

Deep Learning

Transformers and Graph Neural Networks

Abuzar Khan, Yue Jian

11-785, Fall 2022

Part 1

Transformers

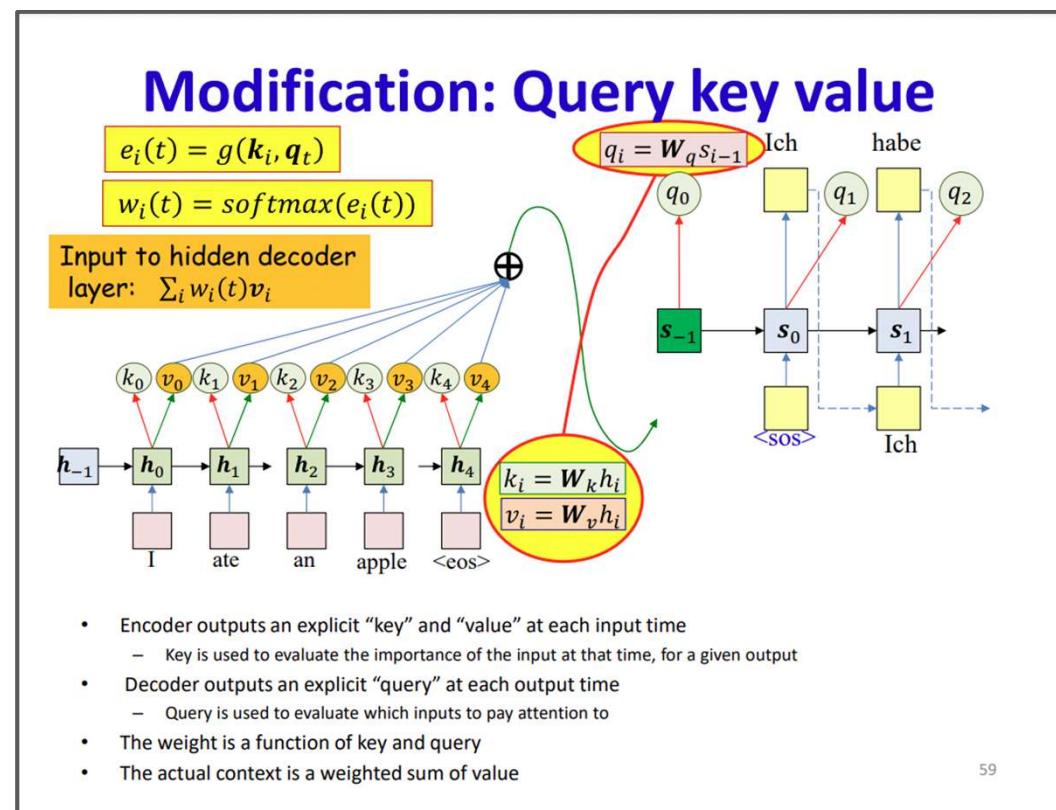
Credit where due

Probably influenced by:

- Louis-Philippe Morency, **11777**
- Graham Neubig, **11711**
- (obviously) Bhiksha Raj, **11-you-know-your-course-number**
- Bunch of youtube videos
- Medium Articles

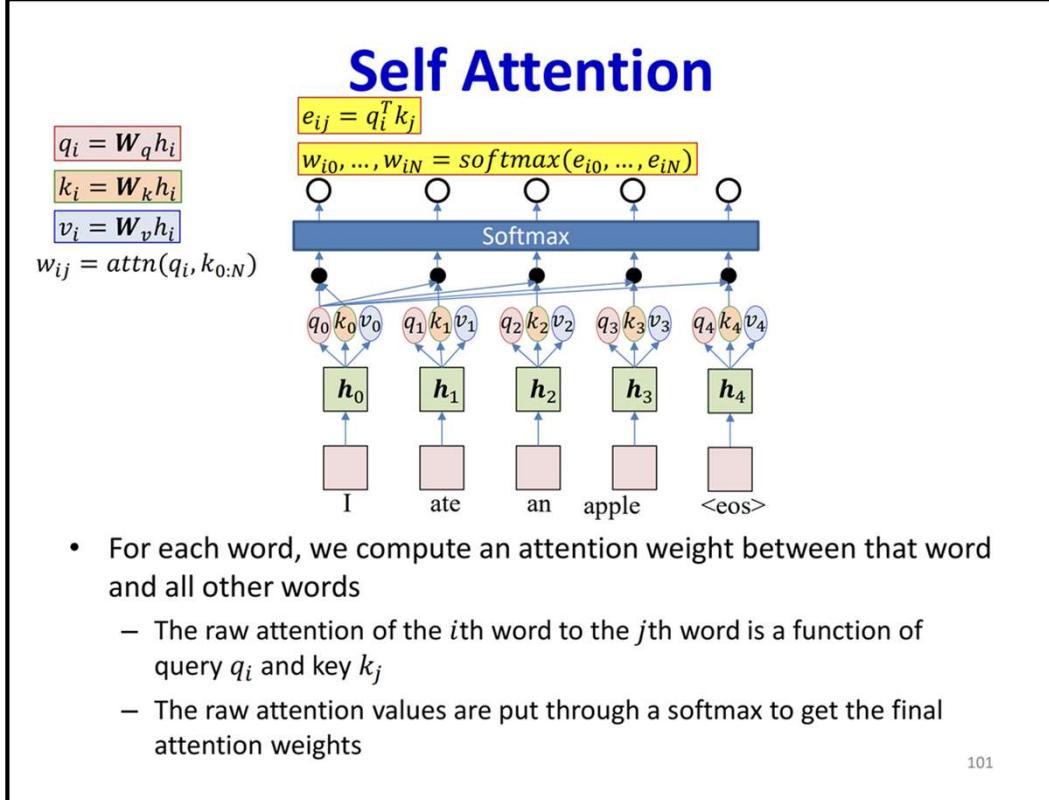
Recall

1. Queries, Keys, and Values



Recall

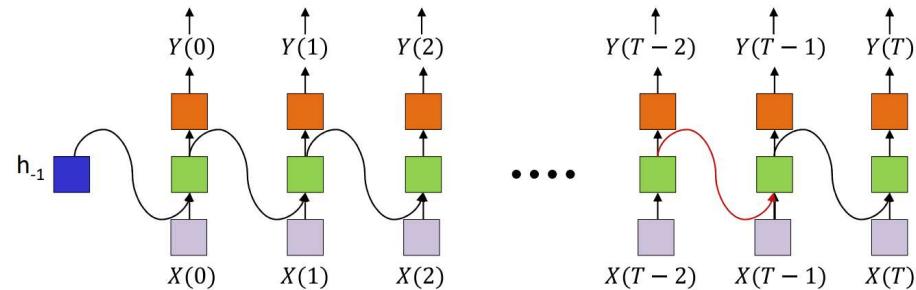
1. Queries, Keys, and Values
2. Self Attention
 - a. Energy Function
 - b. Attention Function



Recall

1. Queries, Keys, and Values
2. Self Attention
 - a. Energy Function
 - b. Attention Function
3. RNNs are slow and sequential
 - a. Attention-based models can be **parallelized!**

Processing order



- Computing $Y(T)$ requires $Y(T - 1)\dots$
- Which requires $Y(T - 2)$, etc...
- RNN inputs must be processed *in order* → slow implementation

Why Transformers

- We want representations that are “dynamic” to context
“I **like** this movie” vs. “I do not **like** this movie”
like should have different representations in both cases
- Vanilla RNNs are **Slow** and have **terrible memory**
- LSTMs and GRUs fix the **memory** problem, but are still **slow** and **sequential**
- CNNs can be **parallelized** but the kernels are static.
- We want **parallelizability**, good **memory**, and **dynamic** computation

Q,K,V in Attention

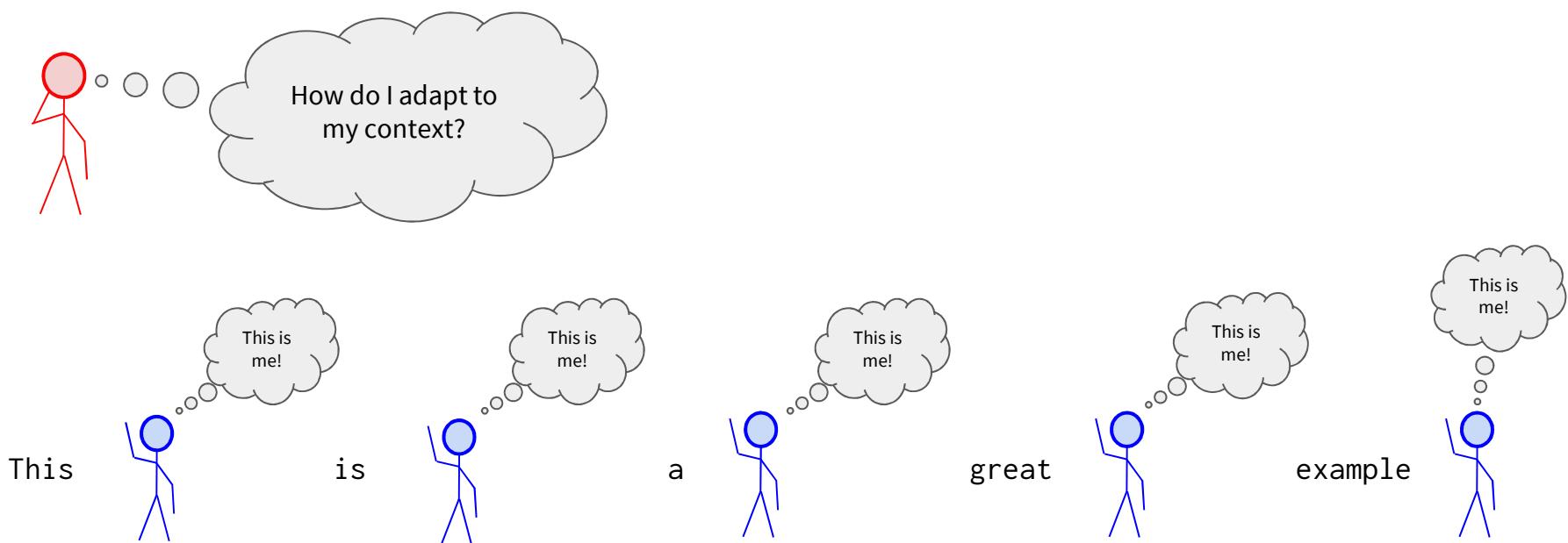
Query: This is what **pays the attention**

Values: These are **paid attention to**

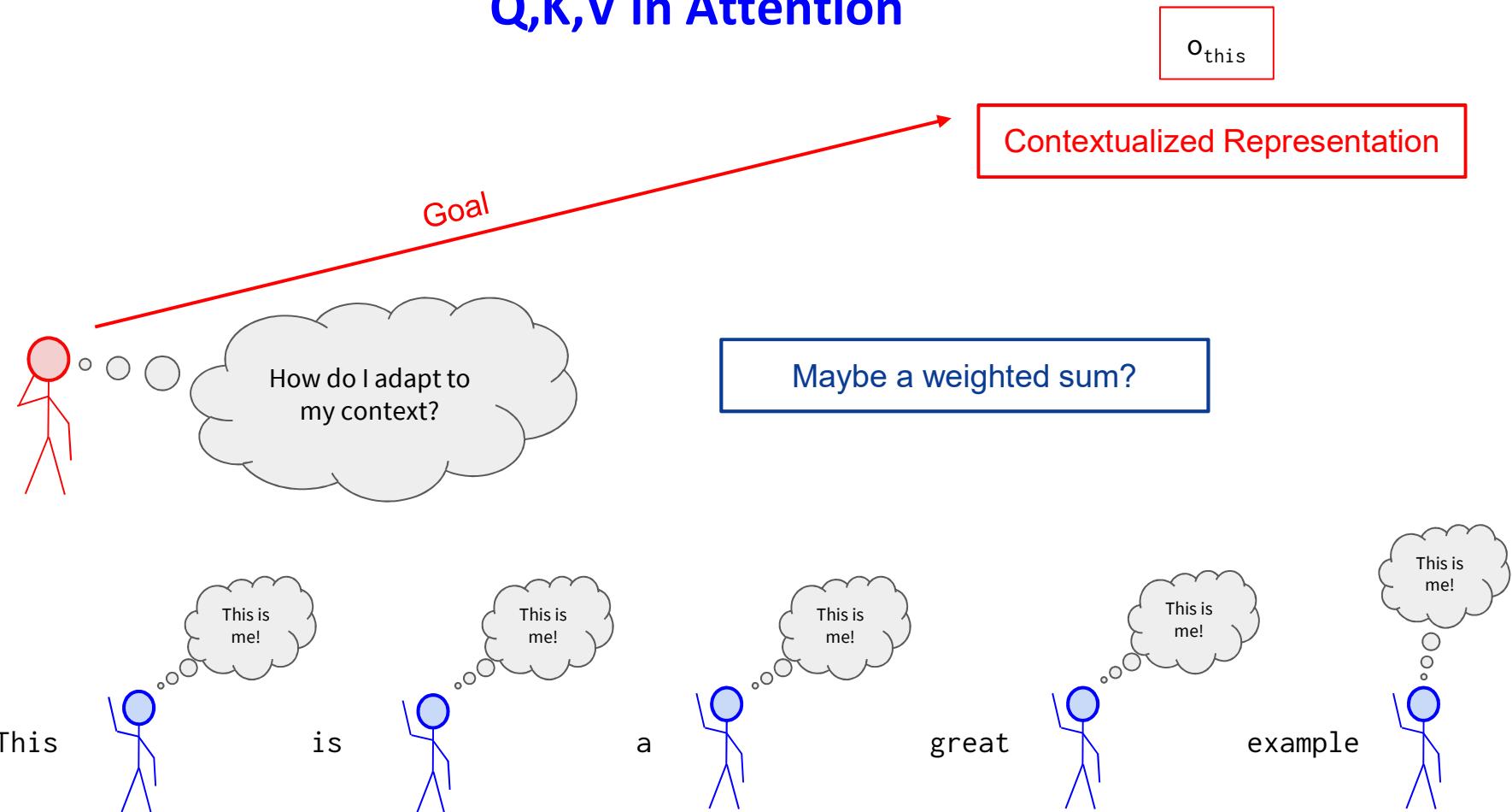
Keys: These help queries figure out **how much attention to pay** to each of the values

Attention Weights: How much attention to pay.

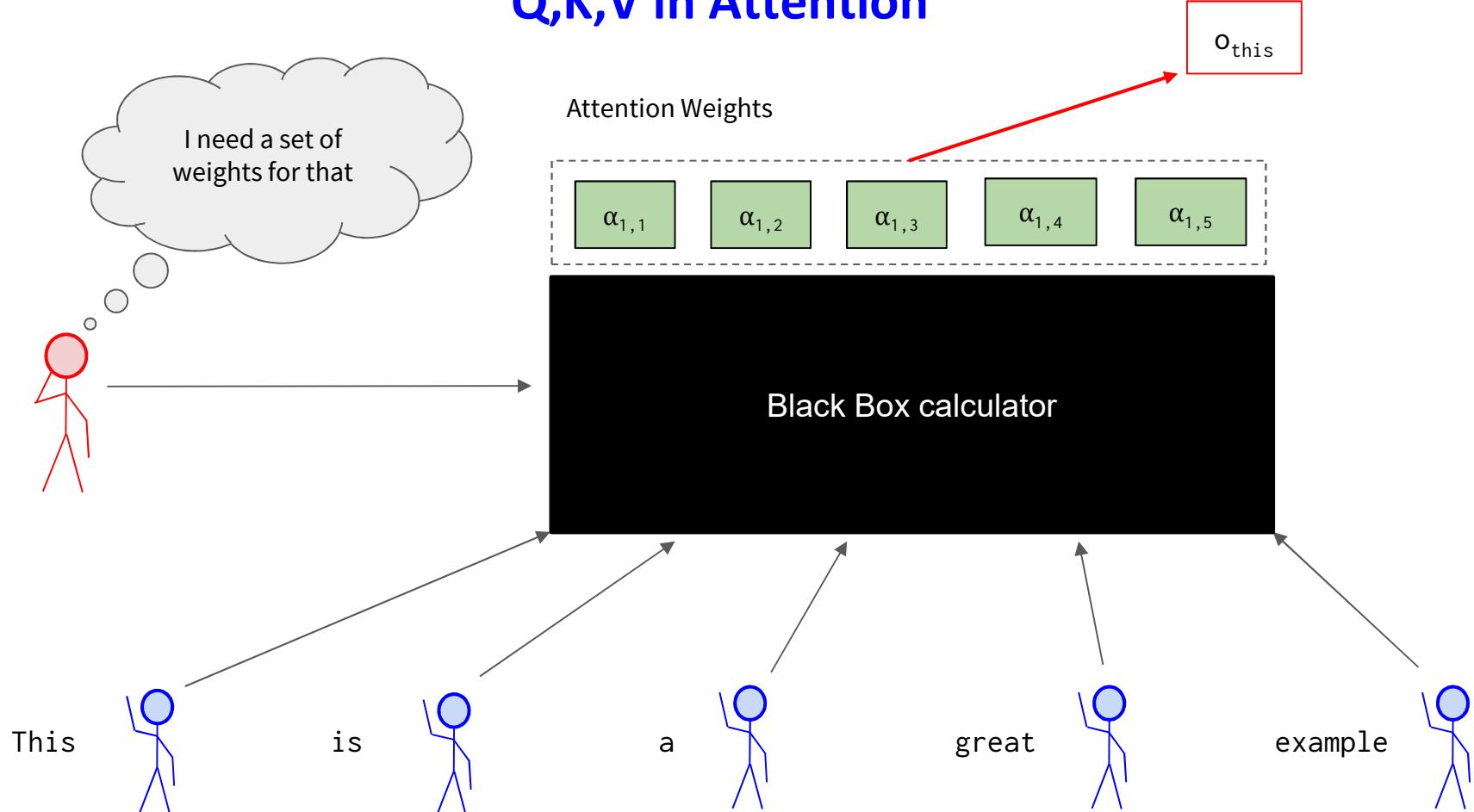
Q,K,V in Attention



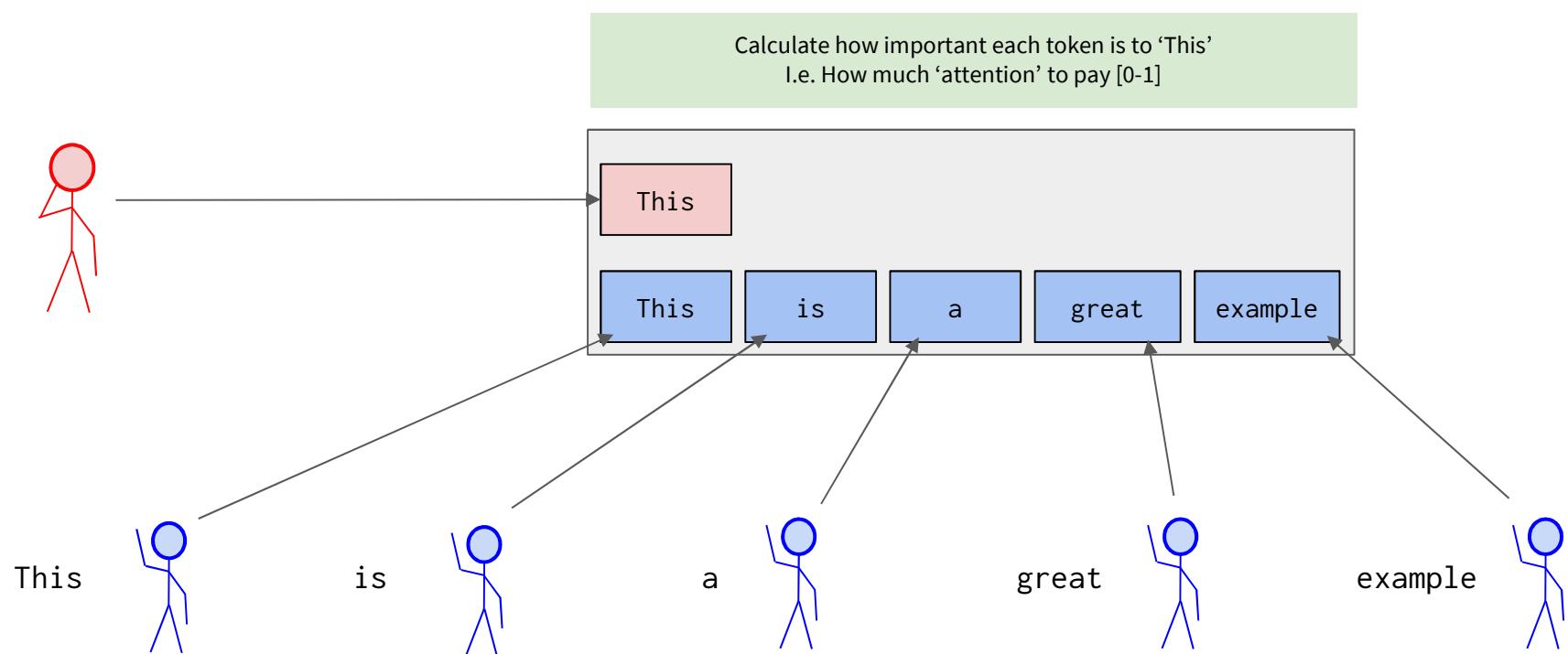
Q,K,V in Attention



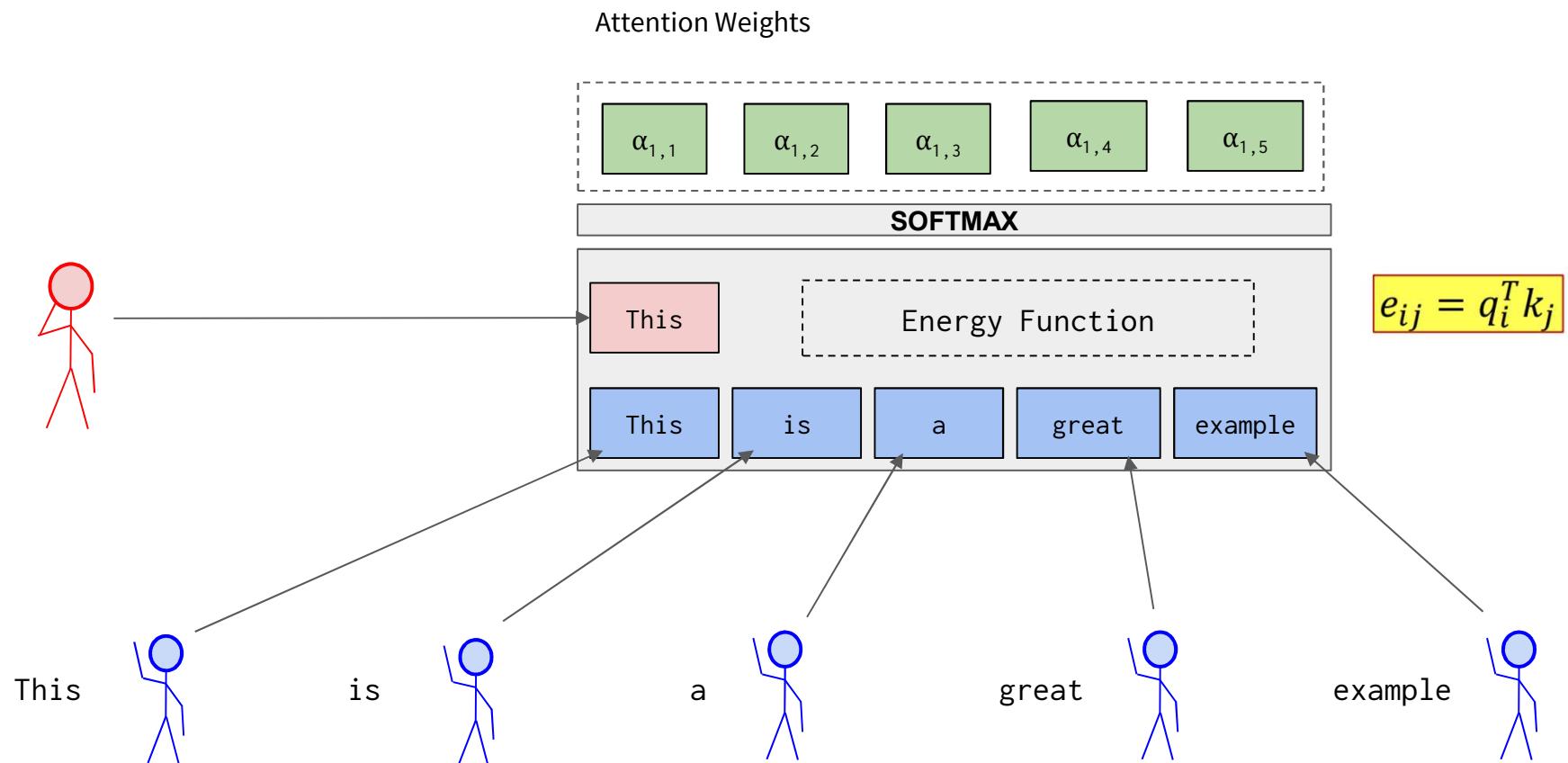
Q,K,V in Attention



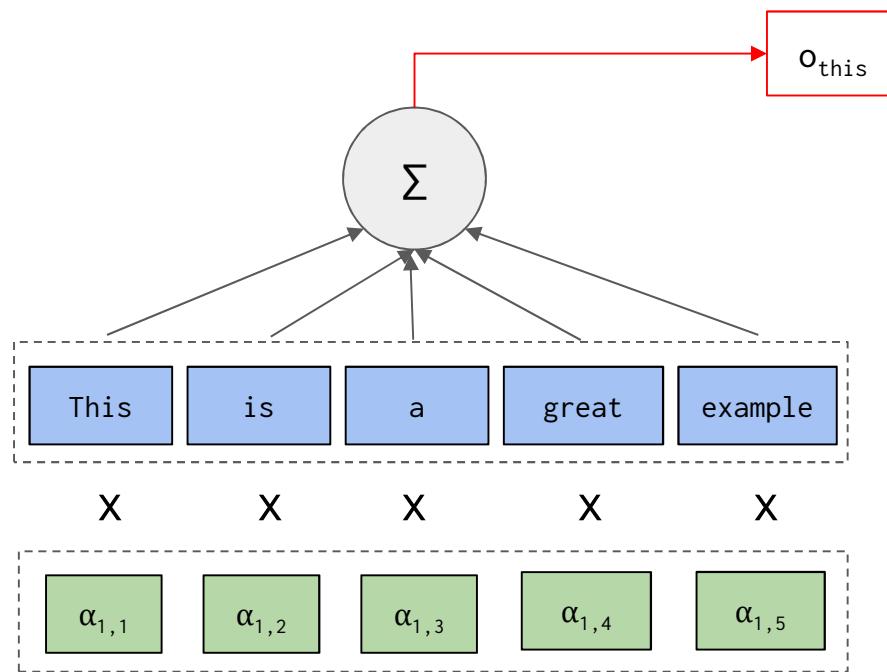
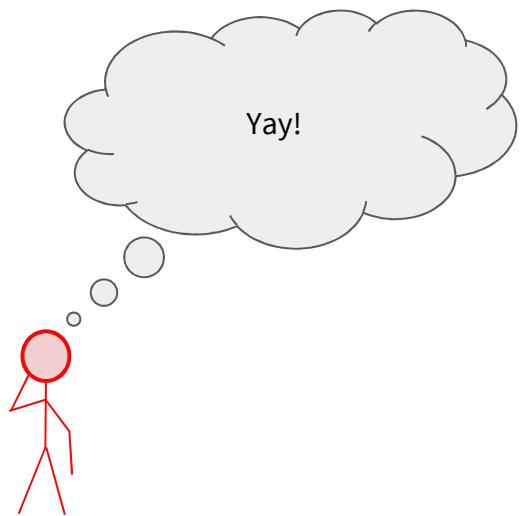
Q,K,V in Attention



