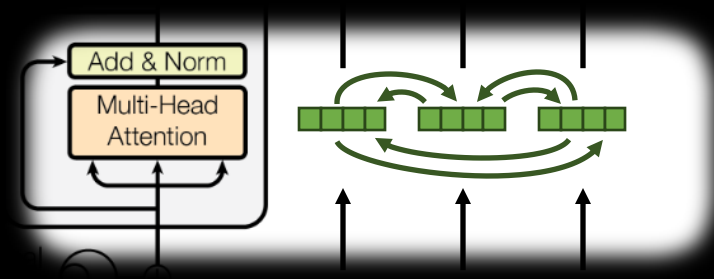


Neural Computation

Attention

Thanks to Constantin Pape and Alex Ecker!



Attention

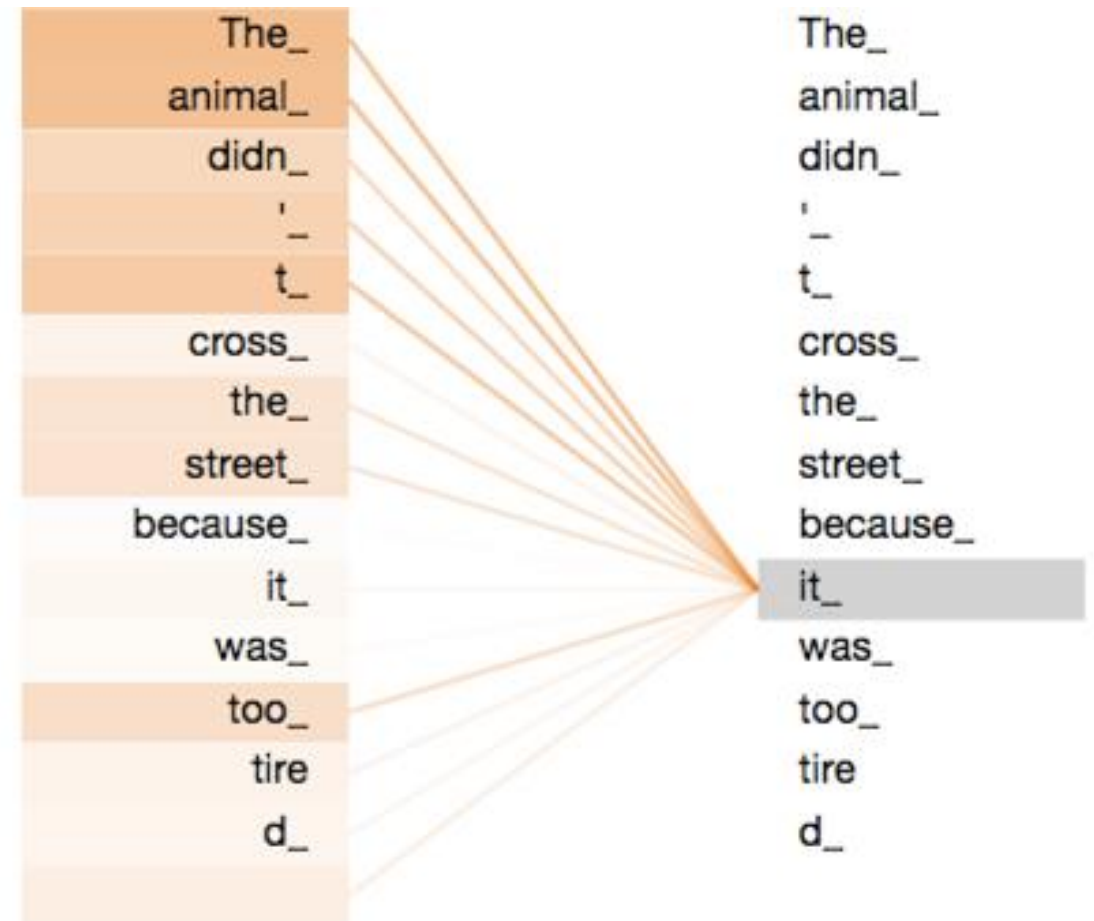
Self-Attention

Self-attention:

- every element in sequence can influence every other element
- learn weighting (“attention”) for each pair of elements

Attention: learnable pairwise weighting that depends on other sequence

Self-Attention: use input sequence for attention



Key-Value Pairs

JSON-Files:

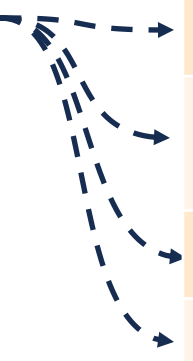
```
{
  "first_name": "John",
  "last_name": "Smith",
  "is_alive": true,
  "age": 27,
  "address": {
    "street_address": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postal_code": "10021-3100"
  },
  "phone_numbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    }
  ]
}
```

Python Dictionary:

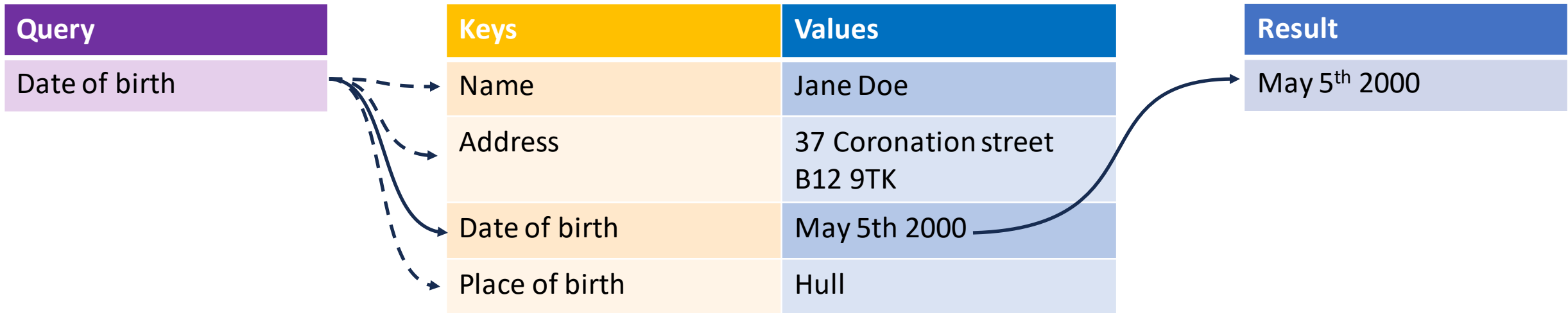
```
thisdict = {
    "brand": "Ford",
    "model": "Mustang",
    "year": 1964
}
```

