



Lecture 09

Transformer based Object Detection

Lin ZHANG, PhD
School of Software Engineering
Tongji University
Fall 2024



Outline

- Transformer in NLP
- Multi-head Attention
- Vision transform (ViT)
- Swin-Transformer
- DEtection TRansformer (DETR)
- Real-time DETR (RT-DETR)



Transformer in NLP

- Transformer was first proposed by Google Brain in the domain of NLP^[1]
 - Transformer encoder is used to transform a set of tokens (vectors) \mathcal{I} into another set of tokens \mathcal{O} , where each element in \mathcal{O} can catch the **global information** of \mathcal{I} ; usually the number of tokens in \mathcal{I} and \mathcal{O} are the same
 - This architecture has since become the foundation for many state-of-the-art NLP models, including BERT and GPT (Generative Pre-trained Transformer)
 - It is now also widely applied in computer vision
- Transformer encoder is composed of basic blocks, including **self-attention**, **positional encoding**, MLP, residual connection and **Layer-norm**
 - A set of vectors is transformed into another different set of vectors; that is why such a structure is called “transformer”

[1] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need[C]//Proc. Adv. Neural Inf. Process. Syst., 2017: 6000-6010. (**Cited by 125704**, Jun. 21, 2024)



Transformer in NLP



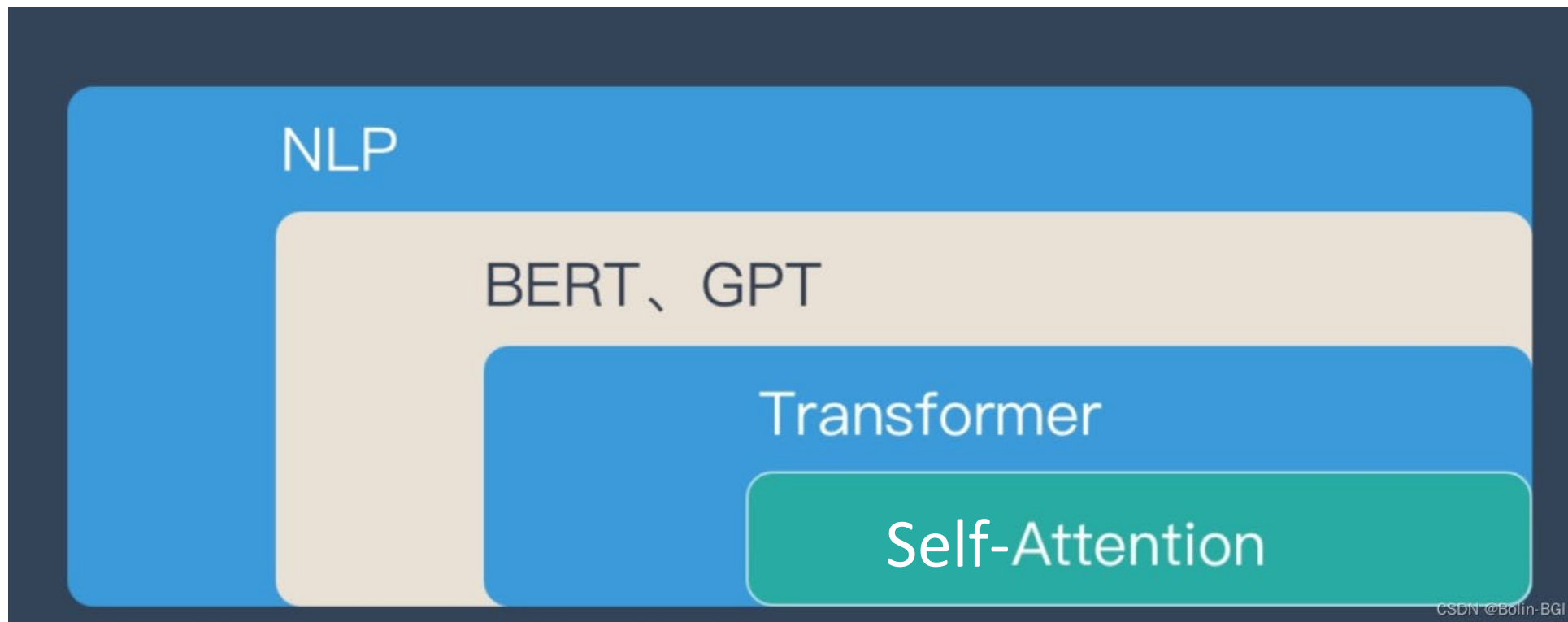
Ashish Vaswani (born in 1986) is a computer scientist. He was a co-founder of Adept AI Labs and a former staff research scientist at Google Brain. Vaswani completed his engineering in Computer Science from BIT Mesra (印度贝拉理工学院, 梅斯拉) in 2002. In 2004, he moved to the US to pursue higher studies at University of Southern California. He did his PhD at the University of Southern California.





Transformer in NLP

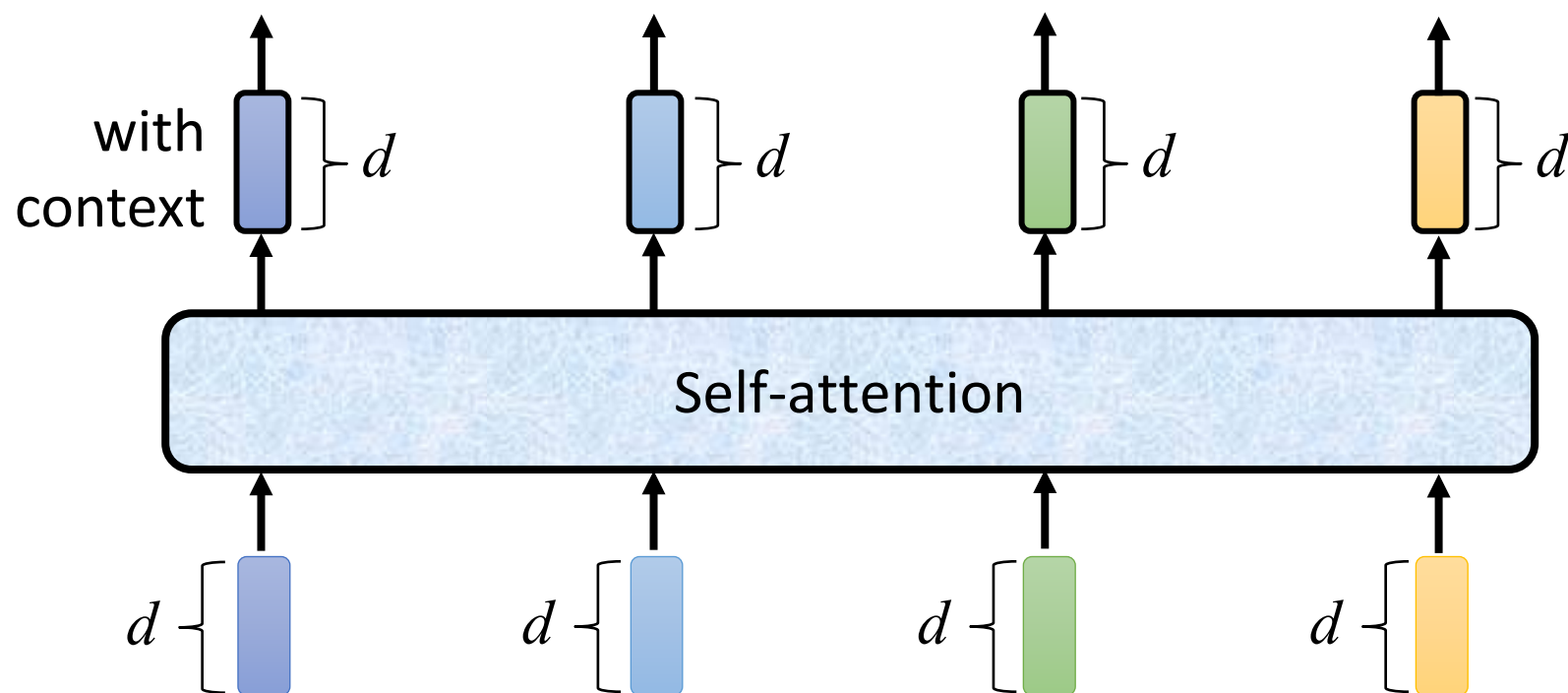
The relationship among several terms, NLP, BERT, GPT, Transformer and self-attention





Self-attention

- It is the core component in Transformer encoder
- Key characteristics of self-attention
 - Input N vectors with dimension d , output N vectors with dimension d
 - Each output vector can capture the **context information of all the input vectors**



Why do we need SA?

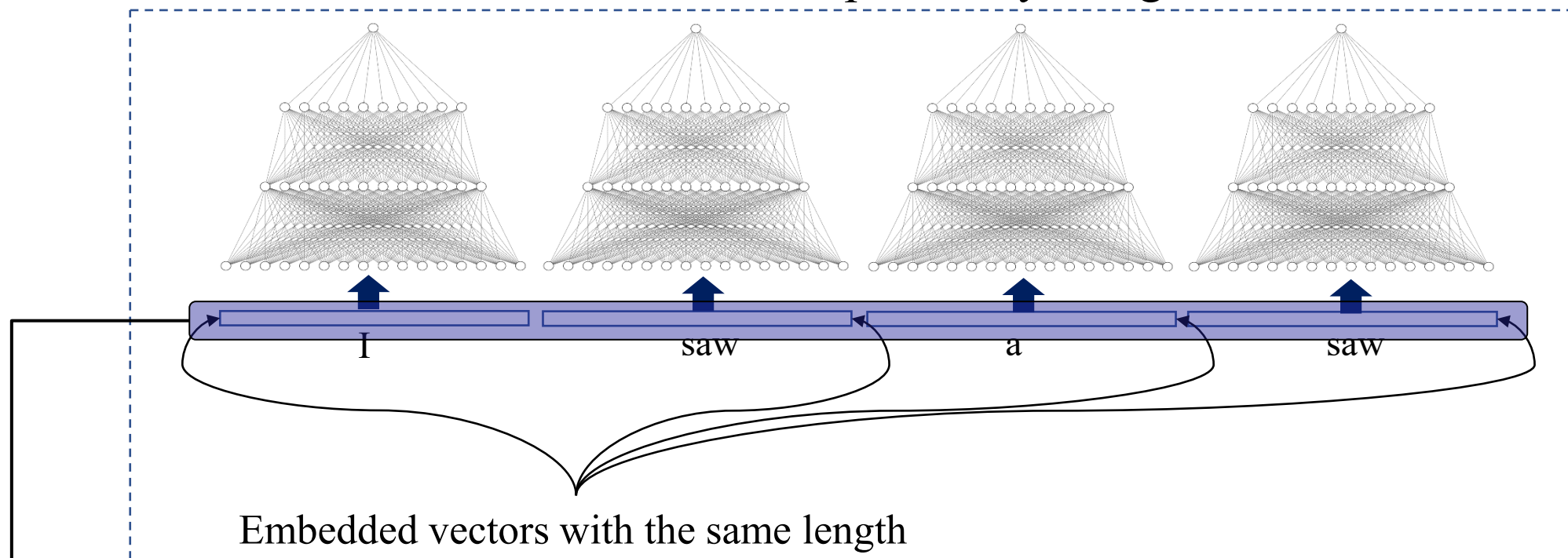
Let's see an example →



Self-attention

Ex: word tagging for an input sentence “I saw a saw”

Naïve solution: deal with each word independently using a network



Is it a good solution? No. Tagging for a word depends on the word itself and also its context

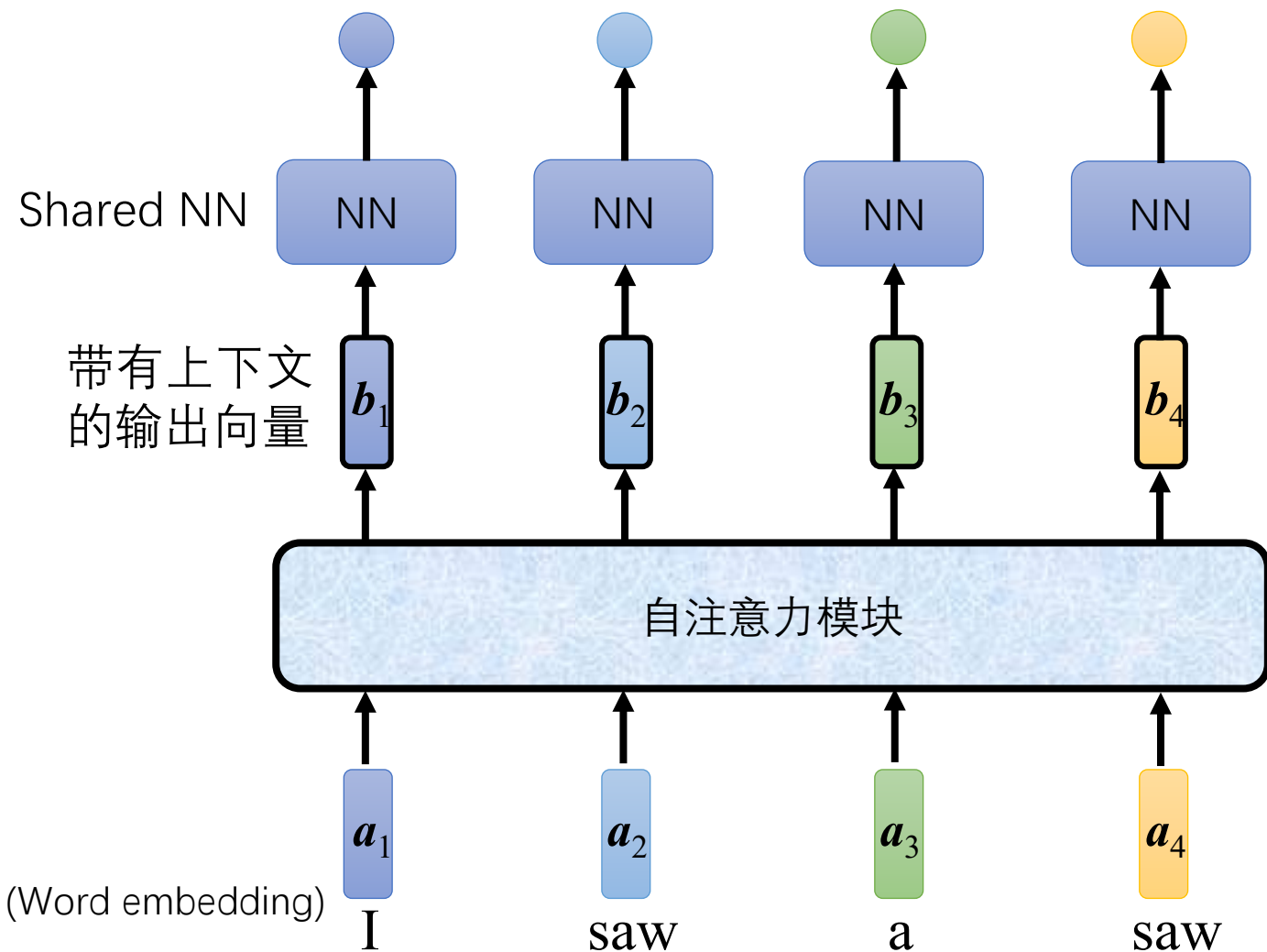
The embedded vectors need to be processed to make them hold the global (contextual) information



Self-attention



Self-attention



- $b_1 \sim b_4$ are computed from $a_1 \sim a_4$ and can be computed in parallel
- Let's see how to compute b_1 in detail; $b_2 \sim b_4$ can be obtained similarly

