



# TRANSFORMERS II

Sandeep Soni

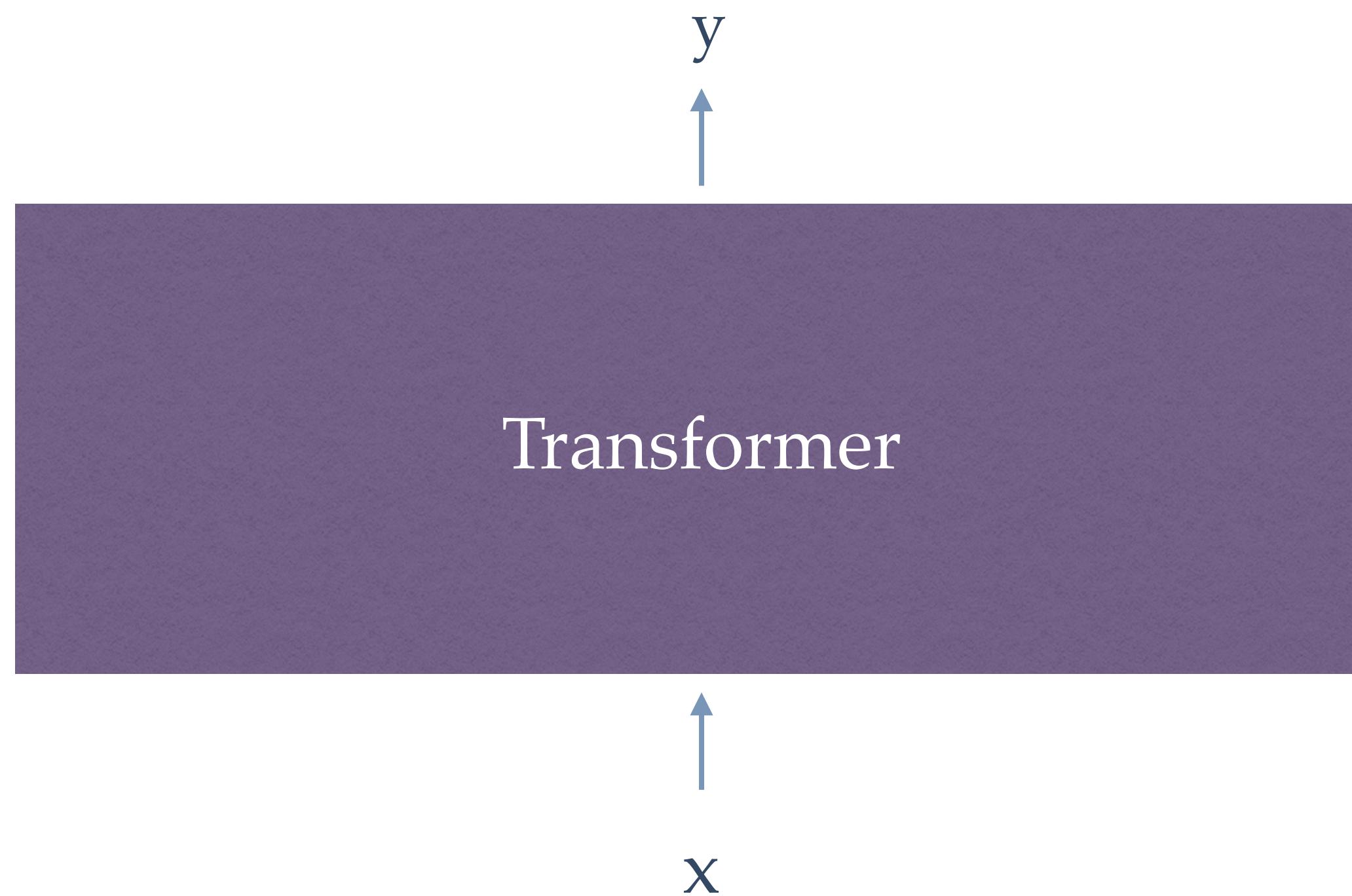
---

03/05/2024

## STORY SO FAR

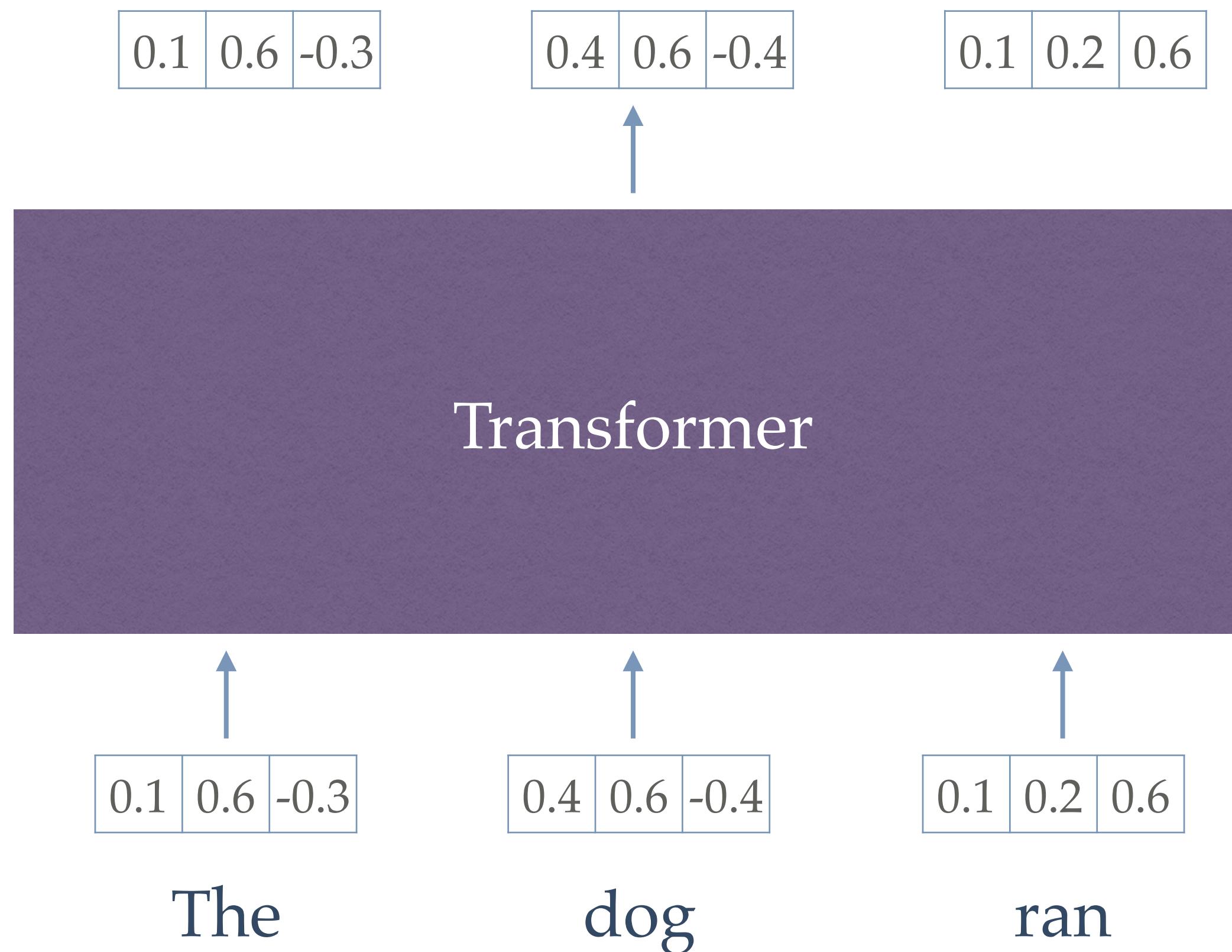
- Language model estimates  $P(x)$  for a sequence  $x$
- Count and normalize language models (e.g., N-gram LMs) assume fixed-length conditioning context
- RNN LM can condition on arbitrarily long context but has sequential computation

