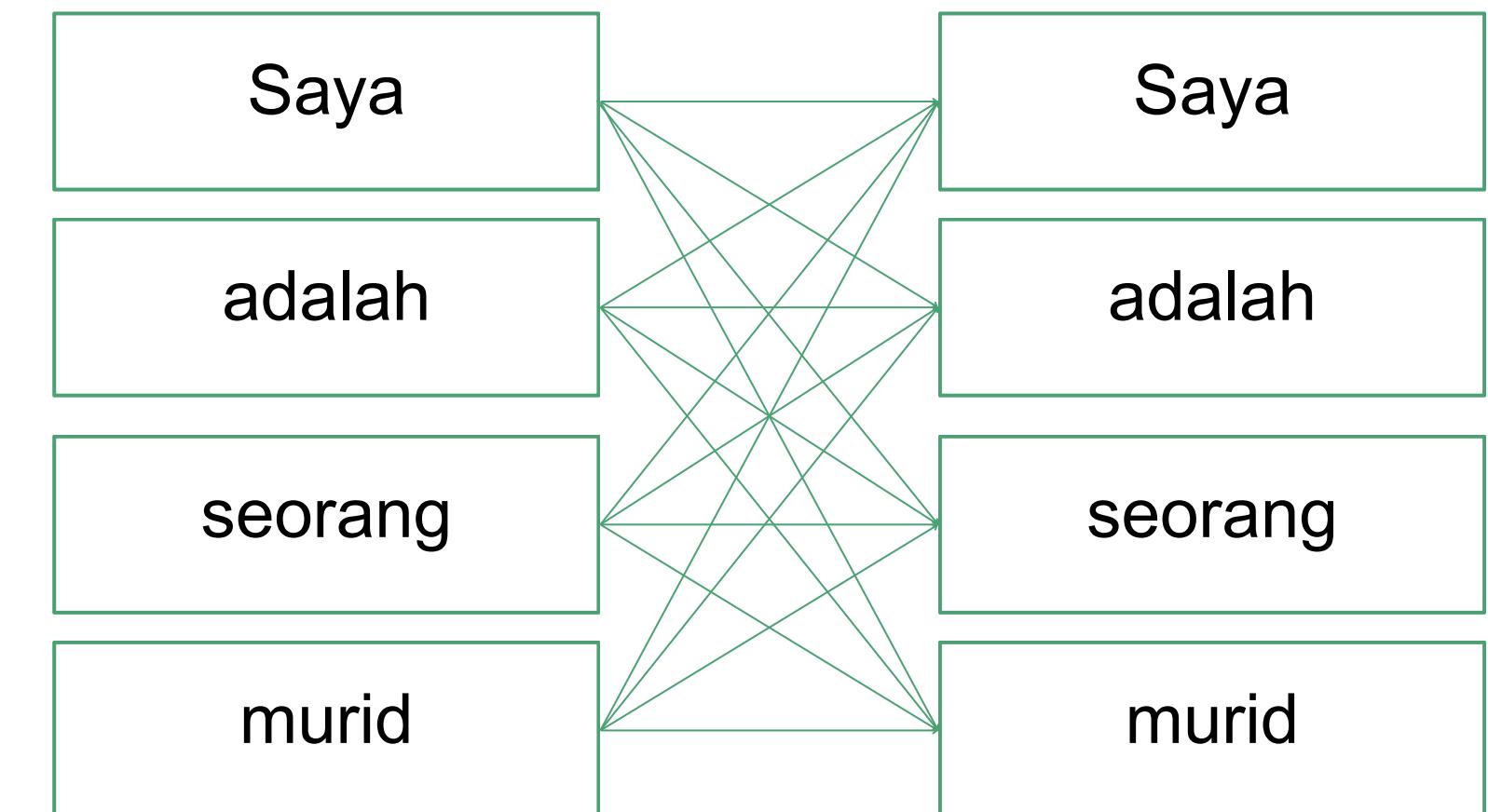
Encoder: Multi-head attention

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words in a sentence

Example

Ari melihat pak **guru** sedang bersama anak **beliau**

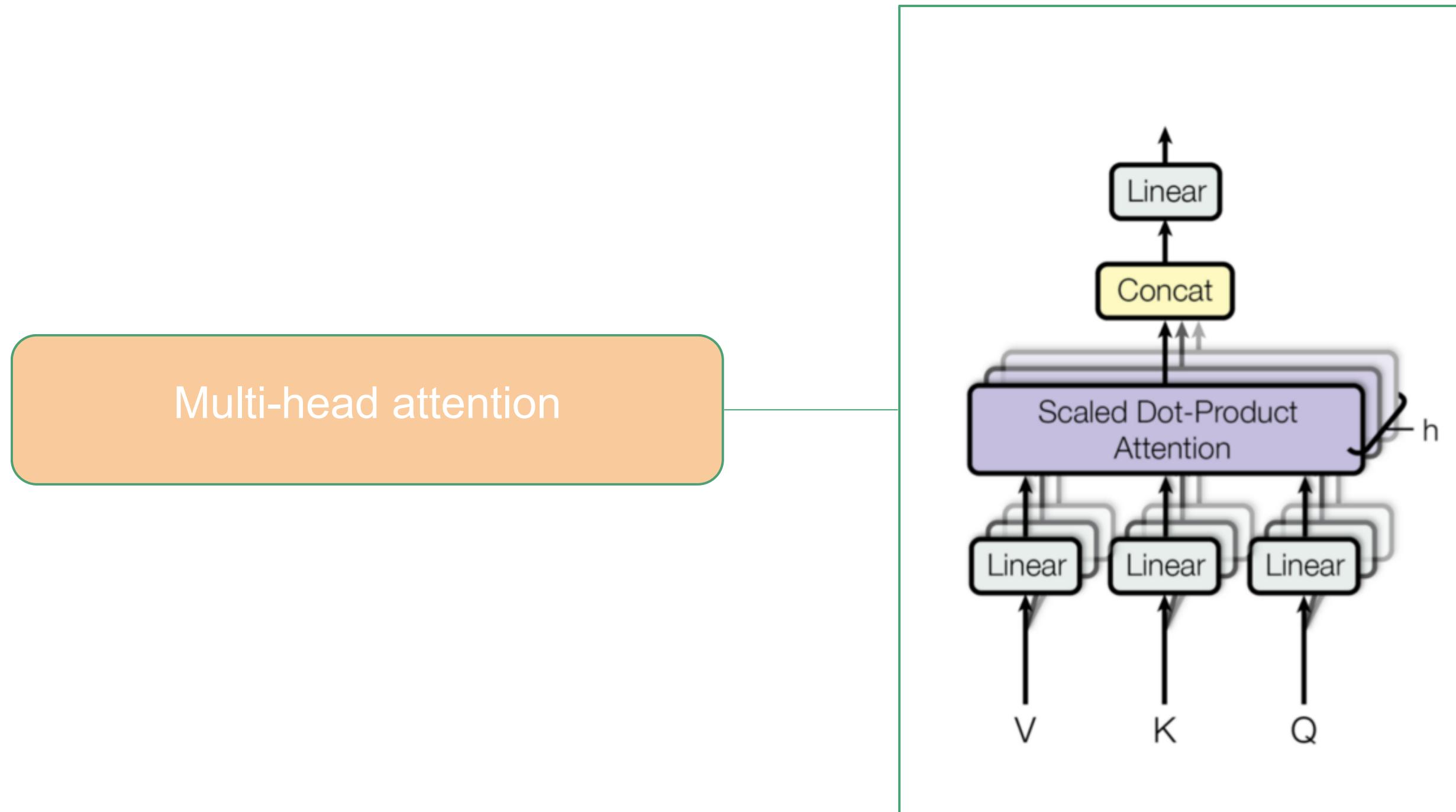


This is done between every words in the sentence. That's why it's also called **multi-headed self-attention**

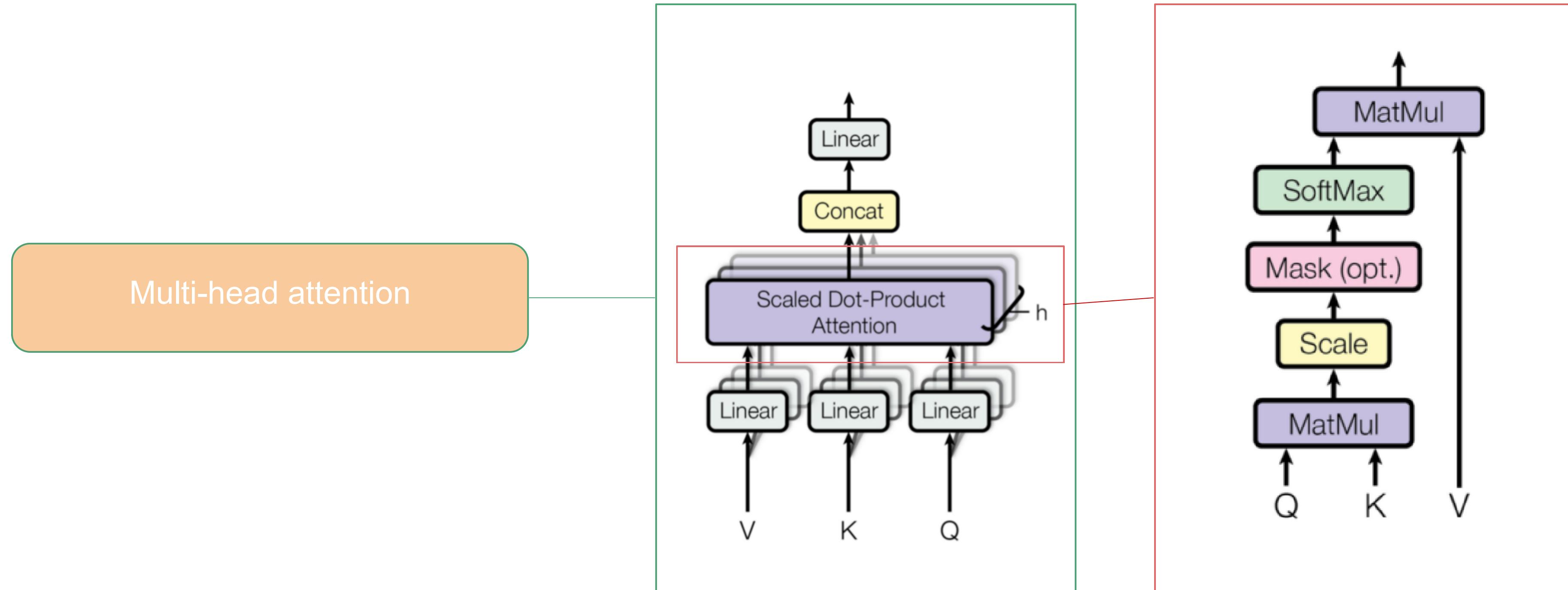
Encoder: Multi-head attention

| | | | | |
|---------|-------|------|------|-------|
| Saya | [0.73 | 0.04 | 0.05 | 0.18] |
| adalah | [0.03 | 0.84 | 0.02 | 0.11] |
| seorang | [0.10 | 0.06 | 0.62 | 0.22] |
| murid | [0.05 | 0.03 | 0.04 | 0.88] |

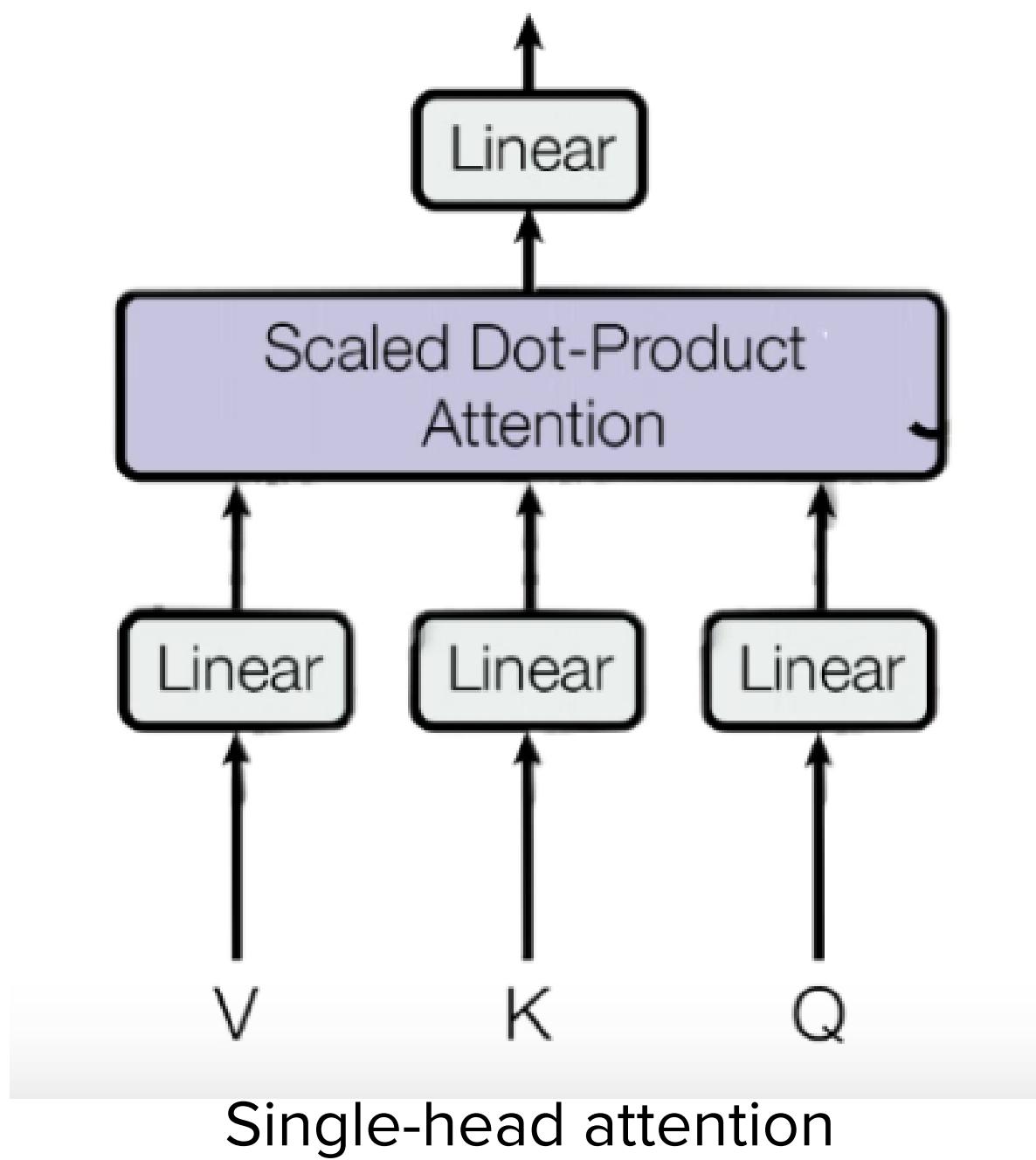
Encoder: Multi-head attention



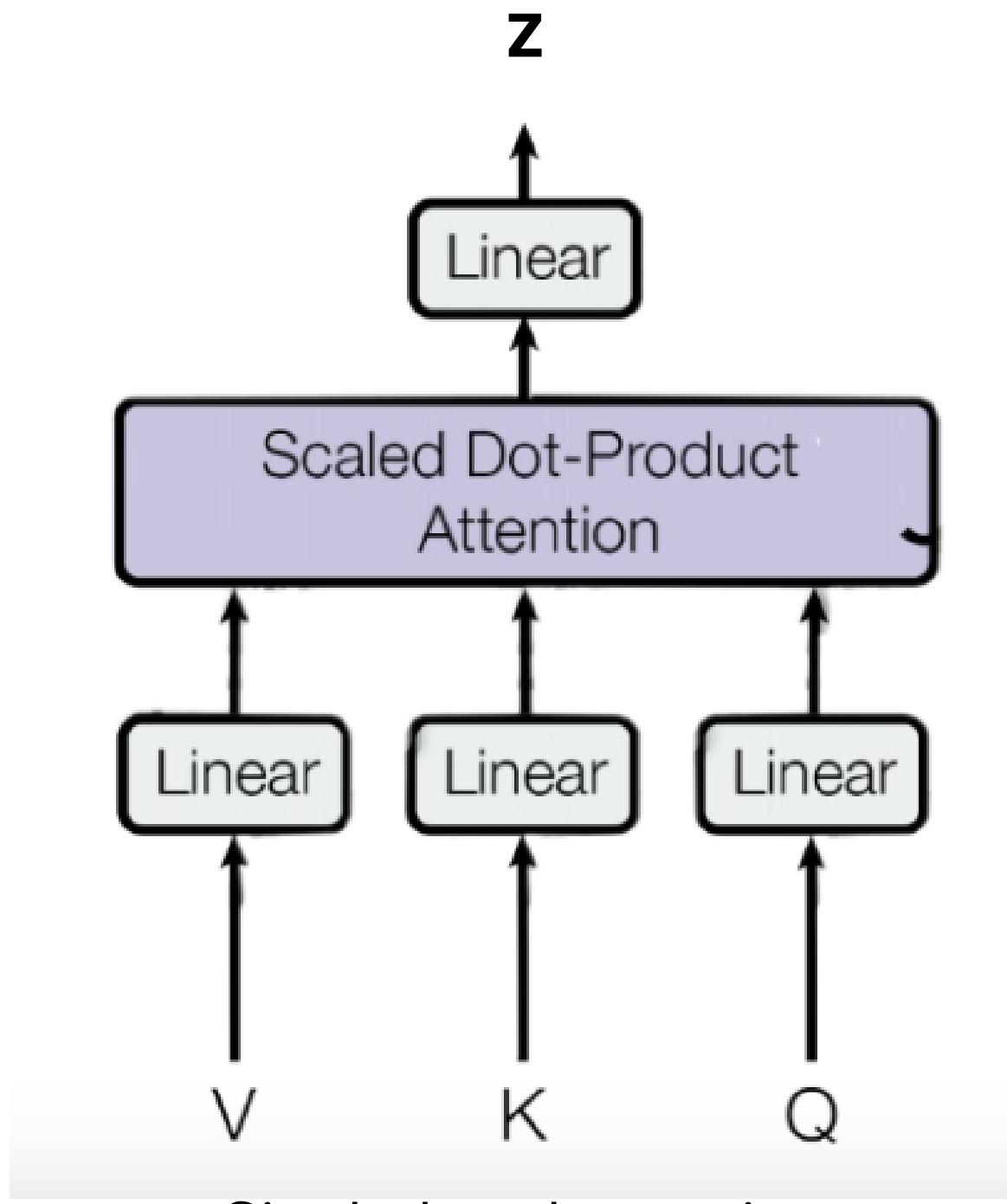
Encoder: Multi-head attention



Encoder: Multi-head attention

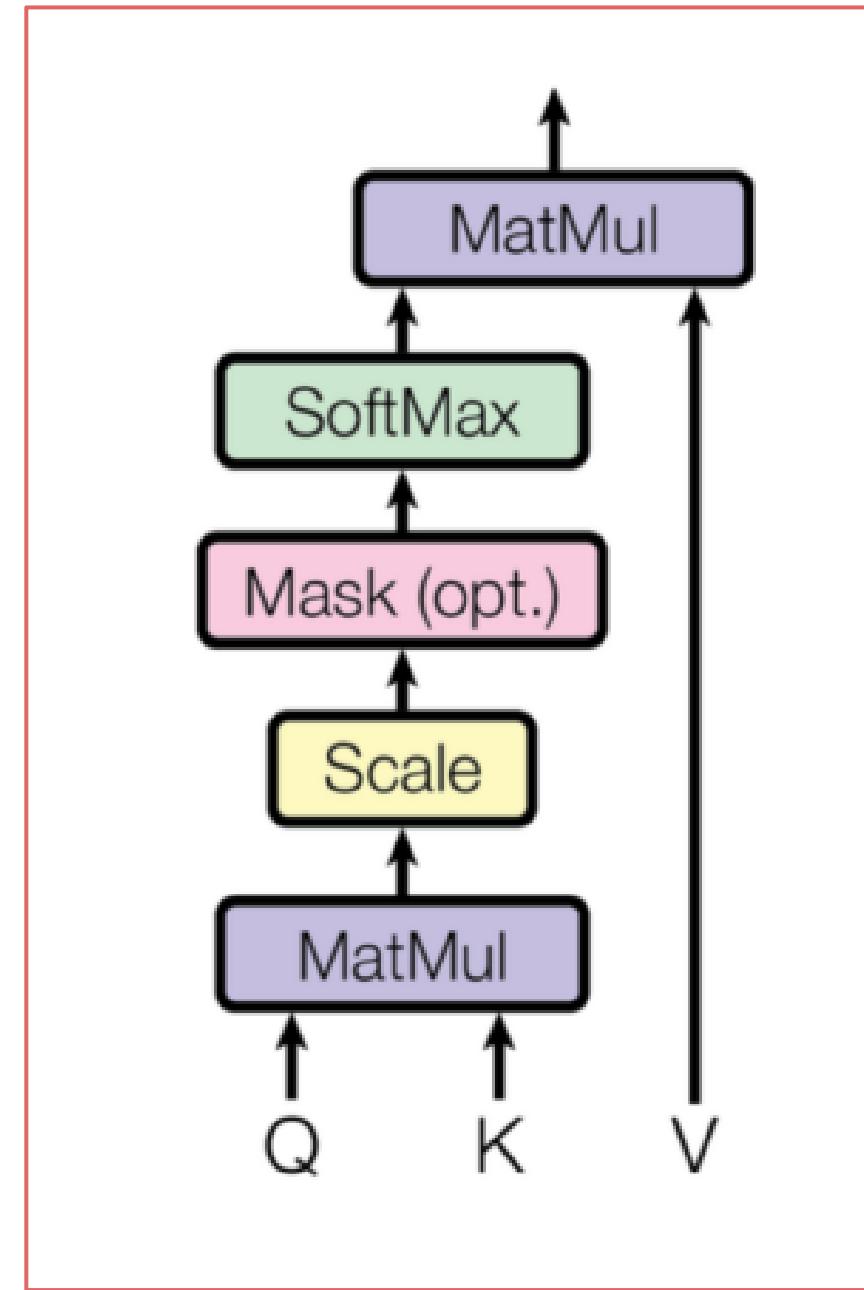
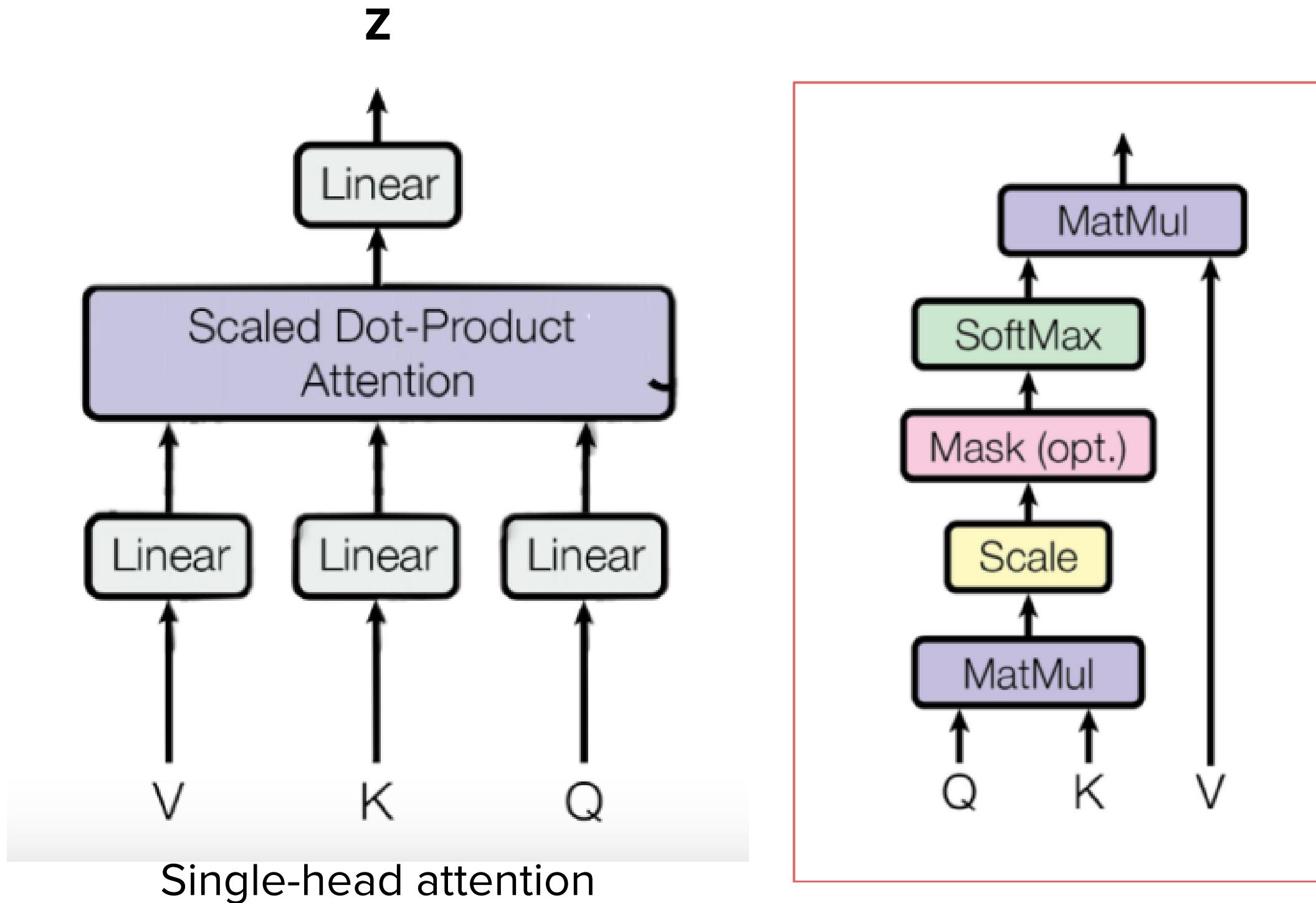


