# Pattern Recognition

# Transformers

林彥宇 教授
**Yen-Yu Lin, Professor**

國立陽明交通大學 資訊工程學系
**Computer Science, National Yang Ming Chiao Tung University**

Some slides are modified from Hong-yi Lee,
Shusen Wang, Sy-Kai Chen, and Shao-Hao Lu

# Outline

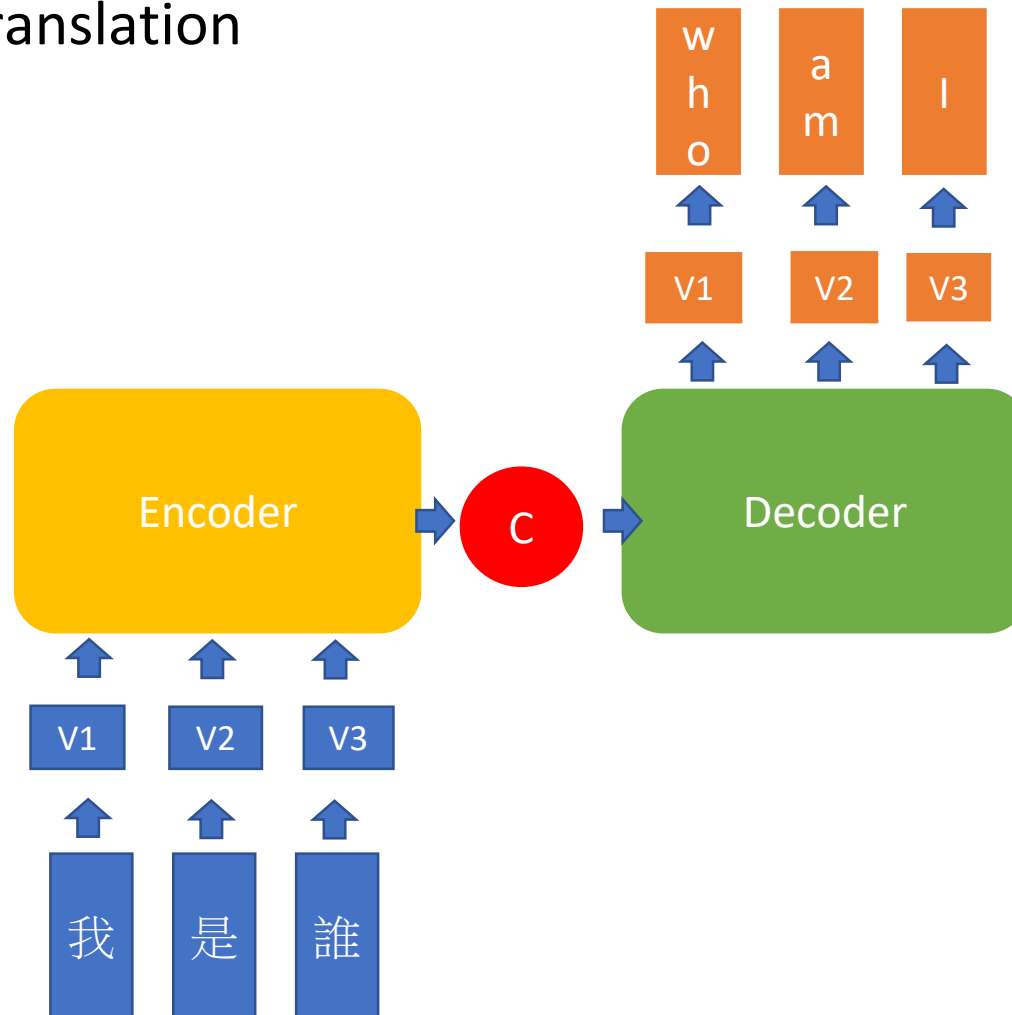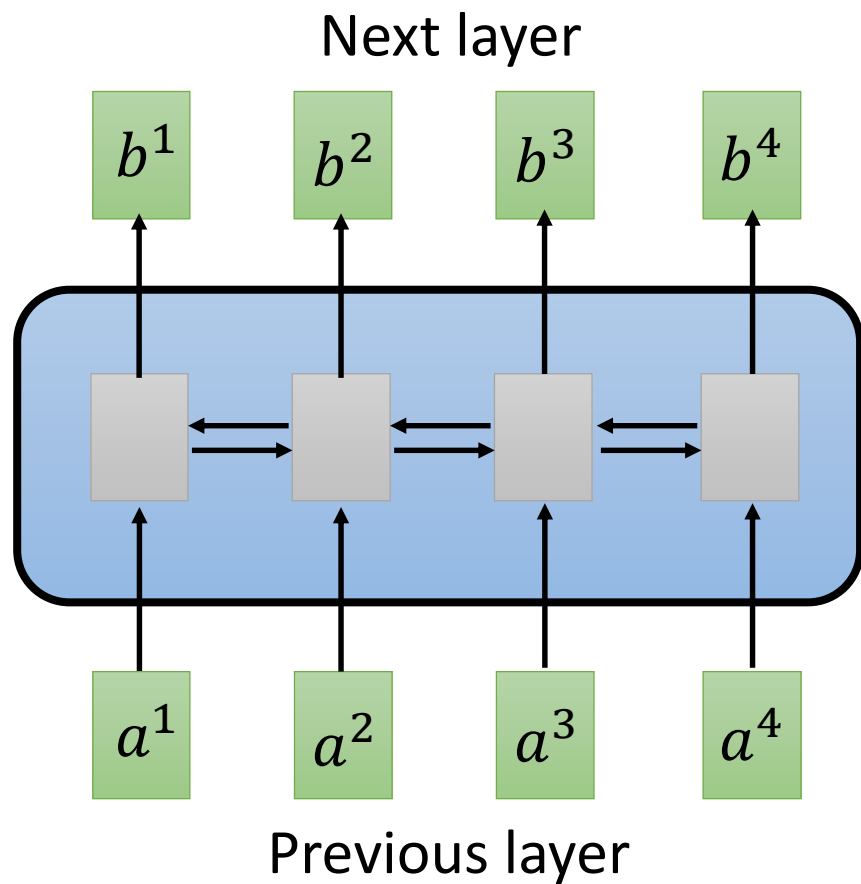- Transformer