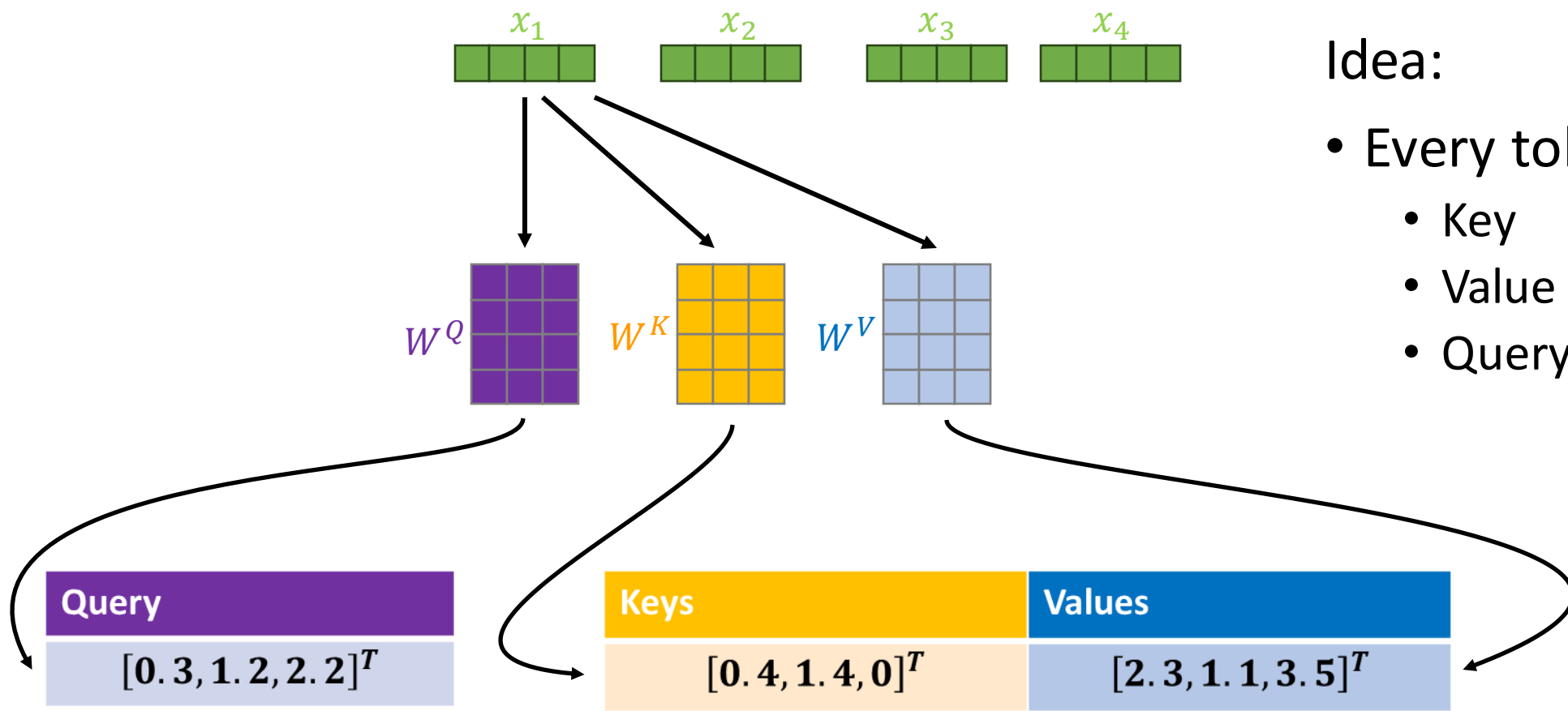
Key-Value Pairs

Query	Keys	Values
Date of birth	Name	Jane Doe
	Address	37 Coronation street B12 9TK
	Date of birth	May 5th 2000
	Place of birth	Hull