Encoder: Multi-head attention

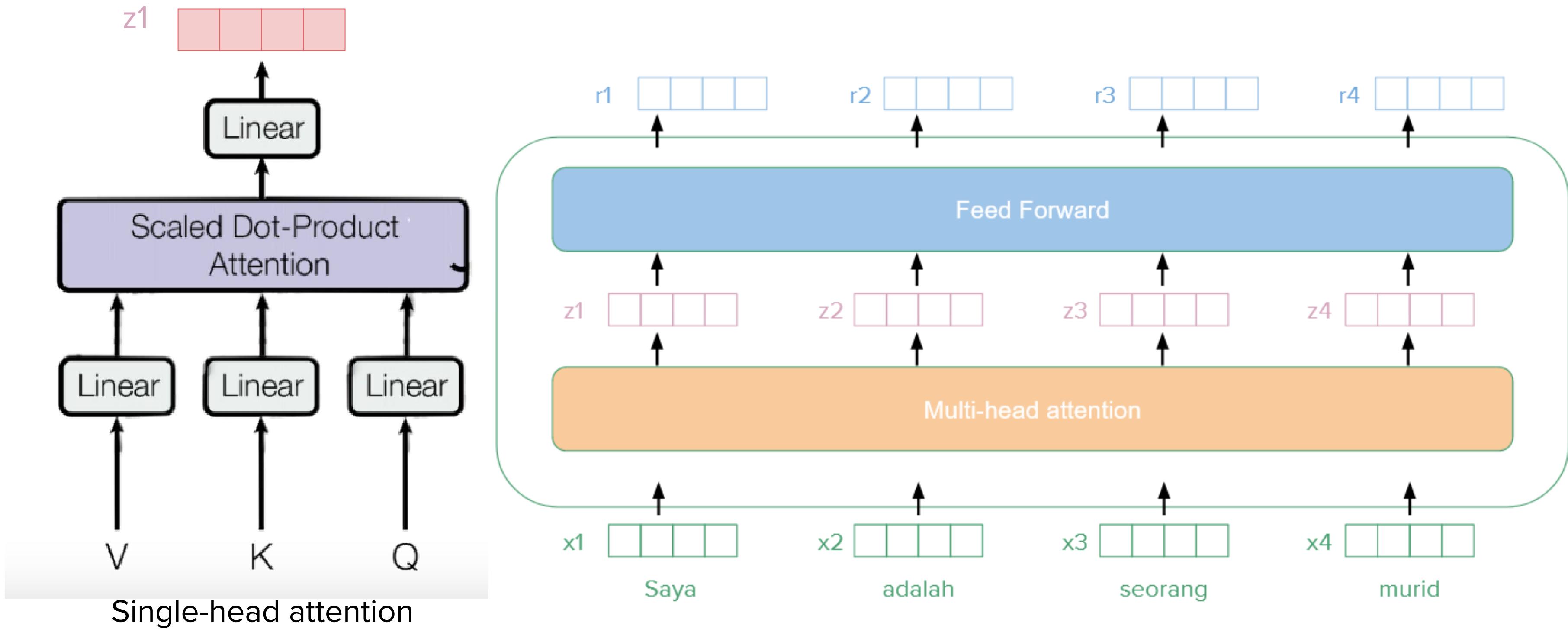


$$Z = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{\text{Dimension of vector } Q, K \text{ or } V}} \right) \cdot V$$

