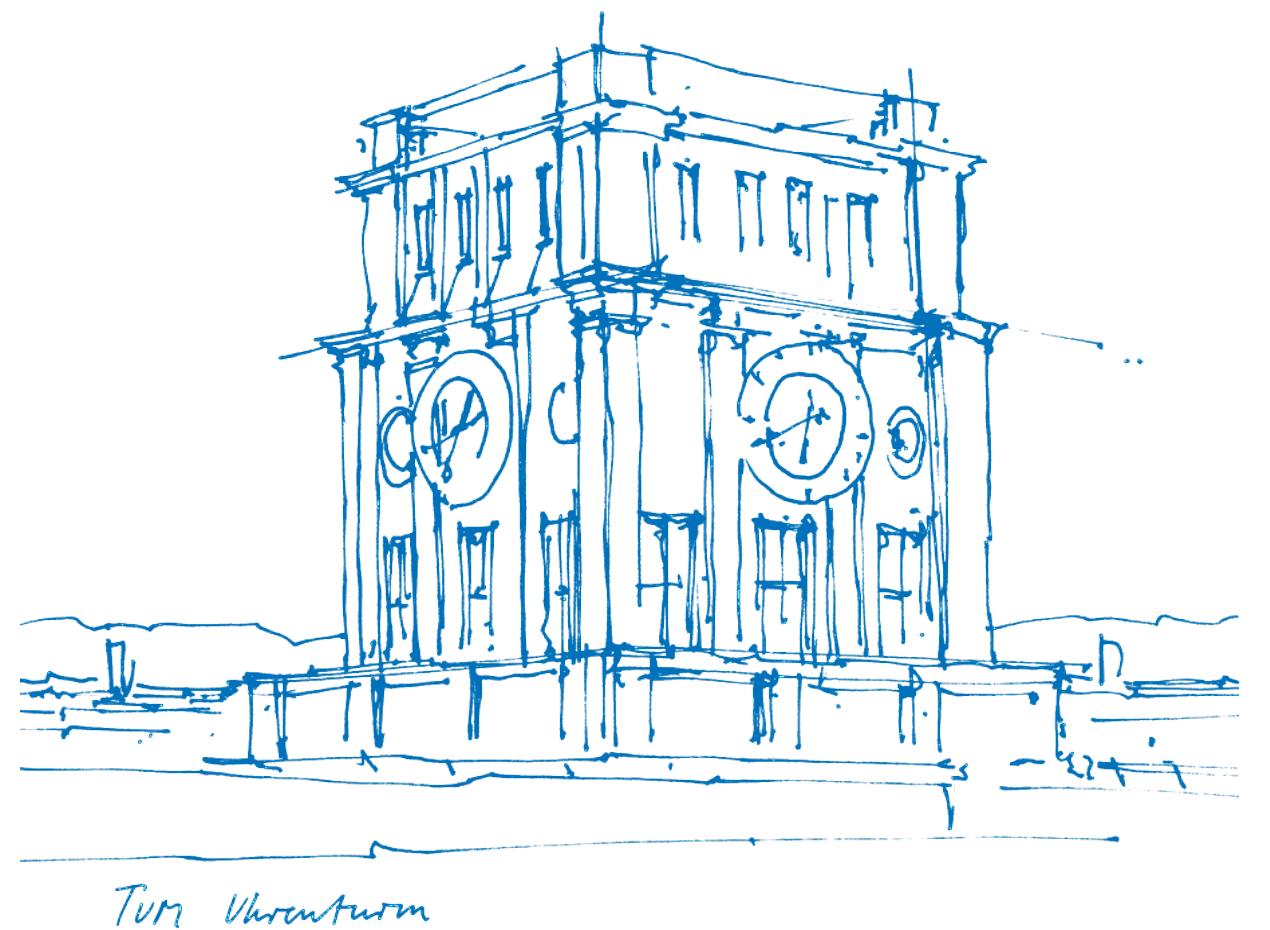


Computer Vision III: Transformers

Nikita Araslanov
19.12.2023

Adapted from:
Prof. Laura Leal-Taixé
<https://dvl.in.tum.de>



Course progress

1. Introduction
2. Object detection 1
3. Object detection 2
4. Single object tracking
5. Multiple object tracking
6. Semantic segmentation
7. Instance & panoptic segmentation
8. Video object segmentation
- 9. Transformers**  We are here
10. Semi-supervised DST (09.01.2024)
11. Unsupervised DST (16.01.2024)

February 6, 16:00:
Q&A interactive lecture (Zoom)
Ask anything about the lecture
material

Reminder:
Please complete course evaluation!
(Deadline: 22.12.2023)

Transformers

Deep Learning Evolution

	CNN Deep Learning	TF Deep Learning
Main idea	Convolution	Attention
Field invented	Computer vision	NLP
Started	NeurIPS 2012	NeurIPS 2017
Paper	AlexNet	Transformers
Conquered vision	Around 2014-2015	Around 2020-2021
Replaced (Augmented)	Traditional ML/CV	CNNs, RNNs

CNNs: A few relevant facts

- CNNs encode local context information 3×3 5×5
- Receptive field increases monotonically with network depth.
 - This enables long-range context aggregation...
 - yet increased depth makes training more challenging.
- We may benefit from more explicit (i.e. in the same layer) non-local feature interactions.

Attention is all you need

Attention Is All You Need

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

Google Brain

lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

Attention is all you need

Attention Is All You Need

Circa 45 thousand
citations in 4.5 years!

Ashish Vaswani*

Google Brain

avaswani@google.com

Noam Shazeer*

Google Brain

noam@google.com

Niki Parmar*

Google Research

nikip@google.com

Jakob Uszkoreit*

Google Research

usz@google.com

Llion Jones*

Google Research

llion@google.com

Aidan N. Gomez* †

University of Toronto

aidan@cs.toronto.edu

Łukasz Kaiser*

Google Brain

lukaszkaiser@google.com

Illia Polosukhin* ‡

illia.polosukhin@gmail.com

Attention is all you need

Attention Is All You Need

Circa 45 thousand
citations in 4.5 years!

Circa 63 thousand
citations in 5.5 years!

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Attention is all you need

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

Circa 45 thousand
citations in 4.5 years!

Circa 63 thousand
citations in 5.5 years!

Circa 101 thousand
citations in 6.5 years!