- ViT

# Outline

- <span style="color:blue">Transformer</span>
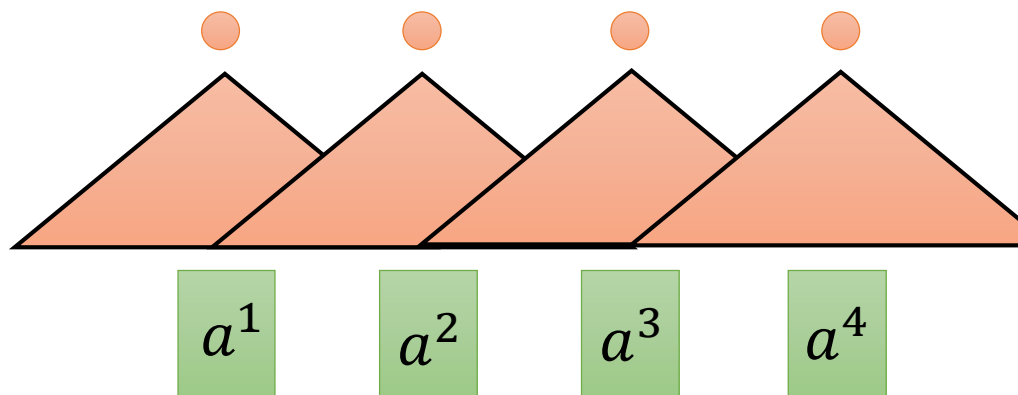  - <span style="color:blue">➤ Most slides are (modified) from Prof. Hung-yi Lee</span>
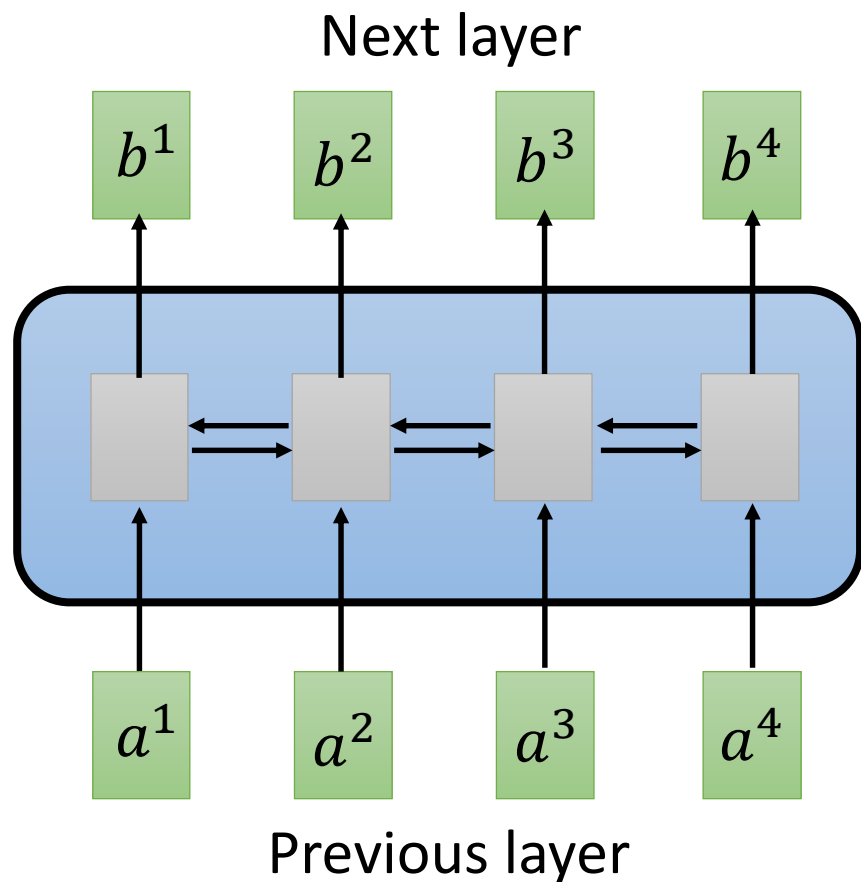
- ViT

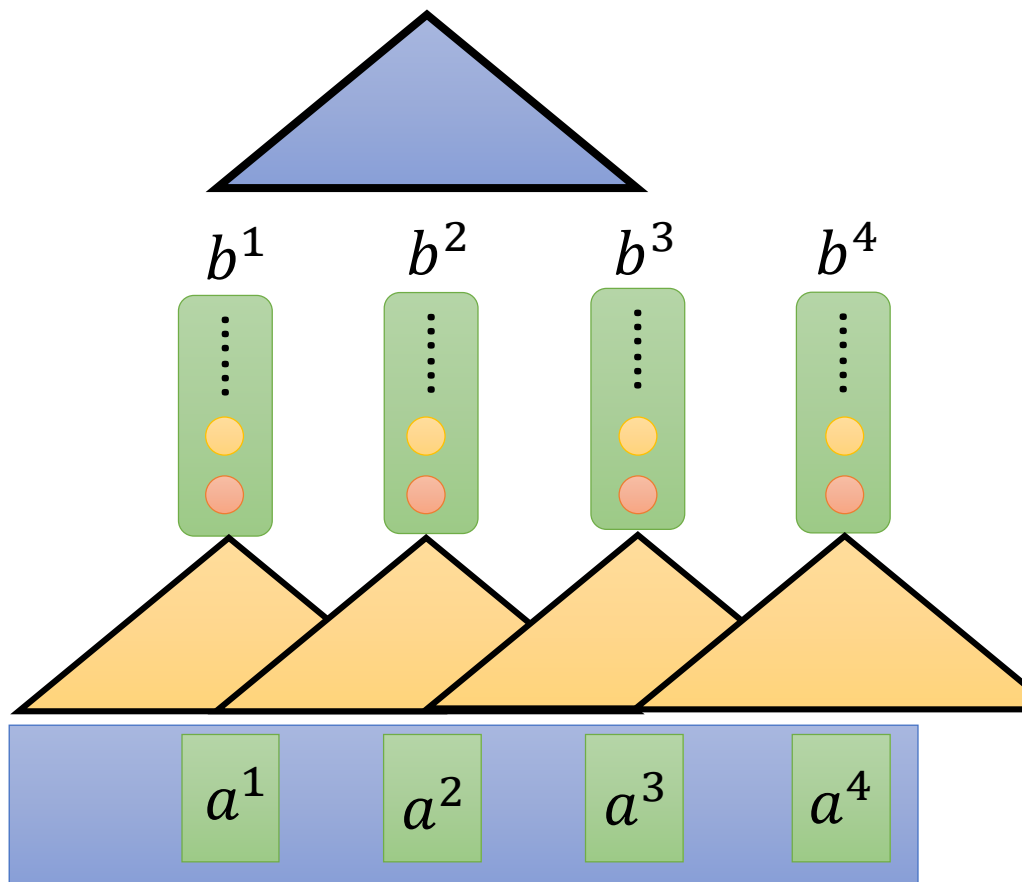# Seq2seq model

- Language translation

# Sequence

Next layer

$b^1$  $b^2$  $b^3$  $b^4$

$a^1$  $a^2$  $a^3$  $a^4$

Previous layer

Hard to parallel !

$a^1$  $a^2$  $a^3$  $a^4$

Using CNN to replace RNN

**Sequence**

Next layer

$b^1$ $b^2$ $b^3$ $b^4$

$a^1$ $a^2$ $a^3$ $a^4$

Previous layer

Hard to parallel

Filters in higher layer can consider longer sequence

$b^1$ $b^2$ $b^3$ $b^4$

$a^1$ $a^2$ $a^3$ $a^4$

Using CNN to replace RNN
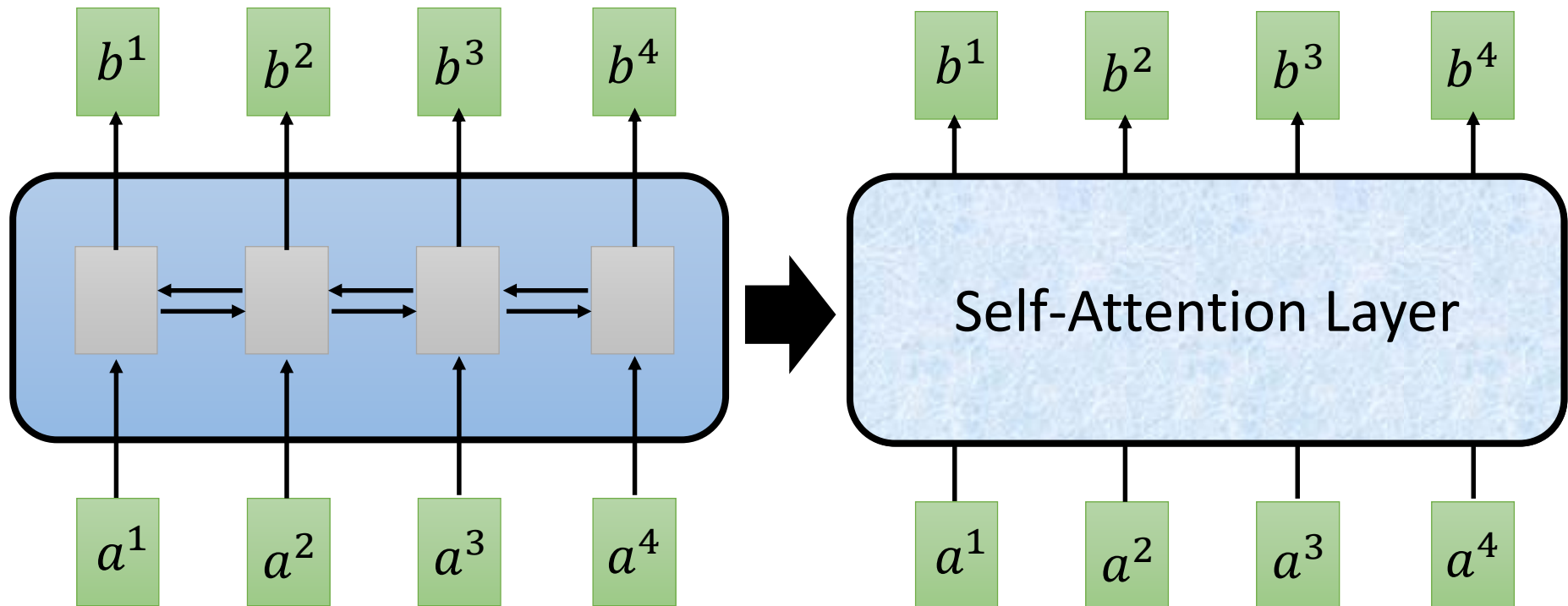(CNN can parallel)