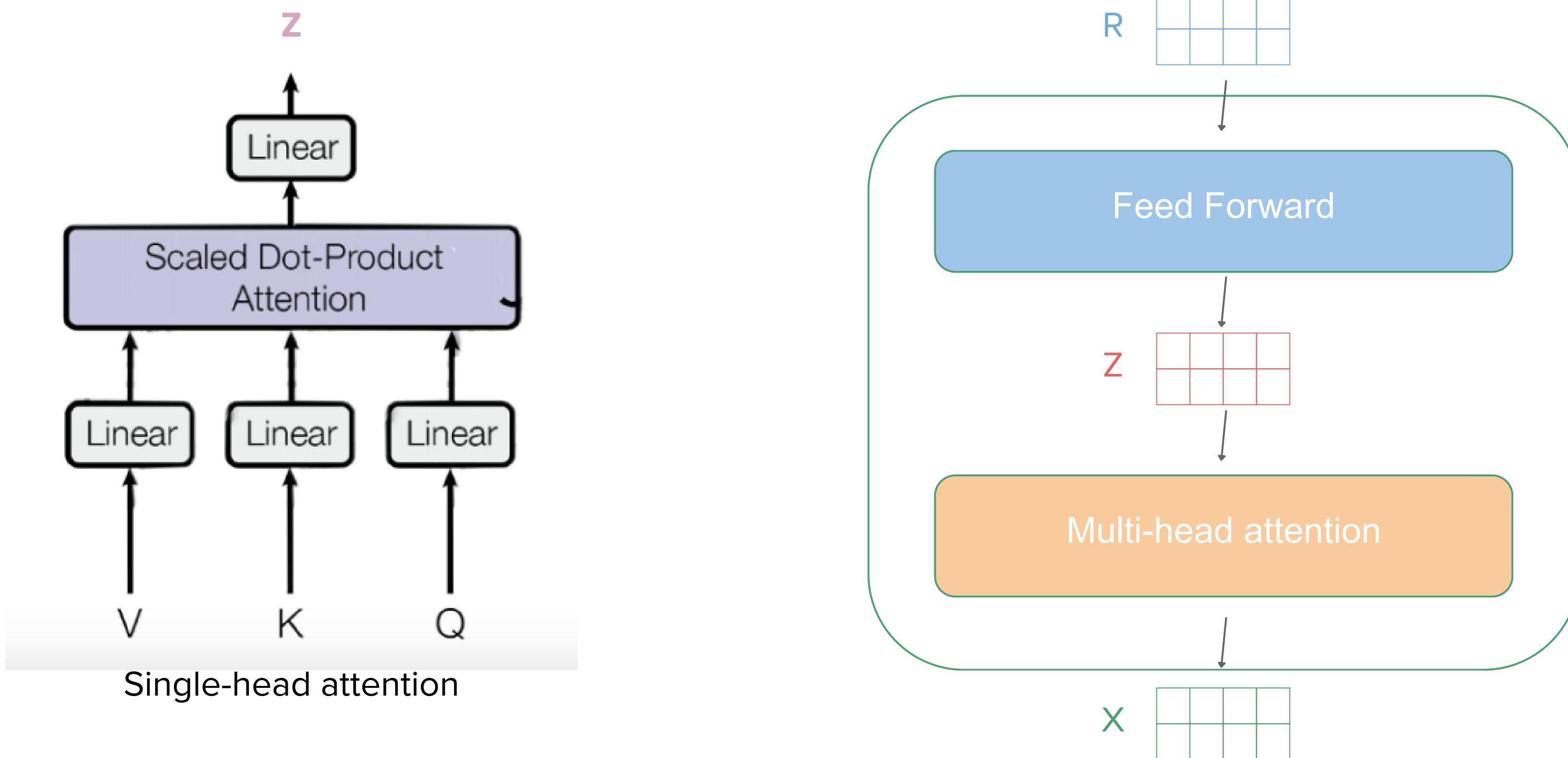
Encoder: Multi-head attention



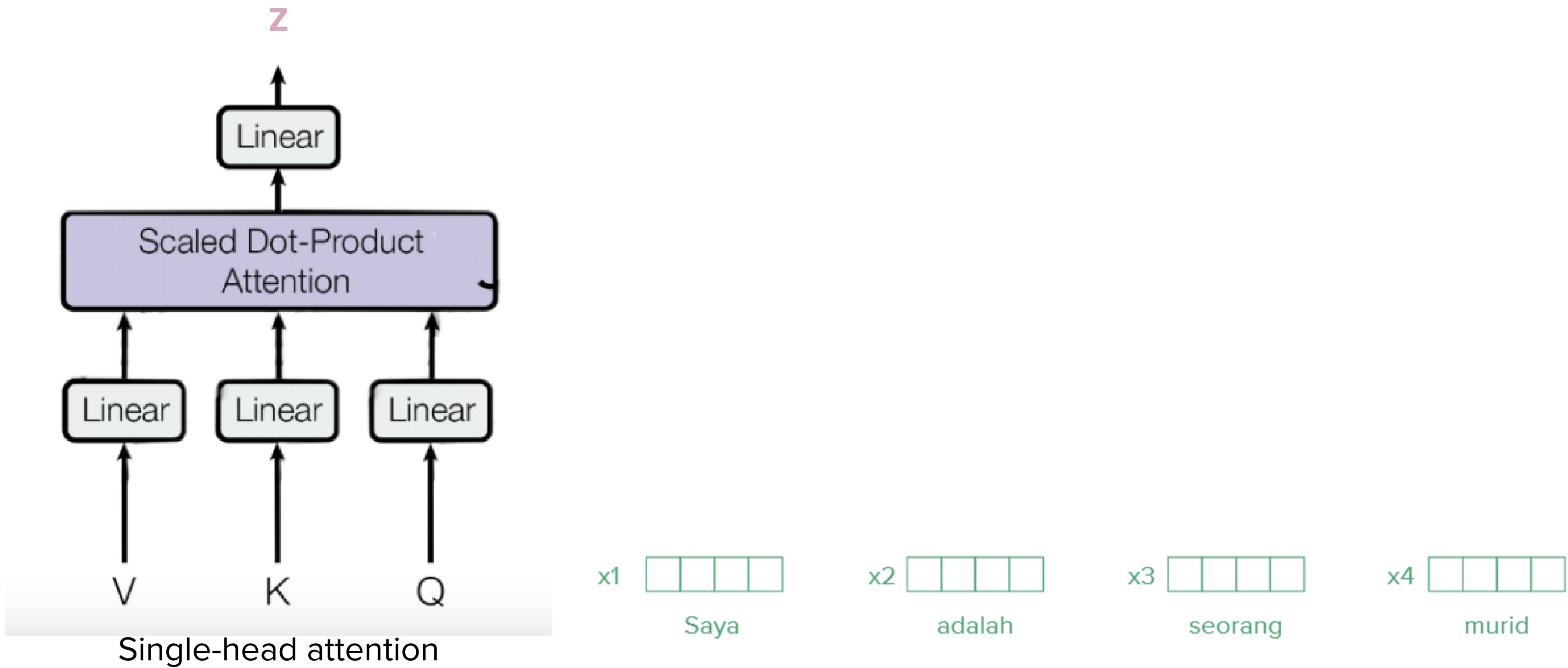
Encoder: Multi-head attention



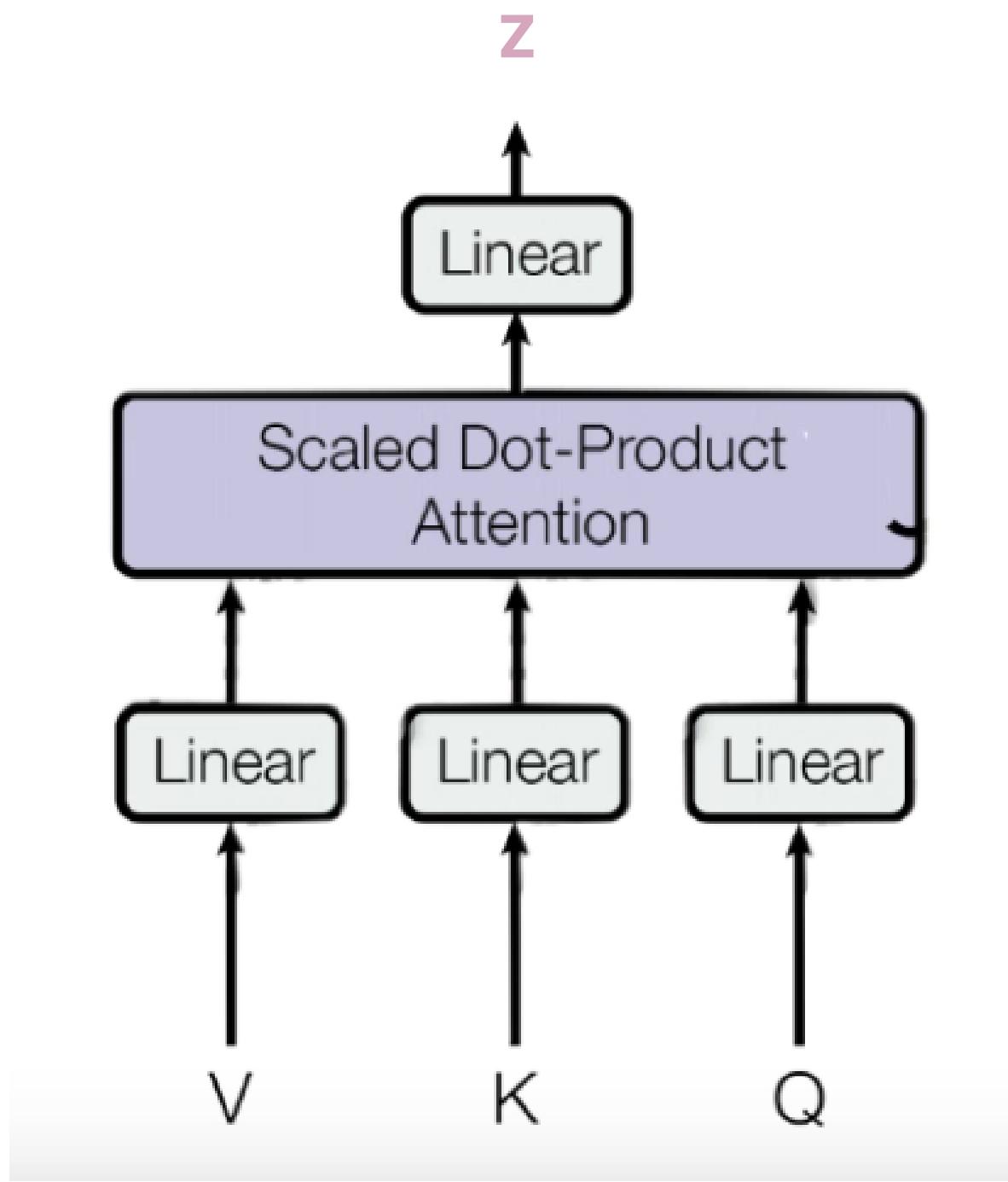
Encoder: Multi-head attention



Encoder: Multi-head attention



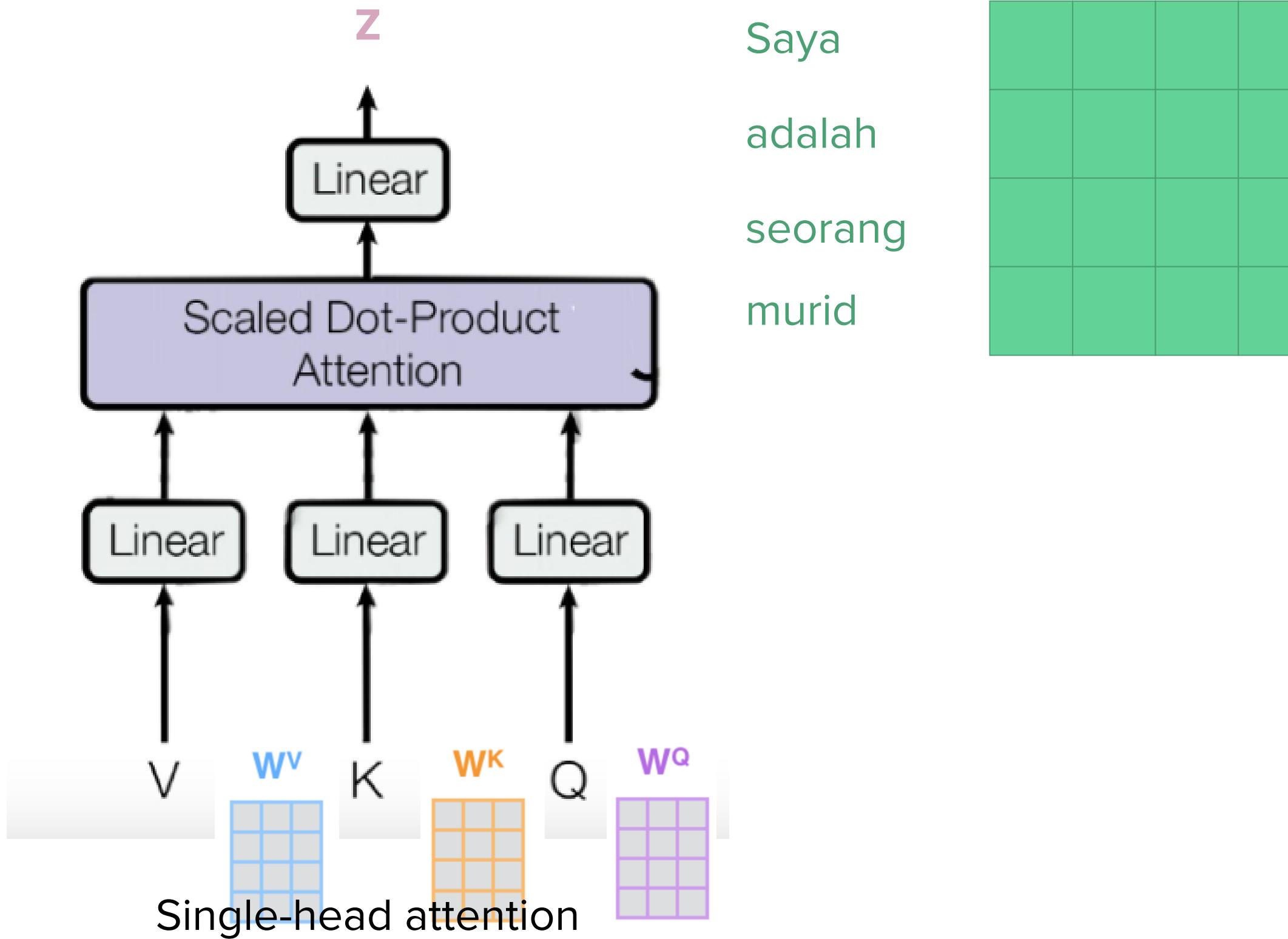
Encoder: Multi-head attention



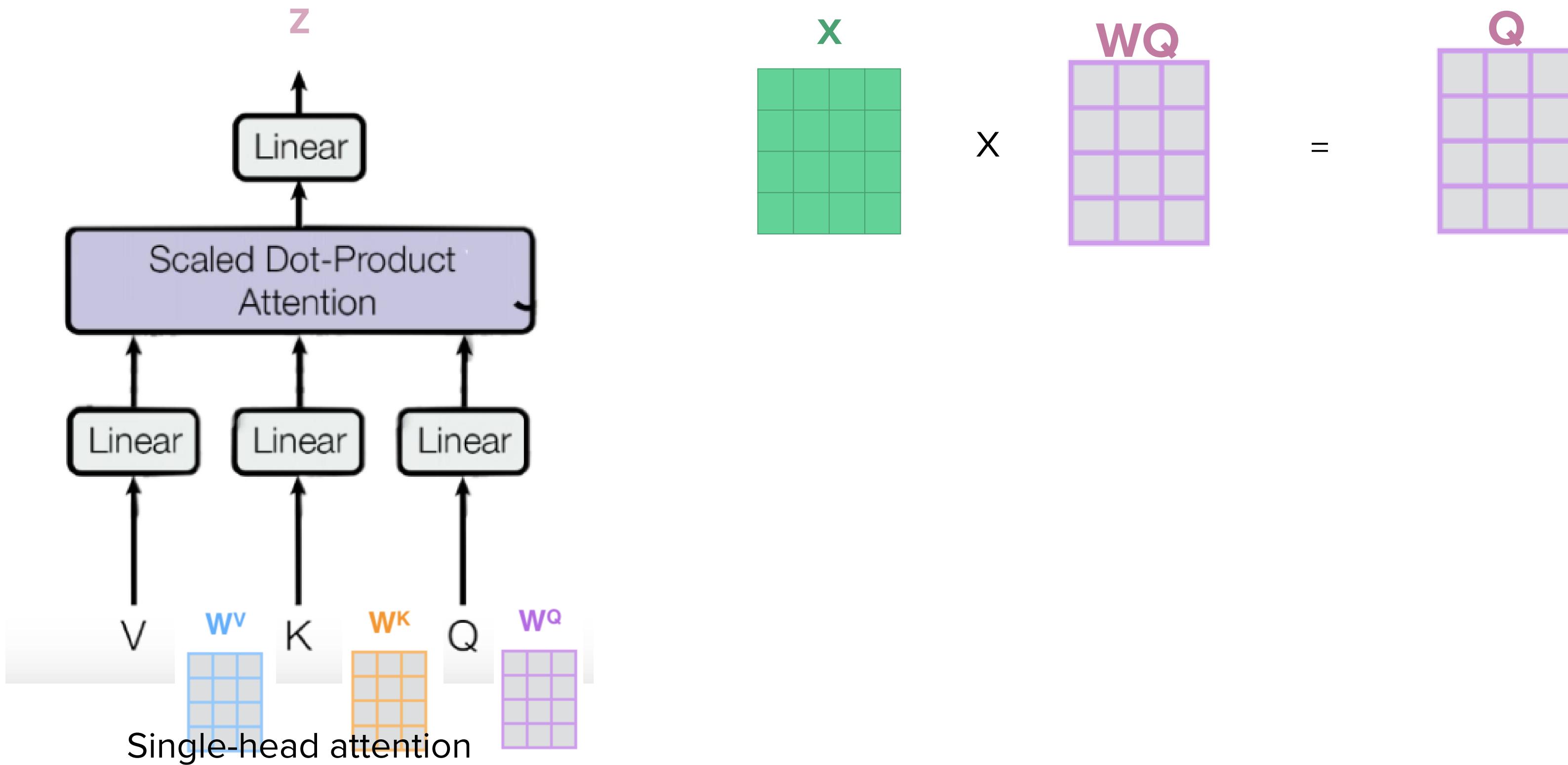
Saya
adalah
seorang
murid

| | | | |
|--|--|--|--|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

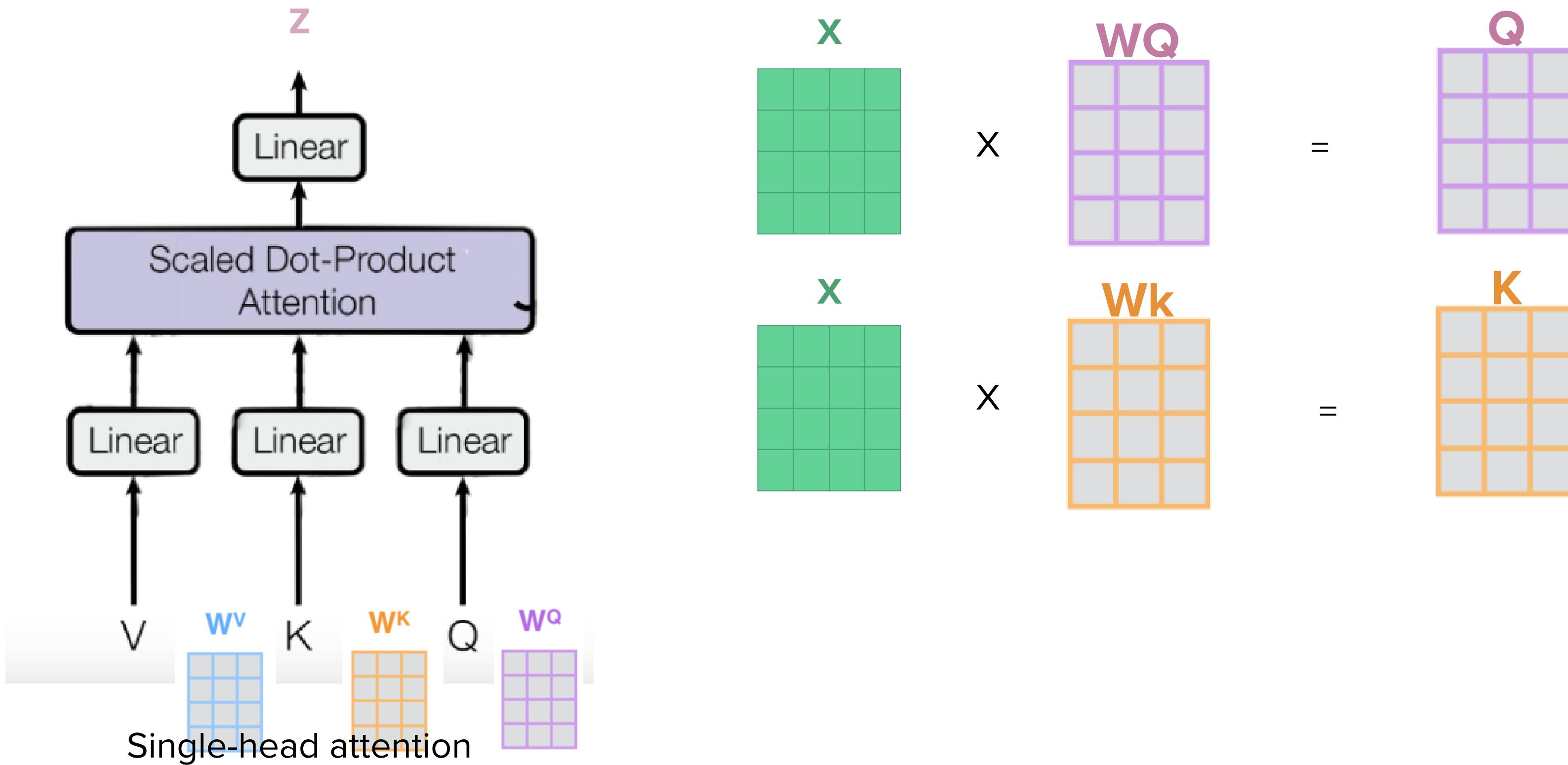
Encoder: Multi-head attention



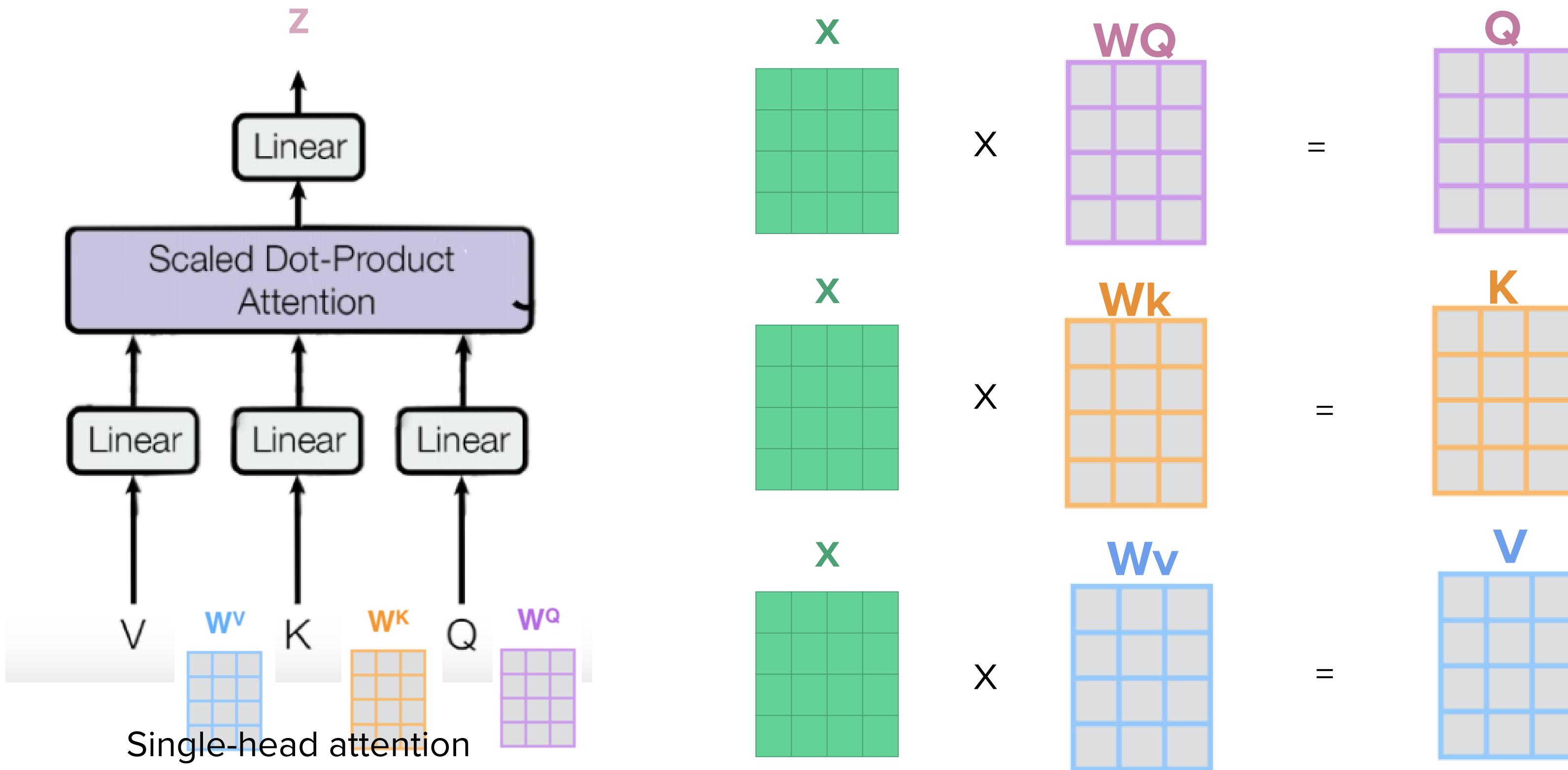
Encoder: Multi-head attention



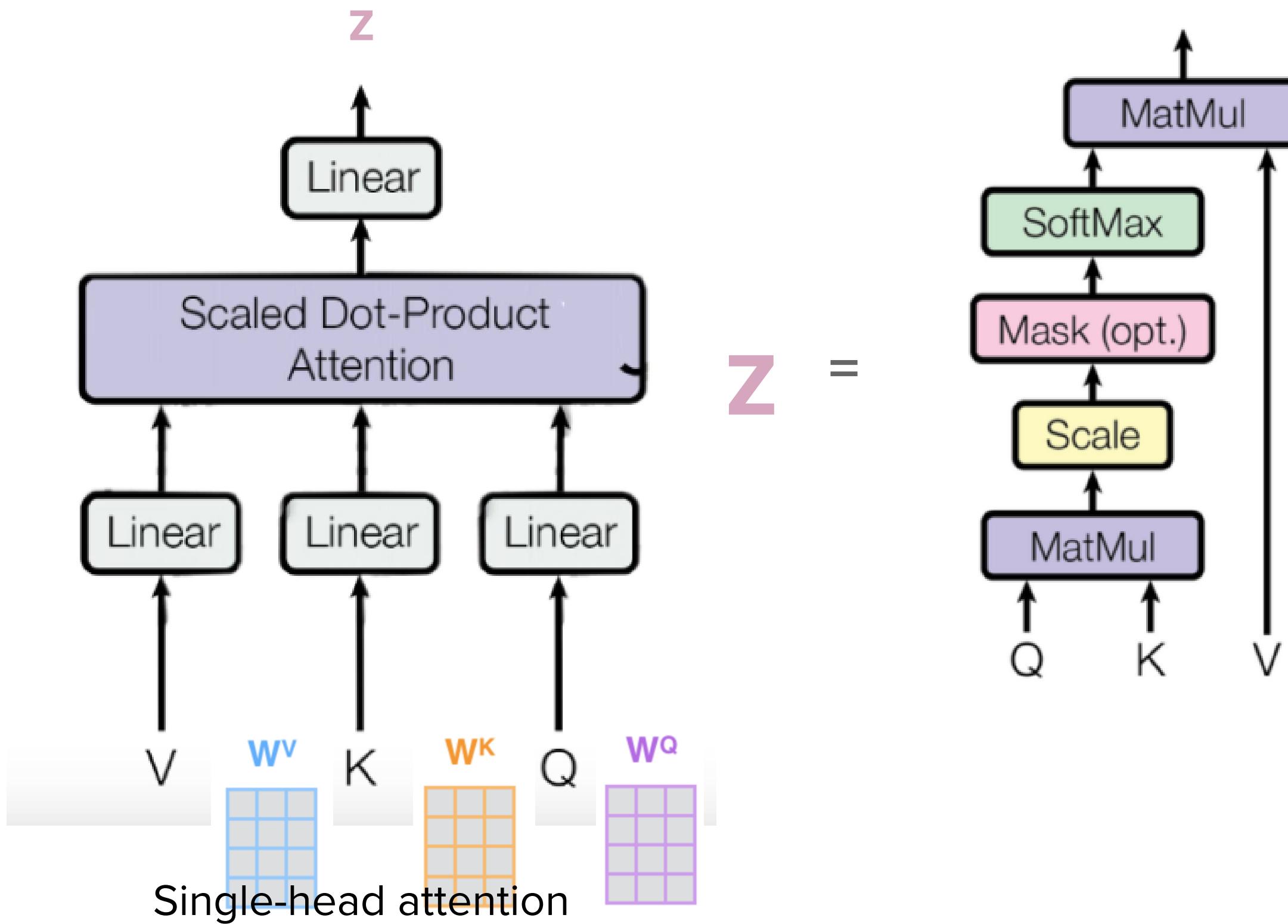
Encoder: Multi-head attention



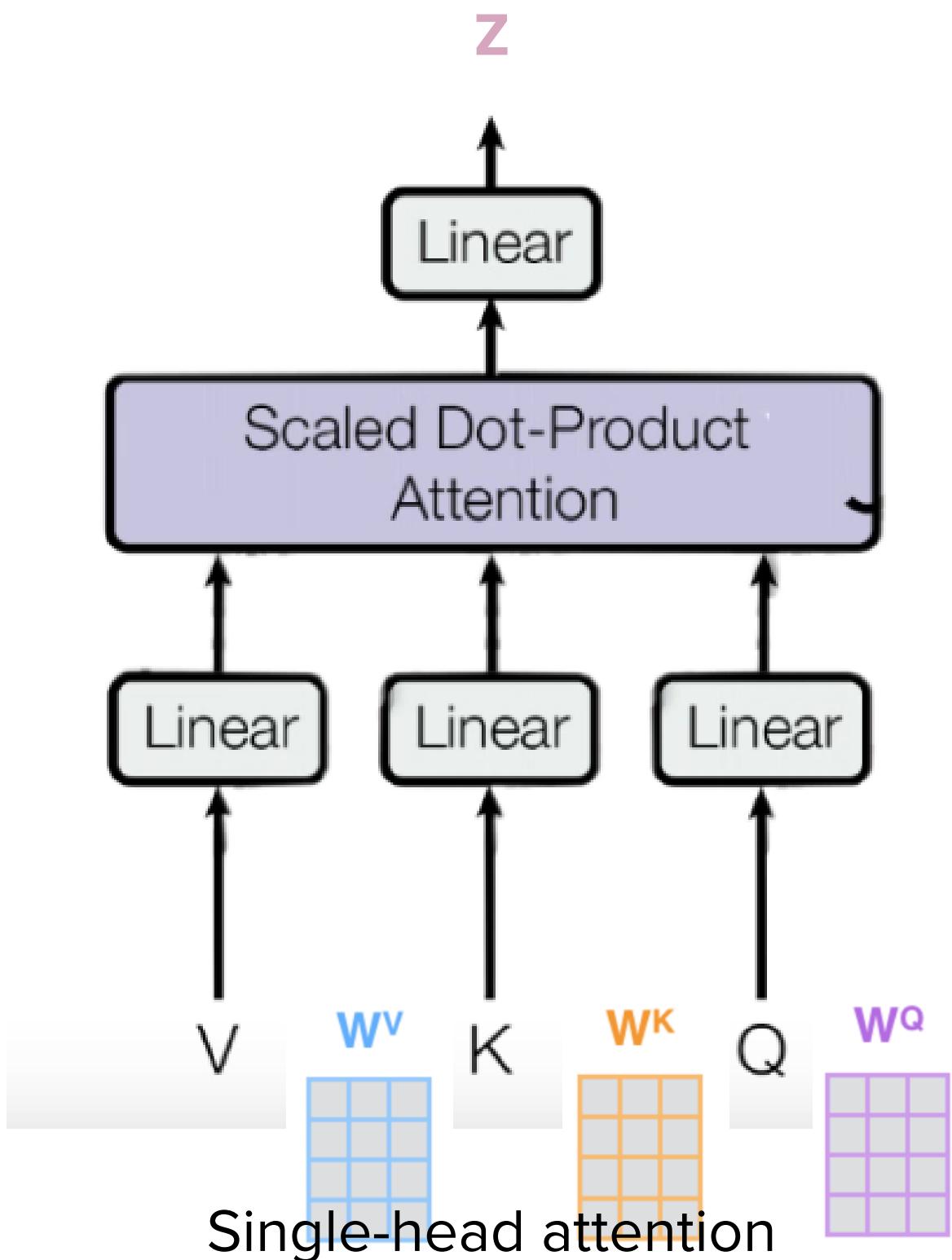
Encoder: Multi-head attention



Encoder: Multi-head attention

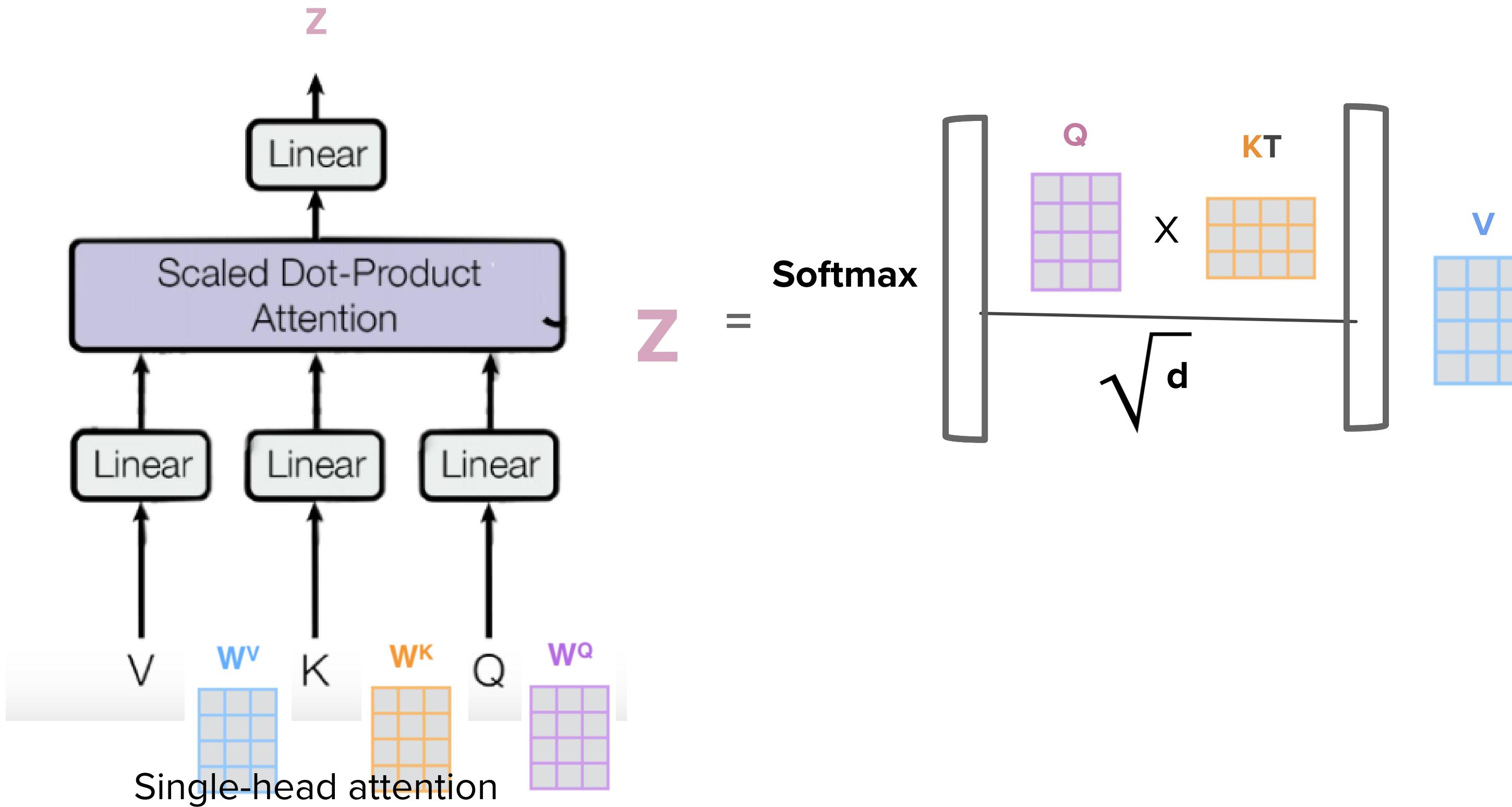


Encoder: Multi-head attention

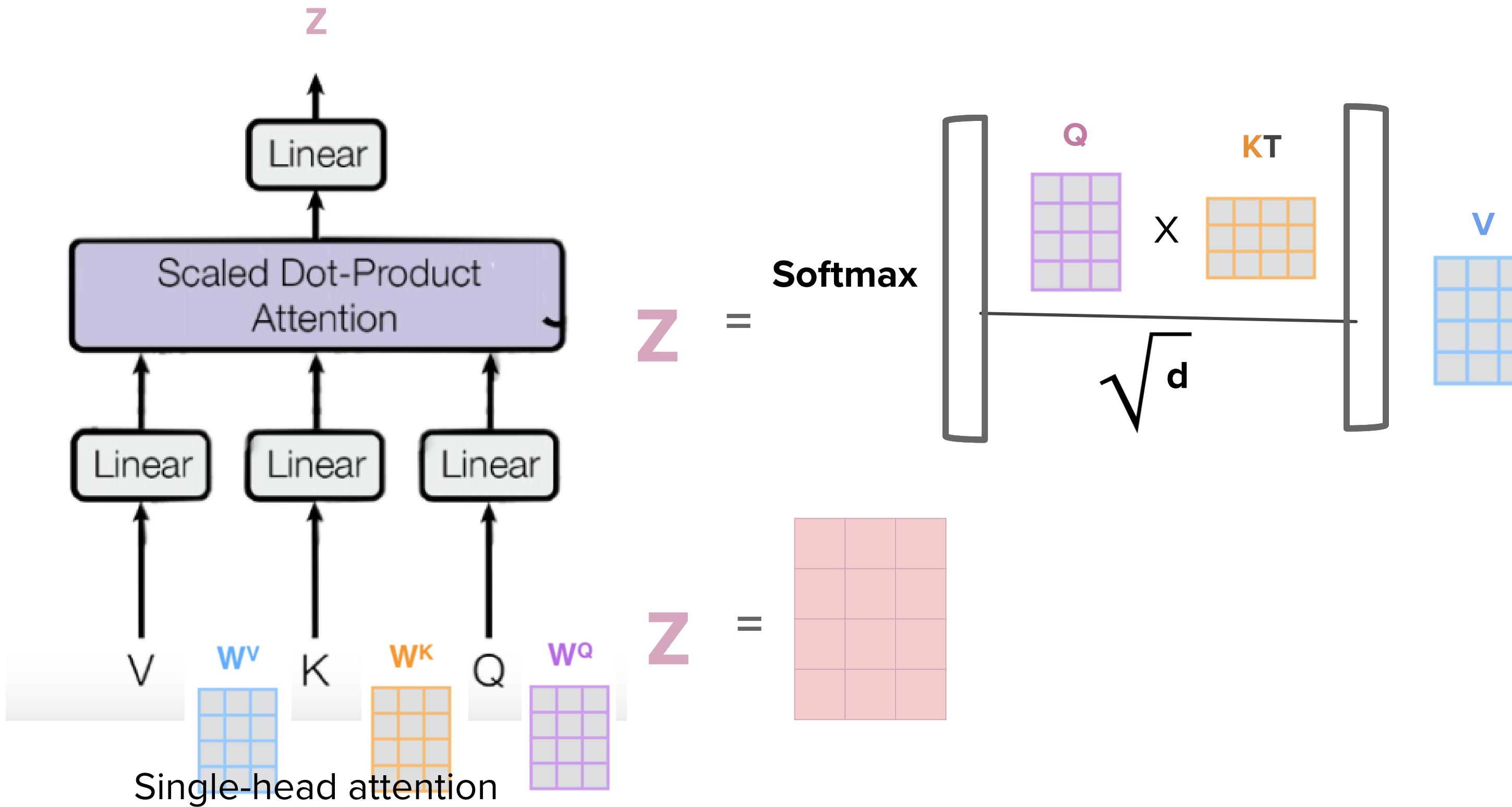


$$Z = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{\text{Dimension of vector } Q, K \text{ or } V}} \right) \cdot V$$