Attention is all you need

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

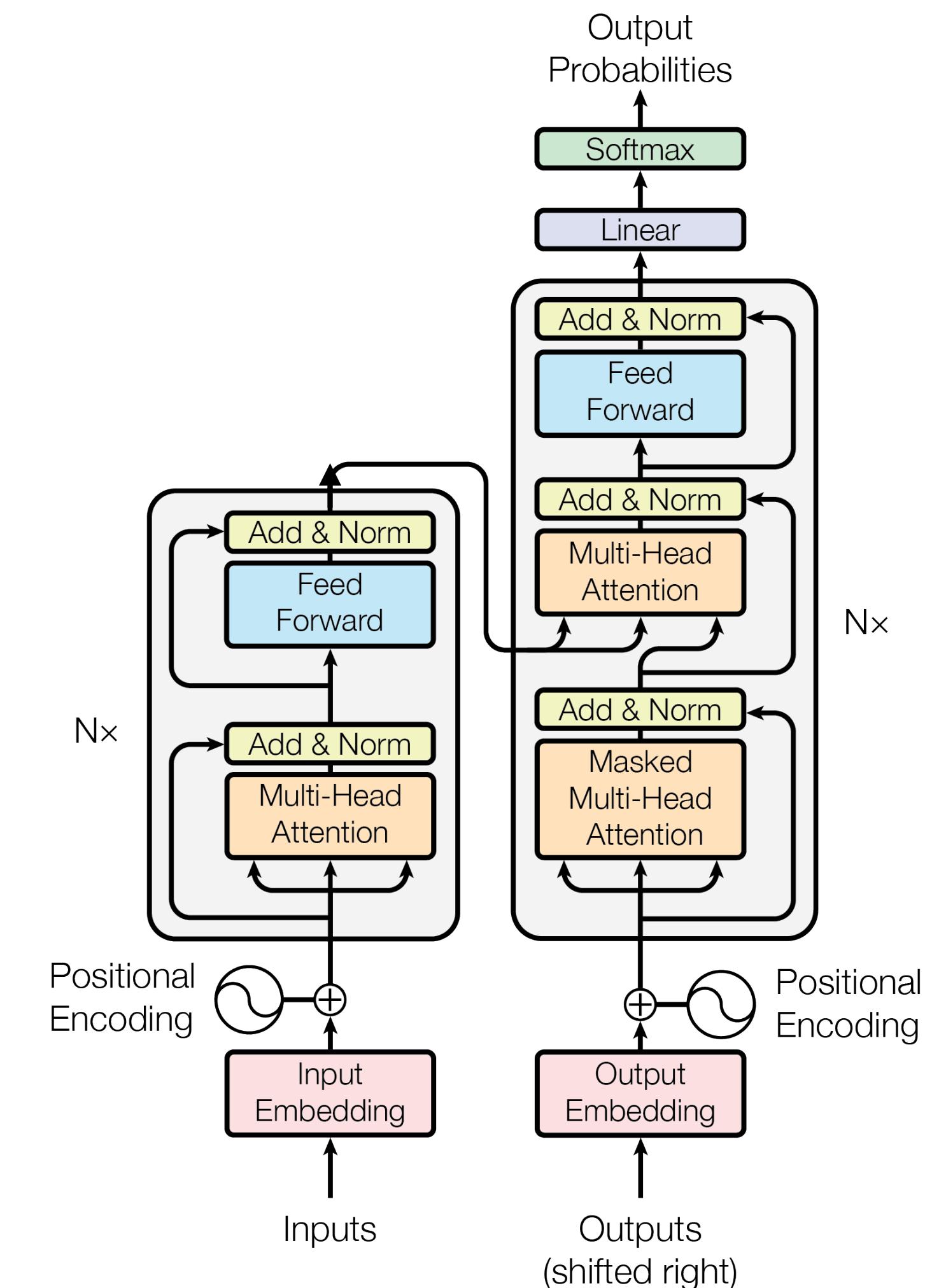
Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com



So, what is so special about it?

Universality

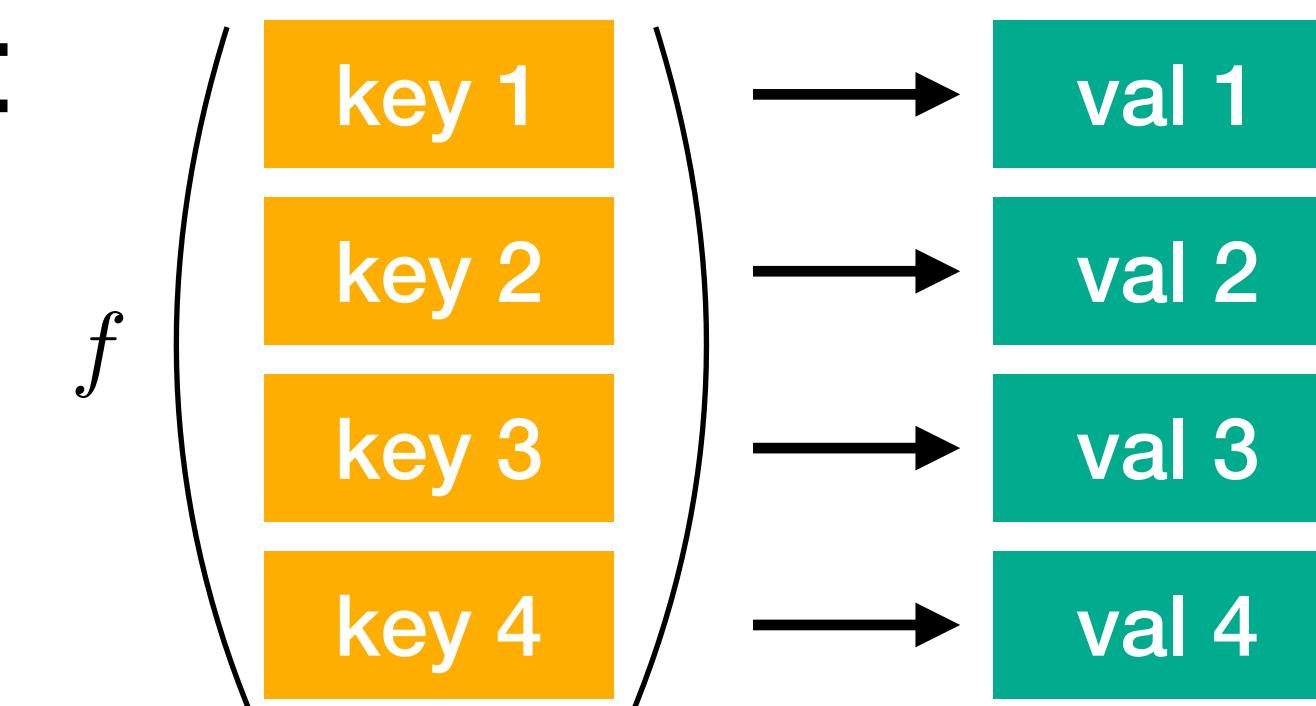
Universality: developing components that work across all possible settings.