# Self-attention

- What is self attention?
  - Self attention is an attention mechanism relating <span style="color:red">different positions of a single sequence</span> in order to compute a representation of the sequence

- Why self attention?
  - Solve the fundamental constraint in recurrent models

# Self-Attention

$b^i$ is obtained based on the whole input sequence.

$b^1, b^2, b^3, b^4$ can be parallelly computed.
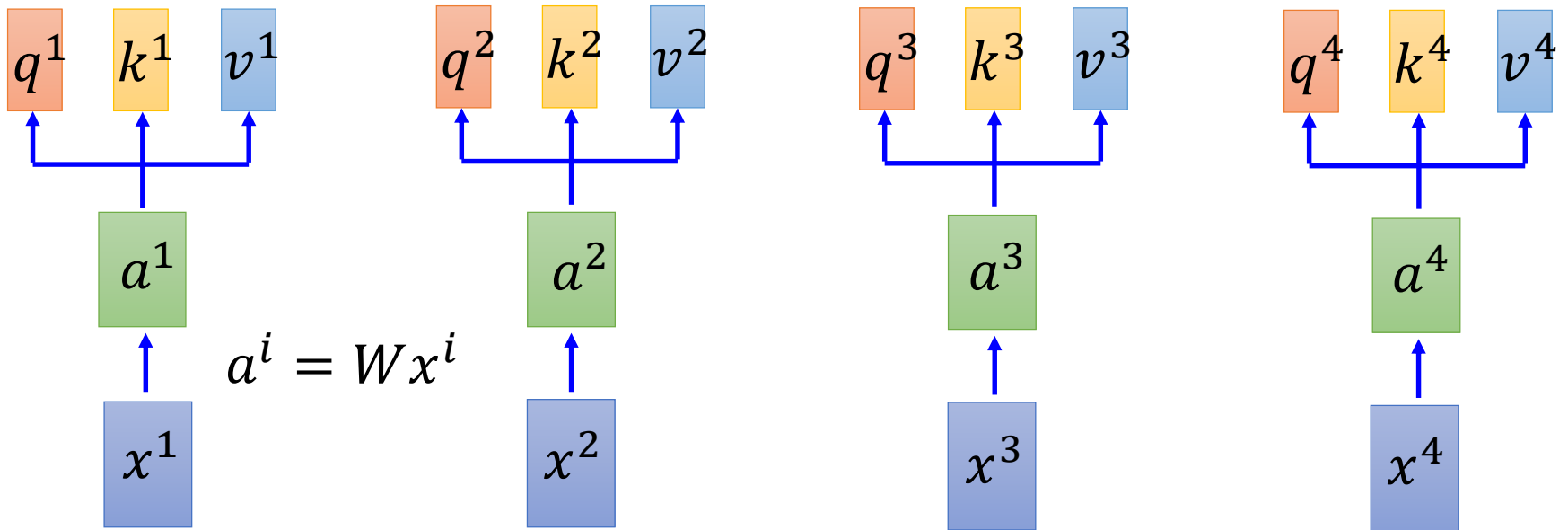
# *Self-attention*

$q$: query (to match others)

$$q^i = W^q a^i$$

$k$: key (to be matched)

$$k^i = W^k a^i$$

$v$: value (information to be extracted)
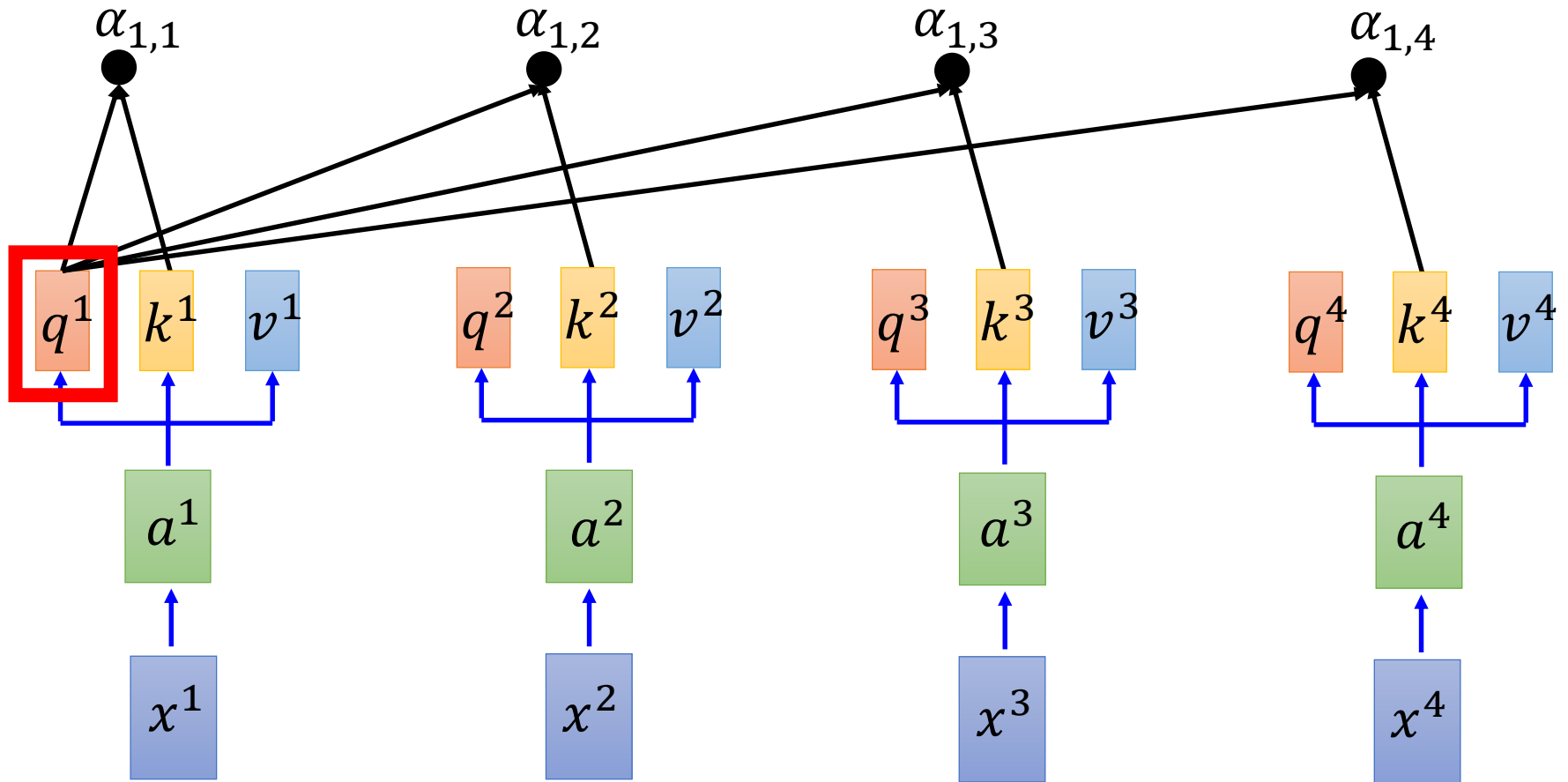
$$v^i = W^v a^i$$

