

# Transformers for Vision

# Transformer

- The de facto neural architecture of many applications
  - GPT: generative pre-training Transformer
  - Gemini: multi-modal LLMs
  - Computer vision
  - Speech recognition



ChatGPT



Gemini



Claude



Llama

# Attention

- Attention compute a weighted average of hidden state vectors.
- Formally:
  - Given  $N$   $k$ -dimensional input vector  $x^1 \dots x^N$  and output vector  $y^1 \dots y^N$
  - Output  $y^i$  is a weighted average of input vector;  $w_{ij}$  is the attention weight from position  $i$  over  $x^j$ 
    - $y^i = \sum_{j=1}^N w_{ij} x^j$

# Formulation of Self-attention

- Given  $N$   $k$ -dimensional input vector  $x^1 \dots x^N$  and output vector  $y^1 \dots y^N$
- Conventional projection layer:
  - Output  $y^i$  is a weighted average of input vector

$$\bullet y^i = \sum_{j=1}^N w_{ij} x^j$$

- Self-attention: self-dependent
  - The weight  $w_{ij}$  depends on input **itself**  $x^j$

$$\bullet w'_{ij} = \sum_k x_k^i x_k^j$$

- Followed by Softmax normalization:  $w_{ij} = \frac{\exp(w'_{ij})}{\sum_j \exp(w'_{ij})}$

# Self-attention with Query, Key, Value

- Self-attention: The weight  $w_{ij}$  depends on input itself
  - $w'_{ij} = \sum_k x_k^i x_k^j$
  - the input  $x^i$  is to matched the input  $x^j$
- More flexible version: add learnable parameters
  - Linear projection layer with learnable weight matrix  $W \in \mathbb{R}^{k \times k}$
  - Query:  $Q^i = W_q x^i$ , to match others
  - Key:  $K^i = W_k x^i$ , to be matched
  - Value:  $V^i = W_v x^i$

# Self-attention with Query, Key, Value

- Q, K, V are linearly projected
  - Query:  $Q^i = W_q x^i$
  - Key:  $K^i = W_k x^i$
  - Value:  $V^i = W_v x^i$
- Compute attention weight based on  $Q^i$  and  $K$ 
  - $A_{ij} = \frac{\exp(Q^i K^j)}{\sum_j \exp(Q^i K^j)}$

# Self-attention with Query, Key, Value

- Compute attention weight based on  $Q^i$  and  $K$

- $A_{ij} = \frac{\exp(Q^i K^j)}{\sum_j \exp(Q^i K^j)}$

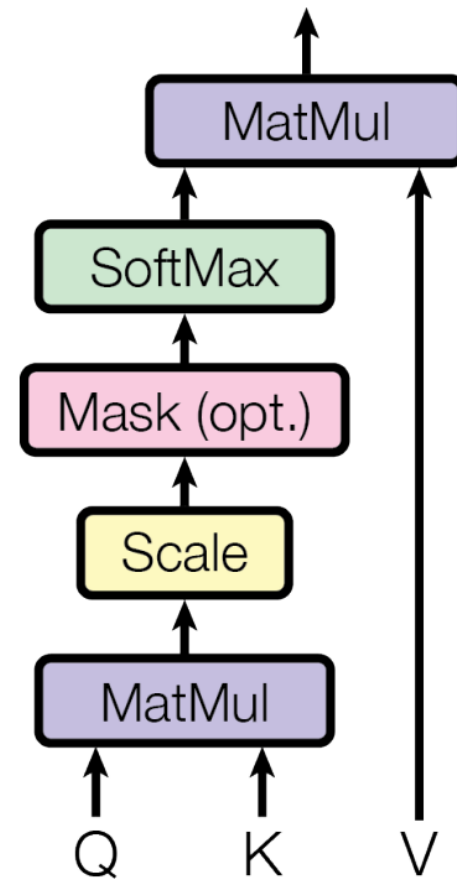
- Aggregated with value vector  $V$

- The new output for i-th position depends on the attention weights and value vectors of all input positions  $j$

- $y^i = \sum_{j=1}^T A_j^i V^j$

# Scaled Dot-product Attention

- Scale factor  $\sqrt{d}$  normalizes the dot product values  $QK^T$ , preventing their variance from becoming overly large.
- Scale dot-product:  $A = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$
- Matrix Multiplication:  $y = AV$





scaled dot – product:

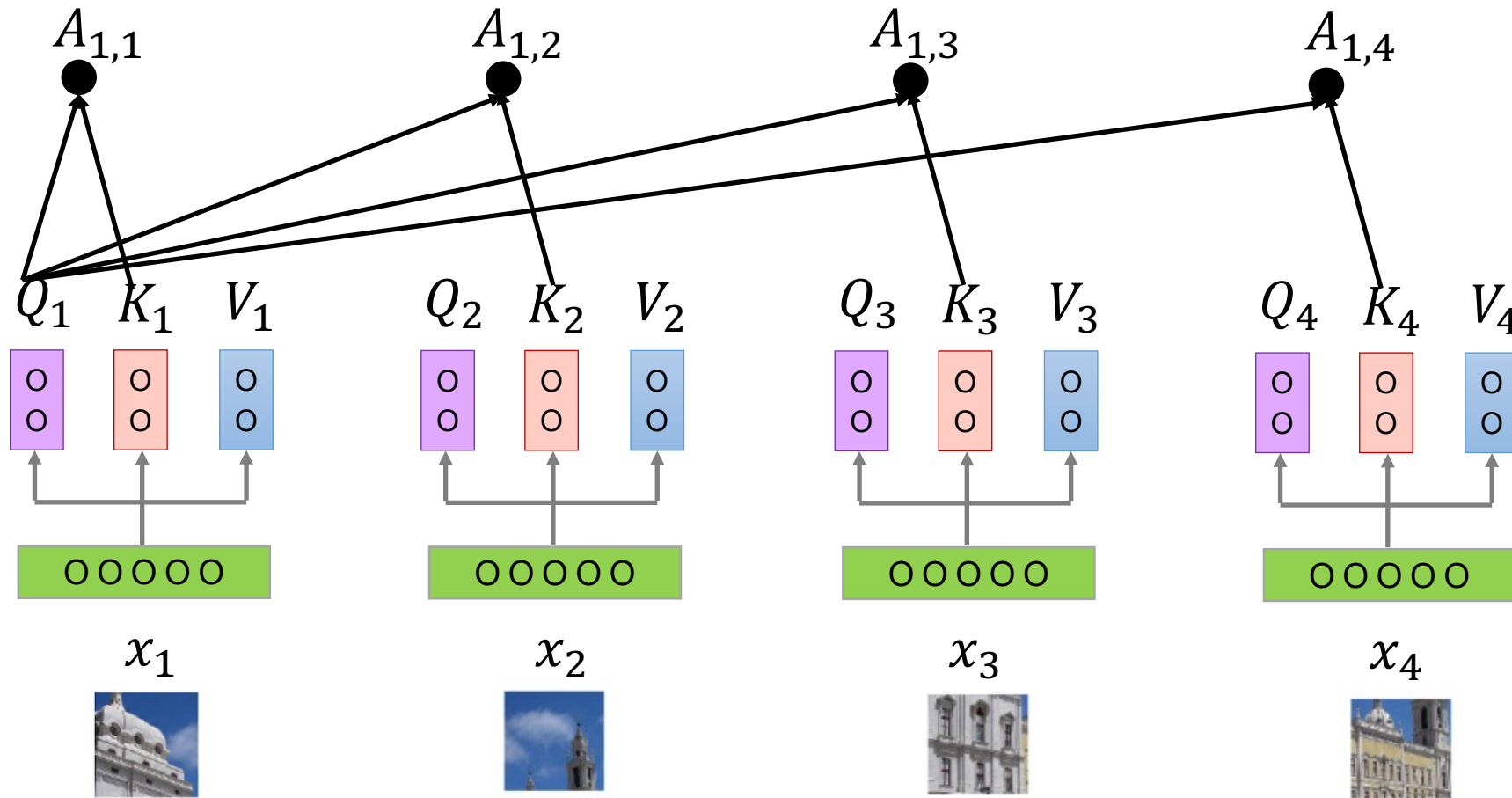
$$A_{1,i} = \frac{QK^T}{\sqrt{d}}$$

$Q$ : query (to match others)

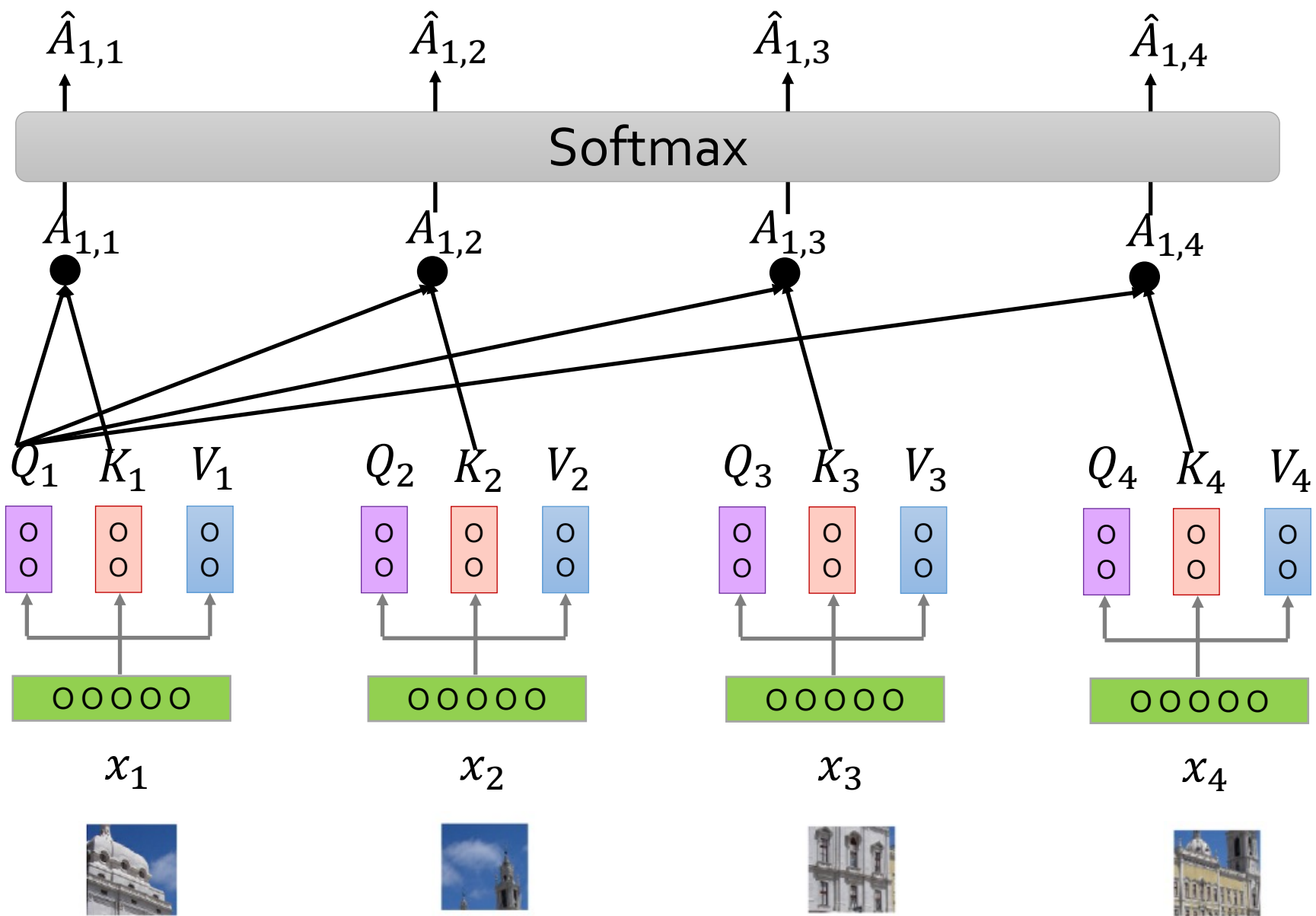
$K$ : key (to be matched)

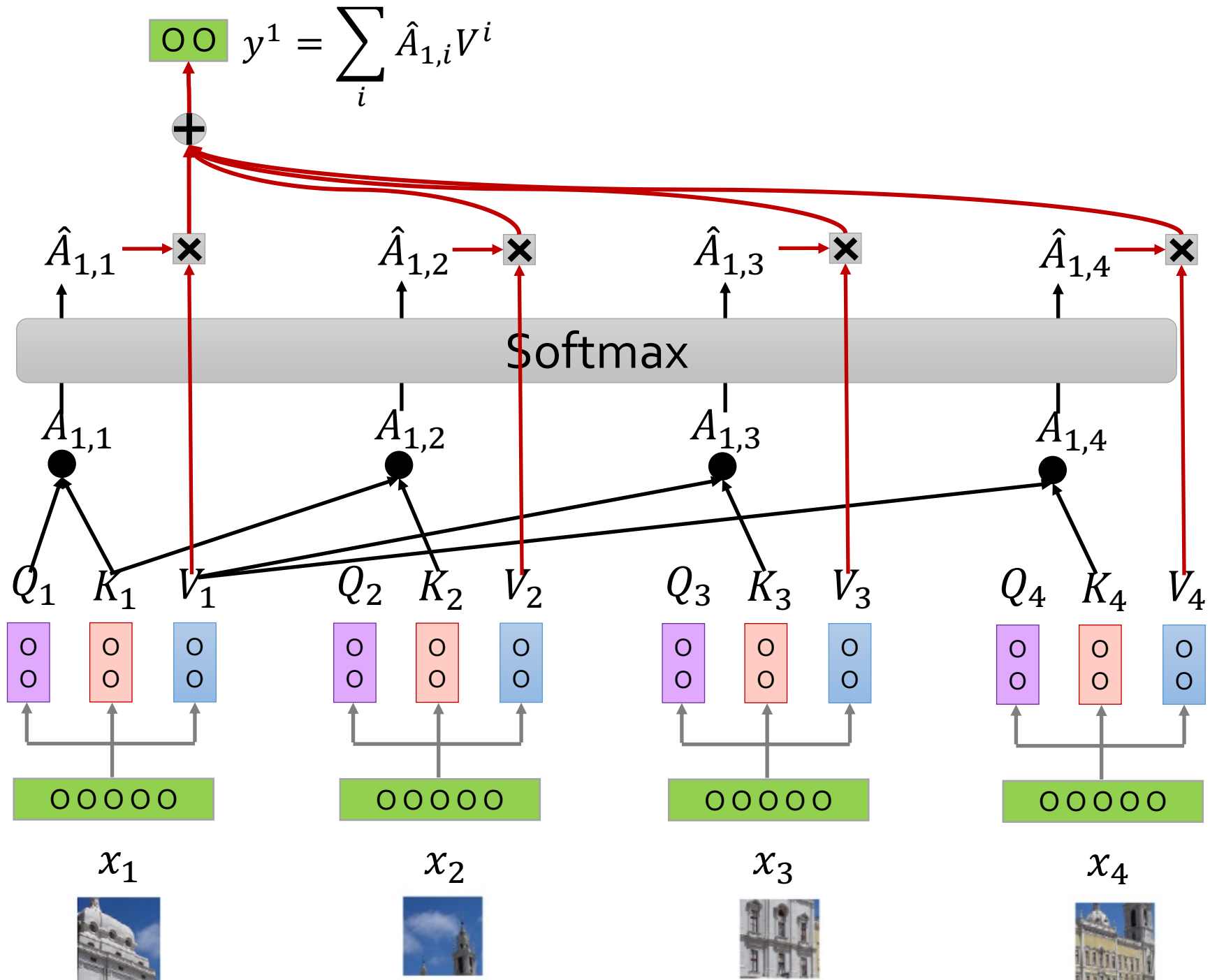
$V$ : value

$x$ : input vector before projection



$$\sigma(z)_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$$





# Quadratic complexity

- $\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right) V$
- Complexity:  $O(n^2 \cdot d)$
- $n$  = sequence length,  $d$  = hidden dimension

# Cross-attention

- Given two different input vector  $x$  and  $m$
- Input to Cross-attention
  - Query:  $Q^i = W_q x^i$
  - Key:  $K^i = W_k m^i$
  - Value:  $V^i = W_v m^i$
- Source vector  $x$  is attended to target vector  $m$
- Compute cross-attention weight based on  $Q^i$  and  $K$ 
  - $A_{ij} = \frac{\exp(Q^i K^j)}{\sum_j \exp(Q^i K^j)}$

# Rethink Cross-attention as Feature Clustering

- (not exactly the same)
- input vector  $x \in \mathbb{R}^{N \times d}$  ( $Q^i = W_q x^i$ ) is the feature centroids with N centroids
- input vector  $m \in \mathbb{R}^{HW \times d}$  ( $K^j = W_k m^j$ ,  $V^i = W_v m^i$ ) is the pixel features
- Cross-attention weight v.s. Soft assignment based on feature distance
  - $A_{ij} = \frac{\exp(Q^i K^j)}{\sum_l \exp(Q^i K^l)}$
- Value aggregation v.s. Centroid update:
  - $O^i = \sum_{j=1}^M A_j^i V^j$

# Positional encoding

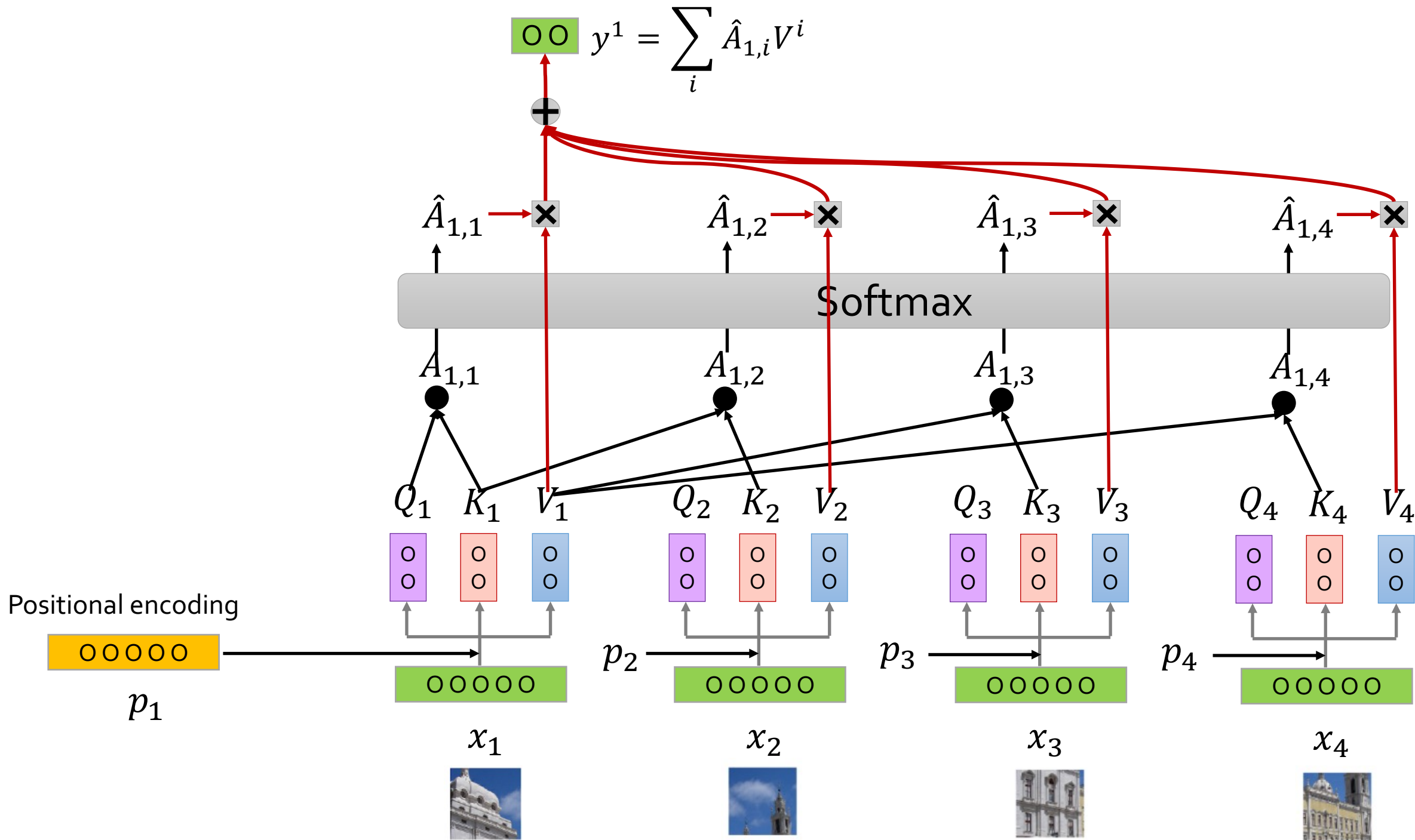
- Motivation:
  - While the attention mechanism is powerful, it doesn't inherently capture the **order of elements** within the sequence.
  - In language, the order of words is crucial for meaning. For example, "The cat sat on the mat" has a different meaning than "The mat sat on the cat."
  - Similarly, in computer vision, the permutation of image pixels matters.