Q,K,V in Attention



Q,K,V in Attention



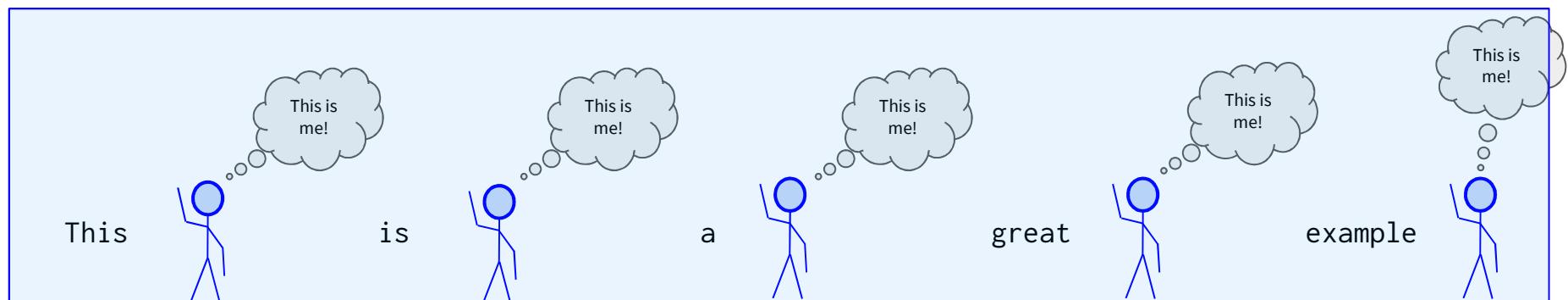
This  is  a  great  example 

Q,K,V in Attention

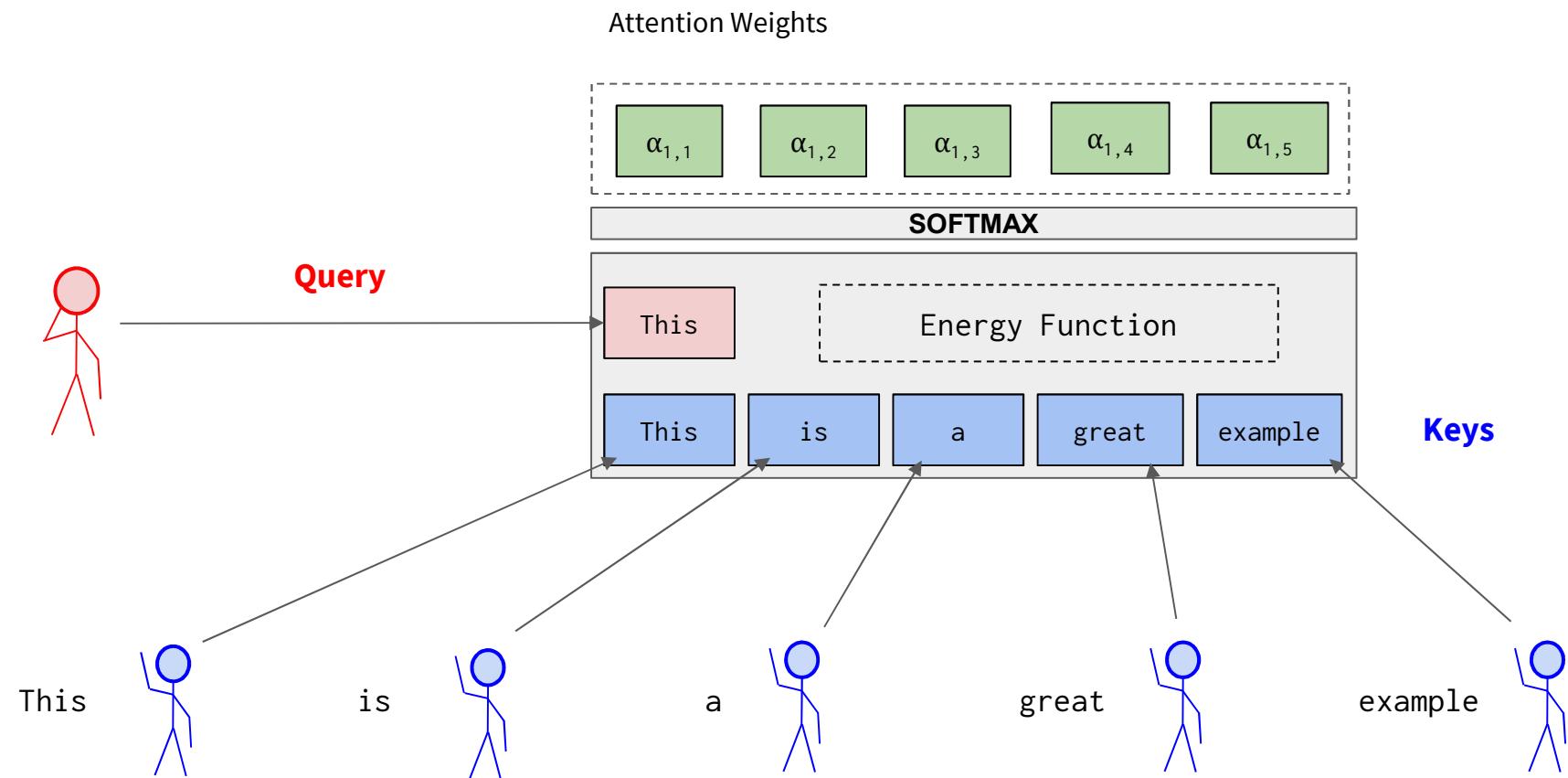
Query



Keys



Q,K,V in Attention



Q,K,V in Attention



Attention Weights

This



is



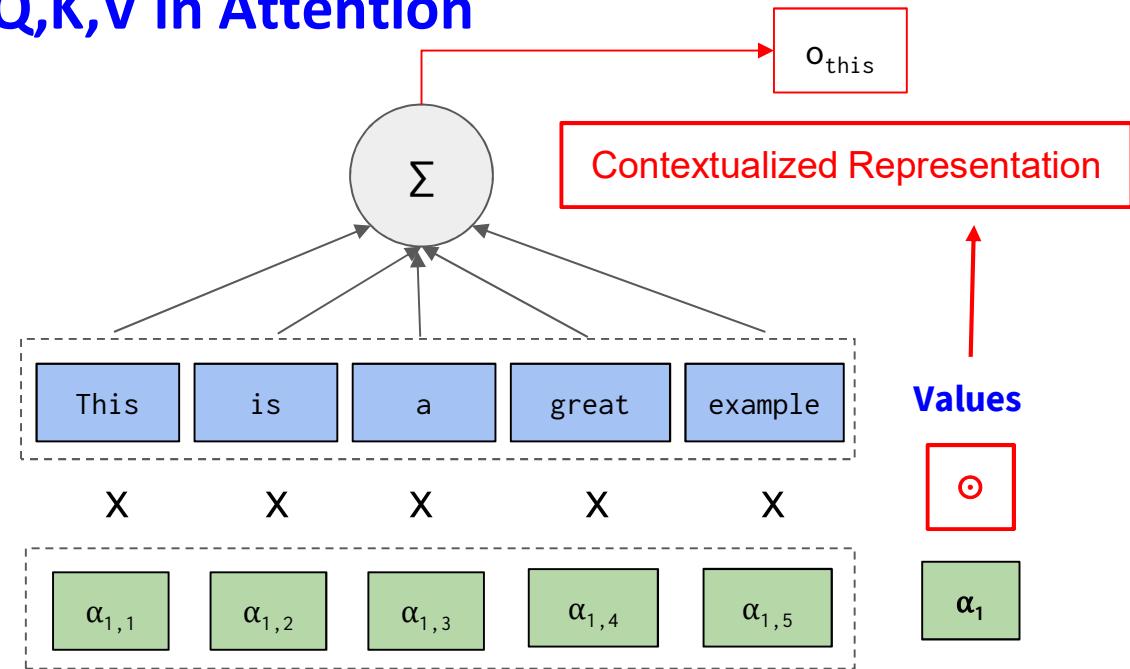
a



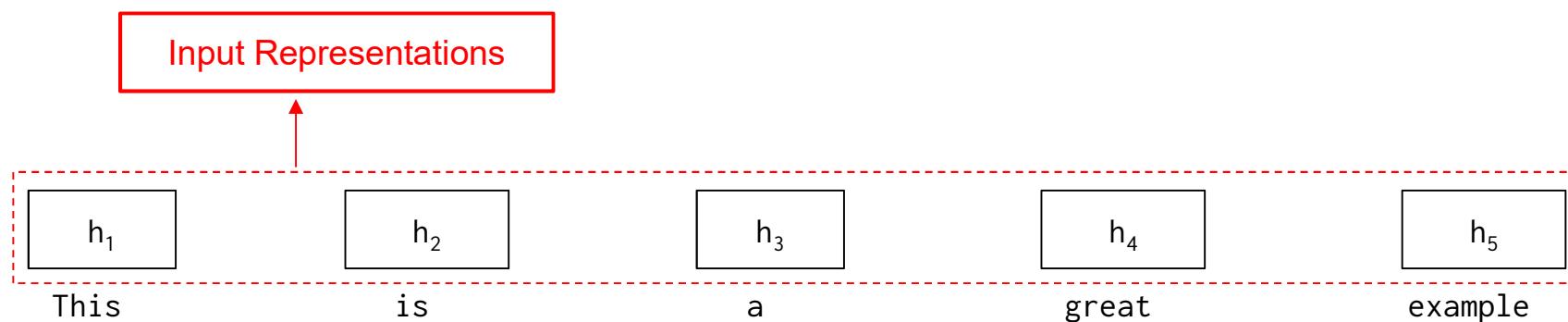
great



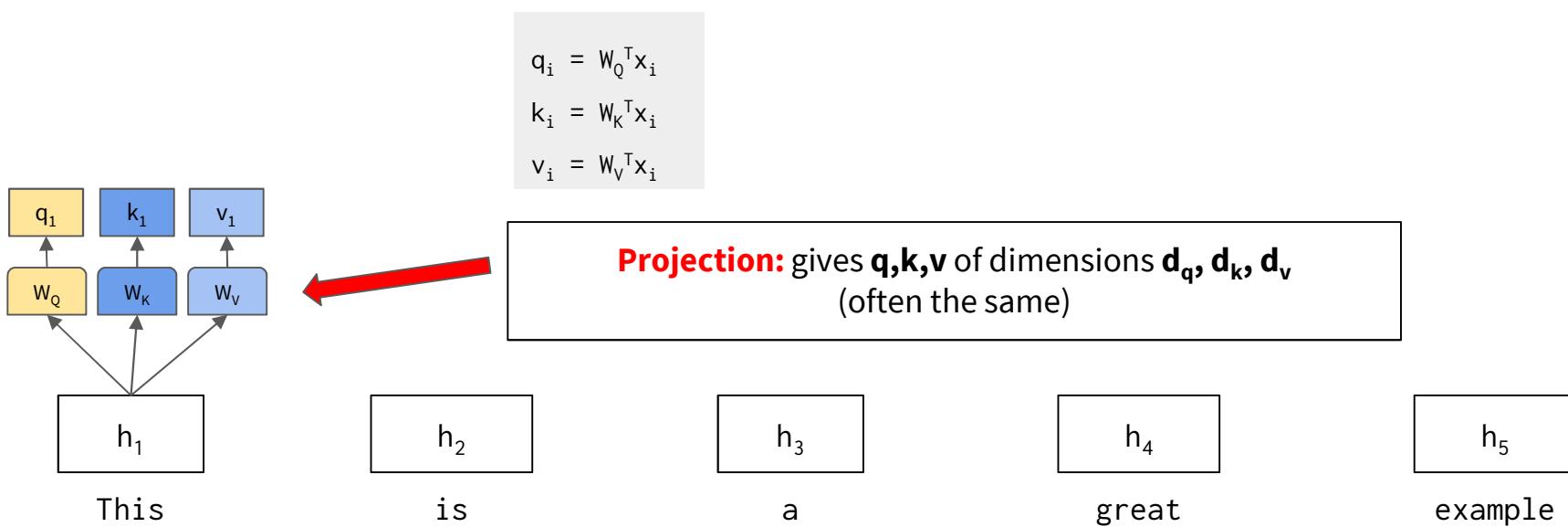
example



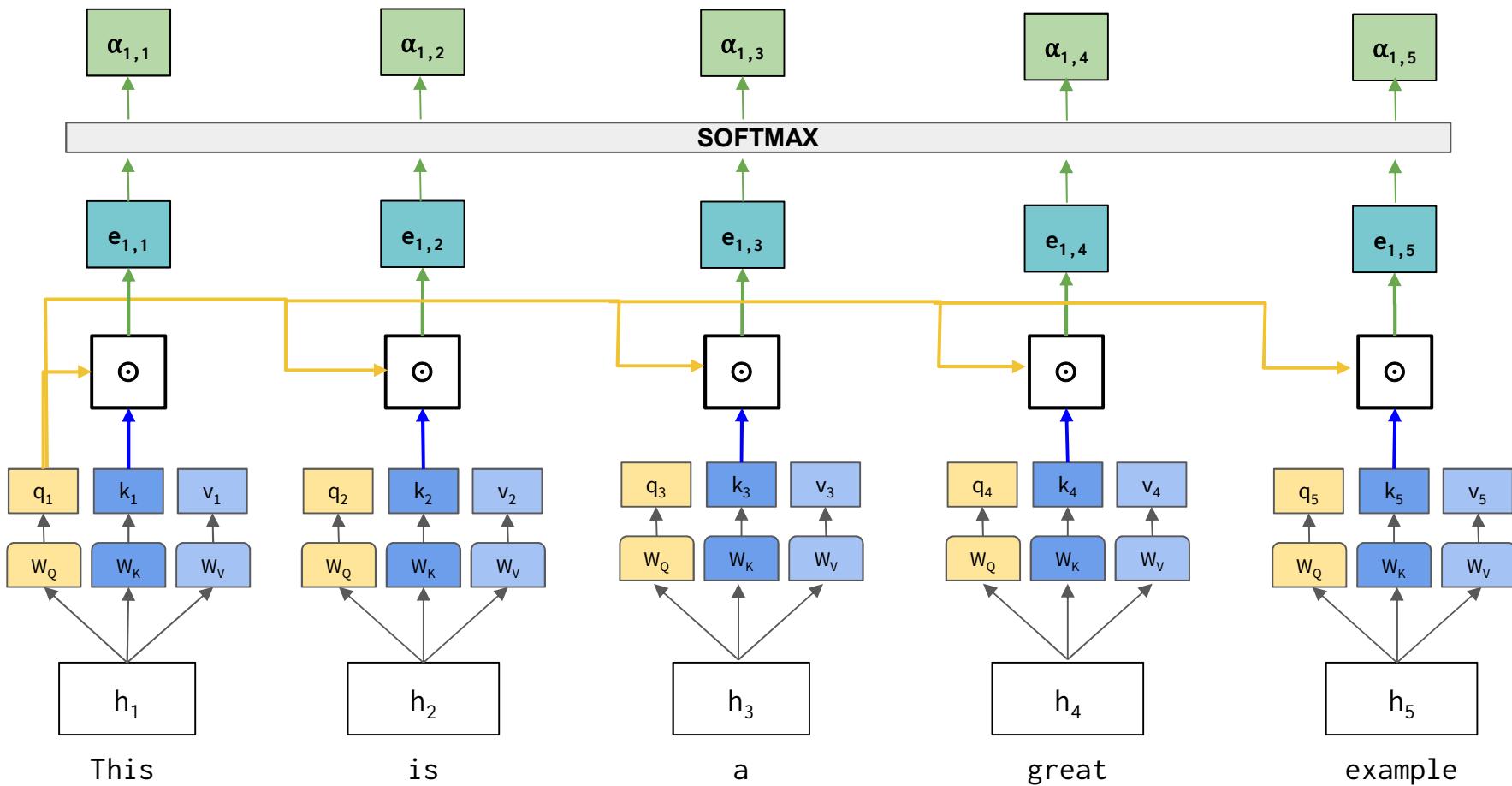
Self Attention

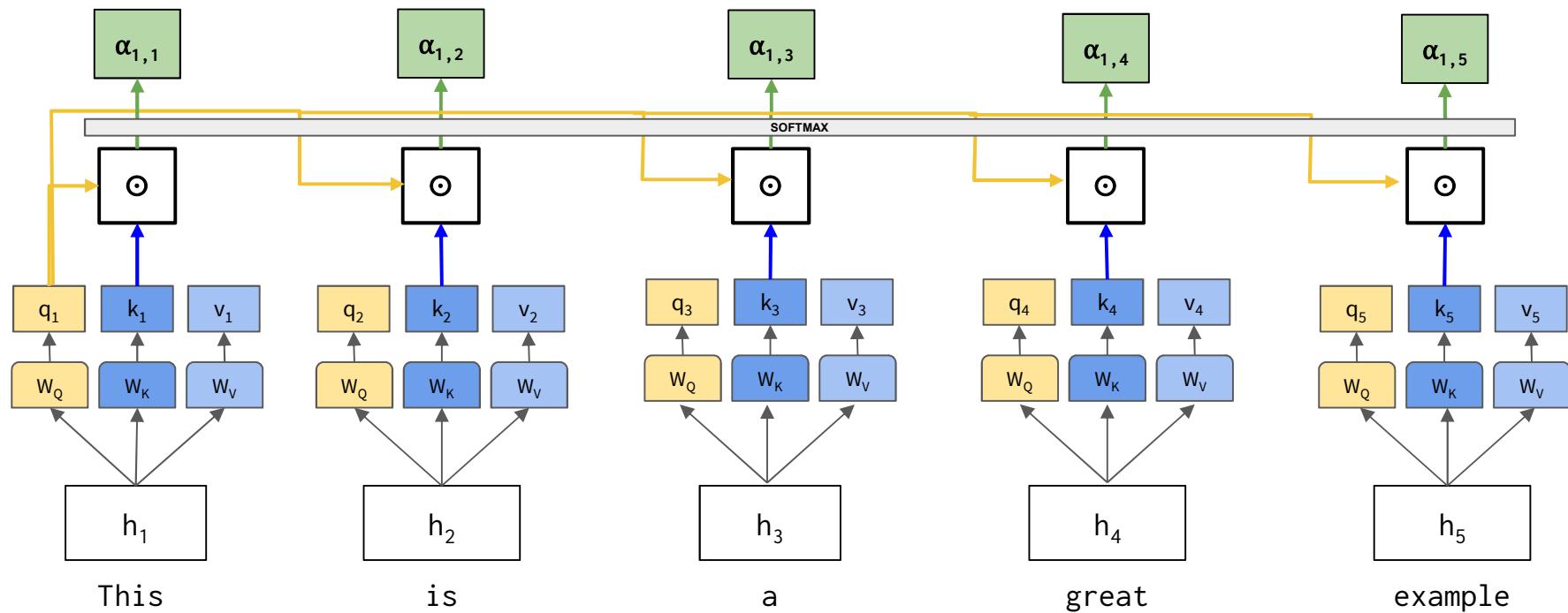


Self Attention

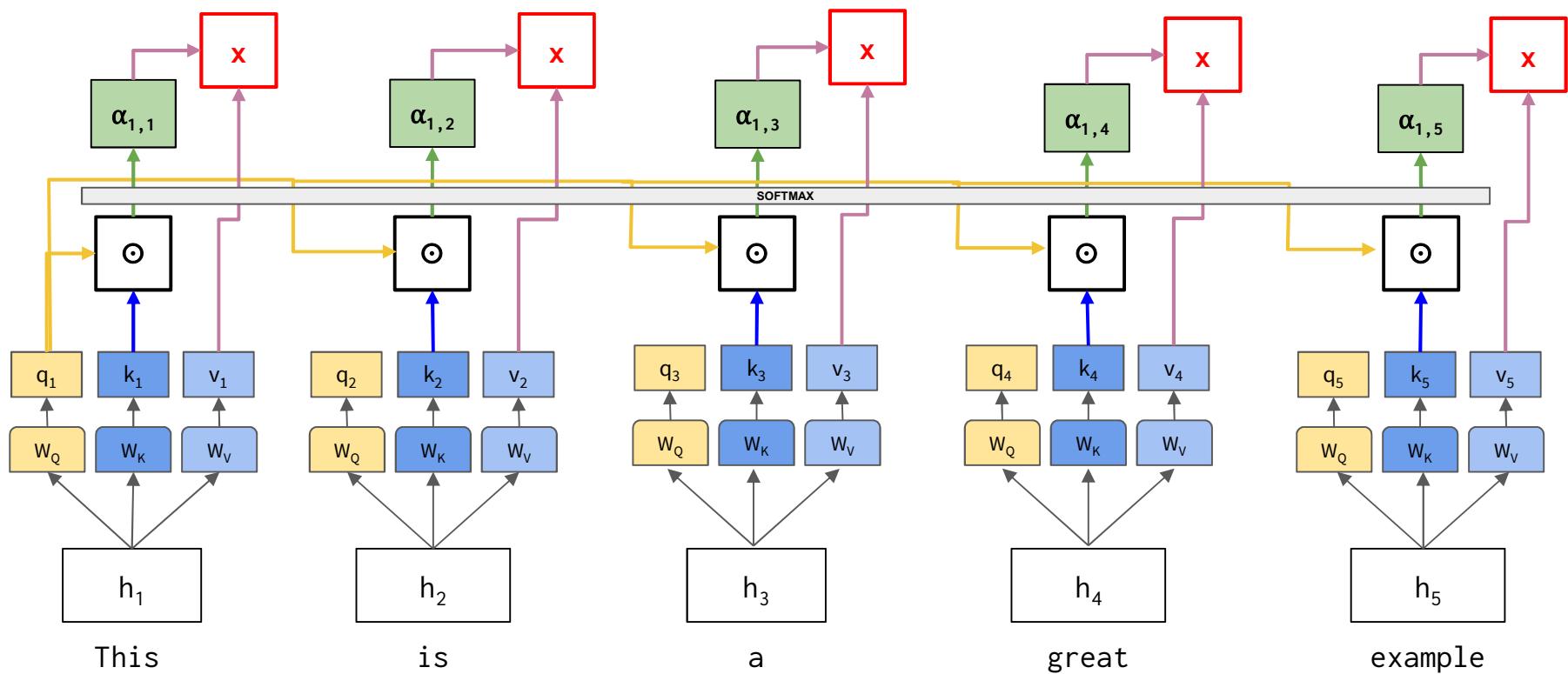


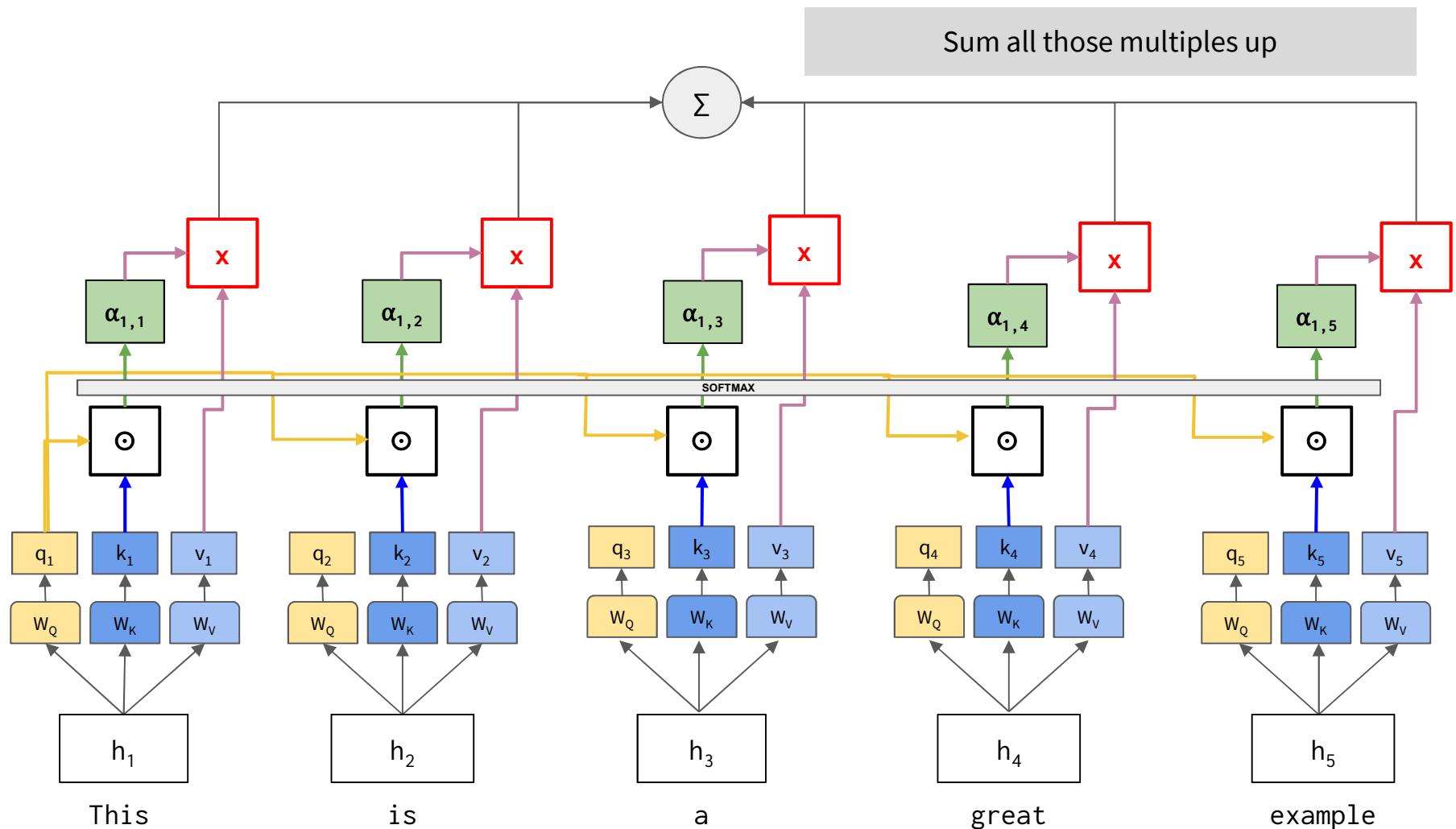
$\alpha_{m,n}$ = How important is token **n** to token **m**'s contextual meaning?