- Consider other modalities (language, speech, etc.);
- Modern deep learning are only partially universal, even in the contexts we have seen so far (detection and segmentation);
- Transformers can be applied to language and images with excellent results in both domains!

[Adapted from: Vaswani et al., ICML '21]

Self-attention: A hash table analogy

Consider a hash table:

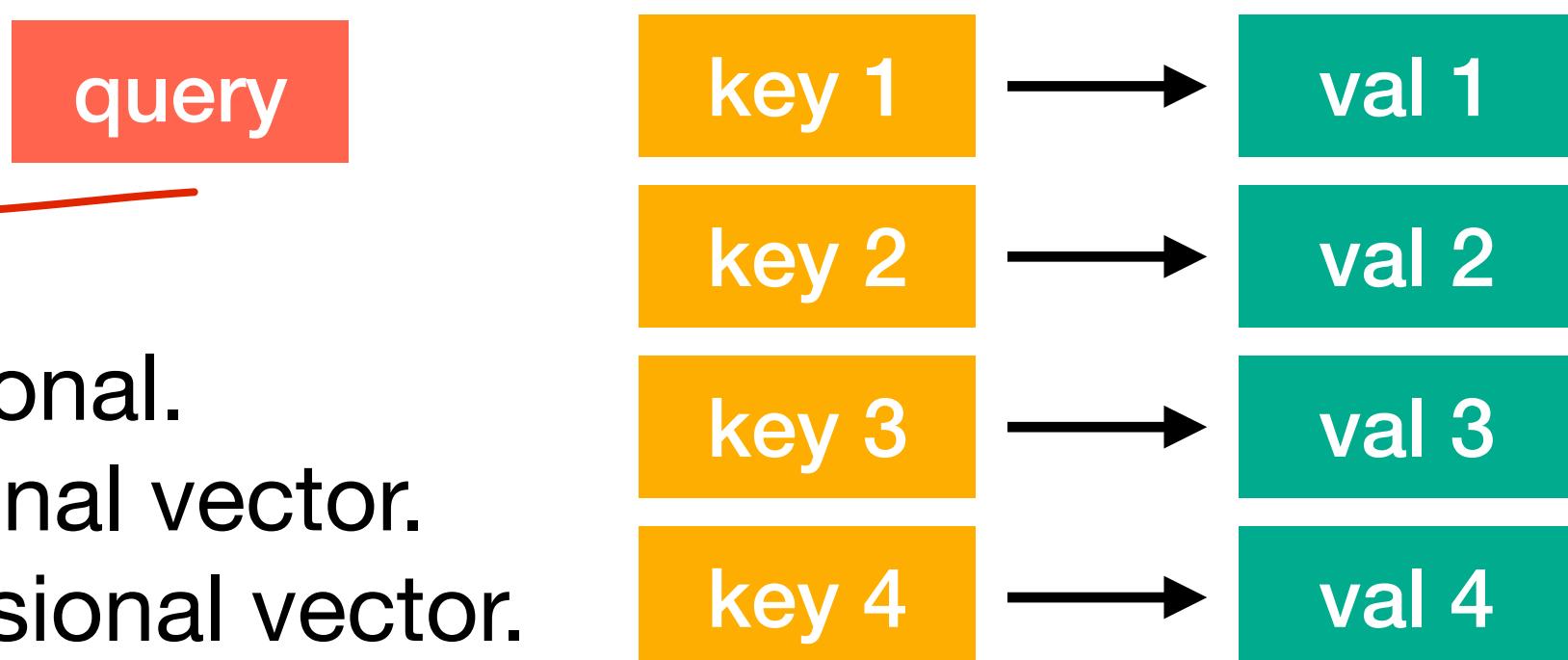


Each key is n-dimensional vector.
Each value is m-dimensional vector.

- A hash table is not a differentiable function (of values w.r.t. keys); it's not even continuous.
 - We can access any value if we know the corresponding key exactly.
 - How can we make it differentiable? – “differentiable soft look-up”

Self-attention: A hash table analogy

Let's say, we have a query vector:



The query is n-dimensional.

Each key is n-dimensional vector.

Each value is m-dimensional vector.

- We can obtain an (approximate) value corresponding to the query:

$$v'(q) = \sum_i w(q, k_i) v_i$$

weighted key

- $w(q, k_i)$ encodes normalised similarity of the query with each key:

$$\sum_i w(q, k_i) = 1$$

Self-attention: A hash table analogy

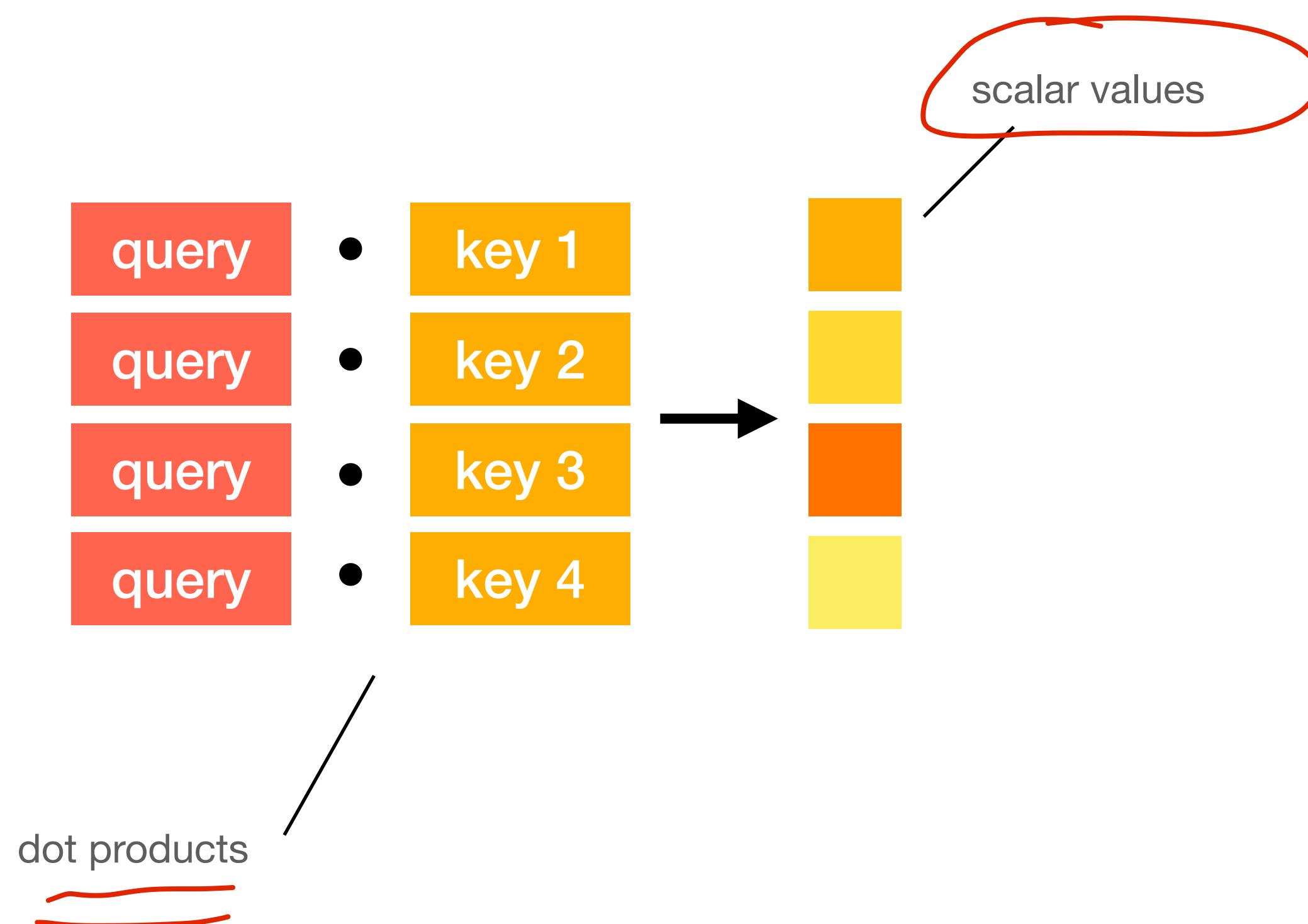
We can use standard operators in our DL arsenal: dot product and softmax