# TRANSFORMERS



- Transform an input sequence of vectors into an output sequence of vectors

# TRANSFORMERS

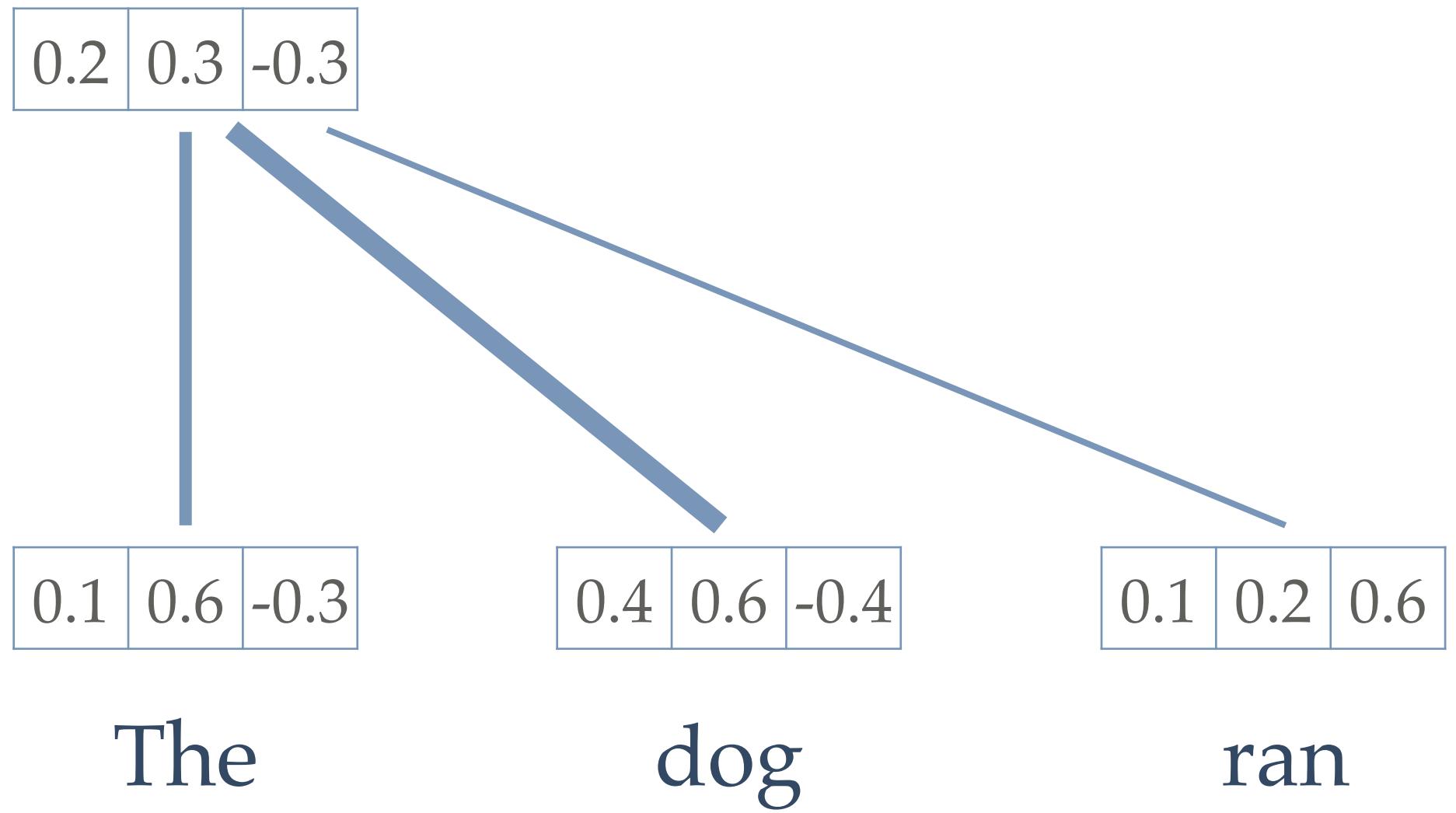


- Transform an input sequence of vectors into an output sequence of vectors

## KEY IDEAS

- A vector associated for every token in a sequence that depends on every other token in the sequence.
- Multiple transformations of vector sequences.
- Parallel computation