# Positional encoding

$$\bar{x}_t = \begin{bmatrix} x_t \\ t \end{bmatrix}$$

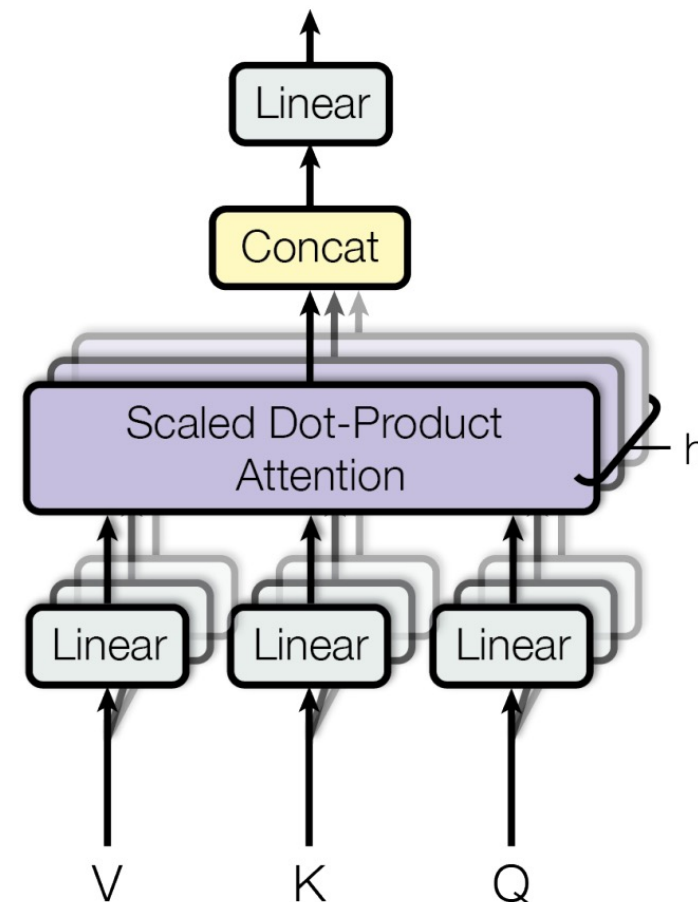
- append  $t$  to the input
  - This is not a great idea, because absolute position is less important than relative position
  - we want to represent position in a way that tokens with similar relative position have similar positional encoding
- More advanced positional encoding (e.g., frequency-based, learning-based, relative positional encoding) are not discussed in this lecture.





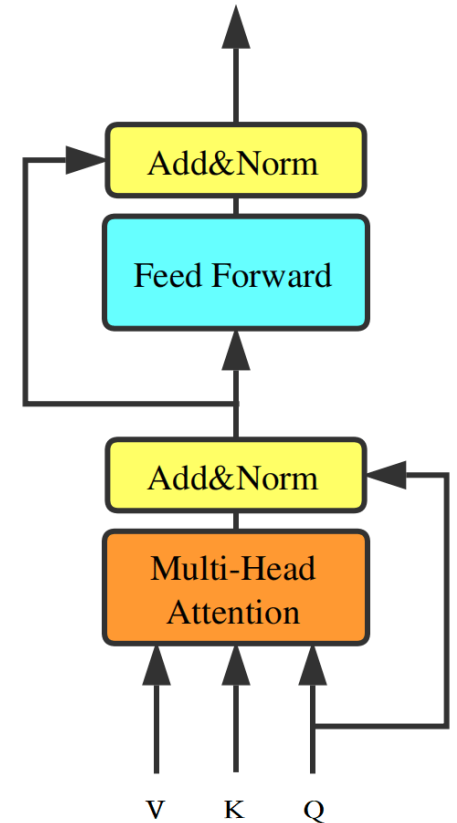
# Multi-head self-attention

- Run  $h$  attention models in parallel on top of different linearly projected versions of  $Q, K, V$ ; concatenate and linearly project the results
- Intuition: enables model to attend to different kinds of information at different positions

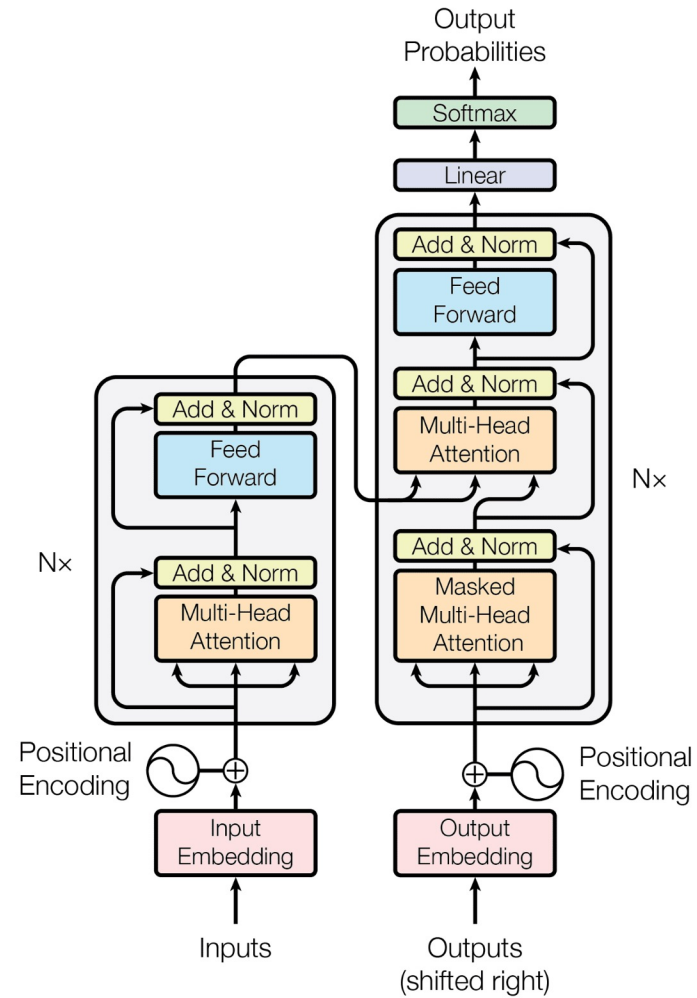


# Transformer block

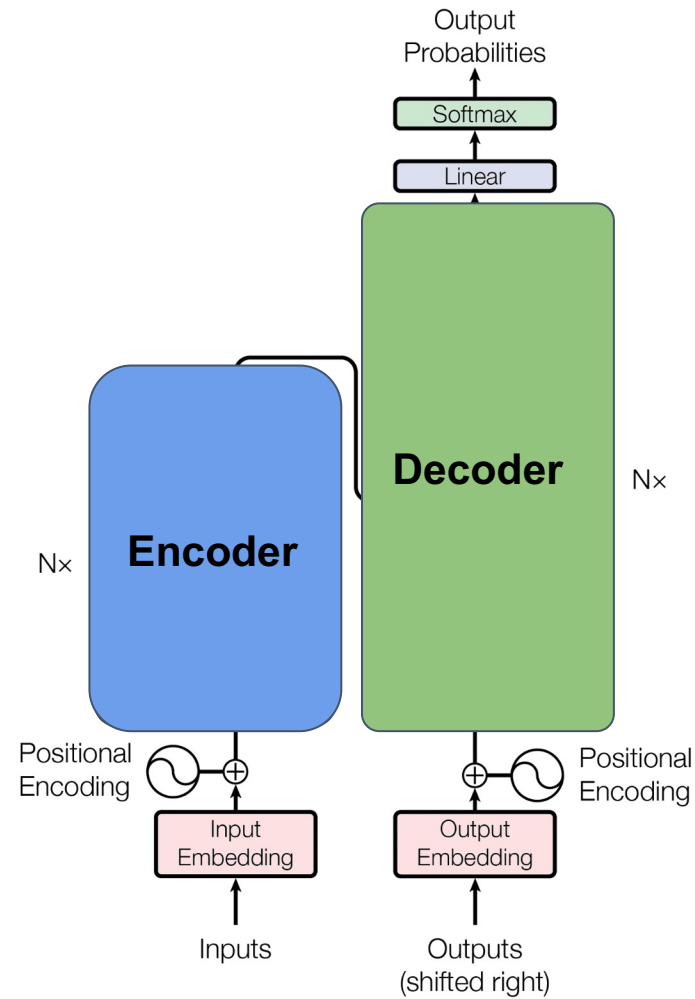
- A Transformer is a sequence of transformer blocks
  - Vaswani et al.:
    - 12 blocks, 512 embedding dimension, 6 attention heads
  - **Multi-Head Attention**: introduced before
  - **Add & Norm**: residual connection followed by layer normalization
  - **Feedforward (Multi-layer perceptron)**: two linear layers with ReLUs in between, applied independently to each vector
- Attention is the only interaction between inputs



# Transformers



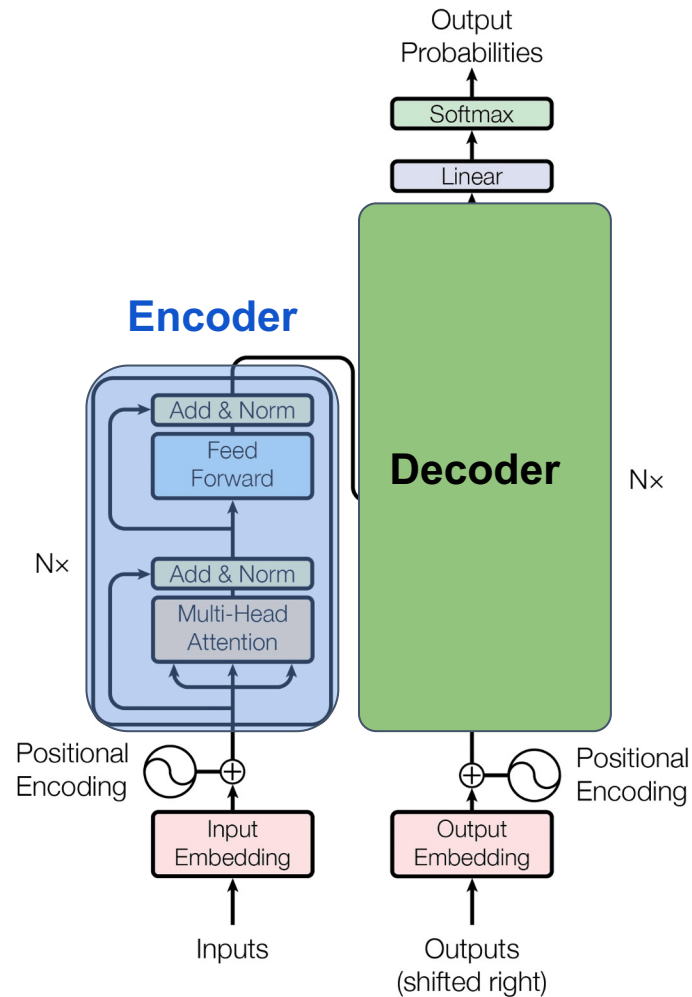
# Transformers



# Transformers

## Transformer Encoder:

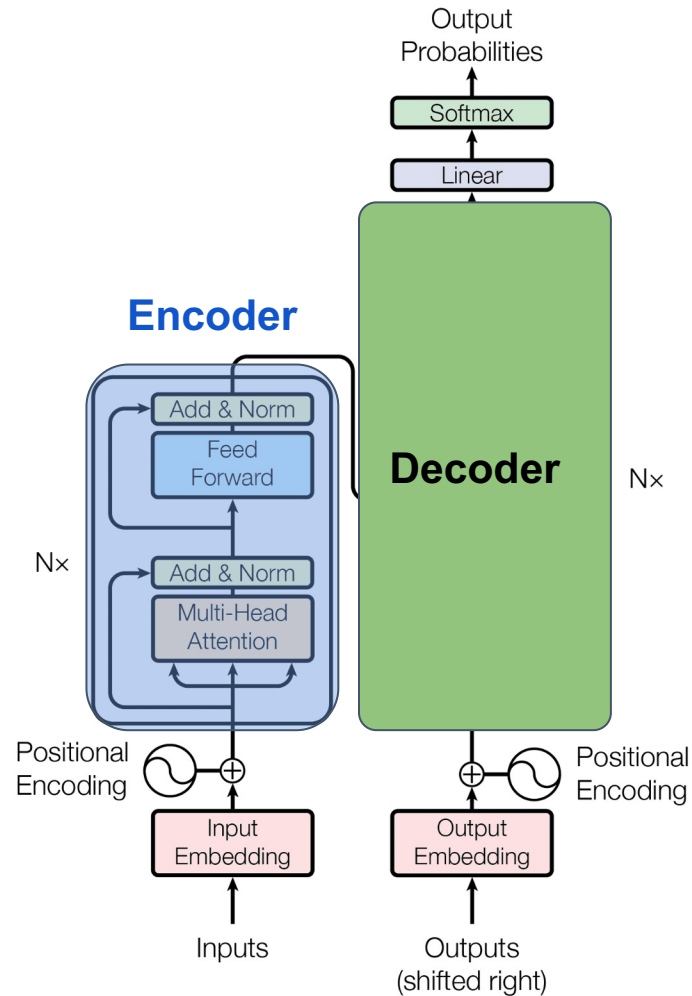
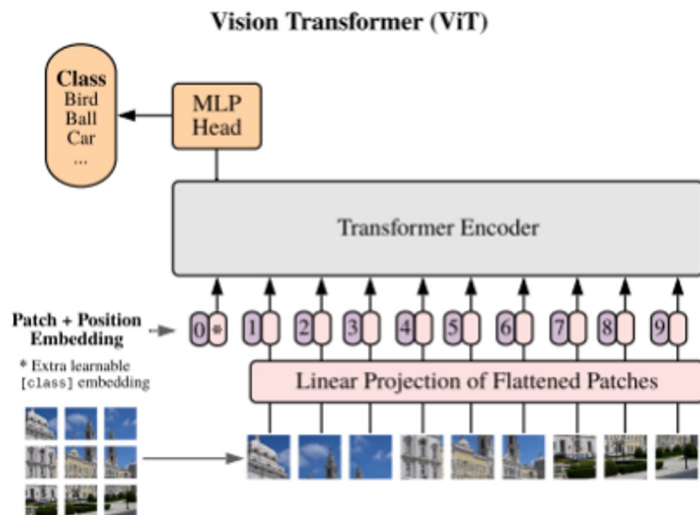
- ❑ Self-attention
- ❑ Serves as strong feature extracting/enhancing modules.



# Transformers

## Transformer Encoder:

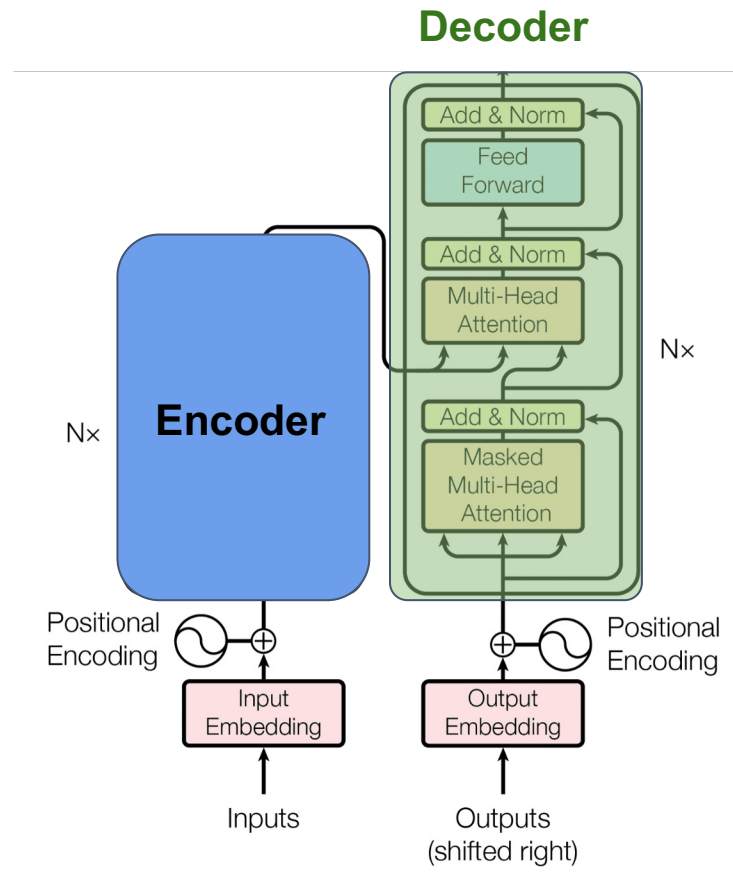
- Vision Transformer (ViT)



[1] Attention Is All You Need. NeurIPS 2017.