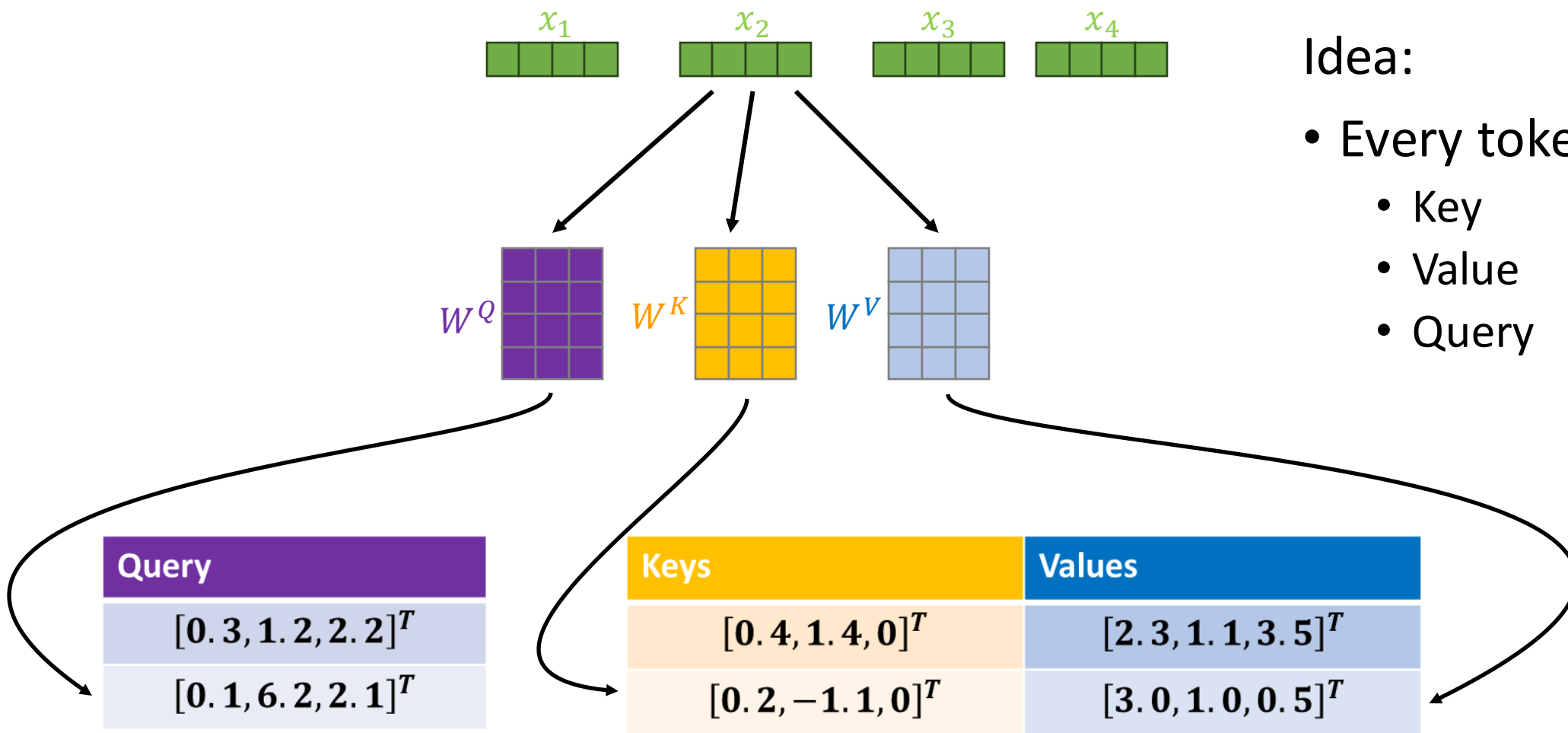
Key-Value Pairs





Idea:

- Every token makes:
 - Key
 - Value
 - Query



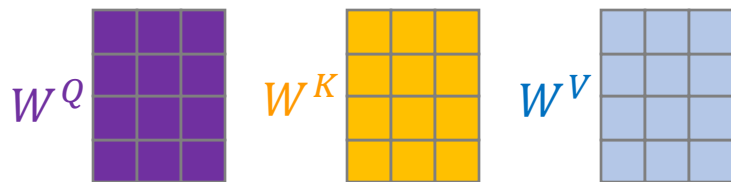
Idea:

- Every token makes:
 - Key
 - Value
 - Query



Idea:

- Every token makes:
 - Key
 - Value
 - Query



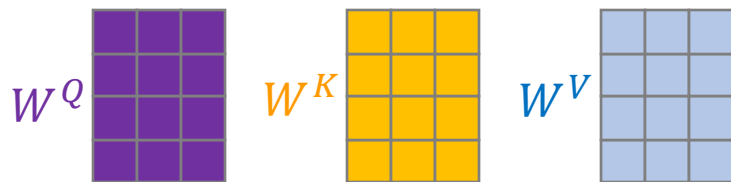
Query
$[0.3, 1.2, 2.2]^T$
$[0.1, 6.2, 2.1]^T$
$[2.1, 2.5, 5.2]^T$

Keys	Values
$[0.4, 1.4, 0]^T$	$[2.3, 1.1, 3.5]^T$
$[0.2, -1.1, 0]^T$	$[3.0, 1.0, 0.5]^T$
$[0.3, 1.2, 4]^T$	$[-2.2, 4.1, 2.5]^T$