# SELF-ATTENTION



- Create a token vector as a weighted combination of input token vectors
- Weights are the importance or attention to other tokens

0.2	0.3	-0.3
-----	-----	------

Value

$$v \in \mathbb{R}^{768} = eW_V$$

$$W_V \in \mathbb{R}^{768 \times 768}$$

0.1	0.2	-0.3
-----	-----	------

Key

$$k \in \mathbb{R}^{96} = eW_K$$

$$W_K \in \mathbb{R}^{768 \times 96}$$

0.3	0.1	-0.5
-----	-----	------

Query

$$q \in \mathbb{R}^{96} = eW_Q$$

$$W_Q \in \mathbb{R}^{768 \times 96}$$

0.1	0.6	-0.3
-----	-----	------

Previous value

$$e \in \mathbb{R}^{768}$$

The

$$a_{1j} = \text{softmax}(\text{dot}(q_1, k_j))$$

0.13

0.83

0.04

Attention weights

0.3

1.9

-0.7

Compatibility scores

$q_1 \cdot k_1$

$q_1 \cdot k_2$

$q_1 \cdot k_3$

0.1	0.6	-0.3
-----	-----	------

The

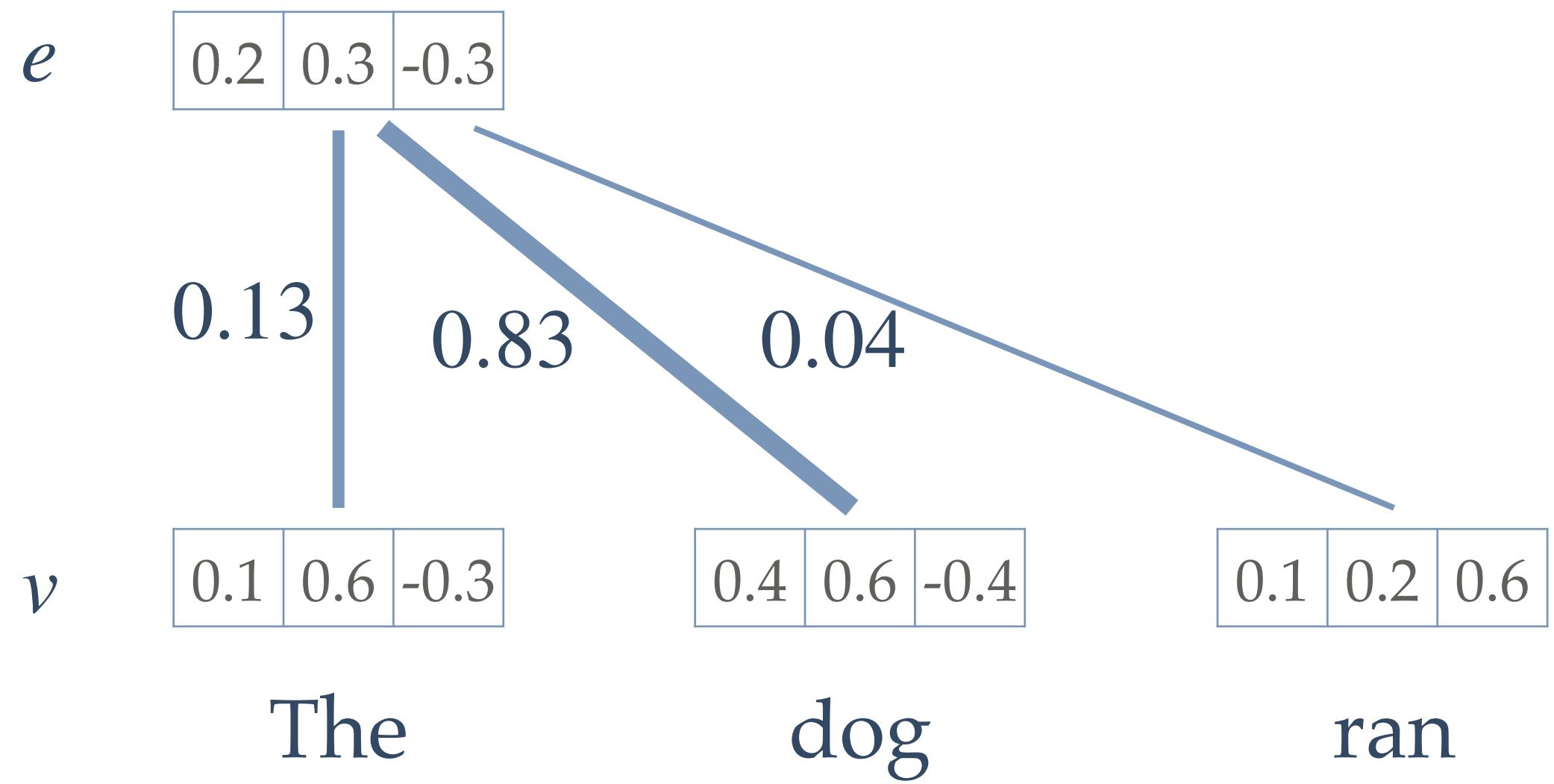
0.4	0.6	-0.4
-----	-----	------

dog

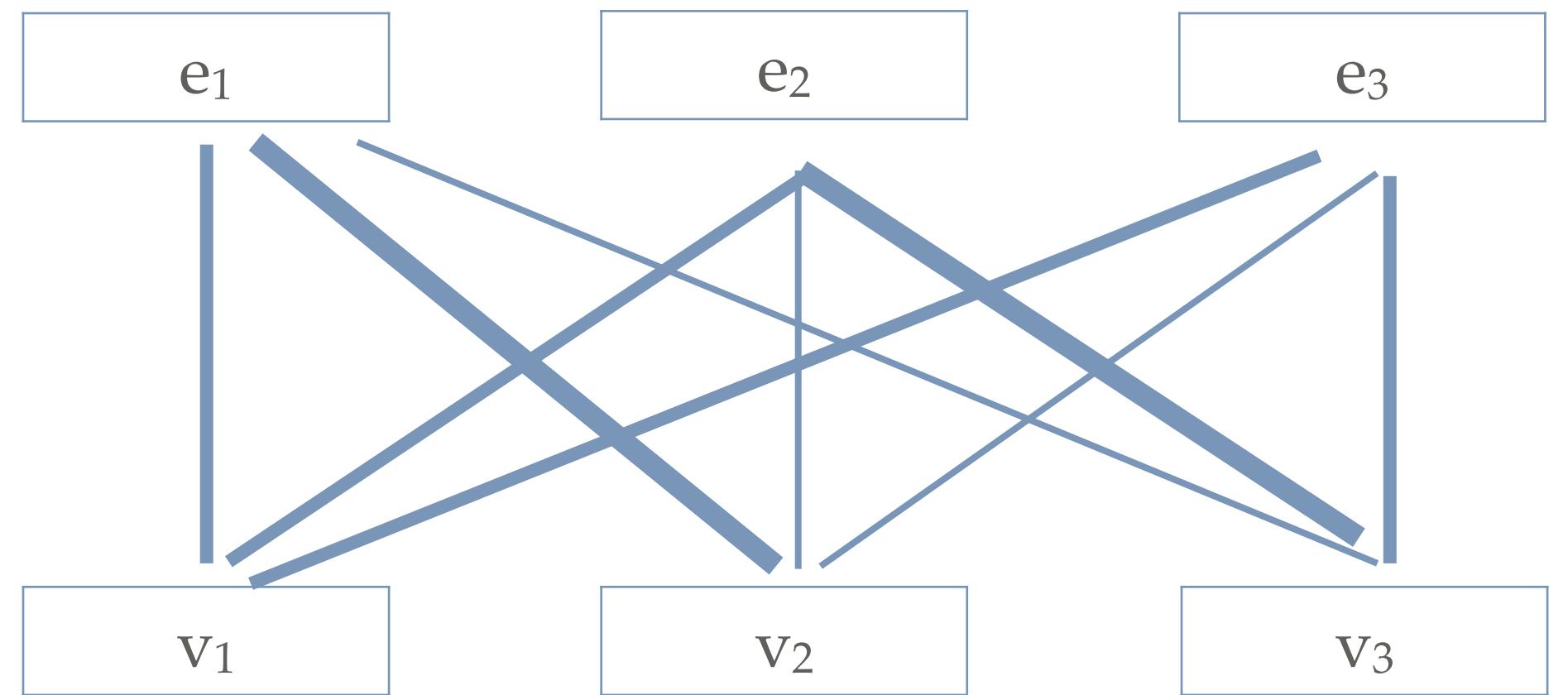
0.1	0.2	0.6
-----	-----	-----

ran

- Query and key vectors are used to calculate compatibility or attention scores.