$$a^i = W x^i$$

# *Self-attention*
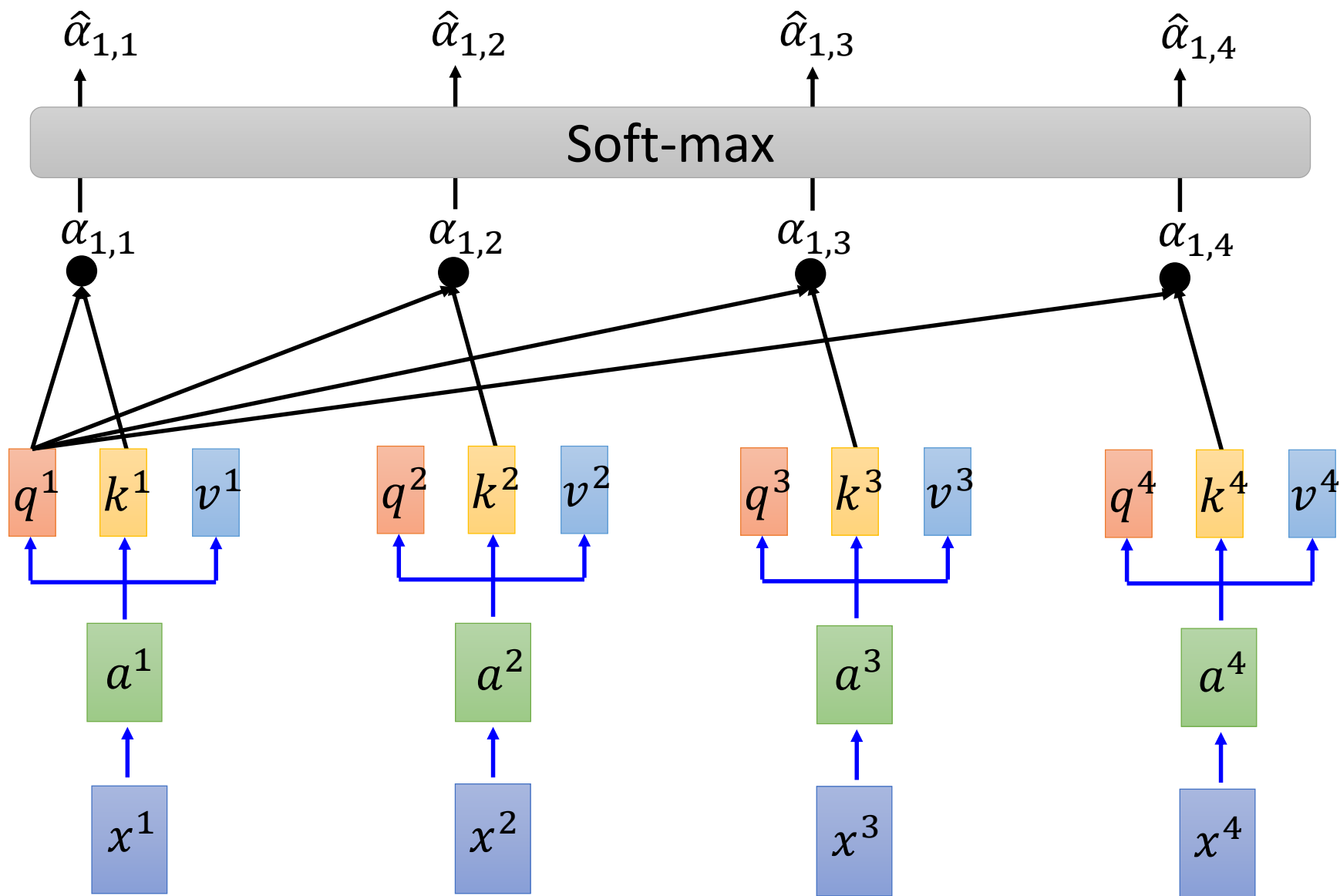
Calculate attention between each query and all keys

d is the dim of $q$ and $k$

Scaled Dot-Product Attention: $\alpha_{1,i} = \underbrace{q^1 \cdot k^i}_{\text{dot product}} / \sqrt{d}$
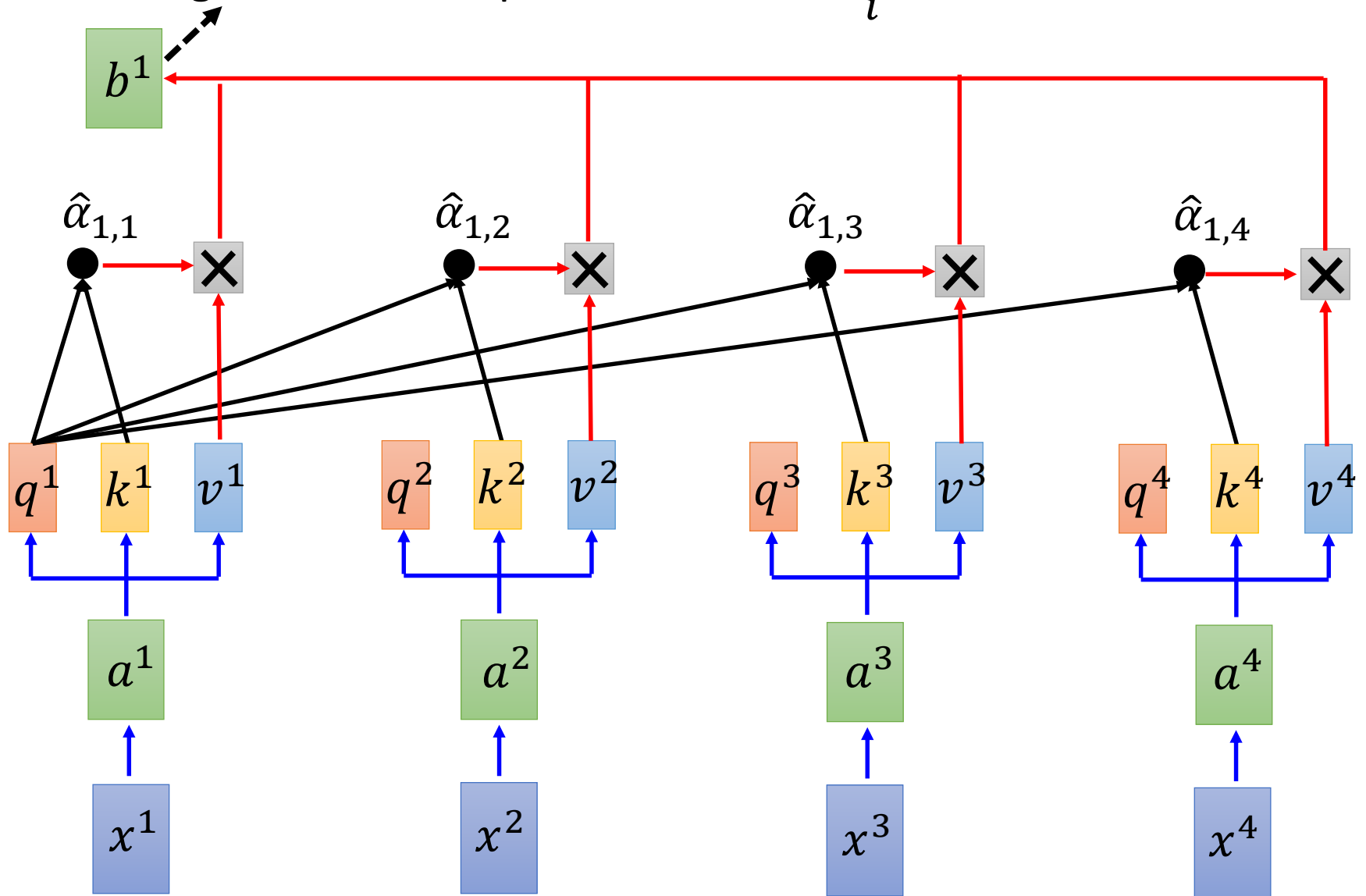
# Self-attention