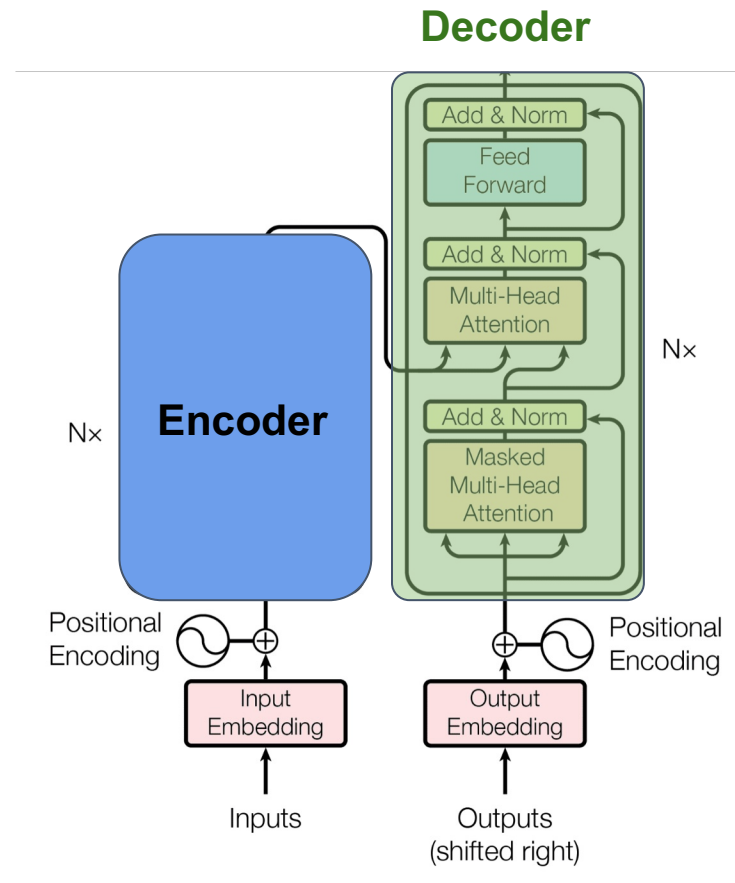
[2] ViT: Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. ICLR 2021.

# Transformers





# Transformers



- ## Transformer Decoder:
- ❑ Cross-attention
  - ❑ Serves as a bridge for sequence-to-sequence translation

# Transformers

➤ Extracting pixel-level representation

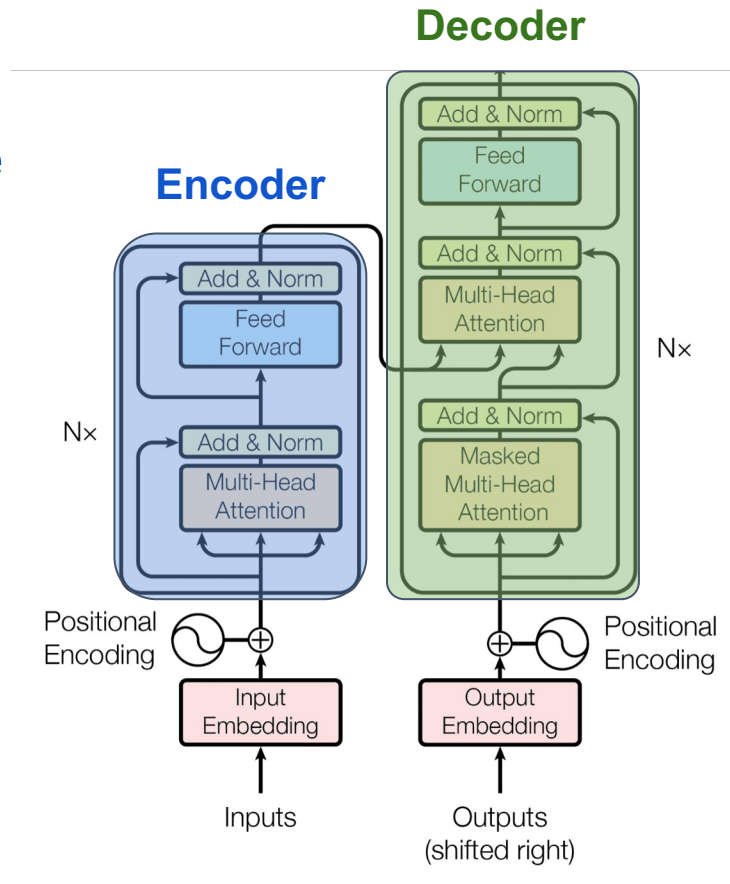
## Transformer Encoder:

- ❑ Self-attention
- ❑ Serves as strong feature extracting/enhancing modules.

➤ Converting to object-level representation

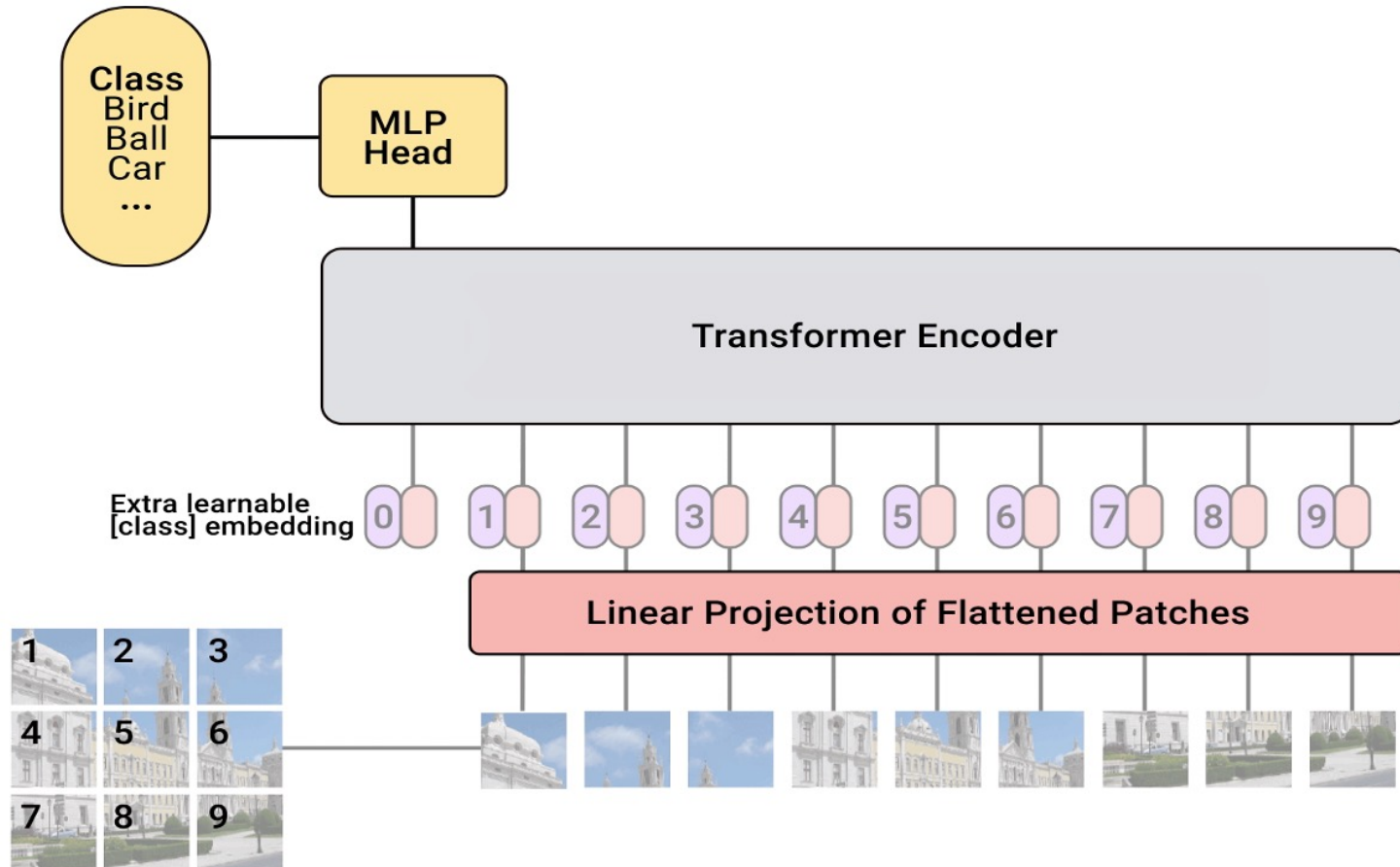
## Transformer Decoder:

- ❑ Cross-attention
- ❑ Serves as a bridge for sequence-to-sequence translation



# Vision Transformers (ViTs)

ViT: the pioneer of Transformer architectures for vision



[1] Dosovitskiy, et al. An image is worth 16x16 words: Transformers for image recognition at scale. ICLR 2021

# Vision Transformers (ViTs)



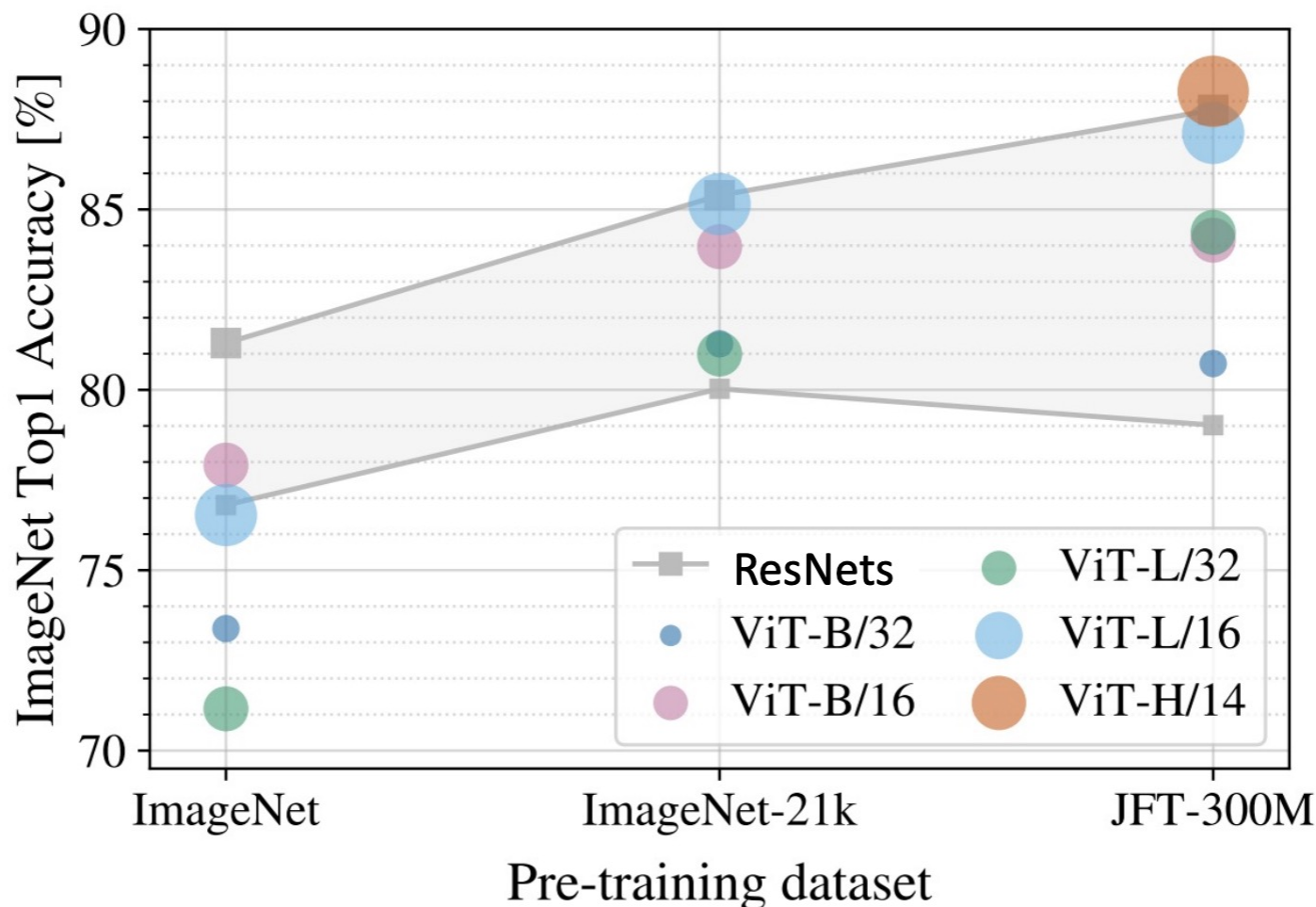
# ViT: overview

- Divide an input image into 196 ( $14 \times 14$ ) small images of size ( $16 \times 16$ )
- Treat it as embedding in NLP
- Use it as an input for traditional transformer encoder (like in BERT)
- Use 12 transformer layers (Norm, Multi-head attention, etc.)
- the last output, use it as input for Dense Layer with 1000
- you have a classification model

# ViT: performance comparison to CNN

JFT-300M is an internal Google dataset with 300M labeled images

If you pretrain on JFT and finetune on ImageNet, large ViTs outperform large ResNets



B = Base  
L = Large  
H = Huge

/32, /16, /14 is patch size; smaller patch size is a bigger model (more patches)

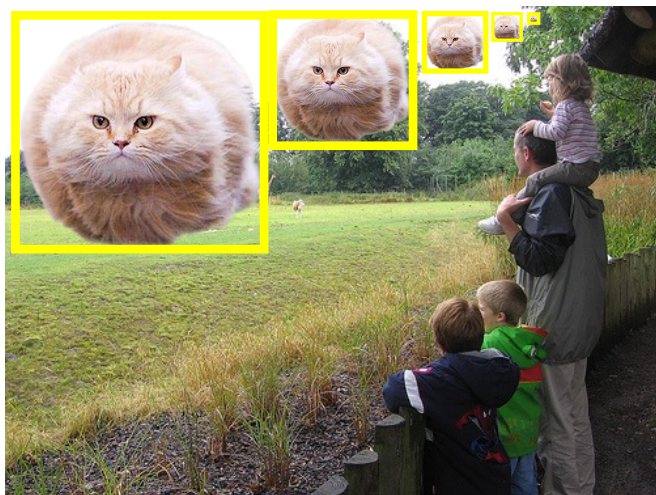
# ViT: strength and weakness

- Strength:
  - Long-range interaction via global attention
  - Strong performance/scalability with large scale of data
- Weakness:
  - Quadratic computational complexity
  - Lacking locality
  - Lacking low-level details due to large patch size (e.g., 16x16)

# Key difference between Visual and Text Signals

**multi-scale**

(scale invariance)



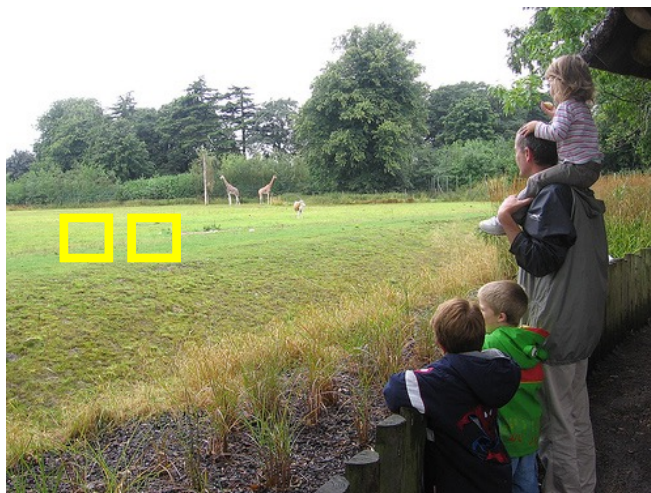
I am a fat cat.

I am a fat fat cat cat. (invalid)

No scale variation

**locality**

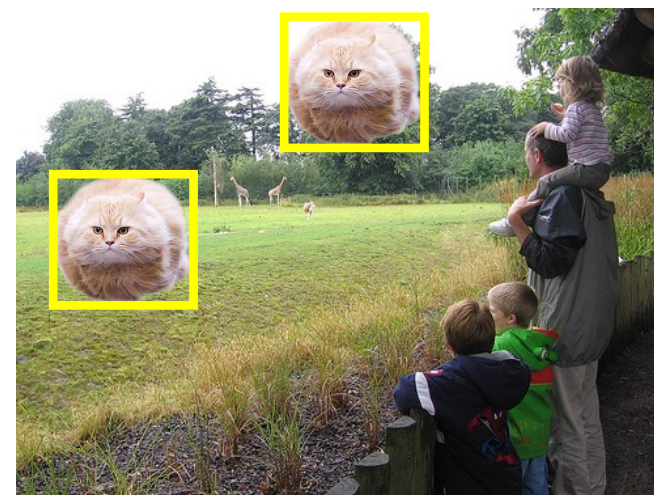
(spatial smoothness)



I like the green grass.

No spatial smoothness

**translation invariance**



I am a fat cat.

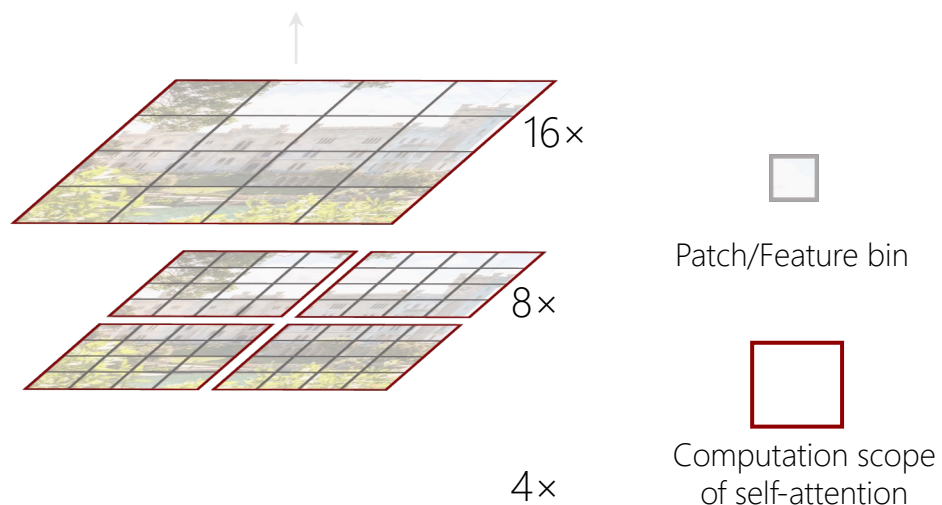
Fat cat is me.

Sensitive to absolute locations



# Hierarchical Vision Transformer

- Swin Transformer: Hierarchical Vision Transformer using Shifted Windows
  - Multi-scale (hierarchical) feature map
  - Use shifted window to get image patches



More friendly to visual signals  
multi-scale / locality / translation invariance

# Hybrid CNN-Transformer Model

- Formulation:
  - Micro-level Hybrid Model
    - Each basic block is designed with convolutional layer and self-attention layers
    - The model stack with several basic blocks repeatedly
  - Macro-level Hybrid Model
    - Each basic block is either pure CNN block or pure Transformer block
    - The model starts with few CNN blocks
    - The model further stacks several Transformer blocks
- Strength
  - Preserve local information due to CNN
  - Preserve translation invariance due to CNN
  - Efficient learning of spatial hierarchy due to CNN
  - Enable global interaction due to Transformer

# Examples

- Hybrid Transformer for visual recognition: Yang et al., MOAT, ICLR 2023
- Hybrid Transformer for medical AI: Chen et al., TransUNet, arXiv 2021
- Hybrid Transformer for vision-language: Chen et al., ViTamin, CVPR 2024

[1] Yang, et al. Moat: Alternating mobile convolution and attention brings strong vision models ICLR 2023

[2] Chen, et al. Transunet: Transformers make strong encoders for medical image segmentation. arXiv preprint arXiv:2102.04306, 2021.