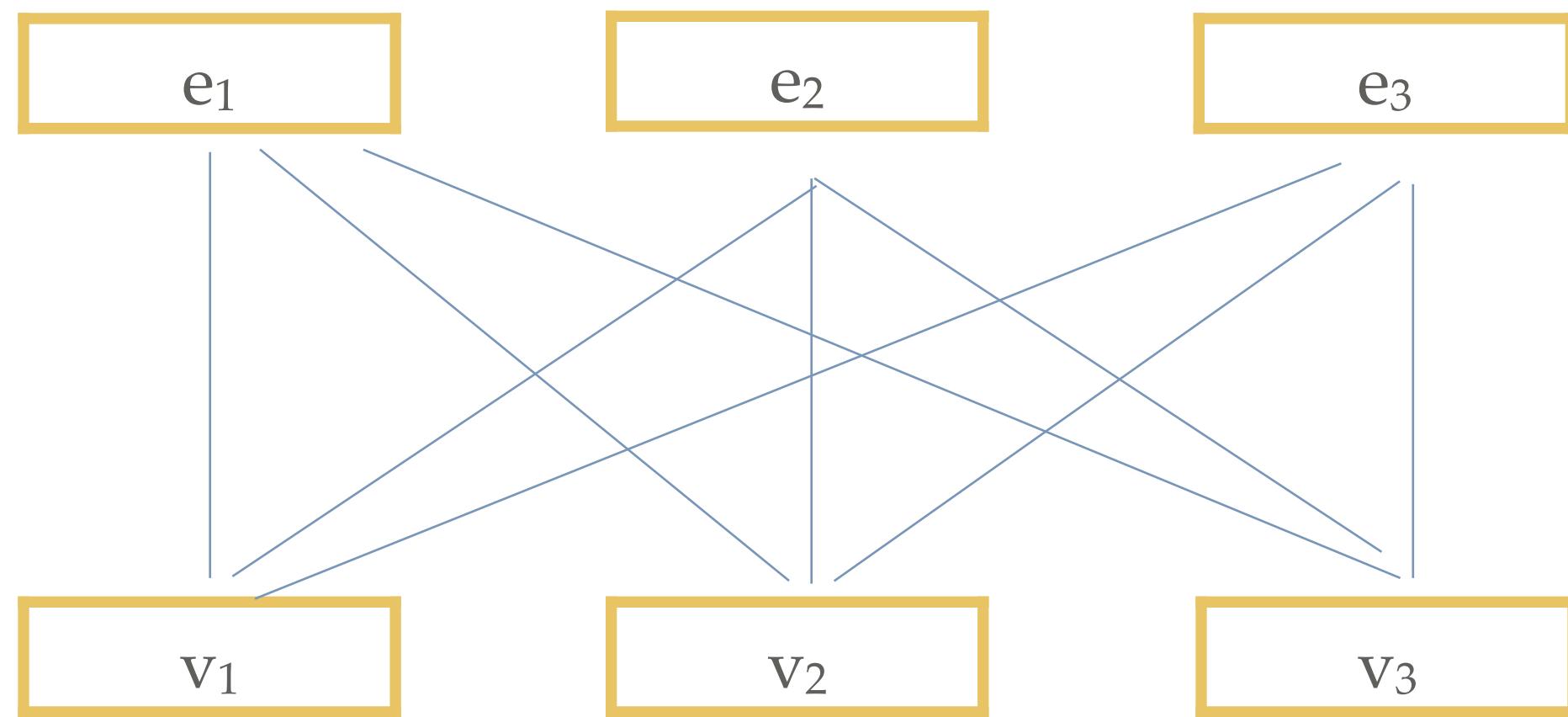


- Attention weights are used to combine value vectors to create a new vector



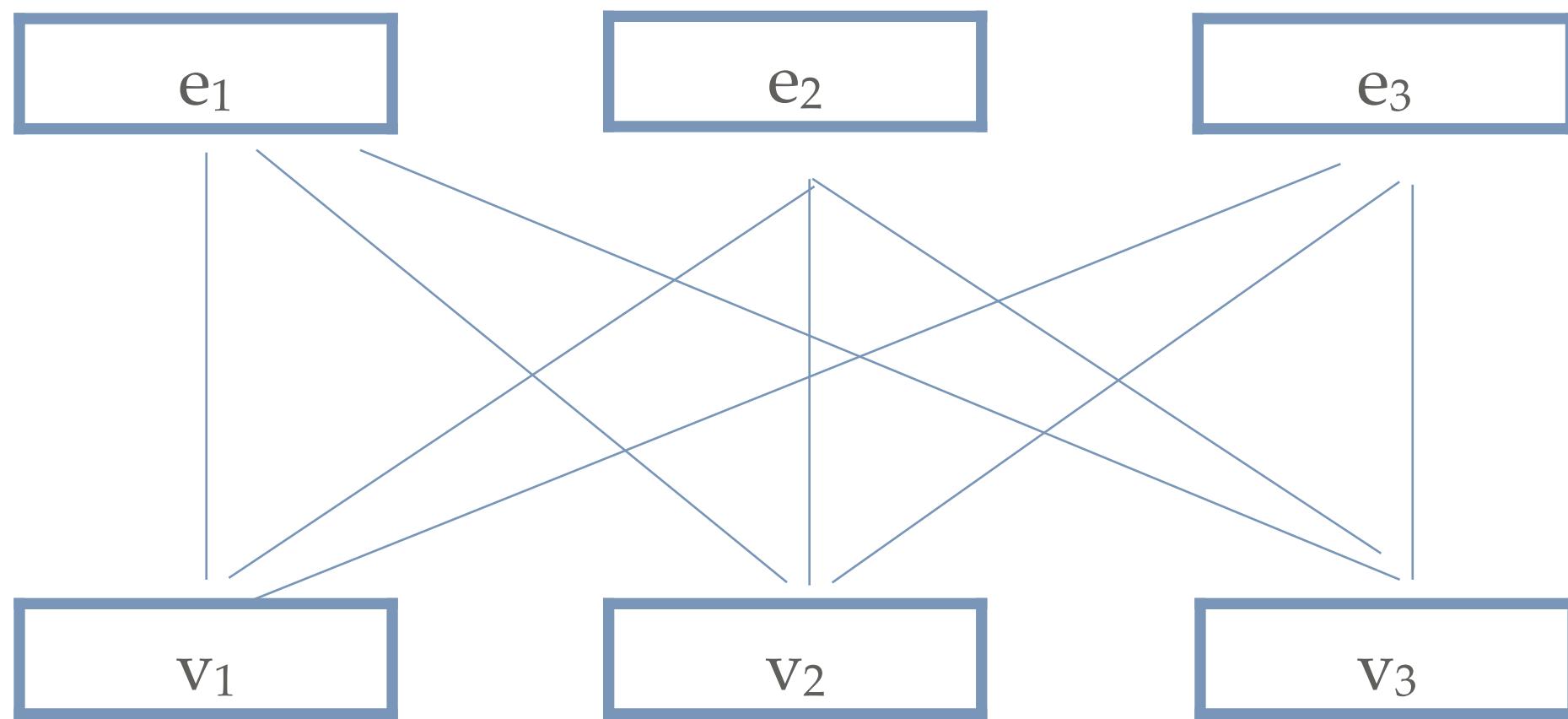
- Compatibility for all the pairs in a sequence is calculated