$$\hat{\alpha}_{1,i} = exp(\alpha_{1,i}) / \sum_j exp(\alpha_{1,j})$$

# *Self-attention*
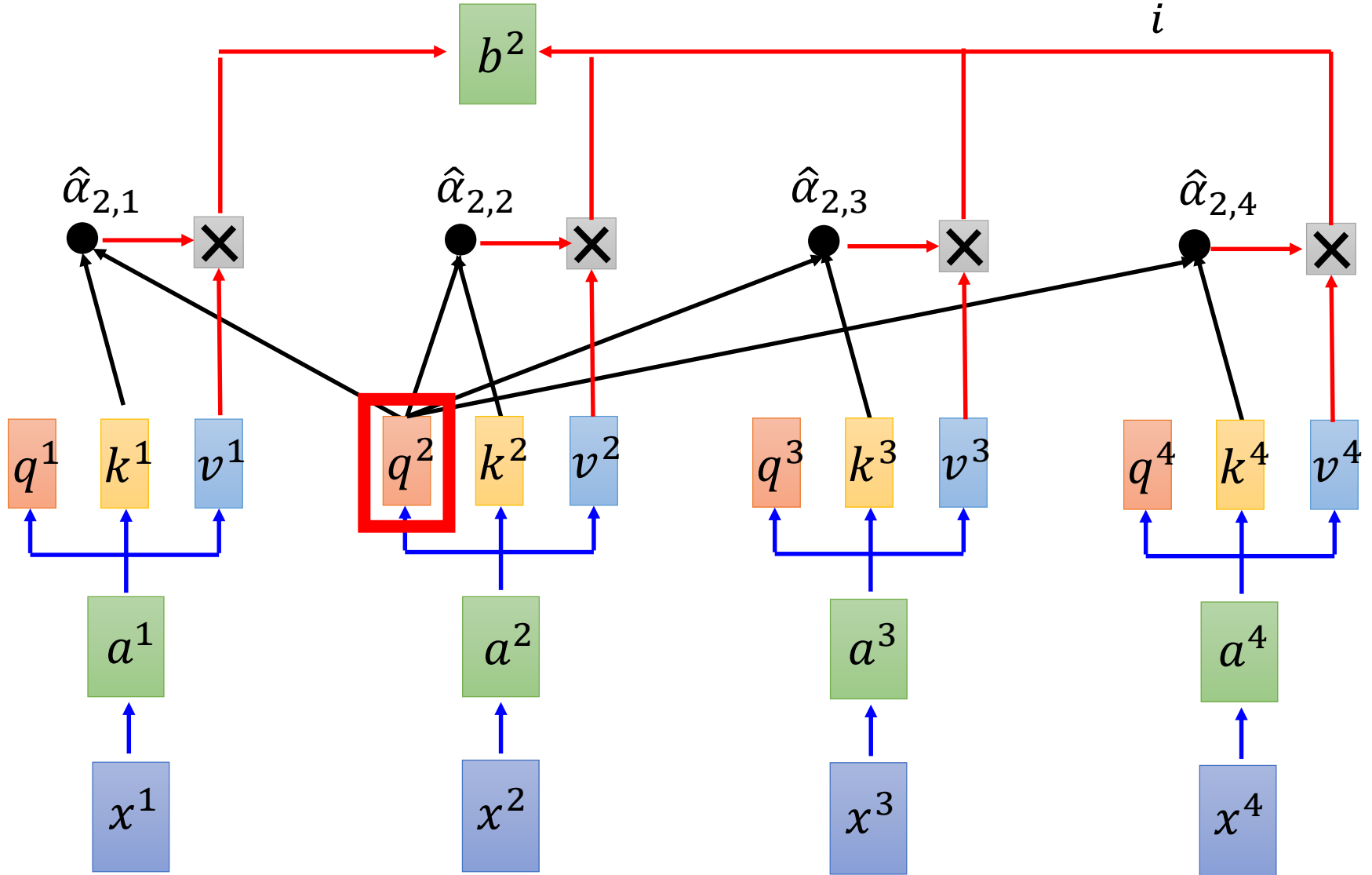
Considering the whole sequence

$$b^1 = \sum_i \hat{\alpha}_{1,i} v^i$$
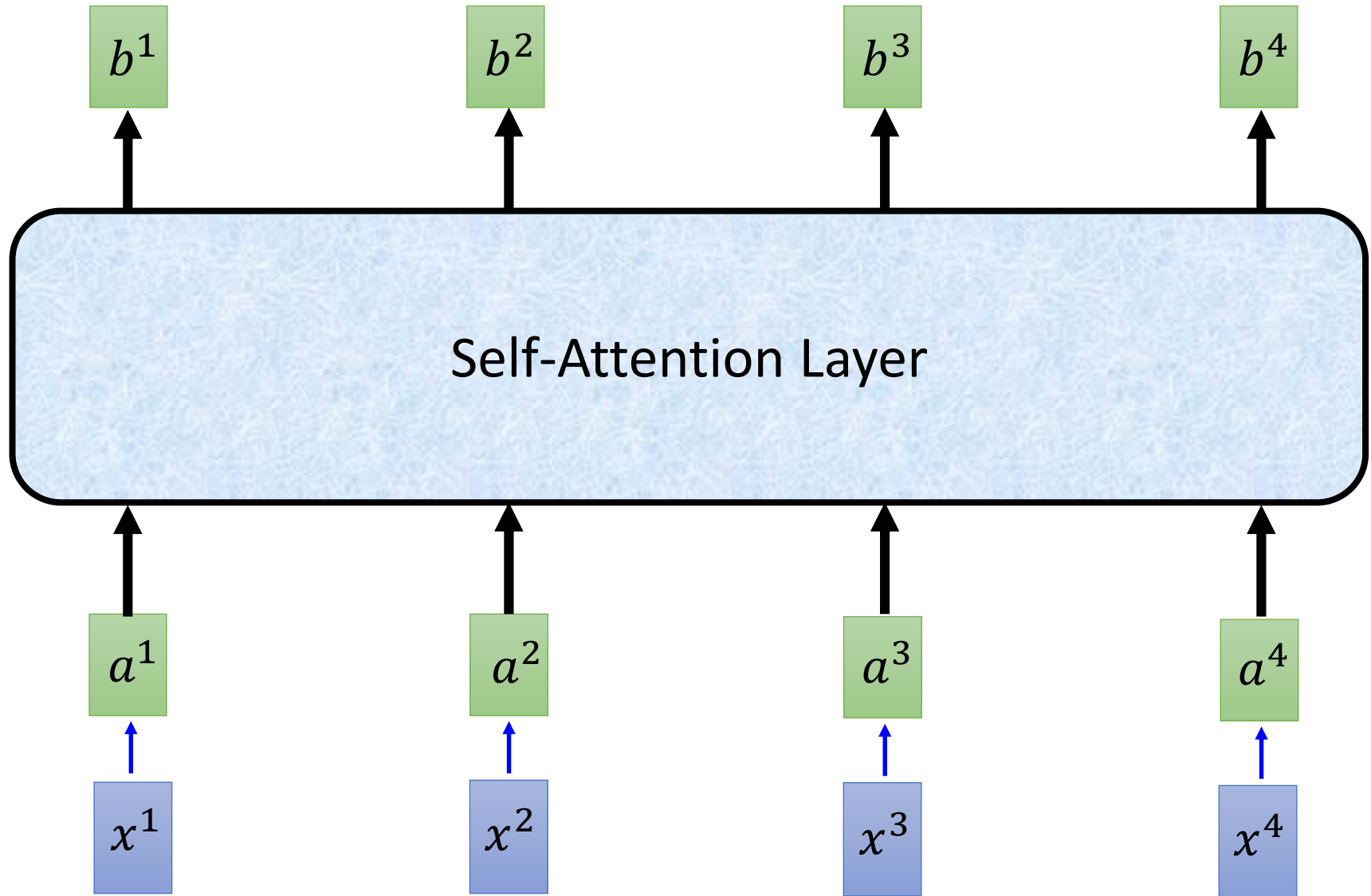
# Self-attention

Calculate attention between each query and all keys

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

# Self-attention

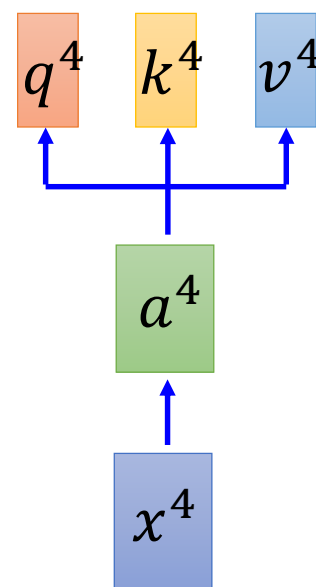$b^1, b^2, b^3, b^4$ can be parallelly computed.