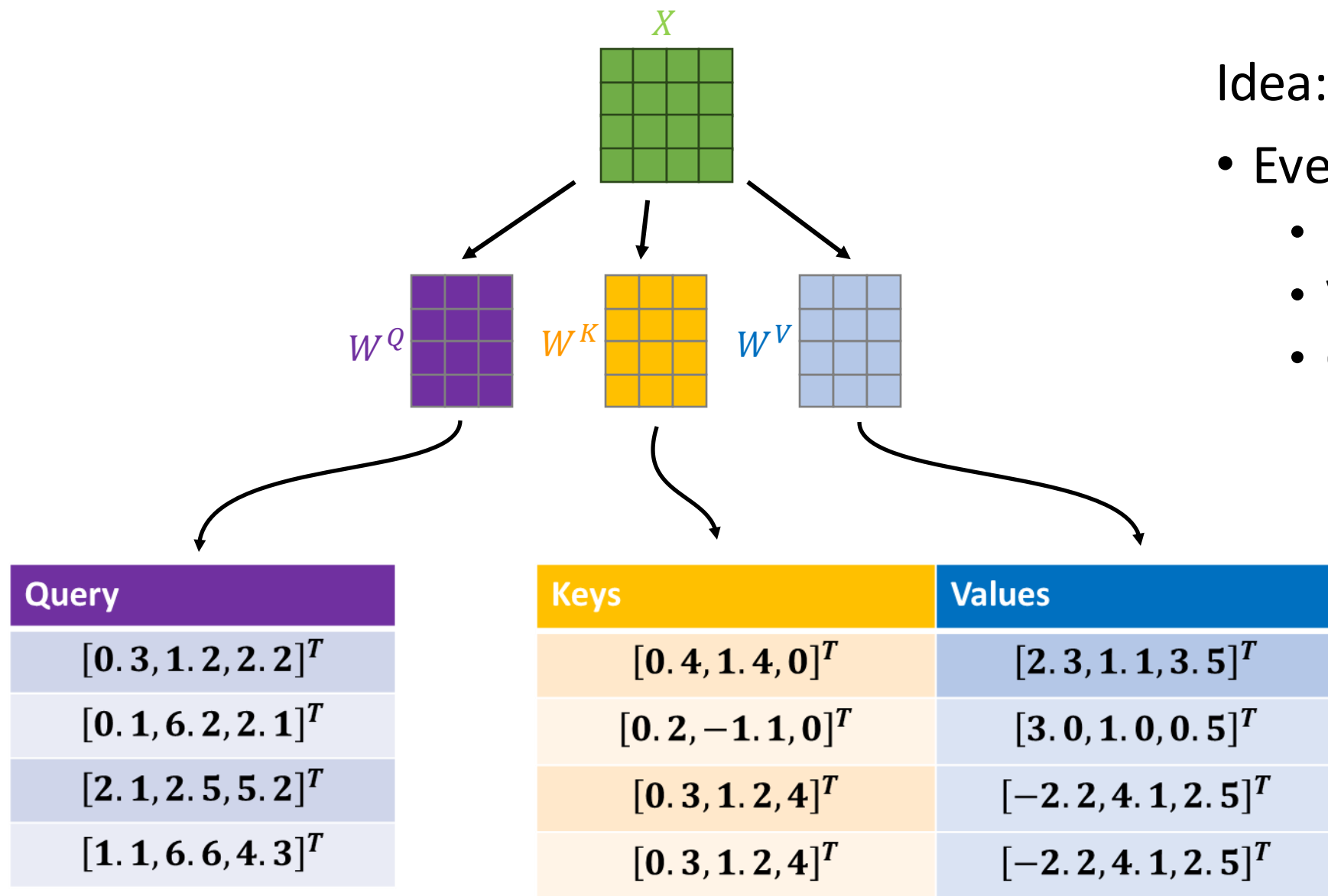
Idea:

- Every token makes:
 - Key
 - Value
 - Query



Query
$[0.3, 1.2, 2.2]^T$
$[0.1, 6.2, 2.1]^T$
$[2.1, 2.5, 5.2]^T$
$[1.1, 6.6, 4.3]^T$

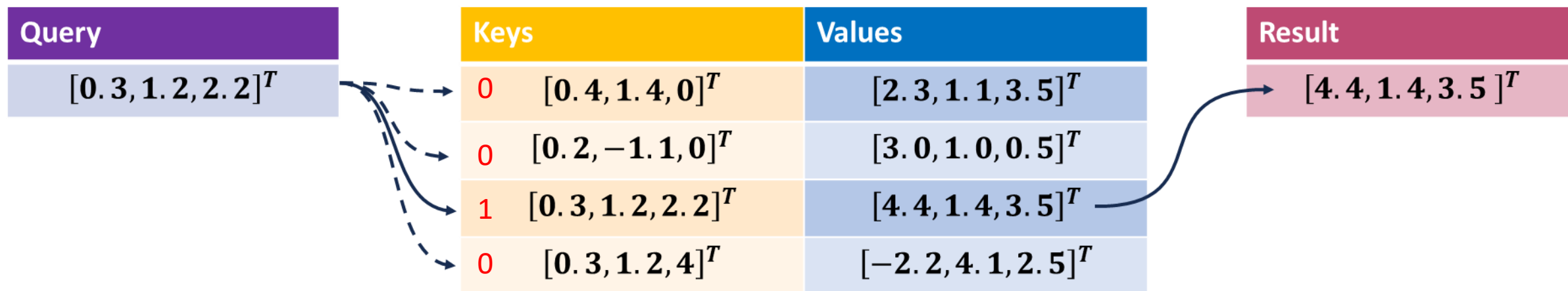
Keys	Values
$[0.4, 1.4, 0]^T$	$[2.3, 1.1, 3.5]^T$
$[0.2, -1.1, 0]^T$	$[3.0, 1.0, 0.5]^T$
$[0.3, 1.2, 4]^T$	$[-2.2, 4.1, 2.5]^T$
$[0.3, 1.2, 4]^T$	$[-2.2, 4.1, 2.5]^T$



Idea:

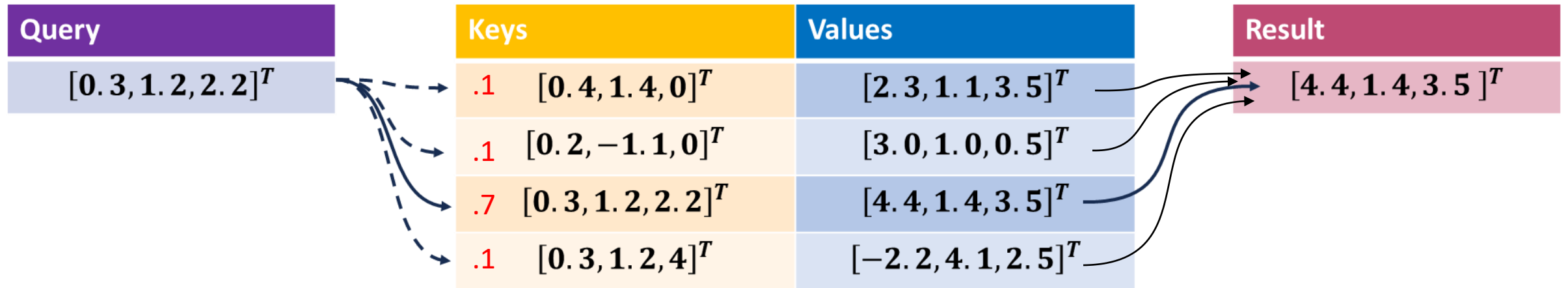
- Every token makes:
 - Key
 - Value
 - Query