Output vectors

$$b_1 \sim b_4 \in \mathbb{R}^{d \times 1}$$

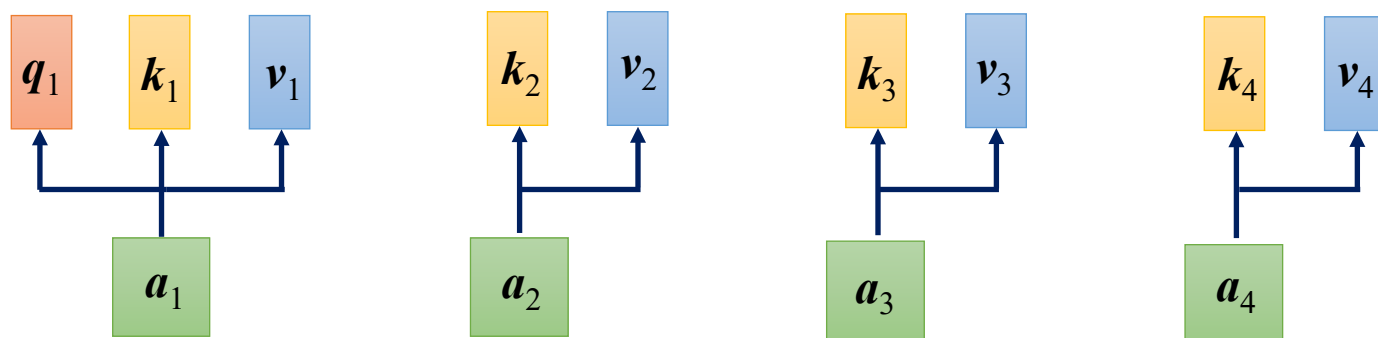
Input vectors

$$a_1 \sim a_4 \in \mathbb{R}^{d \times 1}$$



Self-attention

- ✓ $W_q \in \mathbb{R}^{d \times d}, W_k \in \mathbb{R}^{d \times d}, W_v \in \mathbb{R}^{d \times d}$ are the matrices that need **to be learned by training**
- ✓ q_1 is the **embedded query** vector
- ✓ k_1, k_2, k_3 , and k_4 are the **embedded key** vectors
- ✓ v_1, v_2, v_3 , and v_4 are the **embedded value** vectors



$$q_1 = W_q a_1$$

$$k_1 = W_k a_1$$

$$k_2 = W_k a_2$$

$$k_3 = W_k a_3$$

$$k_4 = W_k a_4$$

$$v_1 = W_v a_1$$

$$v_2 = W_v a_2$$

$$v_3 = W_v a_3$$

$$v_4 = W_v a_4$$

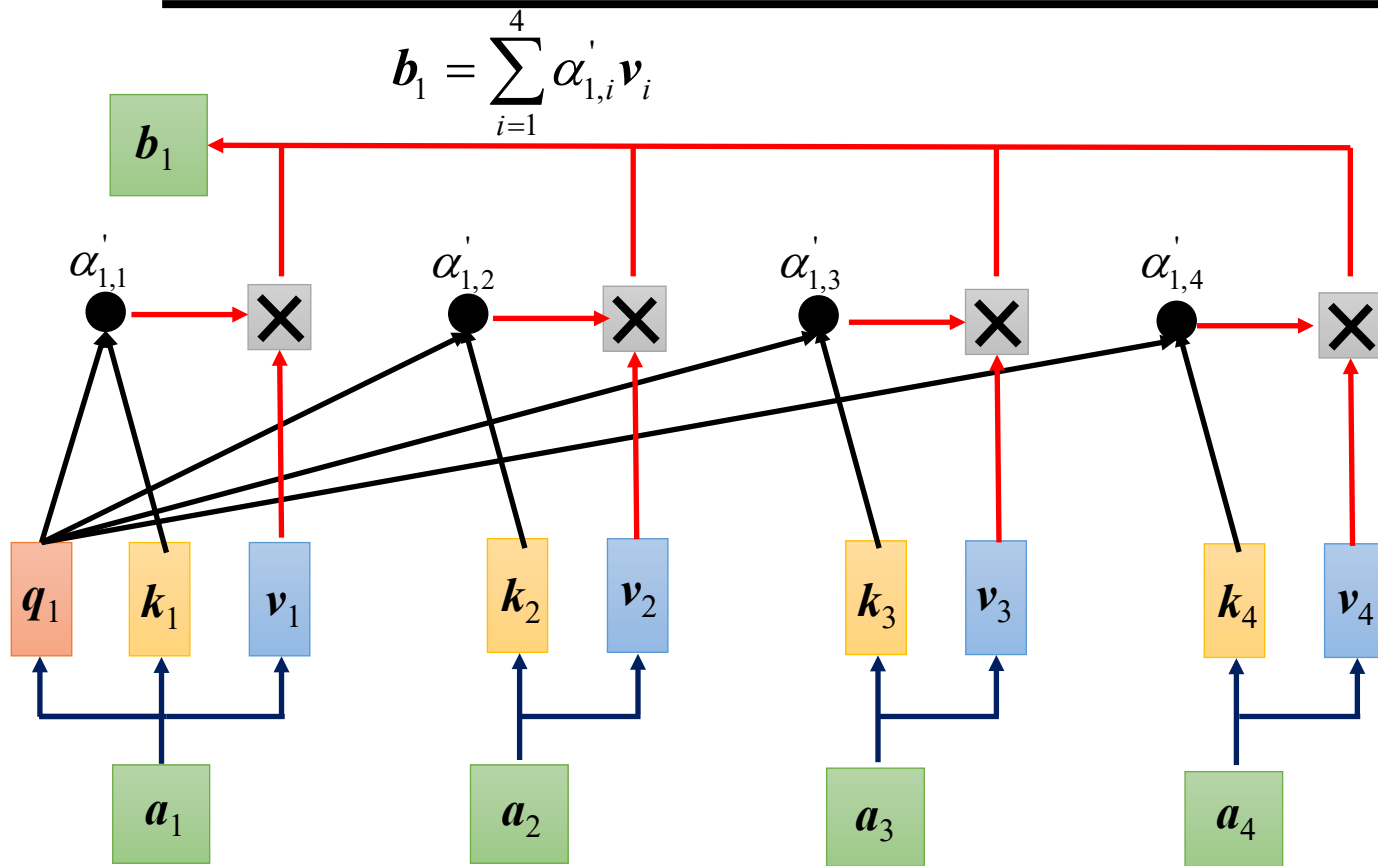
$$q_1 \in \mathbb{R}^{d \times 1}$$

$$k_1, k_2, k_3, k_4 \in \mathbb{R}^{d \times 1}$$

$$v_1, v_2, v_3, v_4 \in \mathbb{R}^{d \times 1}$$



Self-attention



$$q_1 = W_q a_1$$

$$k_1 = W_k a_1$$

$$k_2 = W_k a_2$$

$$k_3 = W_k a_3$$

$$k_4 = W_k a_4$$

$$v_1 = W_v a_1$$

$$v_2 = W_v a_2$$

$$v_3 = W_v a_3$$

$$v_4 = W_v a_4$$

$$\alpha_{1,1} = k_1 \cdot q_1$$

$$\alpha_{1,2} = k_2 \cdot q_1$$

$$\alpha_{1,3} = k_3 \cdot q_1$$

$$\alpha_{1,4} = k_4 \cdot q_1$$

$$\xrightarrow{\text{softmax}} \alpha'_{1,i} = \frac{\exp(\alpha_{1,i} / \sqrt{d})}{\sum_{j=1}^4 \exp(\alpha_{1,j} / \sqrt{d})}$$

$$\alpha_1 \triangleq \begin{pmatrix} \alpha_{1,1} \\ \alpha_{1,2} \\ \alpha_{1,3} \\ \alpha_{1,4} \end{pmatrix}$$

attention scores

softmax

$$\alpha'_1 \triangleq \begin{pmatrix} \alpha'_{1,1} \\ \alpha'_{1,2} \\ \alpha'_{1,3} \\ \alpha'_{1,4} \end{pmatrix}$$

normalized
attention scores

b_1 is the average of $\{v_i\}$ weighted by α'_1 ,

$$b_1 = \sum_{i=1}^4 \alpha'_{1,i} v_i = [v_1 \ v_2 \ v_3 \ v_4] \alpha'_1$$

$b_2 \sim b_4$ can be computed in the same way in parallel



Self-attention

$$q_1 = W_q a_1 \quad q_2 = W_q a_2 \quad q_3 = W_q a_3 \quad q_4 = W_q a_4$$

$$v_1 = W_v a_1 \quad v_2 = W_v a_2 \quad v_3 = W_v a_3 \quad v_4 = W_v a_4$$

$$k_1 = W_k a_1 \quad k_2 = W_k a_2 \quad k_3 = W_k a_3 \quad k_4 = W_k a_4$$

$$\alpha_2 \triangleq \begin{pmatrix} \alpha_{2,1} = k_1 \cdot q_2 \\ \alpha_{2,2} = k_2 \cdot q_2 \\ \alpha_{2,3} = k_3 \cdot q_2 \\ \alpha_{2,4} = k_4 \cdot q_2 \end{pmatrix} \xrightarrow{\text{softmax}} \alpha'_2 \triangleq \begin{pmatrix} \alpha'_{2,1} \\ \alpha'_{2,2} \\ \alpha'_{2,3} \\ \alpha'_{2,4} \end{pmatrix}$$



$$b_2 = \sum_{i=1}^4 \alpha'_{2,i} v_i = [v_1 \ v_2 \ v_3 \ v_4] \alpha'_2$$

$$\alpha_3 \triangleq \begin{pmatrix} \alpha_{3,1} = k_1 \cdot q_3 \\ \alpha_{3,2} = k_2 \cdot q_3 \\ \alpha_{3,3} = k_3 \cdot q_3 \\ \alpha_{3,4} = k_4 \cdot q_3 \end{pmatrix} \xrightarrow{\text{softmax}} \alpha'_3 \triangleq \begin{pmatrix} \alpha'_{3,1} \\ \alpha'_{3,2} \\ \alpha'_{3,3} \\ \alpha'_{3,4} \end{pmatrix}$$



$$b_3 = \sum_{i=1}^4 \alpha'_{3,i} v_i = [v_1 \ v_2 \ v_3 \ v_4] \alpha'_3$$

$$\alpha_4 \triangleq \begin{pmatrix} \alpha_{4,1} = k_1 \cdot q_4 \\ \alpha_{4,2} = k_2 \cdot q_4 \\ \alpha_{4,3} = k_3 \cdot q_4 \\ \alpha_{4,4} = k_4 \cdot q_4 \end{pmatrix} \xrightarrow{\text{softmax}} \alpha'_4 \triangleq \begin{pmatrix} \alpha'_{4,1} \\ \alpha'_{4,2} \\ \alpha'_{4,3} \\ \alpha'_{4,4} \end{pmatrix}$$



$$b_4 = \sum_{i=1}^4 \alpha'_{4,i} v_i = [v_1 \ v_2 \ v_3 \ v_4] \alpha'_4$$

The whole process can be represented in matrix form



Self-attention

Let $I_{d \times 4} \triangleq [a_1 \ a_2 \ a_3 \ a_4]$, $Q_{d \times 4} \triangleq [q_1 \ q_2 \ q_3 \ q_4]$, $K_{d \times 4} \triangleq [k_1 \ k_2 \ k_3 \ k_4]$, $V_{d \times 4} \triangleq [v_1 \ v_2 \ v_3 \ v_4]$, $A_{4 \times 4} \triangleq [\alpha_1 \ \alpha_2 \ \alpha_3 \ \alpha_4]$, $A'_{4 \times 4} = [\alpha'_1 \ \alpha'_2 \ \alpha'_3 \ \alpha'_4]$ and $O_{d \times 4} \triangleq [b_1 \ b_2 \ b_3 \ b_4]$

We have,

$$Q_{d \times 4} = [q_1 \ q_2 \ q_3 \ q_4] = [W_q a_1 \ W_q a_2 \ W_q a_3 \ W_q a_4] = W_q I_{d \times 4}$$

$$K_{d \times 4} = [k_1 \ k_2 \ k_3 \ k_4] = [W_k a_1 \ W_k a_2 \ W_k a_3 \ W_k a_4] = W_k I_{d \times 4}$$

$$V_{d \times 4} = [v_1 \ v_2 \ v_3 \ v_4] = [W_v a_1 \ W_v a_2 \ W_v a_3 \ W_v a_4] = W_v I_{d \times 4}$$

$$\alpha_1 = \begin{pmatrix} \alpha_{1,1} = k_1 \cdot q_1 \\ \alpha_{1,2} = k_2 \cdot q_1 \\ \alpha_{1,3} = k_3 \cdot q_1 \\ \alpha_{1,4} = k_4 \cdot q_1 \end{pmatrix} = \begin{pmatrix} k_1^T q_1 \\ k_2^T q_1 \\ k_3^T q_1 \\ k_4^T q_1 \end{pmatrix} = K^T q_1 \quad \alpha_2 = K^T q_2 \quad \alpha_3 = K^T q_3 \quad \alpha_4 = K^T q_4$$

softmax is applied to each column vector

$$A_{4 \times 4} = [\alpha_1 \ \alpha_2 \ \alpha_3 \ \alpha_4] = [K^T q_1 \ K^T q_2 \ K^T q_3 \ K^T q_4] = K^T [q_1 \ q_2 \ q_3 \ q_4] = K^T Q$$

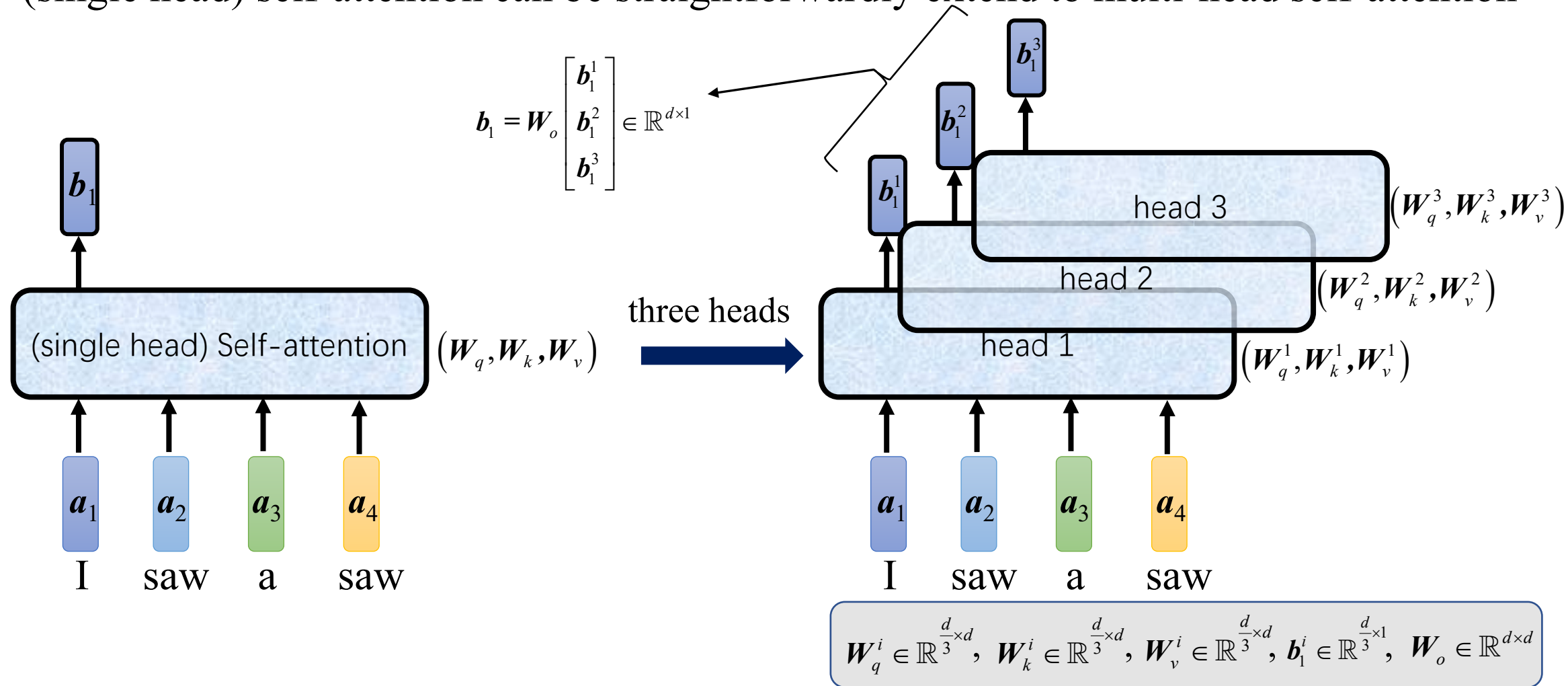
$$A'_{4 \times 4} = [\alpha'_1 \ \alpha'_2 \ \alpha'_3 \ \alpha'_4] = \left[\text{softmax}\left(\frac{\alpha_1}{\sqrt{d}}\right) \text{softmax}\left(\frac{\alpha_2}{\sqrt{d}}\right) \text{softmax}\left(\frac{\alpha_3}{\sqrt{d}}\right) \text{softmax}\left(\frac{\alpha_4}{\sqrt{d}}\right) \right] = \text{softmax}\left(\frac{A}{\sqrt{d}}\right) = \text{softmax}\left(\frac{K^T Q}{\sqrt{d}}\right)$$

$$O = [b_1 \ b_2 \ b_3 \ b_4] = [V_{d \times 4} \alpha'_1 \ V_{d \times 4} \alpha'_2 \ V_{d \times 4} \alpha'_3 \ V_{d \times 4} \alpha'_4] = V_{d \times 4} [\alpha'_1 \ \alpha'_2 \ \alpha'_3 \ \alpha'_4] = V_{d \times 4} A'_{4 \times 4}$$