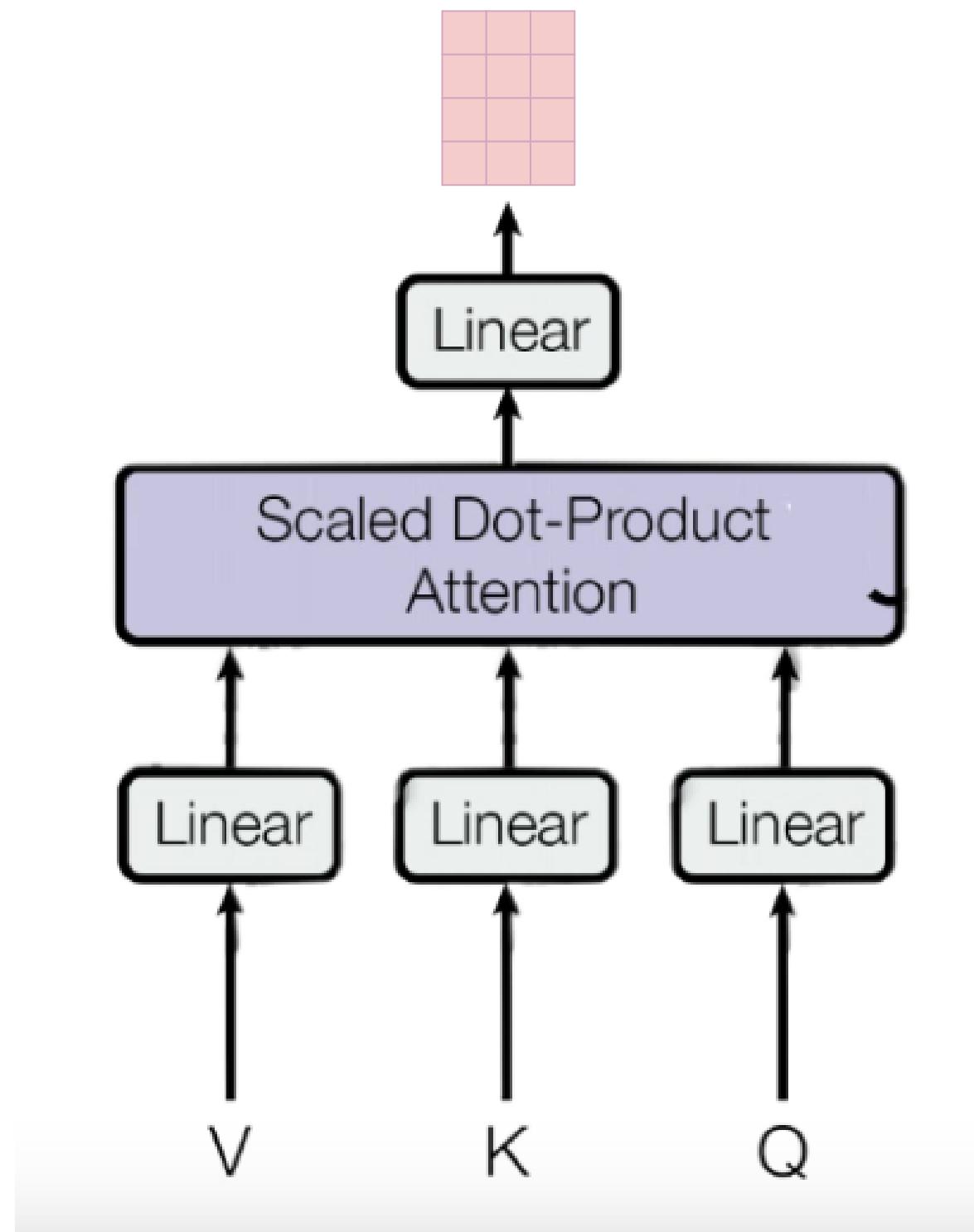
Encoder: Multi-head attention



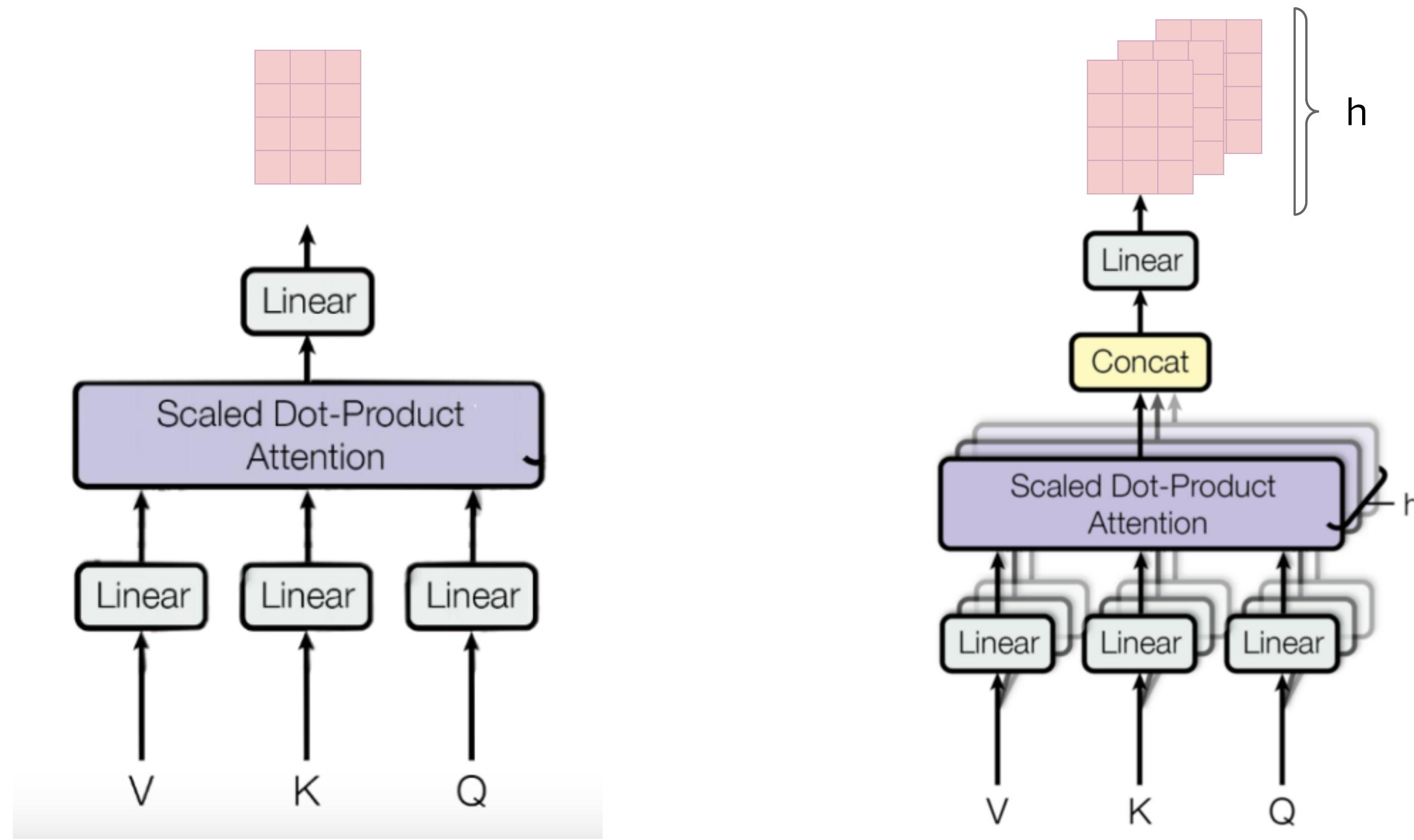
Encoder: Multi-head attention



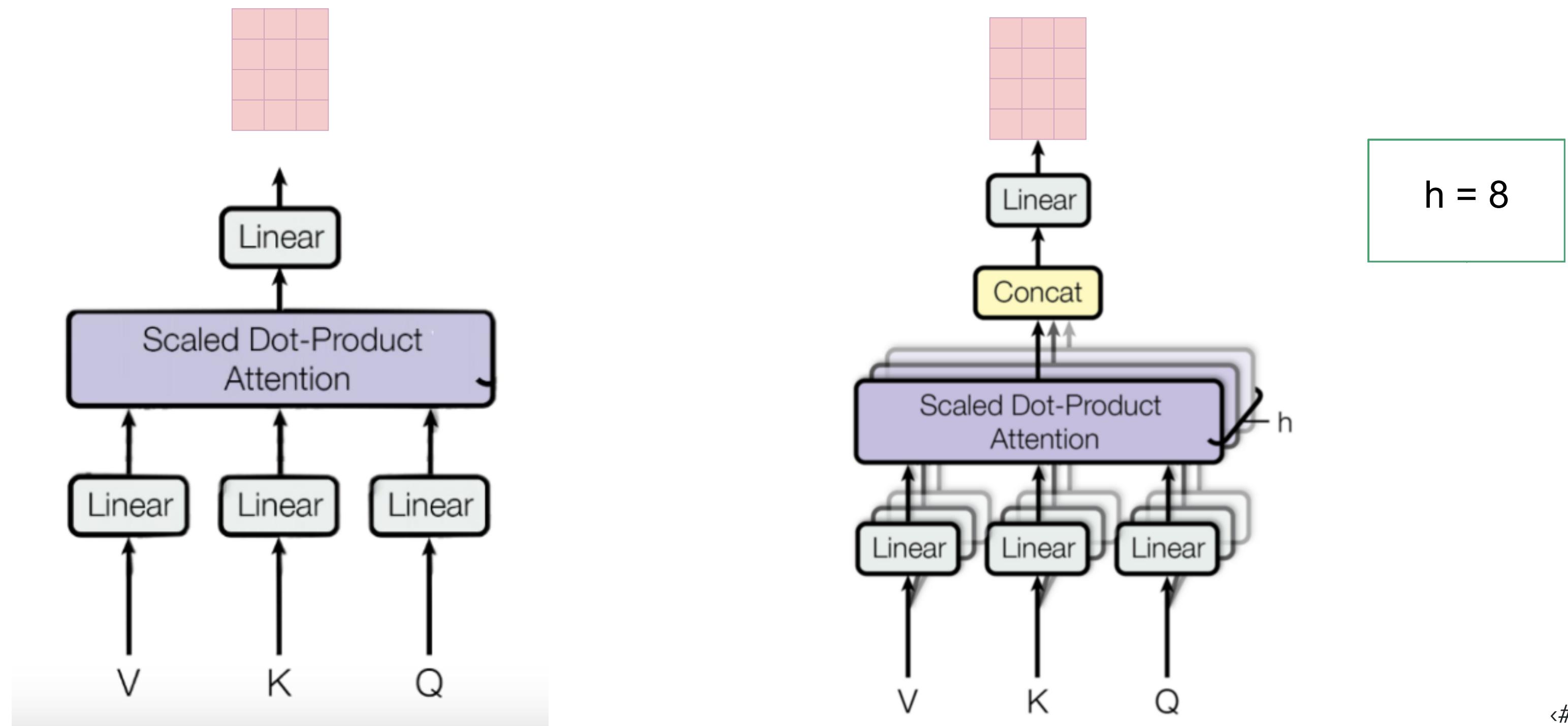
Encoder: Multi-head attention



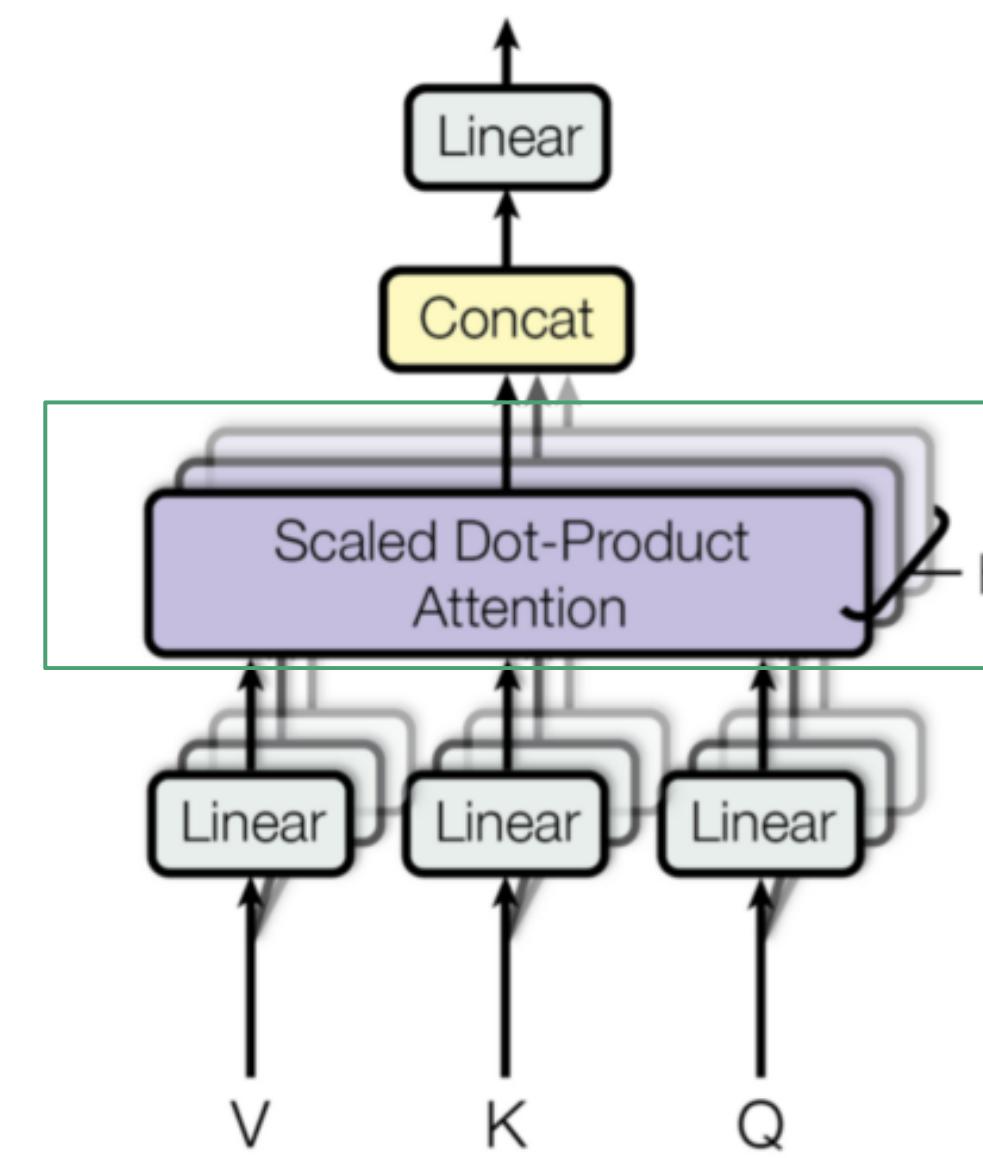
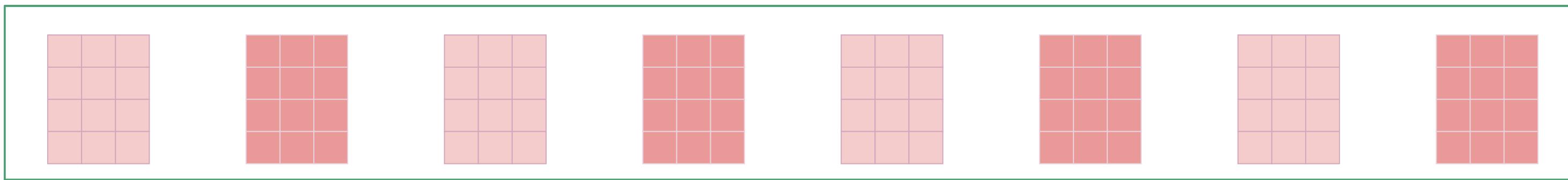
Encoder: Multi-head attention



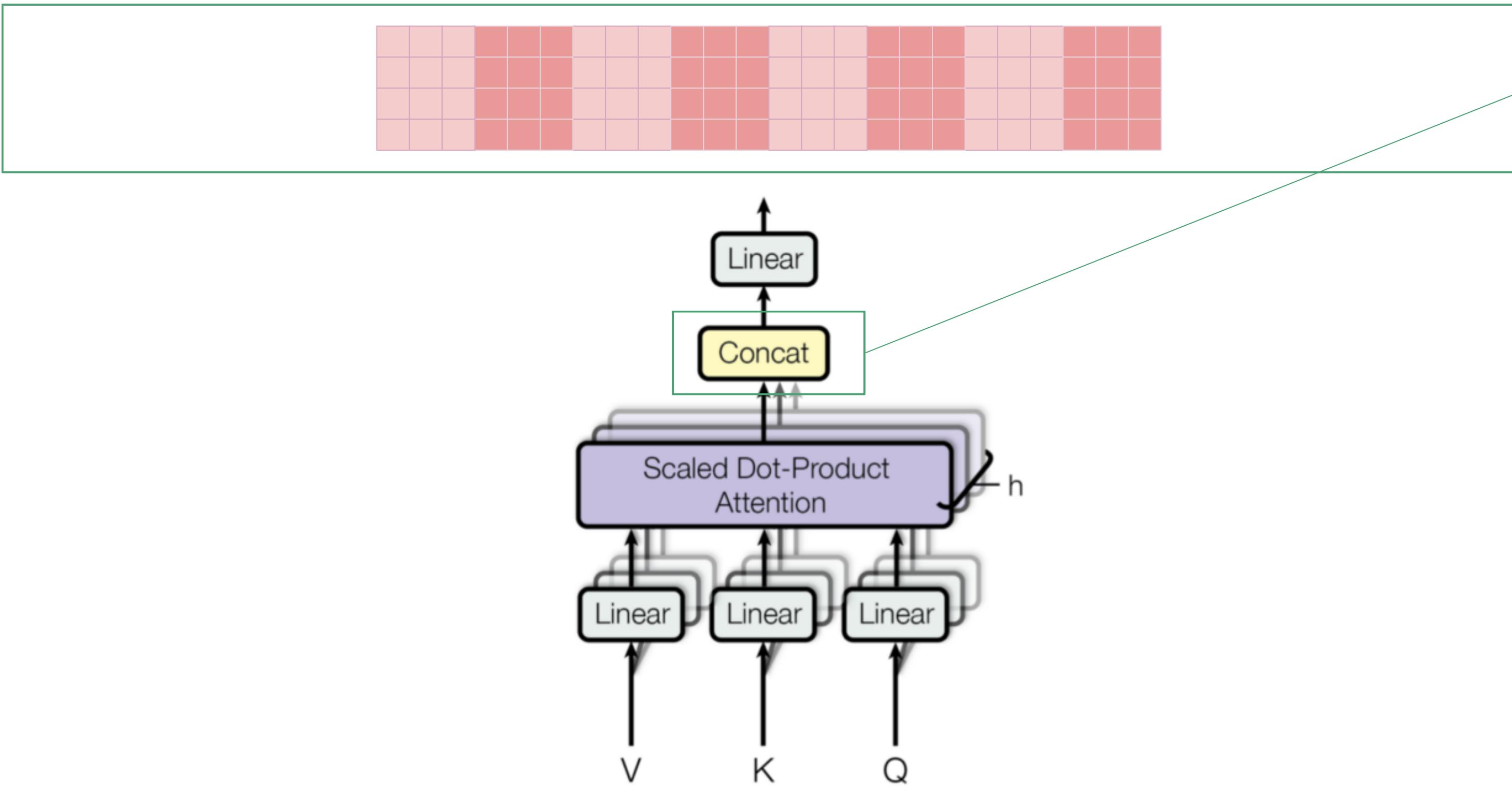
Encoder: Multi-head attention



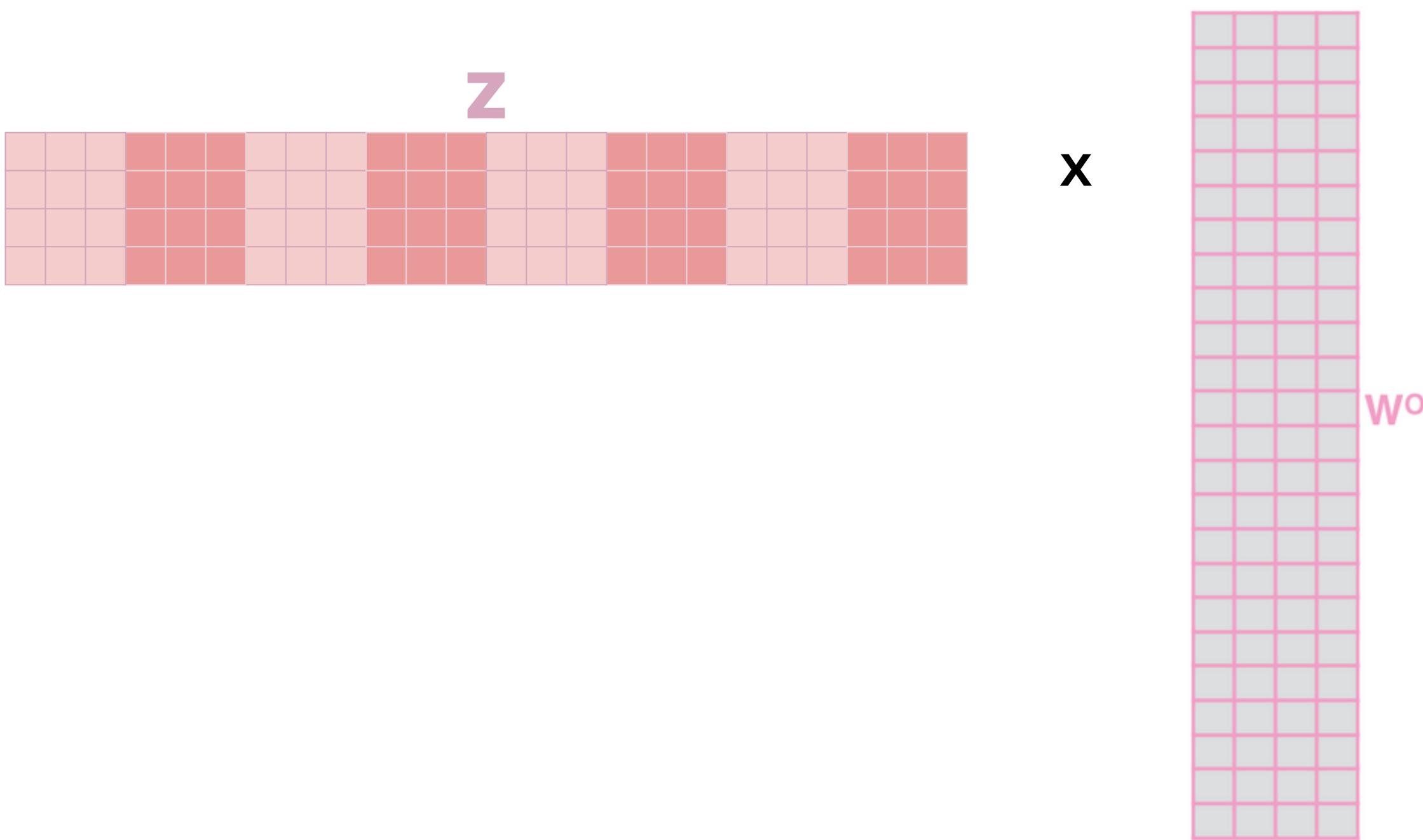
Encoder: Multi-head attention



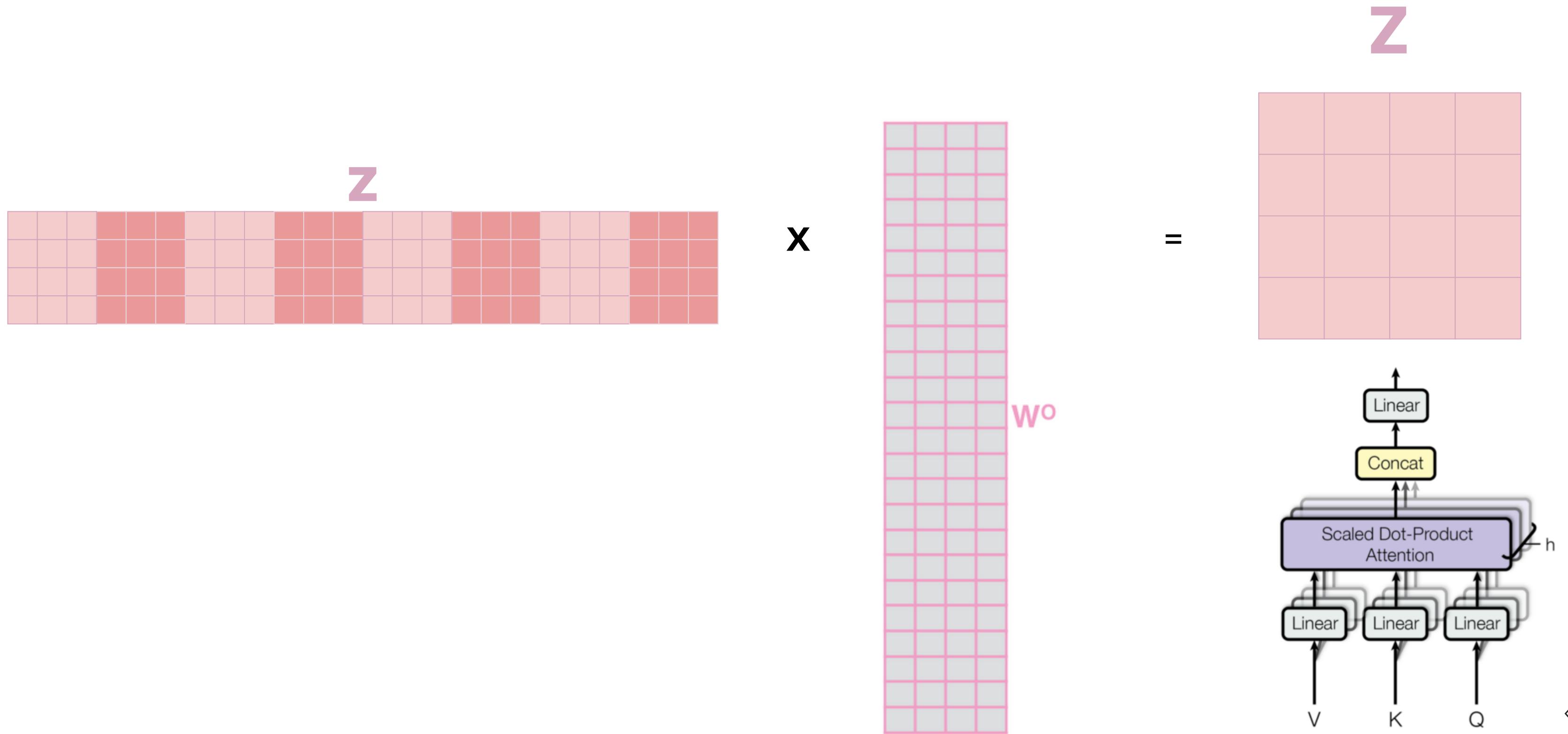
Encoder: Multi-head attention



Encoder: Multi-head attention

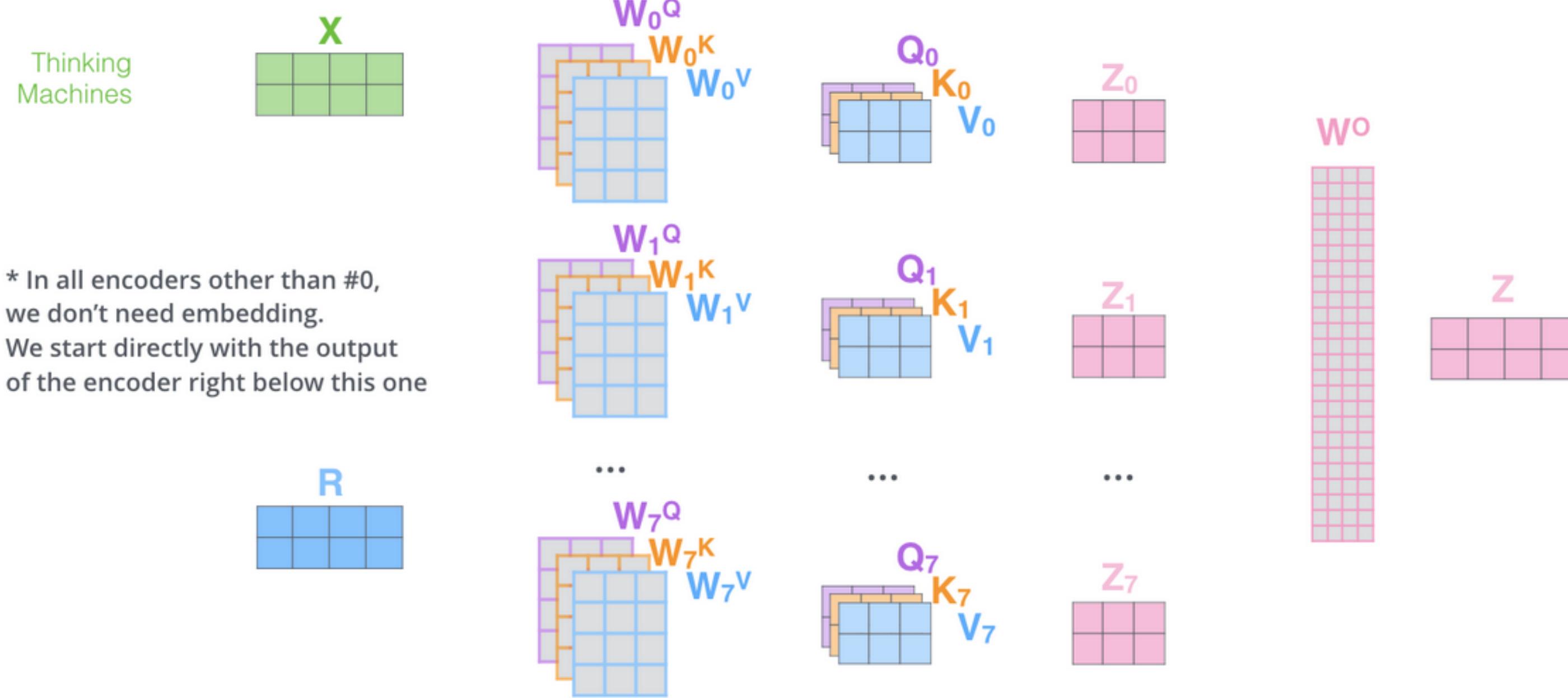


Encoder: Multi-head attention

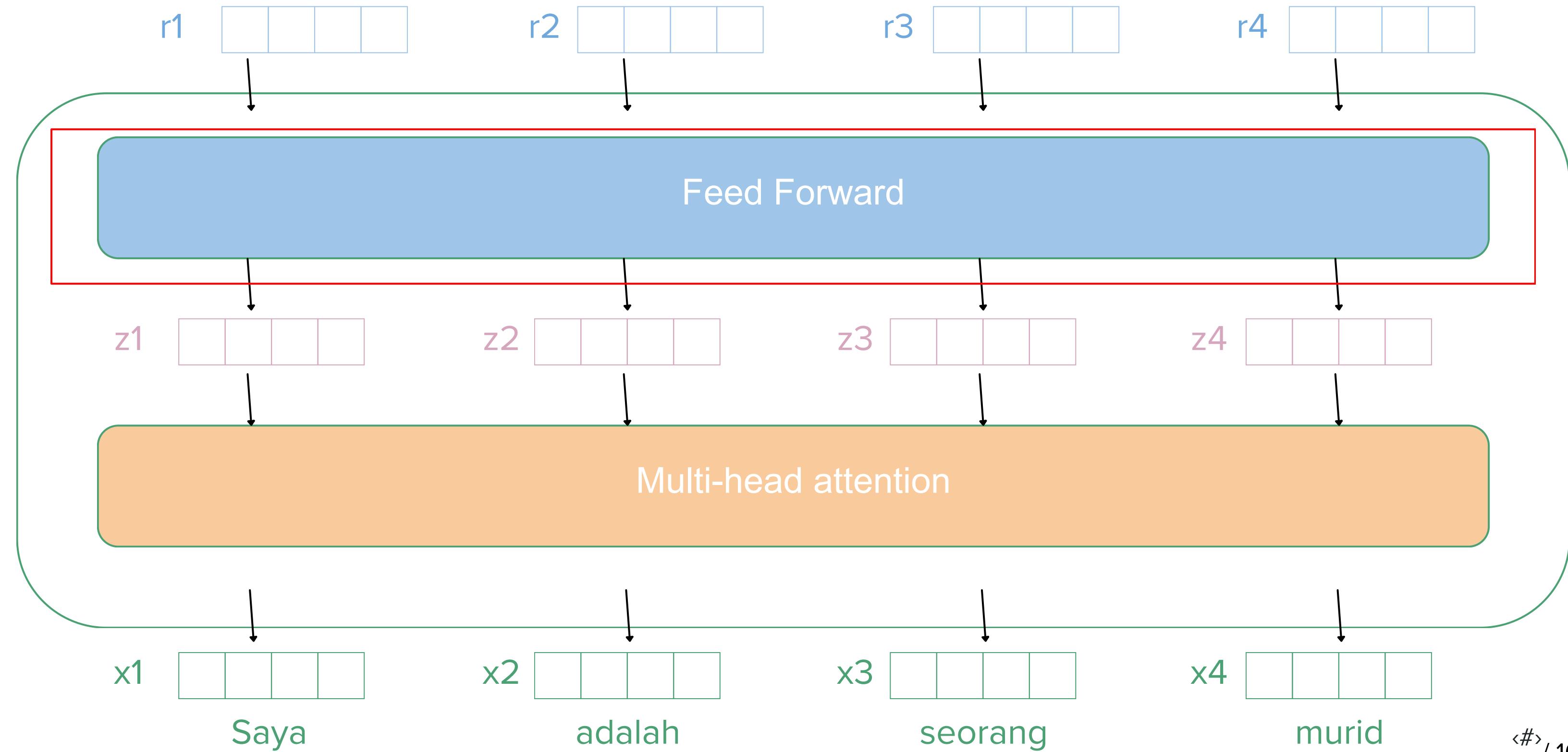


Encoder: Multi-head attention (summary)