[3] Chen, et al. Design scalable vision models in the vision-language era, CVPR 2024

# Hybrid Transformer for Visual Recognition

- Yang et al.  
MOAT ICLR2023

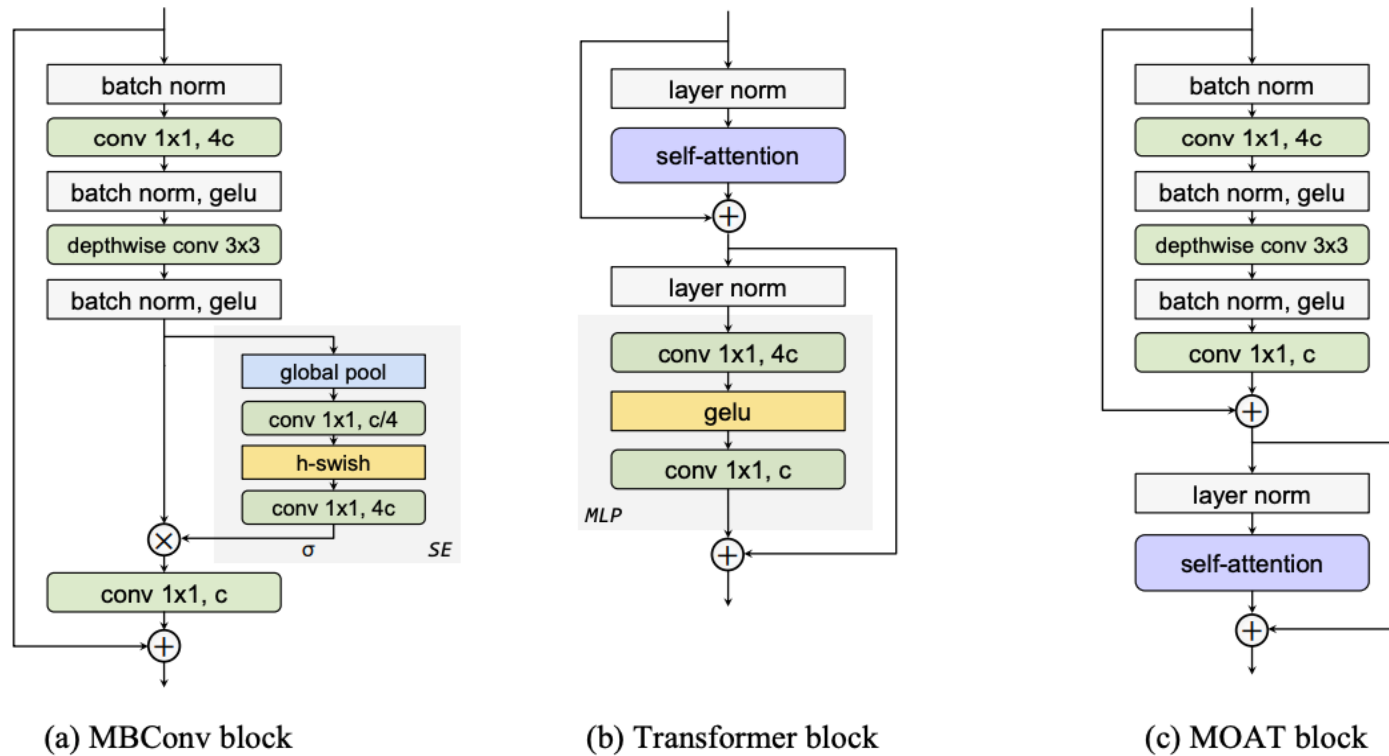
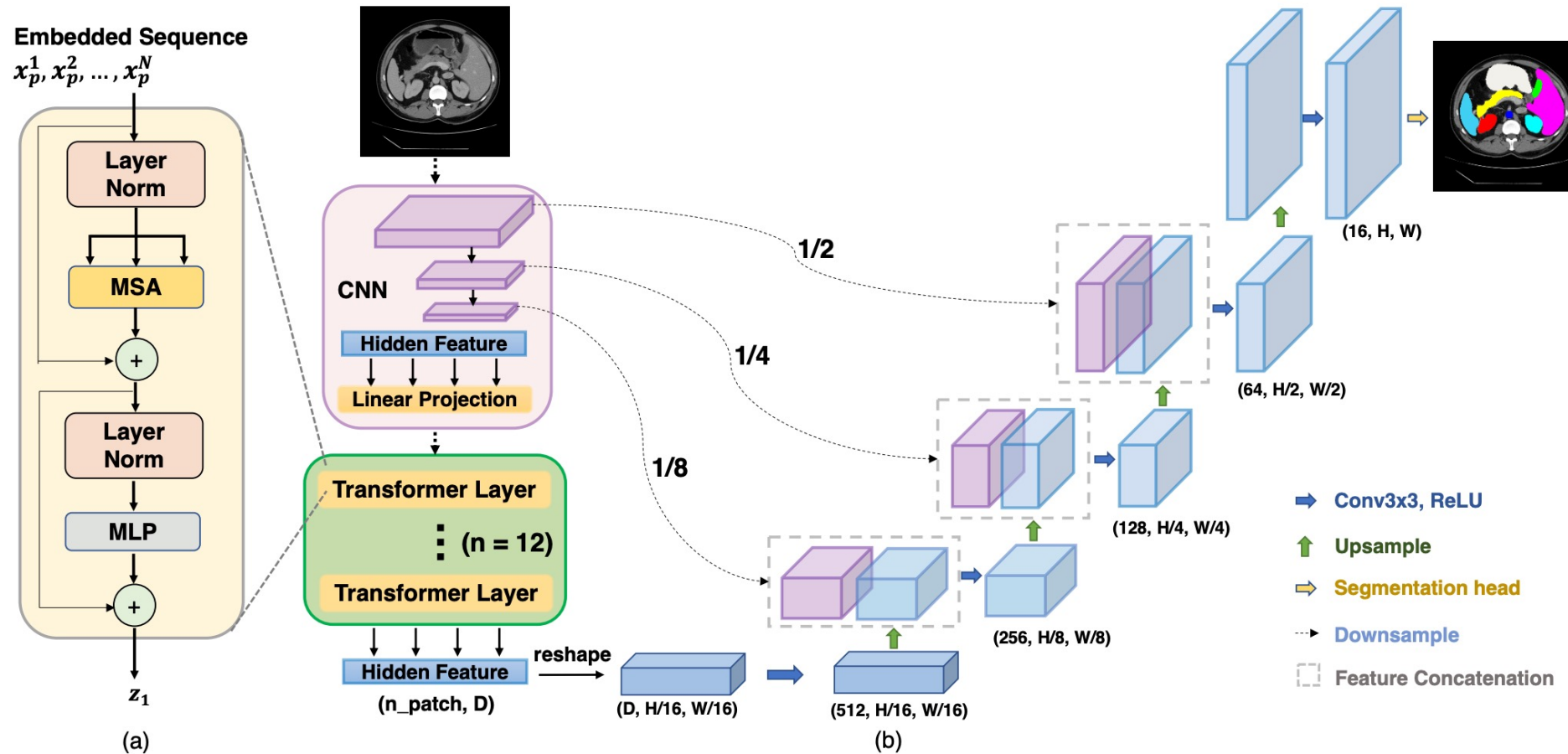


Figure 1: **Block comparison.** (a) The MBConv block (Sandler et al., 2018) employs the inverted bottleneck design with depthwise convolution and squeeze-and-excitation (Hu et al., 2018) applied to the expanded features. (b) The Transformer block (Vaswani et al., 2017) consists of a self-attention module and a MLP module. (c) The proposed MOAT block effectively combines them. The illustration assumes the input tensor has channels  $c$ .

# Hybrid Transformer + UNet: TransUNet

• Chen et al.  
2021



# Hybrid Transformer for Vision-Language

- Chen et al. CVPR 2024

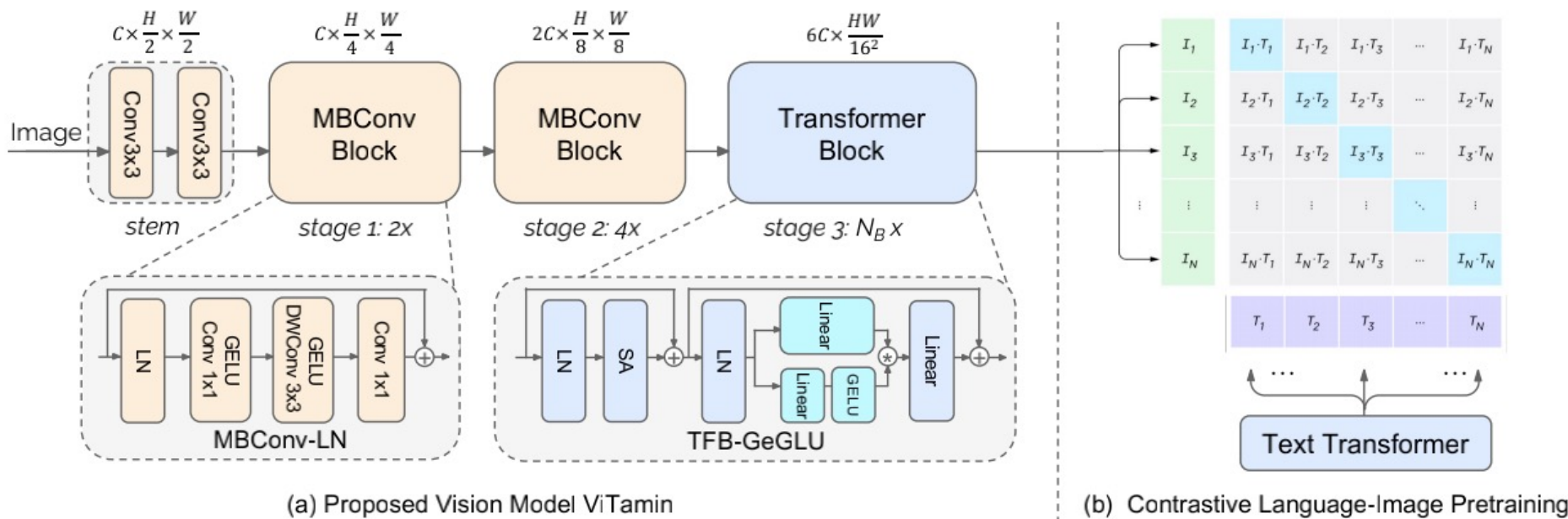


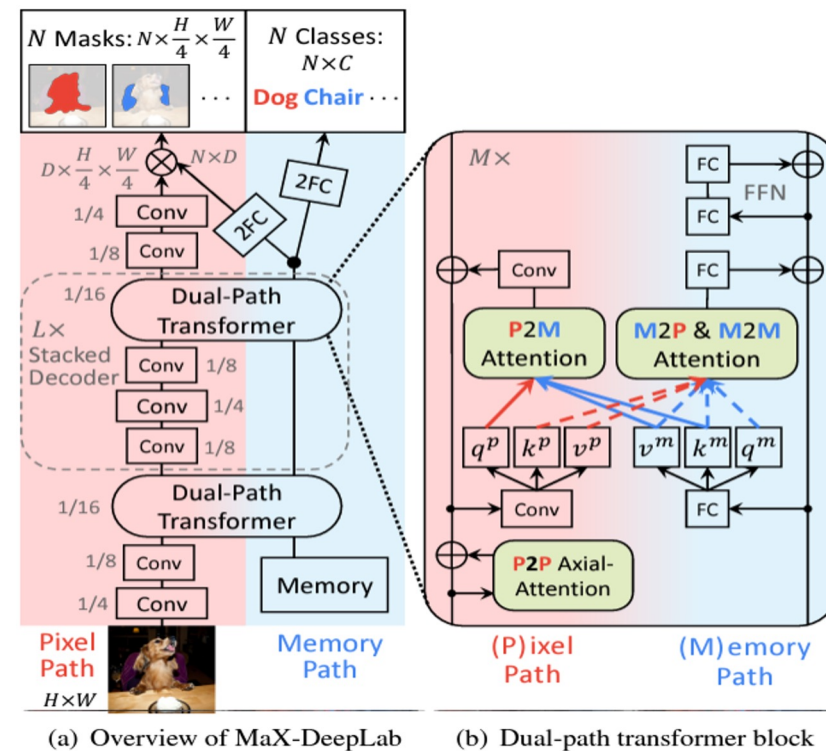
Figure 3. **Overview of ViTamin architecture.** (a) ViTamin begins with a convolutional stem, followed by Mobile Convolution Blocks (MBConv) in stage 1 and 2, and Transformer Blocks (TFB) in stage 3. The 2D input to the stage 3 is flattened to 1D. For the *macro-level* designs, the three-stage layout generates the final feature map with output stride 16, similar to ViT/16 [31]. We set channels sizes for the three stages to be  $(C, 2C, 6C)$ . For the *micro-level* designs, the employed MBConv-LN modifies MBConv [115] by using a single LayerNorm [4]. TFB-GeGLU upgrades TFB’s FFNs [131] (Feed-Forward Networks) with GELU Gated Linear Units [117]. (b) In the CLIP framework, given  $N$  image-text pairs, the vision model’s output  $I_i$  is learned to align with its corresponding text Transformer’s output  $T_i$ . Our text Transformers are the same as OpenCLIP [62]. +: Addition. \*: Multiplication.

# Summary

- Transformers consist of several Transformer blocks with multi-head self-attention layer and feedforward layers.
- It is highly scalable and highly parallelizable
- Faster training, larger models, better performance across vision and language tasks
- Good capability in vision tasks.

# Dual-Path Transformer: MaX-DeepLab

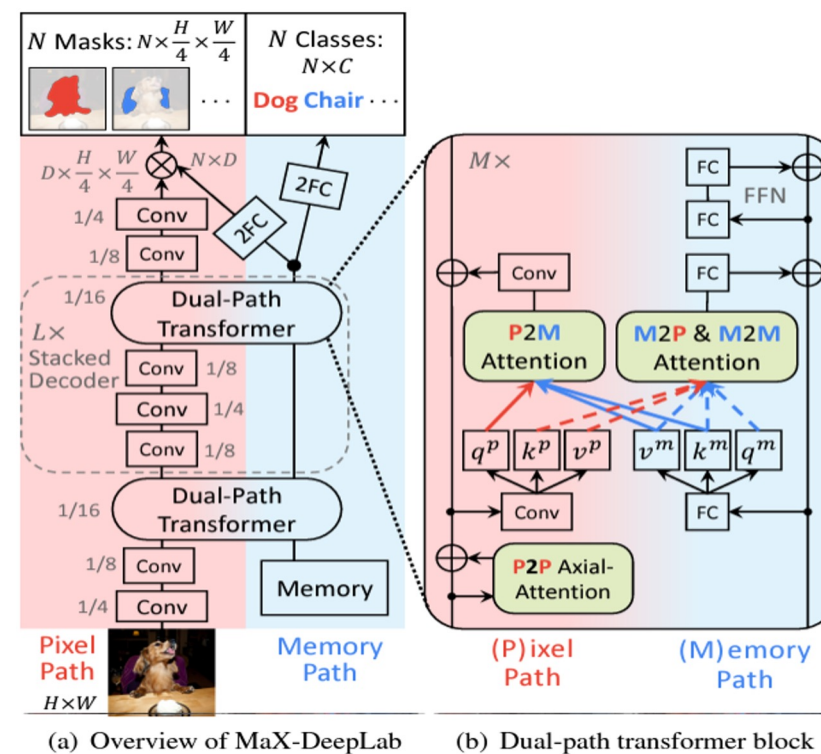
- Dual Path: Pixel Path + Memory Path
  - Memory to store global information
- Self-attention:
  - **P2P**: Pixel-to-pixel self-attention
    - Both query and key are pixel vector itself
  - **M2M**: Memory-to-memory self-attention
    - Both query and key are memory itself
- Cross-attention
  - **P2M**: Pixel-to-memory cross-attention
    - Query is pixel, key is memory
  - **M2P**: Memory-to-pixel cross-attention
    - Query is memory, key is pixel





# Dual-Path Transformer: MaX-DeepLab

- Dual Path: Pixel Path + Memory Path
- Prediction:  $N$  pairs of object class  $p_j^{class} \in \mathbb{R}^1$  and mask  $p_j^{mask} \in \mathbb{R}^{HW}$ 
  - Memory is decoded to  $N$  object class
  - The dot product of Memory and Pixel result in  $N$  object masks

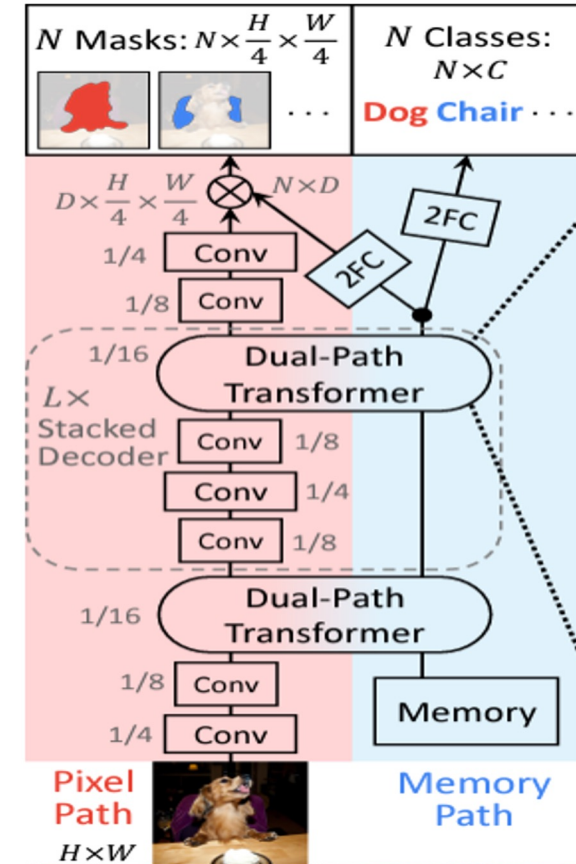


# Hungarian Matching for Prediction-Groundtruth Association

- Goal: match the groundtruth  $y_i$  to a prediction  $p_j$ 
  - Ground-truth:  $M$  pairs of object class  $y_i^{class} \in \mathbb{R}^1$  and mask  $y_i^{mask} \in \mathbb{R}^{HW}$
  - Prediction:  $N$  pairs of object class  $p_j^{class} \in \mathbb{R}^1$  and mask  $p_j^{mask} \in \mathbb{R}^{HW}$ 
    - Memory is decoded to  $N$  object class
    - The dot product of Memory and Pixel result in  $N$  object masks
- Example:
  - Successful match:  $y_i = \{\text{dog}, \text{img}_1\}$ ;  $p_j = \{\text{dog}, \text{img}_1\}$ 
    - Classification accuracy = 1, and mask IoU = 1
  - Failed match:  $y_i = \{\text{dog}, \text{img}_1\}$ ;  $p_j = \{\text{dog}, \text{img}_2\}$  -> low mask IoU
    - Classification accuracy = 1, but mask IoU = 0.1
  - Failed match:  $y_i = \{\text{dog}, \text{img}_1\}$ ;  $p_j = \{\text{cat}, \text{img}_1\}$  -> wrong classification
    - Mask IoU=1, but classification accuracy=0