broadcast the query

query	key 1
query	key 2
query	key 3
query	key 4

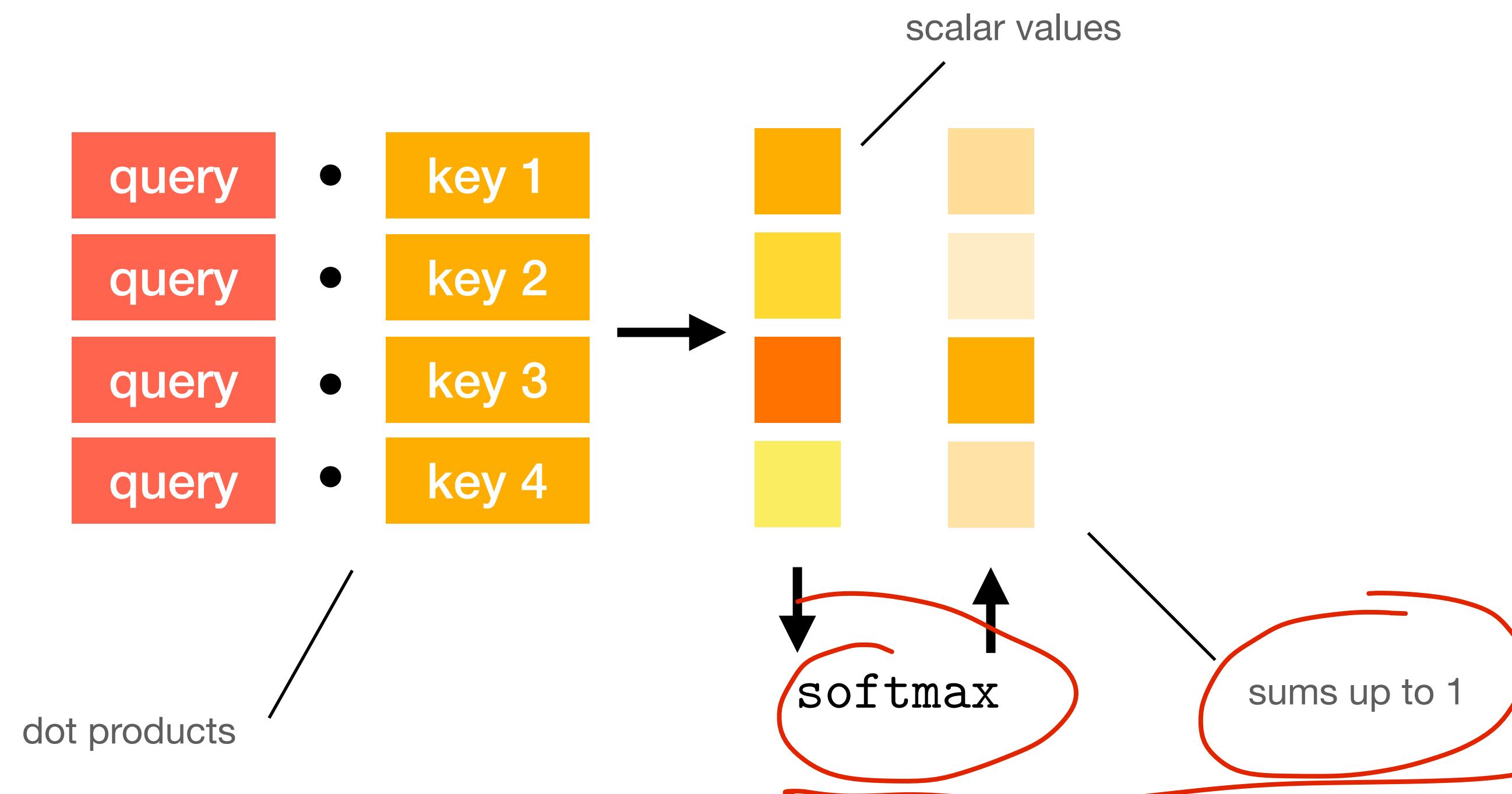
Self-attention: A hash table analogy

We can use standard operators in our DL arsenal: dot product and softmax



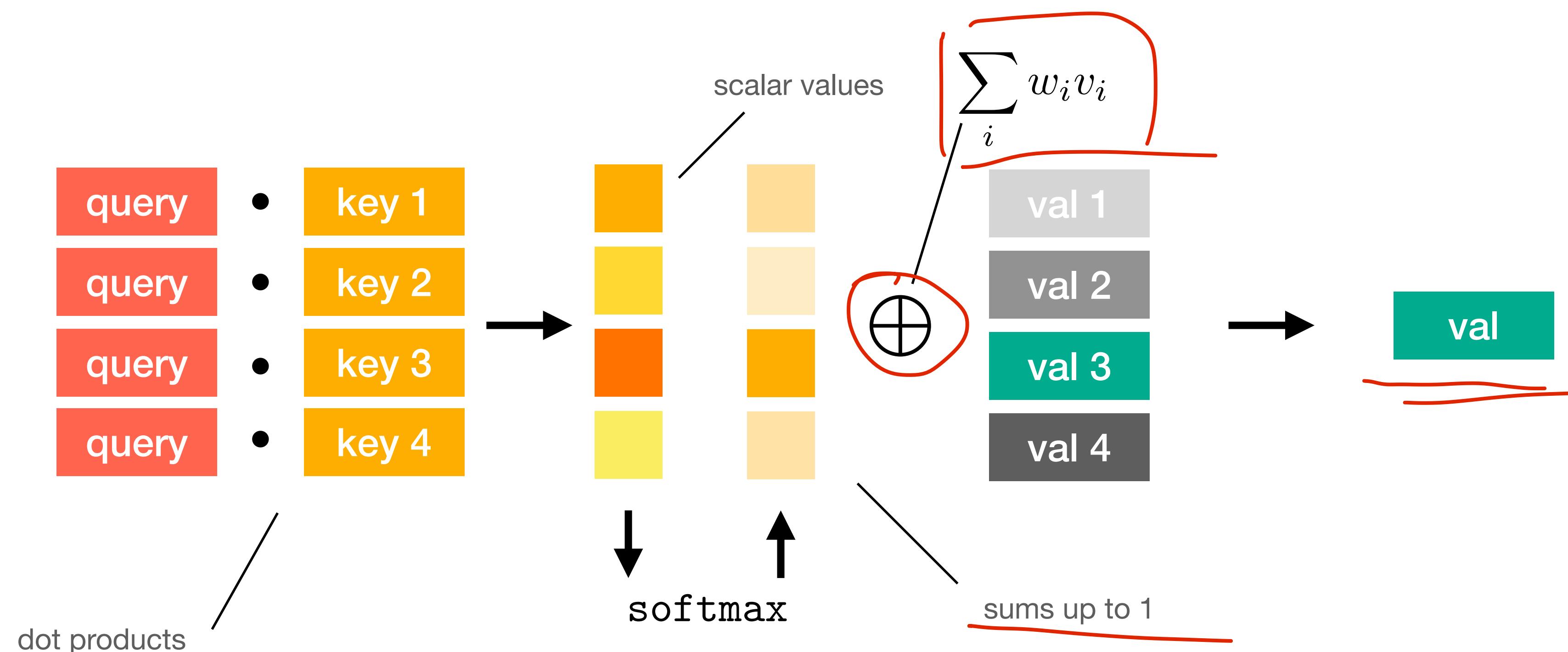
Self-attention: A hash table analogy

We can use standard operators in our DL arsenal: dot product and softmax



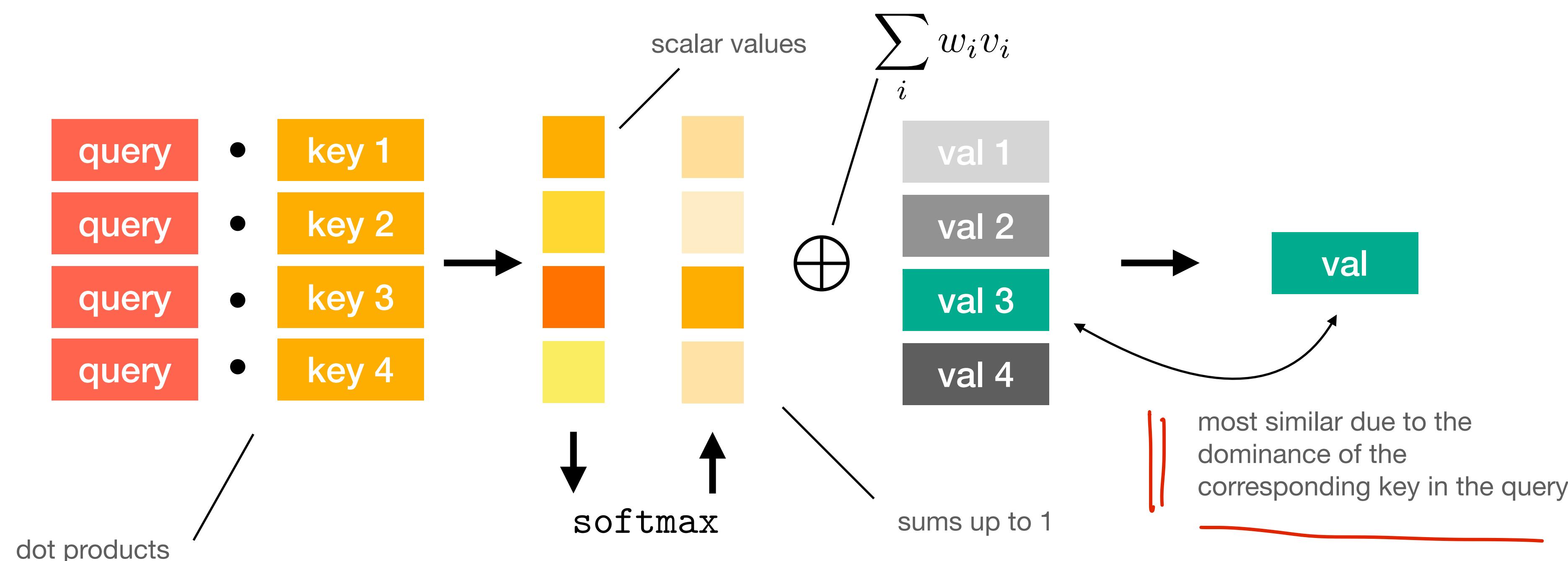
Self-attention: A hash table analogy

We can use standard operators in our DL arsenal: dot product and softmax



Self-attention: A hash table analogy

We can use standard operators in our DL arsenal: dot product and softmax



Self-attention: Summary

Putting it all together. Define:

- a query vector $Q \in \mathbb{R}^{S \times n}$ We have S queries
- a key matrix $K \in \mathbb{R}^{T \times n}$ Our hash table has T (key, value) pairs
- a value matrix $V \in \mathbb{R}^{T \times m}$

Self-attention: $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{n})V$

Self-attention: Scaling factor

$$\text{Self-attention: } \text{Attention}(Q, K, V) = \text{softmax}\left(QK^T / \sqrt{n}\right)V$$

|
scaling factor

$$Q \in \mathbb{R}^{S \times n}$$
$$K \in \mathbb{R}^{T \times n}$$
$$V \in \mathbb{R}^{T \times m}$$

- factor \sqrt{n} eases optimisation especially for large n (QUIZ: Ideas why?)
 - Large n increases the magnitude of the sum (we have n terms);
 - Values with large magnitude inside softmax lead to vanishing gradient.