Query



$$z = \sum_{j=1}^n \mathbf{1}(q = k_j) v_j$$

Relaxed Query



$$\mathbf{z} = \sum_{j=1}^n \text{score}_j \mathbf{v}_j$$

$$\text{score}_j = \text{softmax}\left(\underbrace{\text{similarity}(q, \mathbf{k}_j)}_{\frac{\mathbf{q}^T \mathbf{k}_j}{\sqrt{d_k}}}\right)$$

Self-Attention

Application to sequence: matrix multiplication

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

Q, K, V are computed from **X** (embedding of input sequence) with *learned* weight matrix



Self-Attention

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

Self-Attention

$$\text{Attention}(Q, K, V) = \text{softmax}(\text{ } \begin{array}{|c|c|} \hline \text{red} & \text{red} \\ \hline \text{red} & \text{red} \\ \hline \end{array} \text{ }) \begin{array}{|c|c|c|} \hline \text{blue} & \text{blue} & \text{blue} \\ \hline \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array}$$

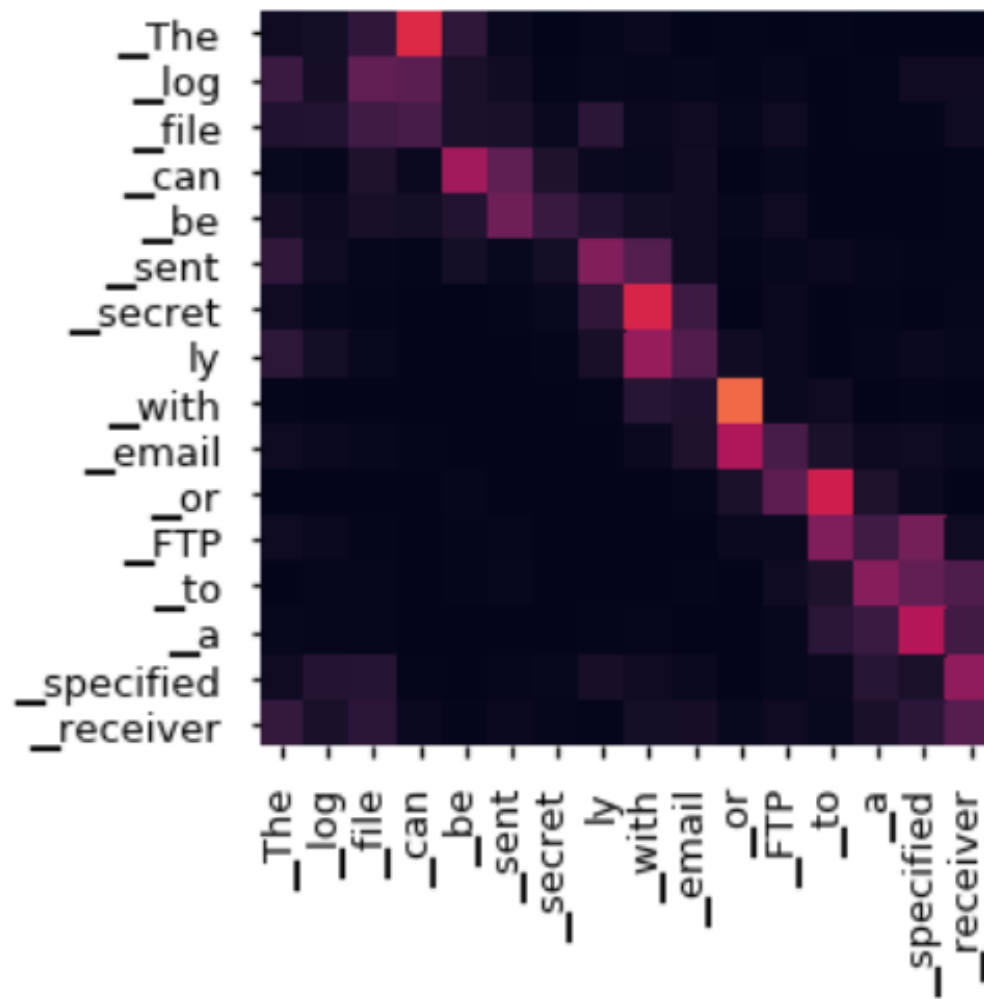
The diagram illustrates the self-attention mechanism. It shows the formula $\text{Attention}(Q, K, V) = \text{softmax}(\text{ } \begin{array}{|c|c|} \hline \text{red} & \text{red} \\ \hline \text{red} & \text{red} \\ \hline \end{array} \text{ }) \begin{array}{|c|c|c|} \hline \text{blue} & \text{blue} & \text{blue} \\ \hline \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array}$. The red grid represents the query (Q) and key (K) matrices, and the blue grid represents the value (V) matrix. The blue 'V' label is positioned above the blue grid.