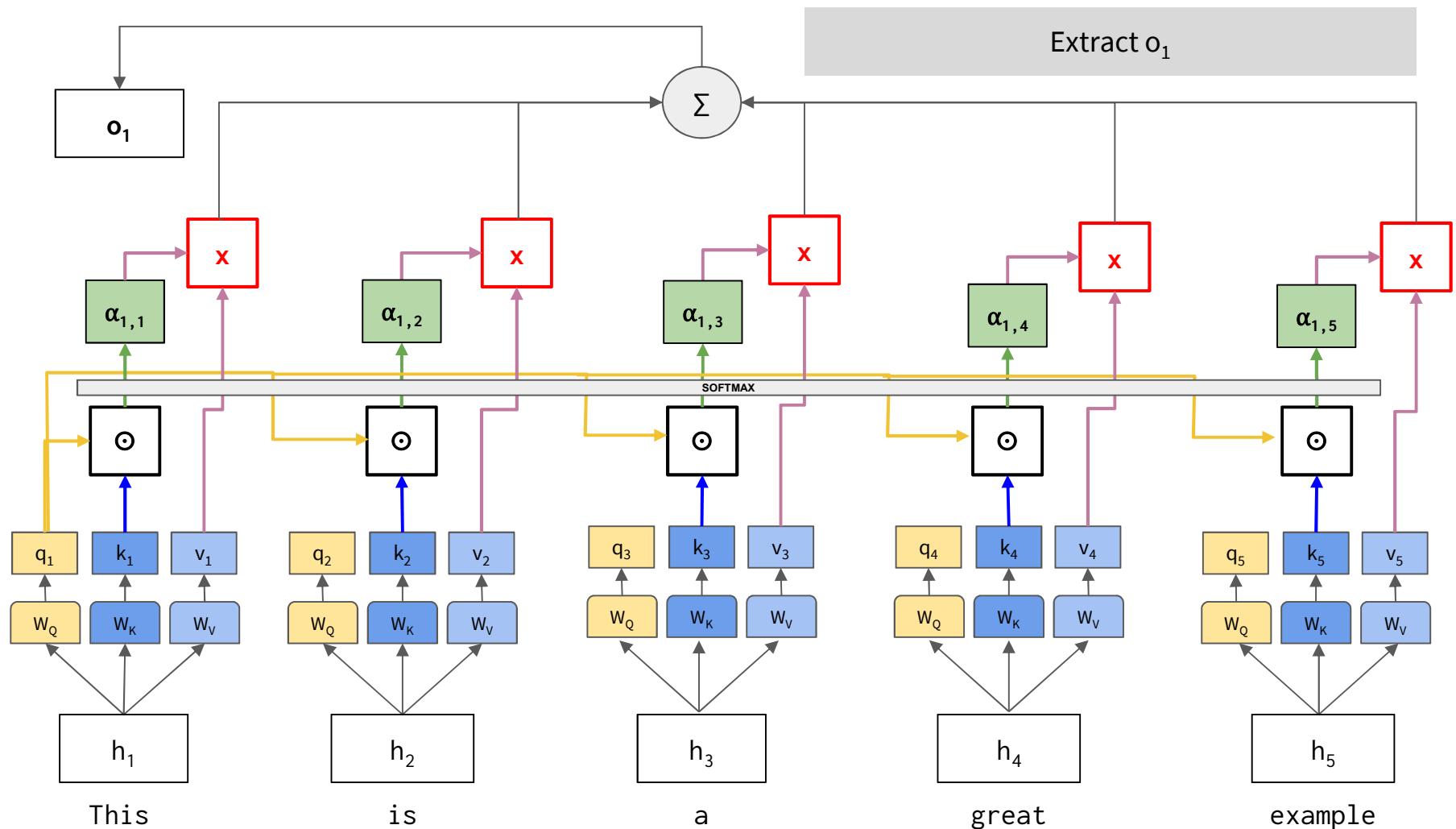


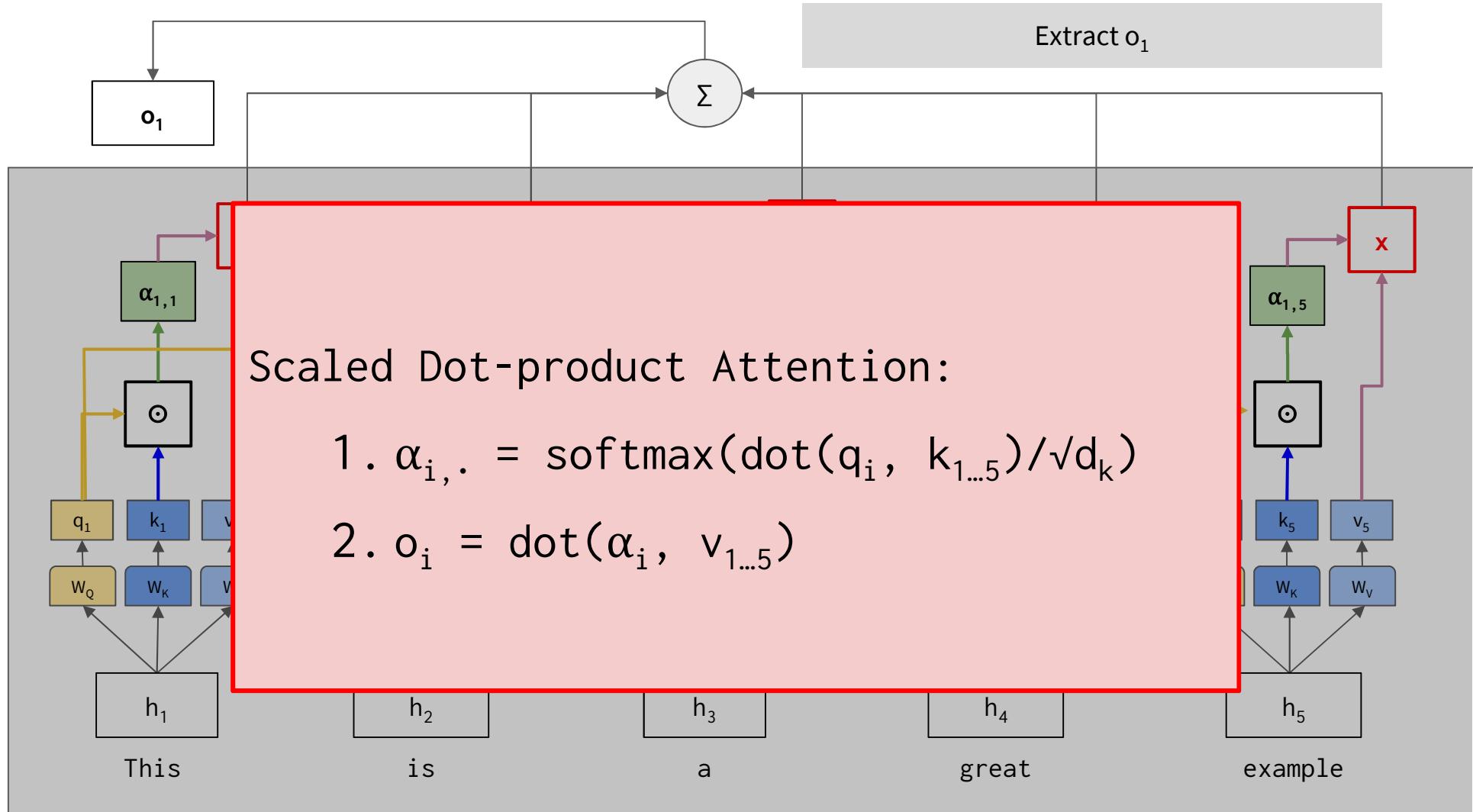


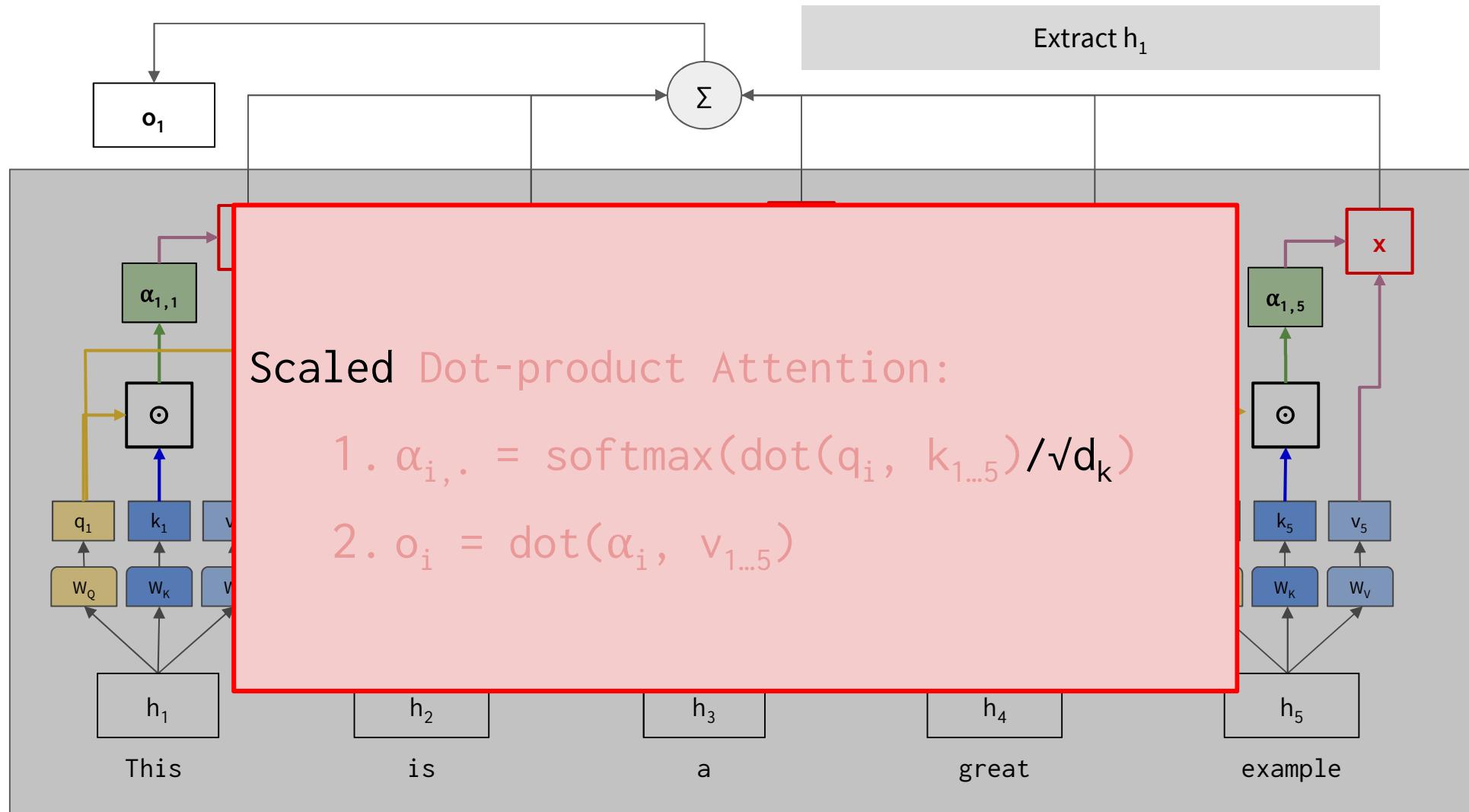
Multiply each $\alpha_{1,i}$ with v_i

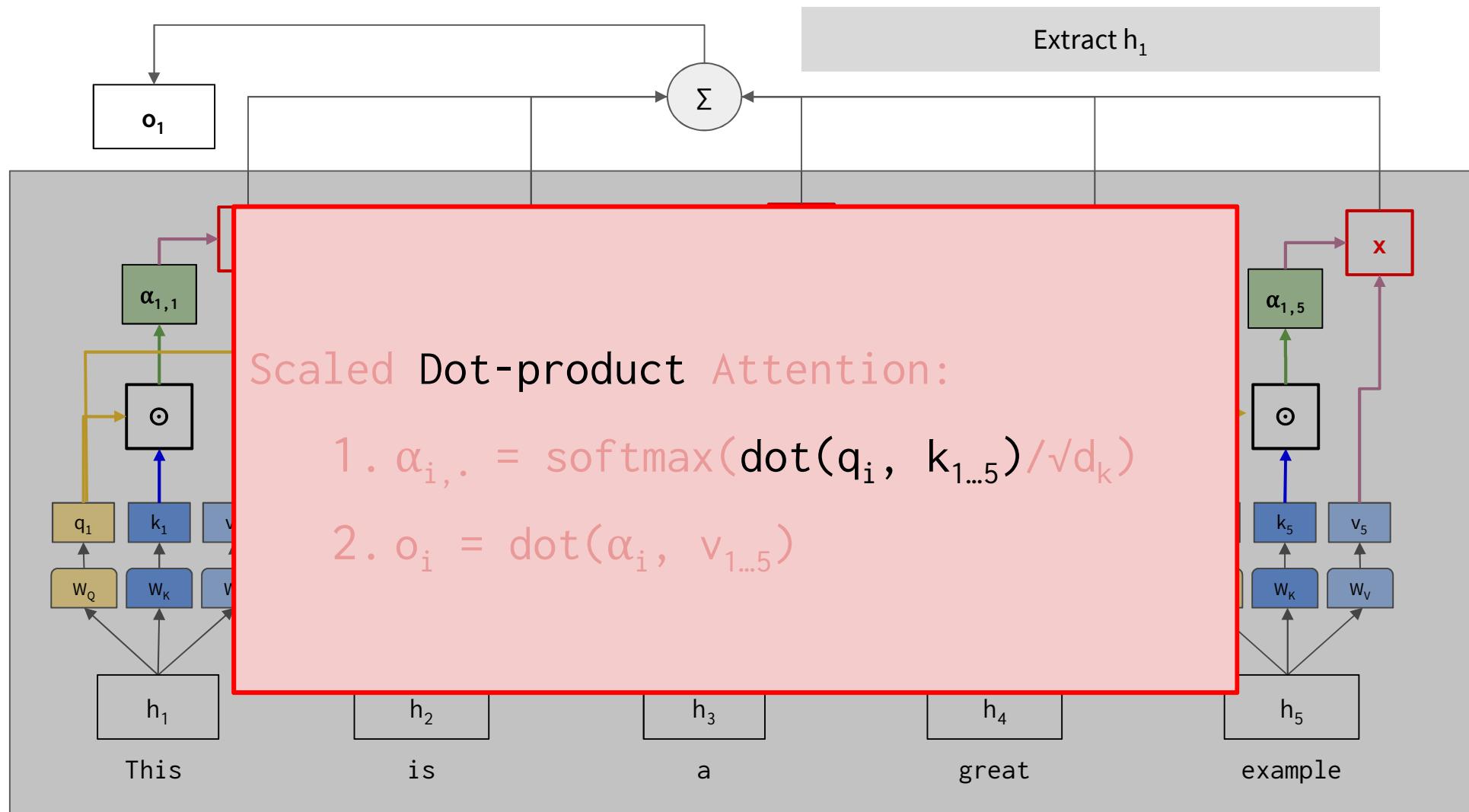


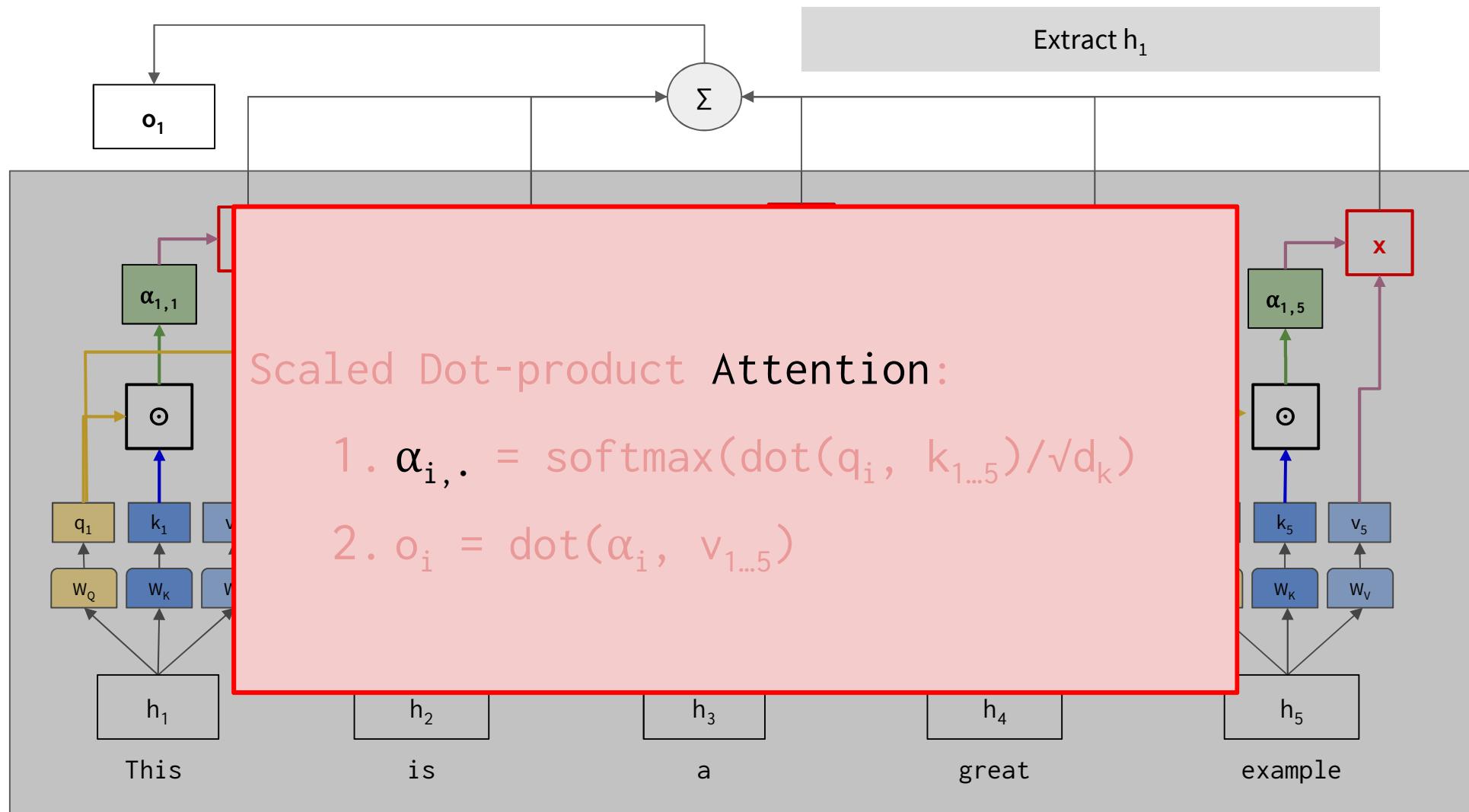


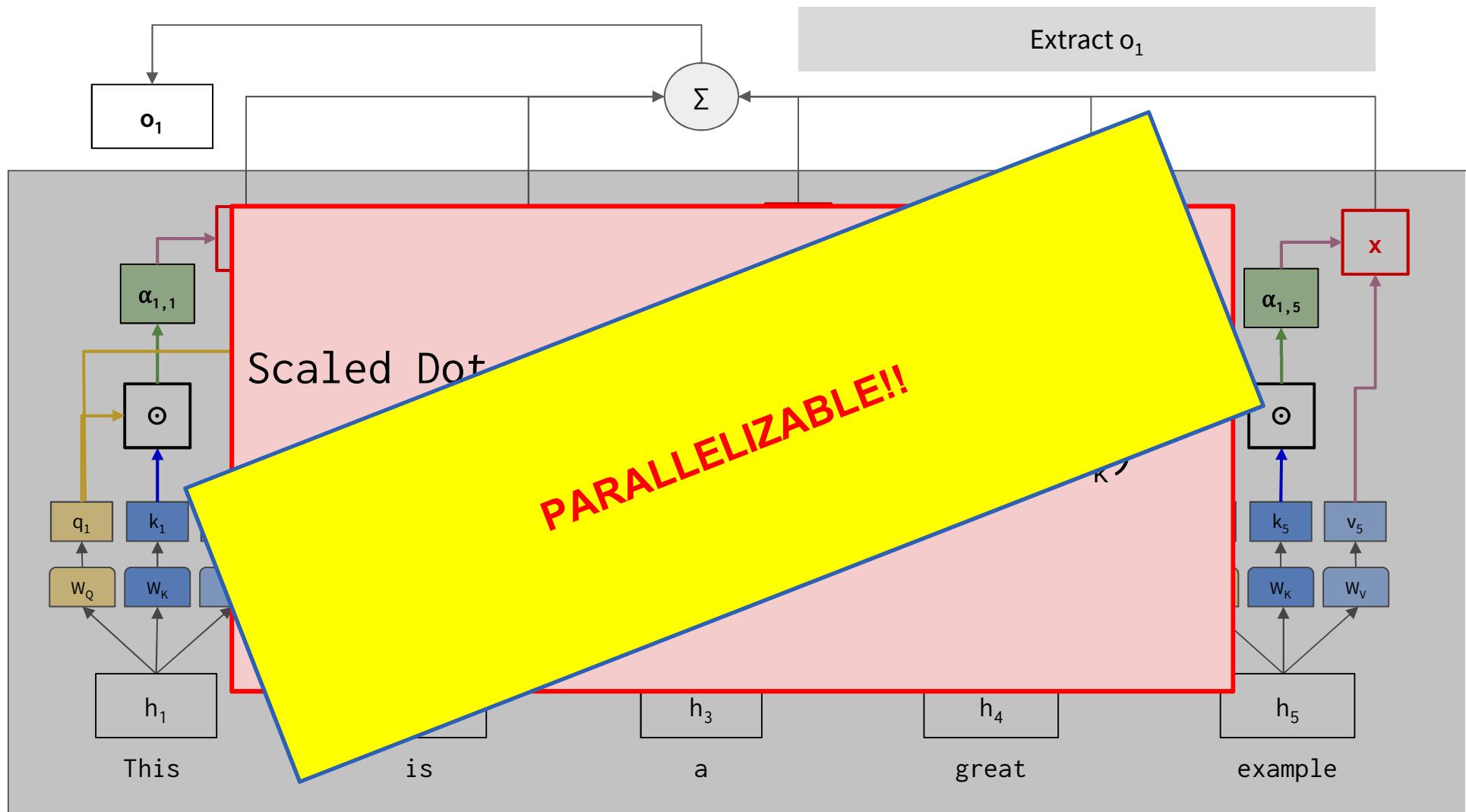








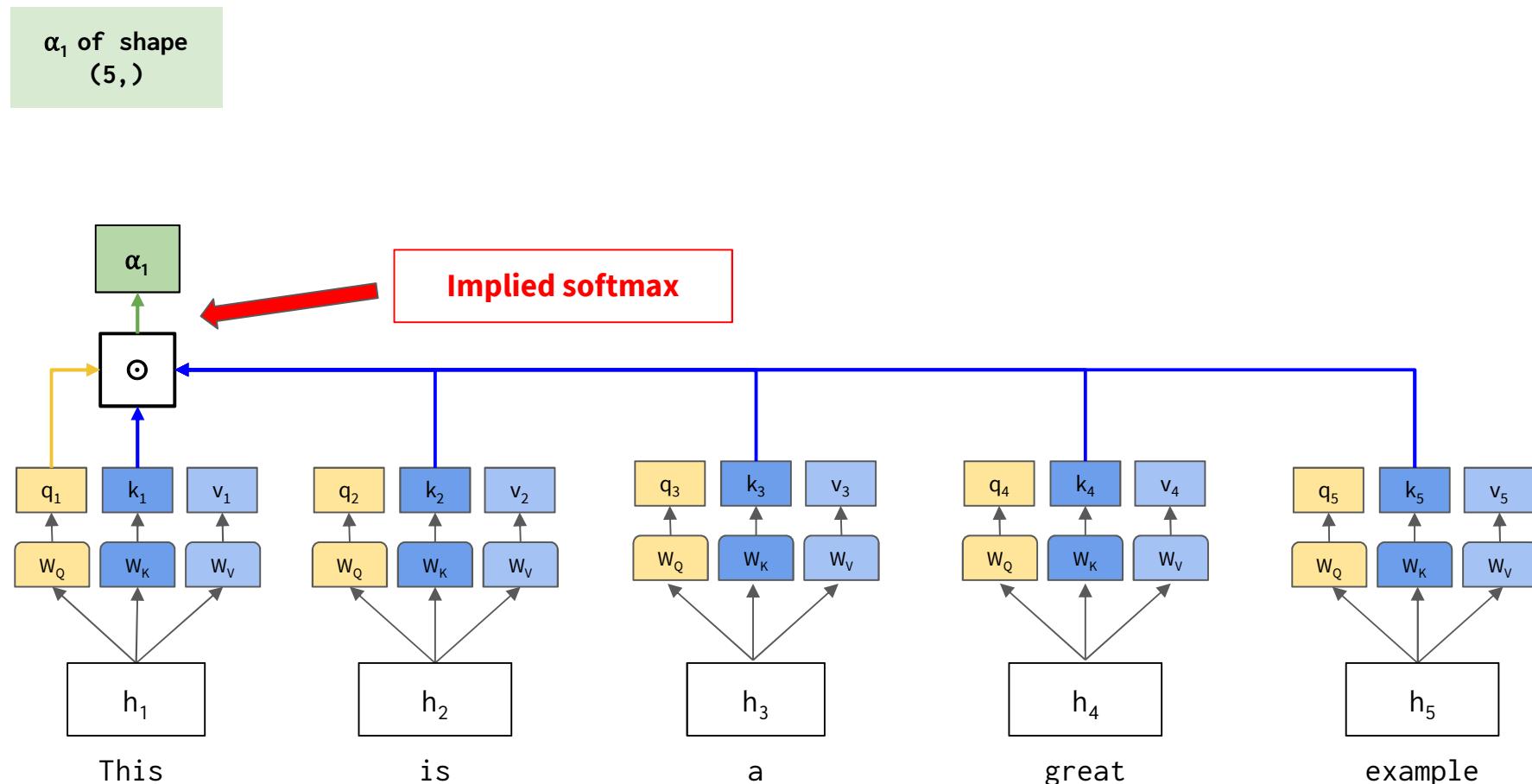




Example

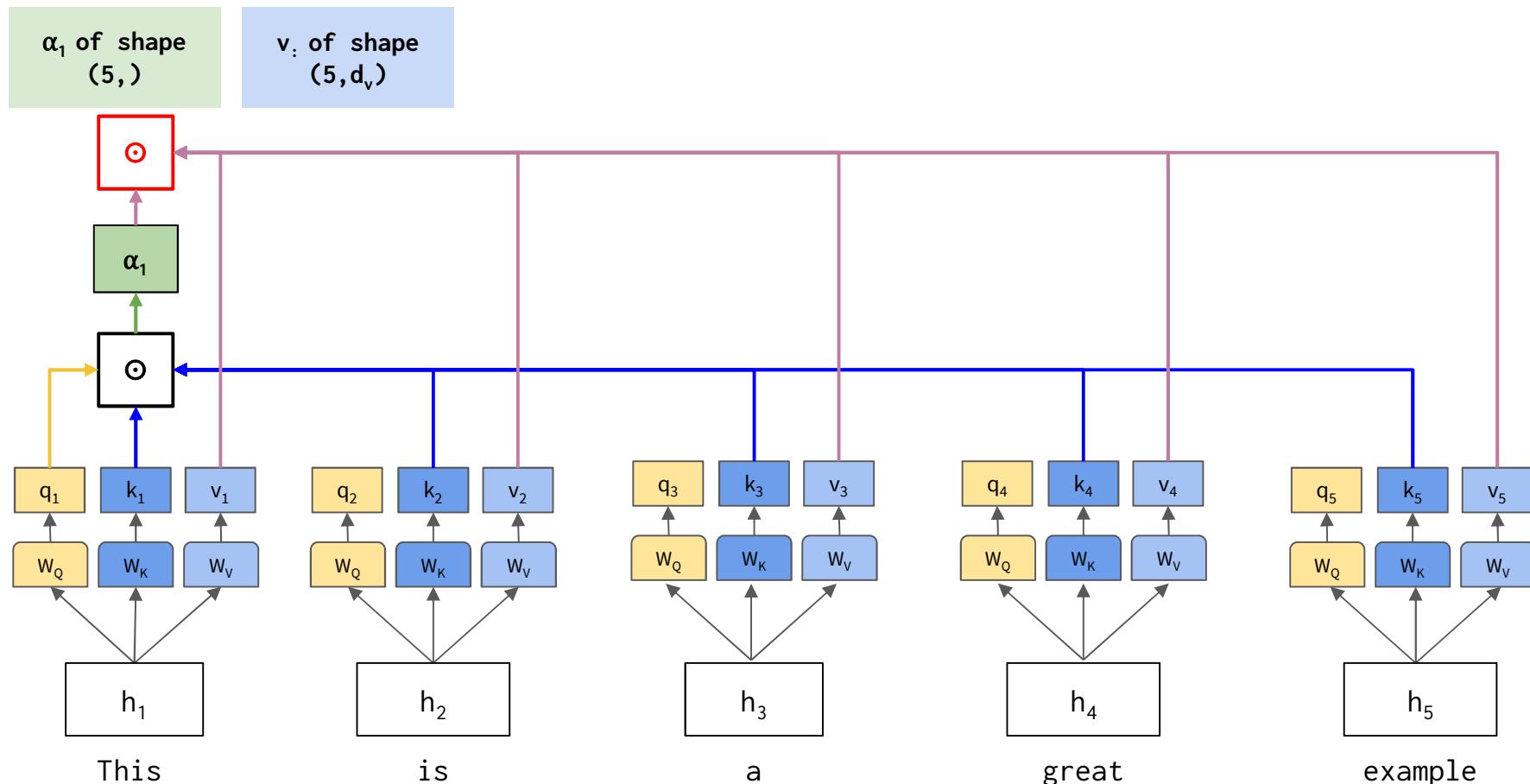
- $q, k, v_1 = [1, 2, 3, 4]$
 $q, k, v_2 = [4, 5, 9, 1]$
 $q, k, v_3 = [6, 2, 1, 4]$
- $e_1 = q_1 k_1^T / \sqrt{4} = 15.0$
 $e_2 = q_1 k_2^T / \sqrt{4} = 22.5$
 $e_3 = q_1 k_3^T / \sqrt{4} = 14.5$
- $\alpha_1 = \text{softmax}(e) = [0.00055, 0.99911, 0.00033]$
- $o_1 = \alpha_1^T V = [3.99, 4.99, 8.99, 1.00]$
 V is the 3×4 matrix of all values