# *Self-attention*

$$q^i = W^q a^i$$

$$k^i = W^k a^i$$

$$v^i = W^v a^i$$

$$\underbrace{q^1 \; q^2 \; q^3 \; q^4}_{Q} = W^q \underbrace{a^1 \; a^2 \; a^3 \; a^4}_{I}$$

$$\underbrace{k^1 \; k^2 \; k^3 \; k^4}_{K} = W^k \underbrace{a^1 \; a^2 \; a^3 \; a^4}_{I}$$

$$\underbrace{v^1 \; v^2 \; v^3 \; v^4}_{V} = W^v \underbrace{a^1 \; a^2 \; a^3 \; a^4}_{I}$$

# Self-attention



$$\alpha_{1,1} = \boxed{k^1}\ \boxed{q^1} \qquad \alpha_{1,2} = \boxed{k^2}\ \boxed{q^1}$$

$$\alpha_{1,3} = \boxed{k^3}\ \boxed{q^1} \qquad \alpha_{1,4} = \boxed{k^4}\ \boxed{q^1}$$
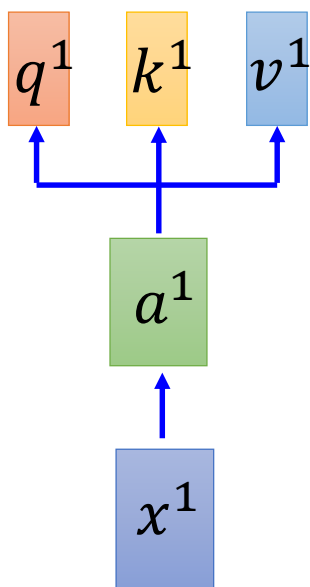
(ignore $\sqrt{d}$ for simplicity)

$$\begin{array}{c} \boxed{\alpha_{1,1}} \\ \boxed{\alpha_{1,2}} \\ \boxed{\alpha_{1,3}} \\ \boxed{\alpha_{1,4}} \end{array} = \begin{array}{c} \boxed{k^1} \\ \boxed{k^2} \\ \boxed{k^3} \\ \boxed{k^4} \end{array} \boxed{q^1}$$
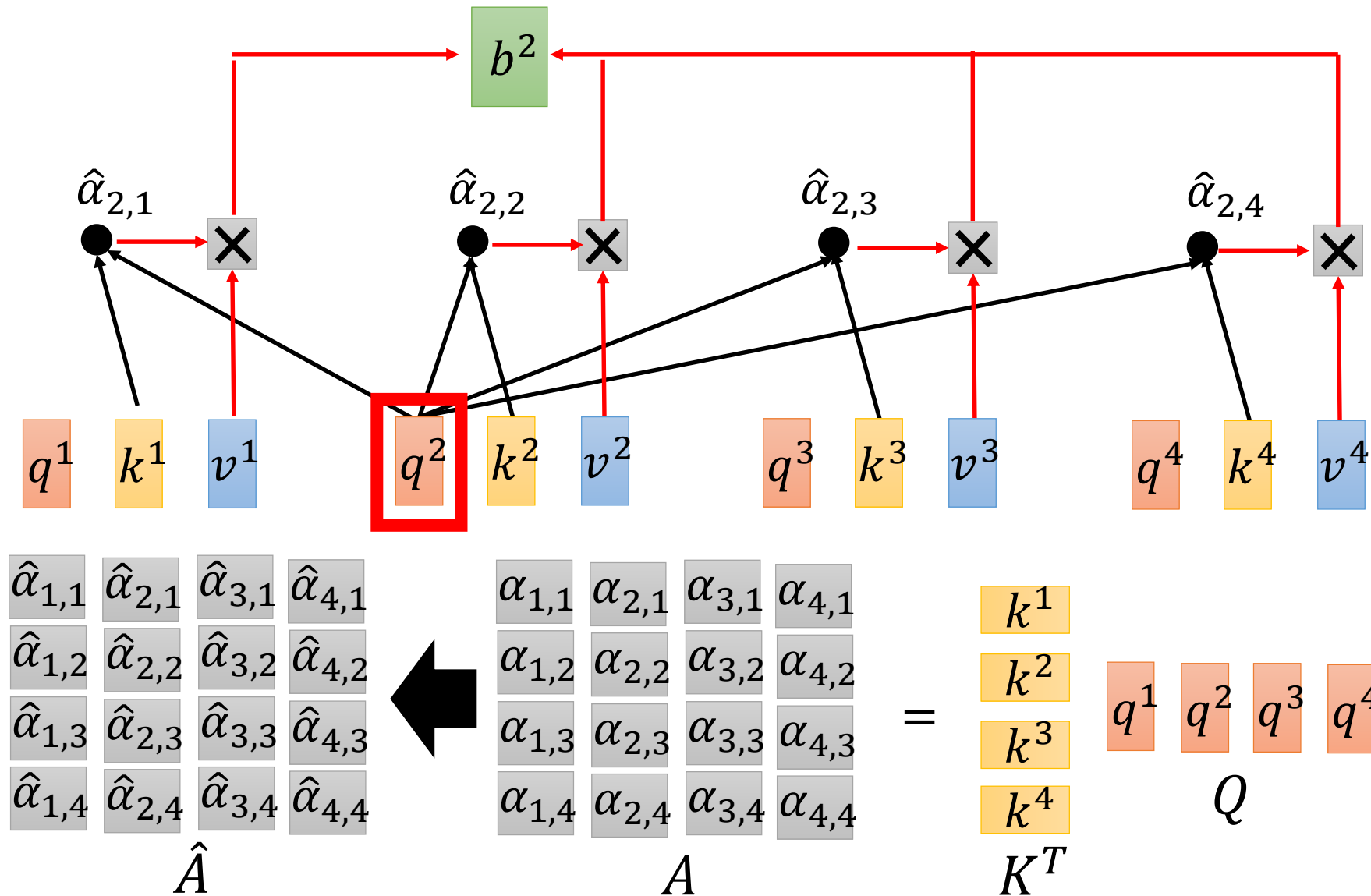
# Self-attention

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

# *Self-attention*

$$b^2 = \sum_i \hat{\alpha}_{2,i} v^i$$

# Self-attention



GPU can accelerate matrix multiplications

# *Multi-head Self-attention*   (2 heads as example)

$$q^{i,1} = W^{q,1}q^i$$

$$q^{i,2} = W^{q,2}q^i$$



$b^{i,1}$

$q^{i,1}$  $q^{i,2}$  $k^{i,1}$  $k^{i,2}$  $v^{i,1}$  $v^{i,2}$  $q^{j,1}$  $q^{j,2}$  $k^{j,1}$  $k^{j,2}$  $v^{j,1}$  $v^{j,2}$

$q^i$  $k^i$  $v^i$  $q^j$  $k^j$  $v^j$

$$q^i = W^q a^i$$

$a^i$  $a^j$

# Multi-head Self-attention   (2 heads as example)

$$q^{i,1} = W^{q,1} q^i$$

$$q^{i,2} = W^{q,2} q^i$$

$b^{i,1}$

$b^{i,2}$

$q^{i,1}$  $q^{i,2}$  $k^{i,1}$  $k^{i,2}$  $v^{i,1}$  $v^{i,2}$   $q^{j,1}$  $q^{j,2}$  $k^{j,1}$  $k^{j,2}$  $v^{j,1}$  $v^{j,2}$

$q^i$   $k^i$   $v^i$      $q^j$   $k^j$   $v^j$

$$q^i = W^q x^i$$

$a^i$

$a^j$

# Multi-head Self-attention    (2 heads as example)

$$b^i = W^O \begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix}$$

$$\begin{bmatrix} b^{i,1} \\ b^{i,2} \end{bmatrix} \longrightarrow b^i$$

$q^{i,1}$  $q^{i,2}$  $k^{i,1}$  $k^{i,2}$  $v^{i,1}$  $v^{i,2}$    $q^{j,1}$  $q^{j,2}$  $k^{j,1}$  $k^{j,2}$  $v^{j,1}$  $v^{j,2}$

$q^i$  $k^i$  $v^i$    $q^j$  $k^j$  $v^j$

$a^i$    $a^j$

# Positional Encoding

- No position information in self-attention.