Multi-head Self-attention (MSA)

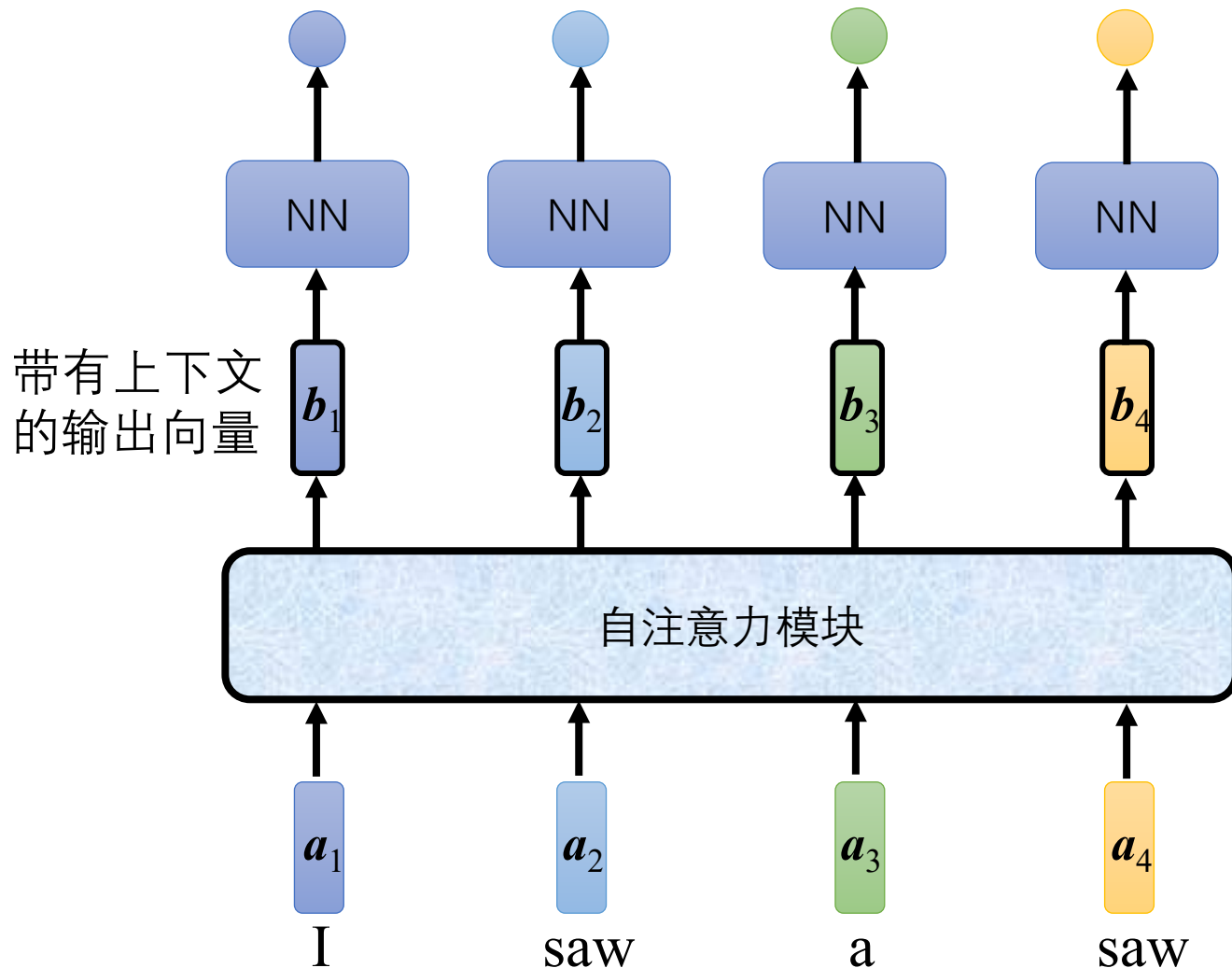
(single head) self-attention can be straightforwardly extend to multi-head self-attention





Positional encoding

- Can the self-attention really solve our word-tagging problem?



- ✓ Using the self-attention, actually b_2 and b_4 are the same!
- ✓ Accordingly, the two “saw”s will be predicted to have the same tagging

What do we miss?

The position information of the vectors in the input sequence

That the two “saw”s have different tagging largely owes to the fact that they have different positions in the sequence



We need to modify the input vectors by embedding their positional information



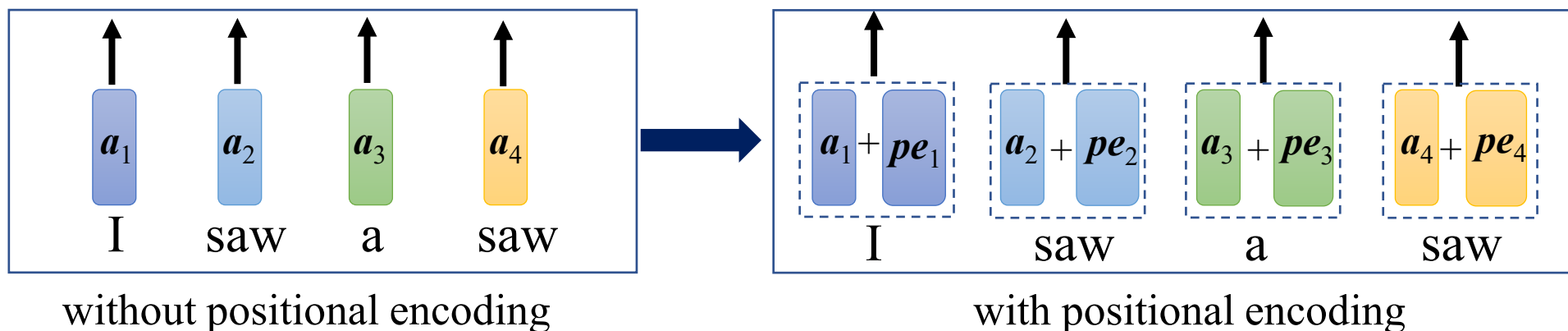
Positional encoding

For each input vector $\mathbf{a}_t \in \mathbb{R}^{d \times 1}$, construct a positional encoding vector $\mathbf{pe}_t \triangleq \left\{ pe_t^{(i)} \right\}_{i=0}^{d-1}$

$$pe_t^{(i)} = \begin{cases} \sin(w_k t), & \text{if } i = 2k \\ \cos(w_k t), & \text{if } i = 2k + 1 \end{cases}$$

where t is the position of this input vector in the sequence, $w_k = \frac{1}{10000^{2k/d}}$, $k=0, 1, 2, \dots, d/2-1$

Positional encoding means that we modify \mathbf{a}_t as $\mathbf{a}_t + \mathbf{pe}_t$





Layer-norm VS batch-norm

- Batch-norm and layer-norm are two strategies for training neural networks faster and more stably
- In Yolov3, we have met batch-norm; in transformer-related structures, in most cases, they use layer-norm
- Batch-norm VS layer-norm
 - **Batch-norm normalizes each feature (channel) independently across the samples in a mini-batch**
 - **Layer-norm normalizes each sample in the mini-batch independently across all features (channels)**
 - As batch-norm is dependent on batch size, it's not effective for small batch sizes

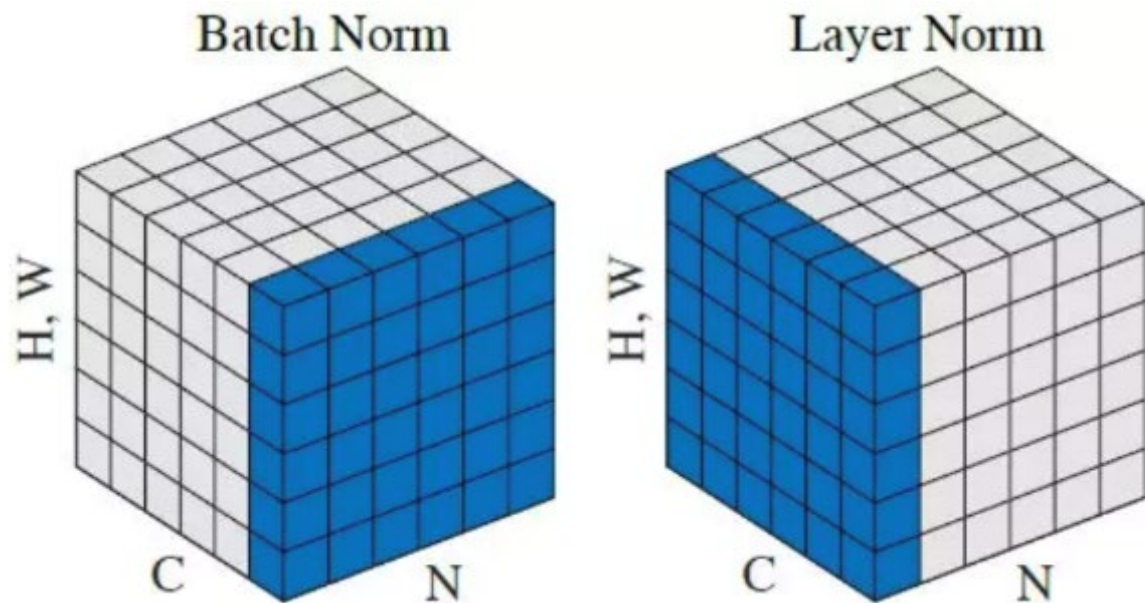


Layer-norm VS batch-norm

They use the same updating formular, but adopt different ways to compute statistics (μ , σ^2),

$$y_i \leftarrow \gamma \frac{x_i - \mu}{\sqrt{\sigma^2 + \varepsilon}} + \beta$$

where γ and β are learnable parameters; for batch-norm, each neuron (channel) has a (γ, β) pair while for layer-norm each layer has a (γ, β) pair

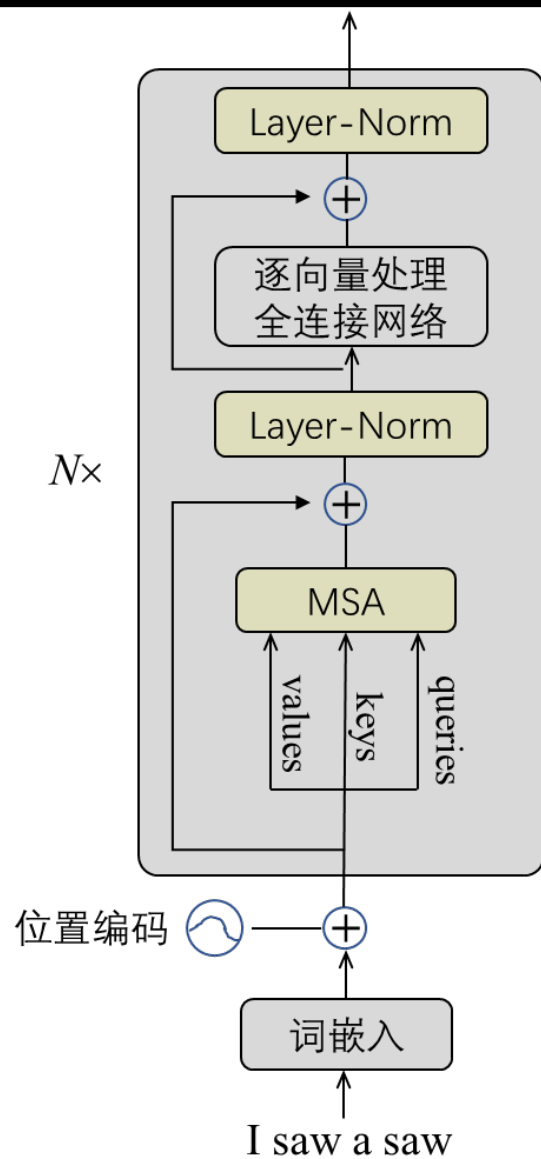


- ✓ C means the feature channels, $H \times W$ defines the instances of a feature by a sample, and N means the number of samples in a mini-batch
- ✓ The normalization is applied to the blue part

For our case, $C=d$, $H \times W=4$, $N=1$



Transformer encoder



- ✓ The inputs are embedded in a sequence of vectors with the same length
- ✓ Then, the embedded vectors are modified with positional encoding
- ✓ Transformer encoder is composed of N blocks with the same structures
- ✓ Each block processes the input vectors by a multi-head self-attention layer, a residual connection layer, a layer-norm operation, a vector-wise fully connected MLP, another residual connection layer, and another layer-norm
- ✓ If the dimension of the input is $d \times N$, the output of a transformer encoder will also be $d \times N$



Outline

- Transformer in NLP
- Multi-head Attention
- Vision transform (ViT)
- Swin-Transformer
- DEtection TRansformer (DETR)
- Real-time DETR (RT-DETR)



Multi-head Attention

Self-attention

Let $I_{d \times 4} \triangleq [a_1 \ a_2 \ a_3 \ a_4]$, $Q_{d \times 4} \triangleq [q_1 \ q_2 \ q_3 \ q_4]$, $K_{d \times 4} \triangleq [k_1 \ k_2 \ k_3 \ k_4]$, $V_{d \times 4} \triangleq [v_1 \ v_2 \ v_3 \ v_4]$, $A_{4 \times 4} \triangleq [\alpha_1 \ \alpha_2 \ \alpha_3 \ \alpha_4]$, $A'_{4 \times 4} = [\alpha'_1 \ \alpha'_2 \ \alpha'_3 \ \alpha'_4]$ and $O_{d \times 4} \triangleq [b_1 \ b_2 \ b_3 \ b_4]$

$$Q_{d \times 4} = [q_1 \ q_2 \ q_3 \ q_4] = [W_q a_1 \ W_q a_2 \ W_q a_3 \ W_q a_4] = W_q I_{d \times 4}$$

query sequence

$$K_{d \times 4} = [k_1 \ k_2 \ k_3 \ k_4] = [W_k a_1 \ W_k a_2 \ W_k a_3 \ W_k a_4] = W_k I_{d \times 4}$$

key sequence

$$V_{d \times 4} = [v_1 \ v_2 \ v_3 \ v_4] = [W_v a_1 \ W_v a_2 \ W_v a_3 \ W_v a_4] = W_v I_{d \times 4}$$

value sequence

In self-attention, the query sequence, the key sequence and the value sequence are actually identical; that is why it is called **self**-attention.