# MULTI-HEAD ATTENTION



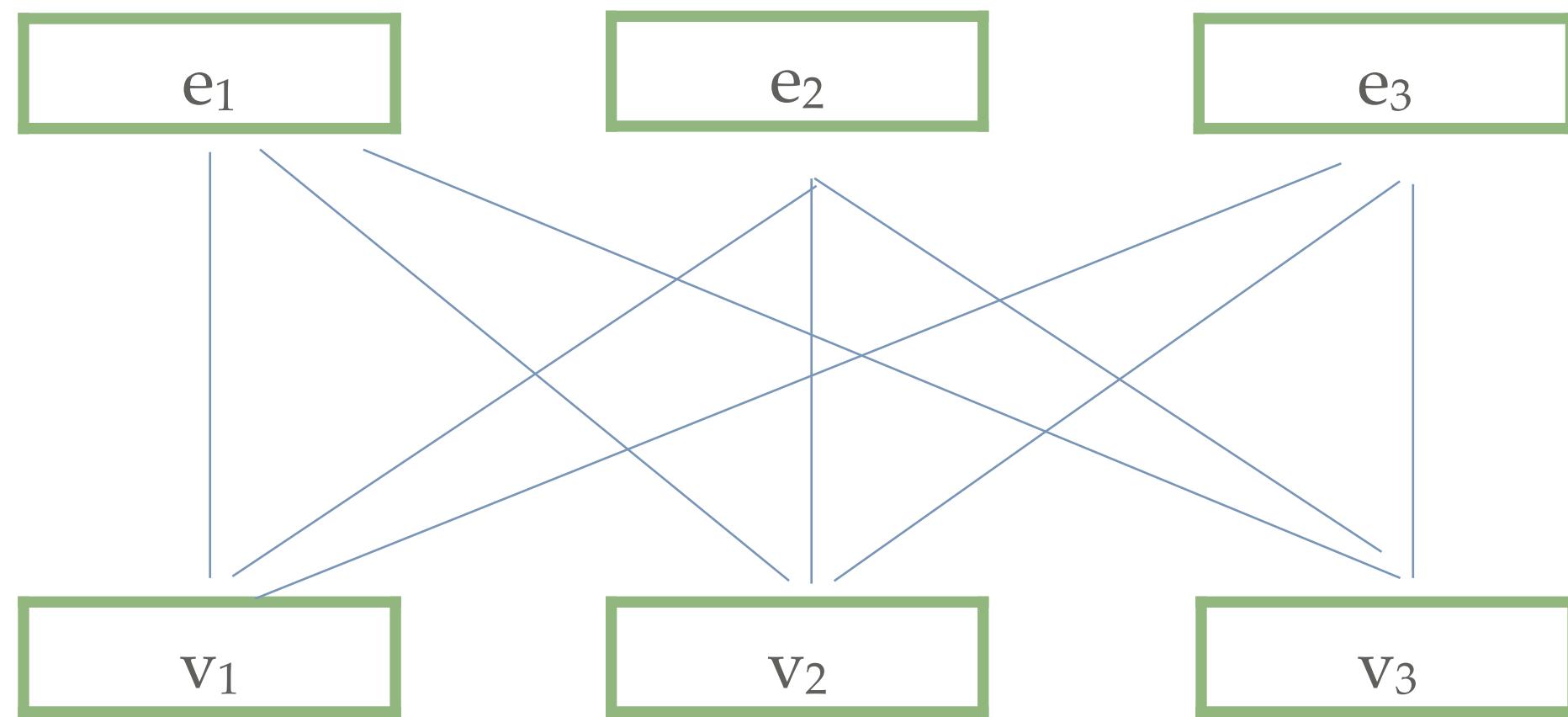
- Each attention head can have its own query, key, and value vectors such that every attention head gives one view of the data

# MULTI-HEAD ATTENTION



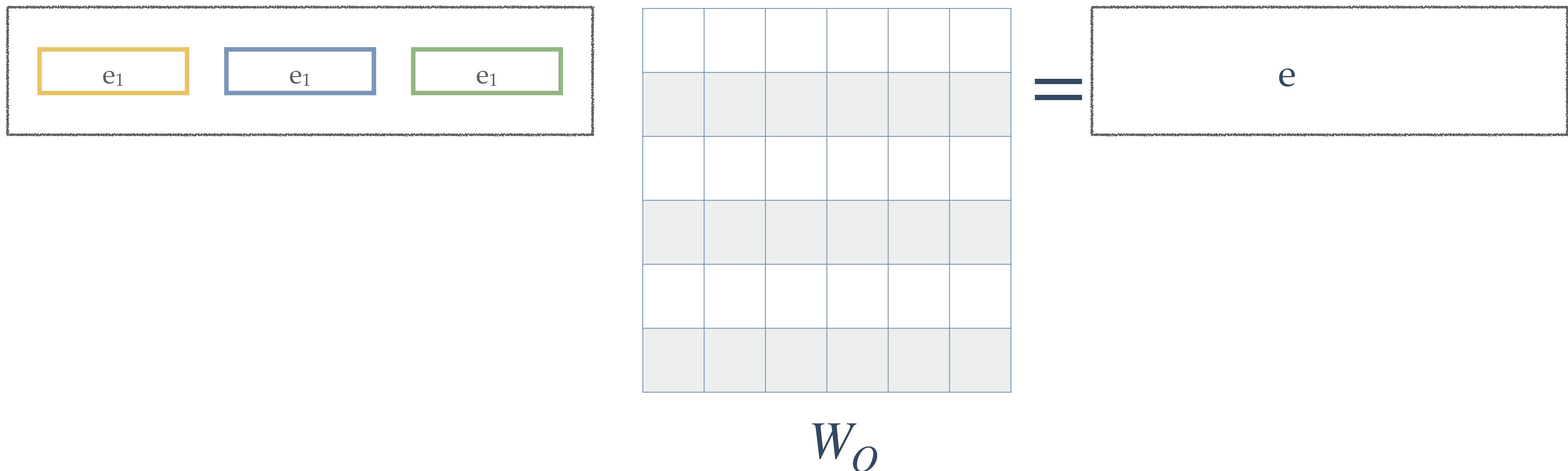
- Each attention head can have its own query, key, and value vectors such that every attention head gives one view of the data