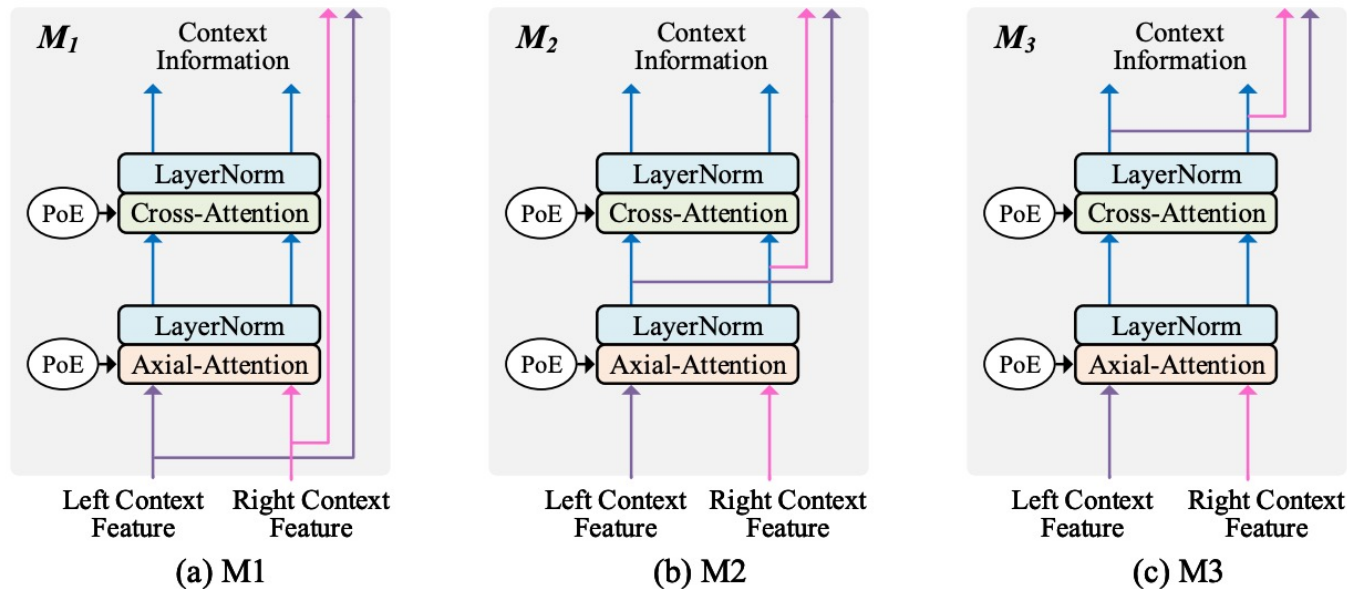
# Hungarian Matching for Prediction-Groundtruth Association

- Goal: match the groundtruth  $y_i$  to a prediction  $p_j$ 
  - Ground-truth:  $M$  pairs of object class  $y_i^{class}$  and mask  $y_j^{mask}$
  - Prediction:  $N$  pairs of object class  $p_i^{class}$  and mask  $p_j^{mask}$ 
    - **Memory** is decoded to  $N$  object class
    - The dot product of **Memory** and **Pixel** result in  $N$  object masks
- Define cost matrix  $C$  of size  $M \times N$ 
  - $C_{i,j} = -accuracy(y_i^{class}, p_j^{class}) - IoU(y_i^{mask}, p_j^{mask})$
- Linear assignment problem to do association
  - Formally, the task is to find an injection  $f: \{1, 2, \dots, m\} \rightarrow \{1, 2, \dots, n\}$  that minimizes the total assignment cost: Minimize  $\sum_{i=1}^M C_{i,f(i)}$

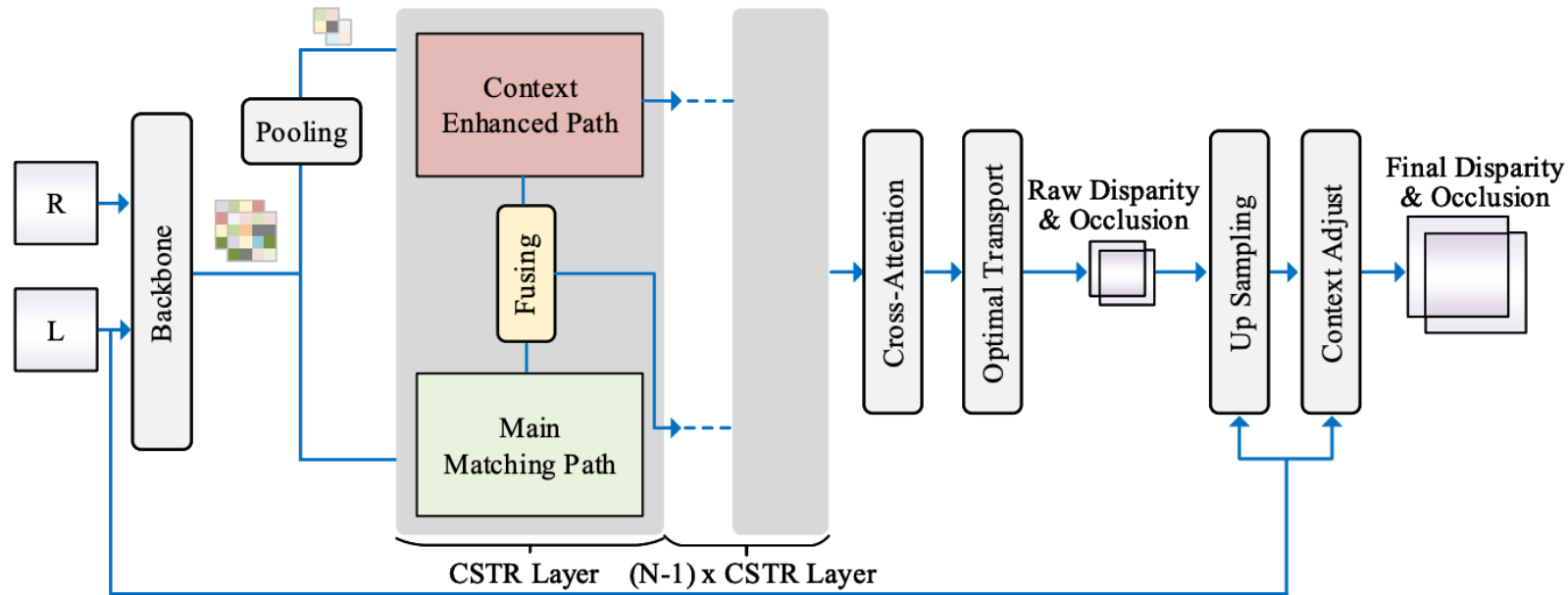


# Cross-attention for Stereo Matching / Fusion

- The context feature of left and right images are fused with cross-attention
- Query is left context features and key is the right context feature, or vice versa
- Reference to context-enhanced stereo Transformer [1] to estimate disparity



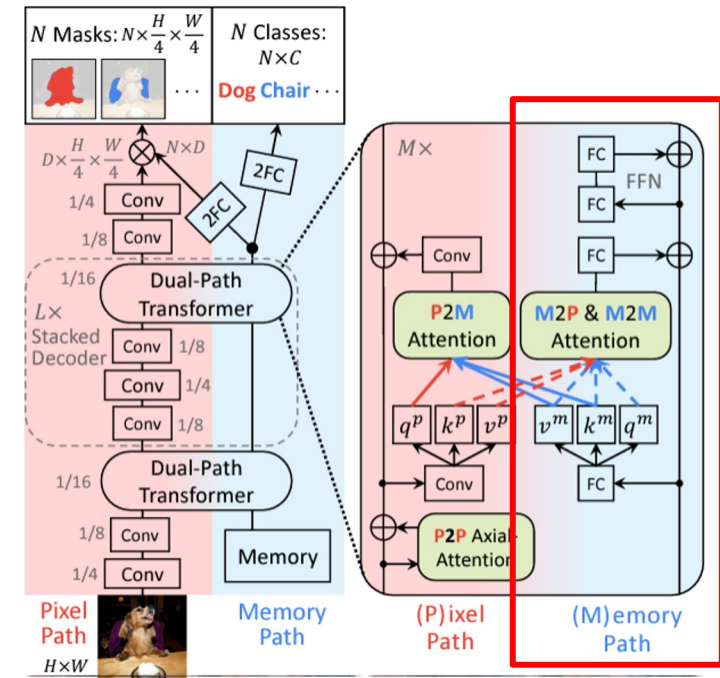
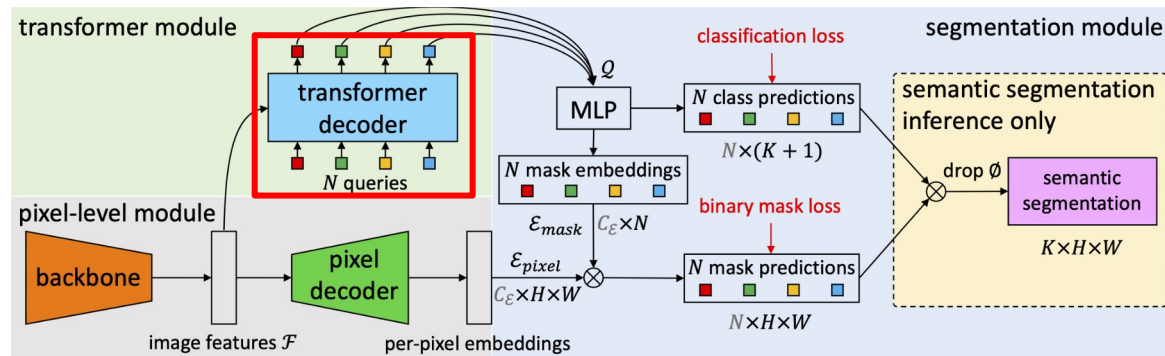
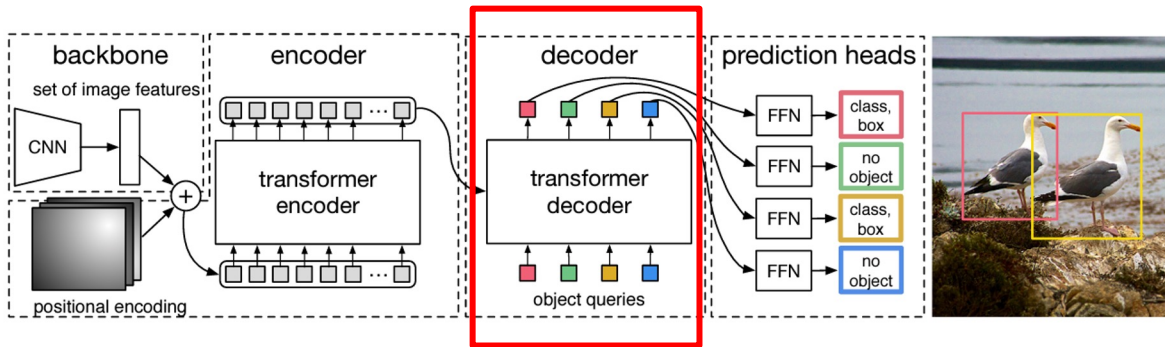
# Context Enhanced Stereo Transformer



**Fig. 3.** CSTR consists of two main components:(1) Context Enhanced Path that extracts long-range context information in low resolution feature. (2) Main Matching Path that use Axial-Attention to enhance context and Cross-Attention to compute raw disparity. Then a learnable Up Sampling block up restore the original scale of disparity and Context Adjustment block refines the disparity with context information across epipolar lines conditioned on the left image.

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Transformer decoder (cross-attention) is crucial in terms of an end-to-end problem formulation:



(a) Overview of MaX-DeepLab

(b) Dual-path transformer block

Carion, et al. End-to-end object detection with transformers. ECCV 2020.

Wang, et al. MaX-DeepLab: End-to-End Panoptic Segmentation With Mask Transformers. CVPR 2021.

Cheng., et al. Per-Pixel Classification is Not All You Need for Semantic Segmentation. NeurIPS 2021.

# **CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation**

Common pipeline:

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Common pipeline:

- Image feature extraction through a backbone



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Common pipeline:

- Image feature extraction through a backbone

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Common pipeline:

- Image feature extraction through a backbone
- Image feature enhancement through transformer encoder (self-attention)

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Common pipeline:

- Image feature extraction through a backbone
- Image feature enhancement through transformer encoder (self-attention)
- learnable object queries aggregate pixel features through transformer decoder (cross-attention) -> box/mask embedding

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Common pipeline:

- Image feature extraction through a backbone
- Image feature enhancement through transformer encoder (self-attention)
- learnable object queries aggregate pixel features through transformer decoder (cross-attention) -> box/mask embedding
- hungarian matching for supervision assignments

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

Common pipeline:

- Image feature extraction through a backbone
- Image feature enhancement through transformer encoder (self-attention)
- **learnable object queries aggregate pixel features through transformer decoder (cross-attention) -> box/mask embedding**
- hungarian matching for supervision assignments

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

## Cross-Attention:

- Affinity logits are computed, with linear projections, between queries ( $\mathbf{Q}^c$ ) and pixels ( $\mathbf{K}^p$ )
- A spatial-wise (HW) softmax is applied to convert the affinity logits map to attention weights
- The attention is multiplied with pixel features, with linear project  
$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$
- The update of queries is added in a residual manner.

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

## Cross-Attention:

- Affinity logits are computed, with linear projections, between queries ( $Q^c$ ) and pixels ( $K^p$ )

- A spatial-wise (HW) softmax is applied to convert the affinity logits map to attention weights

- The attention is multiplied with pixel features, with a linear projection

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

- The update of queries is added in a residual manner.

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

## Cross-Attention:

- Affinity logits are computed, with linear projections, between queries ( $Q^c$ ) and pixels ( $K^p$ )
- A spatial-wise (HW) softmax is applied to convert the affinity logits map to attention weights

- The attention linear project

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

- The update of queries is added in a residual manner.



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

## Cross-Attention:

- Affinity logits are computed, with linear projections, between queries ( $Q^c$ ) and pixels ( $K^p$ )
- A spatial-wise (HW) softmax is applied to convert the affinity logits map to attention weights
- The attention linear project  $\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p$ , pixel features, with
- The update of queries is added in a residual manner.

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

## Cross-Attention:

- Affinity logits are computed, with linear projections, between queries ( $\mathbf{Q}^c$ ) and pixels ( $\mathbf{K}^p$ )
- A spatial-wise (HW) softmax is applied to convert the affinity logits map to attention weights
- The attention is multiplied with pixel features, with a linear projection: 
$$\hat{\mathbf{C}} = \mathbf{C} + \text{softmax}_{HW}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$
- The update of queries is added in a residual manner.

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

## Cross-Attention:

- Affinity logits are computed, with linear projections, between queries ( $Q^c$ ) and pixels ( $K^p$ )
- A spatial-wise (HW) softmax is applied to convert the affinity logits map to attention weights

- The attention is applied to pixel features, with linear project
- $$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

As a result, the queries will be updated and converted to correspond to a specific object in prediction.

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

- Inconsistency of using object queries for mask prediction and updating object queries
- Sparse attention map

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

- Inconsistency of using object queries for mask prediction and updating object queries

mask prediction:

$$\mathbf{Z} = \underset{\text{mask probs}}{\text{softmax}} \underset{\text{affinity logits}}{\mathbf{F} \times \mathbf{C}^T},$$

The equation shows the mask prediction process. A horizontal line is drawn under the expression. Below the line, two red brackets are positioned: one under the  $\text{softmax}$  function and another under the  $\mathbf{F} \times \mathbf{C}^T$  product. The text "mask probs" is centered under the first bracket, and "affinity logits" is centered under the second bracket.

- Sparse attention map

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

- Inconsistency of using object queries for mask prediction and updating object queries

*mask prediction:* 
$$\mathbf{Z} = \underset{N}{\text{softmax}}(\mathbf{F} \times \mathbf{C}^T),$$

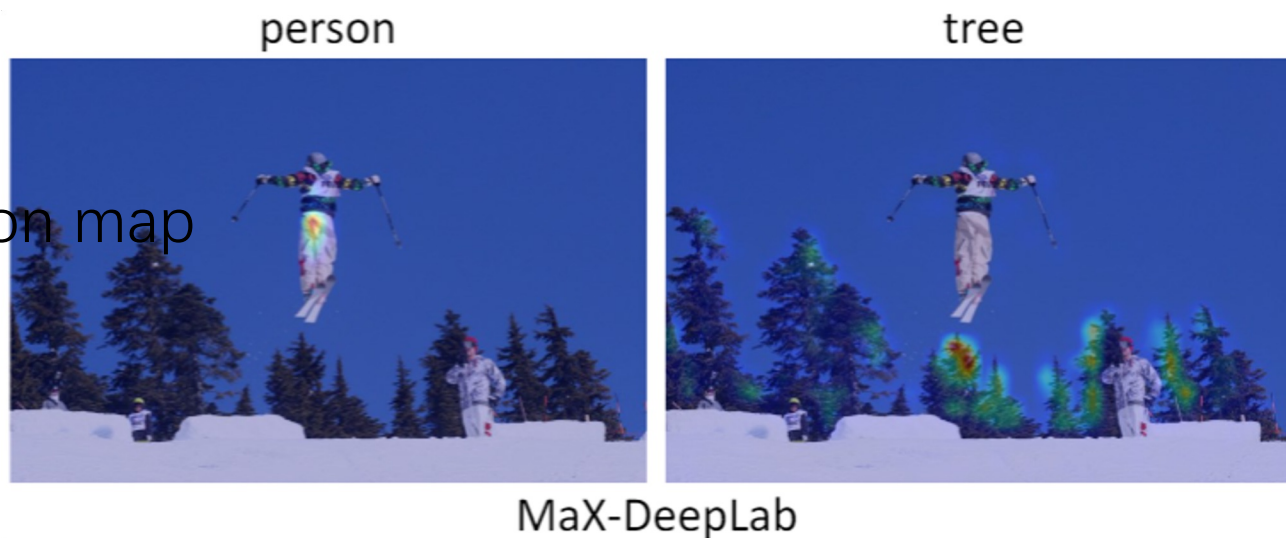
*updating object queries:* 
$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

- Sparse attention map

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

- Inconsistency of using memory queries for mask prediction and updating memory queries

- Sparse attention map



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

Cross-attention -> a clustering process

object queries -> cluster centers

attention map -> clustering assignment

updating object queries -> updating clustering center



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

Cross-attention -> a clustering process

object queries -> cluster centers

attention map -> clustering assignment

updating object queries -> updating clustering center

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}(\mathbf{Q}^c \times (\mathbf{K}^p)^T)} \times \mathbf{V}^p,$$