If the key sequence and the value sequence are the same while the query sequence is different, the self-attention changes to **attention**. In other words, self-attention is a special case of attention

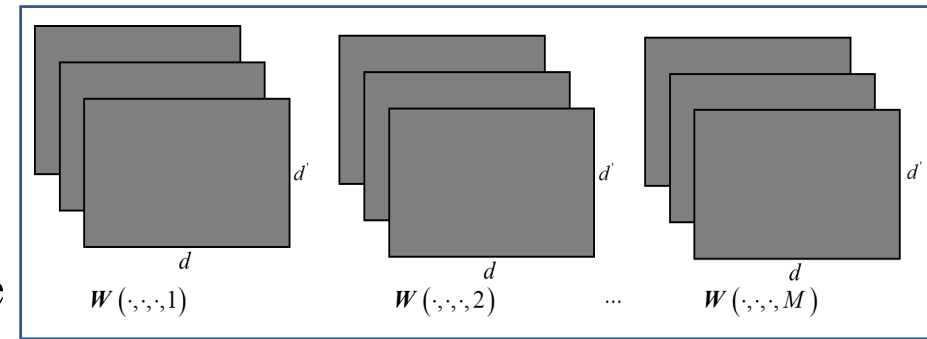


Multi-head Attention

- Multi-head attention is an extension to the multi-head self-attention; their computation frameworks are the same, except that **in multi-head attention, the query sequence is different from the key and value sequences**
- Multi-head attention can be used in **transformer decoders**

A multi-head attention module with M heads can be expressed as a function mh-attn,

$$\text{mh-attn} \left(\underbrace{\mathbf{I}_q}_{d \times N_q}, \underbrace{\mathbf{I}_{kv}}_{d \times N_{kv}}, \underbrace{\mathbf{W}}_{d' \times d \times 3 \times M}, \underbrace{\mathbf{W}_o}_{d \times d} \right) \mapsto \underbrace{\mathbf{O}}_{d \times N_q}$$



$\mathbf{I}_q \in \mathbb{R}^{d \times N_q}$ is the query sequence; $\mathbf{I}_{kv} \in \mathbb{R}^{d \times N_{kv}}$ is the key/value sequence

$\mathbf{W} \in \mathbb{R}^{d' \times d \times 3 \times M}$ is the weight tensor, where $d' = \frac{d}{M}$ is the dimension of the embedded vectors in each single-head

$\mathbf{W}_o \in \mathbb{R}^{d \times d}$ linearly maps the concatenation of the outputs of single-heads to a space of dimension d

$\mathbf{O} \in \mathbb{R}^{d \times N_q}$ is the final output

Remember: each query generates an output vector



Multi-head Attention

Multi-head attention

$$\text{mh-attn} \left(\underbrace{\mathbf{I}_q}_{d \times N_q}, \underbrace{\mathbf{I}_{kv}}_{d \times N_{kv}}, \underbrace{\mathbf{W}}_{d' \times d \times 3 \times M}, \underbrace{\mathbf{W}_o}_{d \times d} \right) \mapsto \underbrace{\mathbf{O}}_{d \times N_q}$$

Similar as multi-head self-attention, the final output of a MHA module is generated by

- ✓ Concatenating the outputs of all the single-heads
- ✓ then perform a linear mapping using \mathbf{W}_o

$$\mathbf{O}'_{d \times N_q} = [\text{attn}(\mathbf{I}_q, \mathbf{I}_{kv}, \mathbf{W}_1); \text{attn}(\mathbf{I}_q, \mathbf{I}_{kv}, \mathbf{W}_2); \dots; \text{attn}(\mathbf{I}_q, \mathbf{I}_{kv}, \mathbf{W}_M)]$$

$$\mathbf{O}_{d \times N_q} = \mathbf{W}_o \mathbf{O}'_{d \times N_q}$$

where $\text{attn}(\mathbf{I}_q, \mathbf{I}_{kv}, \mathbf{W}_h) \in \mathbb{R}^{d' \times N_q}$ is the output of the h^{th} single-head and $\mathbf{W}_h \in \mathbb{R}^{d' \times d \times 3}$ is the h^{th} “slice” of the tensor \mathbf{W} ; $[\cdot]$ denotes the channel-wise concatenation

Single-head

$$\text{attn}(\mathbf{I}_q, \mathbf{I}_{kv}, \mathbf{W}') \in \mathbb{R}^{d' \times N_q}$$

where $\mathbf{W}' = [\mathbf{W}'_1; \mathbf{W}'_2; \mathbf{W}'_3] \in \mathbb{R}^{d' \times d \times 3}$ is the weight tensor for this head; $\mathbf{W}'_1 \in \mathbb{R}^{d' \times d}$, $\mathbf{W}'_2 \in \mathbb{R}^{d' \times d}$, $\mathbf{W}'_3 \in \mathbb{R}^{d' \times d}$ are used to compute embedded queries, keys, and values, respectively

$$\mathbf{Q}_{d' \times N_q} = \mathbf{W}'_1 (\mathbf{I}_q + \mathbf{P}_q), \mathbf{K}_{d' \times N_{kv}} = \mathbf{W}'_2 (\mathbf{I}_{kv} + \mathbf{P}_{kv}), \mathbf{V}_{d' \times N_{kv}} = \mathbf{W}'_3 \mathbf{I}_{kv}$$

Compute the normalized attention scores between \mathbf{q}_i and \mathbf{K} ,

$$\boldsymbol{\alpha}_i \triangleq (\alpha_{i,1}, \alpha_{i,2}, \dots, \alpha_{i,N_{kv}})^T, \quad \alpha_{i,j} = \frac{\exp(\mathbf{k}_j^T \mathbf{q}_i / \sqrt{d'})}{\sum_{k=1}^{N_{kv}} \exp(\mathbf{k}_k^T \mathbf{q}_i / \sqrt{d'})}$$

The i -th output vector of this head is,

$$\text{attn}_i(\mathbf{I}_q, \mathbf{I}_{kv}, \mathbf{W}') = \sum_{j=1}^{N_{kv}} \alpha_{ij} \mathbf{v}_j = \mathbf{V} \boldsymbol{\alpha}_i \in \mathbb{R}^{d' \times 1}$$



Outline

- Transformer in NLP
- Multi-head Attention
- Vision transform (ViT)
- Swin-Transformer
- DEtection TRansformer (DETR)
- Real-time DETR (RT-DETR)



Vision Transformer (ViT)

- Inspired by the great success of transformer in the field of NLP, researchers in CV have started to adopt transformers in CV tasks
- The first image classification framework totally based on transformer is ViT (Vision Transformer) proposed also by researchers of Google in 2021

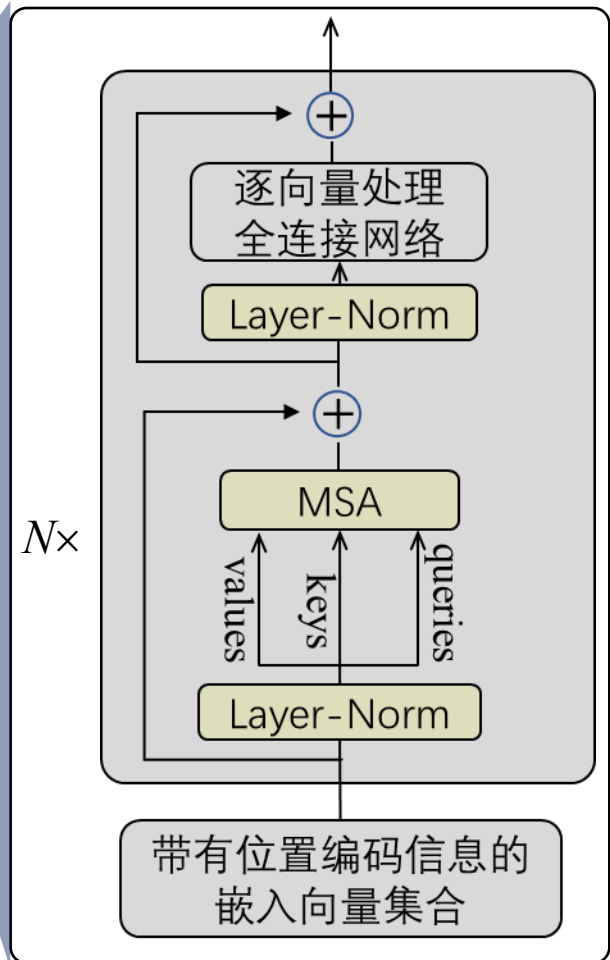
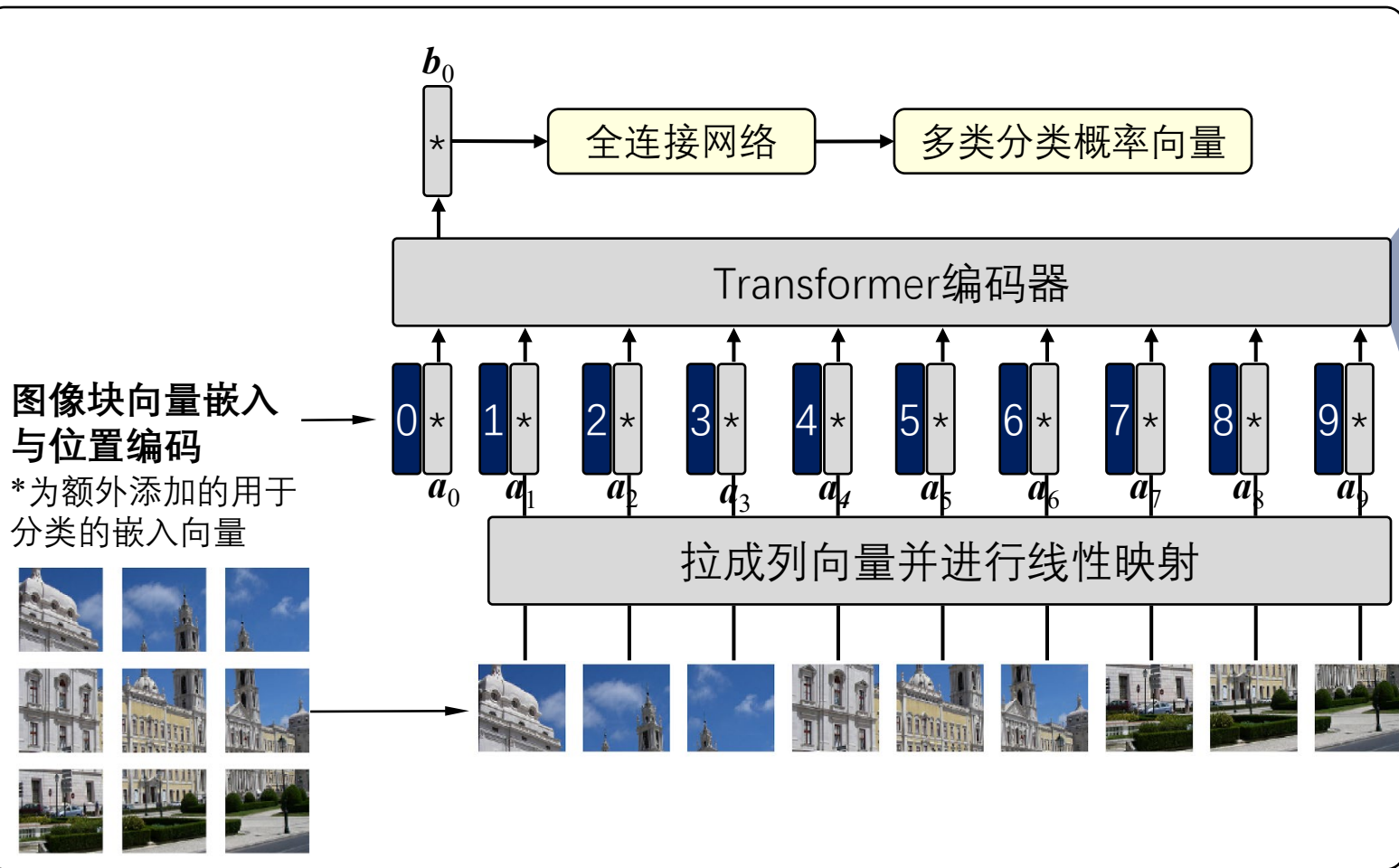


Alexey Dosovitskiy received the M.Sc. and Ph.D. degrees in mathematics (functional analysis) from Moscow State University, Moscow, Russia, in 2009 and 2012, respectively. He is currently a Research Scientist with the Intelligent Systems Laboratory, Intel, Munich, Germany. From 2013 to 2016, he was a Postdoctoral Researcher, with Prof. T. Brox, with the Computer Vision Group, University of Freiburg, Breisgau, Germany, working on various topics in deep learning, including self-supervised learning, image generation with neural networks, motion, and 3-D structure estimation.

[1] DOSOVITSKIY A, BEYER L, KOLESNIKOV A, et al. An image is worth 16x16 words: Transformers for image recognition at scale[C]//Proc. Int'l. Conf. Learning Representations, 2021. (**Cited by 38008**, Jun. 25, 2024)



Vision Transformer (ViT)

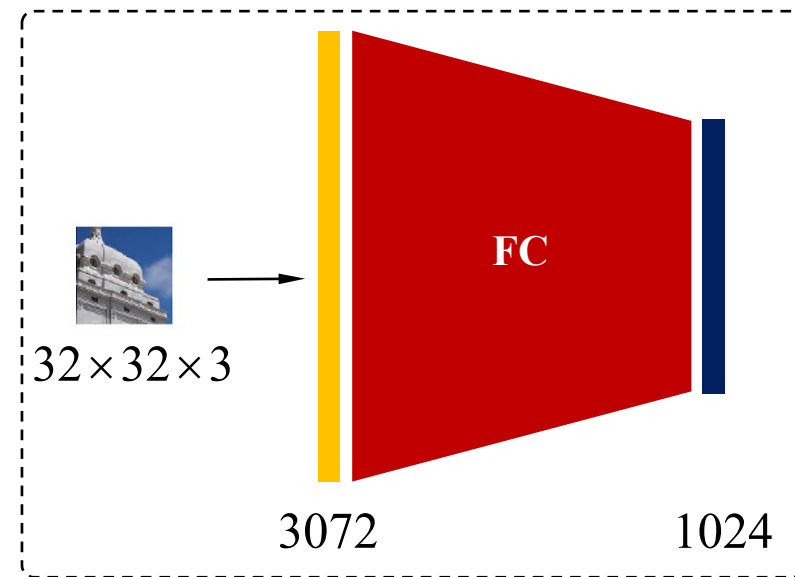




Vision Transformer (ViT)



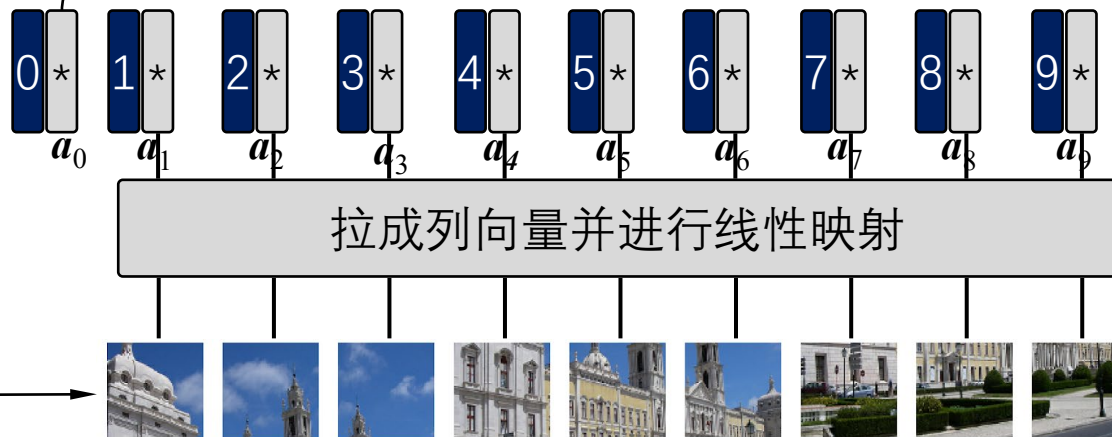
视觉Transformer



Vision Transformer (ViT)

图像块向量嵌入
与位置编码

*为额外添加的用于
分类的嵌入向量



视觉Transformer

✓ A learnable vector \mathbf{a}_0 is added here and is expected to hold the global information of the image for classification; \mathbf{a}_0 can be randomly initialized

With \mathbf{a}_0 , the input vectors can be seen as a matrix,

$$\mathbf{A} \triangleq [\mathbf{a}_0 \ \mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_9] \in \mathbb{R}^{1024 \times 10}$$