# MULTI-HEAD ATTENTION

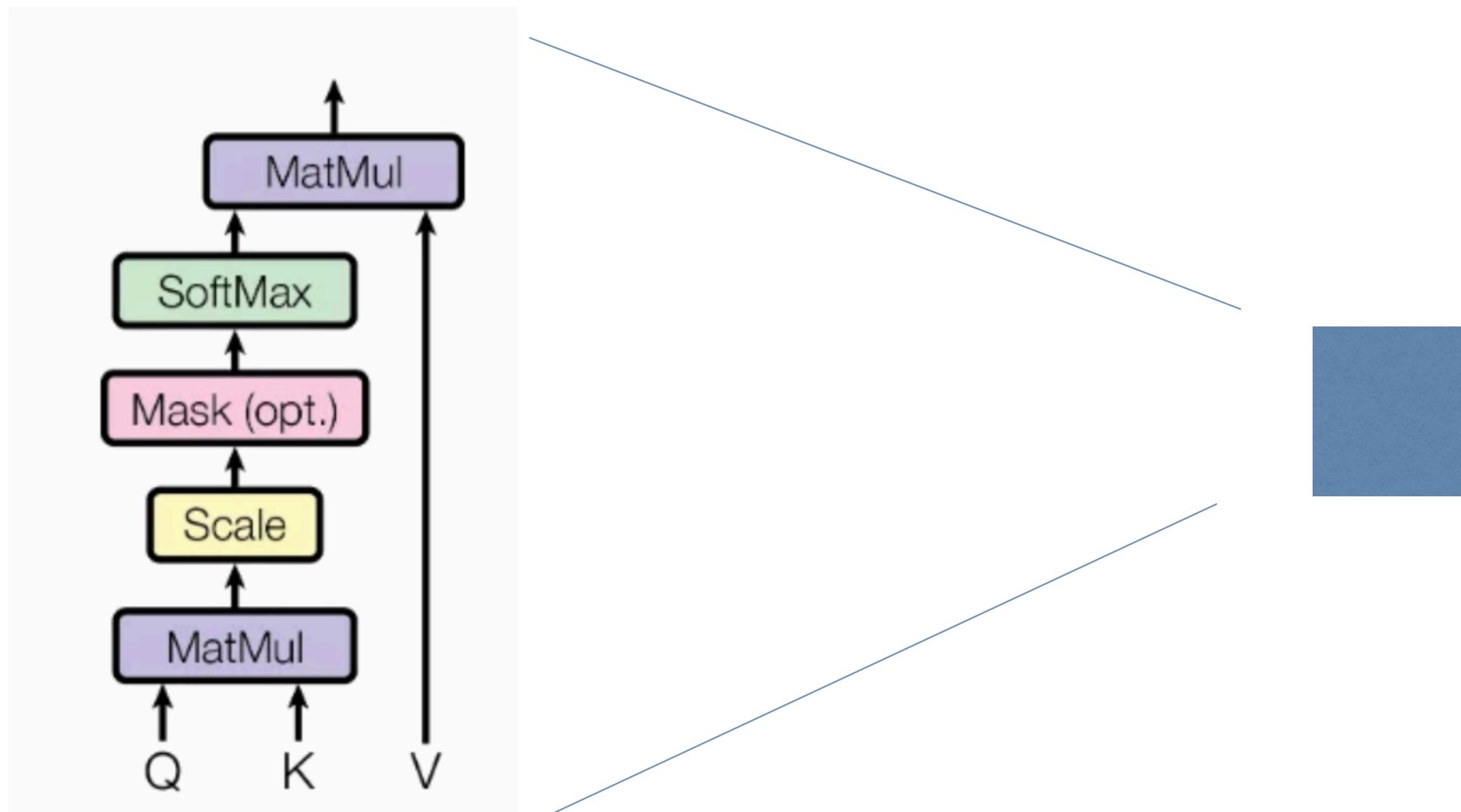


- Each attention head can have its own query, key, and value vectors such that every attention head gives one view of the data

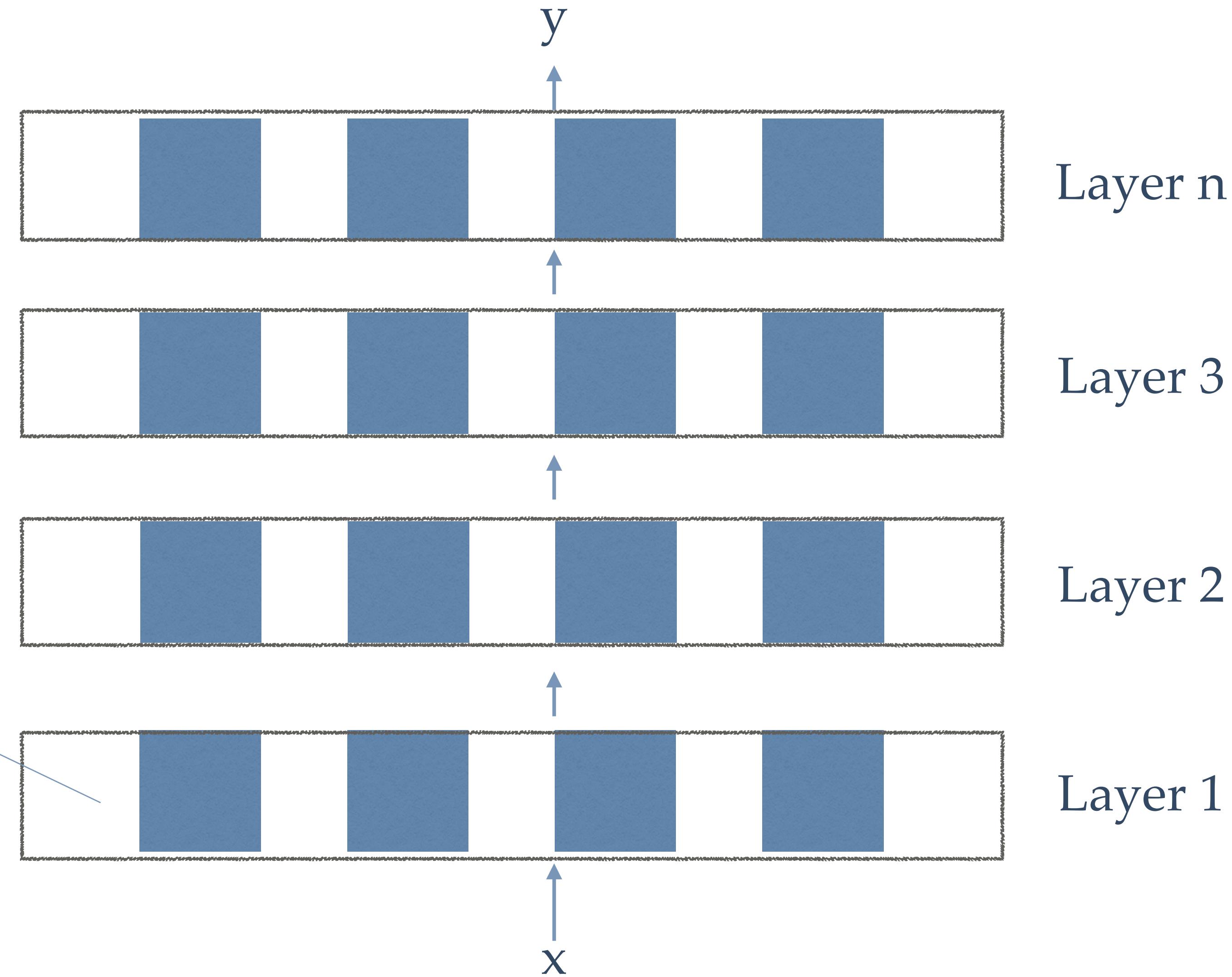
# MULTI-HEAD ATTENTION



# SINGLE HEAD ATTENTION



Self attention units  
within an attention block



# TRANSFORMERS

- There are some more transformations within an attention block (layer normalization, feedforward network)
- To account for the position of each token, positional embeddings are used in place to augment token embeddings