Self-Attention

$$\text{Attention}(Q, K, V) =$$



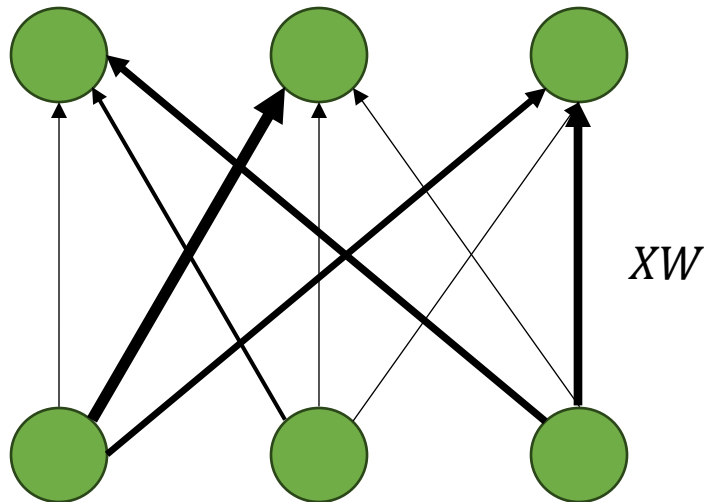
Attention Matrix



Self-Attention vs. Feed-Forward

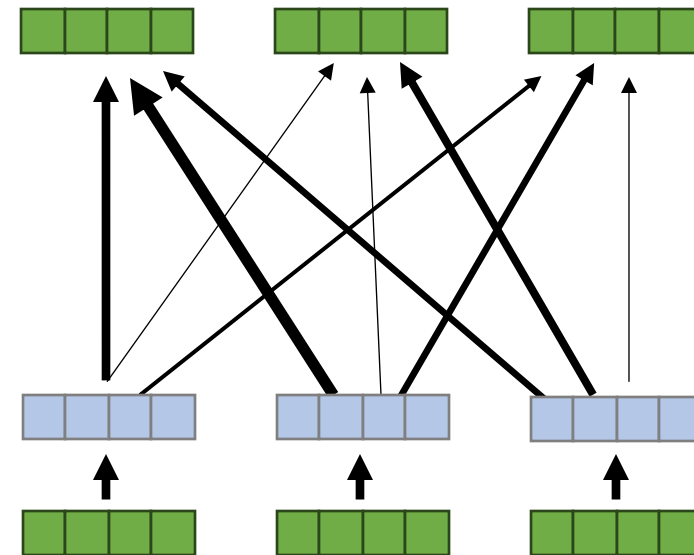
Perceptron / Feed-Forward:

- **Fixed weight matrix**



Self-Attention:

- **Dynamic weight matrix**
- Computed from inputs



$$\text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

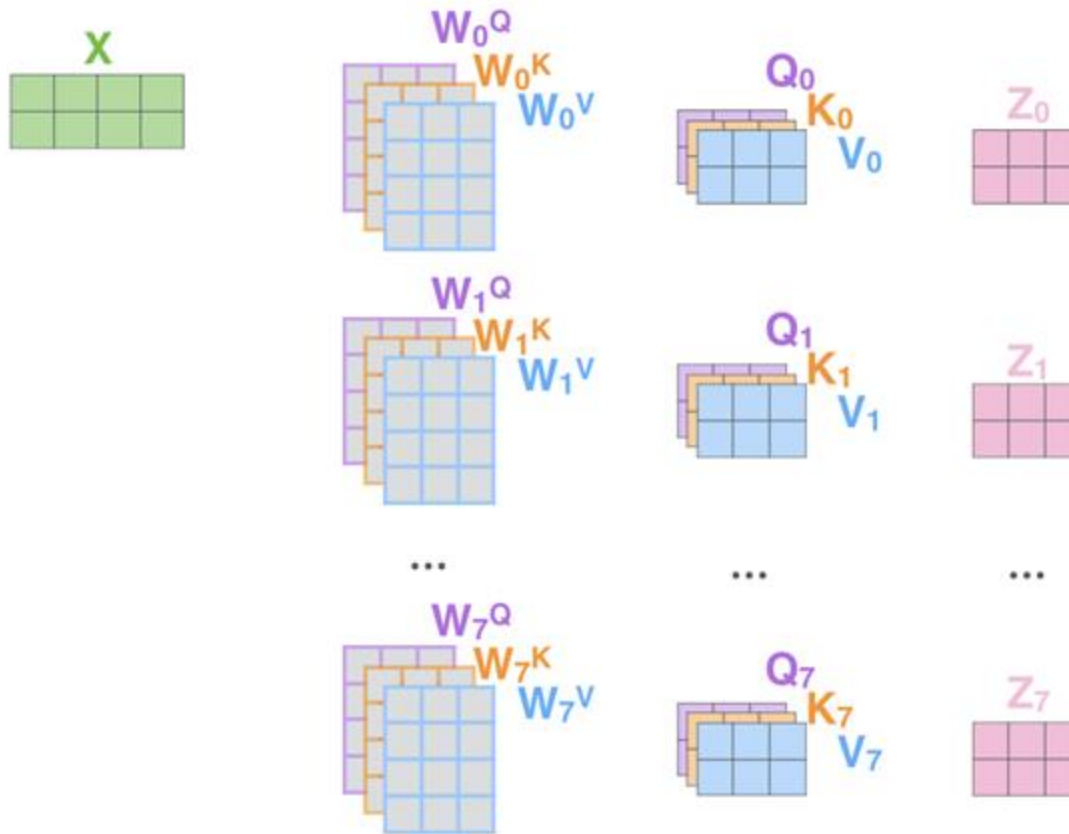
The equation shows the computation of the self-attention output. It includes a purple matrix Q , an orange matrix K^T , and a blue matrix V .