- 1) This is our input sentence* X
- 2) We embed each word*
- 3) Split into 8 heads. We multiply X or R with weight matrices
- 4) Calculate attention using the resulting $Q/K/V$ matrices
- 5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

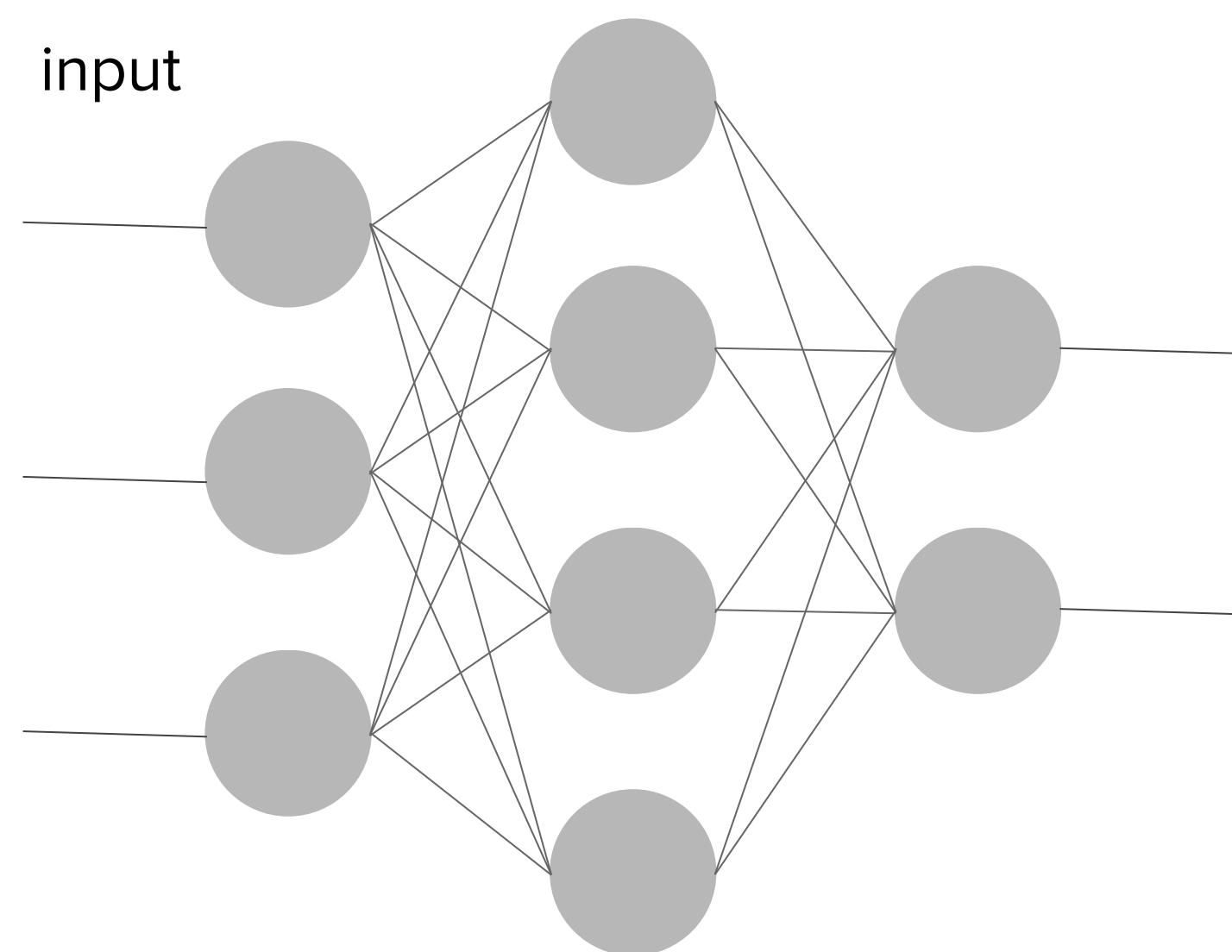
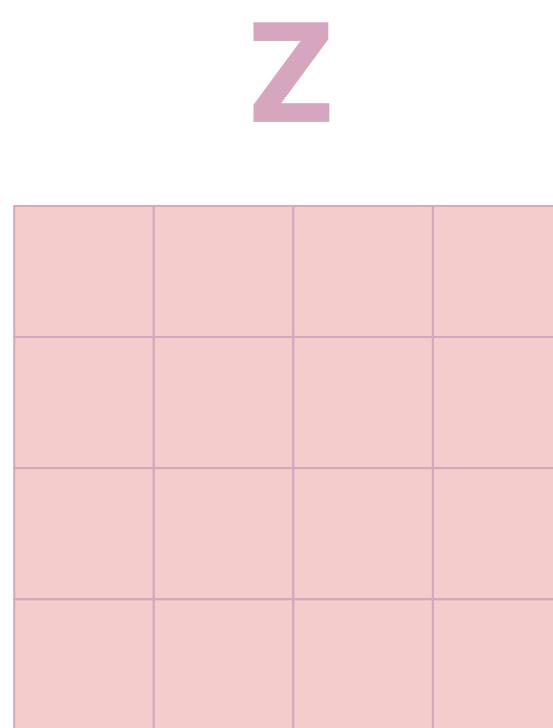


Encoder: Feed Forward

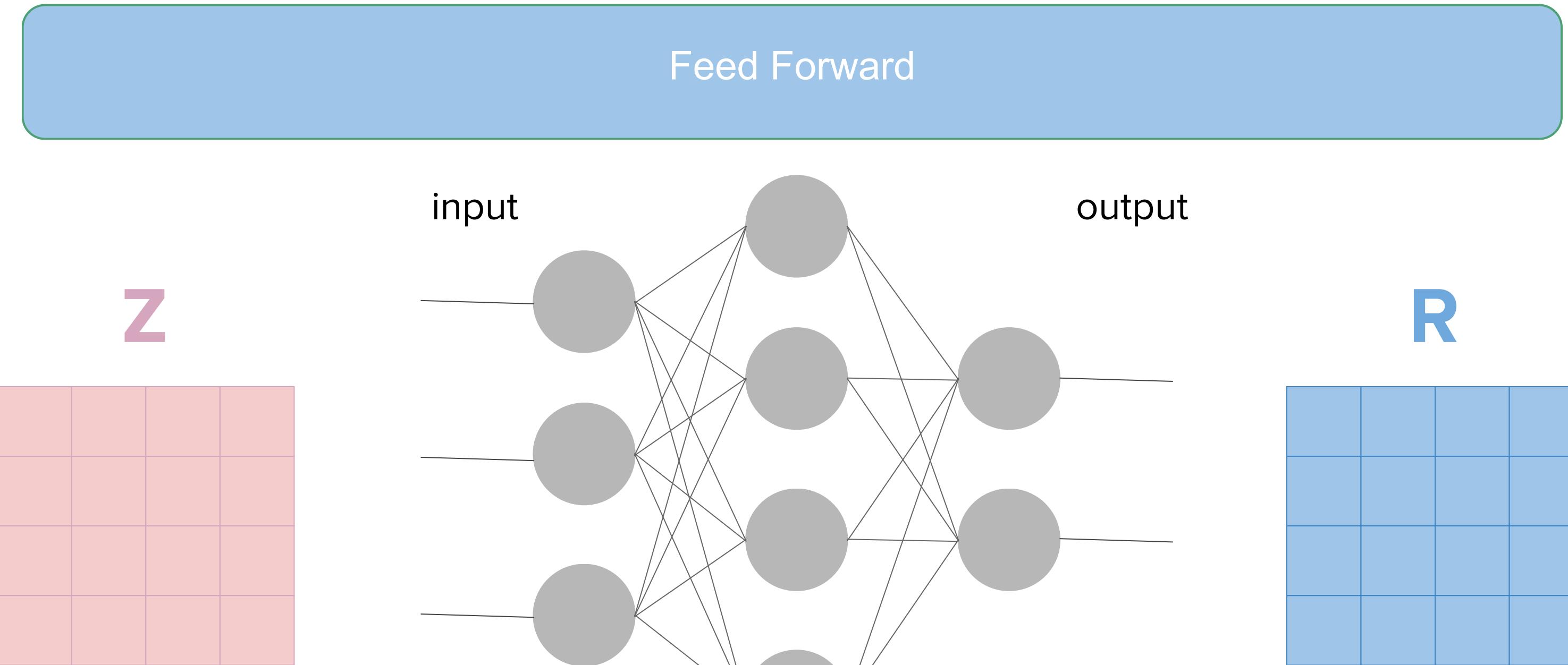


Encoder: Feed Forward

Feed Forward



Encoder: Feed Forward

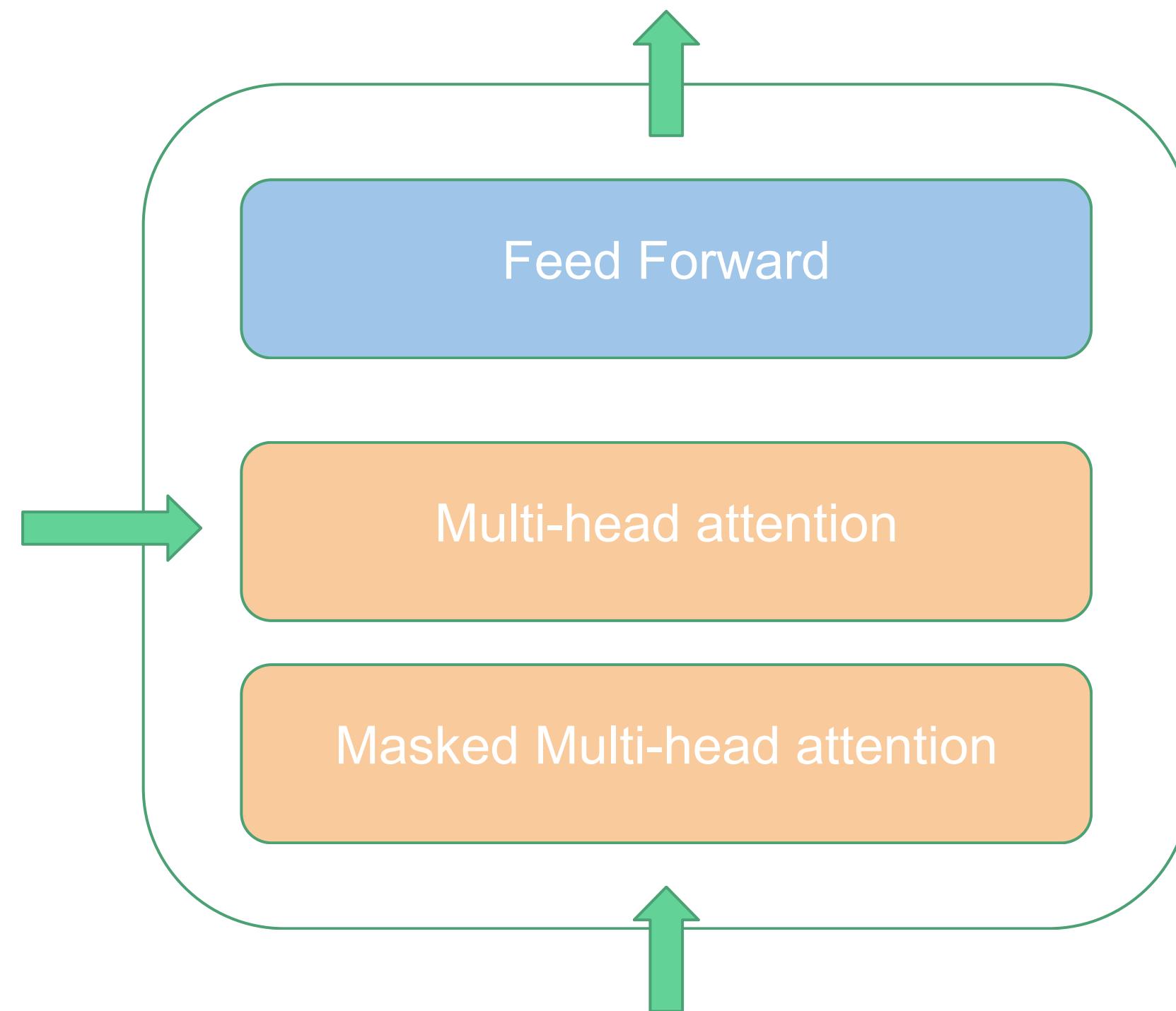


— Any Question Guys ~

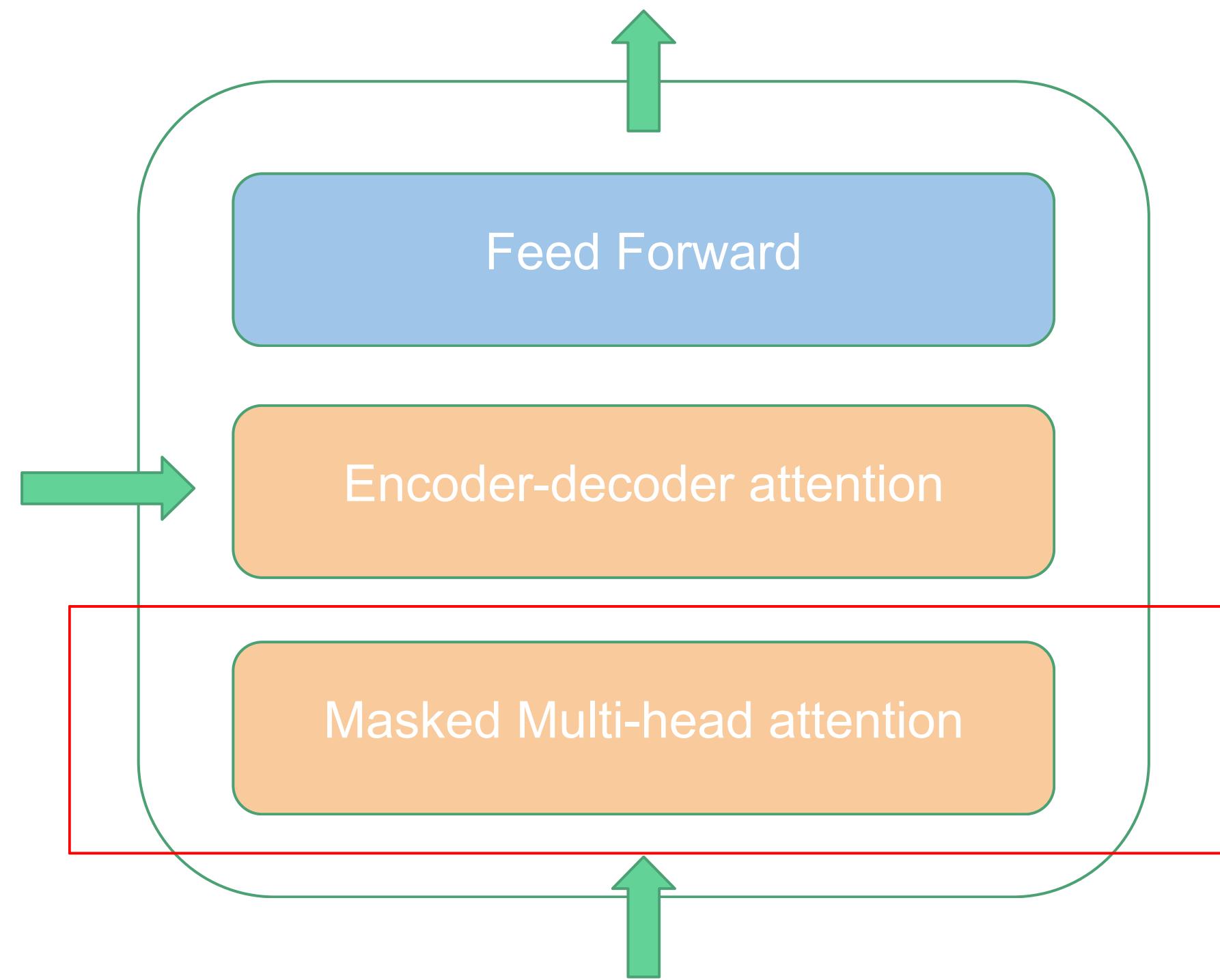


— Transformer's Decoder

Decoder



Decoder

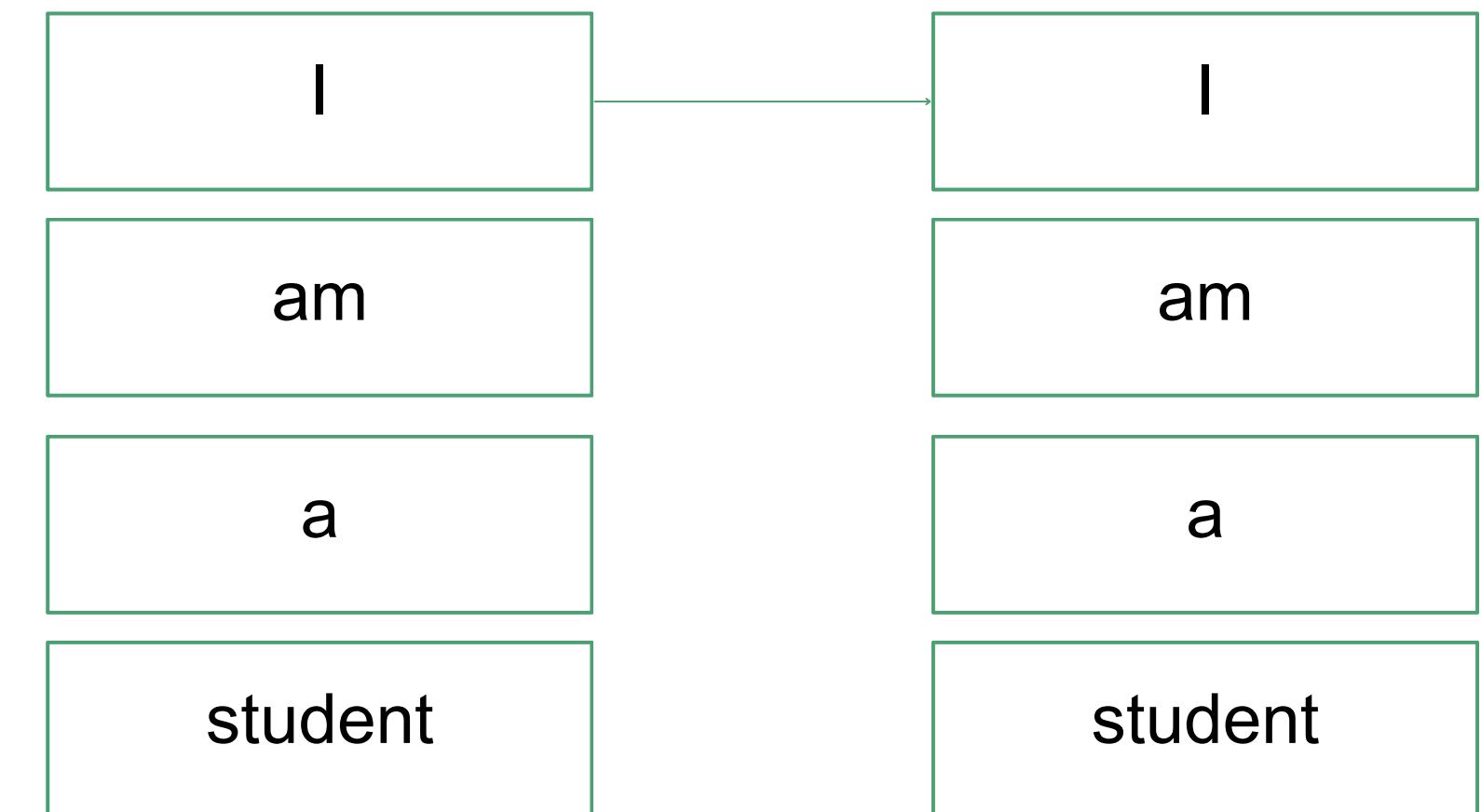


Decoder: Masked Multi-head attention

Very similar with multi-head attention in
encoder

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words and another previous words in a sentence

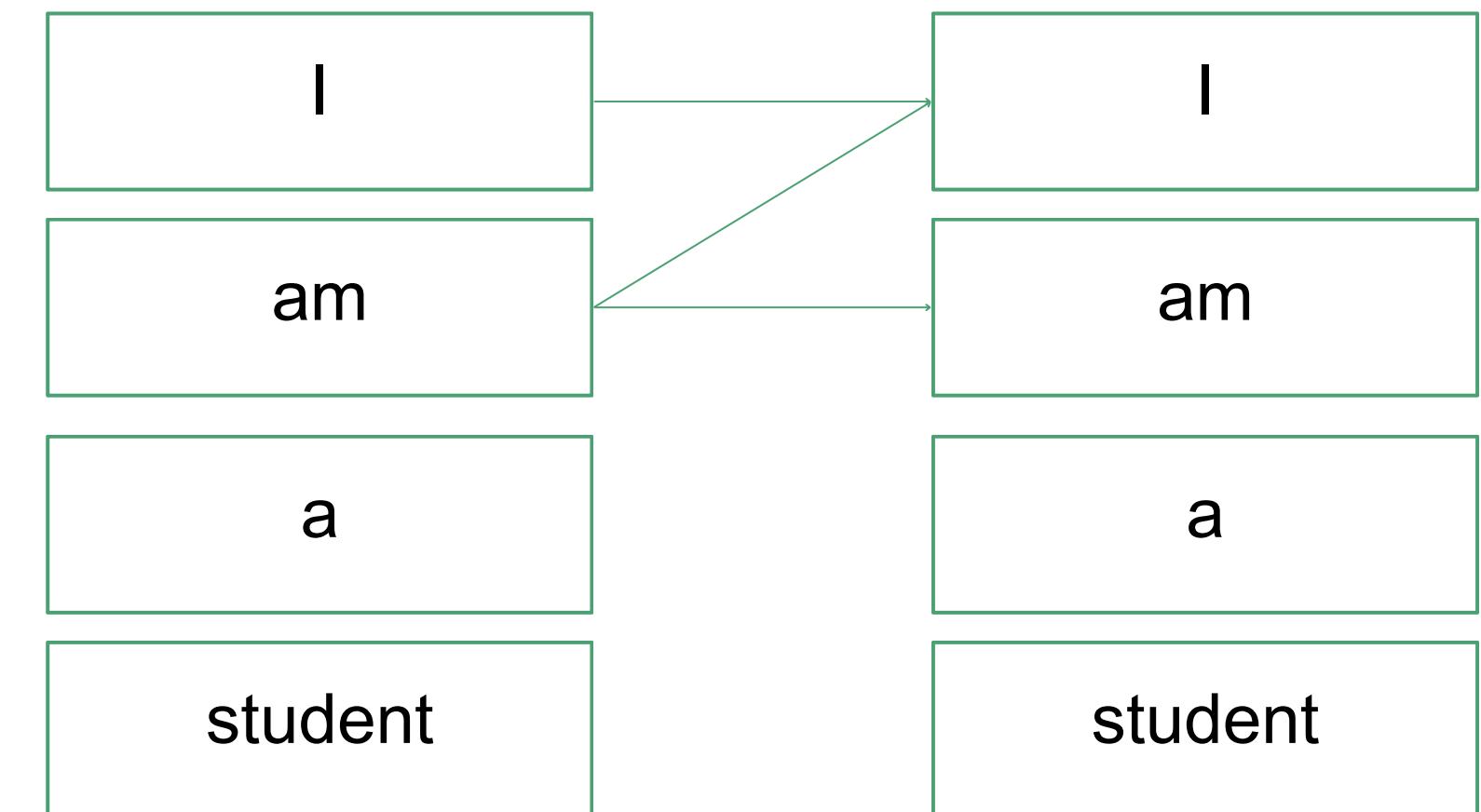


Decoder: Masked Multi-head attention

Very similar with multi-head attention in
encoder

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words and another previous words in a sentence