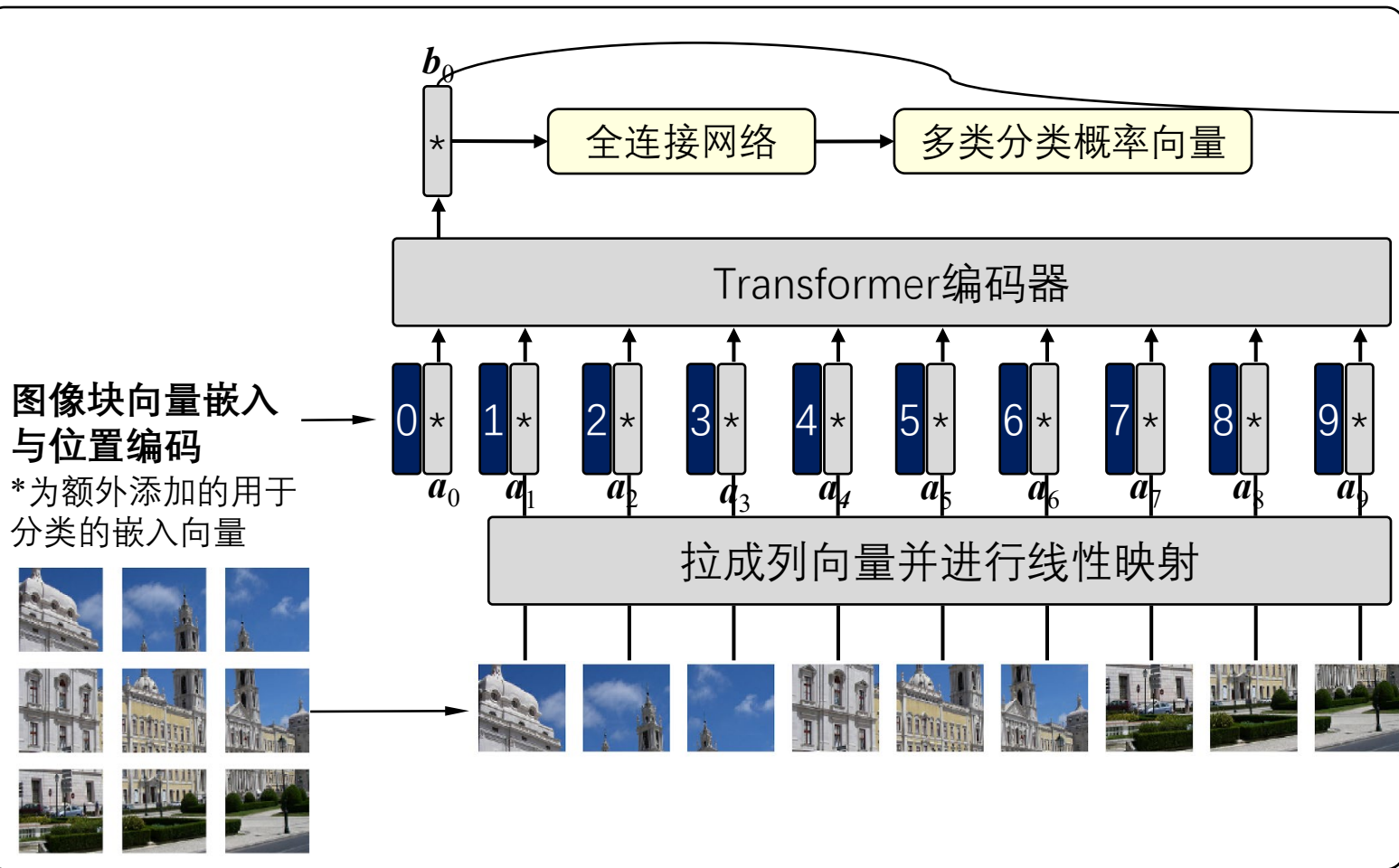
For positional encoding, ViT uses a learnable matrix $\mathbf{PE} \in \mathbb{R}^{1024 \times 10}$

Then, the input vectors with positional encoding form the matrix

$$\mathbf{A} + \mathbf{PE} \in \mathbb{R}^{1024 \times 10}$$



Vision Transformer (ViT)



✓ b_0 corresponds to a_0 ; as it is the output of a transformer encoder, it can hold the global information of all the block vectors and can be trained to perform the classification task

视觉Transformer



Outline

- Transformer in NLP
- Multi-head Attention
- Vision transform (ViT)
- Swin-Transformer
- DEtection TRansformer (DETR)
- Real-time DETR (RT-DETR)



Swin-Transformer

- ViT is a cool idea and demonstrates that transformer has a good potential in CV; however, it has too many parameters and thus is computationally expensive
- Swin-Transformer^[1] proposes several universal and straightforward ideas to improve ViT
 - To reduce computational cost: **self-attention is restricted to a local window**
 - To make the outputs of the self-attention still capture the contextual information: **windows used to compute self-attention are shifted and the self-attention is repeated**; the shifted window partitioning approach introduces connections between neighboring non-overlapping windows in the previous layer
 - To get multi-resolution feature maps (as Yolov3 and Yolov8): **a patch merging mechanism is proposed**; the obtained hierarchical representation is very similar to multi-resolution CNN feature maps

[1] LIU Z, LIN Y, CAO Y, et al. Swin transformer: Hierarchical vision transformer using shifted windows[C]//Proc. IEEE Int'l. Conf. Computer Vision, 2021: 9992-10002. (**Cited by 19016**, July 4, 2024)



Swin-Transformer

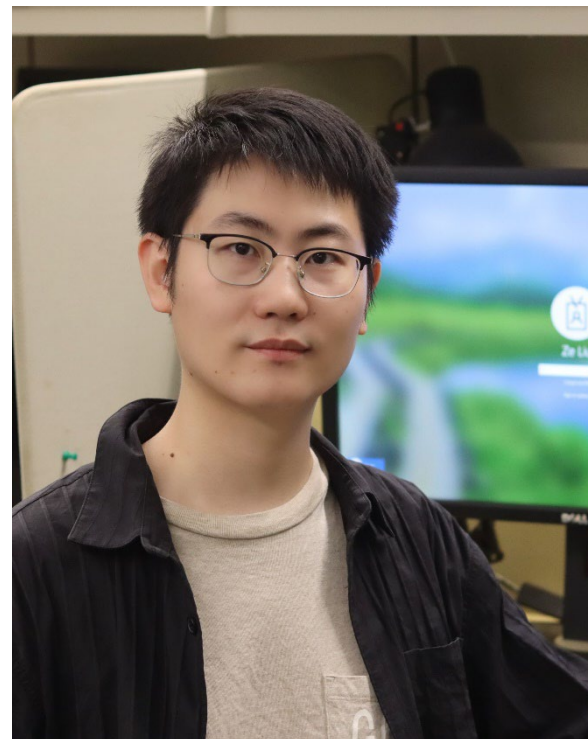
Marr Prize

**Swin Transformer: Hierarchical Vision Transformer
using Shifted Windows**

Ze Liu (USTC), Yutong Lin (Xi'an Jiaotong University),
Yue Cao (Microsoft Research), Han Hu (Microsoft Research Asia),
Yixuan Wei (Tsinghua University), Zheng Zhang (MSRA, Huazhong University of
Science and Technology), Stephen Lin (Microsoft Research),
Baining Guo (MSR Asia)

Session 8 (A/B)

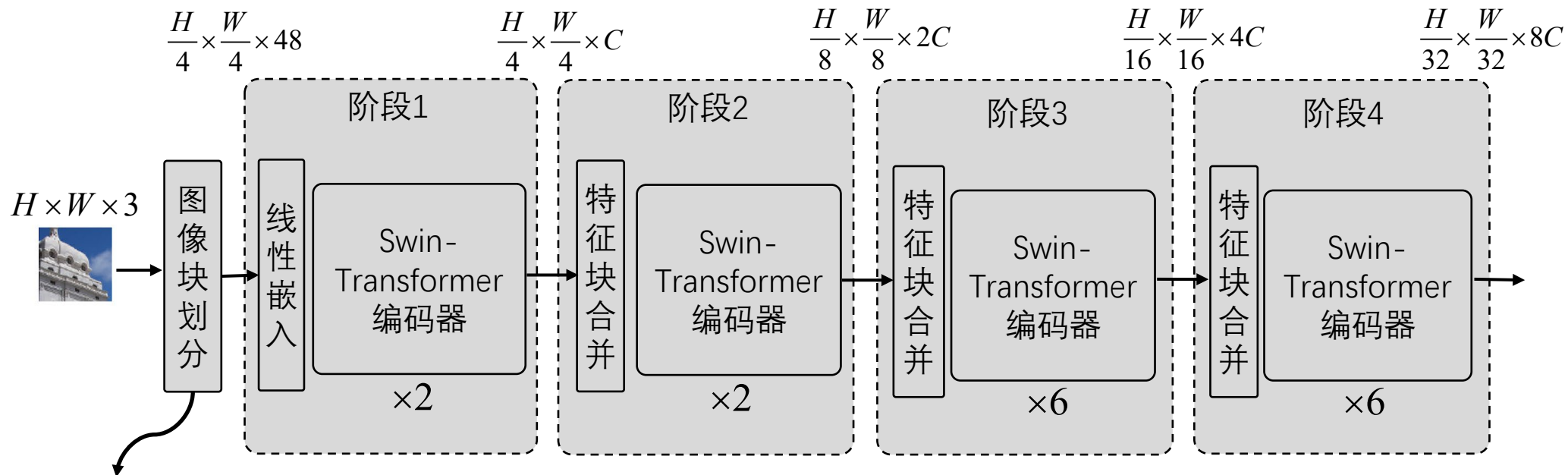
2021 **ICCV** OCTOBER 11-17
VIRTUAL



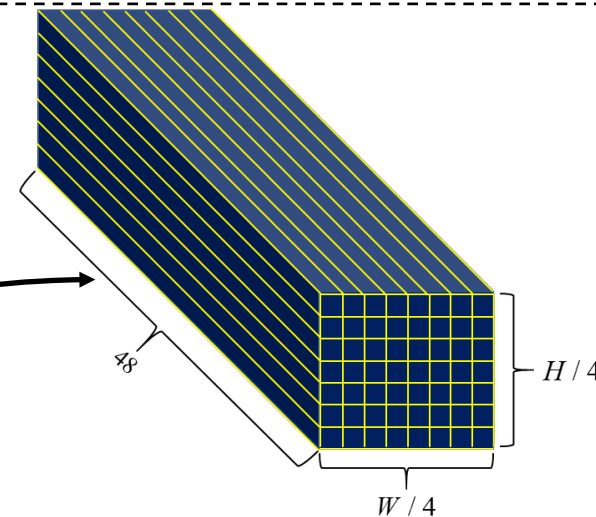
Ze Liu is currently a final-year joint Ph.D. candidate at the University of Science and Technology of China (USTC) and Microsoft Research Asia (MSRA)



Swin-Transformer

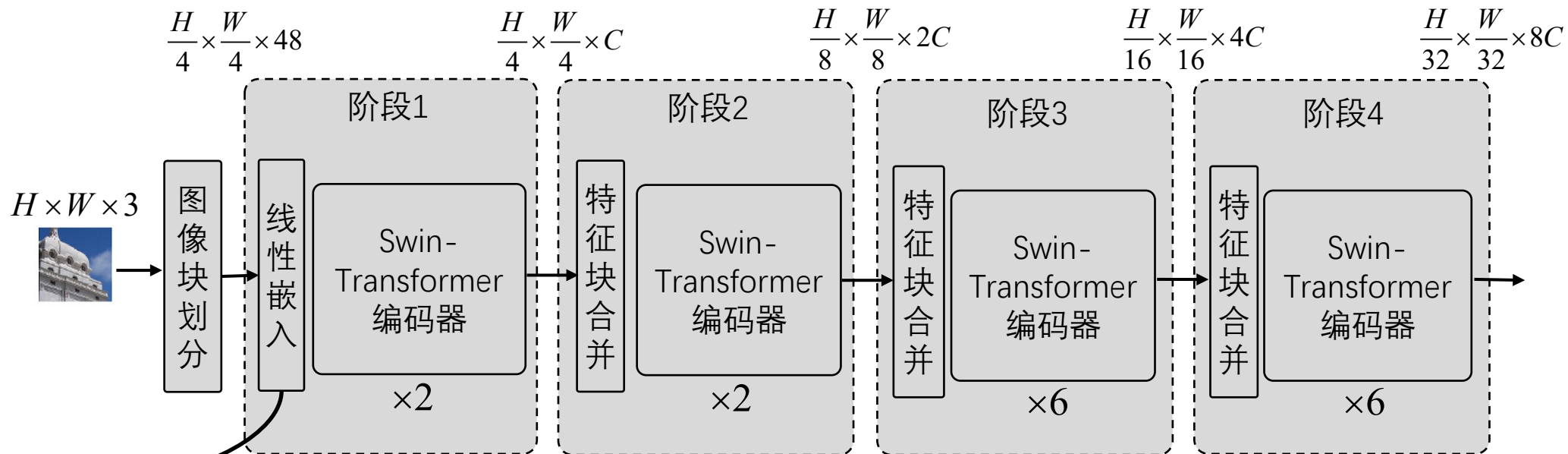


- ✓ Each patch is of the size $4 \times 4 \times 3$, and is stacked into a 48- d vector
- ✓ After the patch partition, the input to the “stage 1” is a feature map $V_0 \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times 48}$

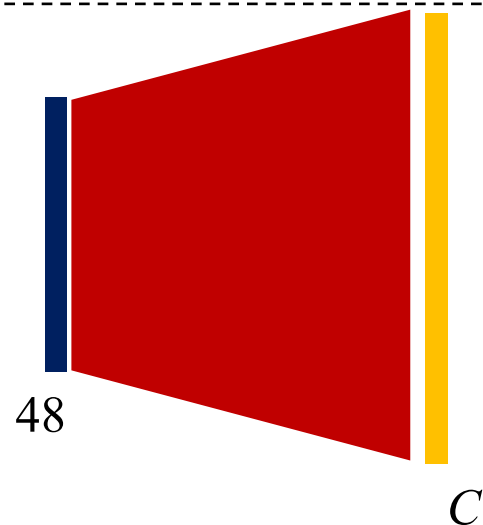




Swin-Transformer

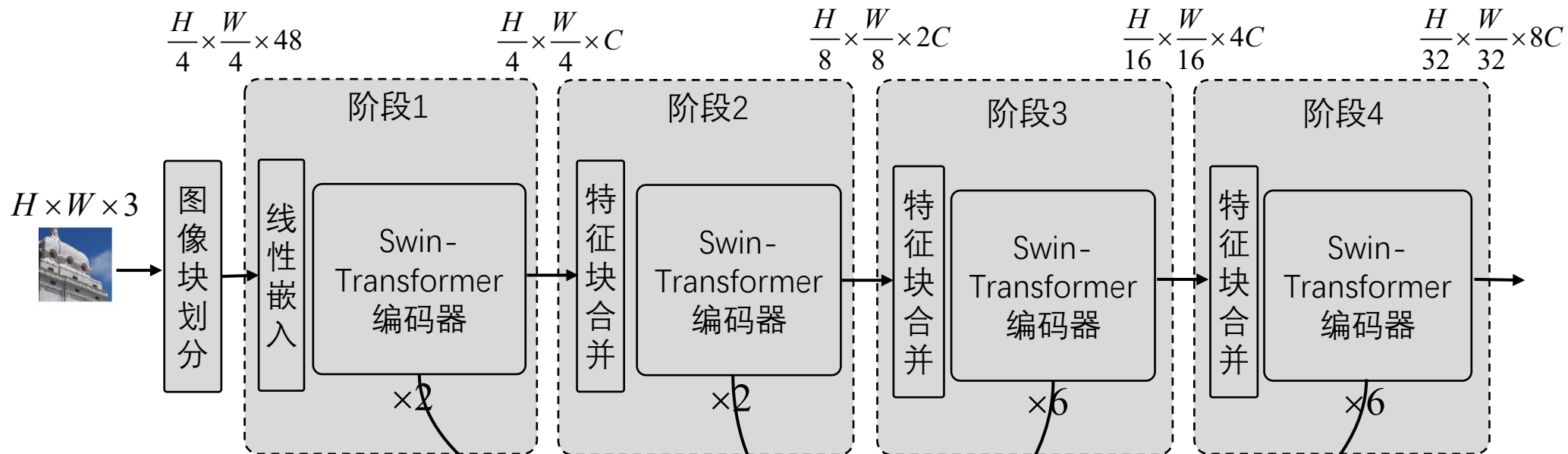


- ✓ Using a fully-connect network to map each 48-d vector in V_0 to a C -d vector
- ✓ After linear embedding, the feature map becomes $V_0^1 \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times C}$
- ✓ Since the swin-transformer encoder does not change the dimension, the output dimension of stage 1 also is $\frac{H}{4} \times \frac{W}{4} \times C$