Self Attention



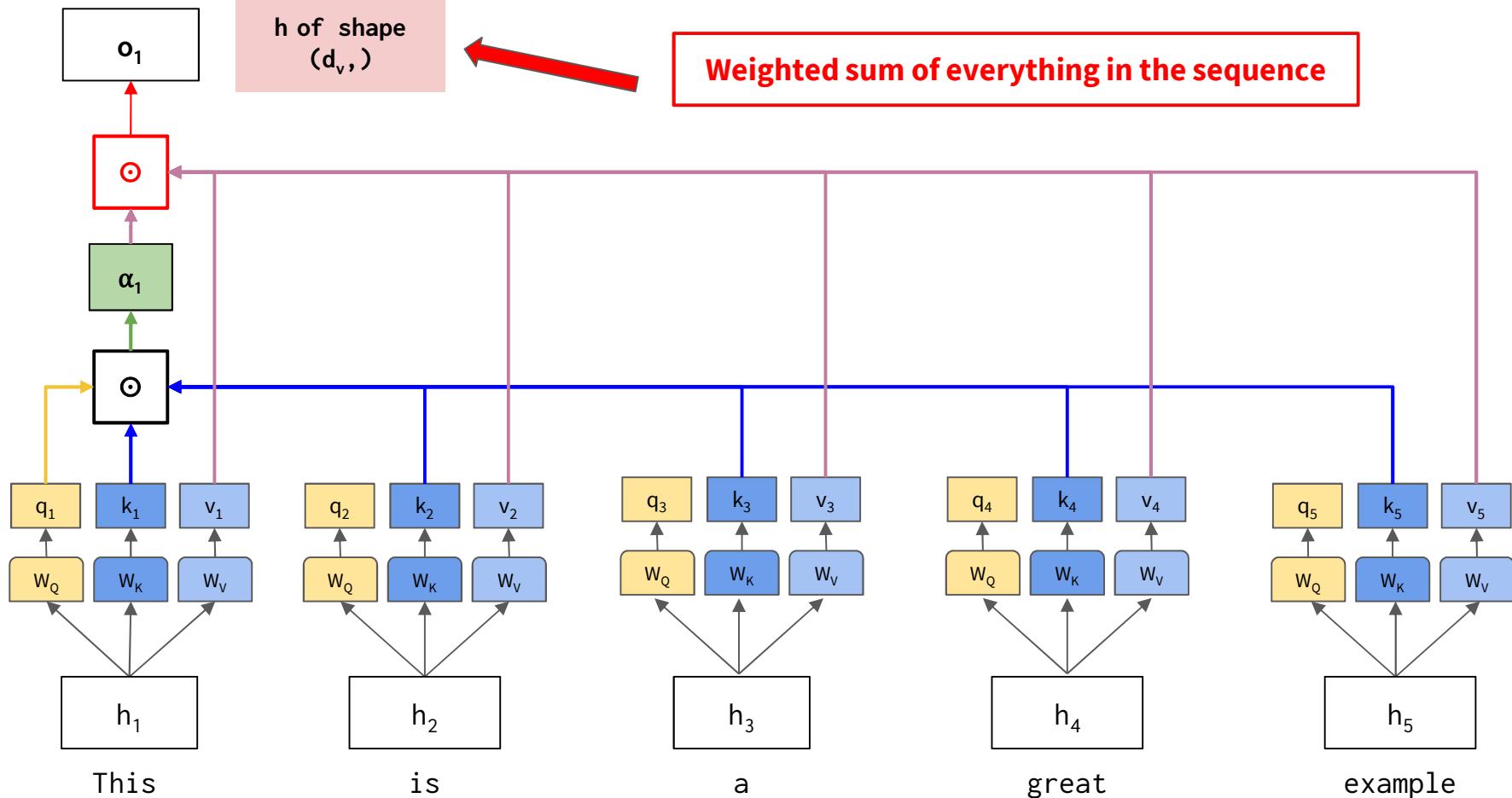
Self Attention

*Implied softmax



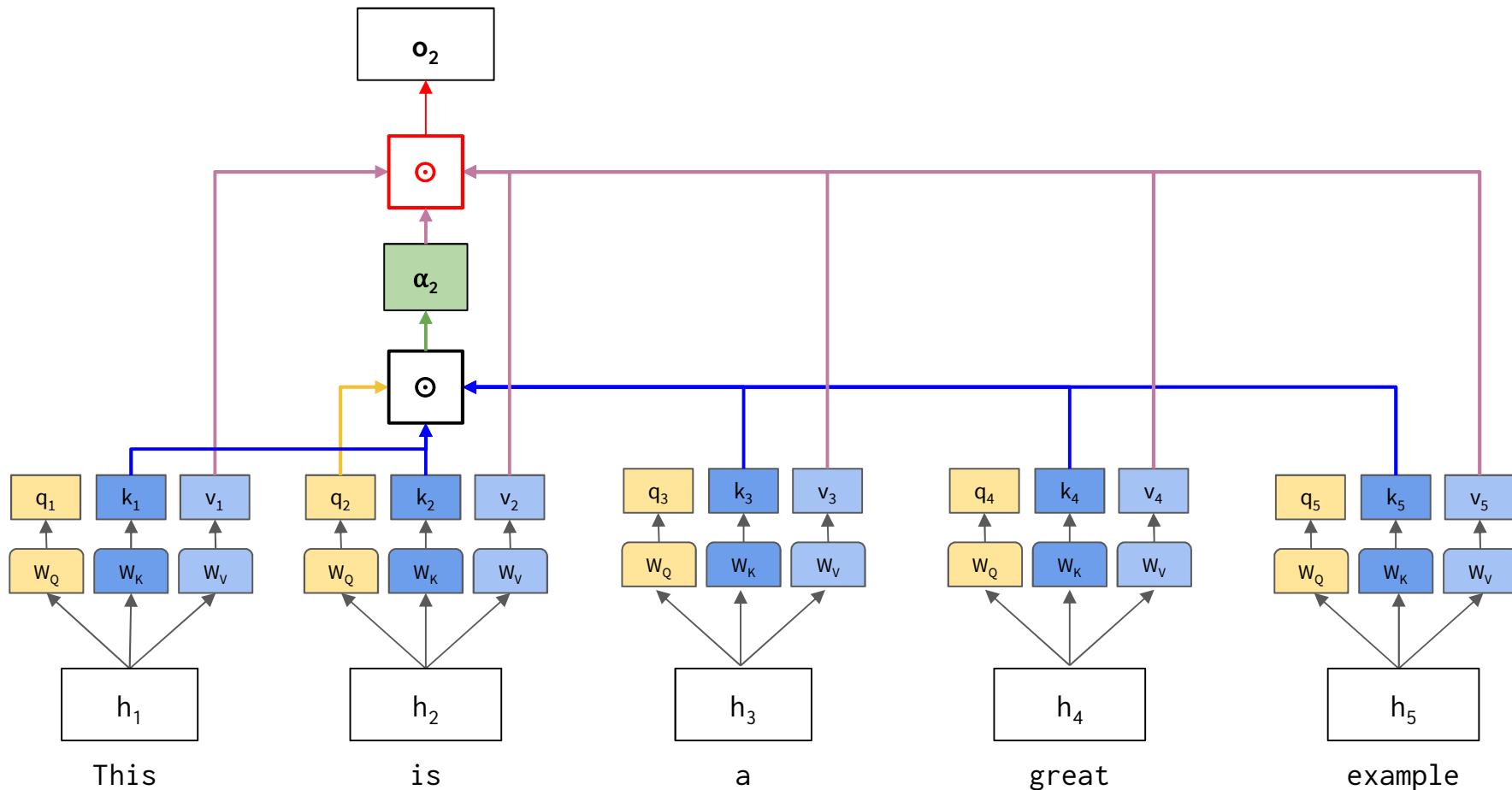
Self Attention

*Implied softmax



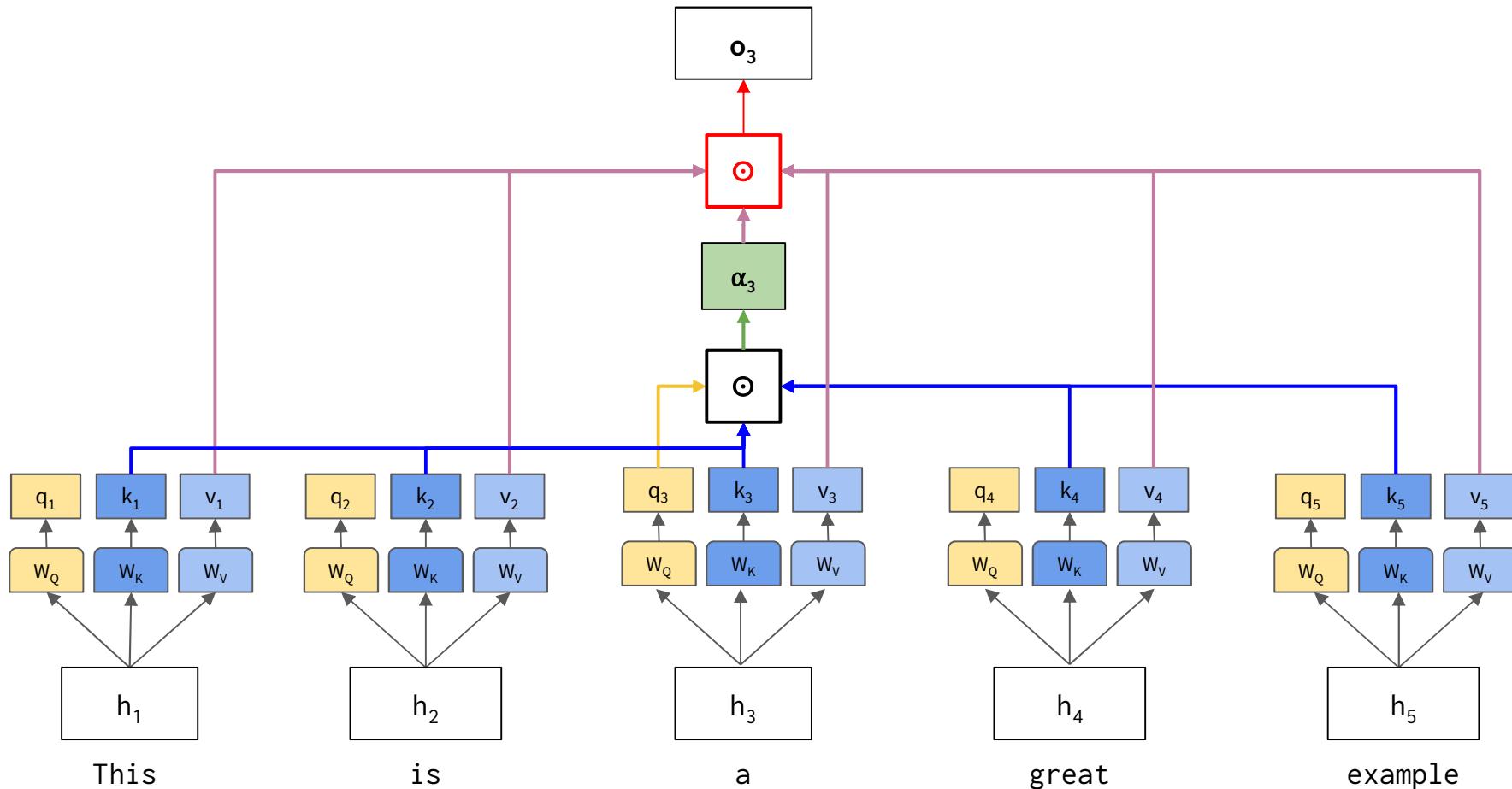
Self Attention

*Implied softmax



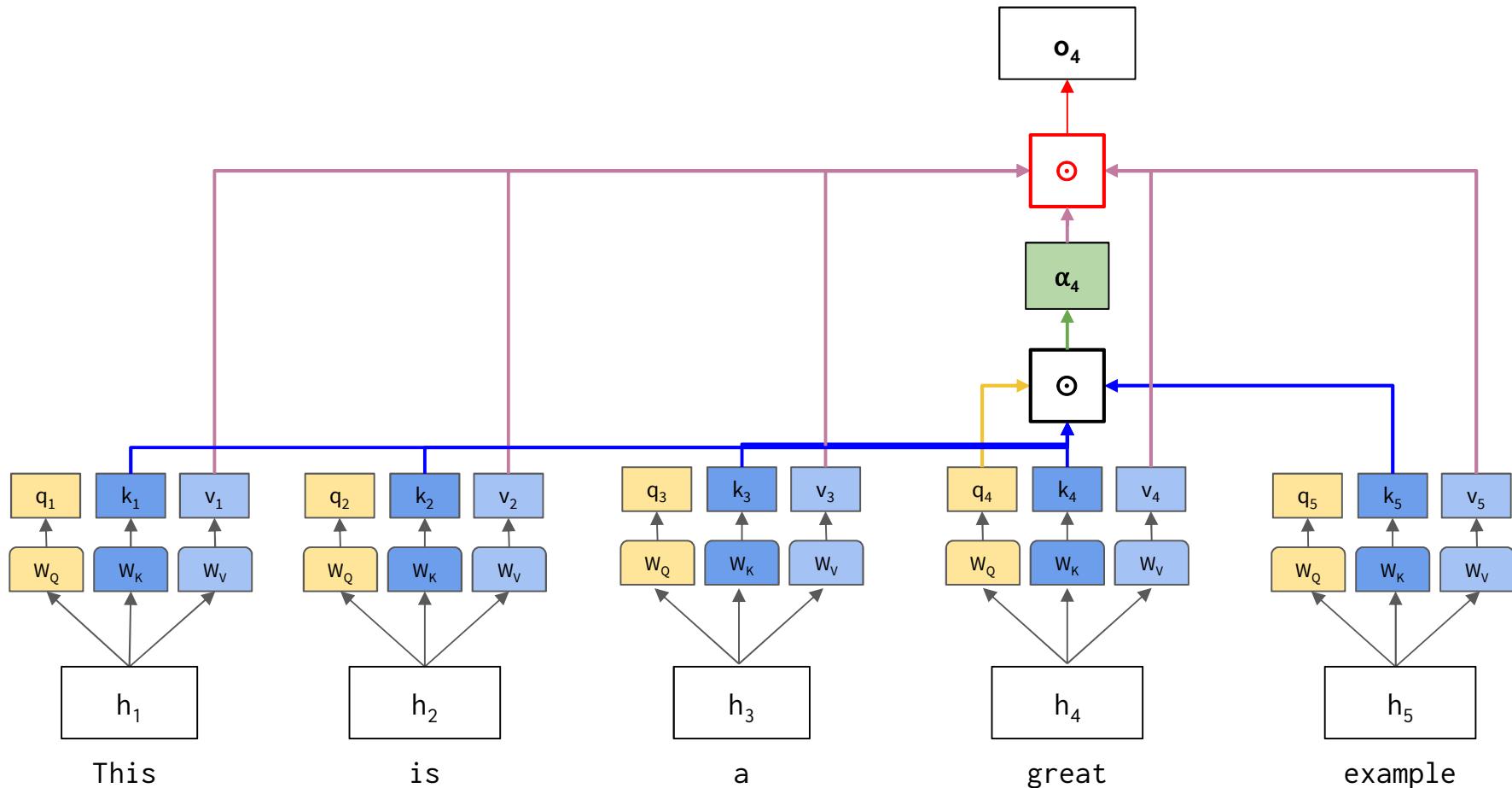
Self Attention

*Implied softmax



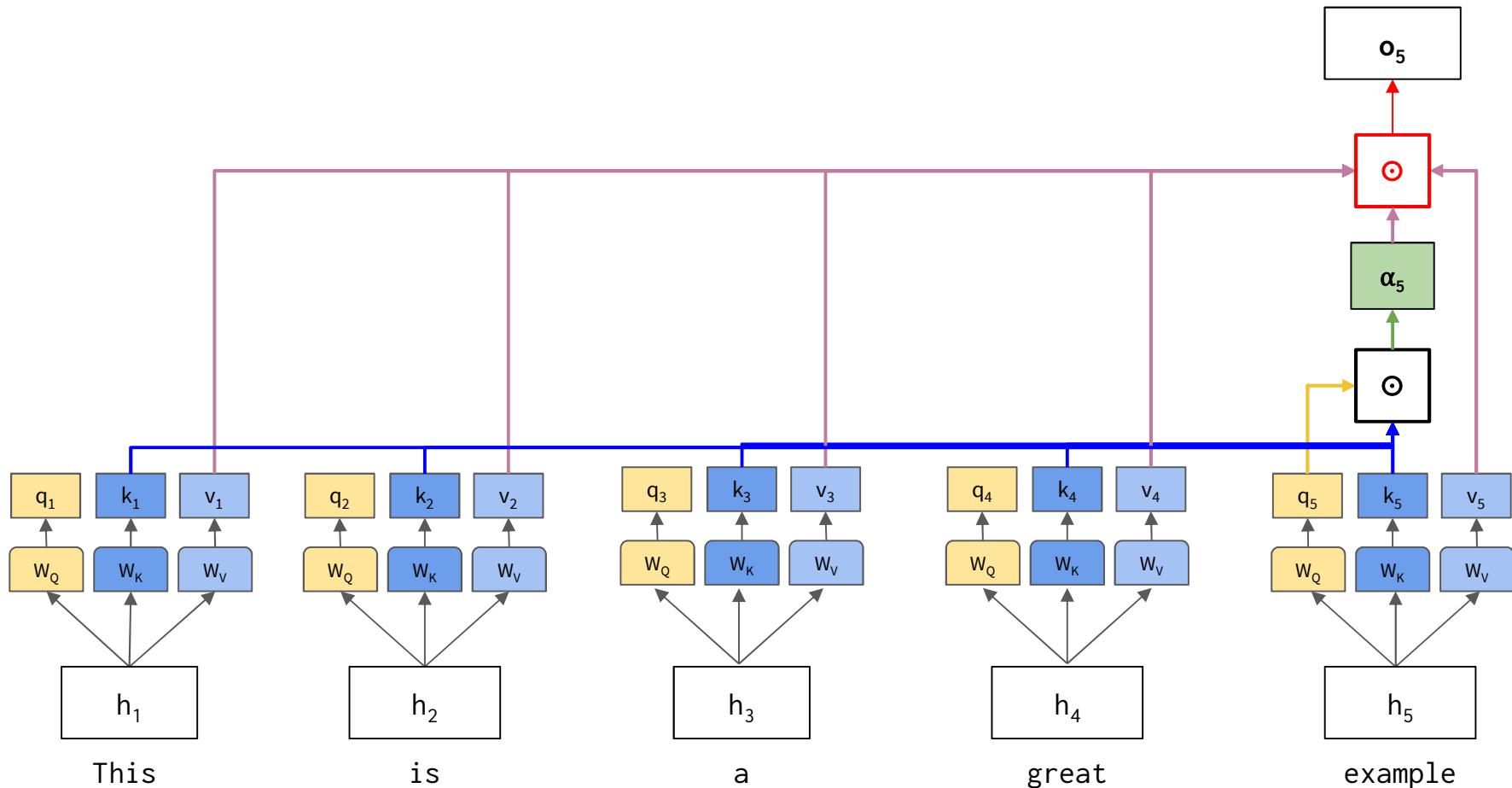
Self Attention

*Implied softmax



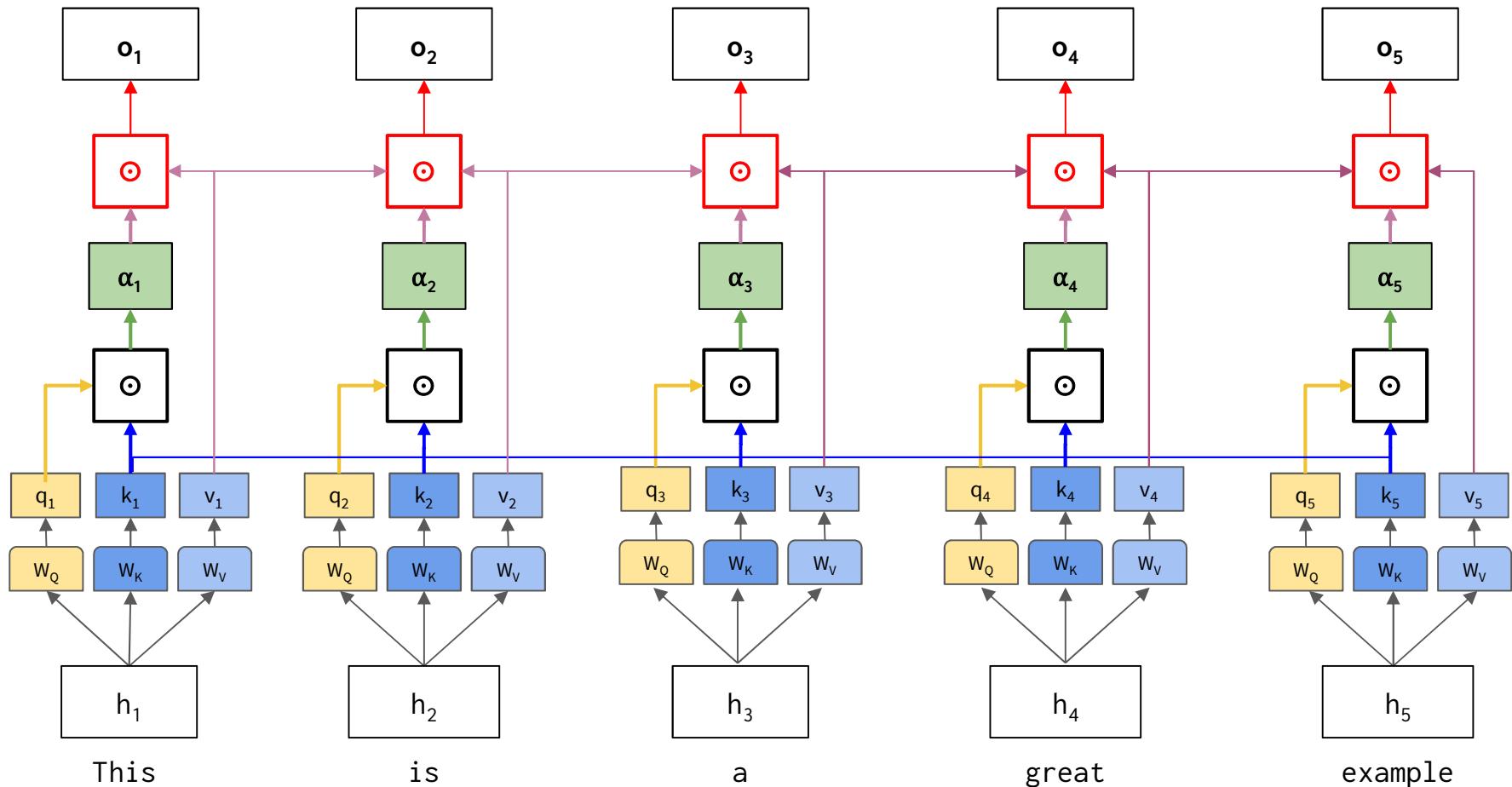
Self Attention

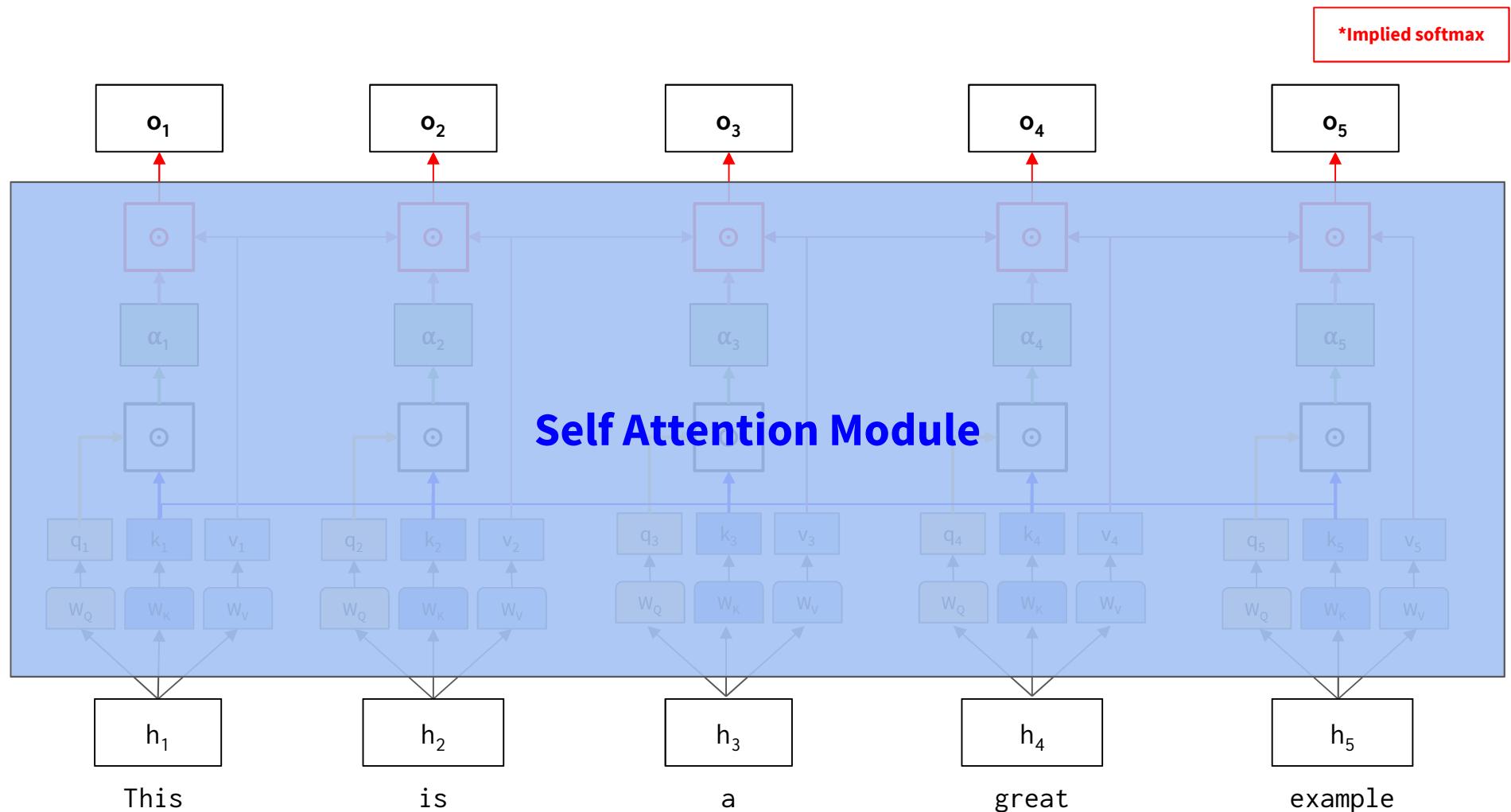
*Implied softmax

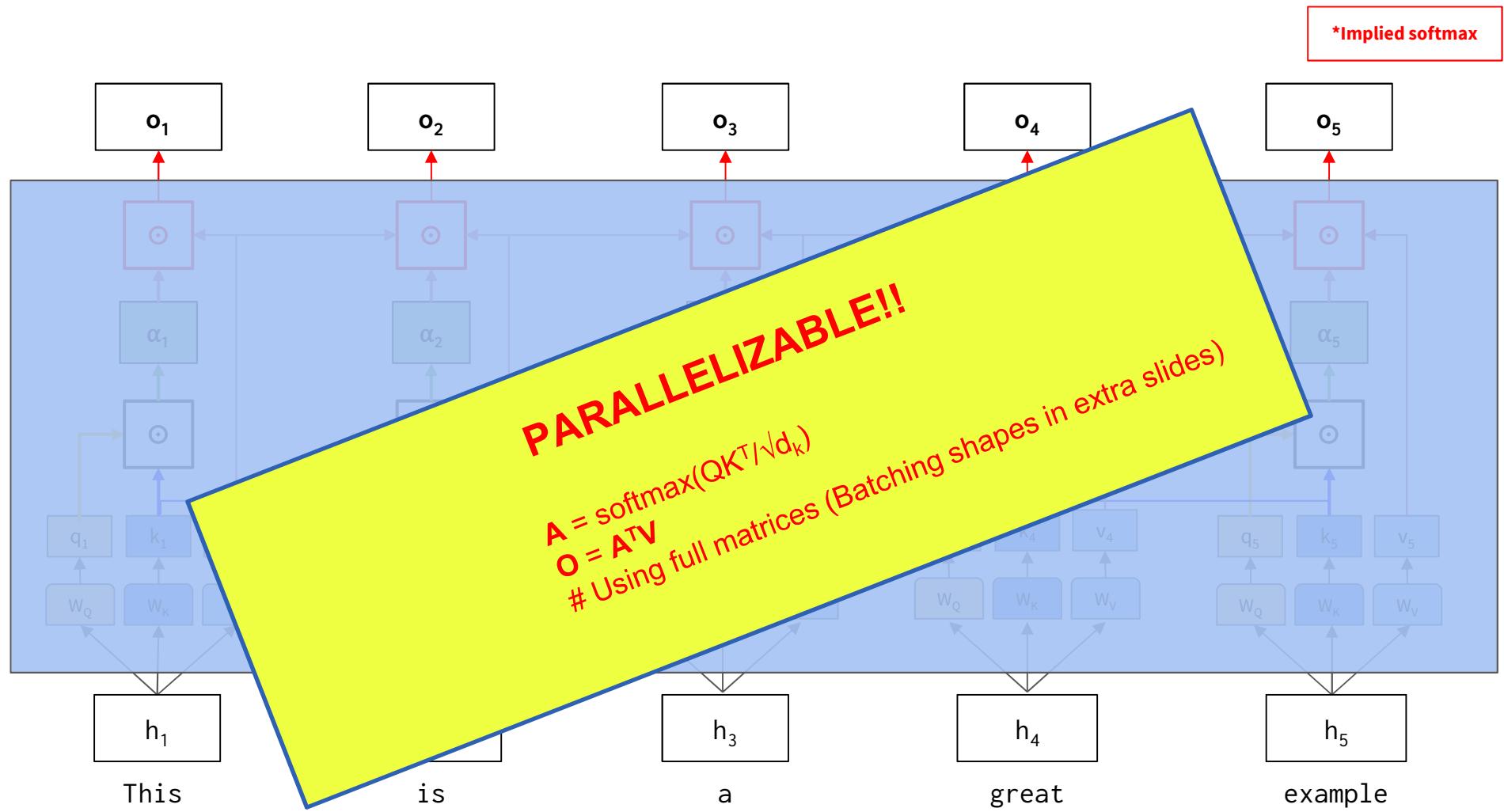


Self Attention

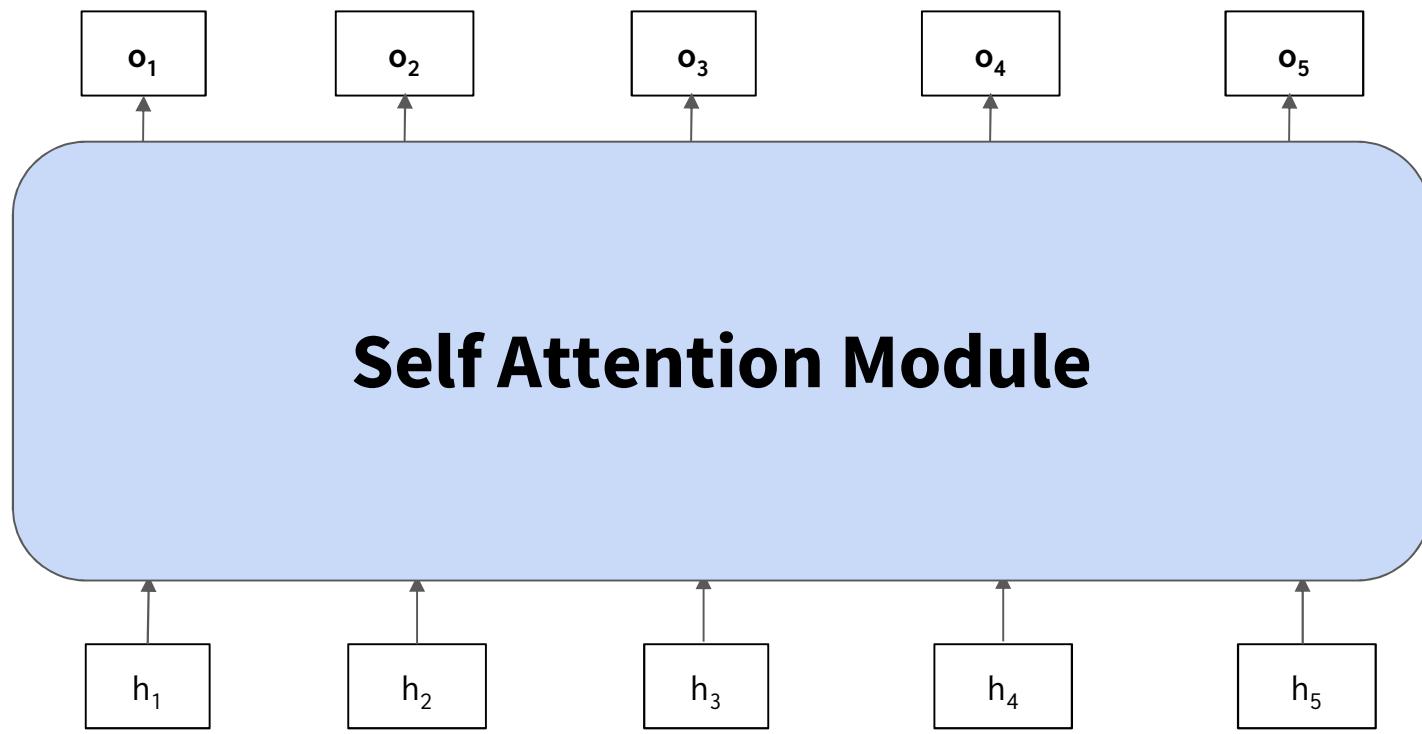
*Implied softmax







Single Headed Self Attention



Poll 1 (@1442)

Which of the following are true about self attention? (Select all that apply)

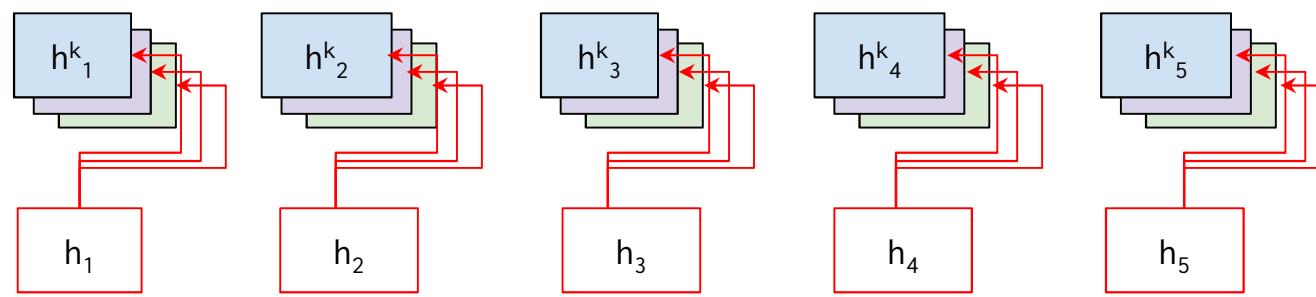
- a. To calculate attention weights for input h_i , you would use key k_i , and all queries
- b. To calculate attention weights for input h_i , you would use query q_i , and all keys
- c. The energy function is scaled to bring attention weights in the range of [0,1]
- d. The energy function is scaled to allow for numerical stability

Poll 1 (@1442)

Which of the following are true about self attention? (Select all that apply)

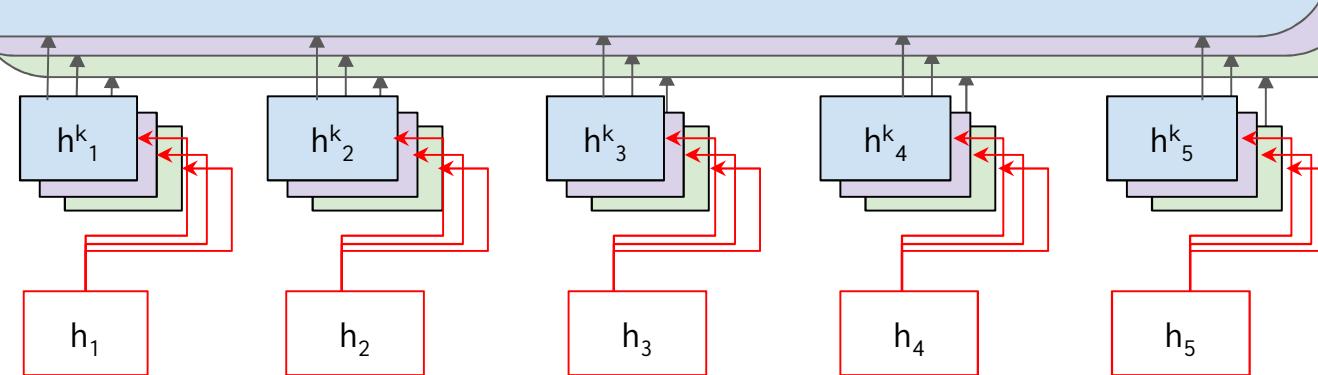
- a. To calculate attention weights for input h_i , you would use key k_i , and all queries
- b. **To calculate attention weights for input h_i , you would use query q_i , and all keys**
- c. The energy function is scaled to bring attention weights in the range of [0,1]
- d. **The energy function is scaled to allow for numerical stability**

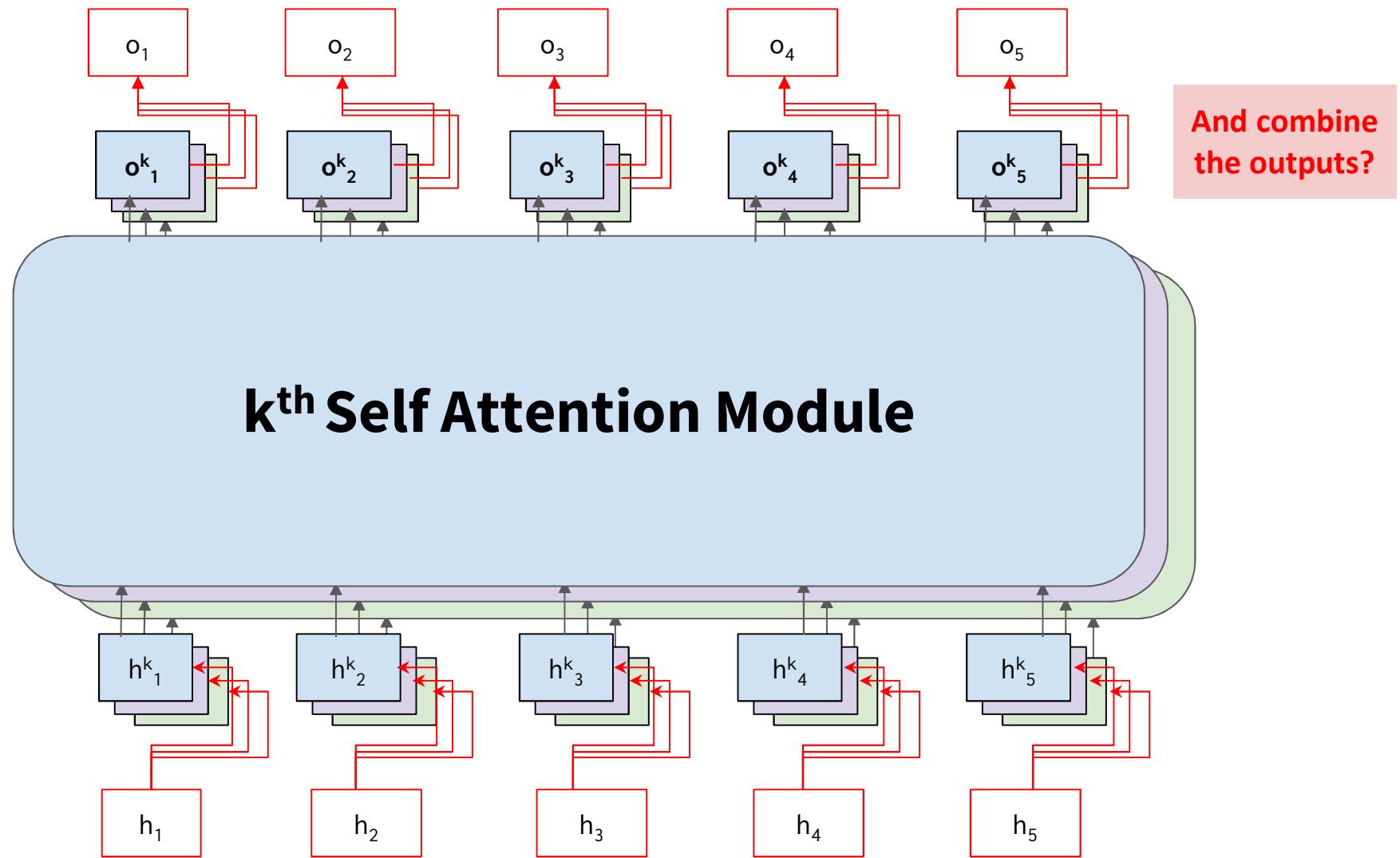
What if we split the input into 'k' sub-inputs?

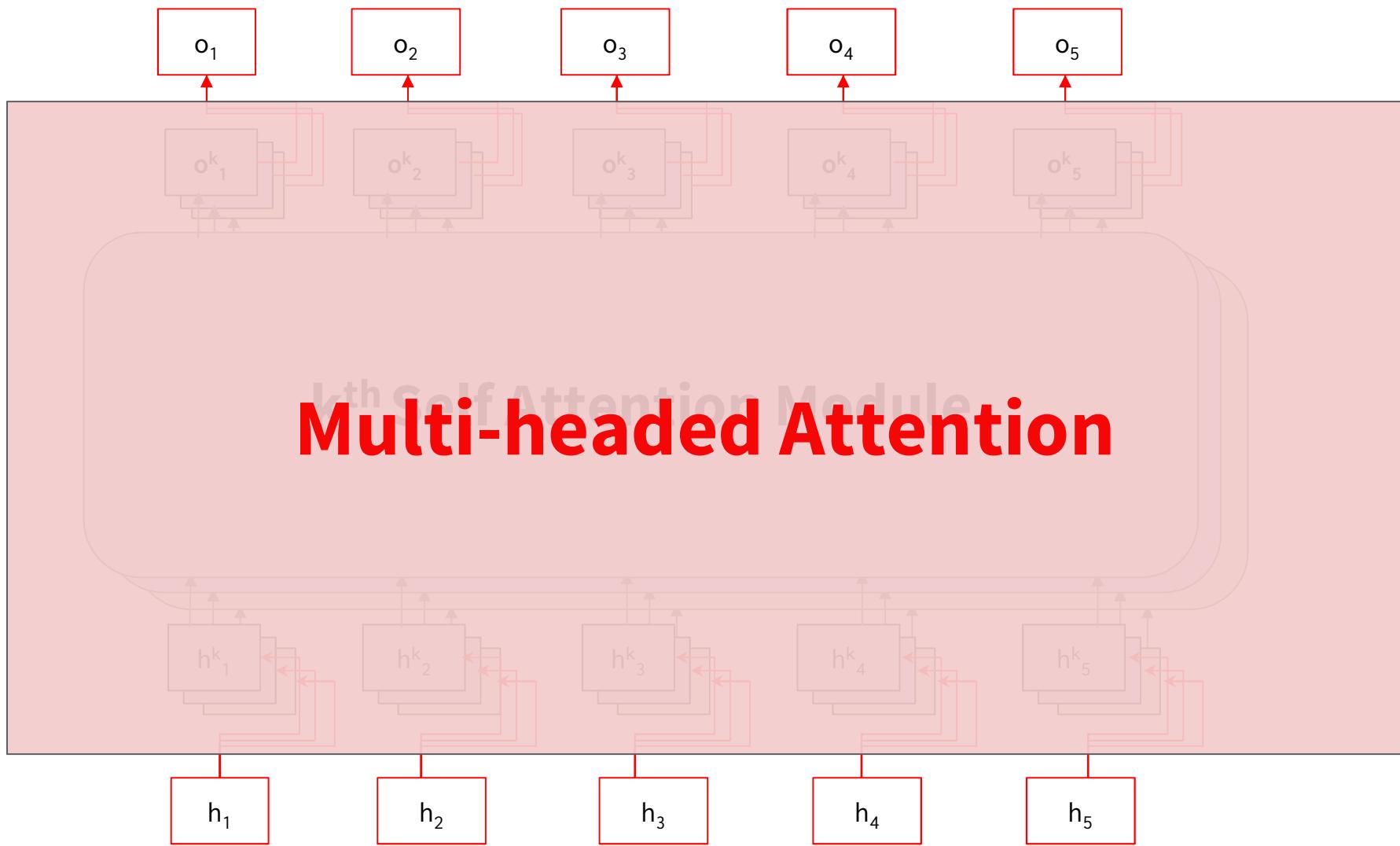


And pass each sub-input into a Self-Attention Module?

k^{th} Self Attention Module







Multi Headed Self Attention

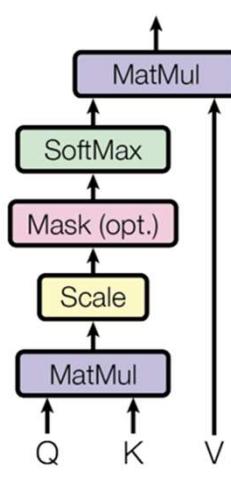
- Split input into **k** parts
- Pass the j^{th} part of **each input** into the j^{th} **attention head**
- Concatenate each of the **k** outputs

Why go through the trouble?

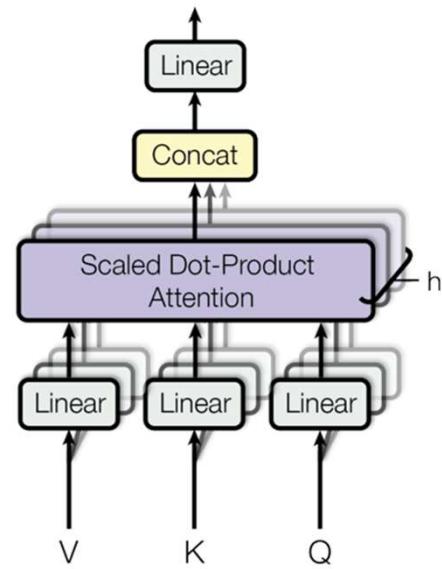
- Each head **could** find a different kind of relation between the tokens
 - Subject-verb, subject-object, verb-modifier, dependency, etc.

Attention is all you need

Scaled Dot-Product Attention

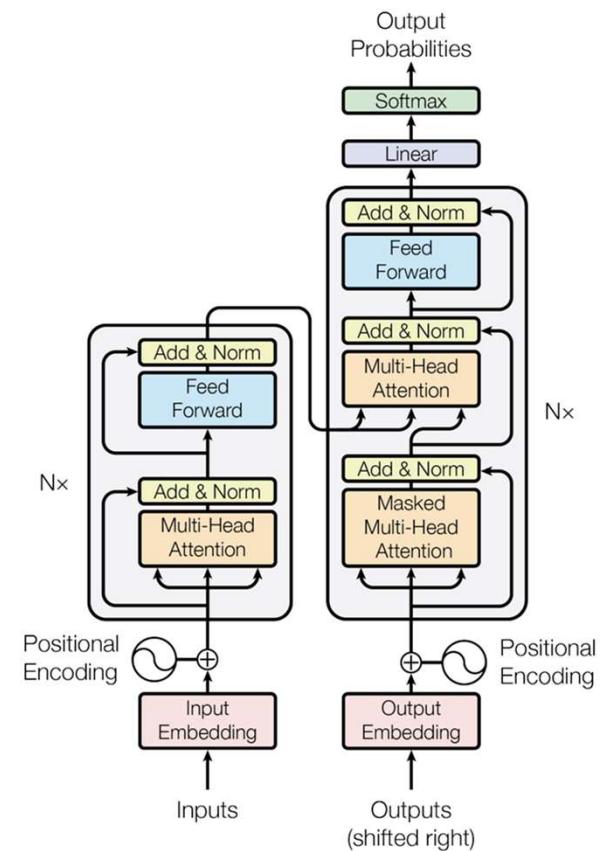


Multi-Head Attention



Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

Breaking down the transformer



Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

Figure 1: The Transformer - model architecture.