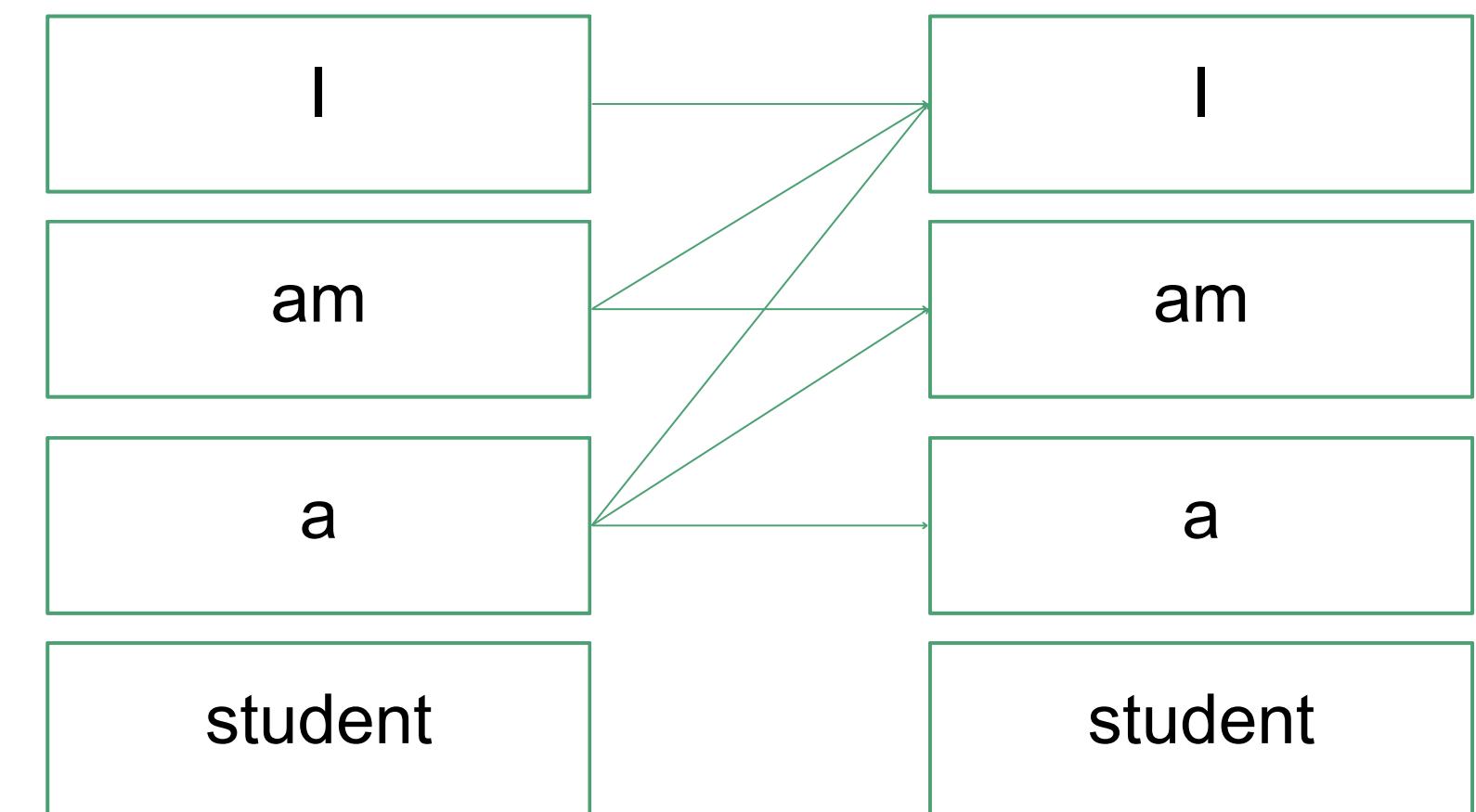


Decoder: Masked Multi-head attention

Very similar with multi-head attention in
encoder

Purpose:

1. Find the focus of the sentence
2. Find the relationship between words and another previous words in a sentence

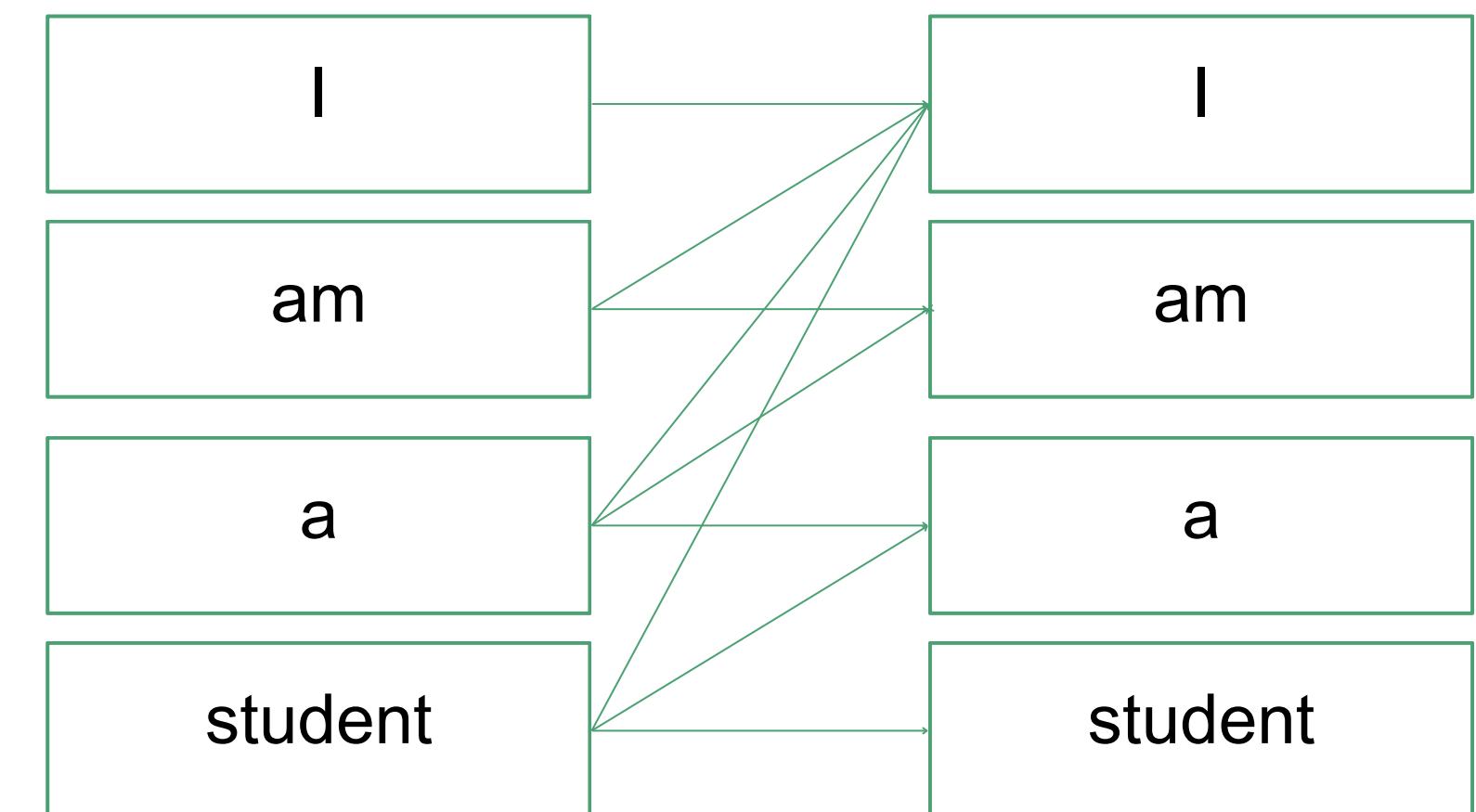


Decoder: Masked Multi-head attention

Very similar with multi-head attention in
encoder

Purpose:

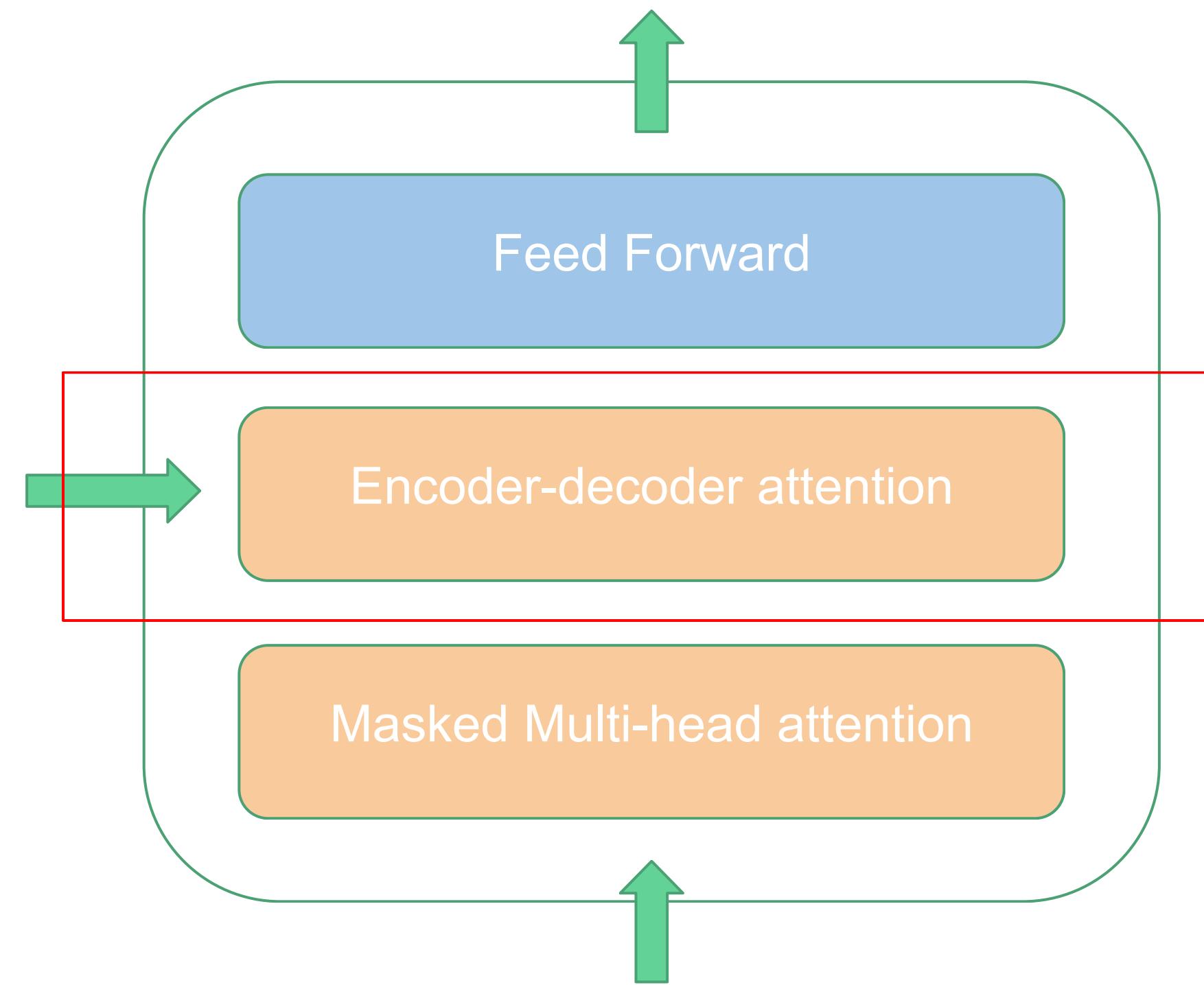
1. Find the focus of the sentence
2. Find the relationship between words and another previous words in a sentence



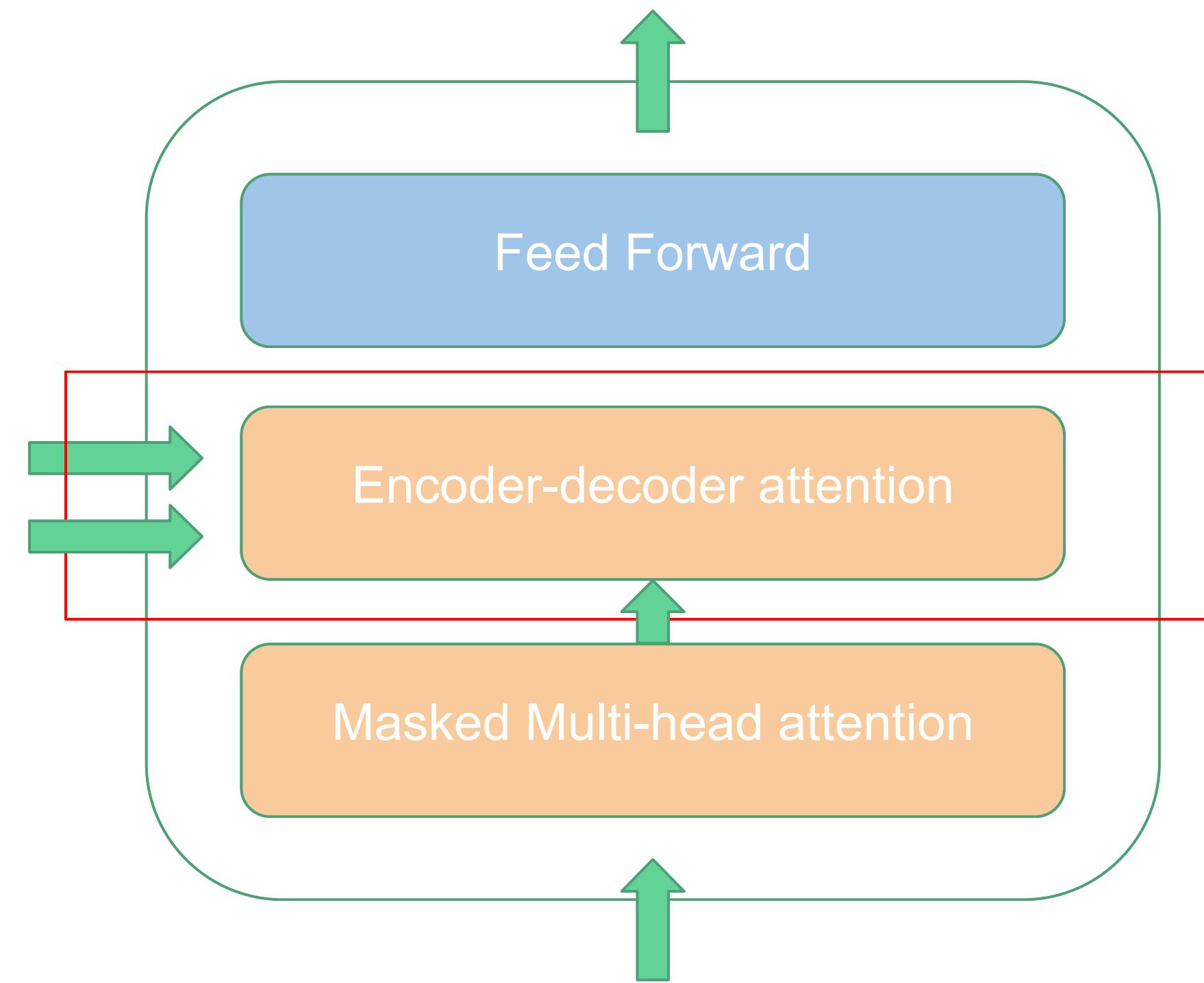
Decoder: Masked Multi-head attention

| | |
|---------|-----------------------|
| I | [1 0 0 0] |
| am | [0.05 0.95 0 0] |
| a | [0.20 0.05 0.75 0] |
| student | [0.05 0.03 0.04 0.88] |