Cross-attention -> a clustering process

object queries -> cluster centers

attention map -> clustering assignment

updating object queries -> updating clustering center

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

Cross-attention -> a clustering process

object queries -> cluster centers

attention map -> clustering assignment

updating object queries -> updating clustering center

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

*Machine Translation:* Each object query corresponds to a word in target language, and it will be assigned to one most affiliated word in source language as its update

*Panoptic Segmentation:* Each object query corresponds to an object in prediction, and it will be assigned to one most affiliated pixel as its update

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

*Machine Translation:* Each object query corresponds to a word in target language, and it will be assigned to one most affiliated word in source language as its update

*Panoptic Segmentation:* Each object query corresponds to an object in prediction, and it will be assigned to one most affiliated pixel as its update

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

*Machine Translation:* Each object query corresponds to a word in target language, and it will be assigned to one most affiliated word in source language as its update

*Panoptic Segmentation:* Each object query corresponds to an object in prediction, and it will be assigned to one most affiliated pixel as its update  
**objects assigned to pixels?**

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

*Machine Translation:* Each object query corresponds to a word in target language, and it will be assigned to one most affiliated word in source language as its update

*Panoptic Segmentation:* Each object query corresponds to an object in prediction, and it will be assigned to one most affiliated pixel as its update

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

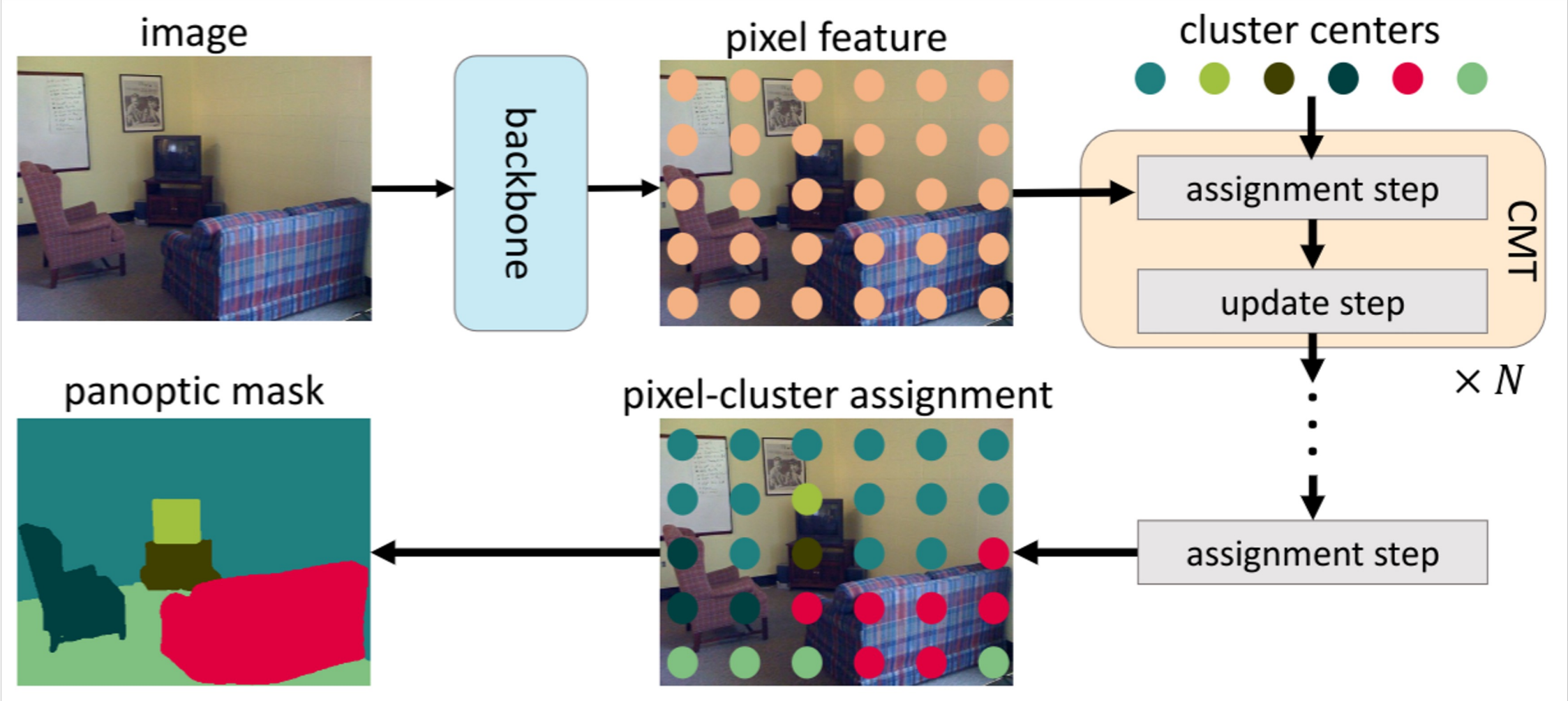
$$\hat{\mathbf{C}} = \mathbf{C} + \left( \underset{N}{\text{softmax}}(\tilde{\mathbf{K}}^p \times (\tilde{\mathbf{Q}}^c)^T) \right)^T \times \mathbf{V}^p.$$

*Panoptic Segmentation:* Each object query corresponds to an object in prediction, and it will be assigned to one most affiliated pixel as its update

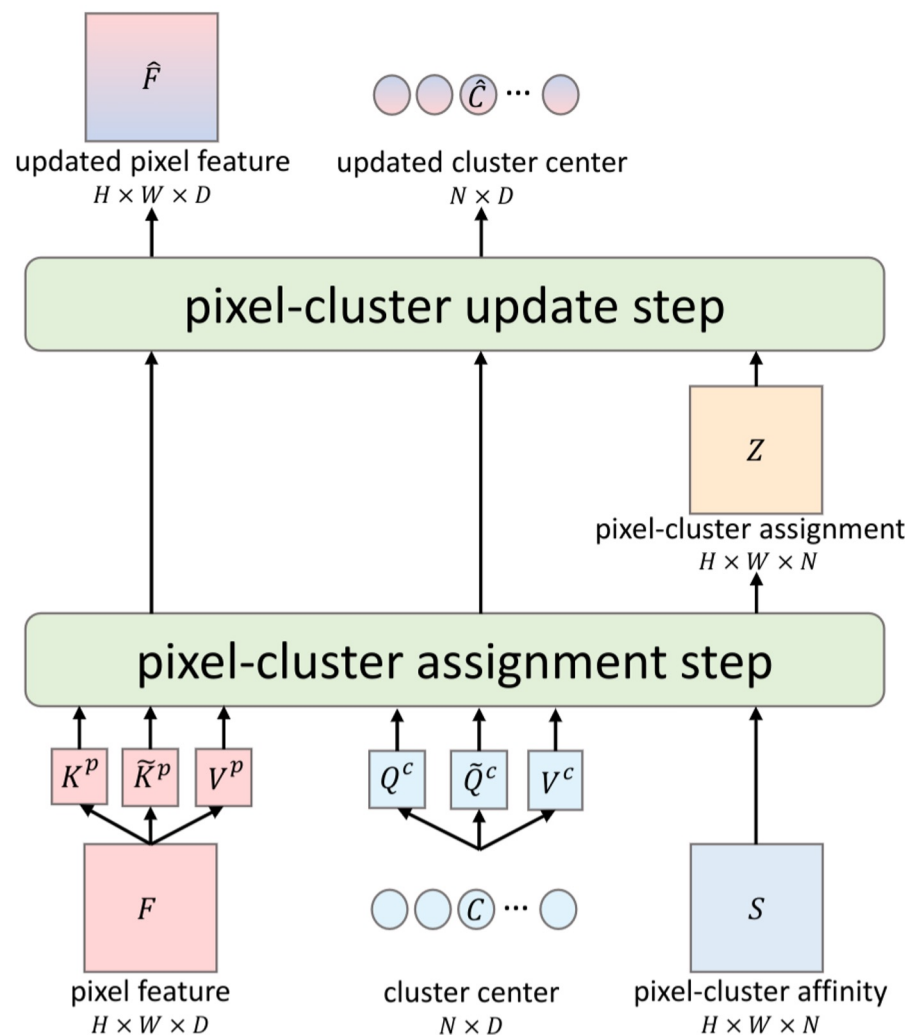
*Panoptic Segmentation:* Each pixel will choose one most affiliated object, all assigned pixels will serve as an update to corresponding object query



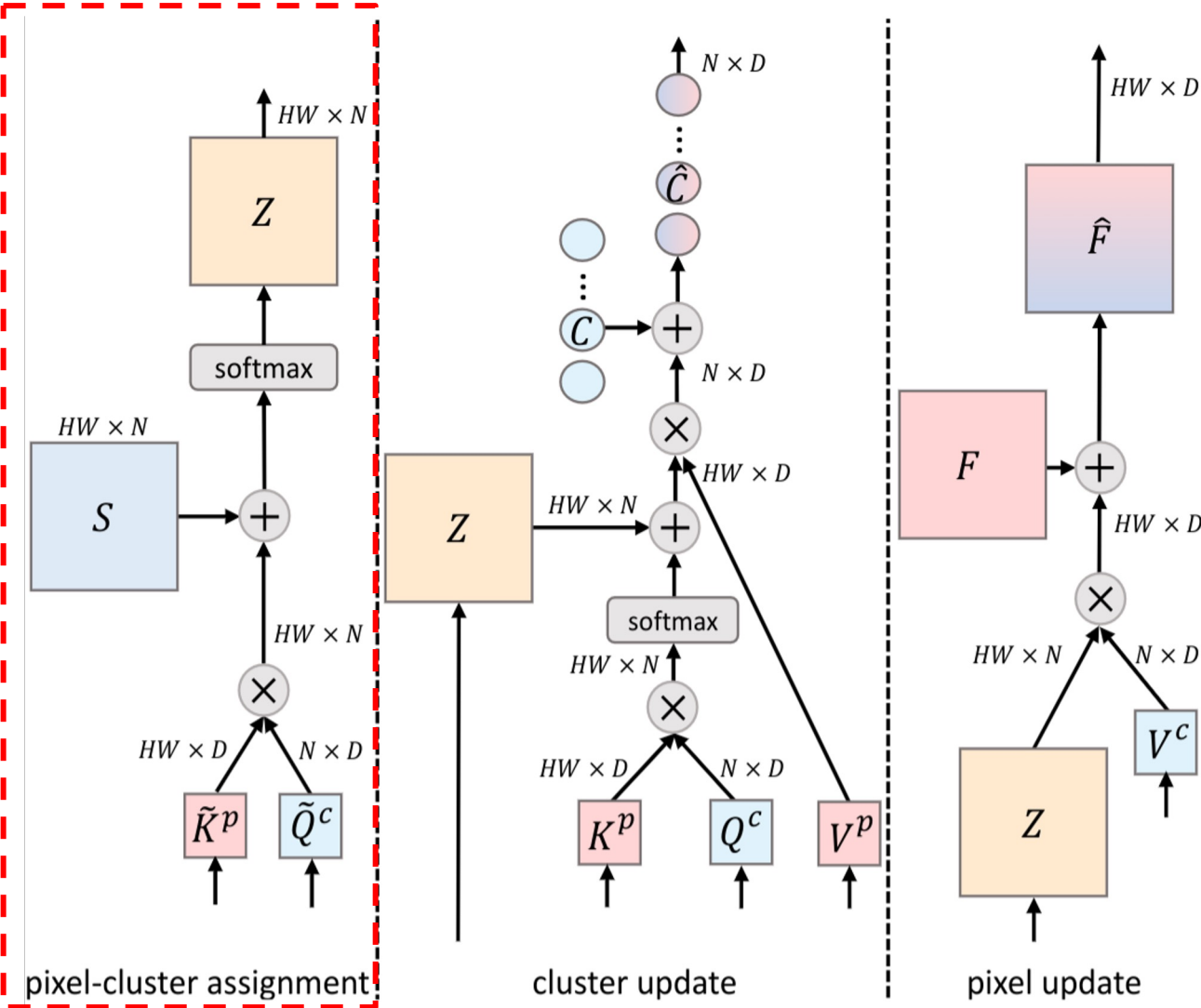
# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation



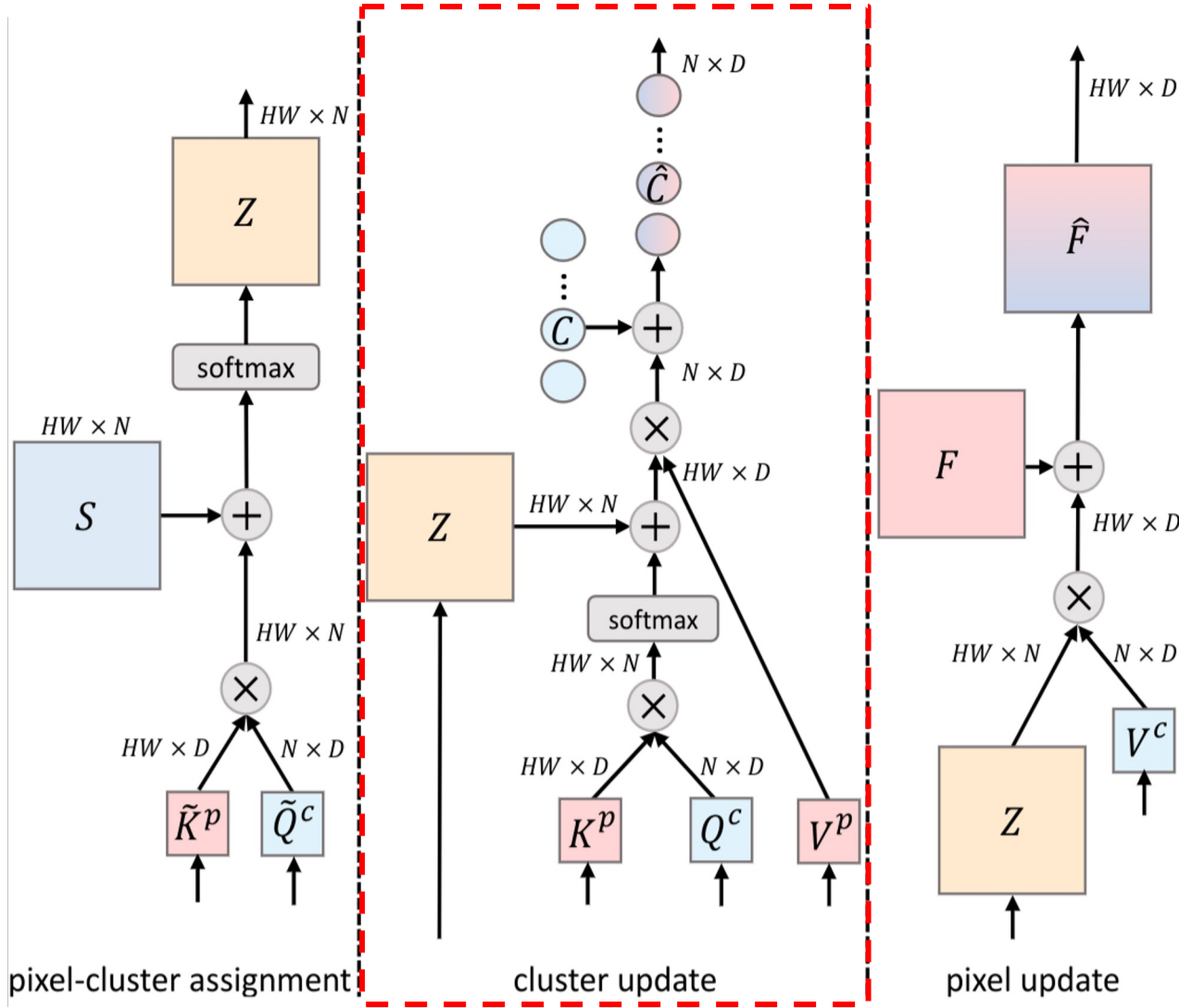
# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation



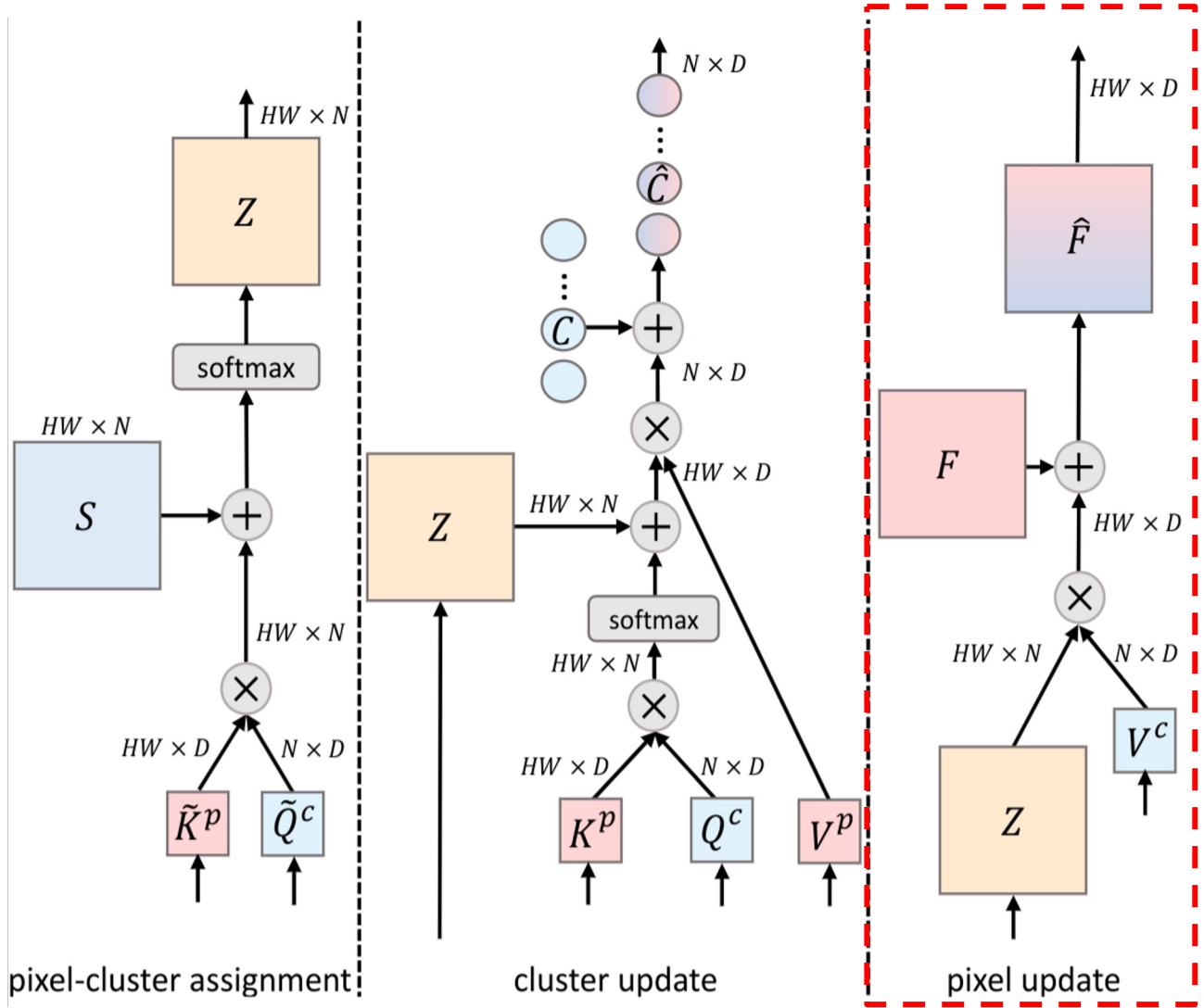
# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation



# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

- Include coordinates into the clustering
- Improve instance discrimination loss in MaX-DeepLab to pixel-wise contrastive loss

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation

We build CMT-DeepLab upon previous SOTA method MaX-DeepLab:

COCO Panoptic:

+4.2% *val* PQ

+4.4% *test* PQ

Cityscapes

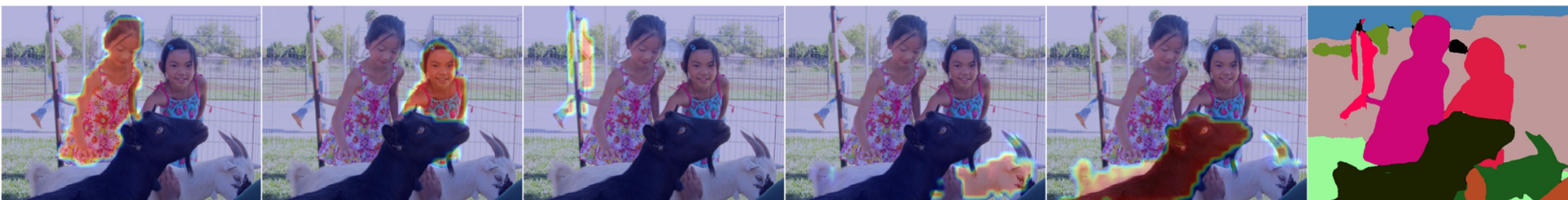
Panoptic:

+2.9% *val* PQ

# CMT-DeepLab: Clustering Mask Transformers for Panoptic Segmentation



MaX-DeepLab



CMT-DeepLab



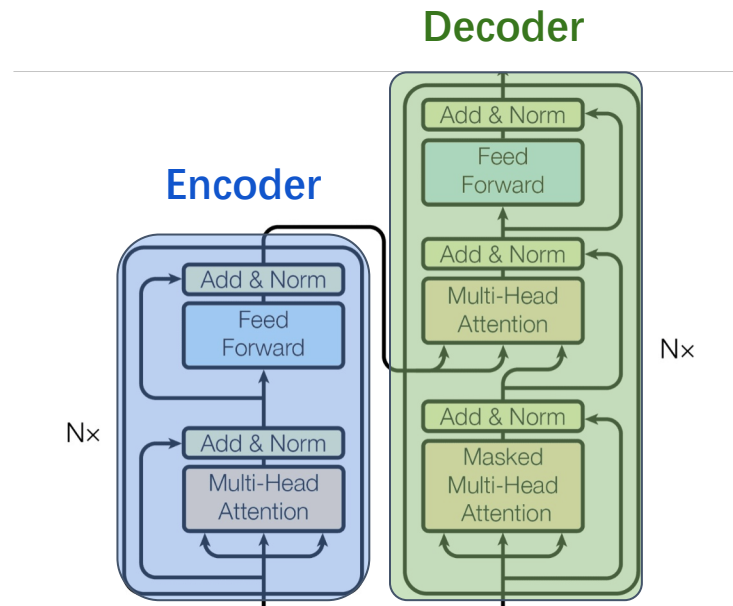
# Summary

- A novel clustering view to better understand and design Transformer modules for object-centric representation learning
- We introduce CMT-DeepLab, which unifies cross-attention and panoptic segmentation from a clustering perspective

# Overview

## Extracting pixel-level representation:

- Vision Transformer (ViT)



## Converting to object-level representation:

- CMT-DeepLab [2] (CVPR 22 *Oral*)
- **kMaX-DeepLab [3] (ECCV 22)**

# *k*-means Mask Transformer

Cross-attention:

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

# *k*-means Mask Transformer

Cross-attention:

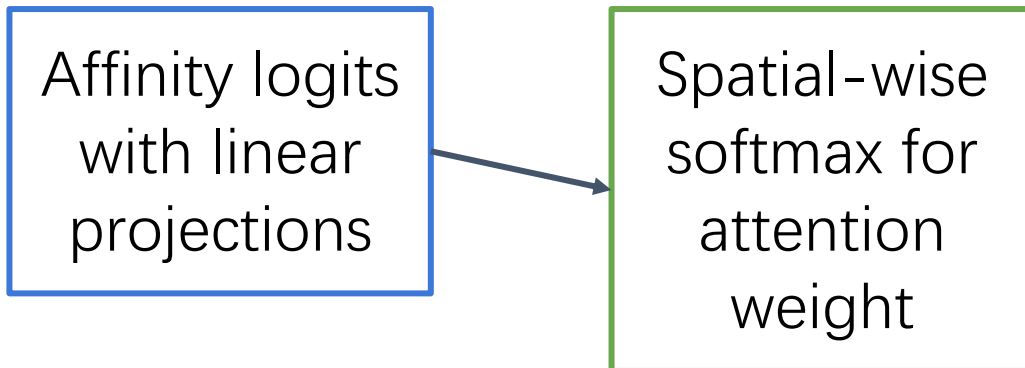
$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

Affinity logits  
with linear  
projections

# *k*-means Mask Transformer

Cross-attention:

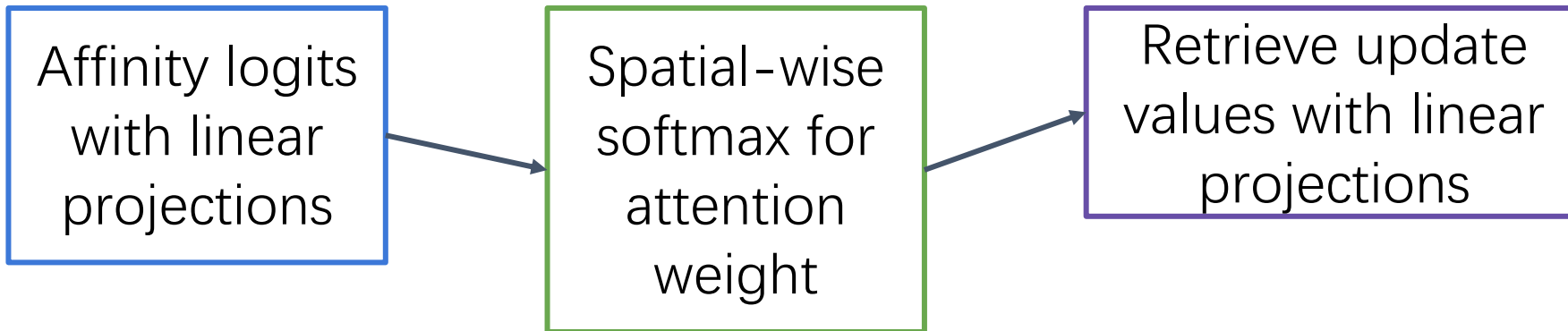
$$\hat{\mathbf{C}} = \mathbf{C} + \text{softmax}_{HW}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$



# *k*-means Mask Transformer

Cross-attention:

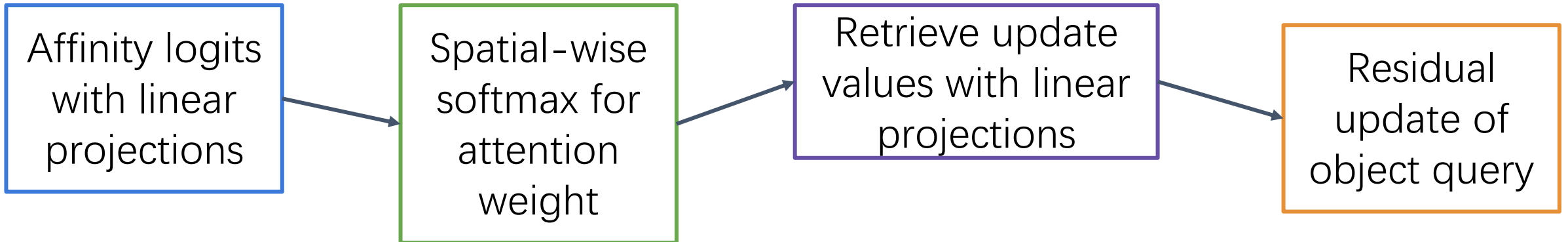
$$\hat{\mathbf{C}} = \mathbf{C} + \text{softmax}_{HW}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$



# *k*-means Mask Transformer

Cross-attention:

$$\hat{\mathbf{C}} = \mathbf{C} + \underset{HW}{\text{softmax}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$



# *k*-means Mask Transformer

k-means clustering algorithm:

$$\mathbf{A} = \underset{N}{\operatorname{argmax}}(\mathbf{C} \times \mathbf{P}^T),$$

$$\hat{\mathbf{C}} = \mathbf{A} \times \mathbf{P},$$

---



# *k*-means Mask Transformer

k-means clustering algorithm:

$$\mathbf{A} = \underset{N}{\operatorname{argmax}}(\mathbf{C} \times \mathbf{P}^T),$$

$$\hat{\mathbf{C}} = \mathbf{A} \times \mathbf{P},$$

---

Affinity logits

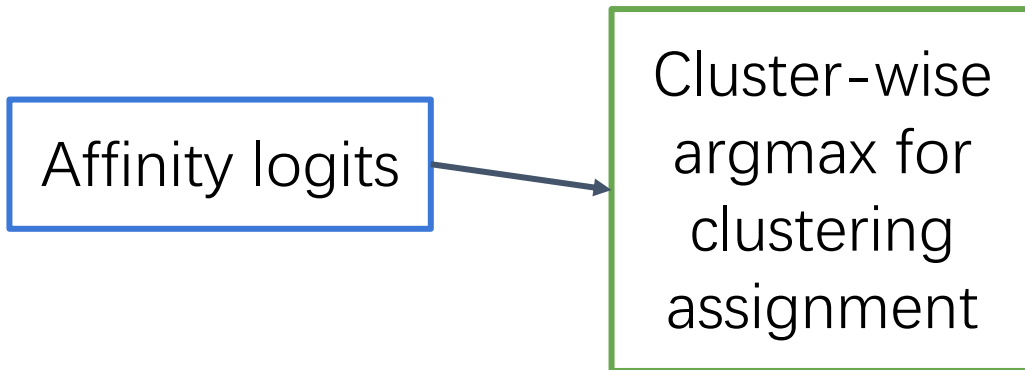
# *k*-means Mask Transformer

k-means clustering algorithm:

$$\mathbf{A} = \underset{N}{\operatorname{argmax}}(\mathbf{C} \times \mathbf{P}^T),$$

$$\hat{\mathbf{C}} = \mathbf{A} \times \mathbf{P},$$

---

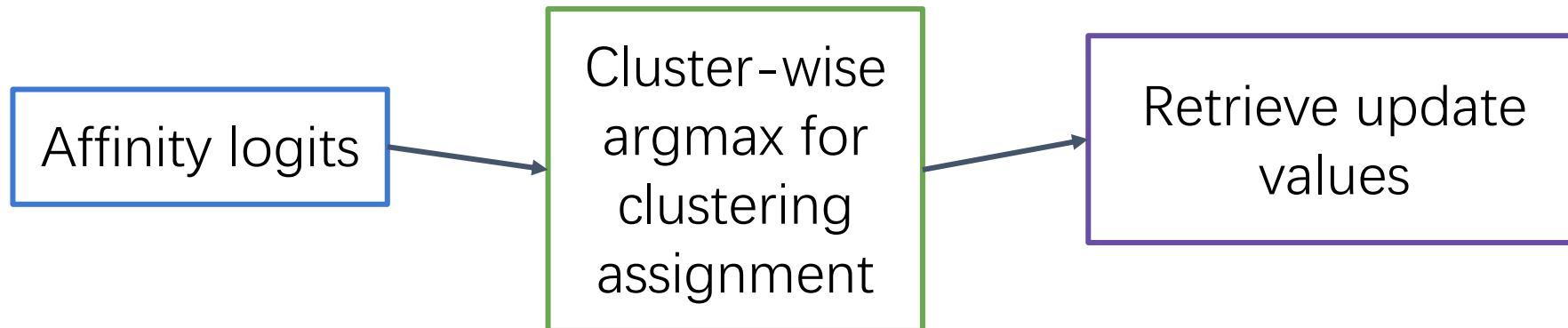


# *k*-means Mask Transformer

k-means clustering algorithm:

$$\mathbf{A} = \underset{N}{\operatorname{argmax}}(\mathbf{C} \times \mathbf{P}^T),$$

$$\hat{\mathbf{C}} = \mathbf{A} \times \mathbf{P},$$

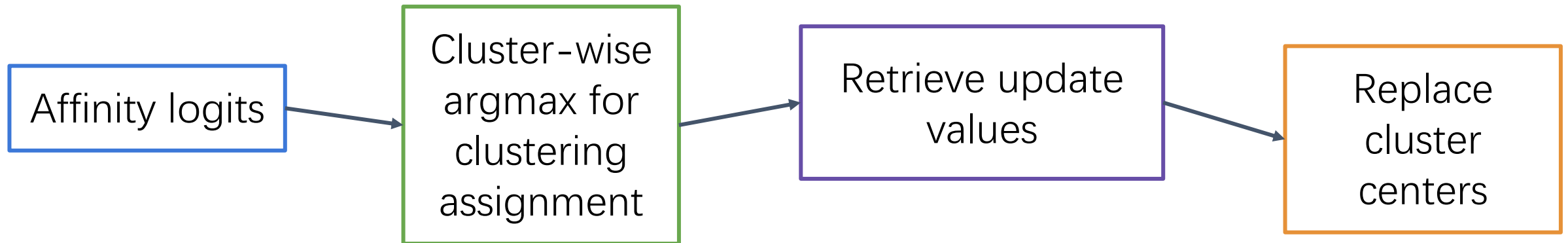


# *k*-means Mask Transformer

k-means clustering algorithm:

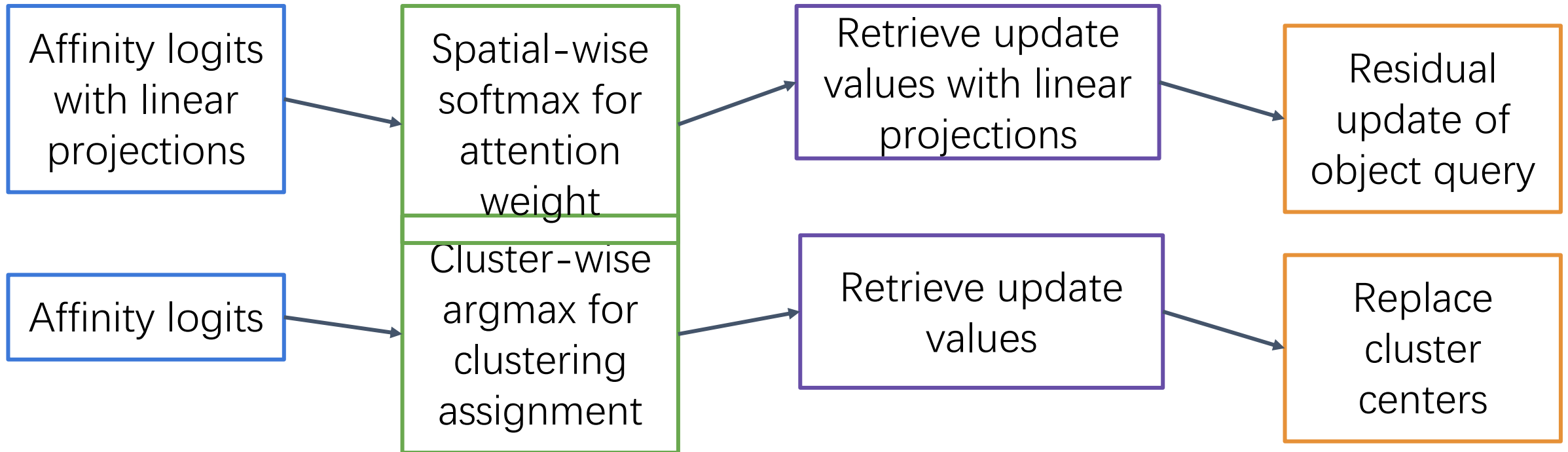
$$\mathbf{A} = \underset{N}{\operatorname{argmax}}(\mathbf{C} \times \mathbf{P}^T),$$

$$\hat{\mathbf{C}} = \mathbf{A} \times \mathbf{P},$$



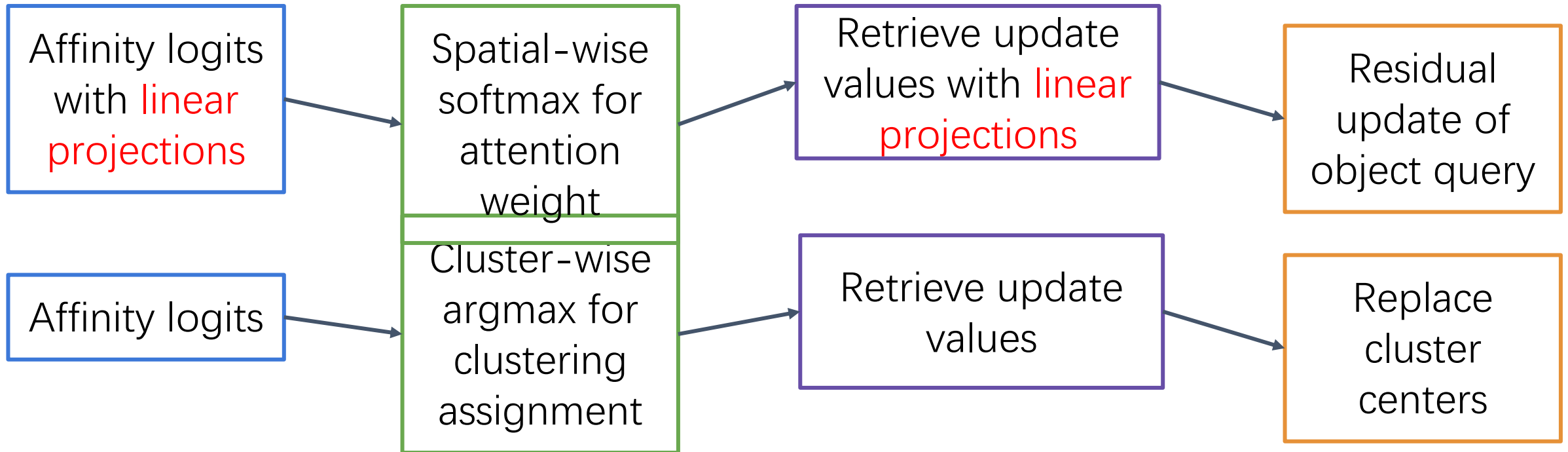
# *k*-means Mask Transformer

cross-attention v.s. k-means clustering algorithm:



# *k*-means Mask Transformer

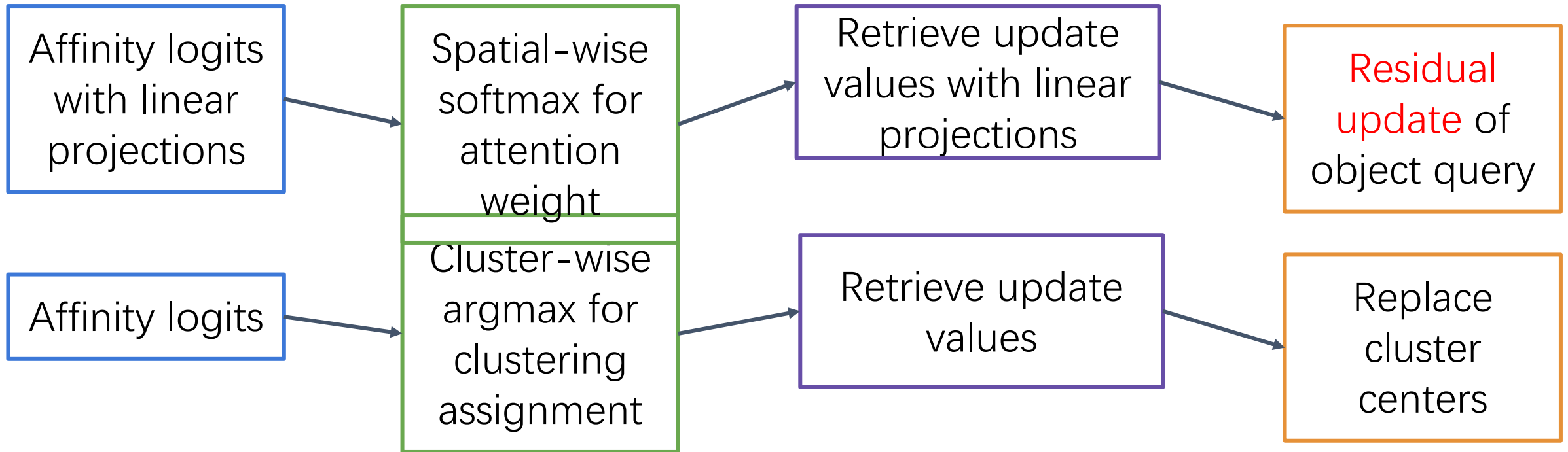
cross-attention v.s. k-means clustering algorithm:



Linear  
projections

# *k*-means Mask Transformer

Cross-attention v.s. *k*-means clustering algorithm:

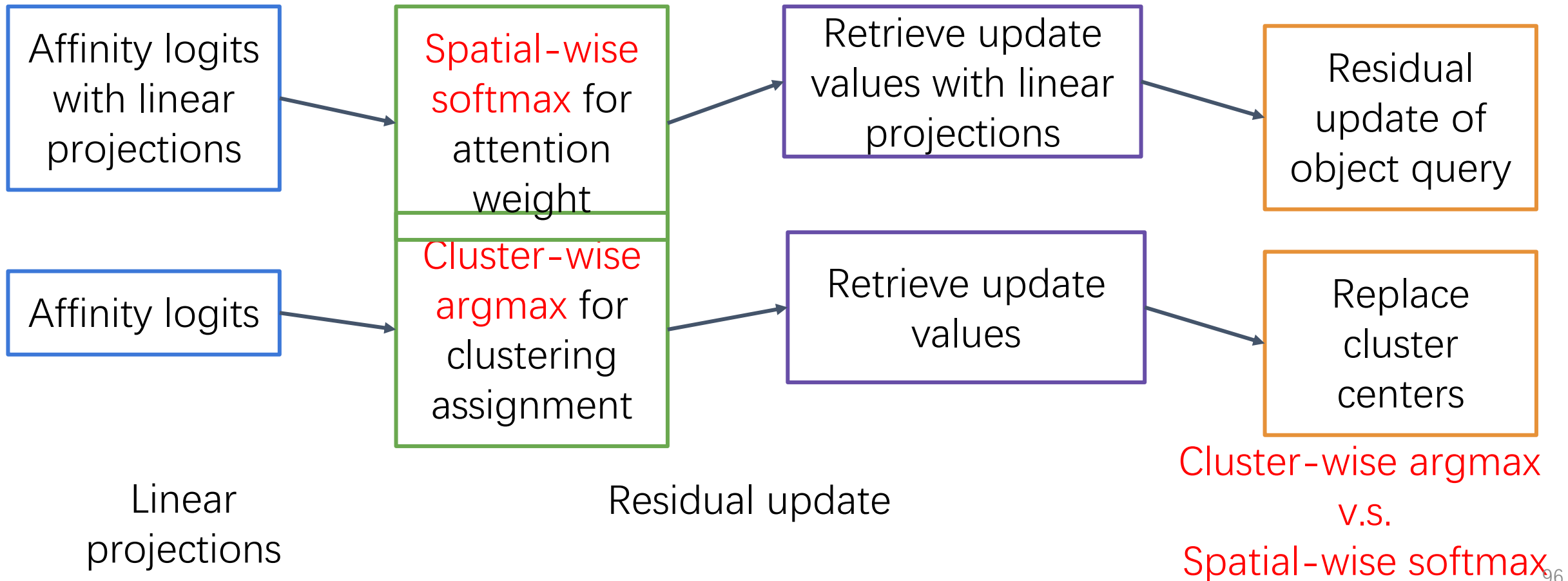


Linear  
projections

Residual update

# *k*-means Mask Transformer

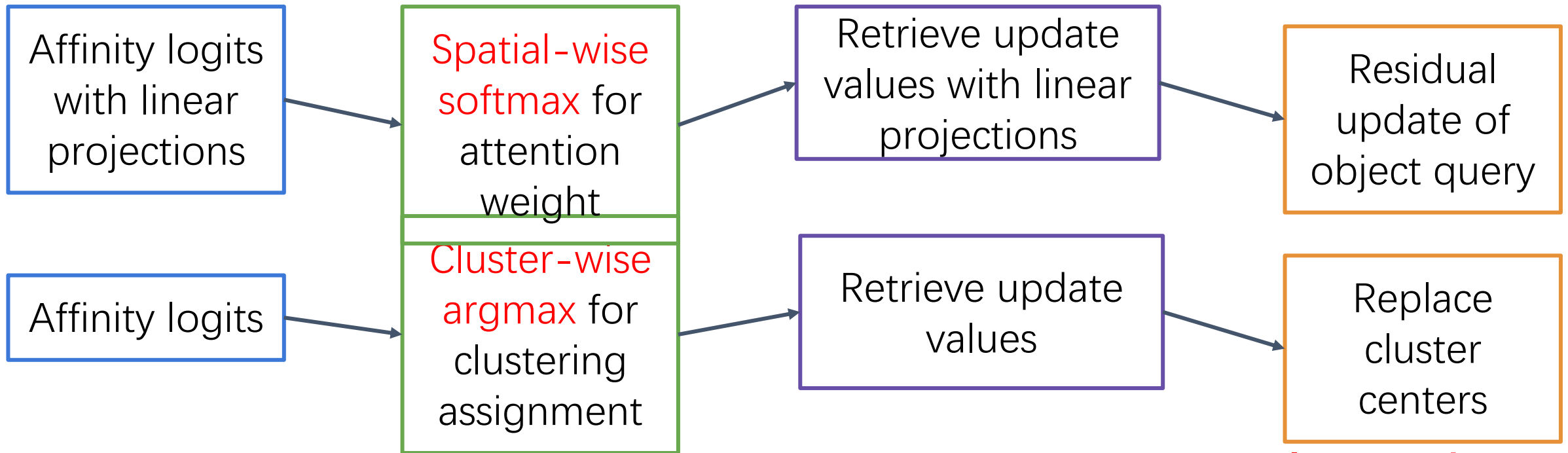
cross-attention v.s. k-means clustering algorithm:





# *k*-means Mask Transformer

cross-attention v.s. k-means clustering algorithm:



Linear projections

Residual update

**Cluster-wise  
argmax  
v.s.  
Spatial-wise**

# *k*-means Mask Transformer

A simple *change* for k-means cross-attention:

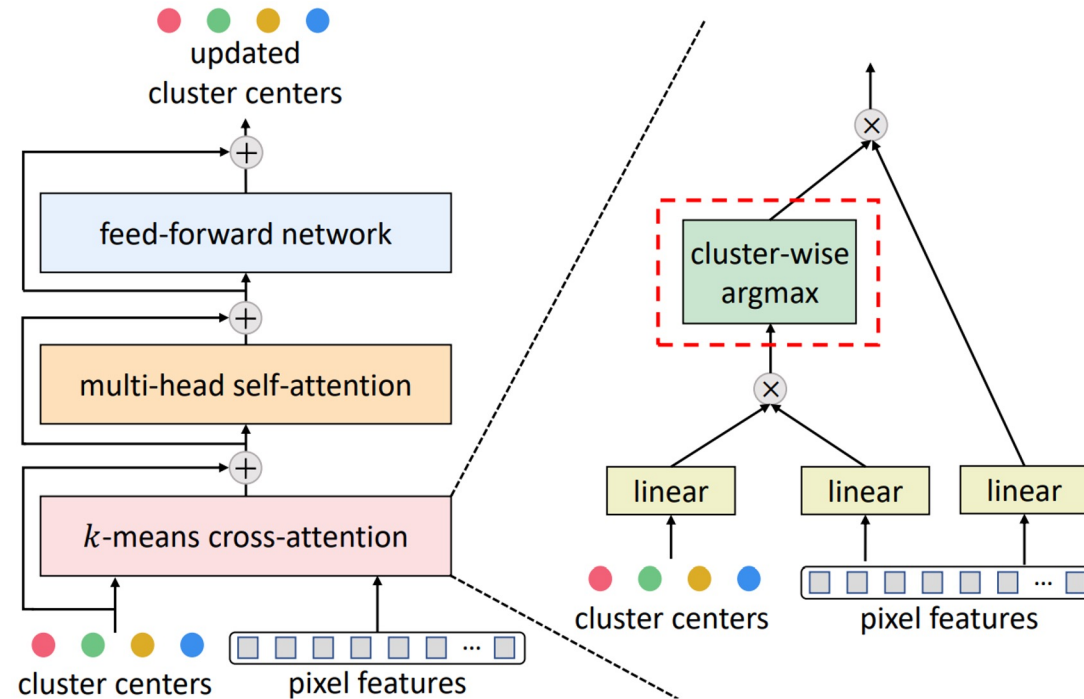
---

$$\hat{\mathbf{C}} = \mathbf{C} + \boxed{\text{softmax}_{HW}}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p,$$

$$\hat{\mathbf{C}} = \mathbf{C} + \boxed{\text{argmax}_N}(\mathbf{Q}^c \times (\mathbf{K}^p)^T) \times \mathbf{V}^p.$$

# *k*-means Mask Transformer

A simple *change* for k-means cross-attention:



# *k*-means Mask Transformer

A simple *change* for k-means cross-attention:

pixel-cluster interaction module	ResNet-50			MaX-S		
	params	FLOPs	PQ	params	FLOPs	PQ
cross-attention [89]	56M	165G	47.5	73M	237G	52.0
dual-path cross-attention [92]	58M	175G	48.0	75M	247G	52.3
<i>k</i> -means cross-attention	57M	168G	52.7	74M	240G	56.1
dual-path <i>k</i> -means cross-attention	59M	176G	53.0	76M	248G	56.2

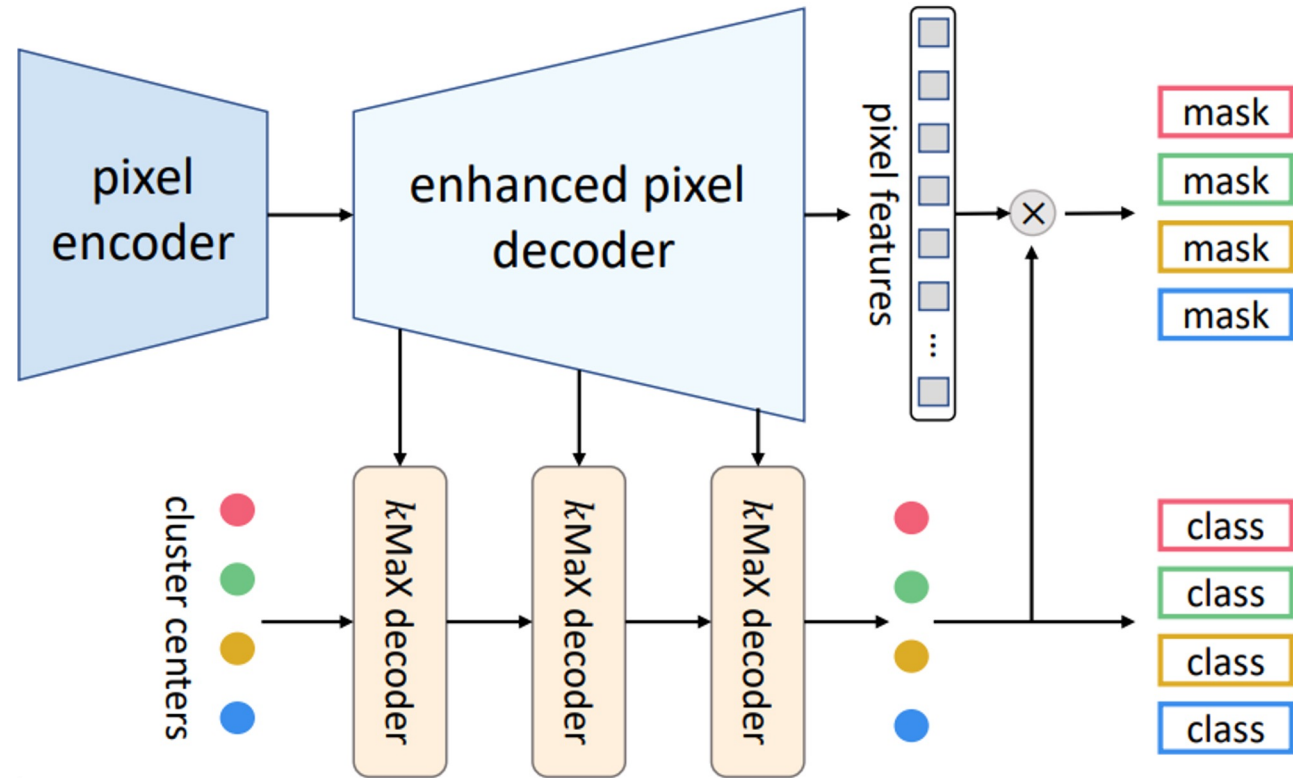
# *k*-means Mask Transformer

A simple change for *k*-means cross-attention:

pixel-cluster interaction module	ResNet-50			MaX-S		
	params	FLOPs	PQ	params	FLOPs	PQ
cross-attention [89]	56M	165G	47.5	73M	237G	52.0
dual-path cross-attention [92]	58M	175G	48.0	75M	247G	52.3
<i>k</i> -means cross-attention	57M	168G	52.7	74M	240G	56.1
dual-path <i>k</i> -means cross-attention	59M	176G	53.0	76M	248G	56.2

**5.2% (47.5% -> 52.7%) PQ improvement with one change and negelectable extra cost**

# *k*-means Mask Transformer



# *k-means* Mask Transformer

Method	params	FLOPs	FPS	PQ
<a href="#">MaX-DeepLab</a>	451M	3692G	-	51.1% (-6.9%)
<a href="#">MaskFormer</a>	212M	792G	5.2	52.7% (-5.3%)
<a href="#">K-Net</a>	-	-	-	54.6% (-3.4%)
<a href="#">CMT-DeepLab</a>	270M	1114G	3.2	55.3% (-2.7%)
<b>kMaX-DeepLab</b>	<b>232M</b>	<b>749G</b>	<b>6.6</b>	<b>58.0%</b>

COCO val/set

# *k*-means Mask Transformer

Method	params	FLOPs	FPS	PQ	AP <sup>mask</sup>	mIoU
<a href="#">Panoptic-DeepLab</a>	47M	548G	5.7	63.0% (-5.4%)	35.3% (-8.7%)	80.5% (-3.0%)
<a href="#">Axial-DeepLab</a>	173M	2447G	-	64.4% (-4.0%)	36.7% (-7.3%)	80.6% (-2.9%)
<a href="#">SWideRNet</a>	536M	10365G	1.0	66.4% (-2.0%)	40.1% (-3.9%)	82.2% (-1.3%)
<b>kMaX-DeepLab</b>	<b>232M</b>	<b>1673G</b>	<b>3.1</b>	<b>68.4%</b>	<b>44.0%</b>	<b>83.5%</b>

Cityscapes val set



# *k*-means Mask Transformer



image



1<sup>st</sup> cluster assignment



2<sup>nd</sup> cluster assignment



3<sup>rd</sup> cluster assignment



4<sup>th</sup> cluster assignment



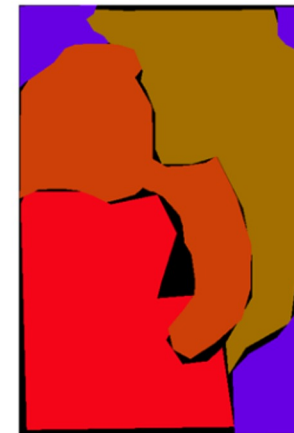
5<sup>th</sup> cluster assignment



6<sup>th</sup> cluster assignment



panoptic prediction

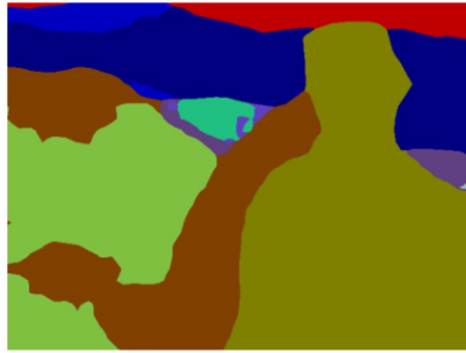


panoptic label

# *k*-means Mask Transformer



image



1<sup>st</sup> cluster assignment



2<sup>nd</sup> cluster assignment



3<sup>rd</sup> cluster assignment



4<sup>th</sup> cluster assignment



5<sup>th</sup> cluster assignment



6<sup>th</sup> cluster assignment



panoptic prediction



panoptic label

# *k*-means Mask Transformer



image



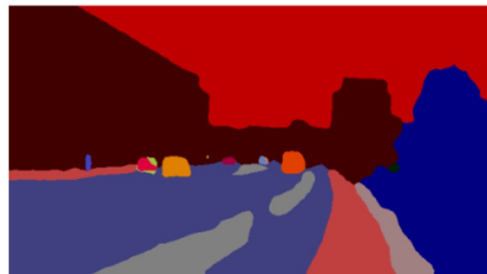
1<sup>st</sup> cluster assignment



2<sup>nd</sup> cluster assignment



3<sup>rd</sup> cluster assignment



4<sup>th</sup> cluster assignment



5<sup>th</sup> cluster assignment



6<sup>th</sup> cluster assignment



panoptic prediction



panoptic label

# Summary

- Discuss the underlying similarity between *cross-attention* and *k-means clustering* algorithm.
- Propose *k-means cross-attention*, which designs cross-attention as a k-means clustering module, leading to better *object-centric representation*.
- A simple change on *activation function* with *SOTA* performance.