Breaking down the transformer

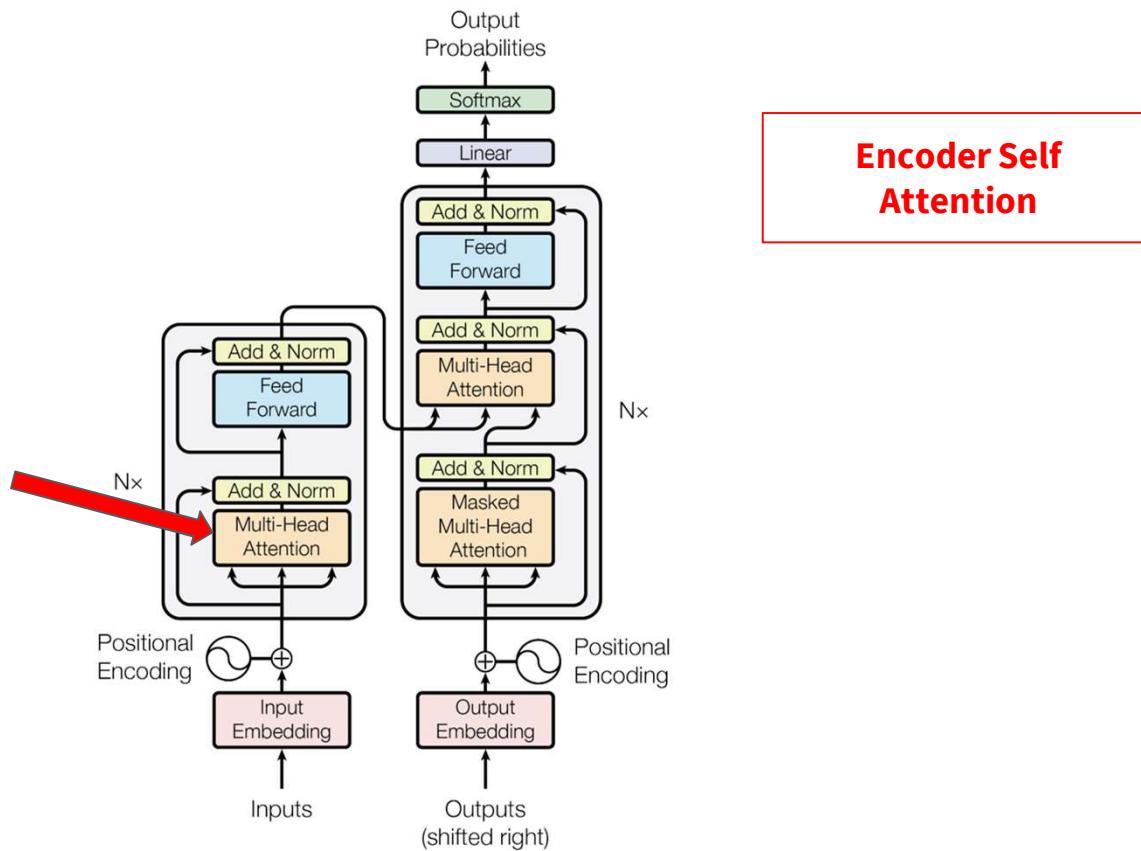
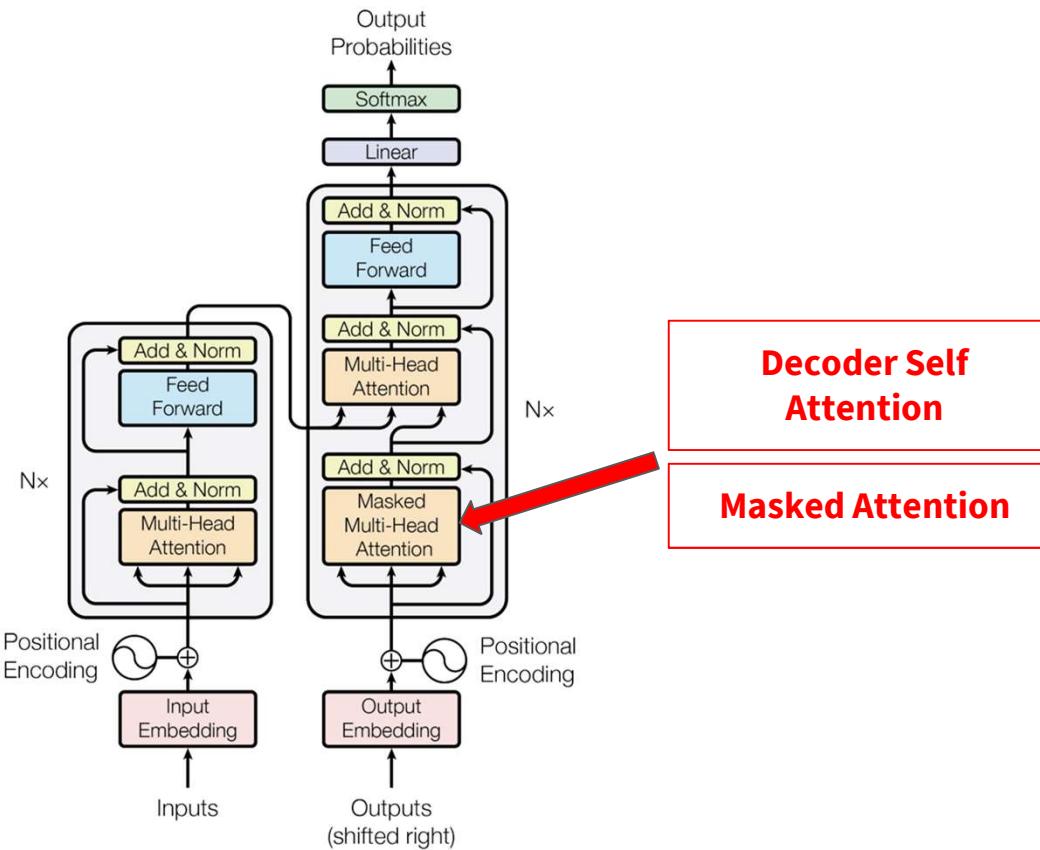


Figure 1: The Transformer - model architecture.

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

Breaking down the transformer



Decoders can be parallelized during **training** (only)

Feed the whole output sequence (outputs) in at once

Need to ensure model doesn't cheat

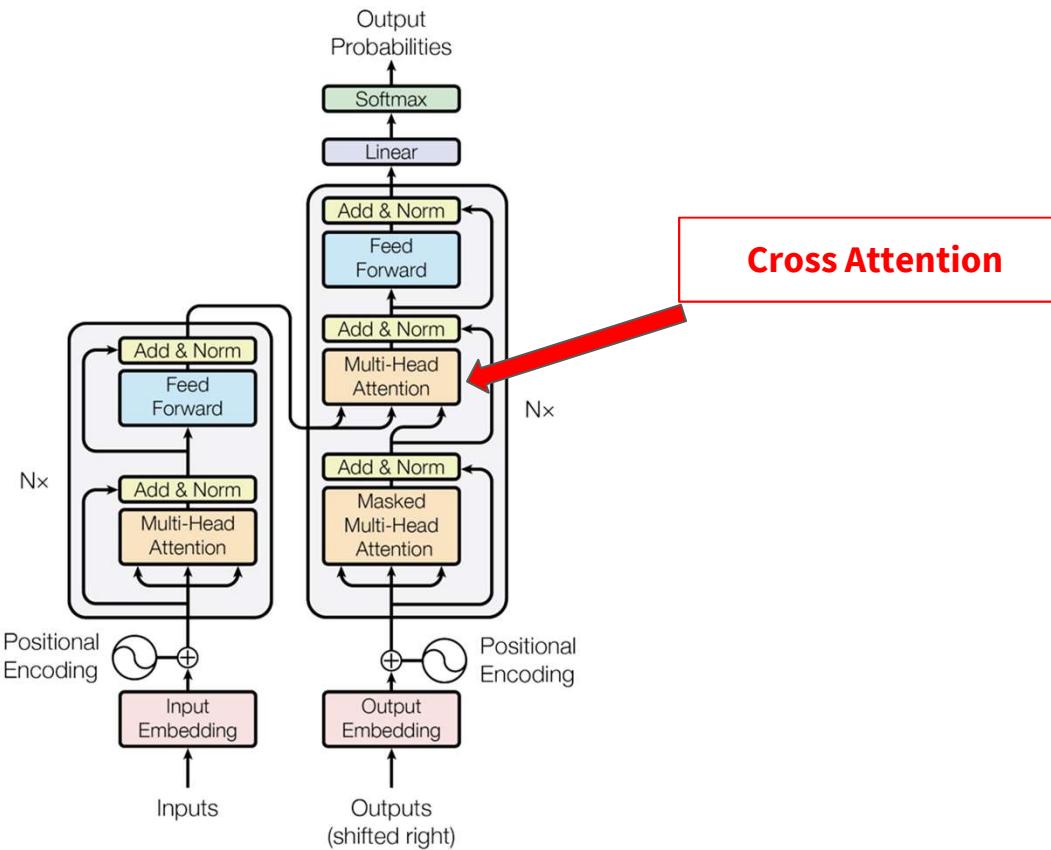
Alter the attention weights to be **0** (set input to softmax to **-inf**) for all times $t' > t$

Ensure **autoregressive property**

Figure 1: The Transformer - model architecture.

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

Breaking down the transformer



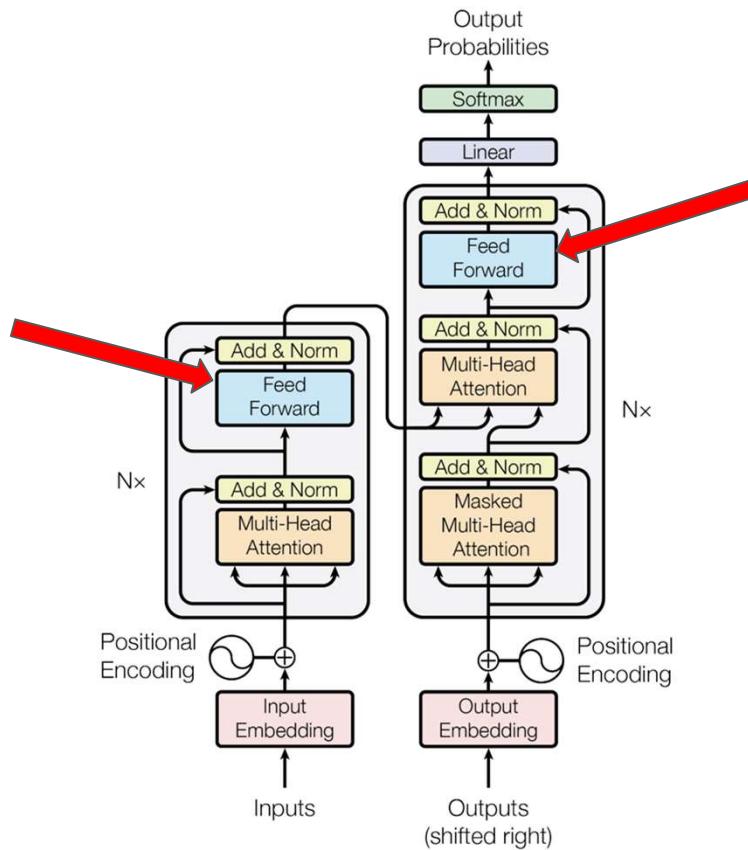
During decoding, the **query** comes from the outputs, **keys** and **values** come from the encoder.

Decoder input “**pays** attention” to the encoder outputs.

Figure 1: The Transformer - model architecture.

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

Breaking down the transformer



Feed Forward Layers

Feed Forward layers allow for high dimensional computations

Simply there to allow the model to capture more information

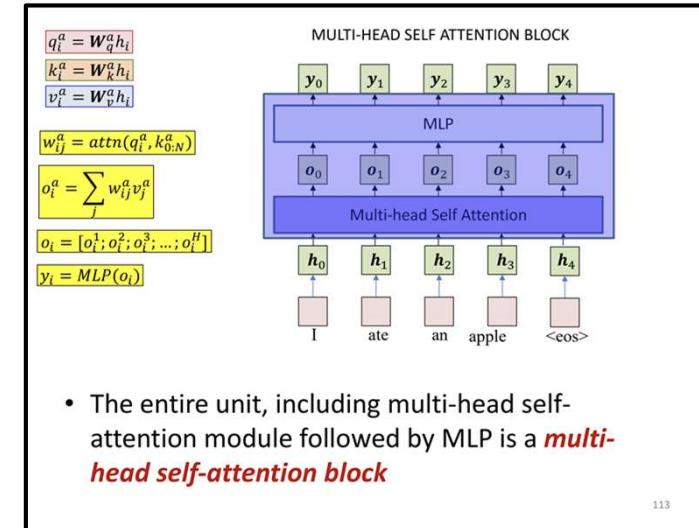
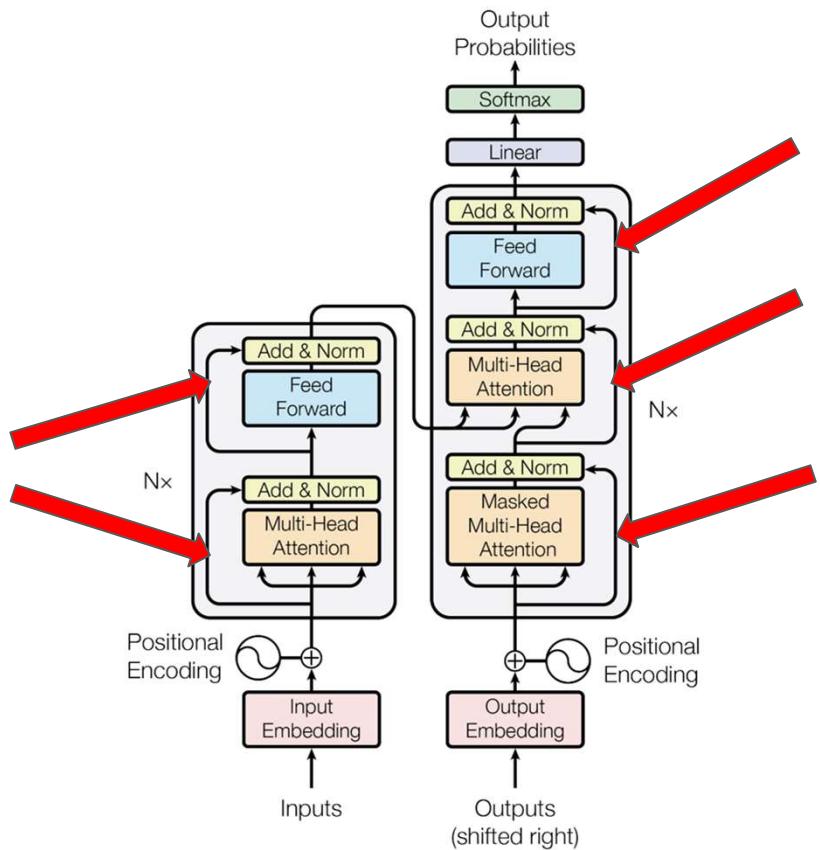


Figure 1: The Transformer - model architecture.

Breaking down the transformer



Residual Connections

Add & Norm:

```
out = LayerNorm(x + Sublayer(x)),
```

Sublayer(x) is whatever layer is below the **Add & Norm**

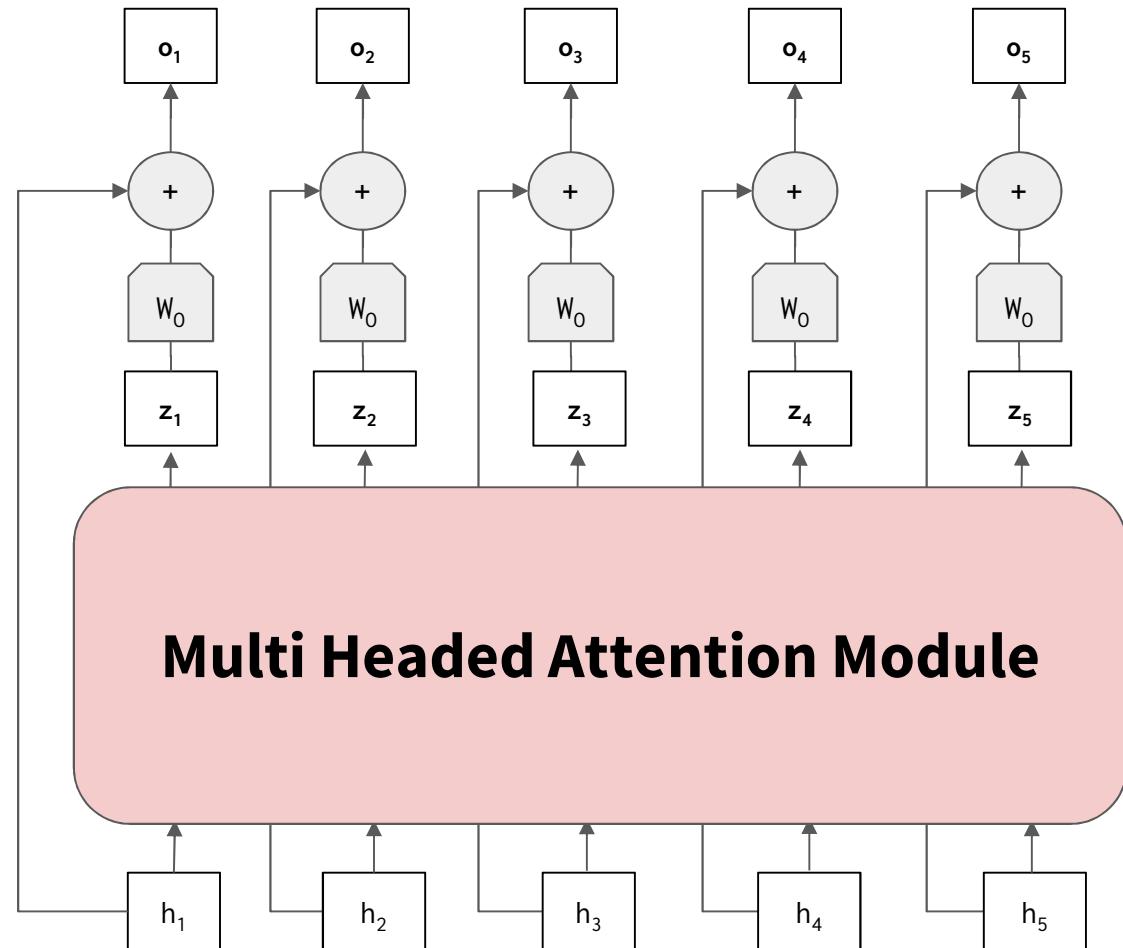
Figure 1: The Transformer - model architecture.

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

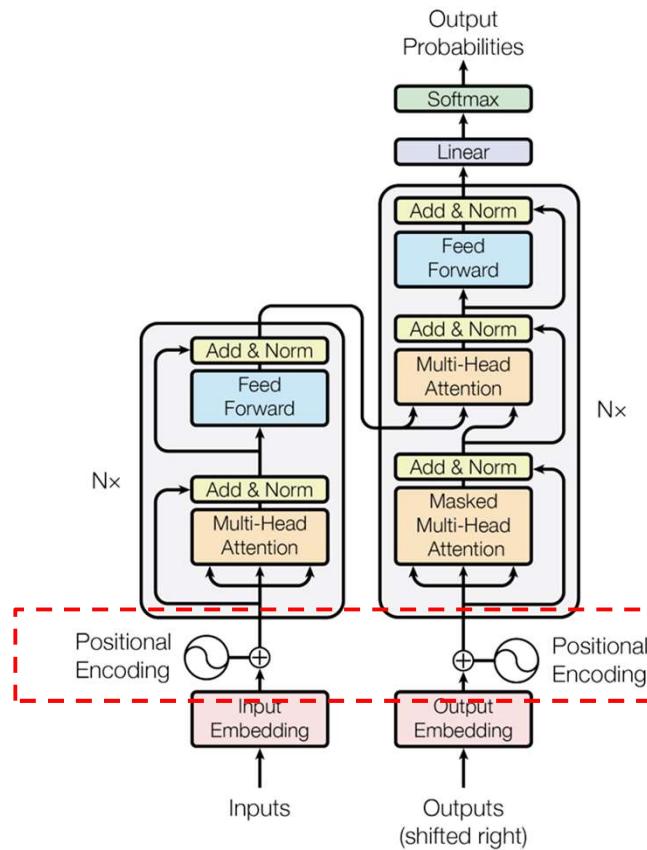
Residual Connection

Transformers are residual machines

"How much do I SHIFT my meaning given my context?"



Breaking down the transformer



Transformers have no inherent notion of order in a sequence.

This notion has to be externally enforced.

Positional Encodings are added to transformer inputs to add information about order.

Positional Encoding

Figure 1: The Transformer - model architecture.

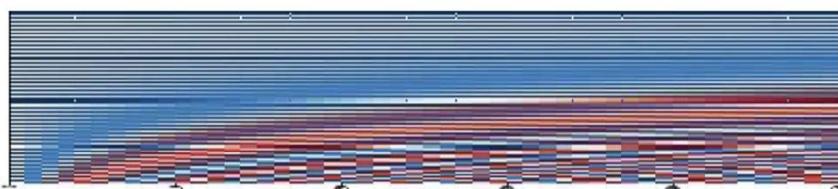
Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).

Recall

Positional encodings as discussed in the last lecture.

Positional Encoding

regenerate

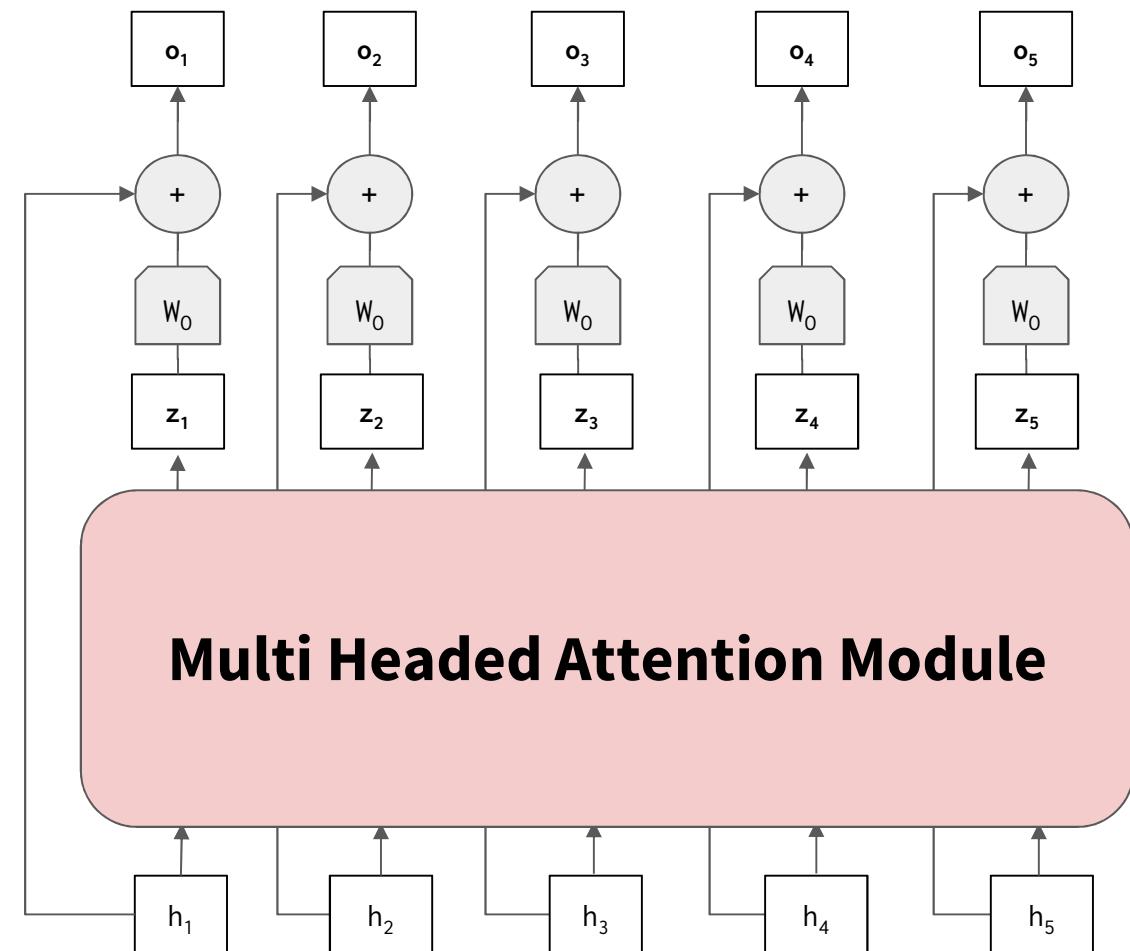


$$P_t = \begin{bmatrix} \sin \omega_1 t \\ \cos \omega_1 t \\ \sin \omega_2 t \\ \cos \omega_2 t \\ \vdots \\ \sin \omega_{d/2} t \\ \cos \omega_{d/2} t \end{bmatrix} \quad \omega_l = \frac{1}{10000^{2l/d}}$$

$$P_{t+\tau} = M_\tau P_t$$
$$M_\tau = \text{diag} \left(\begin{bmatrix} \cos \omega_l \tau & \sin \omega_l \tau \\ -\sin \omega_l \tau & \cos \omega_l \tau \end{bmatrix}, l = 1 \dots d/2 \right)$$

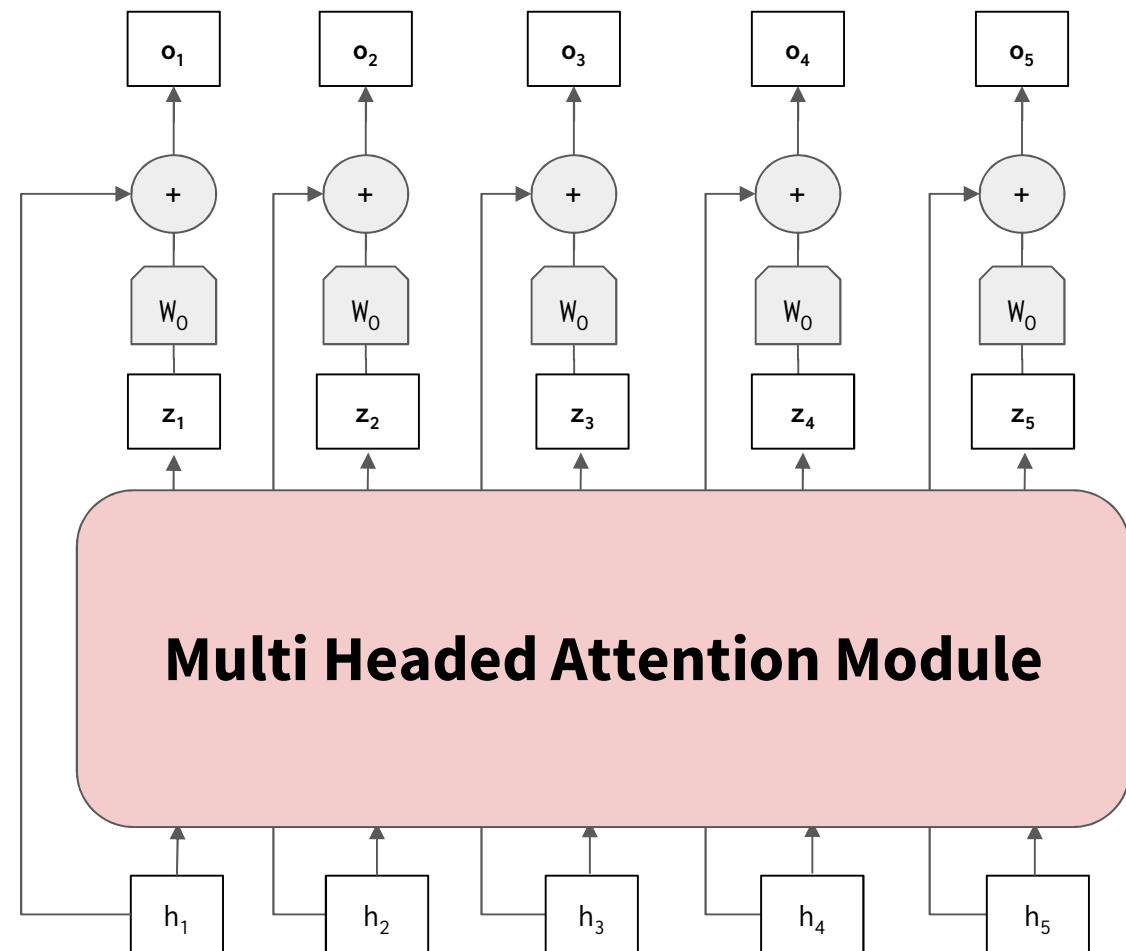
- A vector of sines and cosines of a harmonic series of frequencies
 - Every $2l$ -th component of P_t is $\sin \omega_l t$
 - Every $2l + 1$ -th component of P_t is $\cos \omega_l t$
- Never repeats
- Has the linearity property required

Positional Encoding



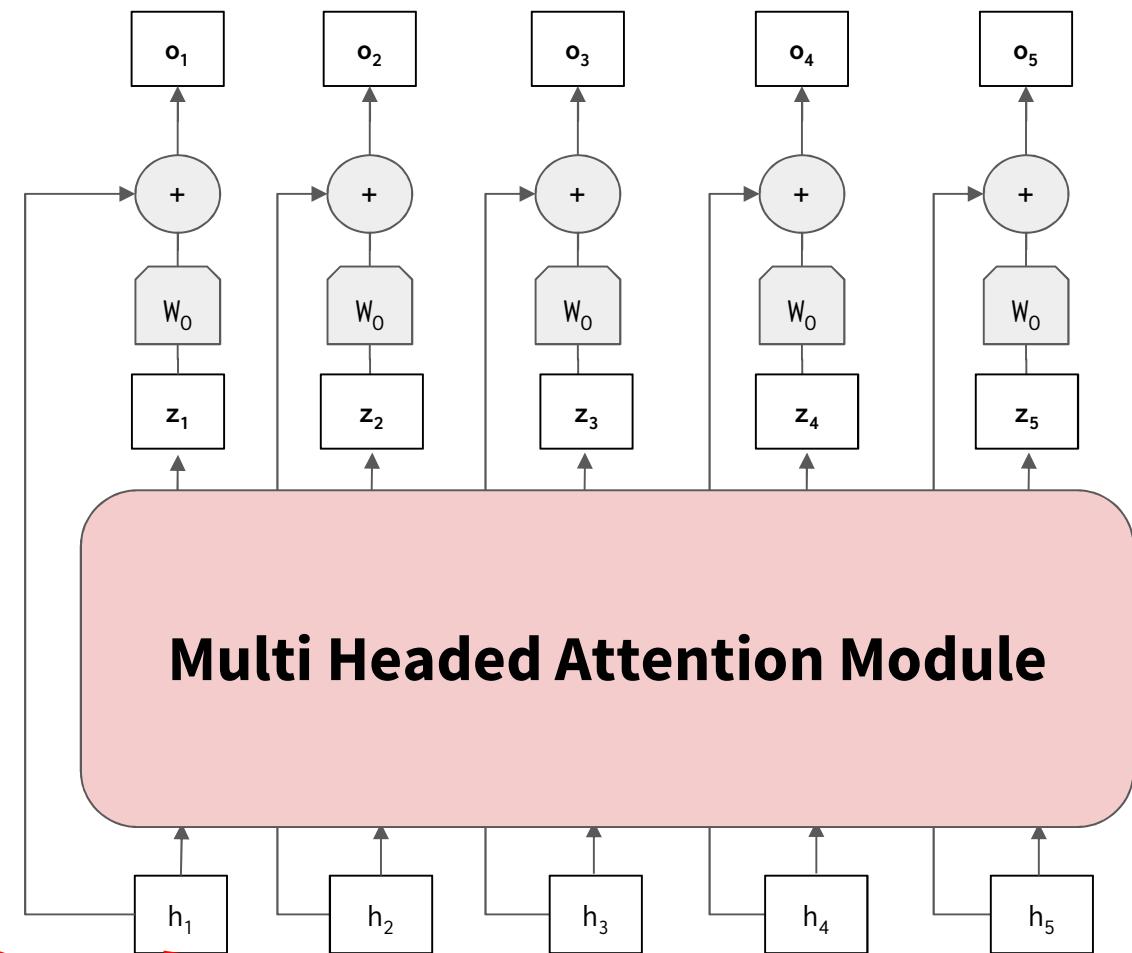
Positional Encoding simplified

```
p1 = [1 0 0 0 0]  
p2 = [0 1 0 0 0]  
...  
p5 = [0 0 0 0 1]
```



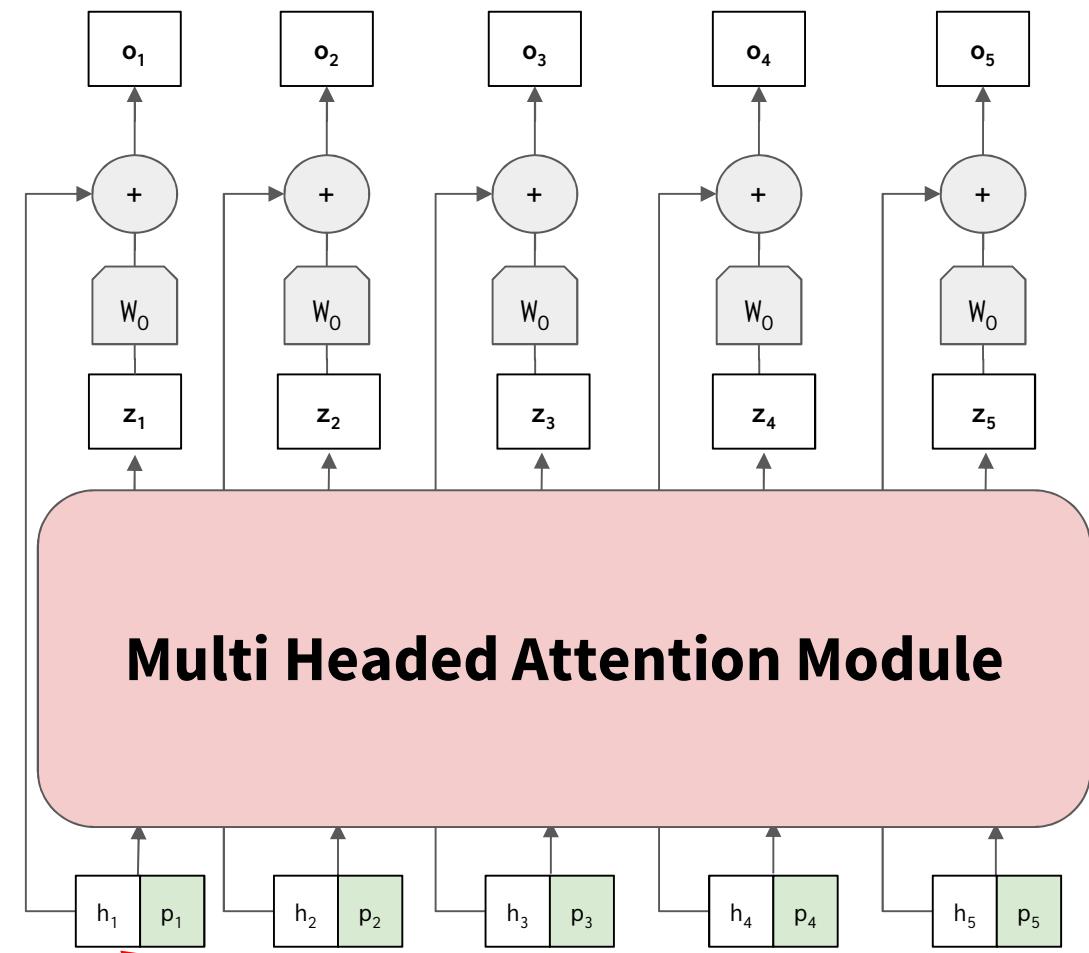
Positional Encoding simplified

```
p1 = [1 0 0 0 0]  
p2 = [0 1 0 0 0]  
...  
p5 = [0 0 0 0 1]
```



Positional Encoding simplified

```
p1 = [1 0 0 0 0]  
p2 = [0 1 0 0 0]  
...  
p5 = [0 0 0 0 1]
```



Poll 2 (@1436)

Which of the following are true about transformers?