Self-attention: Linear projection

Self-attention: $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{n})V$

$Q \in \mathbb{R}^{S \times n}$
 $K \in \mathbb{R}^{T \times n}$
 $V \in \mathbb{R}^{T \times m}$

|
scalar factors

- What is the output dimension? (QUIZ)

Self-attention: Linear projection

$$\text{Self-attention: } \text{Attention}(Q, K, V) = \underbrace{\text{softmax}(QK^T / \sqrt{n})V}_{\substack{| \\ \text{scalar factors}}} \quad Q \in \mathbb{R}^{S \times n} \quad K \in \mathbb{R}^{T \times n} \quad V \in \mathbb{R}^{T \times m}$$

- What is the output dimension? (QUIZ)

$\mathbb{R}^{S \times m}$ “for every query we fetch
the corresponding value”

- We can make the input and the output dimensionality the same.
 - This is important if we want to stack this operation and build deep networks.

Self-attention: Linear projection

Input: $X \in \mathbb{R}^{T \times d}$ “**T tokens** of size d ”

Define 4 linear mappings: $W^Q, W^K \in \mathbb{R}^{d \times n}$ $W^V \in \mathbb{R}^{m \times d}$ $W^O \in \mathbb{R}^{m \times d}$



These are model parameters

Compute Q, K, V :

$$K := XW^K \quad [T \times n] \quad Q := XW^Q \quad [T \times n] \quad V := XW^V \quad [T \times m]$$



Output: $Y := \text{Attention}(Q, K, V)W^O$
 $[T \times d]$



We started with $[T \times d]$ and produced $[T \times d]$ output

Self-attention: Complexity

Self-attention: $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{n})V$

$Q \in \mathbb{R}^{S \times n}$
 $K \in \mathbb{R}^{T \times n}$
 $V \in \mathbb{R}^{T \times m}$

scalar factors

- Memory complexity ($T \gg n, m$): $O(T^2)$

$$\text{softmax}(QK^T / \sqrt{n})V$$

$[T \times T]$

- Runtime complexity: $O(T^2n)$

Multi-head attention

$$\text{Self-attention: } \text{Attention}(Q, K, V) = \text{softmax}\left(QK^T / \sqrt{n}\right)V$$

$Q \in \mathbb{R}^{S \times n}$
 $K \in \mathbb{R}^{T \times n}$
 $V \in \mathbb{R}^{T \times m}$

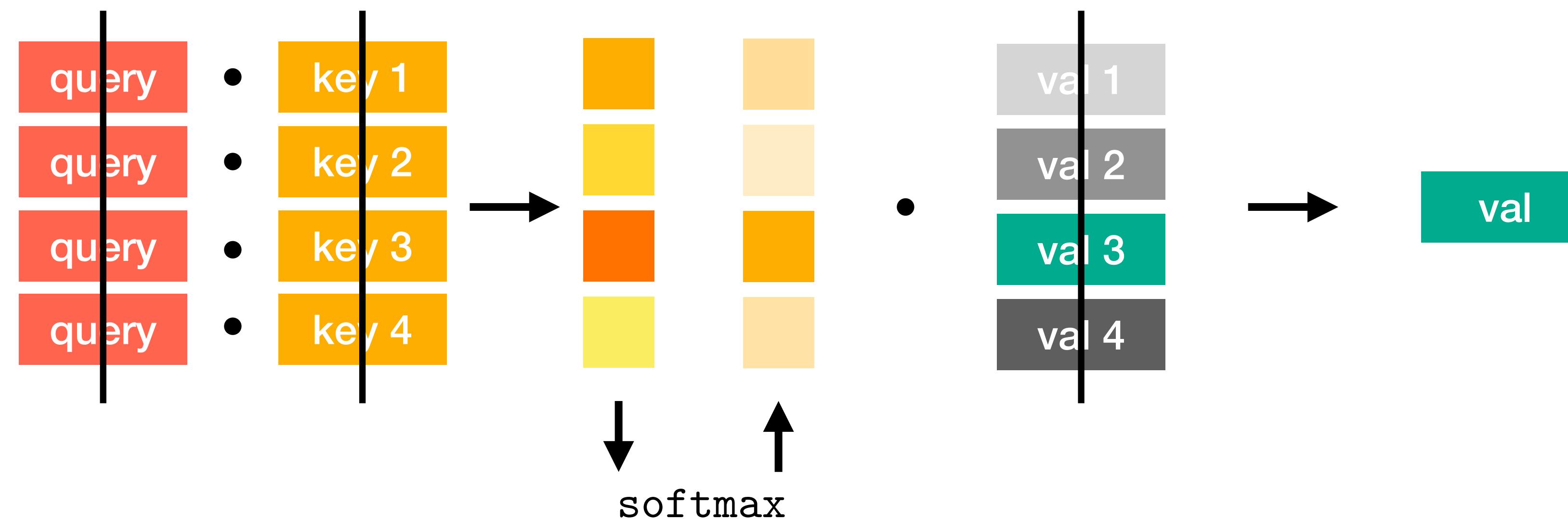
scalar factors

- One softmax – one look-up operation;
- We can extend the same operation to multiple look-ups without increasing the complexity;
- This is called **multi-head** attention.

Multi-head attention

Self-attention: $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{n})V$

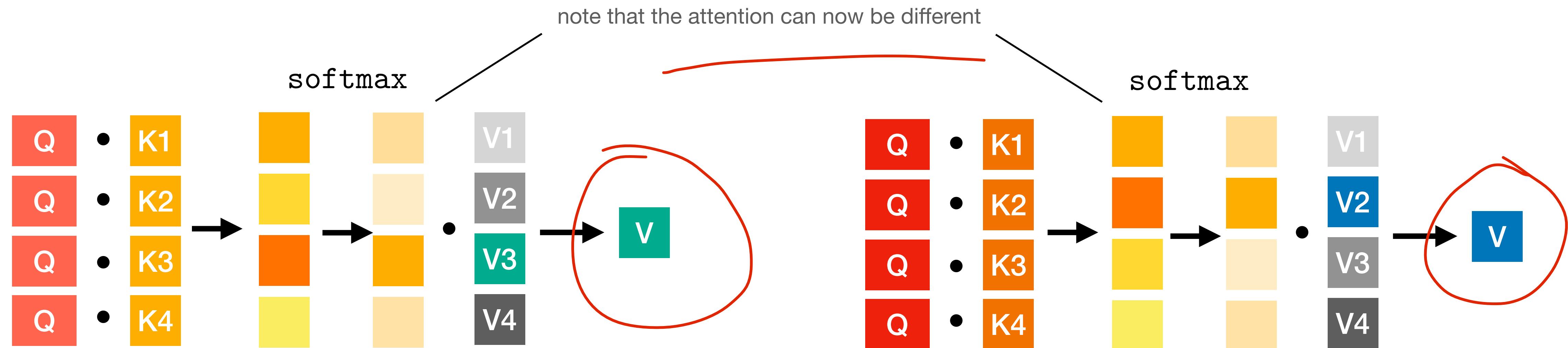
Main idea: chunk Q,K,V vectors



For example, split it in two (e.g. split 128d vectors into 2 64d vectors)