Multi-headed self-attention



Multi-headed self-attention

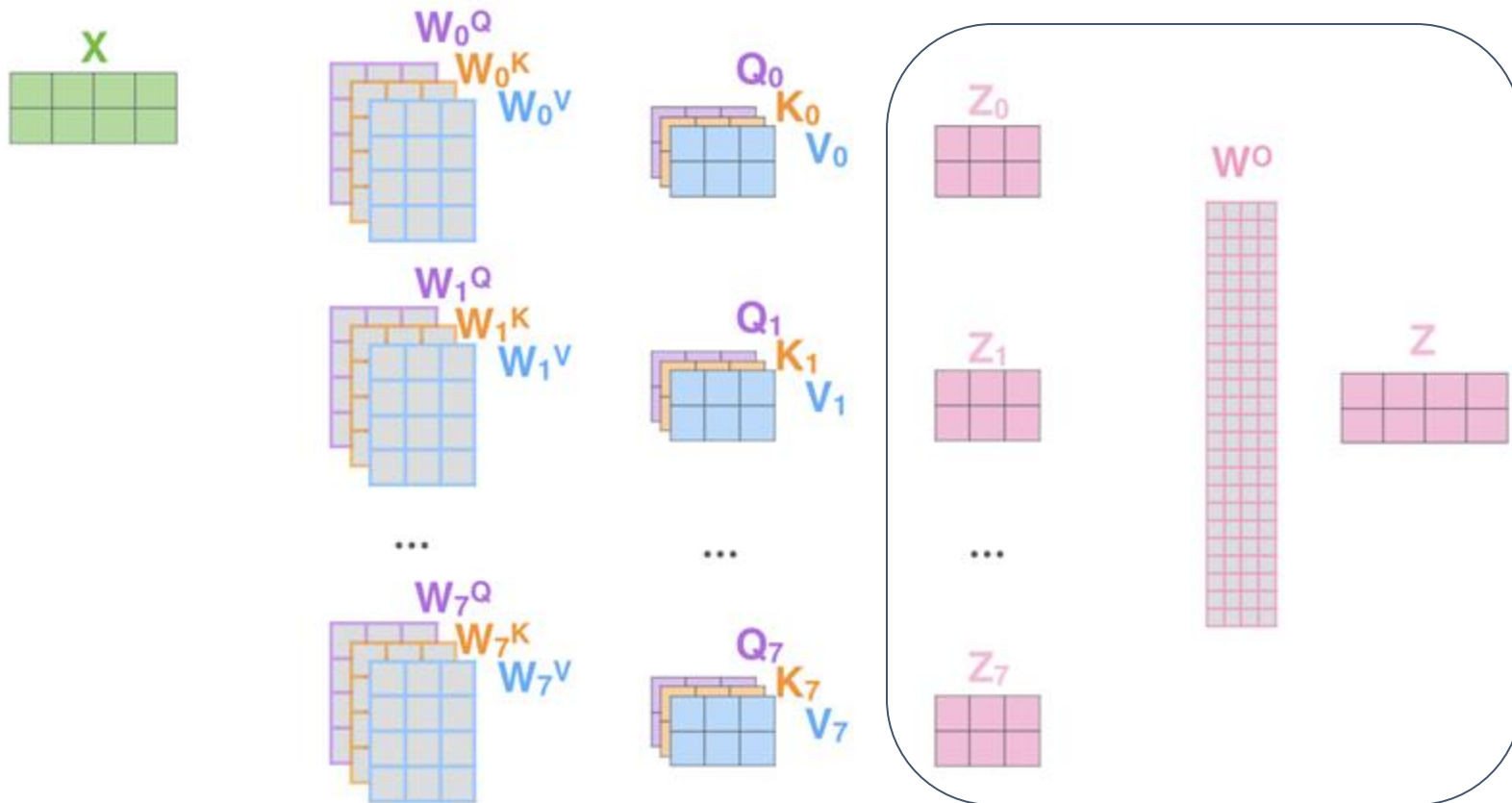
Multiple attention heads for
increased model capacity