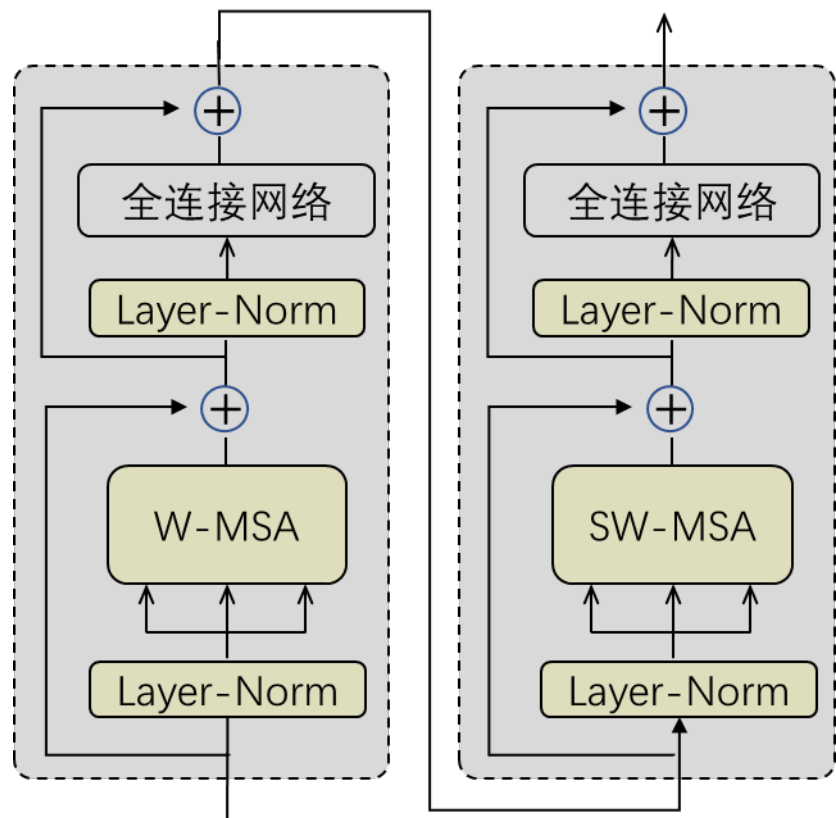




Swin-Transformer

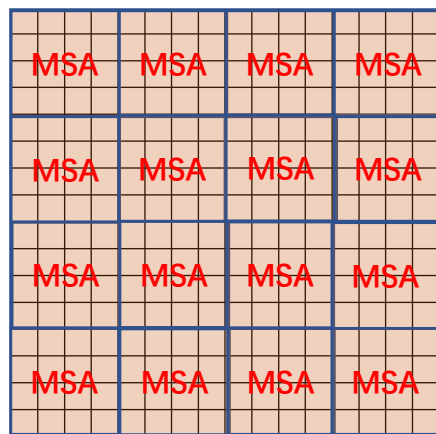


Next page

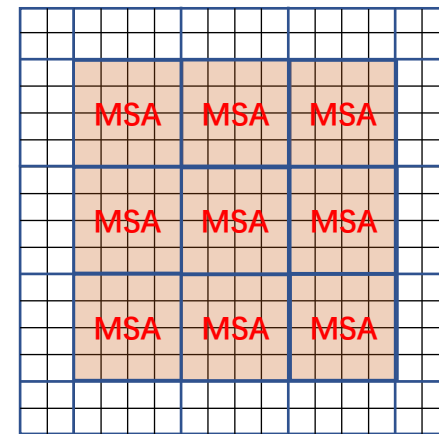


Swin-Transformer Encoders

- ✓ Two swin-transformer encoders are used in pair
- ✓ With respect to structures, they are similar to the transformer encoder of ViT; the only difference is that the swin-transformer encoders compute multi-head self-attention in local windows
- ✓ The windows used in SW-MSA are shifted from windows in W-MSA; the window-shifting mechanism can supplement connections across non-overlapping windows



Window-partition in W-MSA

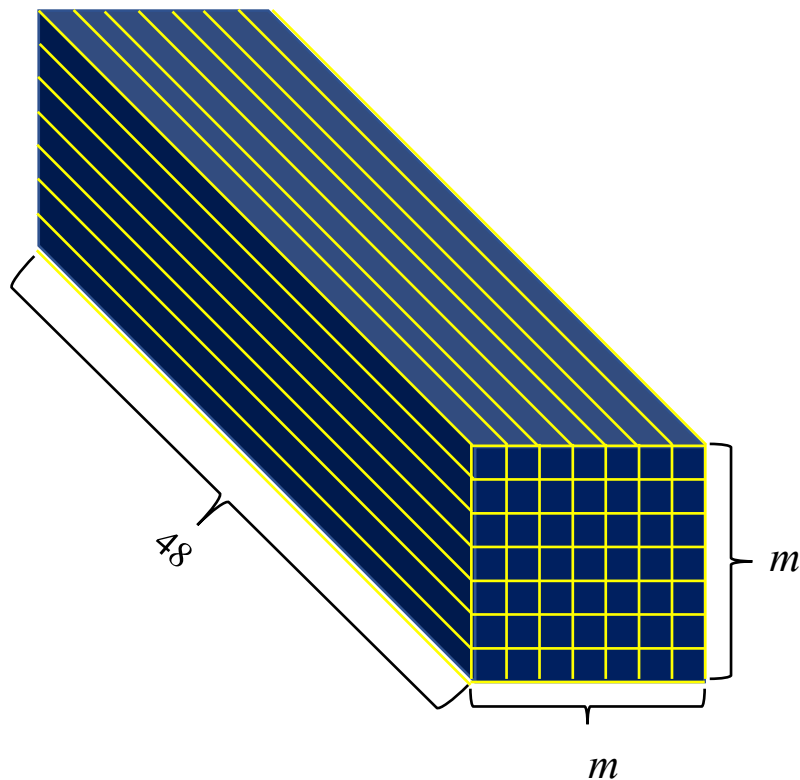


Window-partition in SW-MSA

Note: Each point actually is a vector



Swin-Transformer



How to perform positional encoding?

Using a common positional encoding strategy,

$$\left\{ \begin{array}{c} 48 \\ \left[\begin{array}{c} | \\ | \\ | \\ \vdots \\ | \end{array} \right] \end{array} \right\} + \left\{ \begin{array}{c} \left[\begin{array}{ccc} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{array} \right] \\ 0 \ 1 \ 2 \dots m^2-1 \end{array} \right\} \quad \text{(PE matrix)}$$

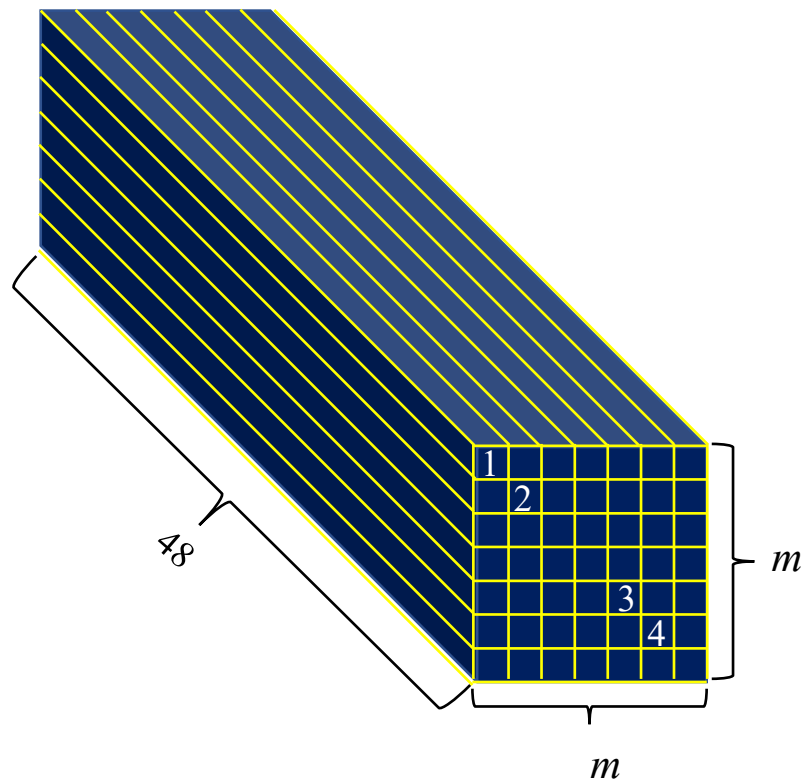
48 \rightarrow Input with PE

- ✓ PE matrix can be pre-fixed (such as sinusoid functions) or can be learned as in the case of ViT
- ✓ Such a positional encoding strategy can be considered as a “1-D” strategy since the vectors are arranged in a line
- ✓ **The “2D” positional relationships of the vectors cannot be well embedded in such a way**

\rightarrow Relative Position Bias



Swin-Transformer



- Relative position bias can encode the relative position relationship between two vectors in a 2D array
E.g., the RPB of “2” to “1” is the same as the RPB of “4” to “3”
- Positional encodings are added to the input vectors; differently, RPBs are added to attention scores

Embedded queries and keys: $\mathbf{Q} \in \mathbb{R}^{d \times m^2}$, $\mathbf{K} \in \mathbb{R}^{d \times m^2}$

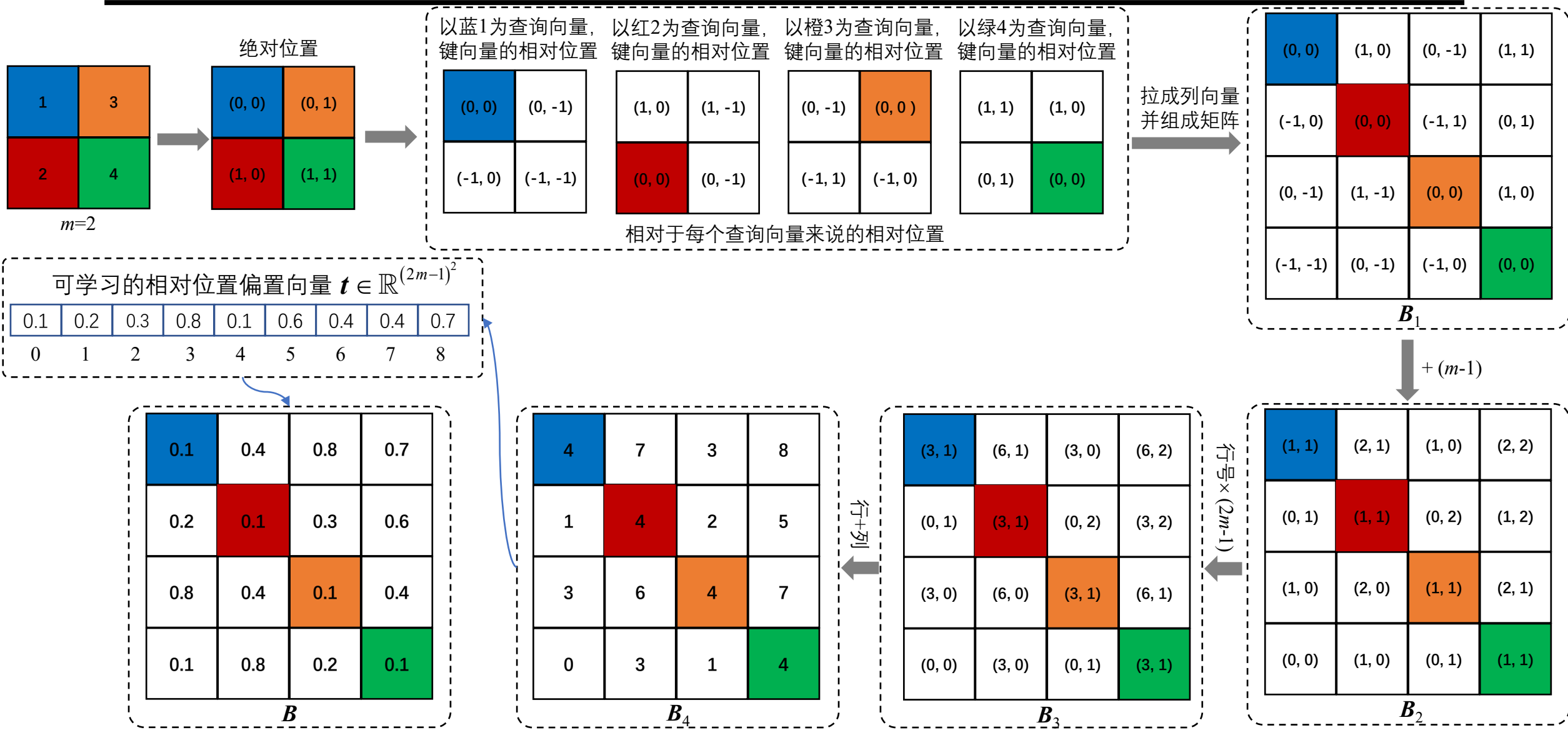
Attention scores,

$$\mathbf{A}_{m^2 \times m^2} = \mathbf{K}^T \mathbf{Q} = \begin{bmatrix} \alpha_{1,1} & \alpha_{2,1} & \cdots & \alpha_{m^2,1} \\ \alpha_{1,2} & \alpha_{2,2} & \cdots & \alpha_{m^2,2} \\ \vdots & \vdots & & \vdots \\ \alpha_{1,m^2} & \alpha_{2,m^2} & \cdots & \alpha_{m^2,m^2} \end{bmatrix}$$

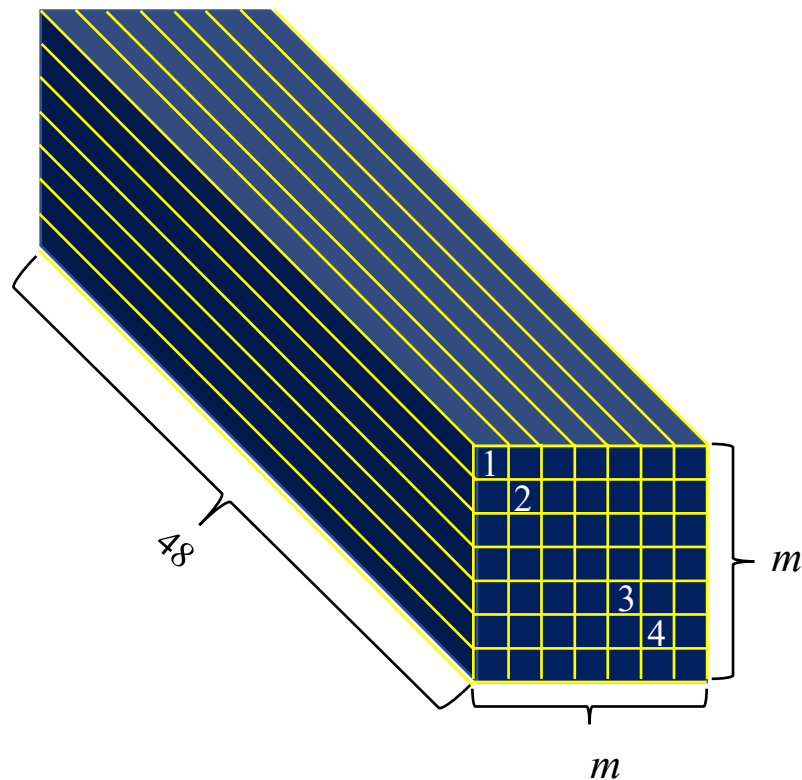
It implies that the RPB matrix \mathbf{B} should be of the dimension $m^2 \times m^2$, and $\mathbf{B}(i, j)$ reflects the relative position of the query \mathbf{q}_i and the key \mathbf{k}_j



Swin-Transformer



Swin-Transformer



- Relative position bias can encode the relative position relationship between two vectors in a 2D array
E.g., the RPB of “2” to “1” is the same as the RPB of “4” to “3”
- Positional encodings are added to the input vectors; differently, RPBs are added to attention scores