Multi-head attention

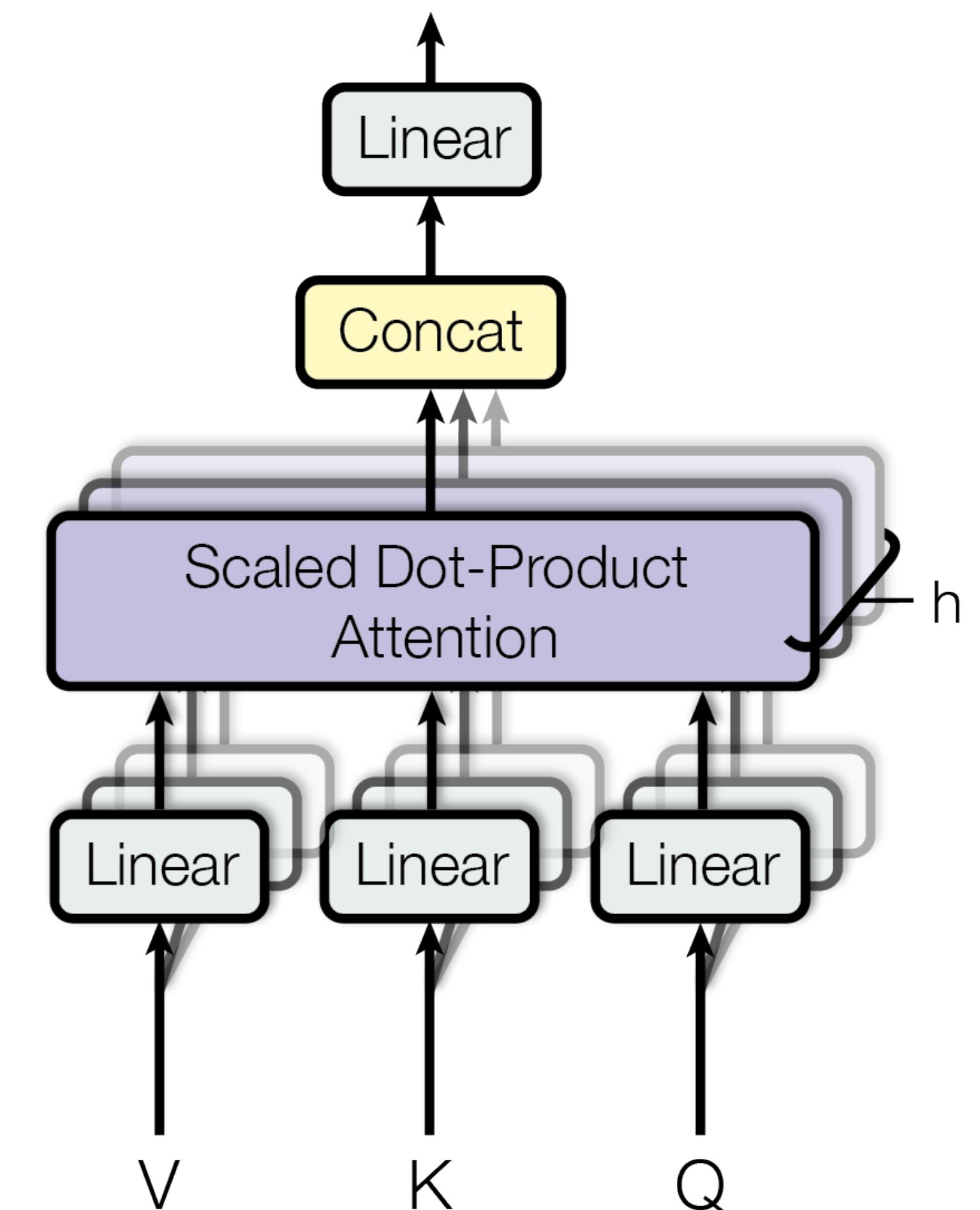
- Now, we have two (Q,K,V) triples, where feature dimension is reduced by a factor of 2.
- Repeat the same process:



- Concatenate v v to obtain the original resolution.

Multi-head attention

Intuition: Given a single query, we can fetch multiple values corresponding to different keys.



[Vaswani et al., 2017]