- a. The attention module tries to calculate the “shift” in meaning of a token given all other tokens in the batch
- b. Transformers can always be run in parallel
- c. Transformer decoders can only be parallelized during training
- d. Positional encodings help parallelize the transformer encoder
- e. Queries, keys, and values are obtained by splitting the input into 3 equal segments
- f. Multiheaded attention helps transformers find different kinds of relations between the tokens
- g. During decoding, decoder outputs function as queries and keys while the values come from the encoder

Poll 2 (@1436)

Which of the following are true about transformers?

- a. **The attention module tries to calculate the “shift” in meaning of a token given all other tokens in the batch**
- b. Transformers can always be run in parallel
- c. **Transformer decoders can only be parallelized during training**
- d. Positional encodings help parallelize the transformer encoder
- e. Queries, keys, and values are obtained by splitting the input into 3 equal segments
- f. **Multiheaded attention helps transformers find different kinds of relations between the tokens**
- g. During decoding, decoder outputs function as queries and keys while the values come from the encoder

Summary (1)

- Roles of Queries, Keys, and Values
 - \mathbf{Q} pay attention to \mathbf{V} according to computation with \mathbf{K}
“Computation” is the attention function.
- **Self** versus **Cross** attention
- Transformers are Residual Machines
- **Positional Encodings:** Transformers have no notion of order - this needs to be explicitly inserted.

Summary (2)

- Transformers' biggest advantage lies in parallelizability and 'omnidirectionality'
- On smaller sequence lengths, RNN-based models might perform better with fewer parameters.

Extra Slides

Few types of energy functions

- MLP

$$e(q, k) = W_2^T(\tanh(W_1^T[q; k]))$$

- Bilinear

$$e(q, k) = (q^T)(W)(k)$$

- Scaled-Dot Product

$$e(q, k) = (q)(k^T) / (s) \# s = \text{scaling factor } (\sqrt{d_k})$$

Batching and shapes

The attention function takes in:

$$q : (B, T, d_q)$$

$$k : (B, T, d_k)$$

$$v : (B, T, d_v)$$

Energy / attention scores:

$$e : (B, T, T) \quad \# \text{Score between each pair of tokens if } e = qk^T/s$$

Output vector:

$$o : (B, T, d_v) \quad \# \text{calculated as softmax}(e)^Tv$$

Part 2

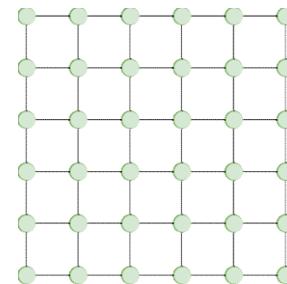
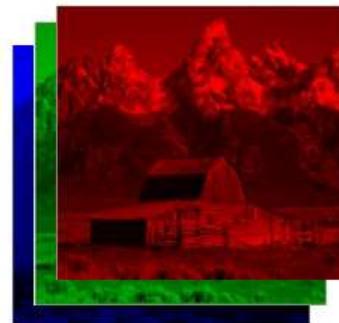
Permutation Invariance & Graph Neural Networks

Revisiting data structure we have met

Sequence data: text/speech



Grid data: image

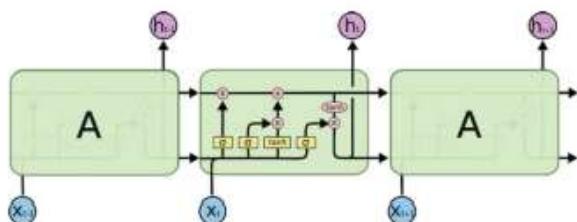


Revisiting data structure we have met

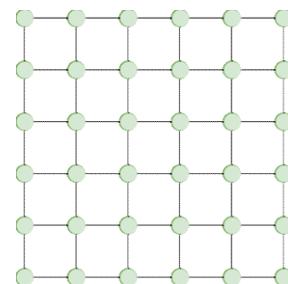
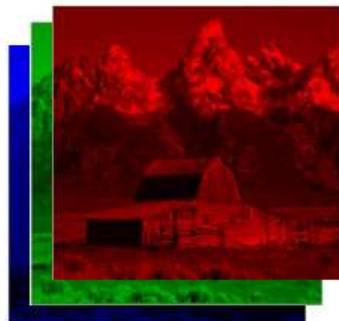
Sequence data: text/speech



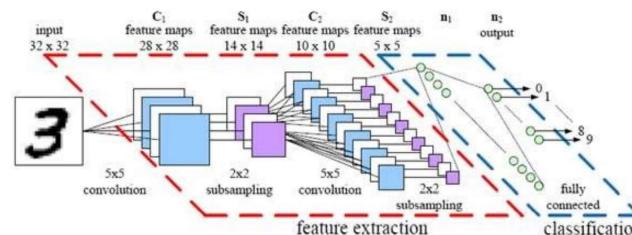
Recurrent Neural Networks



Grid data: image

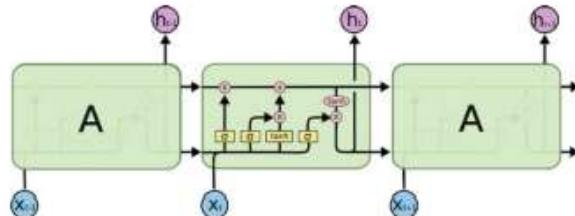


Convolution Neural Networks

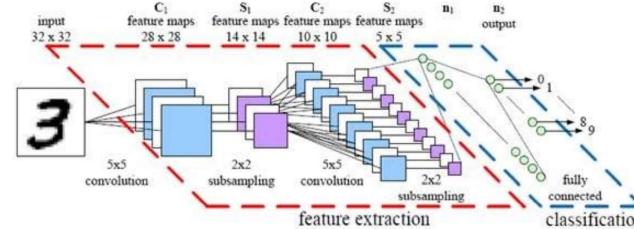
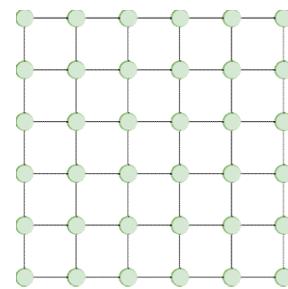
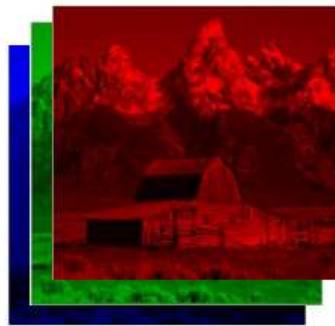


Revisiting data structure we have met

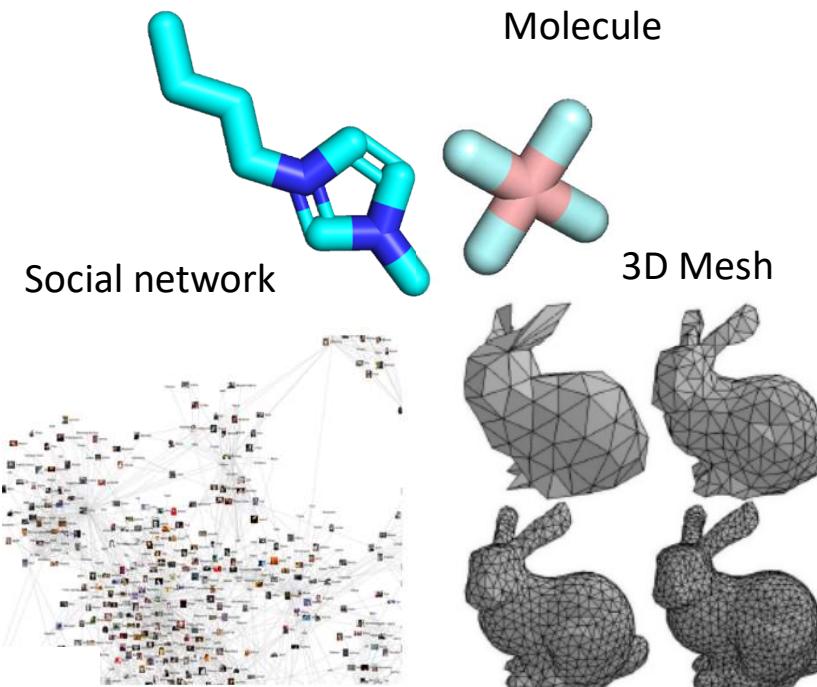
Sequence data: text/speech



Grid data: image

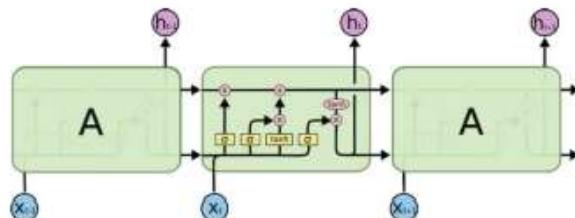


Unstructured data: Molecule, Social Networks
,3D Mesh

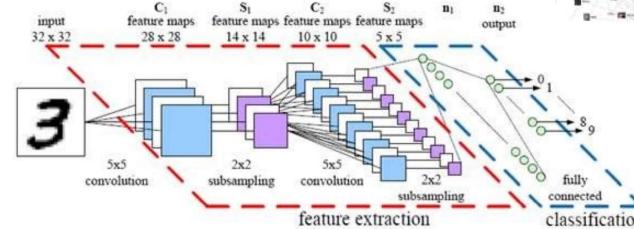
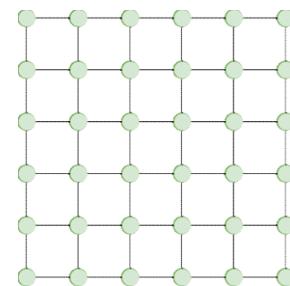
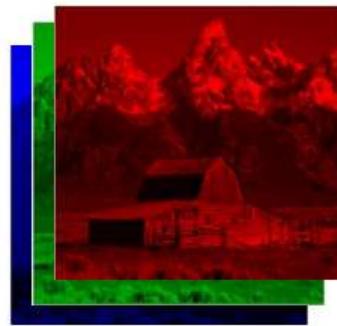


Revisiting data structure we have met

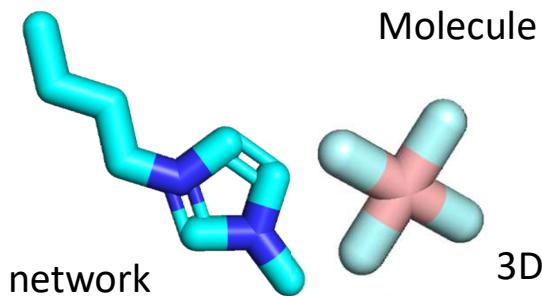
Sequence data: text/speech



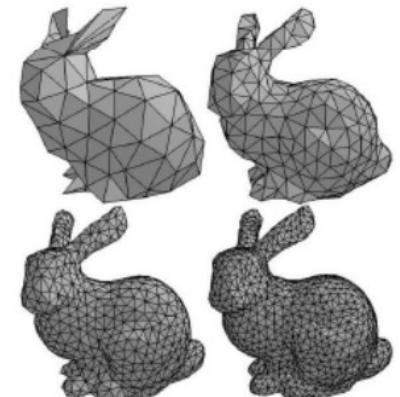
Grid data: image



Unstructured data: Molecule, Social Networks
, 3D Mesh



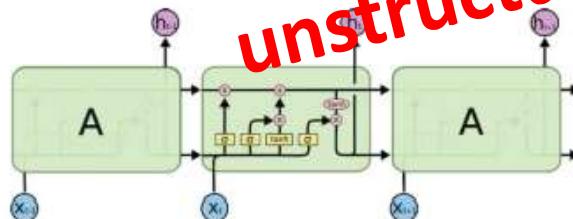
Social network



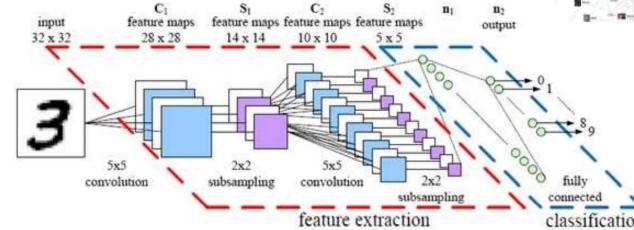
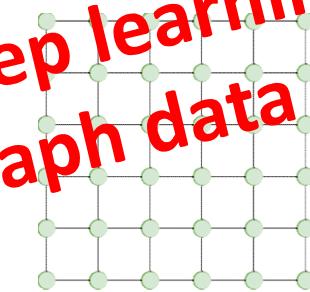
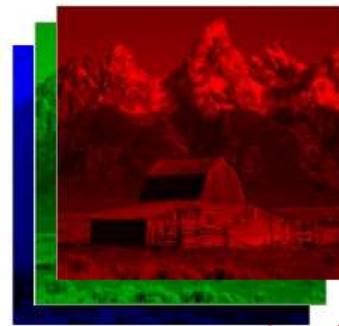
?

Revisiting data structure we have met

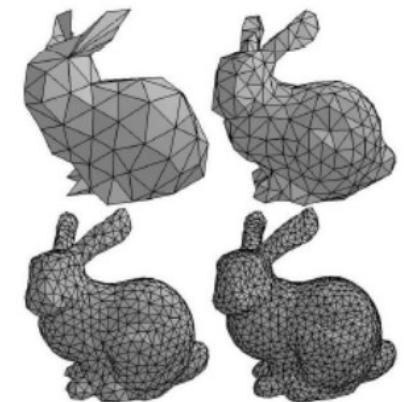
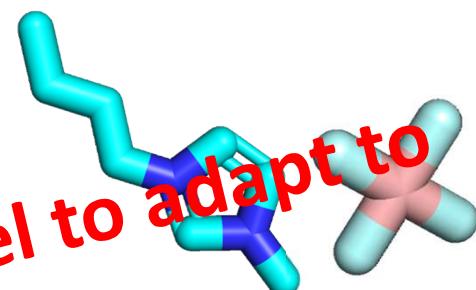
Sequence data: text/speech



Grid data: image



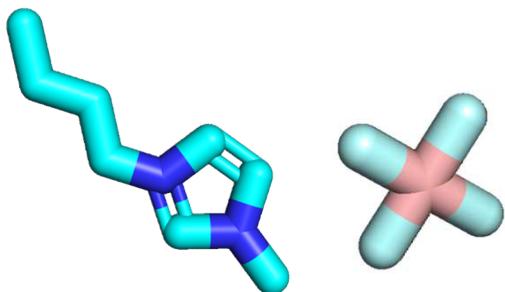
Unstructured data: Molecule, Social Networks
,3D Mesh



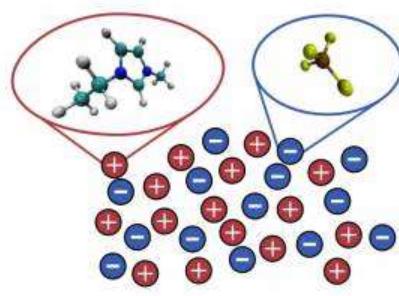
?

Problem Setup: ionic liquid for CO₂ capturing

Ionic liquid molecules

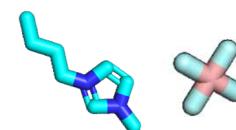


Carbon dioxide (CO₂)



dataset

data	label
IL molecule structure	Solubility(volume percentage)



0.56,
0.119.....

We want to use deep learning model to predict the solubility of ionic liquid based on these molecule data!

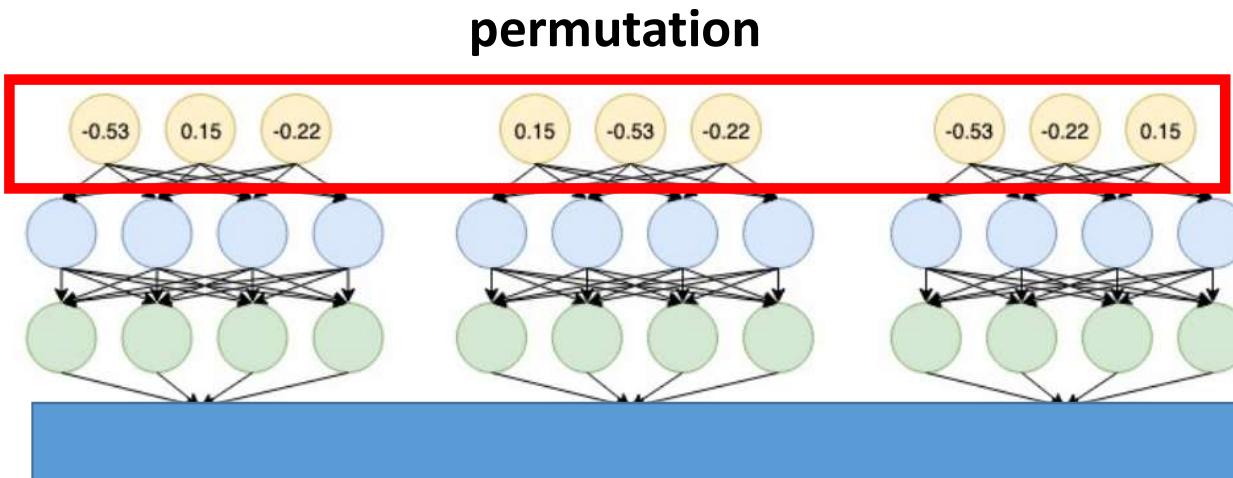
Can we use MLP, CNN or RNN to do this job?

Invariance

Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

Example 1 (Multi-layer Perceptron + Sequence Data):



$f()$: MLP

$g()$: permutation of input order

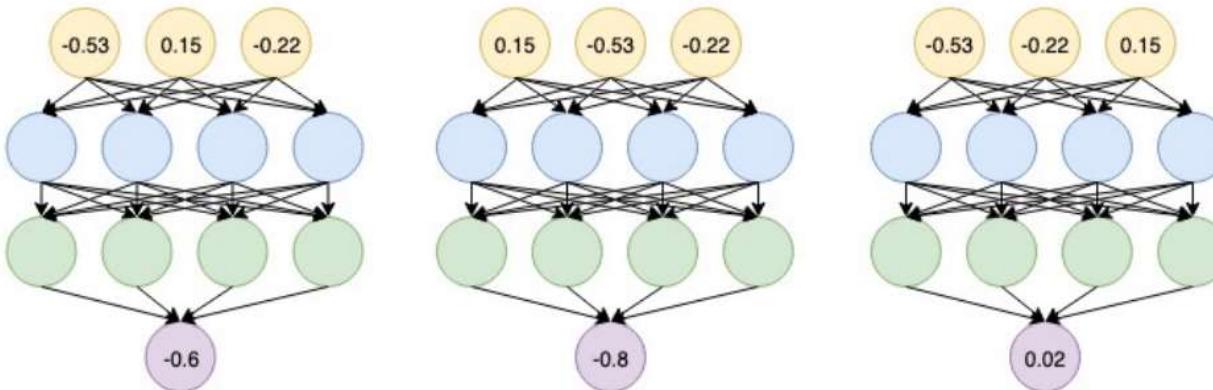
Is $f()$ invariant to $g()$?

Invariance

Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

Example 1 (Multi-layer Perceptron + Sequence Data):



$f()$: MLP

$g()$: permutation of input order

Is $f()$ invariant to $g()$?

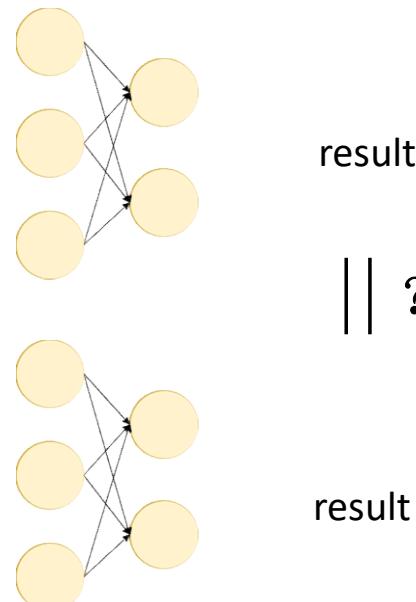
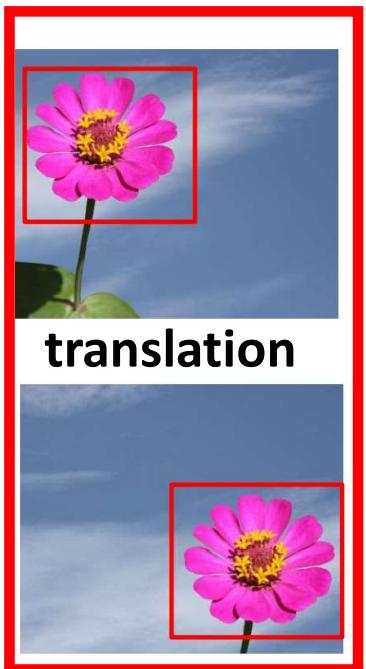
No!

Invariance

Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

Example 2 (MLP + image data):



$f()$: MLP
 $g()$: translation

|| ? Is $f()$ invariant to $g()$?

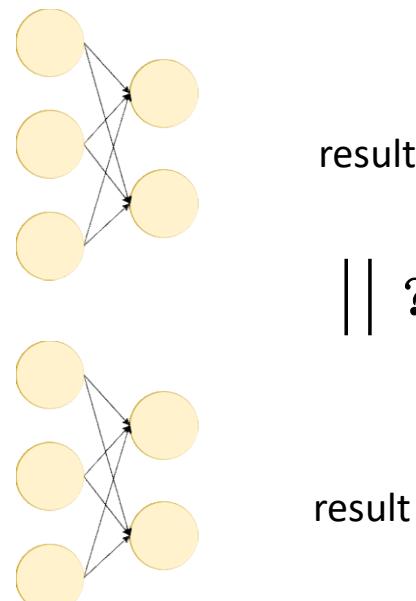
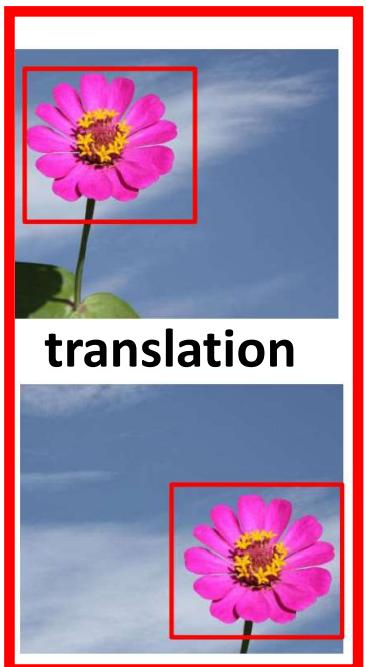
result

Invariance

Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

Example 2 (MLP + image data):



$f()$: MLP
 $g()$: translation

|| ? Is $f()$ invariant to $g()$?

result

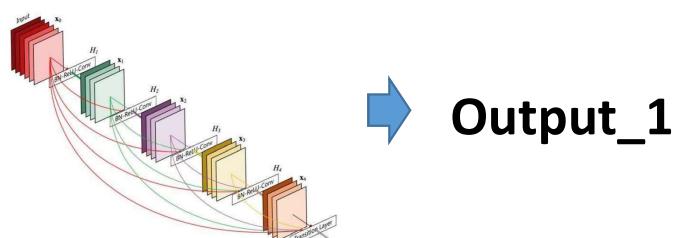
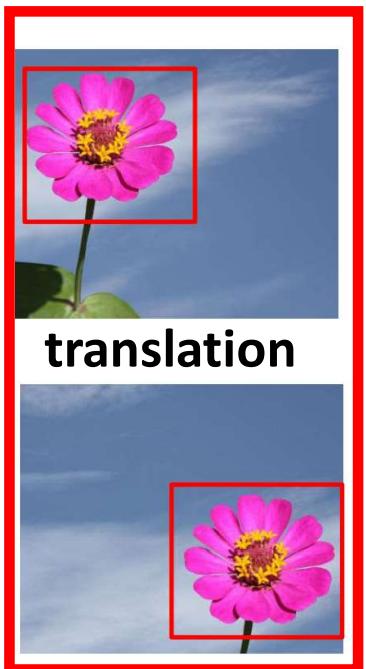
No!

Invariance

Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

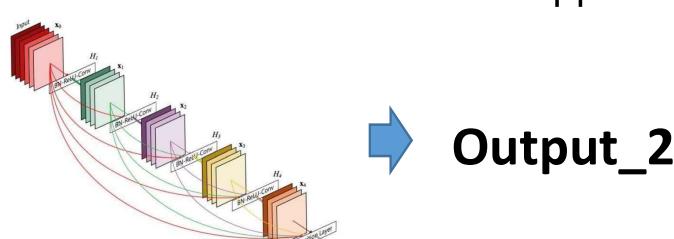
Example 3 (CNN + image data):



→ Output_1

$f()$: CNN
 $g()$:translation

|| ? Is $f()$ invariant to $g()$?



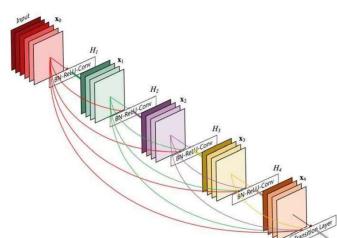
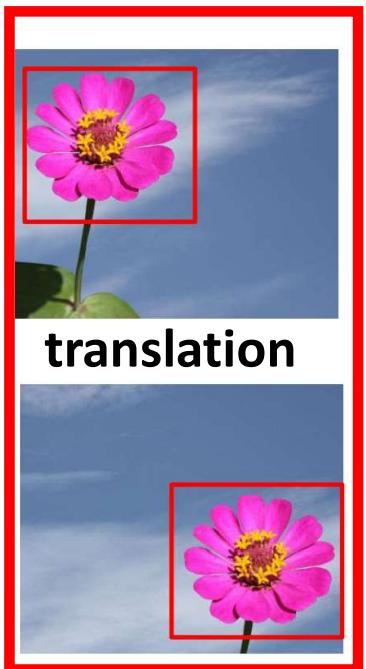
→ Output_2

Invariance

Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

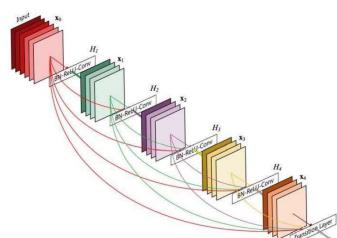
Example 3 (CNN + image data):



→ Output_1

$f()$: CNN
 $g()$: translation

|| ? Is $f()$ invariant to $g()$?



→ Output_2

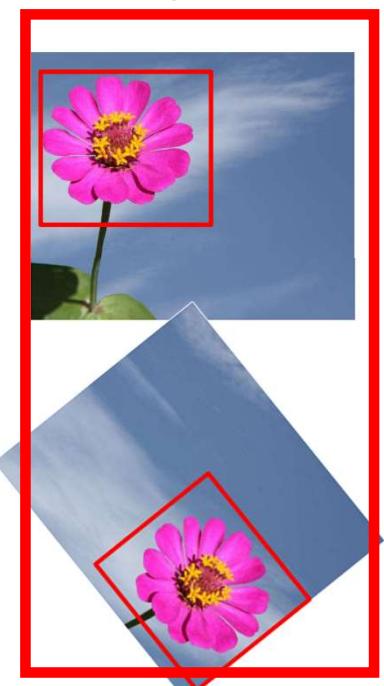
Yes !

Invariance

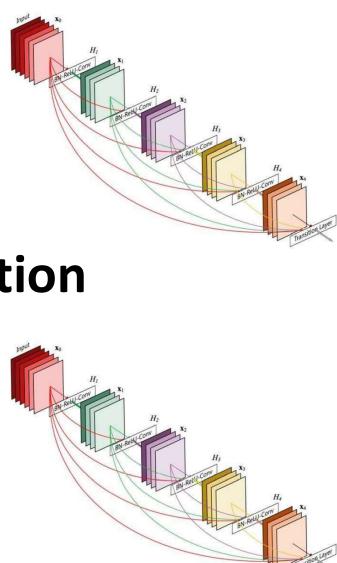
Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

Example 4 (CNN + image data):



rotation



Output_1

$f()$: CNN
 $g()$: rotation

|| ? Is $f()$ invariant to $g()$?

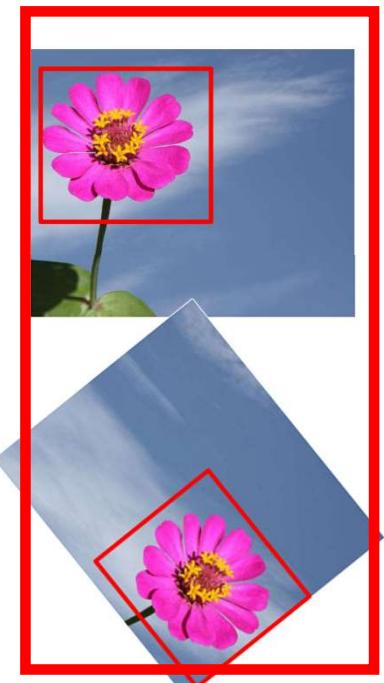
Output_2

Invariance

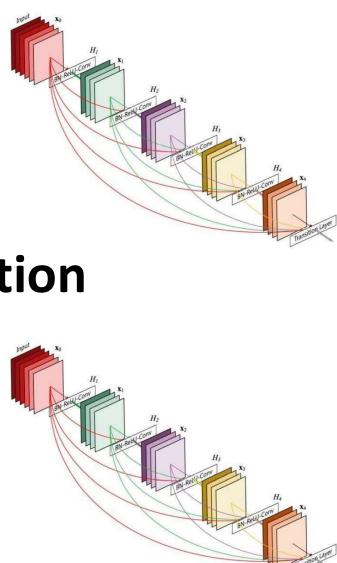
Let's define a mapping from X to Y : $Y = f(X)$

- If there exist a mapping function of $g()$, such that $f(g(X)) = f(X) = Y$
- we say $f()$ is invariant to $g()$

Example 4 (CNN + image data):



rotation



Output_1

$f()$: CNN
 $g()$: rotation

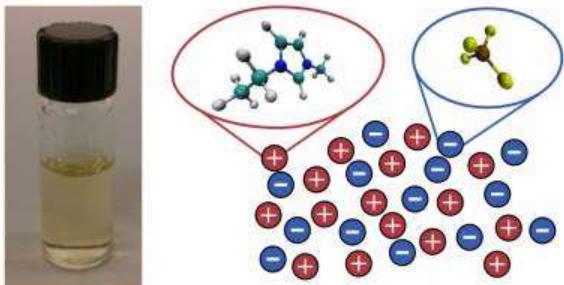
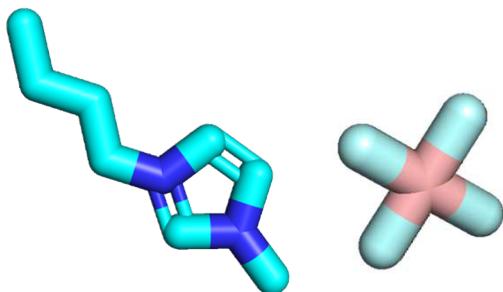
|| ? Is $f()$ invariant to $g()$?

Output_2

No !

Problem Setup: ionic liquid for CO₂ capturing

Ionic liquid molecules

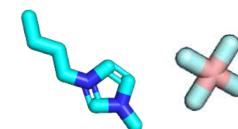


Carbon dioxide (CO₂)



dataset

data	label
IL molecule structure	Solubility(volume percentage)



0.56,
0.119.....

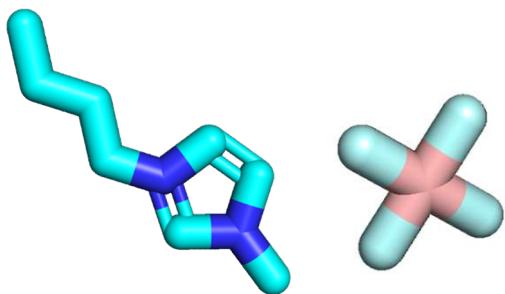
We want to use deep learning model to predict the solubility of ionic liquid based on these data!

- (1) How many ways can you come up with?
- (2) Can we use MLP, CNN or RNN to do this job?
- (3) What are the desired properties of the model we used for molecule?
 - permutation invariant?
 - translation invariant?
 - rotation invariant?

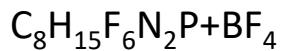
Possible solution for molecule properties prediction

Use MLP to solve the problem

Ionic liquid molecules

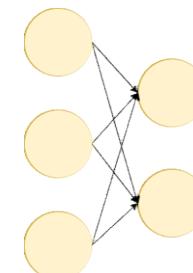


Sequence data: Empirical Formula



C8H15F6N2P+BF4	
C	8
H	15
N	2
P	1
F	6
B	4

Feed sequence to MLP

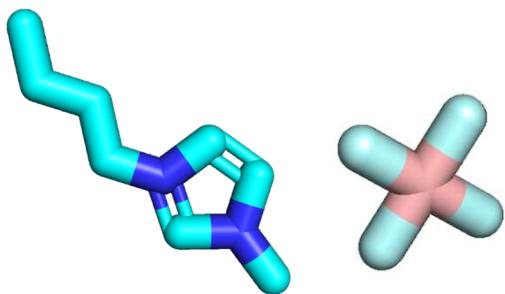


Result

Possible solution for molecule properties prediction

Use MLP to solve the problem

Ionic liquid molecules



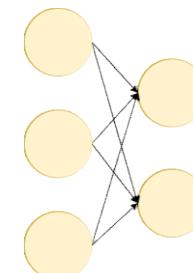
Sequence data: Empirical Formula

$\text{C}_8\text{H}_{15}\text{F}_6\text{N}_2\text{P}+\text{BF}_4$



C	8
H	15
N	2
P	
F	6
B	
	1
	4
	1

Feed sequence to MLP

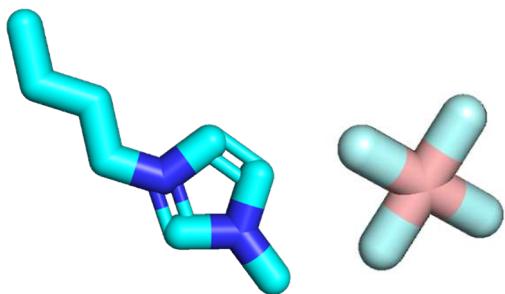


Result

Possible solution for molecule properties prediction

Use MLP to solve the problem

Ionic liquid molecules

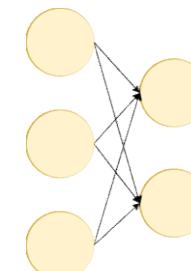


Sequence data: Empirical Formula

$\text{C}_8\text{H}_{15}\text{N}_2\text{F}_6\text{P}+\text{BF}_4$

C	8
H	15
N	2
P	
F	6
B	1 4

Feed sequence to MLP



Result?