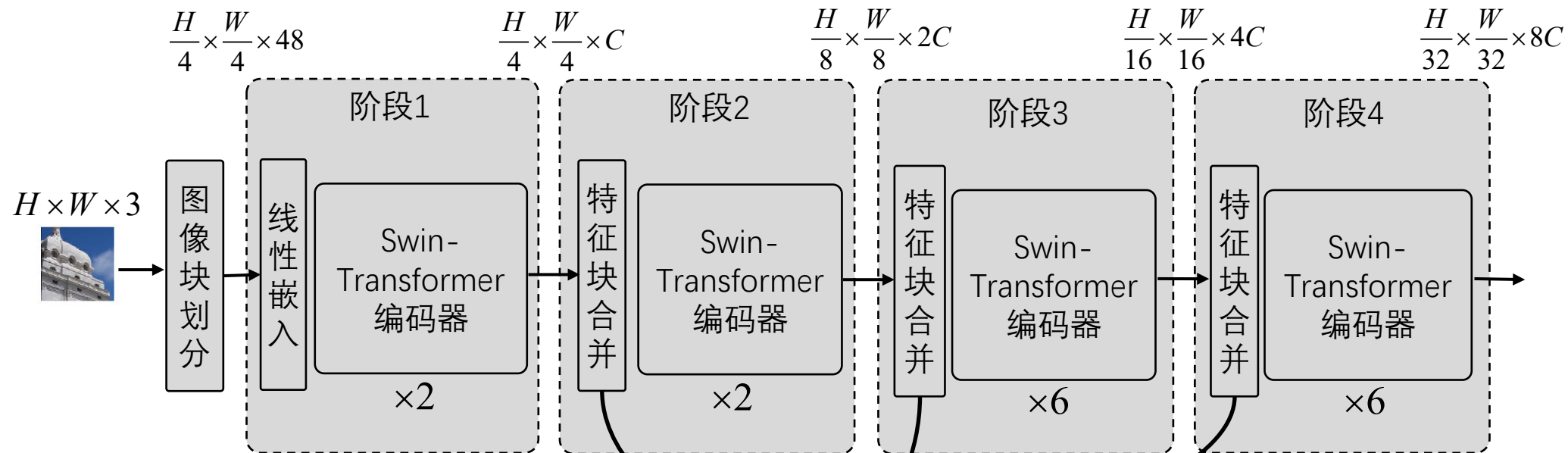
$$A'_{m^2 \times m^2} = \text{softmax} \left(\frac{K^T Q}{\sqrt{d}} + B \right)$$

$$= \text{softmax} \left[\begin{array}{c} \begin{array}{cccc} 1 & 2 & 3 & 4 \\ \alpha_{1,1} & \alpha_{2,1} & \alpha_{3,1} & \alpha_{4,1} \\ \alpha_{1,2} & \alpha_{2,2} & \alpha_{3,2} & \alpha_{4,2} \\ \alpha_{1,3} & \alpha_{2,3} & \alpha_{3,3} & \alpha_{4,3} \\ \alpha_{1,4} & \alpha_{2,4} & \alpha_{3,4} & \alpha_{4,4} \end{array} + \begin{array}{cccc} 0.1 & 0.4 & 0.8 & 0.7 \\ 0.2 & 0.1 & 0.3 & 0.6 \\ 0.8 & 0.4 & 0.1 & 0.4 \\ 0.1 & 0.8 & 0.2 & 0.1 \end{array} \end{array} \right]$$

\sqrt{d}



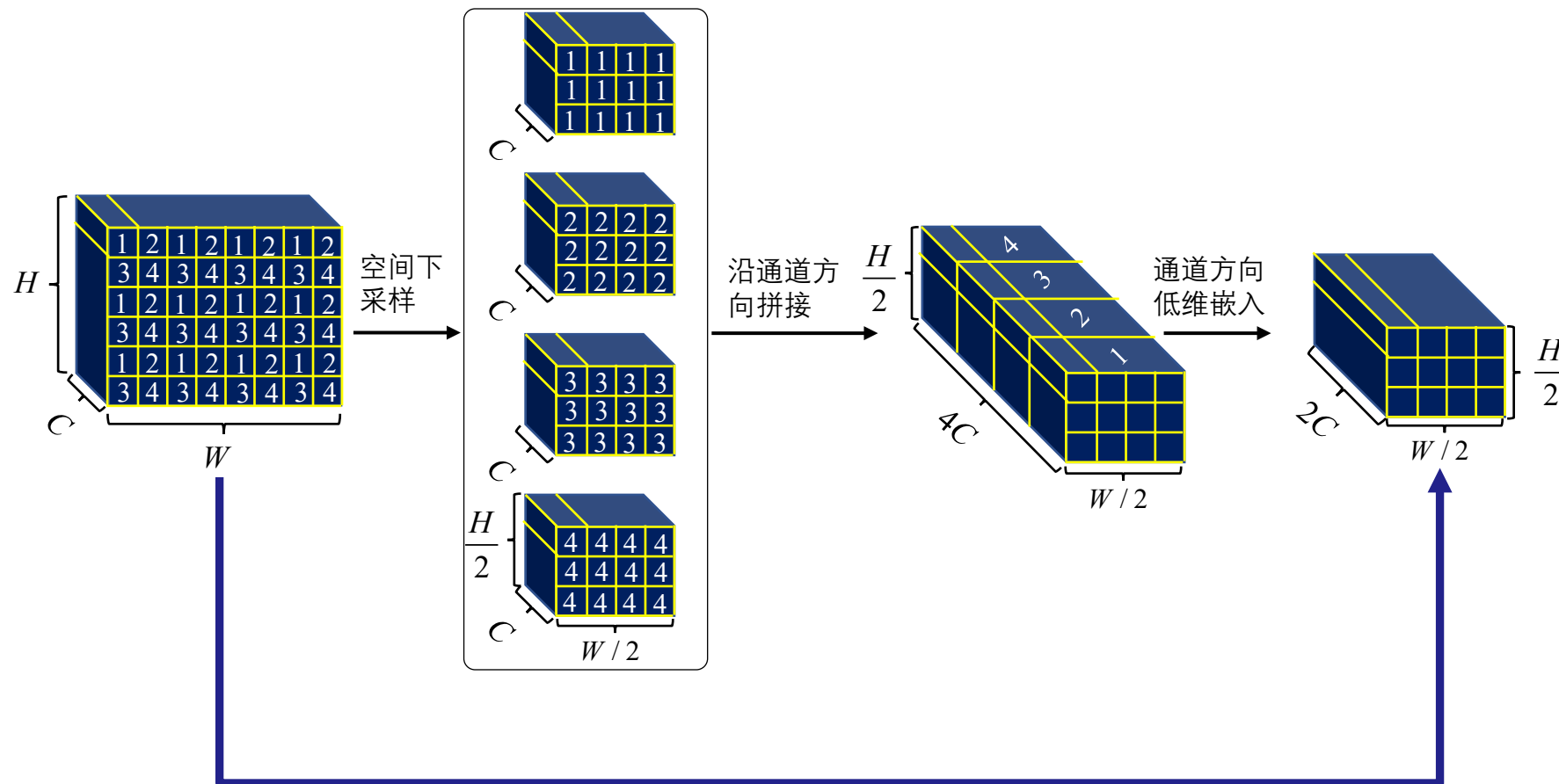
Swin-Transformer



Next page



Swin-Transformer



- ✓ Half the spatial resolution
- ✓ Double the channels



Swin-Transformer

- Play with Swin-Transformer
 - Swin-Transformer is a backbone network architecture
 - It outputs 4 feature maps with four different spatial resolutions; these four feature maps can be used similar as CNN feature maps; thus with different task heads, it can be used for object detection, object classification, semantic segmentation, etc.
 - The authors provide pre-trained Swin-Transformer models at four different scales

	downsp. rate (output size)	Swin-T	Swin-S	Swin-B	Swin-L
stage 1	4× (56×56)	concat 4×4, 96-d, LN	concat 4×4, 96-d, LN	concat 4×4, 128-d, LN	concat 4×4, 192-d, LN
		$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 96, head 3} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 96, head 3} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 128, head 4} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 192, head 6} \end{bmatrix} \times 2$
stage 2	8× (28×28)	concat 2×2, 192-d, LN	concat 2×2, 192-d, LN	concat 2×2, 256-d, LN	concat 2×2, 384-d, LN
		$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 192, head 6} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 192, head 6} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 256, head 8} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 384, head 12} \end{bmatrix} \times 2$
stage 3	16× (14×14)	concat 2×2, 384-d, LN	concat 2×2, 384-d, LN	concat 2×2, 512-d, LN	concat 2×2, 768-d, LN
		$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 384, head 12} \end{bmatrix} \times 6$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 384, head 12} \end{bmatrix} \times 18$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 512, head 16} \end{bmatrix} \times 18$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 768, head 24} \end{bmatrix} \times 18$
stage 4	32× (7×7)	concat 2×2, 768-d, LN	concat 2×2, 768-d, LN	concat 2×2, 1024-d, LN	concat 2×2, 1536-d, LN
		$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 768, head 24} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 768, head 24} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 1024, head 32} \end{bmatrix} \times 2$	$\begin{bmatrix} \text{win. sz. } 7 \times 7, \\ \text{dim 1536, head 48} \end{bmatrix} \times 2$



Outline

- Transformer in NLP
- Multi-head Attention
- Vision transform (ViT)
- Swin-Transformer
- DEtection TRansformer (DETR)
- Real-time DETR (RT-DETR)

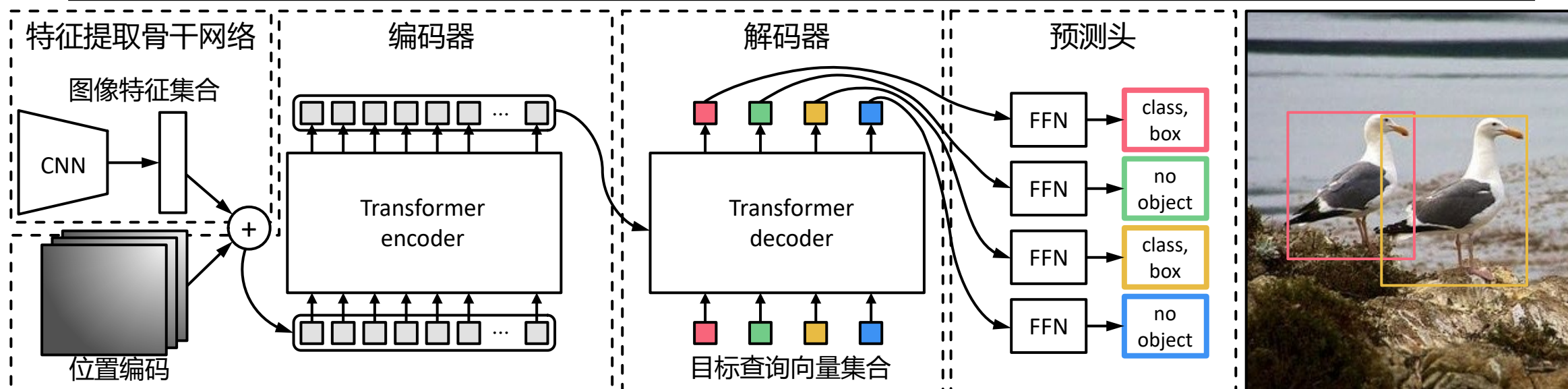


DEtection TRansformer (DETR)

- Swin-transformer actually is a feature extractor, while DETR adopts a **transformer encoder-decoder architecture** to perform object detection
- DETR^[1] is a milestone in the field of object detection
 - A transformer encoder-decoder architecture is adopted
 - It does not depend on anchor-boxes or anchor-points
 - It does not need the NMS (non-maximum suppression) post-processing
 - It creatively makes use of transformer decoder to fulfill the object detection task
 - It models the matching problem between the ground-truths and the detected bounding-boxes as a bipartite matching problem

[1] CARION N, Massa F, SYNNAEVE G, et al. End-to-end object detection with transformers[C]//Proc. Euro. Conf. Computer Vision, 2020: 213-229.

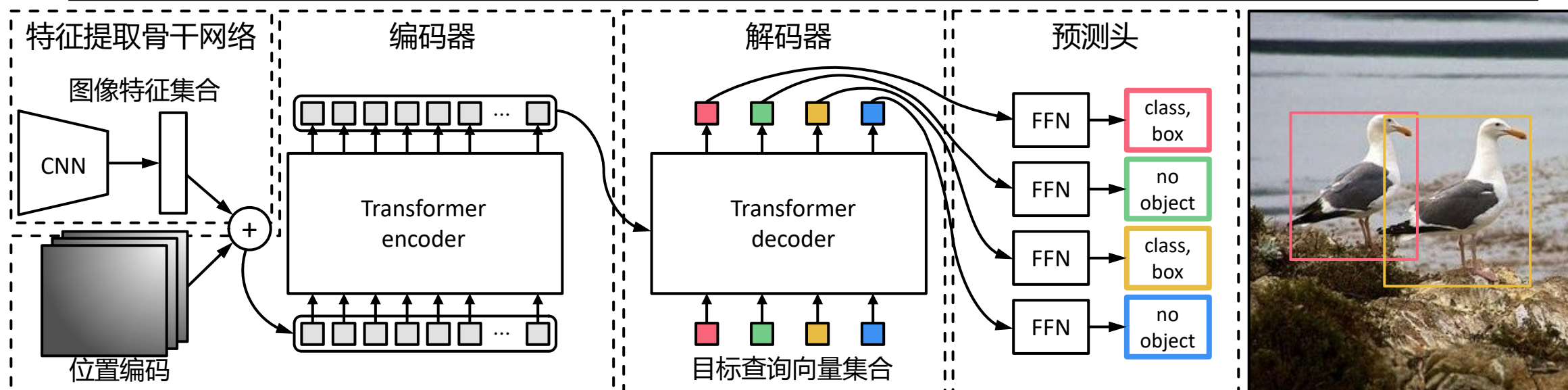
DETR—For Inference



Given an image $I \in \mathbb{R}^{800 \times 1088 \times 3}$, DETR will use the following steps to detect objects from it

- ✓ Use a CNN to extract the feature map F from I , and F 's resolution is 1/32 of I 's. $F \in \mathbb{R}^{25 \times 34 \times 2048}$
- ✓ Linearly map the channels of F to 256 using 1×1 convolutions, $F' \in \mathbb{R}^{25 \times 34 \times 256}$
- ✓ Flatten F' to get the input to the transform encoder, $F_I \in \mathbb{R}^{256 \times 850}$
- ✓ Encode F_I using 6 consecutive encoder blocks
- ✓ Decode the object queries. DETR's decoder has 6 consecutive decoder blocks; each decode block takes its precedent block's output as query sequences and takes the output of the corresponding encoder block as key sequences and value sequences; the input to the first decoder block is object queries $O \in \mathbb{R}^{256 \times 100}$. Denote $P \in \mathbb{R}^{256 \times 100}$ as the output of the last decoder block.