# LARGE LANGUAGE MODELS

- Most modern LLMs are based on the transformer architecture (e.g., BERT, GPT, etc)
- Many layers, high dimensional vector representations, and large corpora for pretraining are all hallmarks of contemporary LLMs

# VARIATIONS OF LANGUAGE MODELING

Masked LM:

$$P(x) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

Encoder-Decoder LM:

$$P(y) = \prod_{i=1}^n P(y_i | y_1, \dots, y_{i-1}, x)$$

Causal LM:

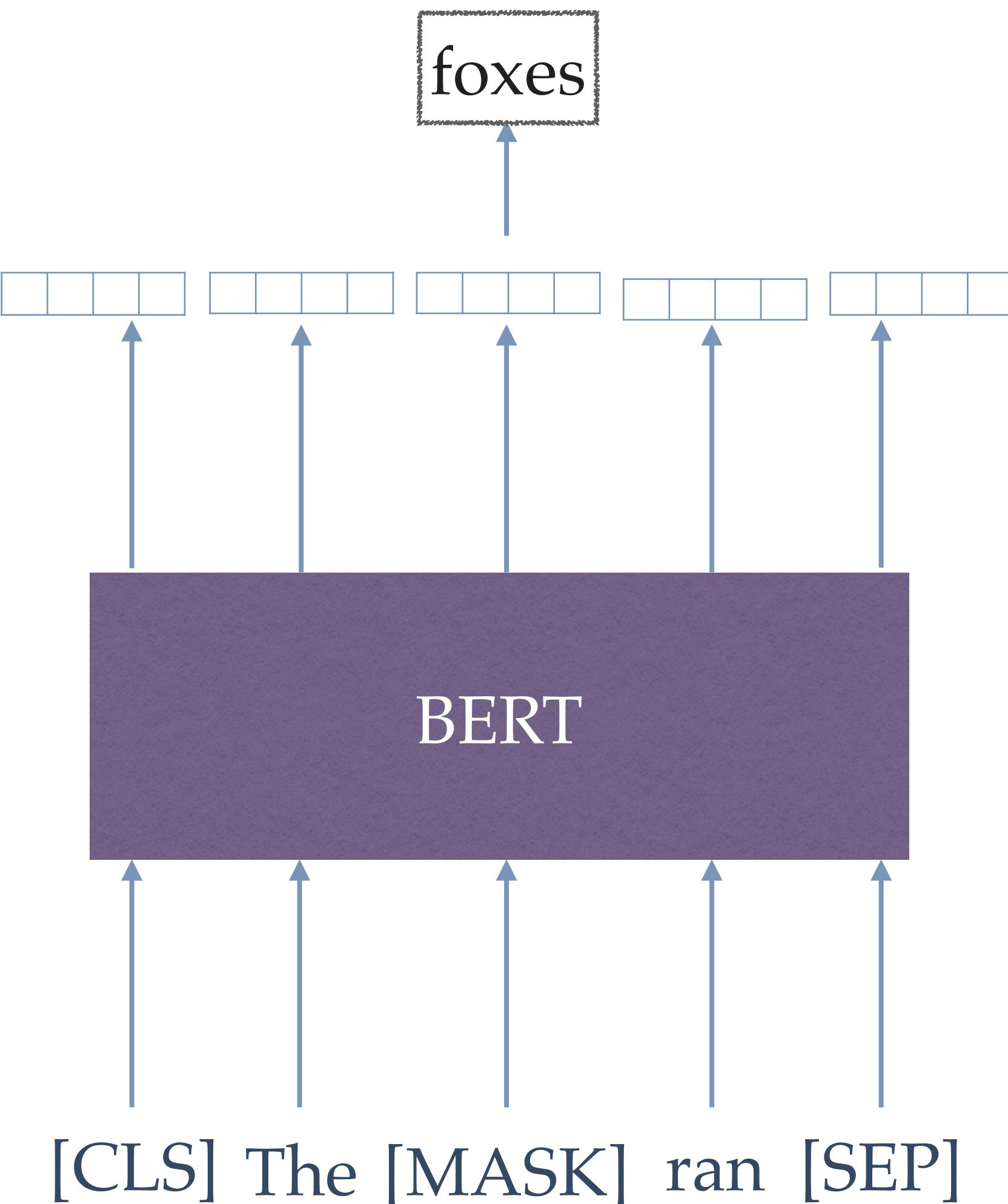
$$P(x) = \prod_{i=1}^n P(x_i | x_1, x_2, \dots, x_{i-1})$$

# MASKED LANGUAGE MODELING

Fill in the blank by  
using the surrounding  
context

The \_\_\_\_\_ ran

$$P(w_t | \neg w_t)$$



BERT is a  
language  
model  
trained to  
predict the  
missing  
word

# BERT

- Masked LM

$$P(x) = \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

- The \_\_\_\_\_ ran →



- Transformer based model
  - trained on 512 sized contexts
  - 12 or 24 transformer blocks
  - 768 or 1024 sized representations
  - pretrained on billions of tokens from wikipedia and books

# T5

- Encoder-Decoder LM: 
$$P(y) = \prod_{i=1}^n P(y_i | y_1, \dots, y_{i-1}, x)$$
- The ~~foxes ran~~ in the forest ~~after~~ seeing the humans



- Pretrained on 750GB of web text

# GPT

- Causal LM:
$$P(x) = \prod_{i=1}^n P(x_i | x_1, x_2, \dots, x_{i-1})$$
- Also called as left-to-right or autoregressive language modeling
- Fill in the blank but always predict the next word in the sequence (e.g., The foxes \_\_\_\_)
- GPT-3 has 175B parameters and trained on 570GB of web text, books, wikipedia