MLP is not permutation invariant!
Not good for molecule with unstructured data

Possible solution for molecule properties prediction

Another possible solution

Ionic liquid molecules

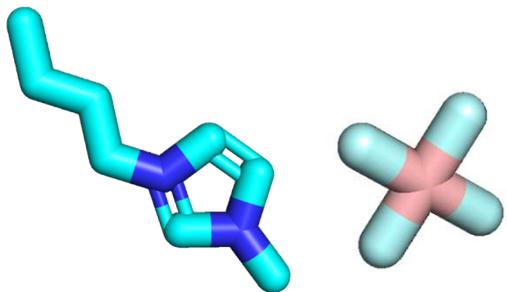
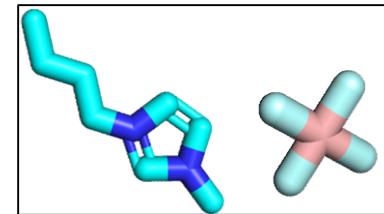
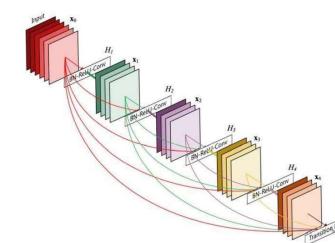


Image data: image for molecule structure



Feed image to CNN



Result?

Possible solution for molecule properties prediction

Another possible solution

Ionic liquid molecules

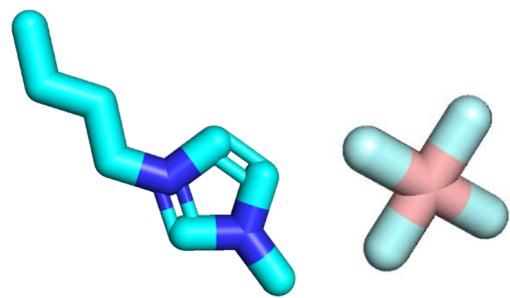
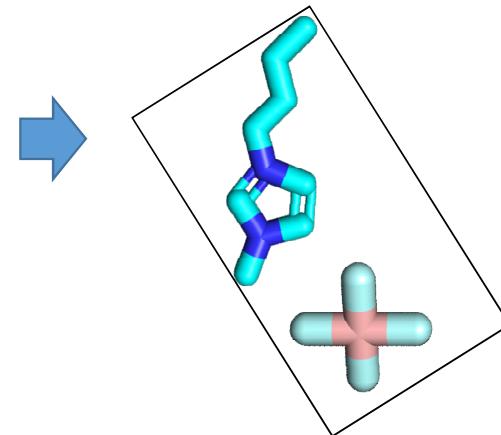
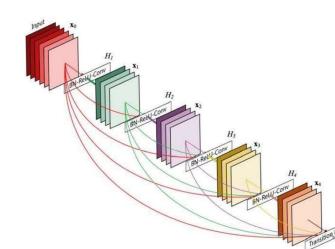


Image data: image for molecule structure



Rotate image

Feed image to CNN



Result?

Possible solution for molecule properties prediction

Another possible solution

Ionic liquid molecules

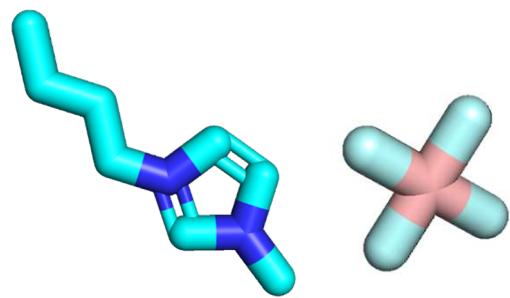
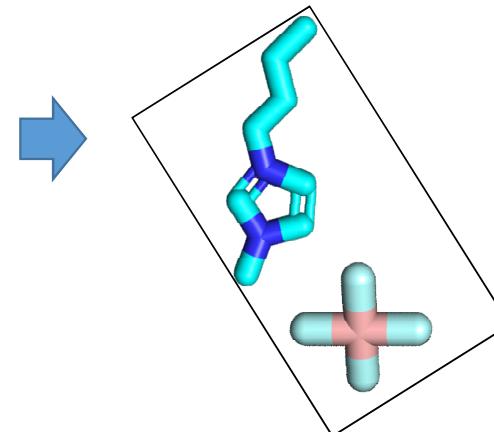
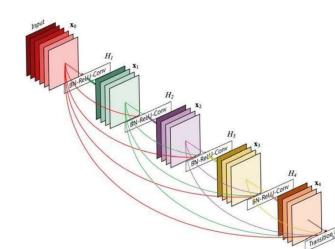


Image data: image for molecule structure



Rotate image

Feed image to CNN



Result?

CNN is not rotational invariant!
Not good for molecule with unstructured data

Story so far

- We want to use deep learning to do prediction for unstructured graph like data such as molecule
 - MLP is not permutation nor translation invariant
 - CNN is not permutation invariant nor rotation invariant
- Those draw backs make MLP and CNN not a good candidates for processing molecule data
- We want a model with a strong invariant property

Graph: Definition

A graph is defined as a pair $G = (V, E)$,

- where **V** is a set whose elements are called **nodes**(vertices),
- and **E** is a set of paired vertices, whose elements are called **edges**.

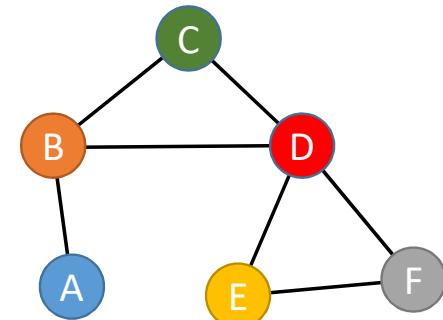
Example:

$$G = (V, E)$$

$$V = \{A, B, C, D, E, F\}$$

$$E = \{(A, B), (B, C), (C, D), (B, D), (C, D), (D, E), (D, F), (E, F)\}$$

Undirected Graph

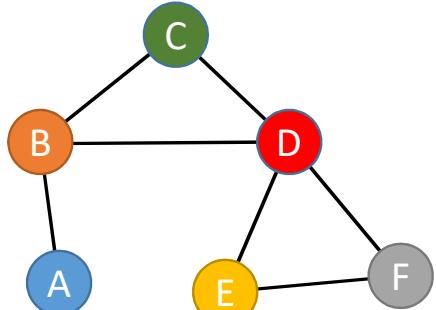


Graph: Matrix representation for graph

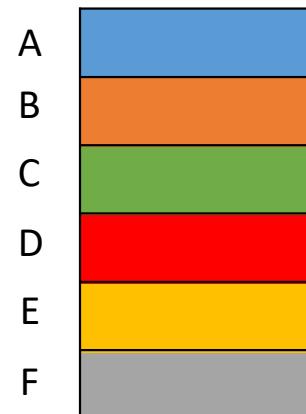
A graph is defined as a pair $G = (V, E)$,

- where **V** is a set whose elements are called **nodes**(vertices),
- and **E** is a set of paired vertices, whose elements are called **edges**.

Undirected Graph



Node information Matrix ($N \times F$)



Node information

Adjacency Matrix ($N \times N$)
6 6

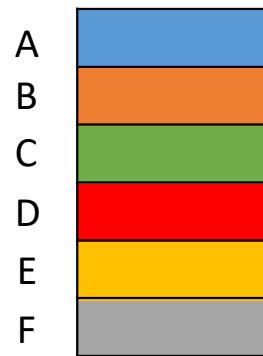
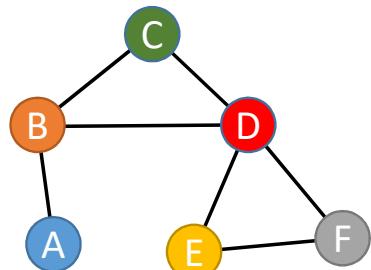
	A	B	C	D	E	F
A		1	1	1	0	0
B	1		1	1	0	0
C	1	1		1	0	0
D	1	1	1		1	1
E	0	0	0	1		1
F	0	0	0	1	1	

Connectivity information

Graph: No canonical order of nodes

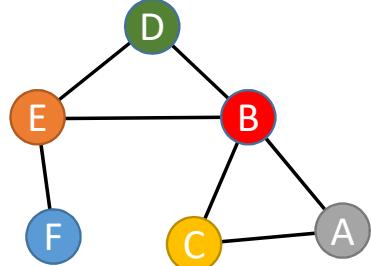
Graph do not have canonical order of the nodes!

Order plan 1



	A	B	C	D	E	F
A		blue				
B	blue		blue	blue		
C		blue		blue	blue	
D			blue		blue	blue
E				blue		blue
F					blue	

Order plan 2

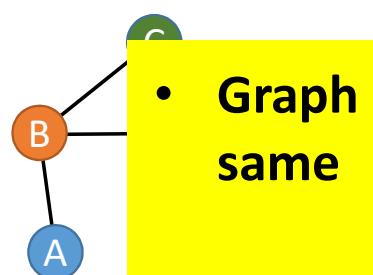


	A	B	C	D	E	F
A		blue	blue			
B	blue			blue	blue	
C		blue				
D					blue	
E					blue	blue
F						blue

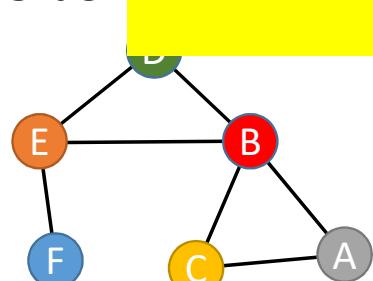
Graph: No canonical order of nodes

Graph do not have canonical order of the nodes!

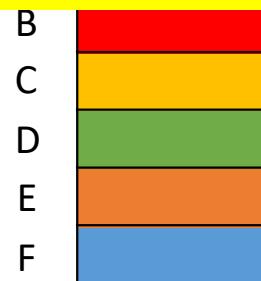
Order plan 1



Order



A	B	C	D	E	F
A	1	0	0	0	0
B	0	1	0	0	0
C	0	0	1	0	0
D	0	0	0	1	0
E	0	0	0	0	1
F	0	0	0	0	0

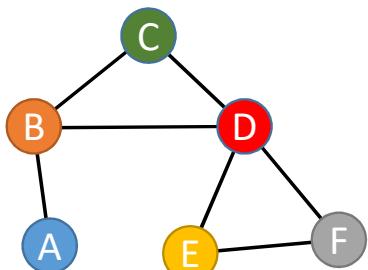


B	C	D	E	F
B	1	0	0	0
C	0	1	0	0
D	0	0	1	0
E	0	0	0	1
F	0	0	0	0

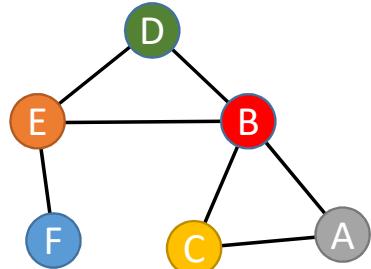
Graph: Permutation invariance

Desired properties for GNN

Order plan 1



Order plan 2

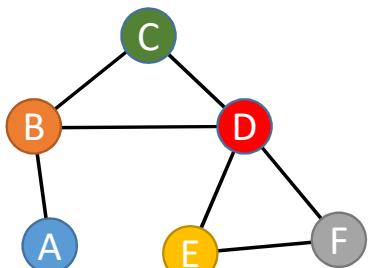


- Consider we use deep learning model to learn a function $f()$ that map a graph $G = (V, E)$ to a vector R^d
- We can write $f: f(V, E) \rightarrow R^d$ (V is node feature matrix, E is represent by adjacency matrix)
- What we want? $f(V_1, E_1) = f(V_2, E_2)$

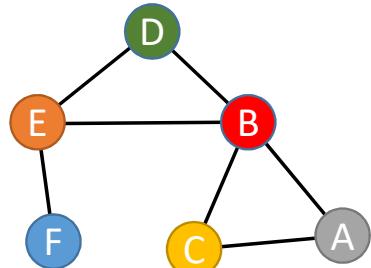
Graph: Permutation invariance

Desired properties for GNN

Order plan i

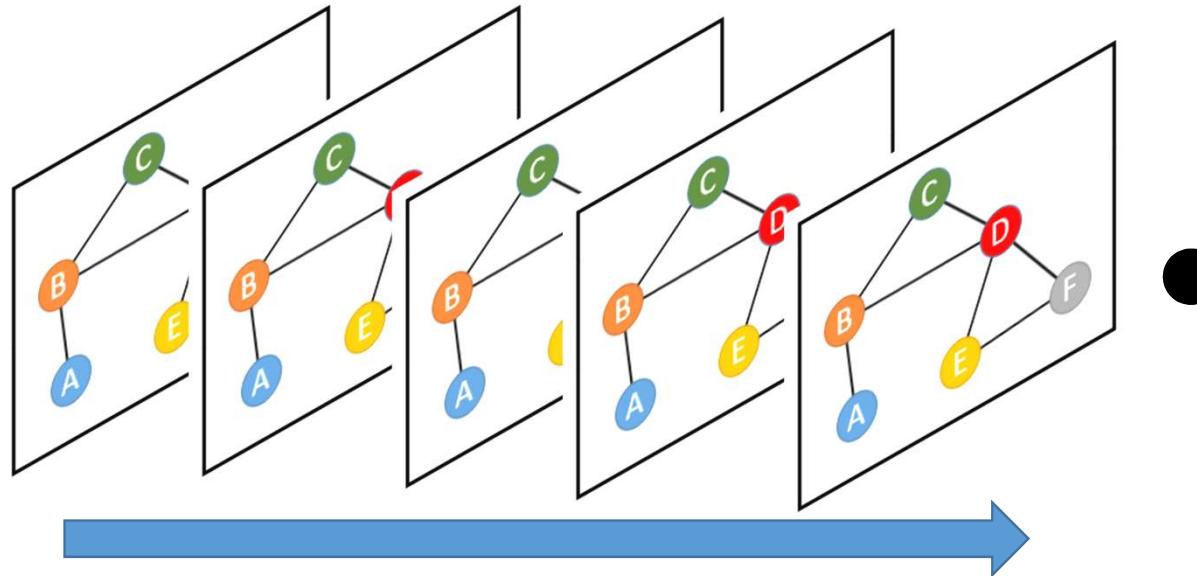


Order plan j



- Consider we use deep learning model to learn a function $f()$ that map a graph $G = (V, E)$ to a vector R^d
- We can write $f: f(V, E) \rightarrow R^d$ (V is node feature matrix, E is represent by adjacency matrix)
- For a graph with m node , there are $m!$ order plans
- Then if $f(V_i, E_i) = f(V_j, E_j)$ works for every pair of order plans i, j
- We formally define that function f is permutation invariant for the graph represented with V , E matrix₂₉

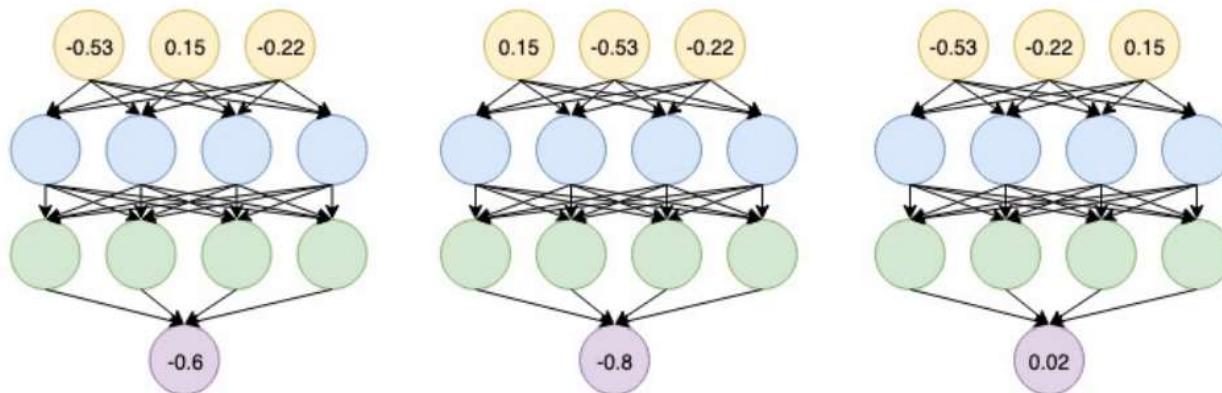
Graph: Permutation invariance



GNN consist of a series of permutation invariant layers!

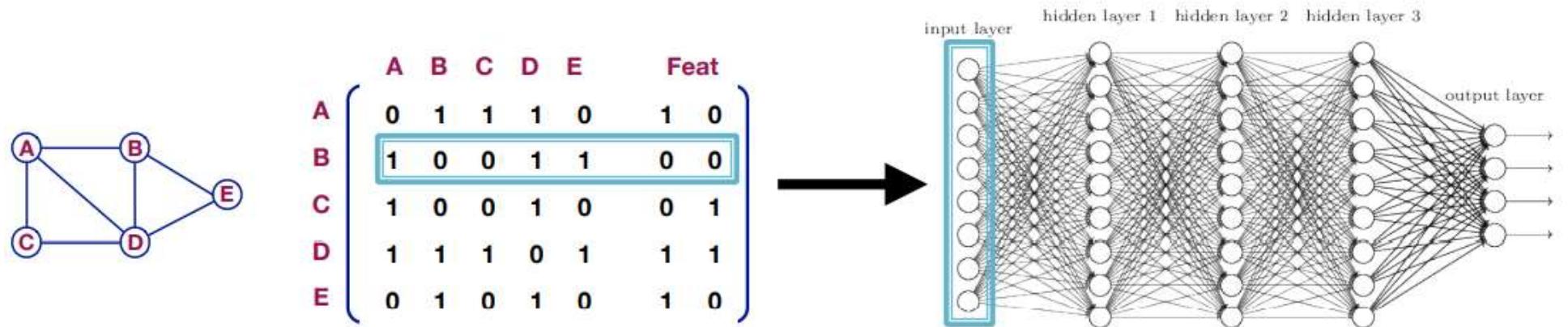
Graph: Revisit MLP

We have discussed that MLP is not permutation invariant



Graph: Revisit MLP

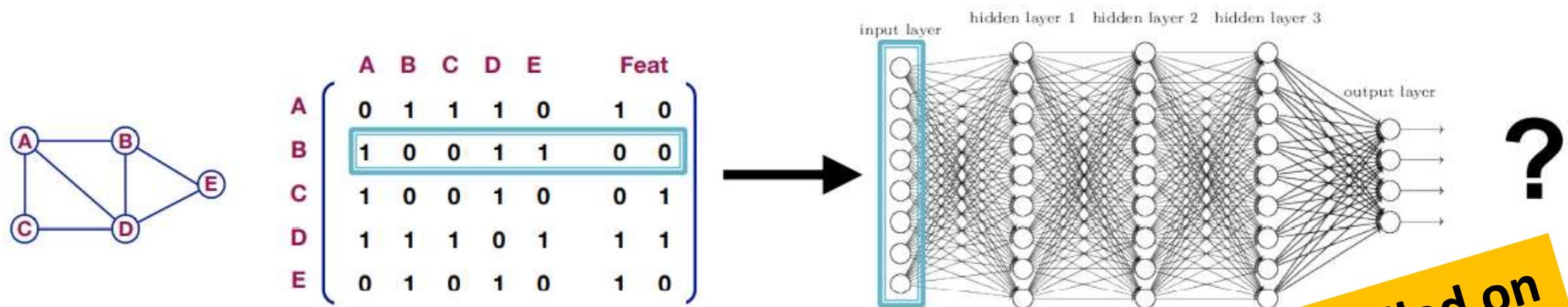
We have discussed that MLP is not permutation invariant
This also hold when we use graph data as input



Change the order plan will change the sequence order
and thus produce a different result!

Graph: Revisit MLP

We have discussed that MLP is not permutation invariant
This also hold when we use graph data as input



Change the order plan will change the sequence
and thus produce a different result!

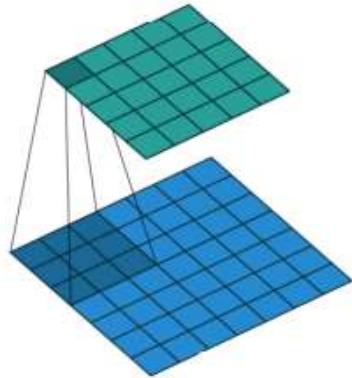
Naïve MLP failed on
graph data!!

Story so far

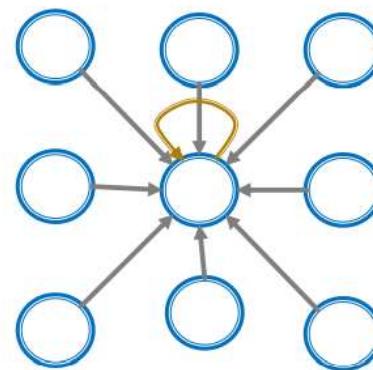
- Graph can be represented by jointly use a feature matrix and a adjacency matrix
- Graph representation does not have canonical order of node
- Permutation invariant is what we want to design a GNN

A single layer of GNN: Graph Convolution

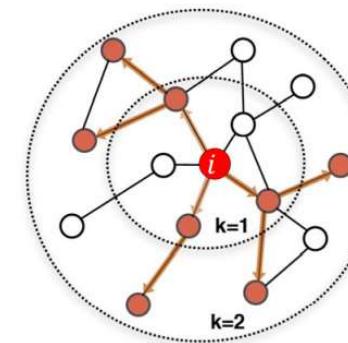
Key idea: Node's neighborhood defines a computation graph



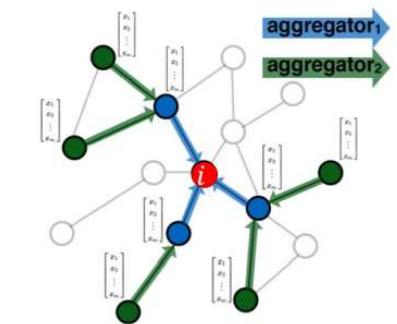
CNN: pixel convolution



CNN: pixel convolution



GNN: graph convolution

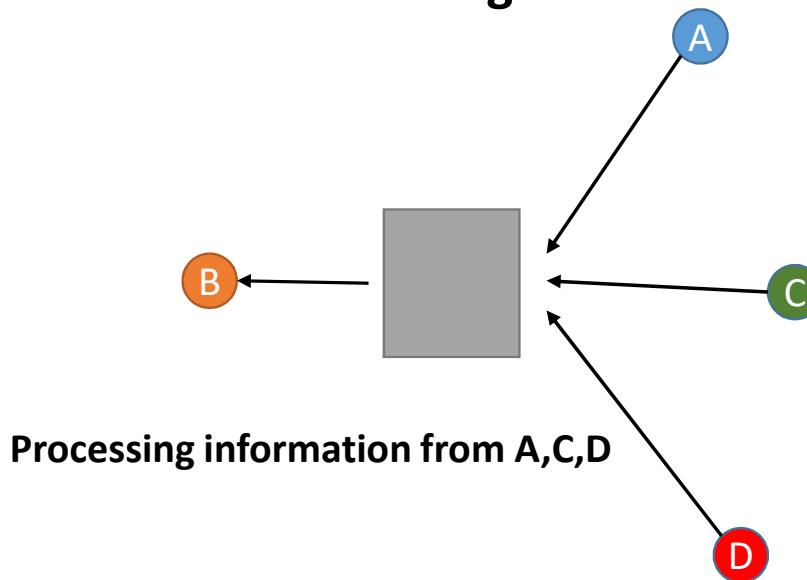
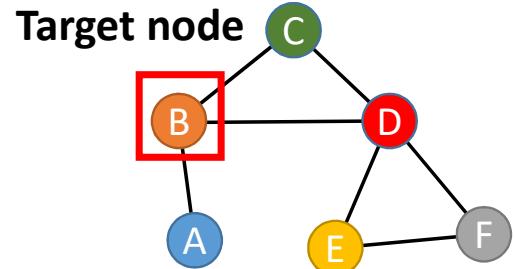


- Learning a node feature by propagating and aggregating neighbor information!
- Node embedding can be defined by local network neighborhoods!

A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods

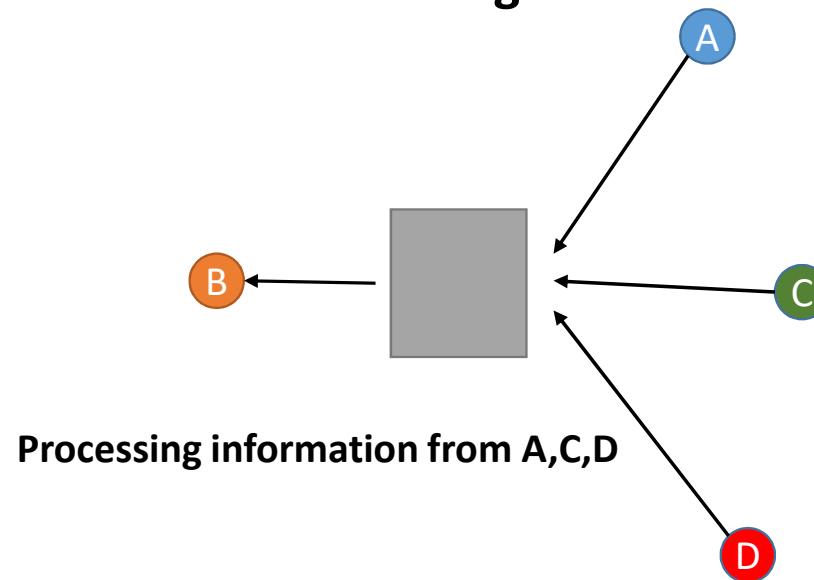
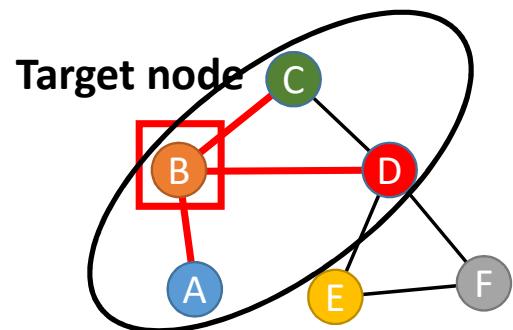
Considering 1 step of feature aggregation
of the nearest neighbor



A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods

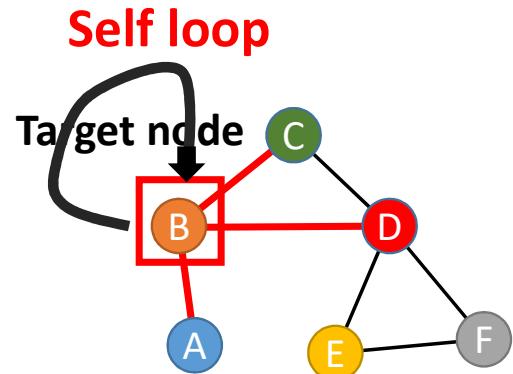
Considering 1 step of feature aggregation
of the nearest neighbor



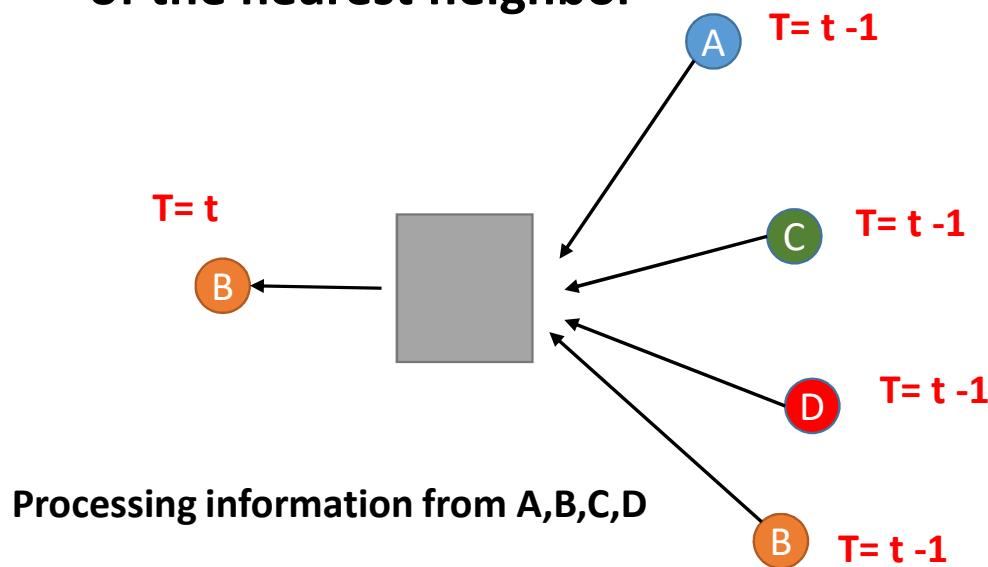
Now B have the information from it's first nearest neighbors

A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods



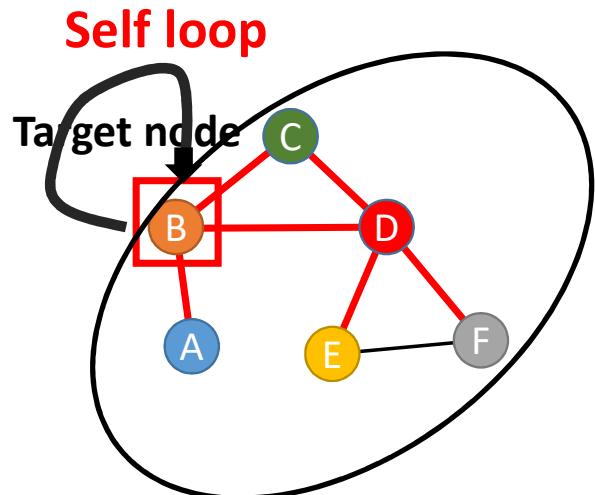
Considering 1 step of feature aggregation
of the nearest neighbor



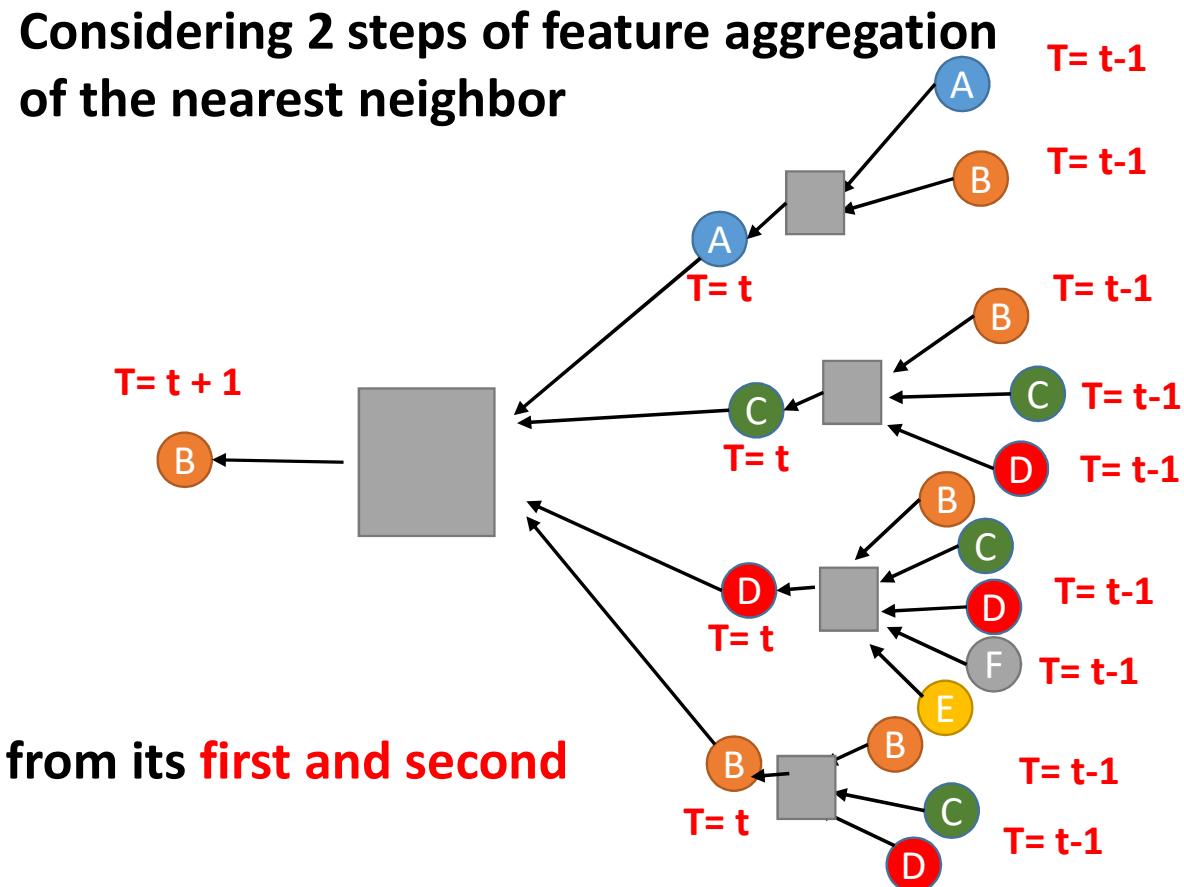
Also we don't want to lose information from B itself

A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods



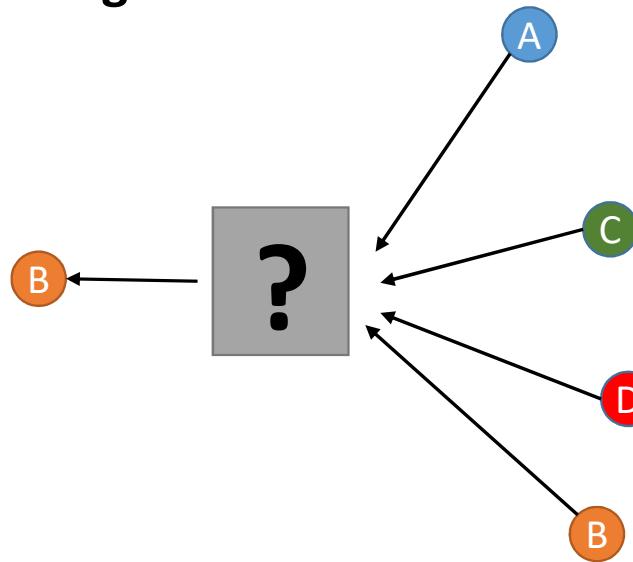
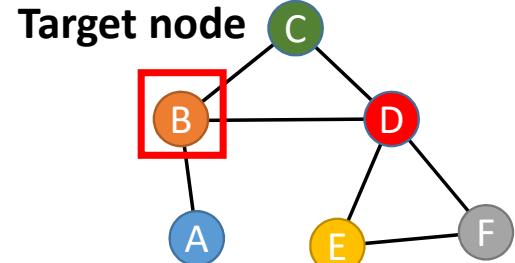
Now B have the information from its **first and second** nearest neighbors



A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods

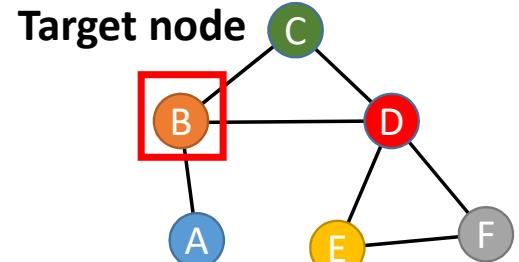
How to process and mix the information from neighbor?