DETR—For Inference

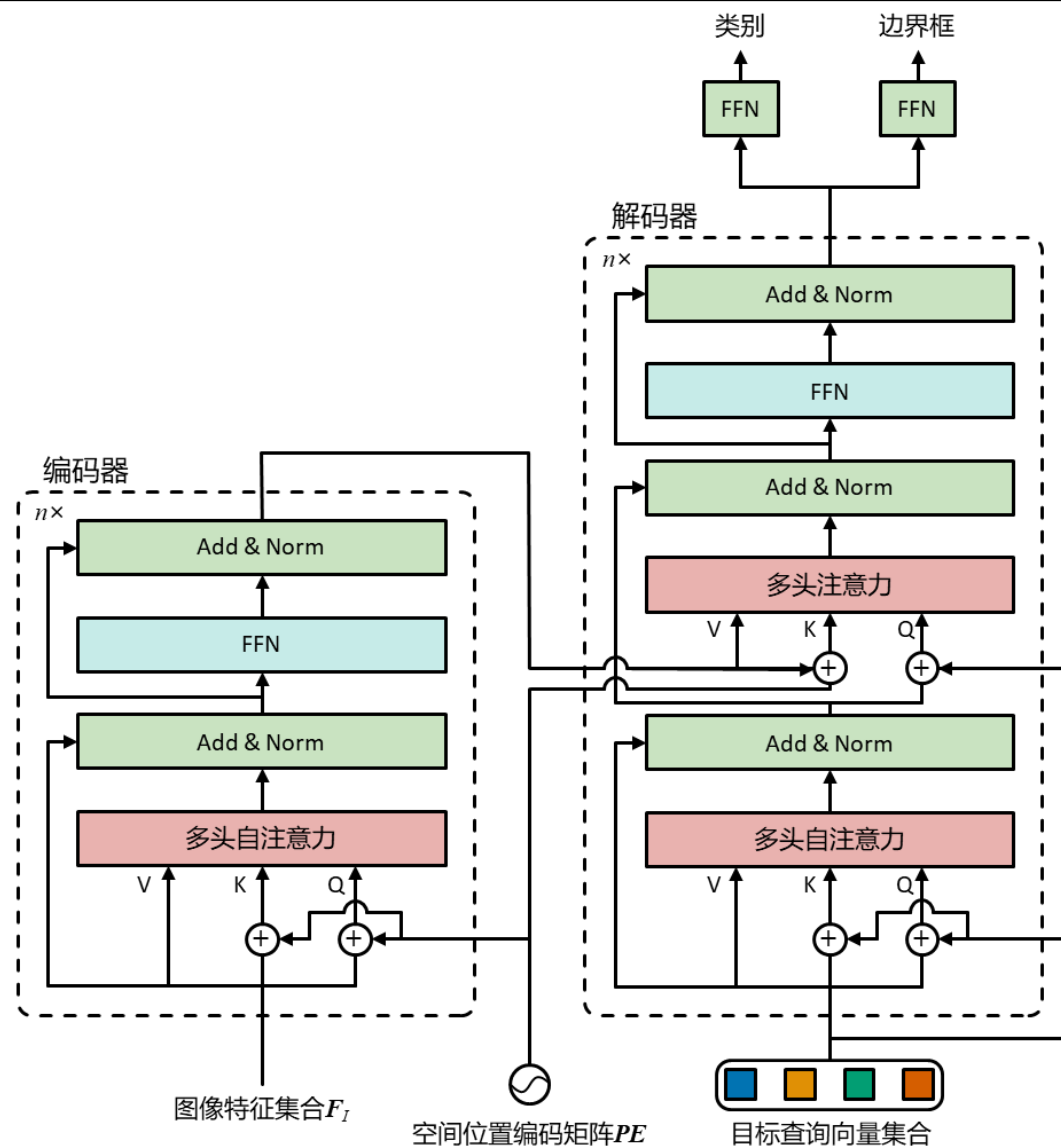


Given an image $I \in \mathbb{R}^{800 \times 1088 \times 3}$, DETR will use the following steps to detect objects from it

- ✓ Interpretate object detection results from $P \in \mathbb{R}^{256 \times 100}$. Column p_i corresponds to one detected object: p_i will be fed into two fully-connected forward networks—one outputs the classification probabilities and one regresses the bounding box (a $4d$ vector, two for the center position and two for the size). It needs to be noted that: if the DETR model is trained to detect objects of 20 classes, the classification probability vector will be of 21d with the last node denoting the probability of this object to be “ Φ ”, i.g., it is not a valid object



DETR—For Inference



The detailed transformer encoder-decoder architecture used in DETR



DETR—For Training

The key problem is: how to compute the loss for a given training image I ?

Denote I 's ground-truth by the set \mathbf{Y} ; denote the detected objects from I by $\hat{\mathbf{Y}} = \left\{ \hat{\mathbf{y}}_i \right\}_{i=1}^N$

Usually N is set to be very large (and by default it is set as $N = 100$), and thus it is valid to suppose that N is larger than the number of ground-truth objects in \mathbf{Y} .

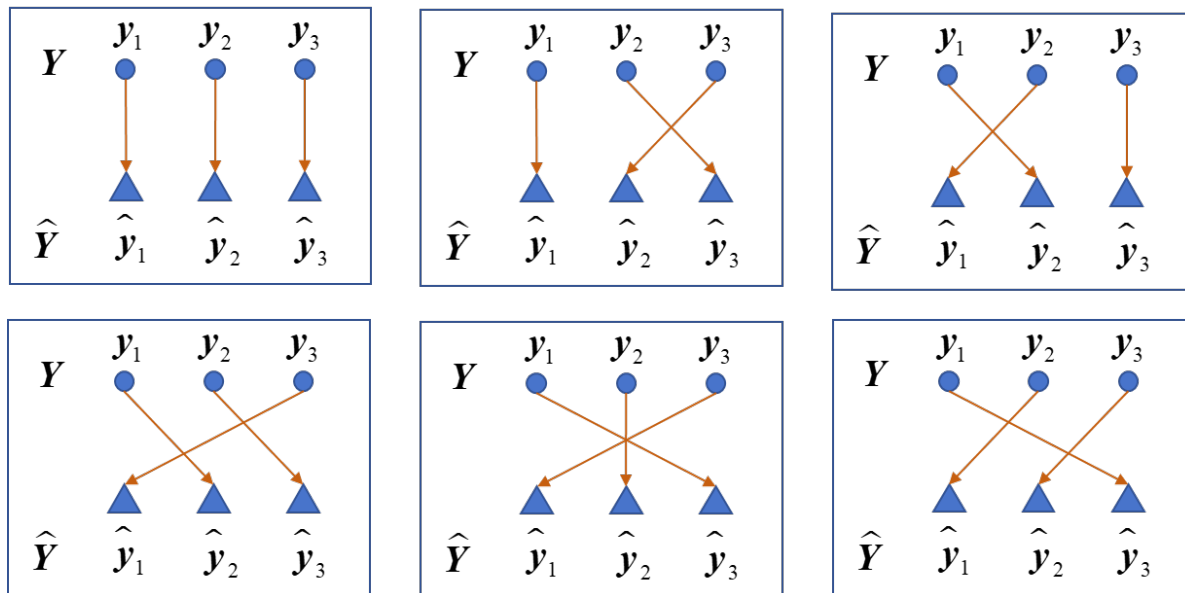
We can extend \mathbf{Y} to make it also have N objects by injecting objects with class labels “ Φ ”

To compute the loss between \mathbf{Y} and $\hat{\mathbf{Y}}$, there are two steps.



DETR—For Training

Step 1: find the optimal matching between \mathbf{Y} and $\hat{\mathbf{Y}}$



$$\sigma^* = \arg \min_{\sigma \in \Omega_N} \sum_{i=1}^N l_{\text{match}} \left(\mathbf{y}_i, \hat{\mathbf{y}}_{\sigma(i)} \right) \quad (1)$$

where Ω_N is the permutation set of N numbers $\{1, 2, \dots, N\}$

The ground-truth object is represented as $\mathbf{y}_i = \begin{pmatrix} c_i \\ \mathbf{b}_i \end{pmatrix}$

where c_i is its class label and \mathbf{b}_i is its bounding-box

$$l_{\text{match}} \left(\mathbf{y}_i, \hat{\mathbf{y}}_{\sigma(i)} \right) = -\mathbf{1}_{\{c_i \neq \phi\}} \hat{p}_{\sigma(i)}(c_i) + \mathbf{1}_{\{c_i \neq \phi\}} l_{\text{box}} \left(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)} \right)$$

where $\hat{p}_{\sigma(i)}(c_i)$ is the probability of the predicted object $\hat{\mathbf{y}}_{\sigma(i)}$ belonging to class c_i , and $\hat{\mathbf{b}}_{\sigma(i)}$ is its bounding-box

The problem (1) is called the bipartite matching problem in graph theory and can be solved by the classic Hungarian algorithm



DETR—For Training

Step 2: Having found the optimal matching σ^* between \mathbf{Y} and $\hat{\mathbf{Y}}$, we compute their loss as,

$$l(\mathbf{Y}, \hat{\mathbf{Y}}) = \sum_{i=1}^N \left[-\log \hat{p}_{\sigma^*(i)}(c_i) + \mathbf{1}_{\{c_i \neq \phi\}} l_{\text{box}}(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma^*(i)}) \right]$$

where,

$$l_{\text{box}}(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)}) = \lambda_{\text{iou}} l_{\text{iou}}(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)}) + \lambda_{L_1} \|\mathbf{b}_i - \hat{\mathbf{b}}_{\sigma(i)}\|_1$$
$$l_{\text{iou}}(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)}) = 1 - \left(\frac{|\mathbf{b}_i \cap \hat{\mathbf{b}}_{\sigma(i)}|}{|\mathbf{b}_i \cup \hat{\mathbf{b}}_{\sigma(i)}|} - \frac{|B(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)}) \setminus (\mathbf{b}_i \cup \hat{\mathbf{b}}_{\sigma(i)})|}{|B(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)})|} \right)$$

where $B(\mathbf{b}_i, \hat{\mathbf{b}}_{\sigma(i)})$ is the least bounding-box enclosing \mathbf{b}_i and $\hat{\mathbf{b}}_{\sigma(i)}$



Outline

- Transformer in NLP
- Multi-head Attention
- Vision transform (ViT)
- Swin-Transformer
- DEtection TRansformer (DETR)
- Real-time DETR (RT-DETR)

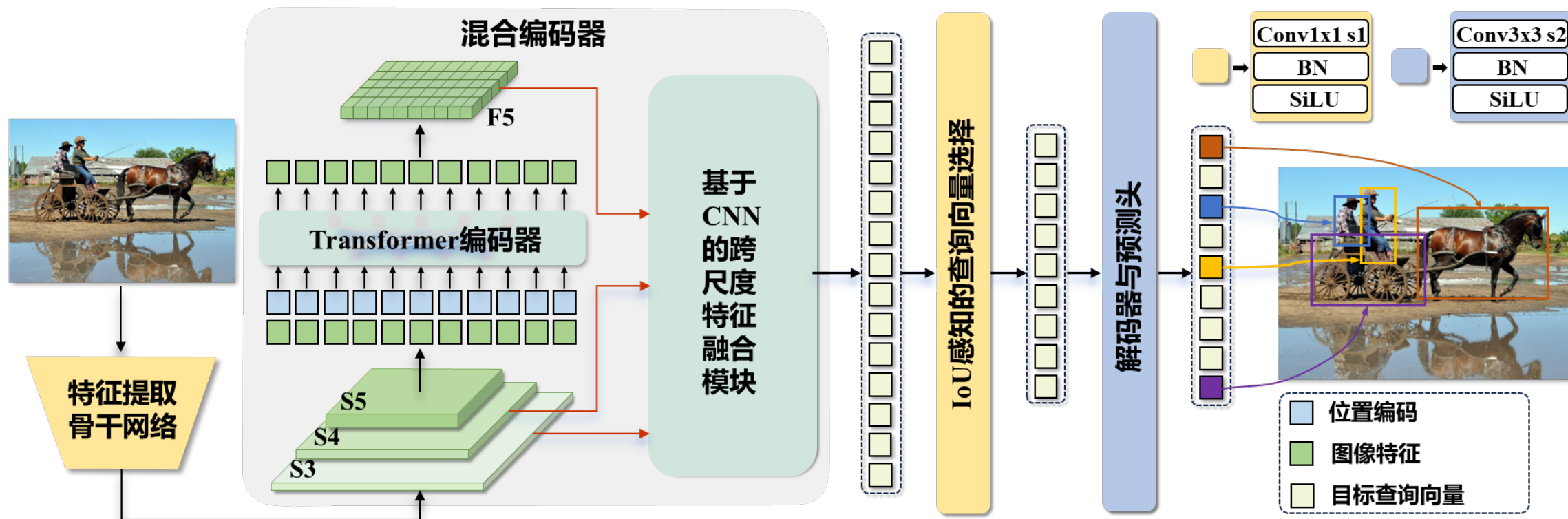


RT-DETR

- RT-DETR^[1] improves upon DETR and it was proposed by scholars from Peking University and Baidu
 - It runs much faster than DETR **and actually it is a real-time end-to-end object detector**
 - The authors designed a hybrid encoder by decoupling intra-scale interaction and cross-scale fusion to improve speed
 - The authors used an IoU-aware query selection mechanism to better initialize the query vectors
 - When computing the loss between the ground-truth and the predicted object, it uses **varifocal loss** (with varifocal loss, if a detected object has higher classification confidence, it will also have a more precise bounding-box)

[1] ZHAO Y, LV W, XU S, et al. DETRs beat YOLOs on real-time object detection[C]//Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2024.

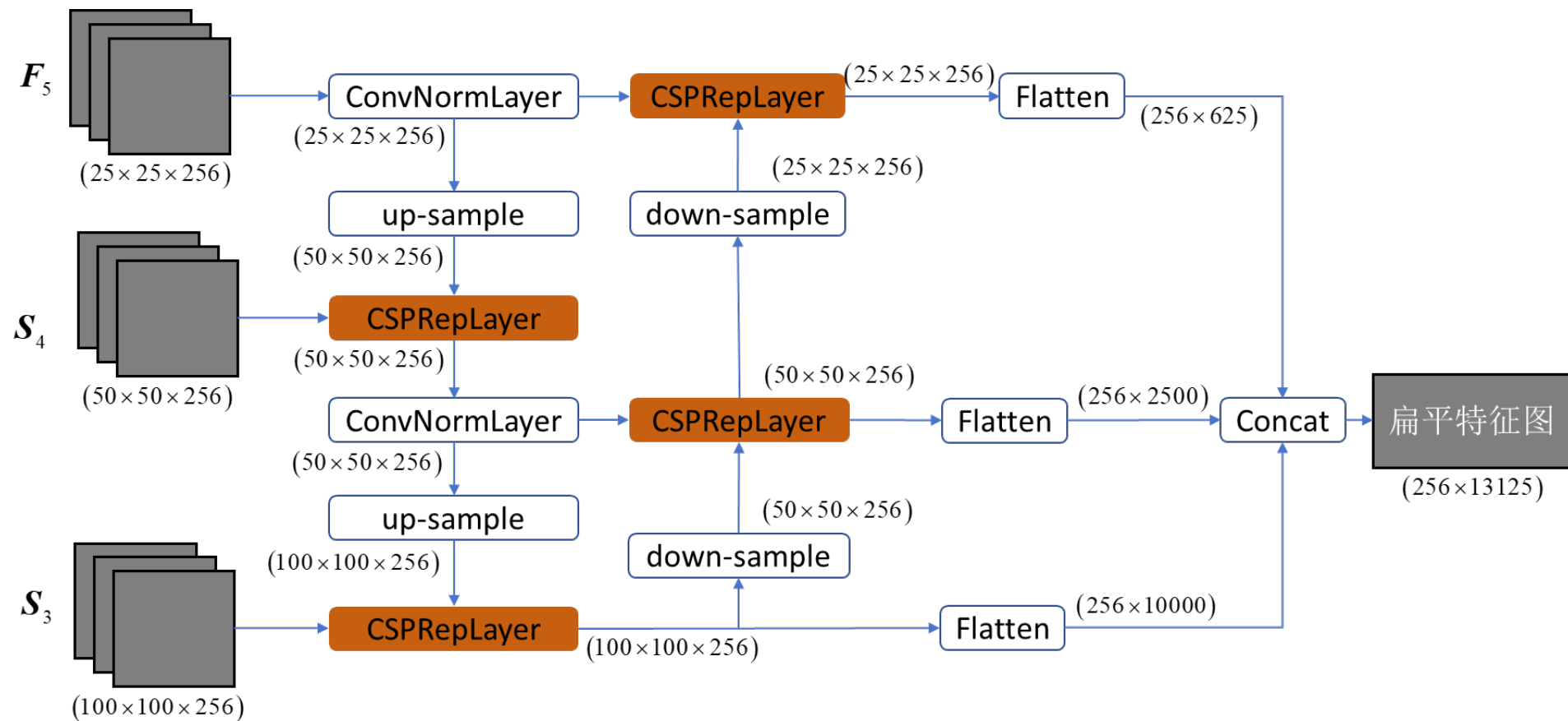
- The general architecture of RT-DETR





RT-DETR

- The CNN-based cross-scale feature fusion module





RT-DETR

Live DEMO