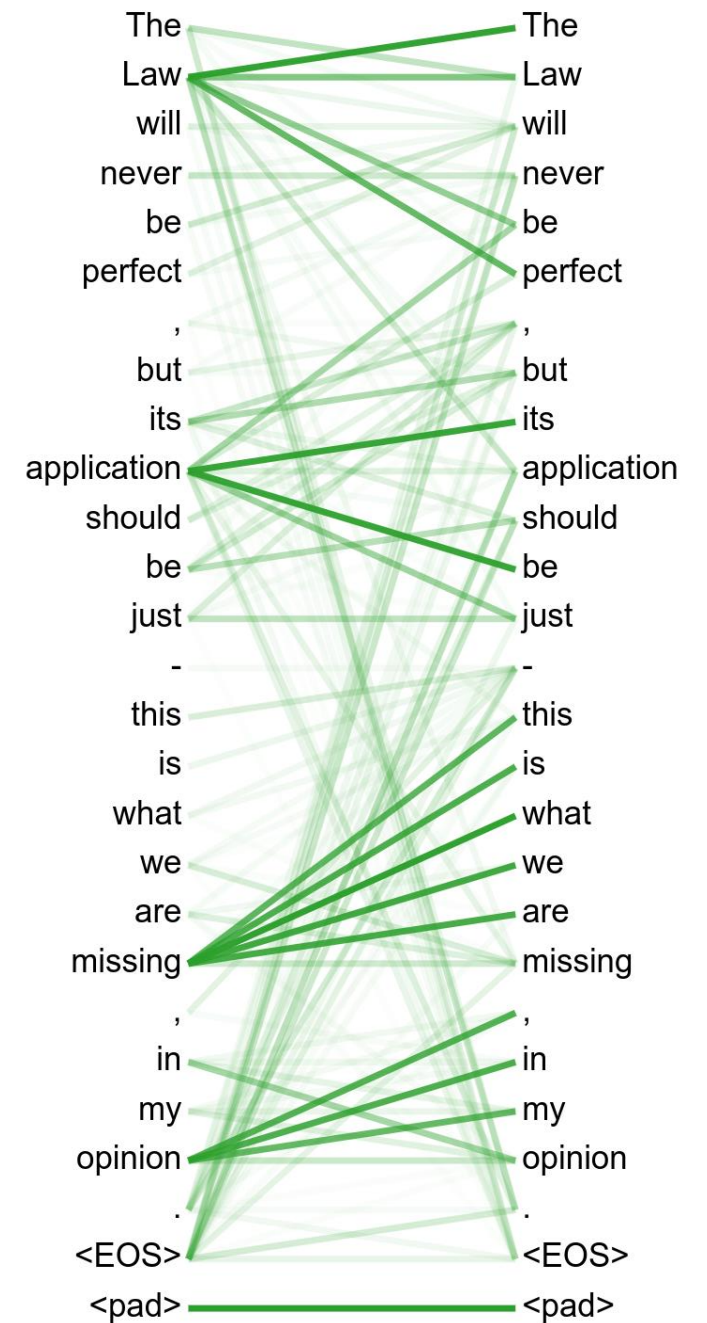
Multi-headed self-attention

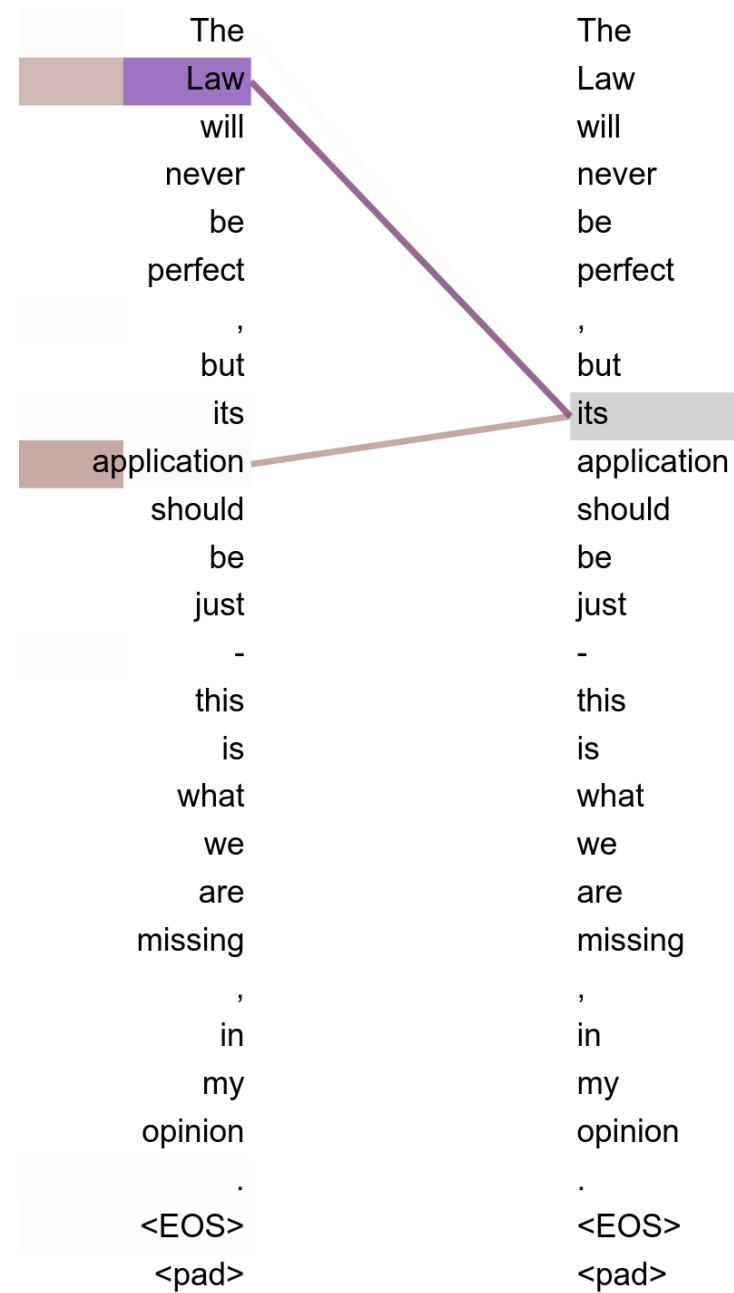
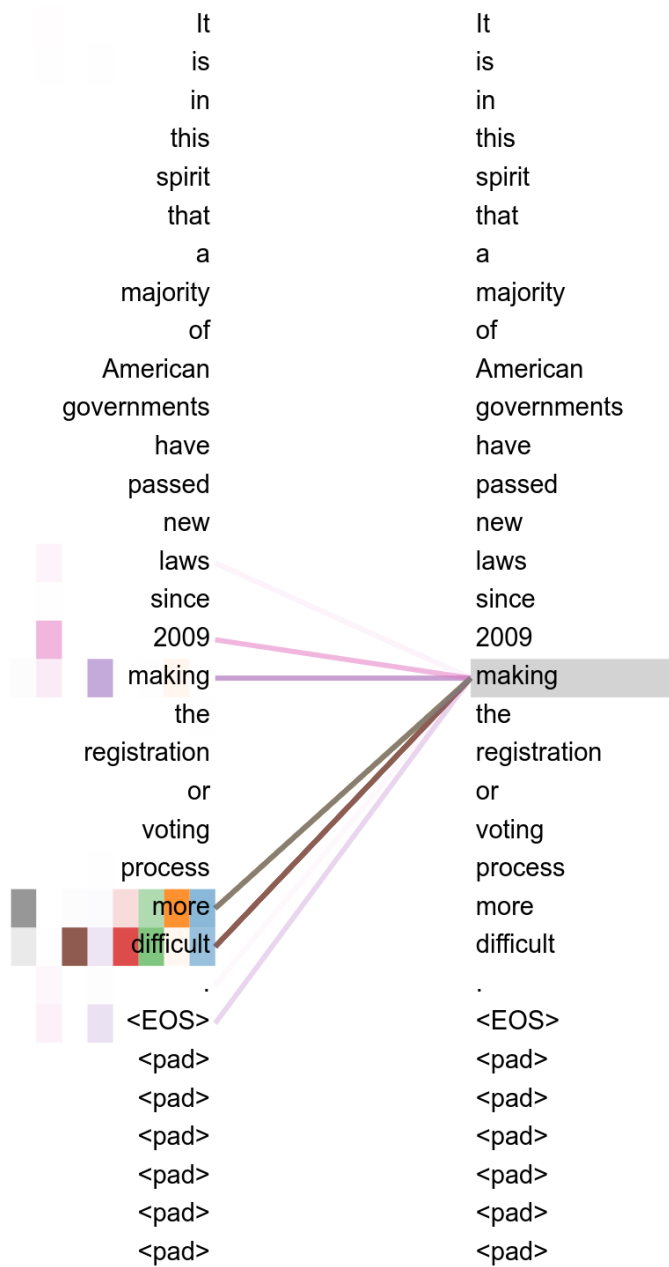
Multiple attention heads for
increased model capacity

Combine with additional feed-
forward layer



Concat z_i outputs, project to z
with learned weight matrix





Next Video: the Decoder