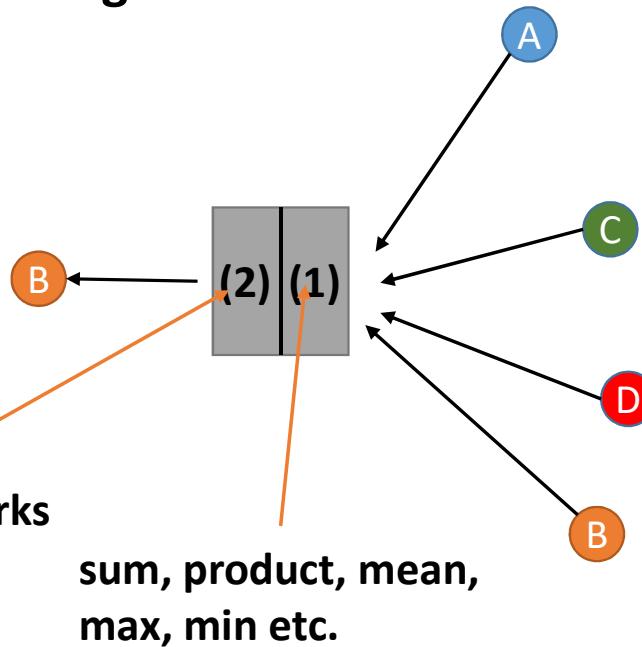
A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods

How to process and mix the information from neighbor?



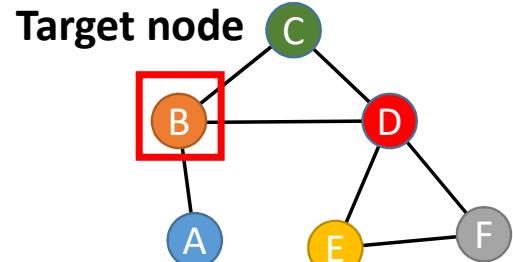
Apply Neural Networks



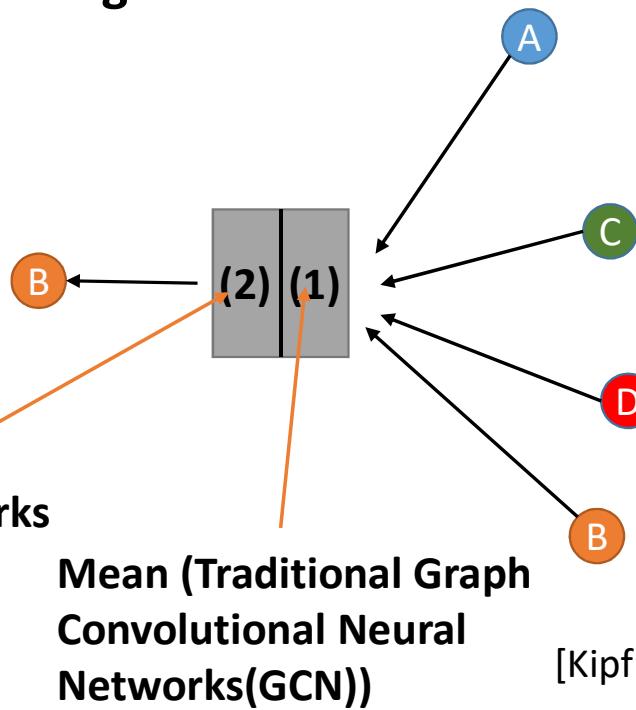
A single layer of GNN: Graph Convolution

Key idea: Generate node embedding based on local network neighborhoods

How to process and mix the information from neighbor?



Apply Neural Networks

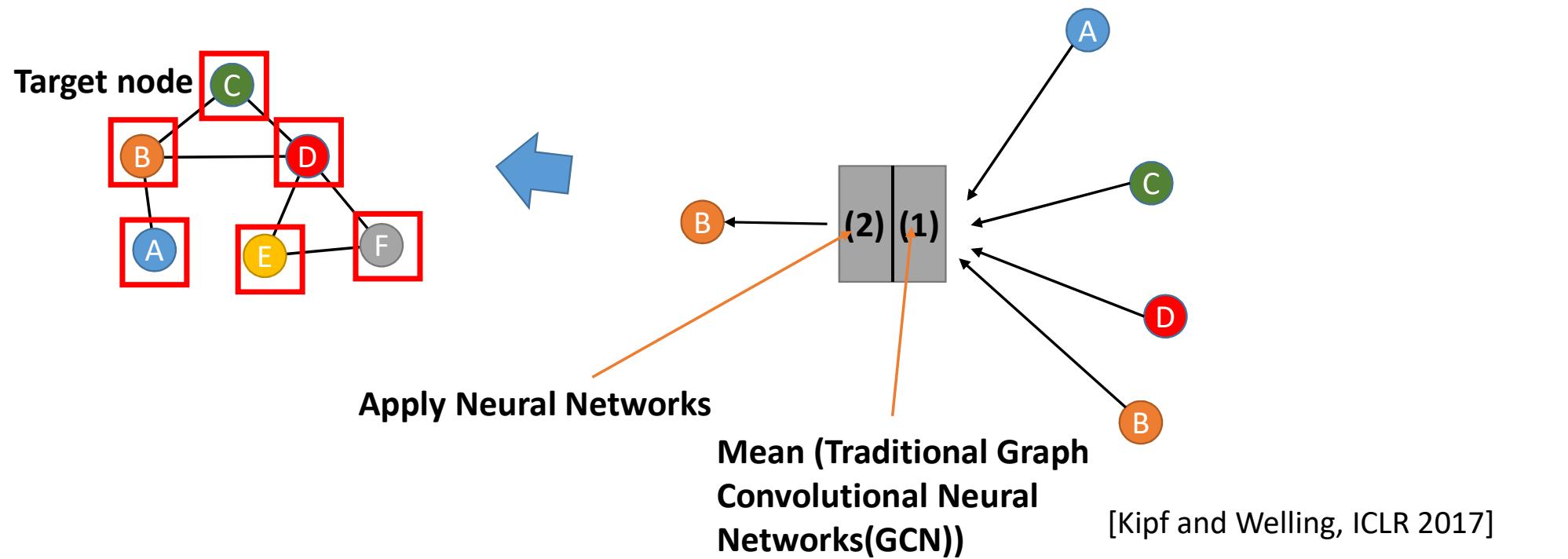


[Kipf and Welling, ICLR 2017]

A single layer of GNN: Graph Convolution

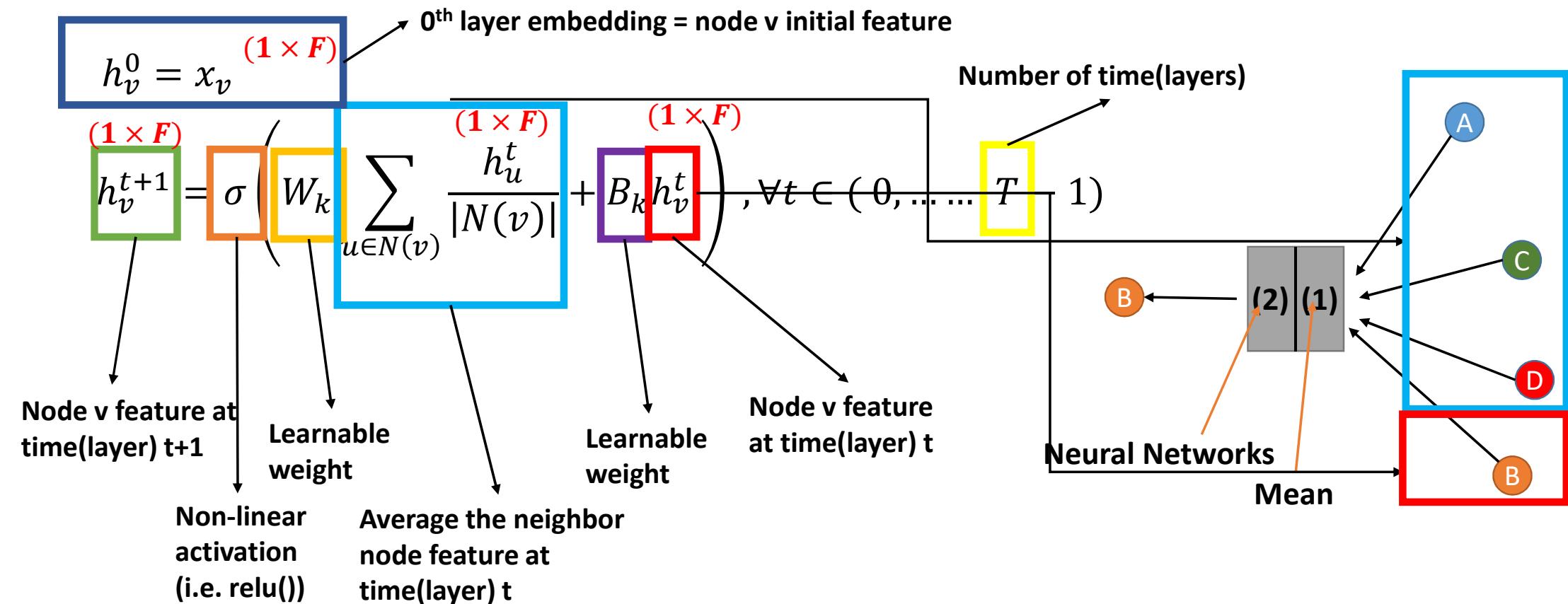
Key idea: Generate node embedding based on local network neighborhoods

During a single Graph Convolution layer, we apply the feature aggregation to every node in the graph at the same time (T)



A single layer of GNN: Graph Convolution-Forward

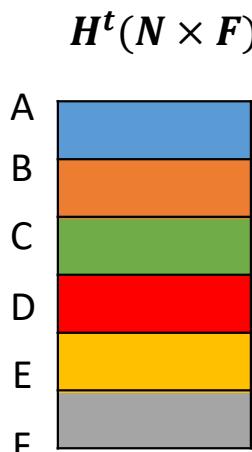
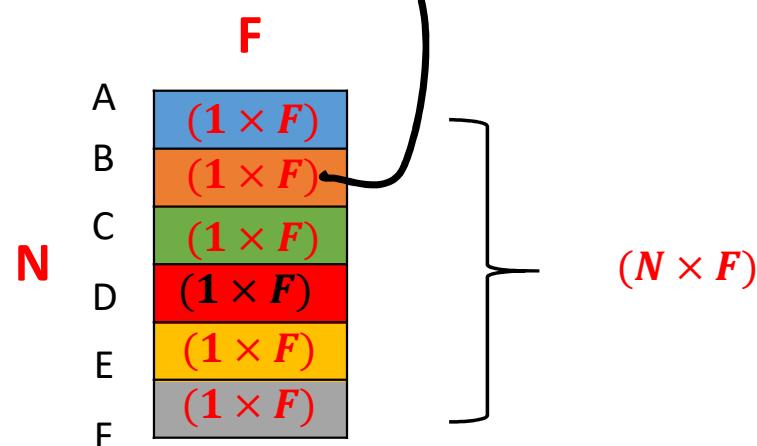
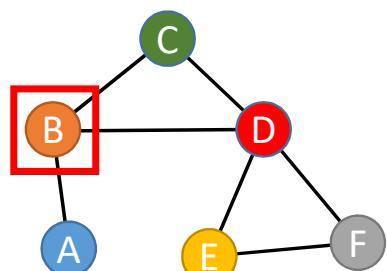
Math for a single layer of graph convolution



A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

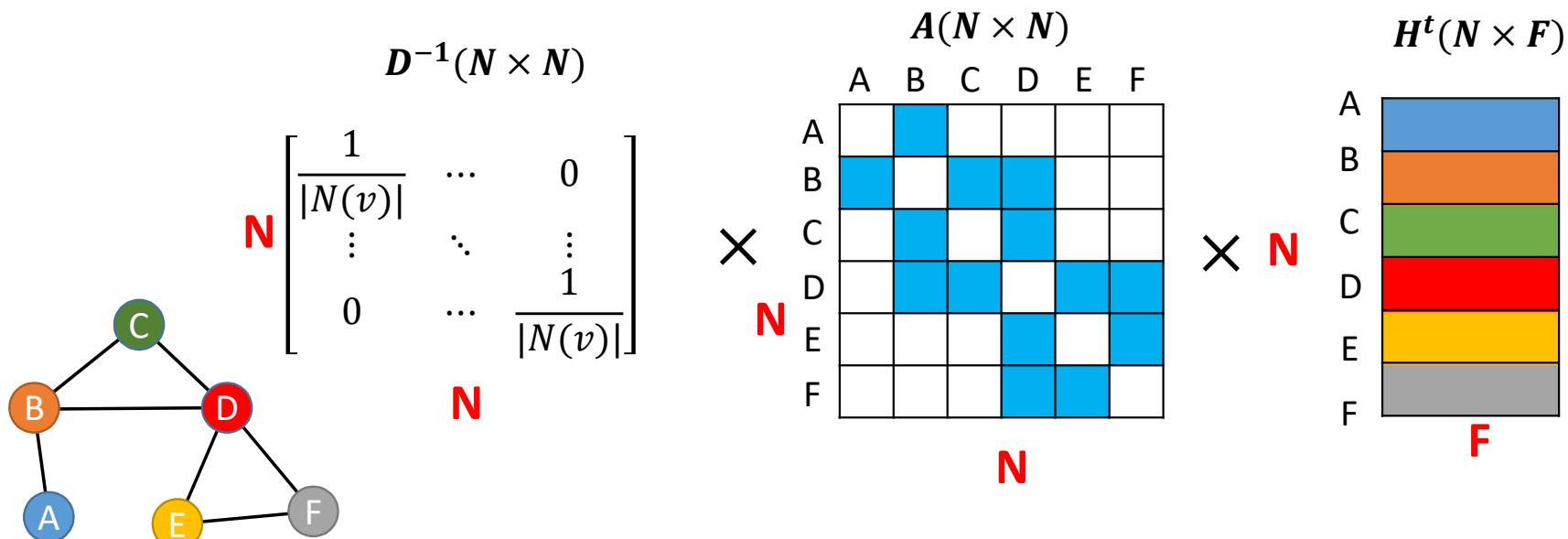


We stack multiple $h_v^t (1 \times F)$ together into $H^t (N \times F)$

A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

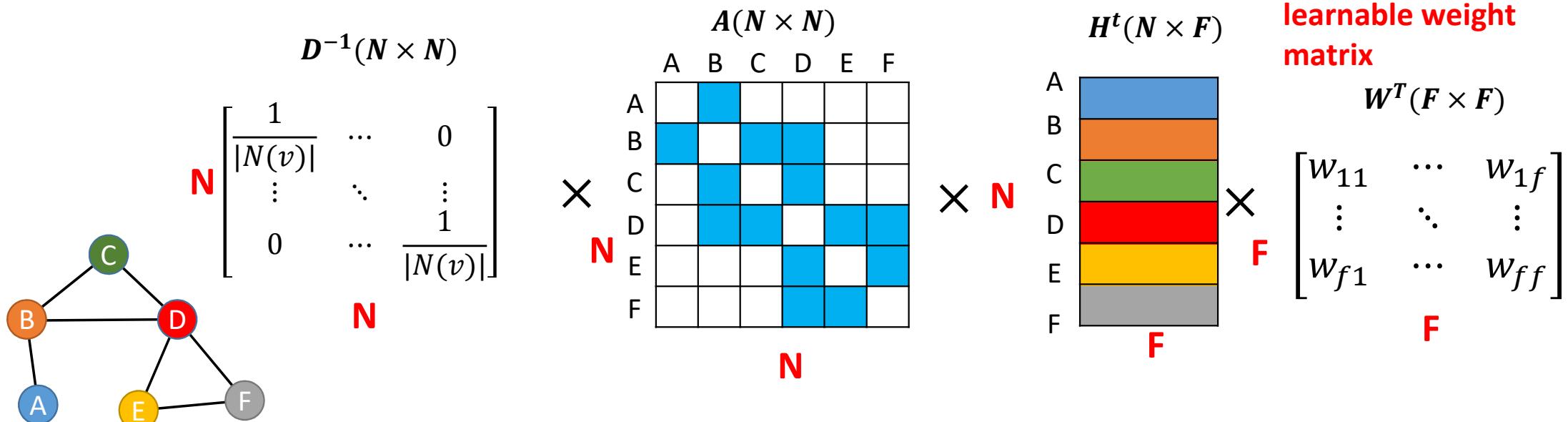
$$h_v^{t+1} = \sigma \left(W_k \left(\sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} \right) + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$



A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \left[\begin{array}{c|cc} & (1 \times F) & (1 \times F) \\ \hline W_k & \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} & + B_k h_v^t \end{array} \right] , \forall t \in (0, \dots, T-1) \right)$$



Noted that W^T is a learnable weight matrix

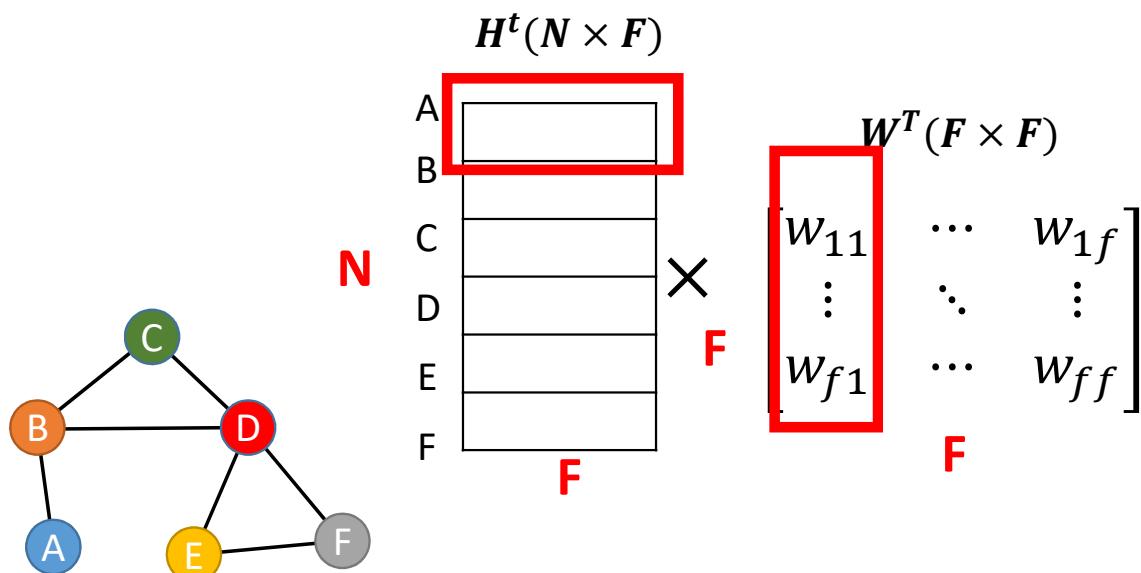
$W^T(F \times F)$

$$\begin{bmatrix} w_{11} & \cdots & w_{1f} \\ \vdots & \ddots & \vdots \\ w_{f1} & \cdots & w_{ff} \end{bmatrix} \times F$$

A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$



Why put W^T on the right hand side of H^t ?

Why not left? With a shape of $(N \times N)$?

A single layer of GNN: Graph Convolution-Forward

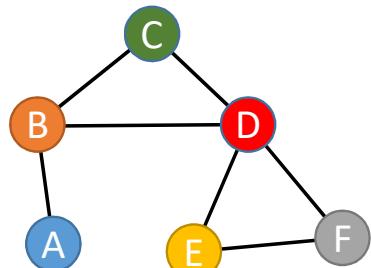
Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

$$\mathbf{N} \quad \begin{matrix} & W^T(N \times N) \\ \begin{matrix} W_{11} & \cdots & W_{1n} \\ \vdots & \ddots & \vdots \\ W_{n1} & \cdots & W_{nn} \end{matrix} & \times \end{matrix} \quad \mathbf{N} \quad \begin{matrix} H^t(N \times F) \\ \begin{matrix} A \\ B \\ C \\ D \\ E \\ F \end{matrix} \end{matrix}$$

Like this?

Seems like nothing goes wrong, the result matrix shape is still $(N \times F)$?



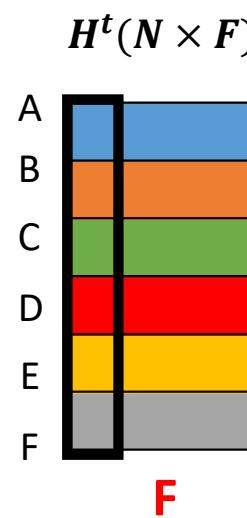
What happen if we still put W on the left hand site?

A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

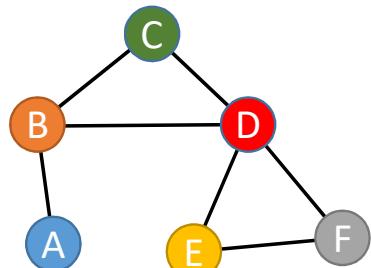
$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

$$\begin{matrix} & W^T(N \times N) \\ N & \begin{bmatrix} W_{11} & \cdots & W_{1n} \\ \vdots & \ddots & \vdots \\ W_{n1} & \cdots & W_{nn} \end{bmatrix} \times N \\ & N \end{matrix}$$



Seems like nothing goes wrong, the result matrix shape is still $(N \times F)$?

No, it's wrong, because we are still mixing information among different nodes, which has the same function with adjacent matrix, feature within node does not receive any mixing

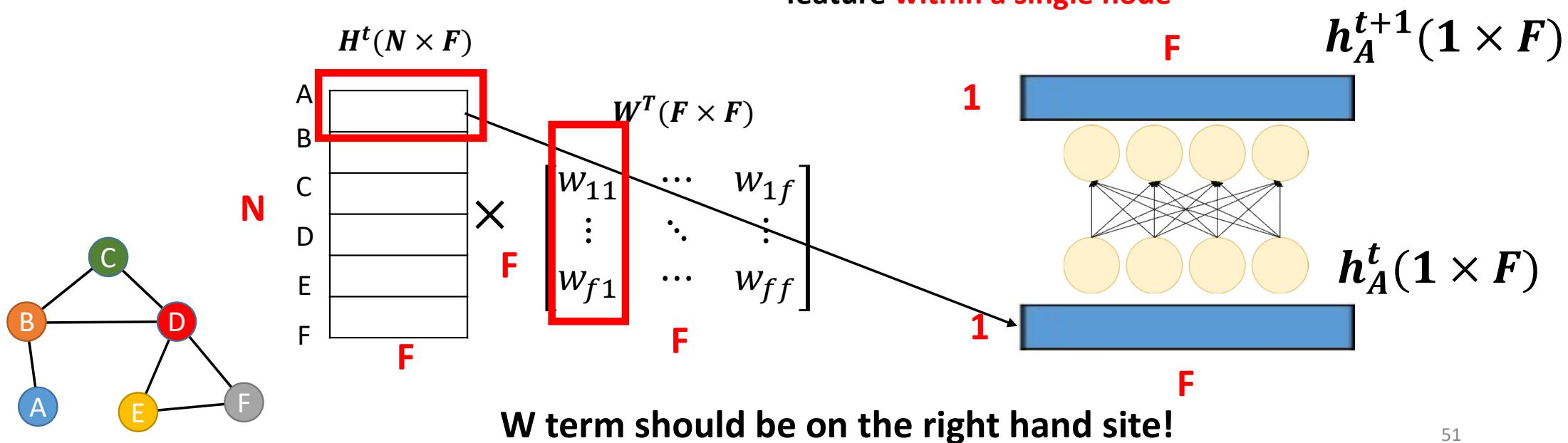


A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

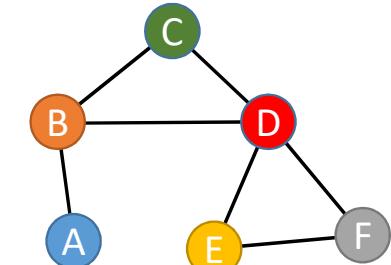
Learnable weight is used to mix information along the feature **within a single node**



A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

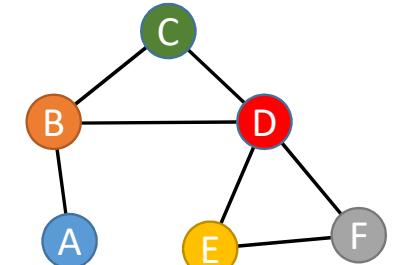


$$\begin{array}{c}
 \mathbf{D}^{-1}(N \times N) \\
 \mathbf{N} \begin{bmatrix} \frac{1}{|N(v)|} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \frac{1}{|N(v)|} \end{bmatrix} \times \mathbf{N} \quad \mathbf{A}(N \times N) \\
 \begin{array}{c} \mathbf{A} \quad \mathbf{B} \quad \mathbf{C} \quad \mathbf{D} \quad \mathbf{E} \quad \mathbf{F} \\ \hline \mathbf{A} & & & & & \\ \mathbf{B} & & & & & \\ \mathbf{C} & & & & & \\ \mathbf{D} & & & & & \\ \mathbf{E} & & & & & \\ \mathbf{F} & & & & & \end{array} \times \mathbf{N} \quad \mathbf{H}^t(N \times F) \\
 \mathbf{N} \quad \mathbf{N} \quad \mathbf{F} \quad \mathbf{F} \\
 \mathbf{H}^{t+1}(N \times F) \\
 \mathbf{W}^T(F \times F) \\
 \mathbf{F} \begin{bmatrix} w_{11} & \cdots & w_{1f} \\ \vdots & \ddots & \vdots \\ w_{f1} & \cdots & w_{ff} \end{bmatrix} = \mathbf{N} \quad \mathbf{F} \quad \mathbf{F}
 \end{array}$$

A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$



Noted that B^T is a learnable weight matrix

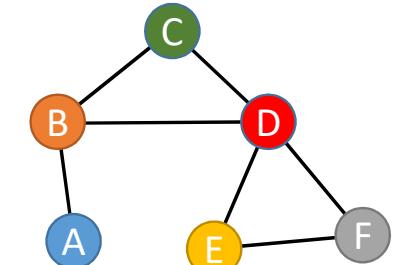
$A'(N \times N)$ Identity matrix	$H_v^t (N \times F)$	$B^T (F \times F)$	$H_v^{t+1} (N \times F)$																																																	
<table border="1" style="border-collapse: collapse; width: 100px;"> <tr><td style="width: 20px;"></td><td style="width: 20px;">A</td><td style="width: 20px;">B</td><td style="width: 20px;">C</td><td style="width: 20px;">D</td><td style="width: 20px;">E</td><td style="width: 20px;">F</td></tr> <tr><td>A</td><td style="background-color: #00FFFF;"></td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>B</td><td></td><td style="background-color: #FF8C00;"></td><td></td><td></td><td></td><td></td></tr> <tr><td>C</td><td></td><td></td><td style="background-color: #3CB371;"></td><td></td><td></td><td></td></tr> <tr><td>D</td><td></td><td></td><td></td><td style="background-color: #FF0000;"></td><td></td><td></td></tr> <tr><td>E</td><td></td><td></td><td></td><td></td><td style="background-color: #FFDAB9;"></td><td></td></tr> <tr><td>F</td><td></td><td></td><td></td><td></td><td></td><td style="background-color: #A9A9A9;"></td></tr> </table>		A	B	C	D	E	F	A							B							C							D							E							F							\times	\times	$=$
	A	B	C	D	E	F																																														
A																																																				
B																																																				
C																																																				
D																																																				
E																																																				
F																																																				
N	F	F	F																																																	

Self loop adjacent matrix is a diagonal matrix!

A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$



Now let's rewrite the scalar form above into matrix form

$$H^{t+1} = \sigma(D^{-1} A H^t W^T + A' H_v^t B^T)$$

Non-Linear Activation
Aggregating neighbor node feature Aggregating self node feature

A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

$$H^{t+1} = \sigma(D^{-1} \hat{A} H^{t'} W'^T)$$

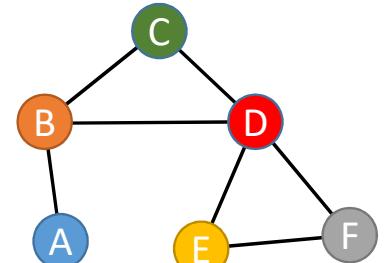
	$A(N \times N)$					
	A	B	C	D	E	F
A						
B						
C						
D						
E						
F						

	$A'(N \times N)$					
	A	B	C	D	E	F
A						
B						
C						
D						
E						
F						

+

	$\hat{A}(N \times N)$					
	A	B	C	D	E	F
A						
B						
C						
D						
E						
F						

=



A single layer of GNN: Graph Convolution-Forward

Matrix form for a single layer of graph convolution

$$h_v^{t+1} = \sigma \left(W_k \sum_{u \in N(v)} \frac{h_u^t}{|N(v)|} + B_k h_v^t \right), \forall t \in (0, \dots, T-1)$$

$$H^{t+1} = \sigma(D^{-1} \hat{A} H^t W'^T)$$

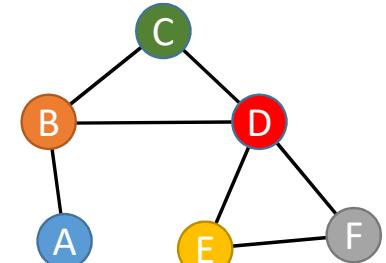
	A($N \times N$)					
A		B	C	D	E	F
B						
C						
D						
E						
F						

	A'($N \times N$)					
A		B	C	D	E	F
B						
C						
D						
E						
F						

+

	\hat{A}($N \times N$)					
A		B	C	D	E	F
B						
C						
D						
E						
F						

=



($N \times N$)

Forward equation
for GCN

Story so far

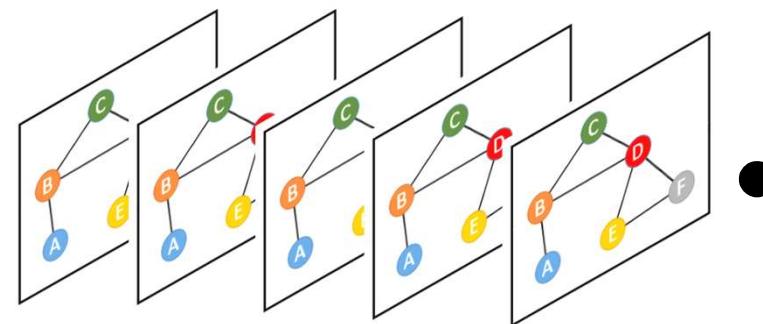
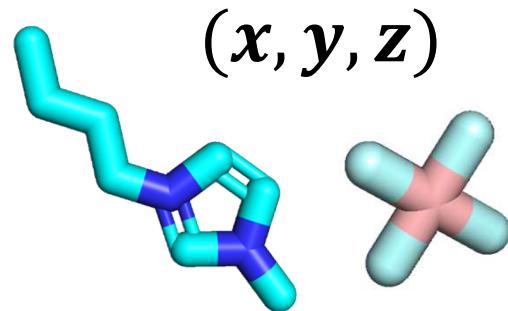
- Node's neighborhood defines a computation graph
- Graph Convolution layer forward

Reference

- Kipf, T.N. and Welling, M., 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*.
- Stanford CS 224 W

Graph: Revisit MLP

Actually, if we include spatial information into GNN node feature
GNN will become neither translation invariant nor rotation invariant



But the permutation invariance for GNN still hold

This problem leads to another topic on GNN which is Equivariant Graph Neural Network, but we don't have time to discuss on that in the lecture