- Original paper: each position has a unique positional vector $e^i$ (not learned from data)

# Seq2seq with Attention

# *Transformer*

Using Chinese to English translation as example

machine    learning

Encoder       Decoder



機 器 學 習        \<BOS\>    machine

# *Transformer*

$b'$ → Layer Norm →

$b'$

+

$b$

⋮

$a$

Layer Norm:
https://arxiv.org/abs/1607.06450
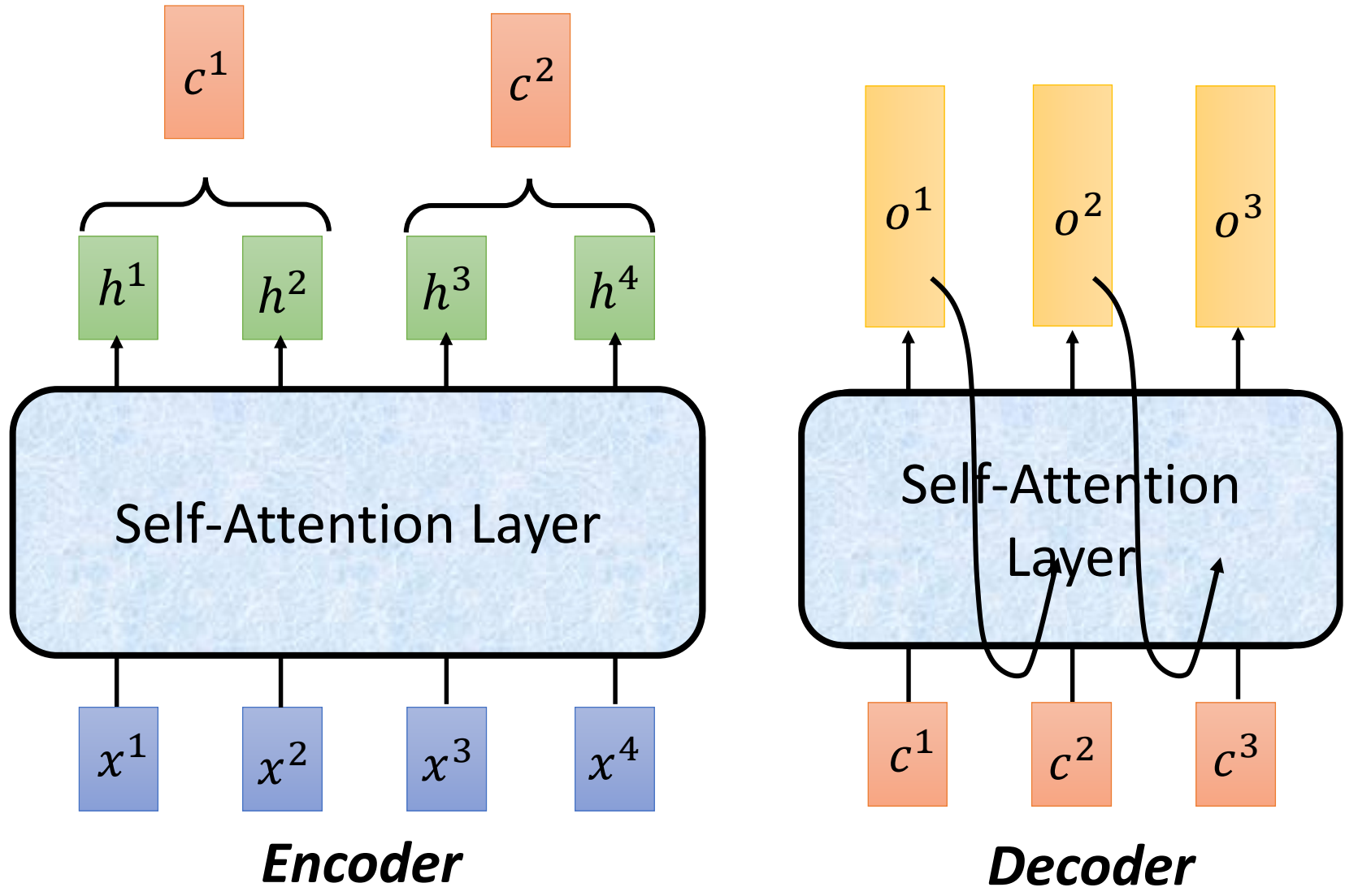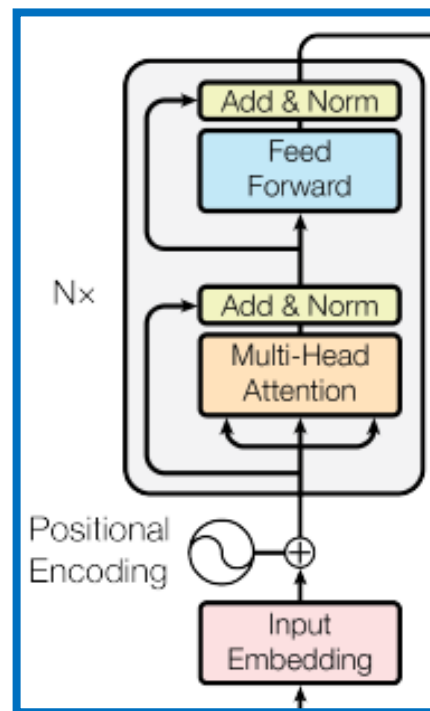
Batch Norm:
https://www.youtube.com/watch?v=BZh1ltr5Rkg

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Batch Size

$\mu = 0, \sigma = 1$ Batch

$\mu = 0, \sigma = 1$ Layer

attend on the input sequence

*Masked*: attend on the generated sequence

$b^1$  $b^2$  $b^3$  $b^4$

Self-Attention Layer

$a^1$  $a^2$  $a^3$  $a^4$

# Outline

- Transformer

- ViT

# Vision Transformer (ViT)

**[Dosovitskiy et al., ICLR'21]**

- Self-attention-based architectures, e.g., Transformers, have become the model of choice in natural language processing
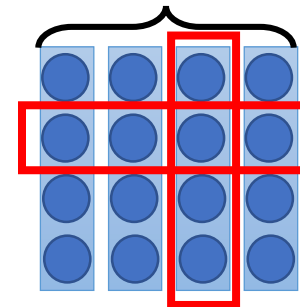
- CNNs, e.g., ResNet, were the state-of-the-art models for object recognition

- Vision transformer (ViT) beats CNNs by a small margin, if the dataset for pre-training is sufficiently large

- ViT is developed mainly based on Transformer
  - ➢ Encoder: linear embedding, multi-head attention, layer normalization, …
  - ➢ Position encoding

# Object recognition

# An image as a set of patches

- To apply Transformer to object recognition with minor modifications, we split an image into patches

- Patches are specified by the user

  ➢ Overlapping or non-overlapping

  ➢ Patch size

# Linear embeddings and positional encoding of patches

- We embed each patch by linear projection $\mathbf{E}$
- 1D learnable positional encoding is used and added to linear embedding