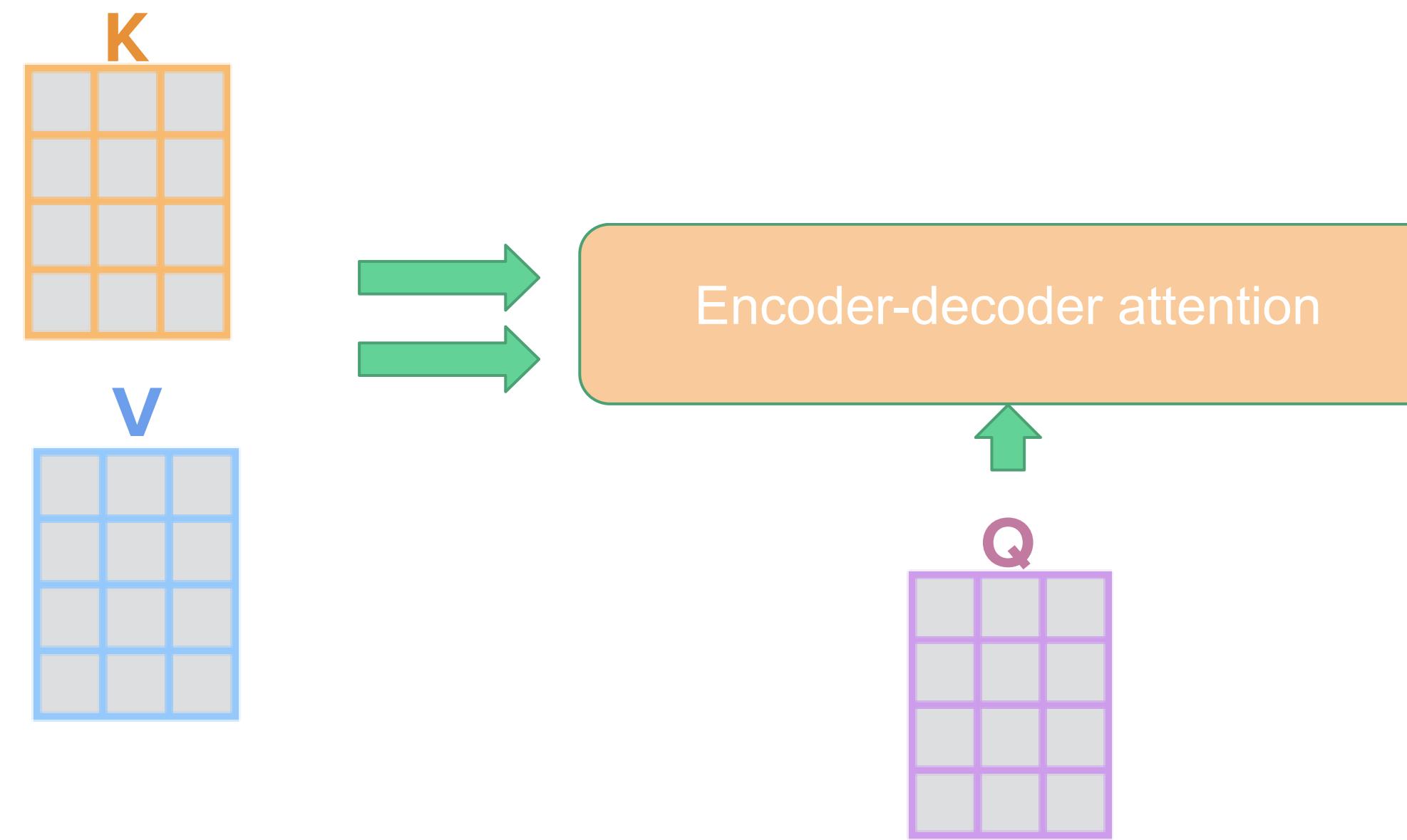
Decoder: Encoder-decoder attention



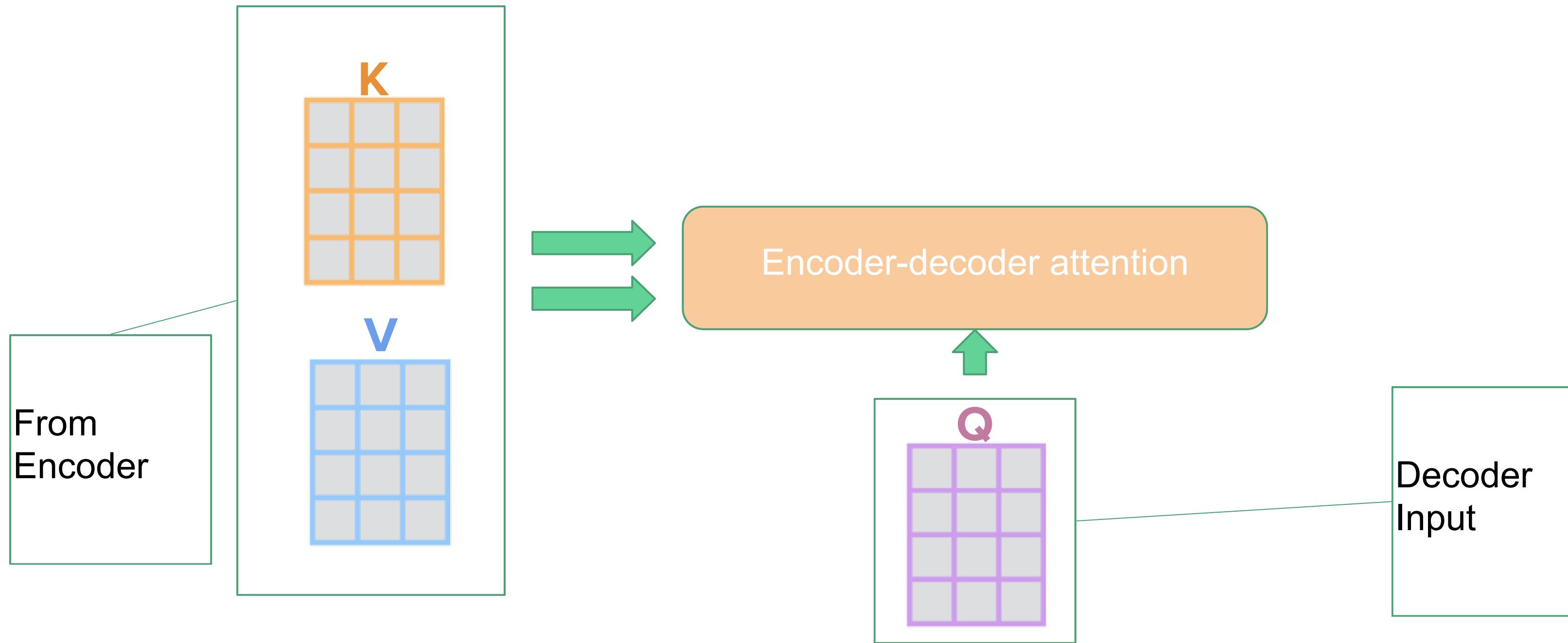
Decoder: Encoder-decoder attention



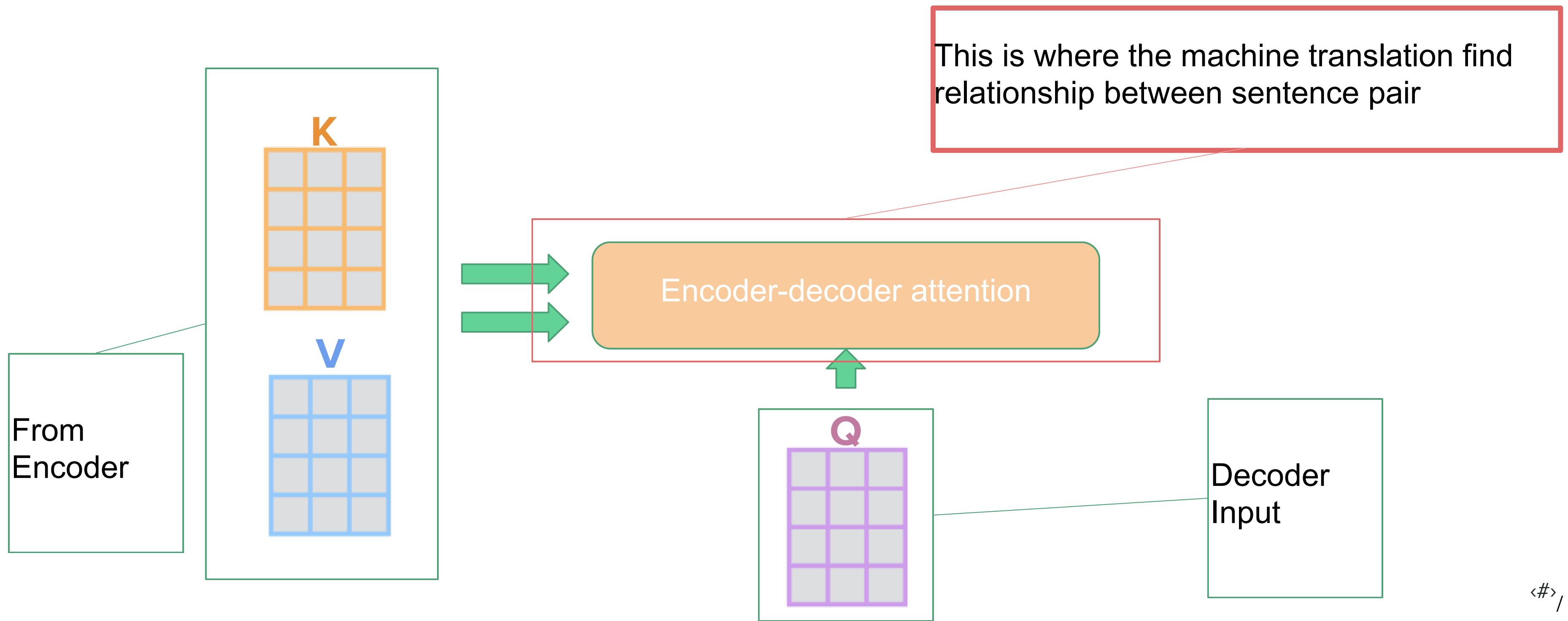
Decoder: Encoder-decoder attention



Decoder: Encoder-decoder attention



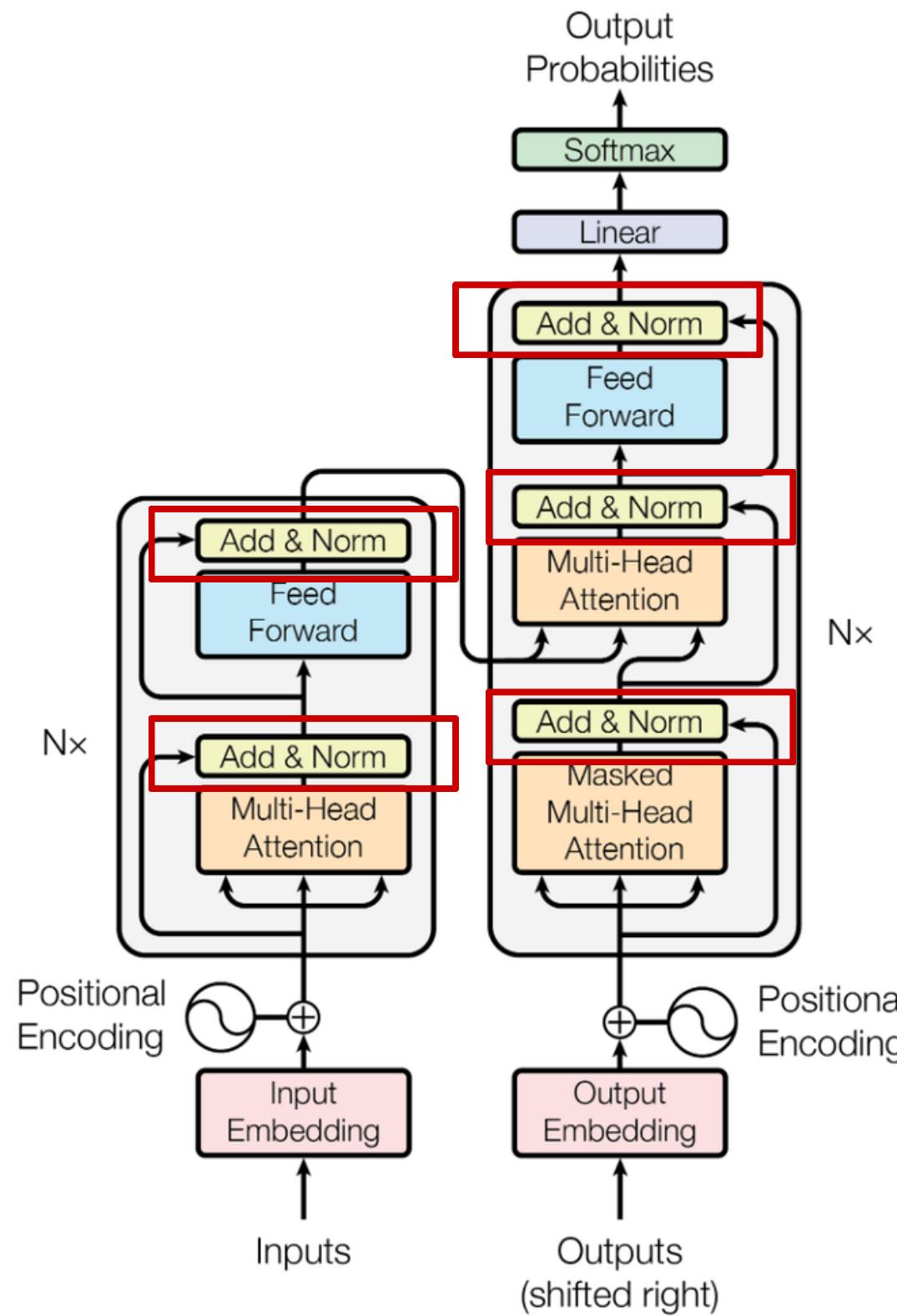
Decoder: Encoder-decoder attention



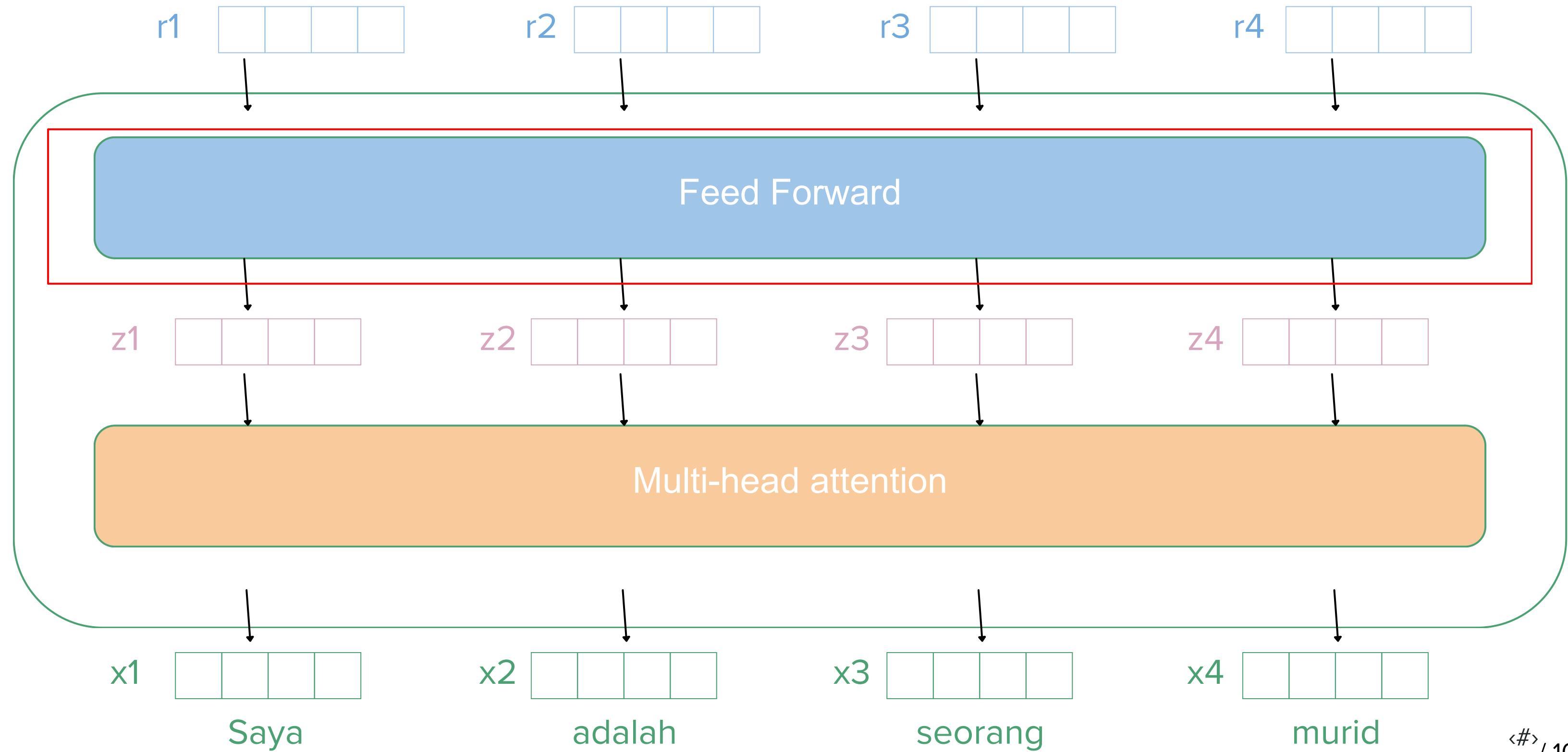
— Any Question Guys ~

— Transformers Add & Norm

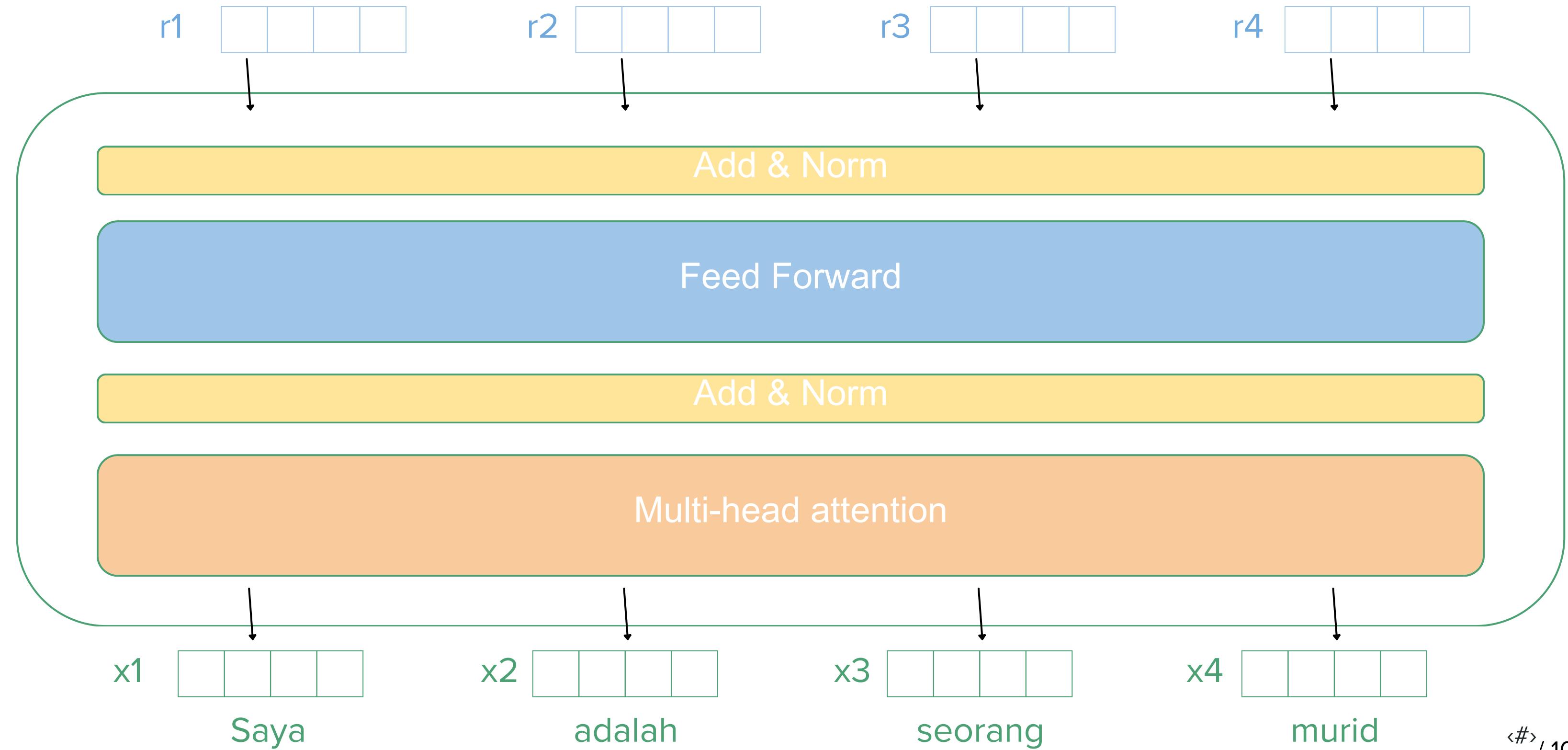
Add & Norm



Add & Norm



Add & Norm



Add & Norm

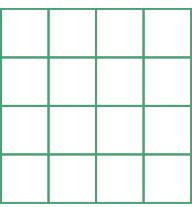
x1
Saya

x2
adalah

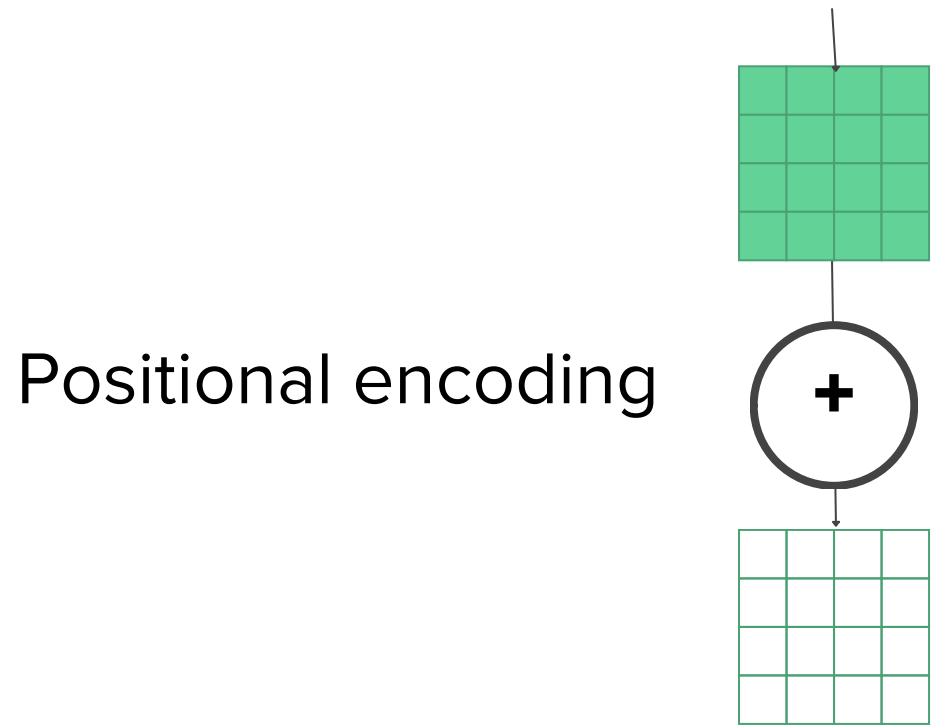
x3
seorang

x4
murid

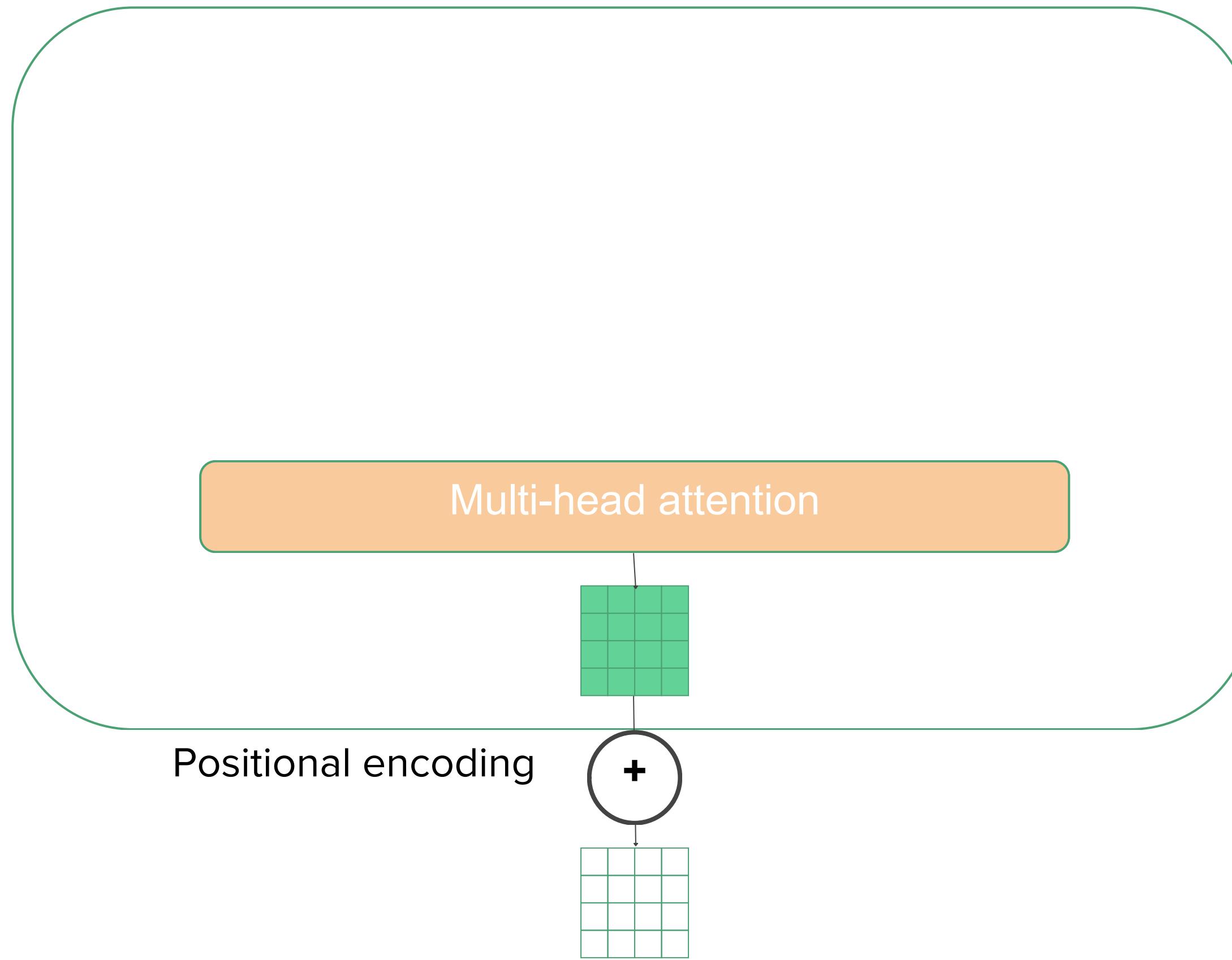
Add & Norm



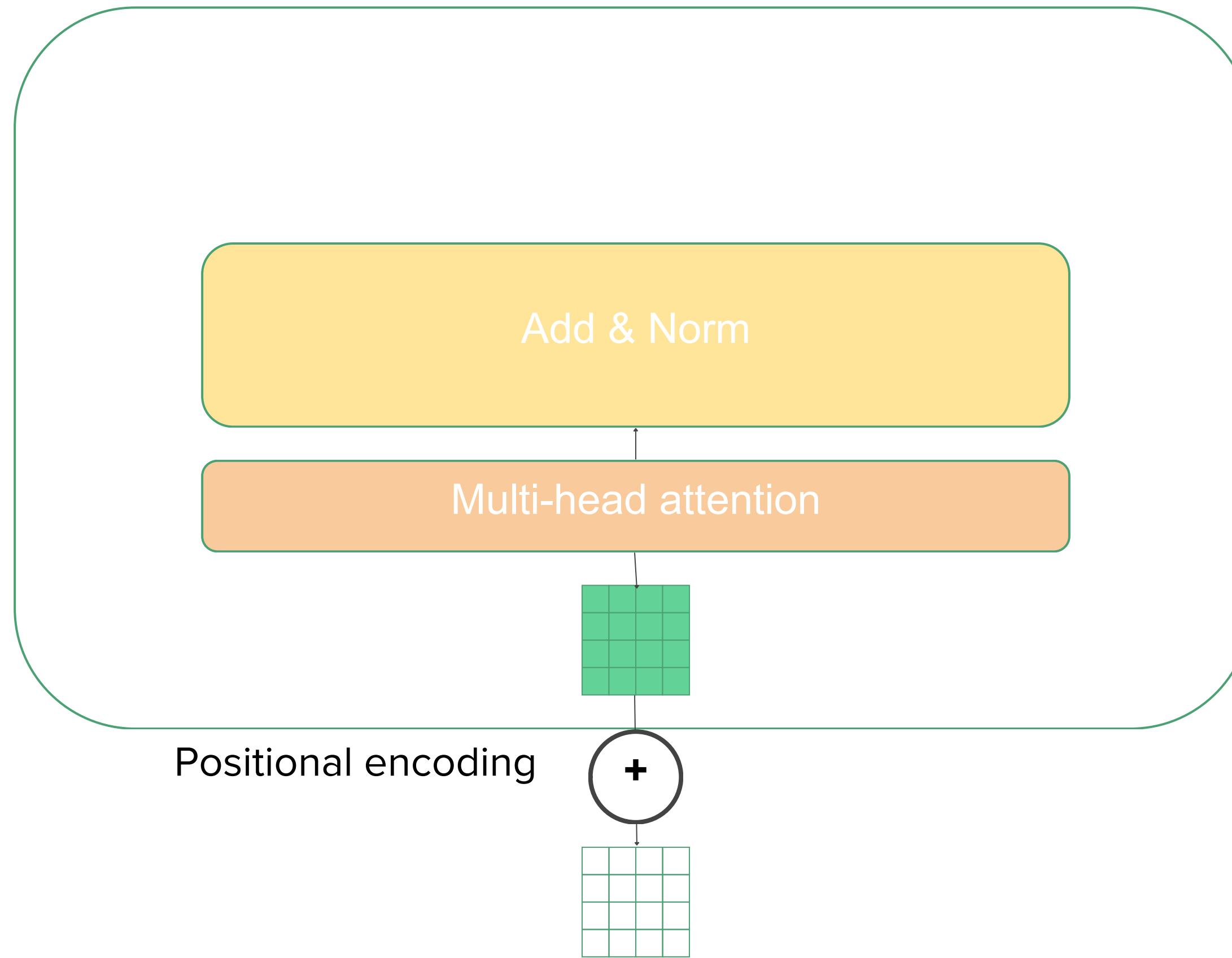
Add & Norm



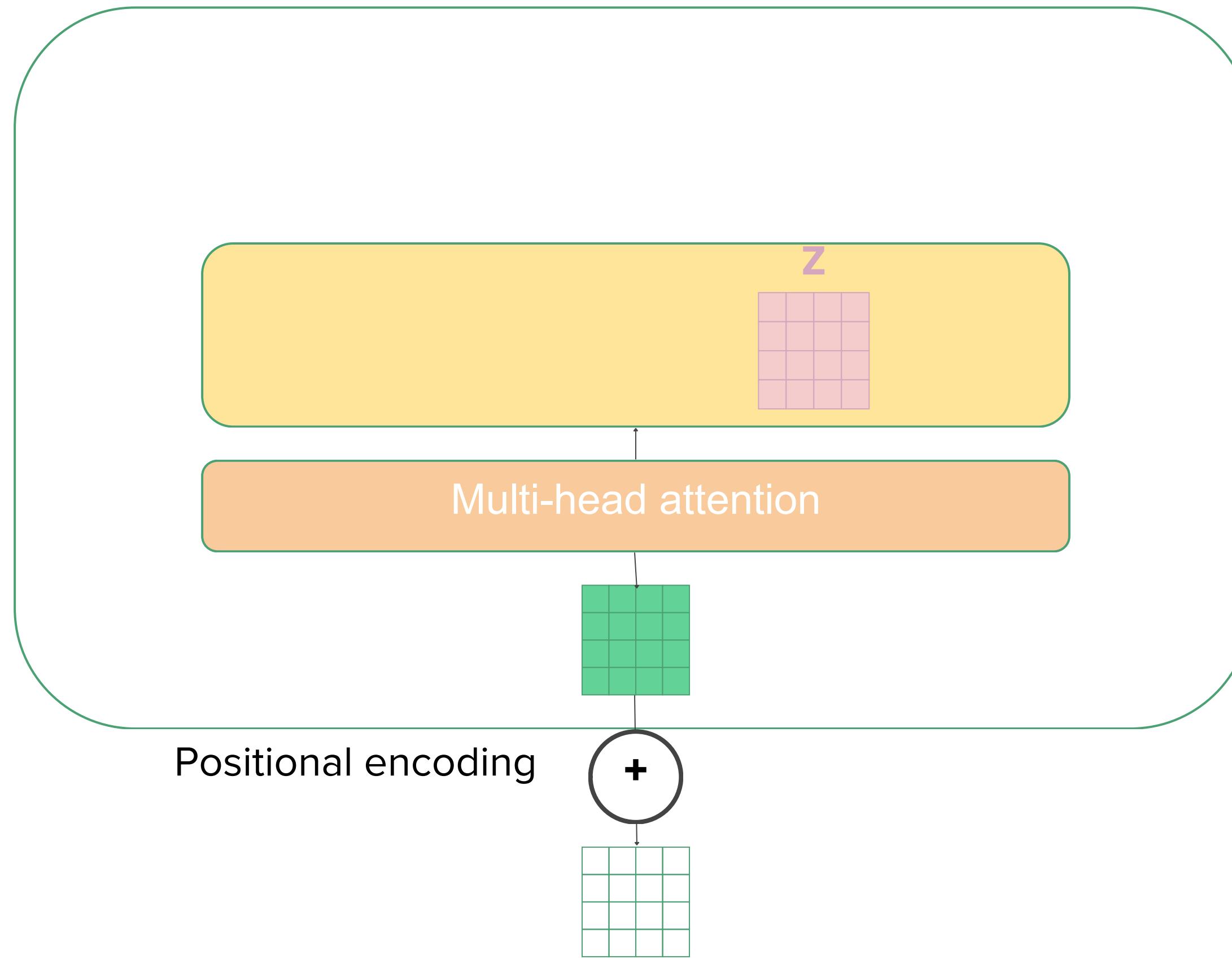
Add & Norm



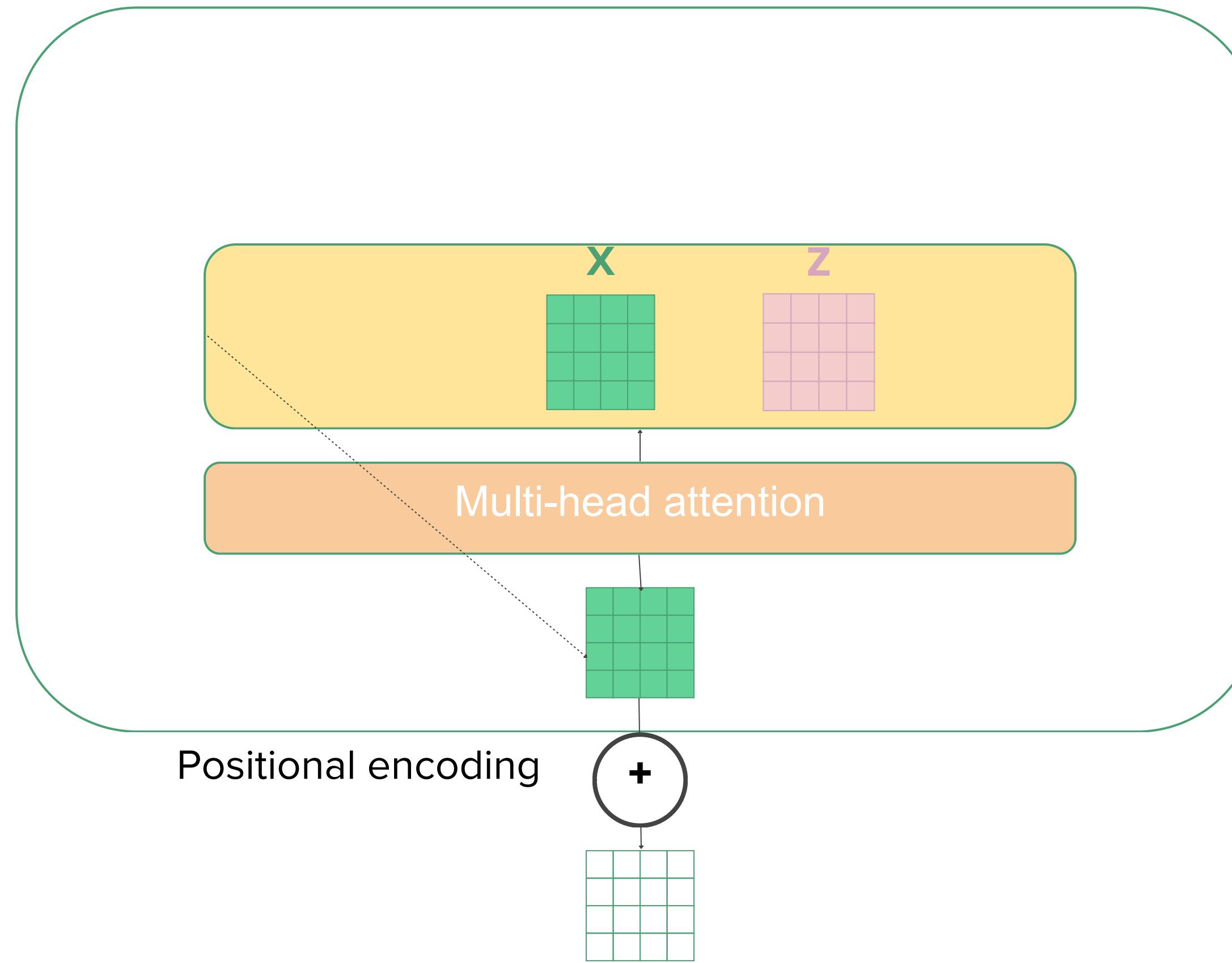
Add & Norm



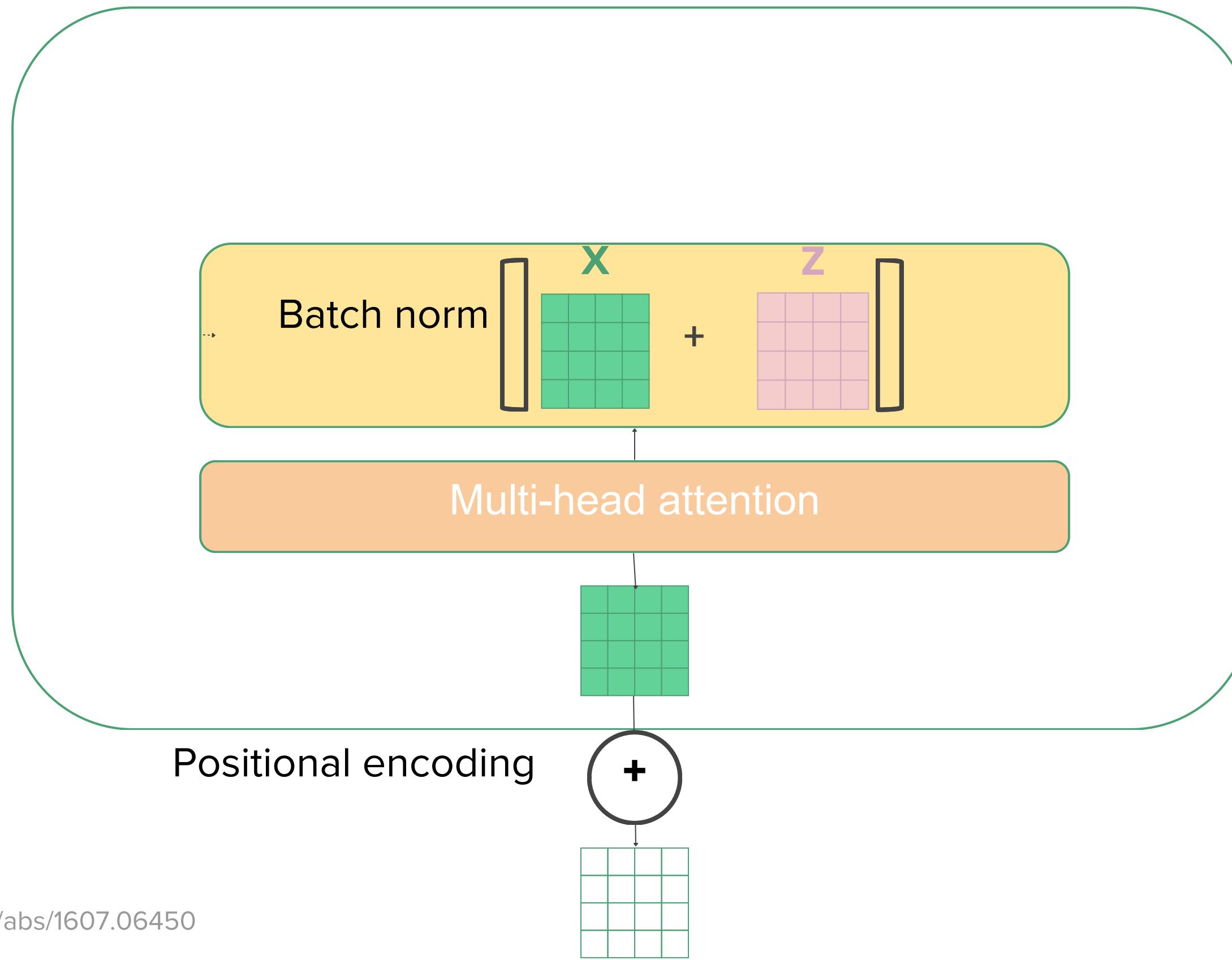
Add & Norm



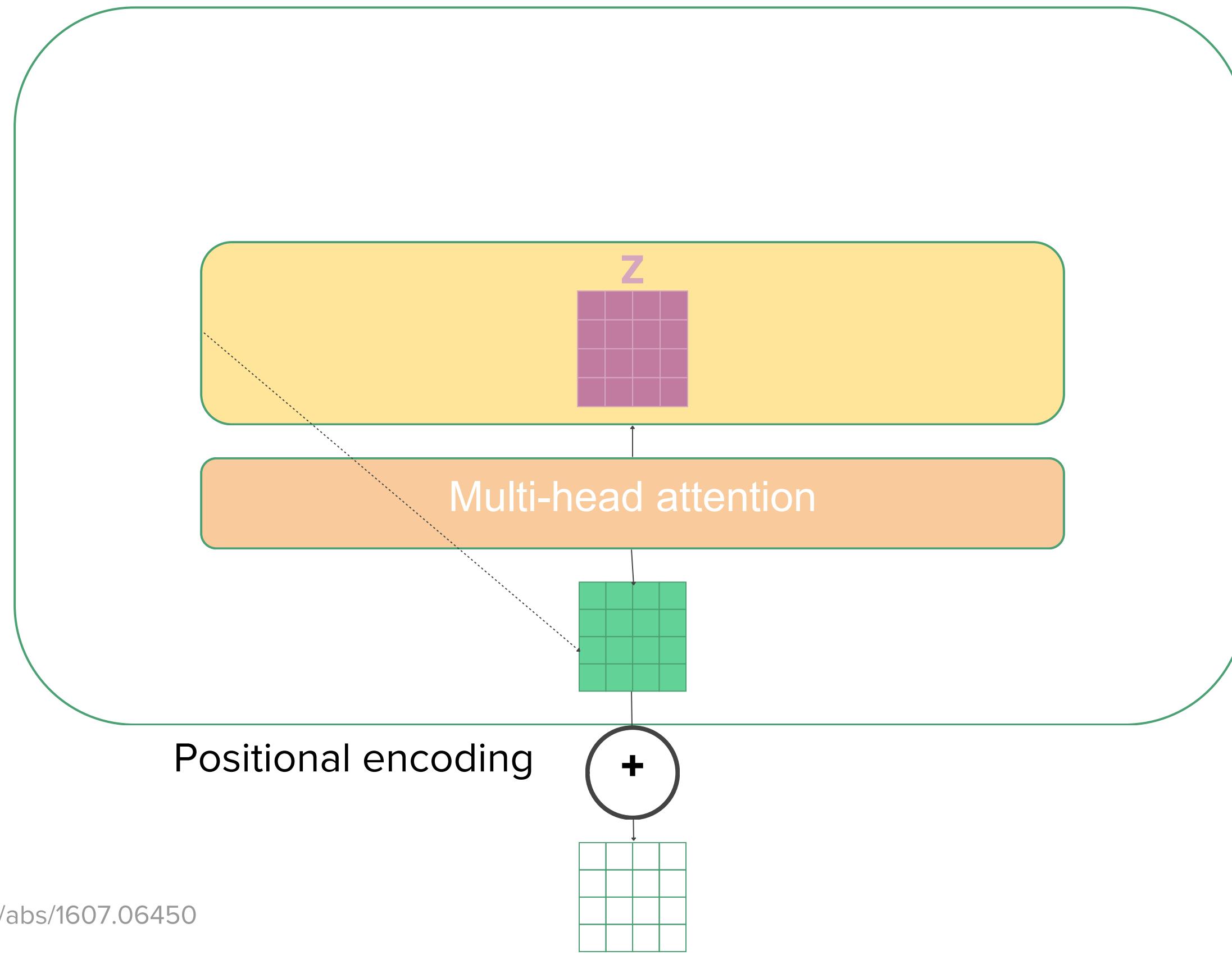
Add & Norm



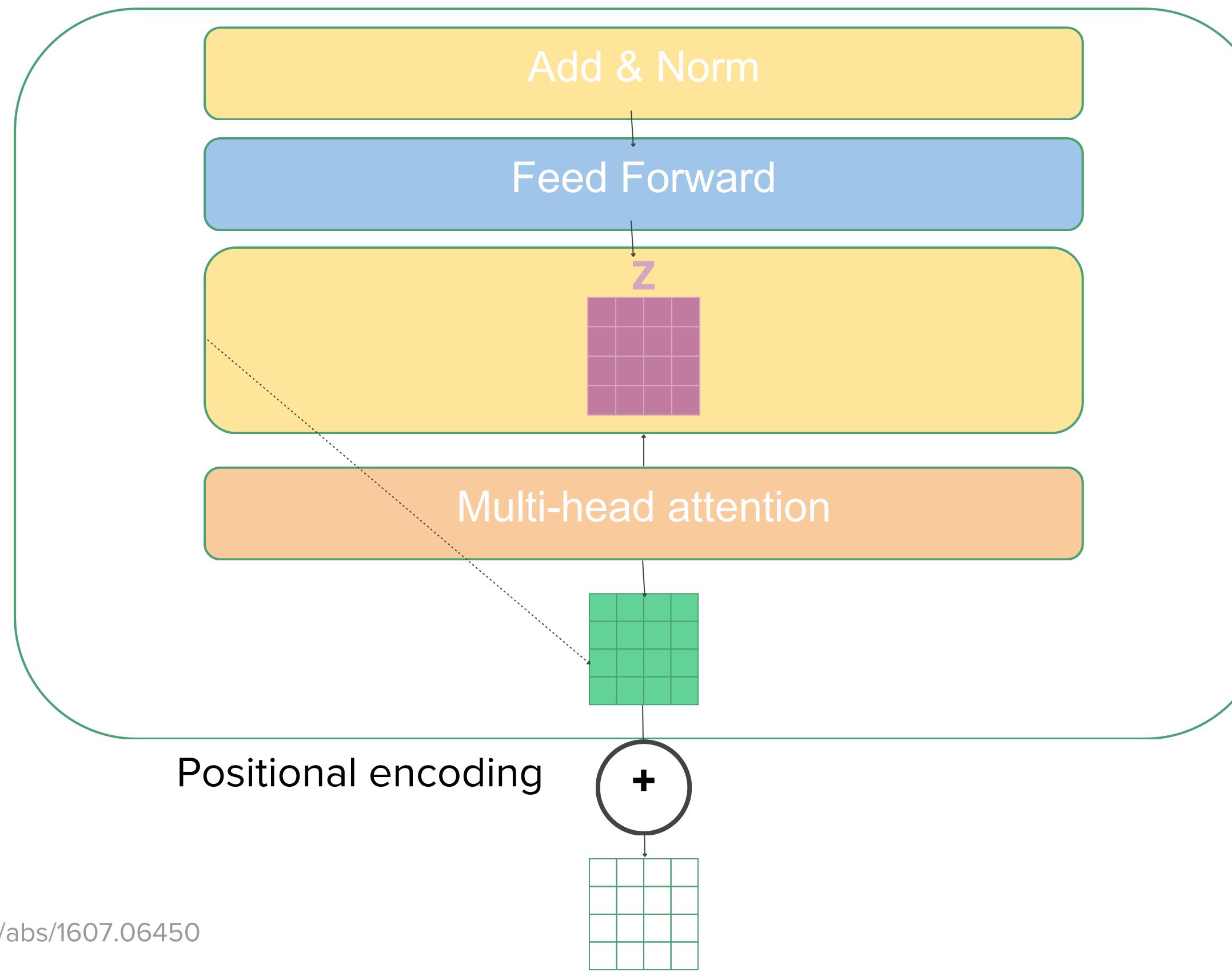
Add & Norm



Add & Norm



Add & Norm



— Any Question Guys ~

— Transformers Performance

Transformers: Performance

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

| Model | BLEU | | Training Cost (FLOPs) | |
|---------------------------------|-------|-------|-----------------------|---------------------|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [18] | 23.75 | | | |
| Deep-Att + PosUnk [39] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [38] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [9] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [32] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [38] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9] | 26.36 | 41.29 | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | | $3.3 \cdot 10^{18}$ |
| Transformer (big) | 28.4 | 41.8 | | $2.3 \cdot 10^{19}$ |

Transformers: Performance

But, <https://arxiv.org/abs/1808.03867> performs better for seq2seq task

The screenshot shows a Cornell University logo and the text "Cornell University" on a black header bar. To the right, it says "the Sim". Below the header, the URL "arXiv.org > cs > arXiv:1808.03867" is on a red bar, along with a "Search..." input field and "Help | Advanced" links. The main content area has a grey header "Computer Science > Computation and Language". Below it, the text "[Submitted on 11 Aug 2018 (v1), last revised 1 Nov 2018 (this version, v3)]" is in blue. The title "Pervasive Attention: 2D Convolutional Neural Networks for Sequence-to-Sequence Prediction" is in large bold black font. The authors' names "Maha Elbayad, Laurent Besacier, Jakob Verbeek" are in blue. A summary of the paper's content follows, mentioning 2D convolutional neural networks for sequence-to-sequence prediction, outperforming state-of-the-art encoder-decoder systems. At the bottom, there are sections for "Comments", "Subjects", "Cite as", and a link to "arXiv:1808.03867v3 [cs.CL]".