# Class token

- A special token, called class token, is included
- Class token does not belong to any patch
- It is represented by a vector that is learnable
- During N blocks, it gathers information from other tokens, and its output vector is used to derive a classifier



**CLS**

**E**  **E**  **E**

•  •  •

**E**

# Overview

1. Split an image into fixed-size patches (16 x 16), linearly embed each of them, and then add CLS token (learnable)
2. Add positional encoding
3. Transformer encoder operation
4. Use [cls] token to do the classification

# Transformer encoder



$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \, \mathbf{x}_p^1 \mathbf{E}; \, \mathbf{x}_p^2 \mathbf{E}; \cdots; \, \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{pos},$$

$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1},$$

$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell,$$

$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0)$$

# ViT vs. Transformer

1. VIT only uses the encoder part
2. Split image into patches and embed them as tokens.
   (Word vectors as tokens in Transformer)
3. Add one additional learnable token: [CLS] token
   (It servers as conceptual summary about the whole image tokens)
4. Different positional encoding methods

# ViT vs. Transformer

- 2. Split an image into patches and embed them as tokens



unroll

embedding

# ViT vs. Transformer

- 3. Add one additional learnable token: [CLS] token
  - ➢ Original transformer is designed for NLP (sec2sec task), so its output is also a sequence (a series of word vectors).
  - ➢ ViT is used for classification task (only 1 output vector). It would be unfair to choose one particular token for image prediction. Hence, one additional learnable token is included as conceptual summary about the whole image tokens and for final prediction



**Vision Transformer (ViT)**

[CLS] token

# ViT vs. Transformer

- 4. Different positional encoding methods
  - ➢ Positional encoding in the original transformer is given by using sine and cosine functions
  - ➢ Positional encoding in ViT is learnable
    - ◆ Cosine similarity among the learned position encoding



Position embedding similarity

# Pre-training, fine-tuning, and testing

# Experiments: Network configurations

- Three variants of ViT of different sizes

| Model | Layers | Hidden size $D$ | MLP size | Heads | Params |
|---|---|---|---|---|---|
| ViT-Base | 12 | 768 | 3072 | 12 | 86M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632M |

➢ ViT-L/16: the Large variant with 16 x 16 input patch size

➢ ViT-H/14: the Huge variant with 14 x 14 input patch size

# Competing methods and datasets

- Competing methods
  - ➤ BiT (Big Transfer): A variant of ResNet
  - ➤ Noisy Student: A variant of EfficientNet

- Datasets for pre-training
  - ➤ ImageNet: 1.3M images of 1K classes (small)
  - ➤ ImageNet-21K: 14M images of 21K classes (medium)
  - ➤ JFT: 300M images of 18K classes (large)
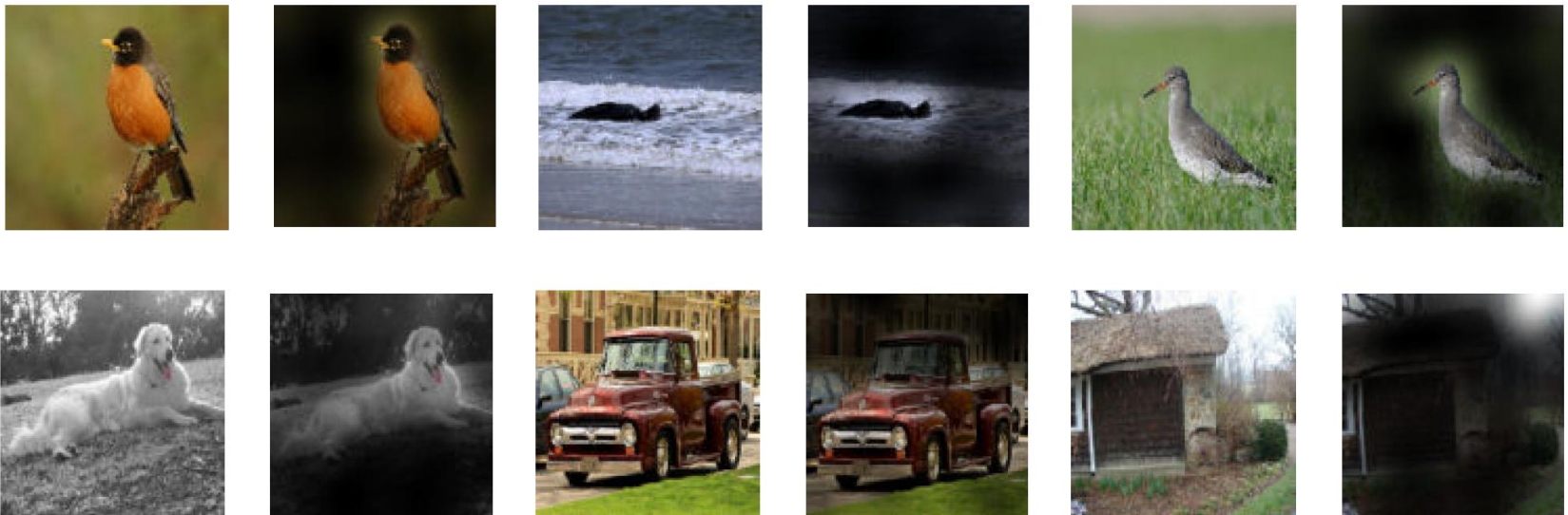
National Chiao Tung University

# Experimental results: Comparison with SoTA

- With the model pre-trained on ImageNet-21K, ViT is comparable with (or slightly underperforms) ResNet

- With the model pre-trained on JFT, ViT slightly outperforms ResNet, while taking less computational resources for pre-training

|  | Ours-JFT (ViT-H/14) | Ours-JFT (ViT-L/16) | Ours-I21k (ViT-L/16) | BiT-L (ResNet152x4) | Noisy Student (EfficientNet-L2) |
|---|---|---|---|---|---|
| ImageNet | $\mathbf{88.55} \pm 0.04$ | $87.76 \pm 0.03$ | $85.30 \pm 0.02$ | $87.54 \pm 0.02$ | $88.4/88.5^*$ |
| ImageNet ReaL | $\mathbf{90.72} \pm 0.05$ | $90.54 \pm 0.03$ | $88.62 \pm 0.05$ | $90.54$ | $90.55$ |
| CIFAR-10 | $\mathbf{99.50} \pm 0.06$ | $99.42 \pm 0.03$ | $99.15 \pm 0.03$ | $99.37 \pm 0.06$ | — |
| CIFAR-100 | $\mathbf{94.55} \pm 0.04$ | $93.90 \pm 0.05$ | $93.25 \pm 0.05$ | $93.51 \pm 0.08$ | — |
| Oxford-IIIT Pets | $\mathbf{97.56} \pm 0.03$ | $97.32 \pm 0.11$ | $94.67 \pm 0.15$ | $96.62 \pm 0.23$ | — |
| Oxford Flowers-102 | $99.68 \pm 0.02$ | $\mathbf{99.74} \pm 0.00$ | $99.61 \pm 0.02$ | $99.63 \pm 0.03$ | — |
| VTAB (19 tasks) | $\mathbf{77.63} \pm 0.23$ | $76.28 \pm 0.46$ | $72.72 \pm 0.21$ | $76.29 \pm 1.70$ | — |
| TPUv3-core-days | 2.5k | 0.68k | 0.23k | 9.9k | 12.3k |

國立交通大學
National Chiao Tung University

# Attention derived by ViT

- Attention maps are inferred by using Attention Rollout [Abnar and Zuidema, ACL'20]
- ViT attends to regions semantically relevant to classification

# References

- Transformer
  - ➢ A. Vaswani et al.: Attention is All you Need, in NIPS 2017.
- ViT
  - ➢ A. Dosovitskiy et al.: An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale, in ICLR 2021.

# **Thank You for Your Attention!**

**Yen-Yu Lin (林彥宇)**

Email: lin@cs.nctu.edu.tw
URL: https://www.cs.nctu.edu.tw/members/detail/lin