Comments: Accepted at CoNLL 2018

Subjects: **Computation and Language (cs.CL)**

Cite as: [arXiv:1808.03867 \[cs.CL\]](https://arxiv.org/abs/1808.03867)
 (or [arXiv:1808.03867v3 \[cs.CL\]](https://arxiv.org/abs/1808.03867v3) for this version)

Transformers Development

Transformers- XL

<https://arxiv.org/abs/1901.02860>

Two new contribution
from vanilla transformers:

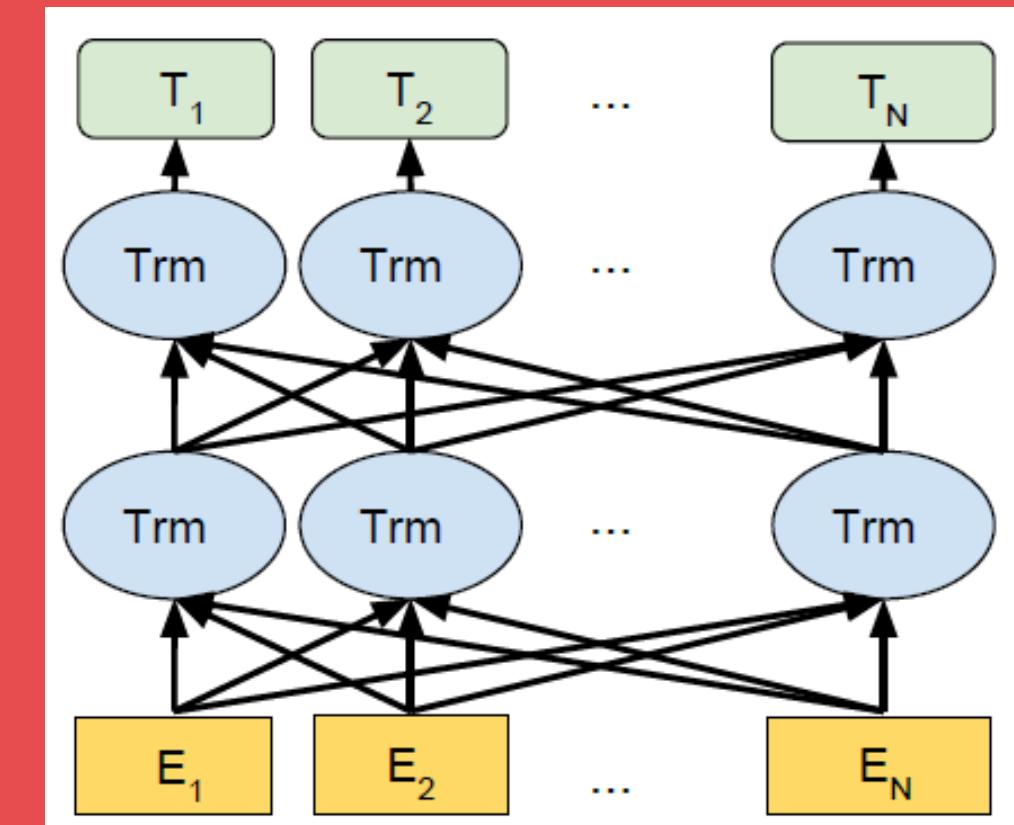
1. Recurrence
mechanism
2. Relative positional
embedding

BERT

Pre-training of Deep
Bidirectional Transformers
for Language
Understanding

Main keys:

1. Masked Language Models (MLM)
2. Next Sentence Prediction

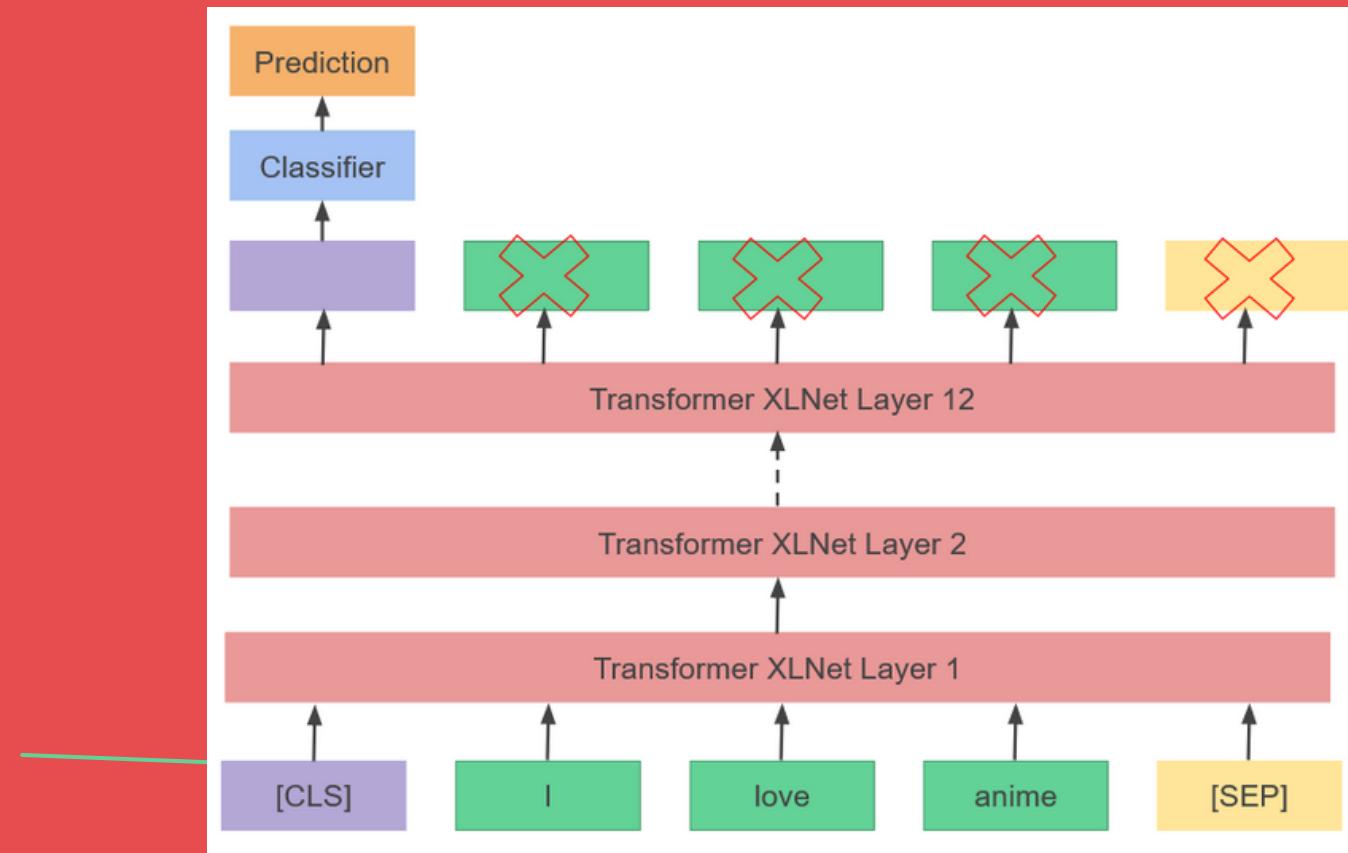


XL-Net

XLNet: Generalized
Autoregressive Pretraining
for Language Understanding

Main keys:

1. Permutation Language Modelling
2. Two-Stream Self-Attention



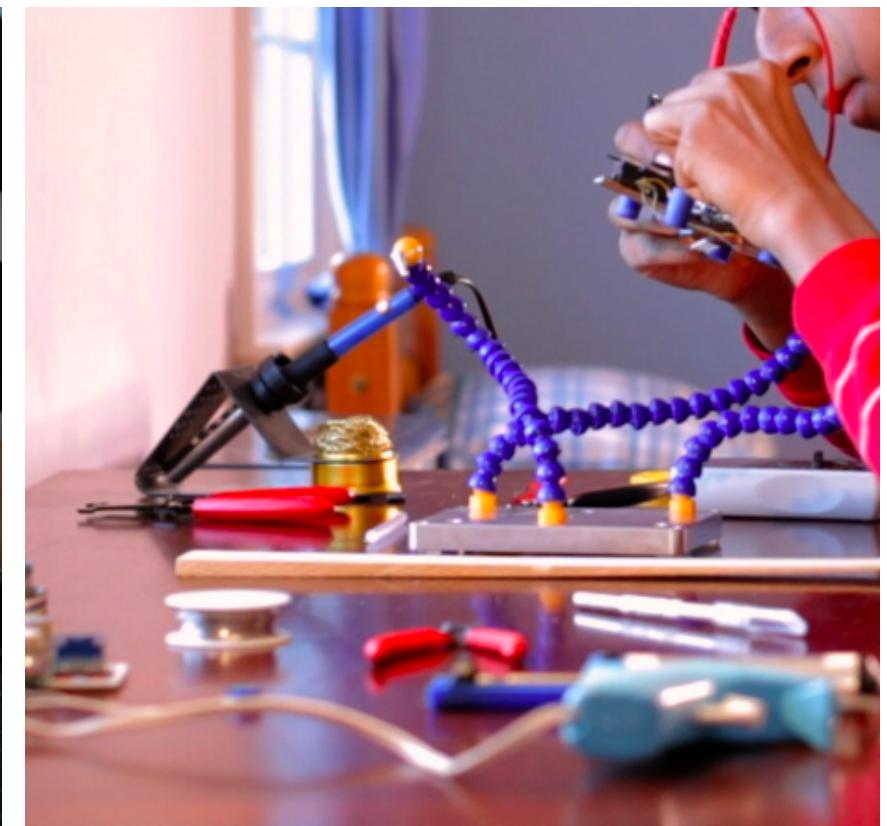
LET'S CODE!

Google Colaboratory
Machine translation using
Transformers:

https://colab.research.google.com/drive/1CoT6zSiFy0IrCEIFJyxeUwcu4peuPK_j?usp=sharing

Transformers from scratch:

<https://www.tensorflow.org/tutorials/text/transformer>



Experiment

- We will use OpenNMT library
- 10.000 English-German sentence pair

<https://opennmt.net/OpenNMT-tf/index.html>

The screenshot shows the documentation for OpenNMT-tf version 2.11. The left sidebar has a dark red background and contains navigation links for 'GETTING STARTED' (Quickstart, Installation) and 'CONFIGURATION' (Model, Parameters, Data, Vocabulary, Tokenization, Embeddings, Alignments). The main content area has a white background and features a header with the 'OpenNMT-tf' logo and version '2.11'. It includes a search bar labeled 'Search docs', a breadcrumb trail 'Docs » Index', and a welcome message: 'Welcome to the OpenNMT-tf documentation! The documentation describes how to install, configure, and use the project.' Below this, there is a note about reporting issues on GitHub and a section titled 'Index' with sub-sections 'Getting started' and 'Configuration', each with a list of links.

Docs » Index

Welcome to the OpenNMT-tf documentation! The documentation describes how to install, configure, and use the project.

If you find a problem with an information presented on this website, please [open an issue on the GitHub repository](#).

Index

Getting started

- [Quickstart](#)
- [Installation](#)

Configuration

- [Model](#)
- [Parameters](#)
- [Data](#)
- [Vocabulary](#)
- [Tokenization](#)
- [Embeddings](#)
- [Alignments](#)

Reference

1. Paper: <https://arxiv.org/abs/1706.03762>
2. Transformers illustration: <http://jalammar.github.io/illustrated-transformer/>
3. <https://www.youtube.com/watch?v=TQQlZhBC5ps>
4. <https://www.youtube.com/watch?v=KN3ZL65Dze0&t=133s>
5. https://kazemnejad.com/blog/transformer_architecture_positional_encoding/

— Any Question Guys ~

Terimakasih!