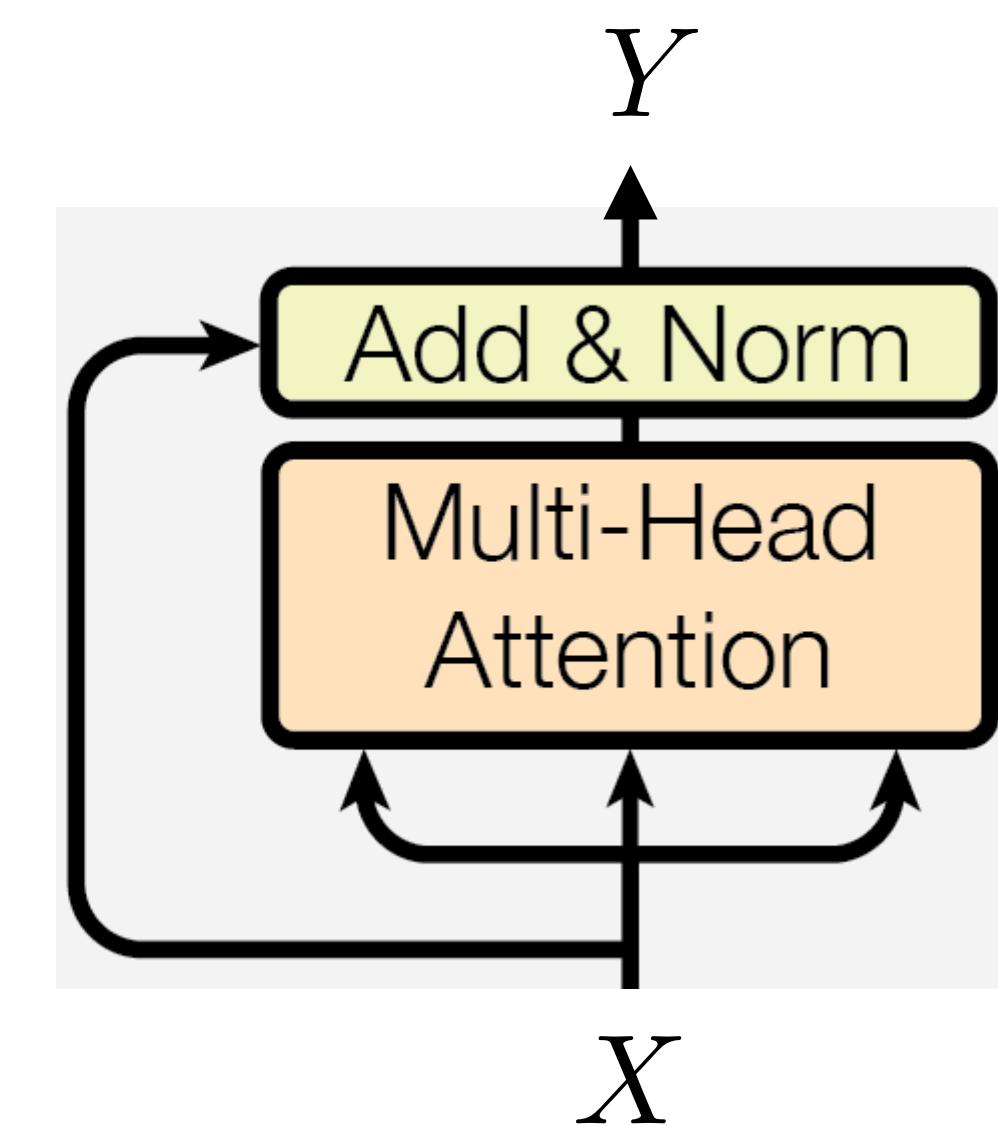
Multi-head attention

- We can design self-attention with arbitrary number of heads.
 - condition: feature dimensionality of Q,K and V should be divisible by it.
- With Z heads, a single query can fetch up to Z different values.
 - we can model more complex interactions between the tokens.
- No increase in runtime complexity – actually faster in wall time.
 - QUIZ: Why?

Self-attention: Normalisation

Further improvements:

- Layer normalisation (Ba et al., 2016)
 - normalise each token feature w.r.t. its own mean and standard deviation.
- Residual connection:
 - add the original token feature before the normalisation.



$$Y := \text{LayerNorm}(X + \text{MHA}(X))$$

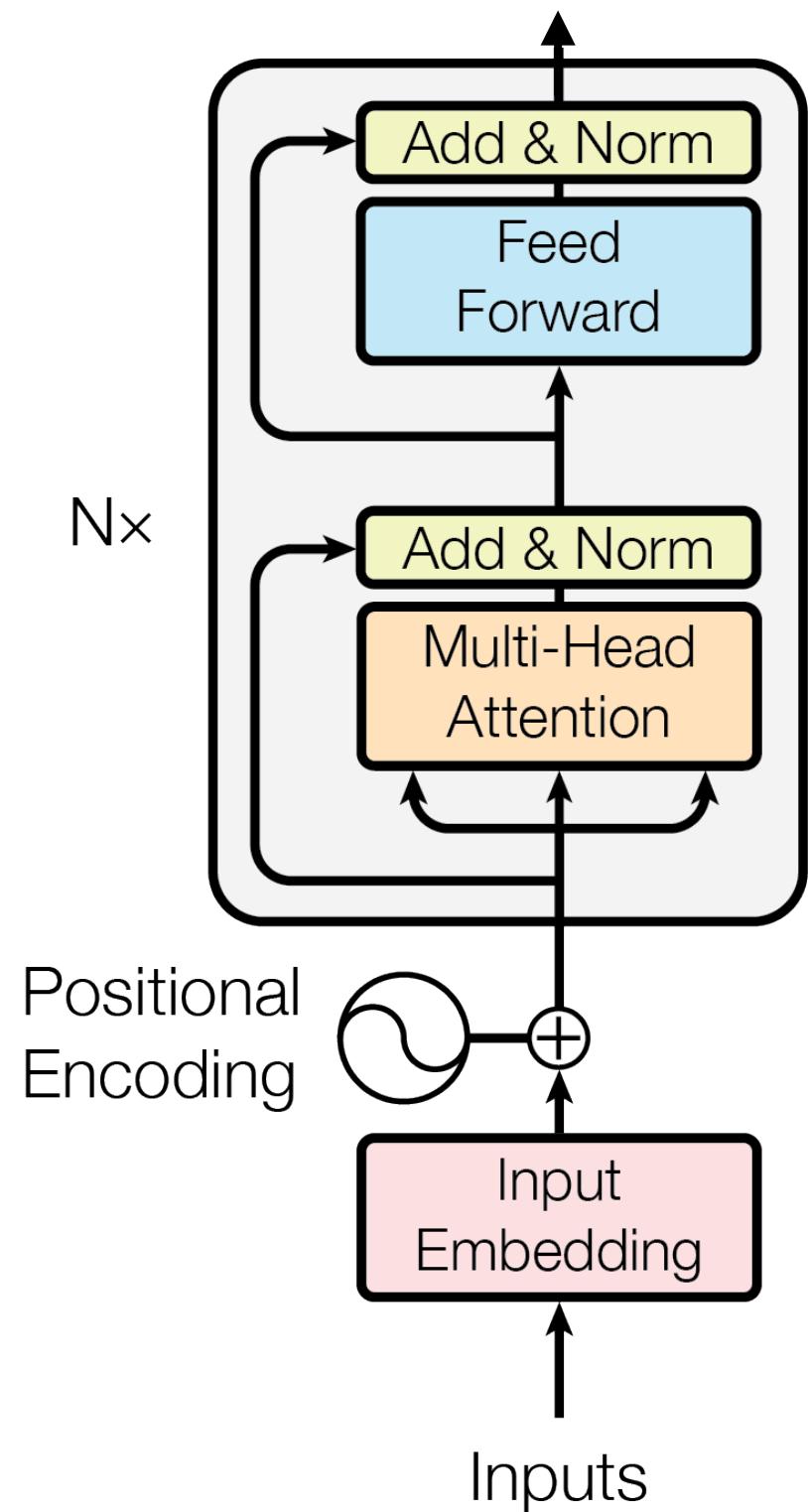
[Vaswani et al., 2017]

Self-attention: Recap

- Self-attention is versatile:
 - arbitrary number of tokens;
 - design choices: the number of heads, feature dimensionality.
- Self-attention computation scales quadratically w.r.t. tokens,
 - since we need to compute the pairwise similarity matrix.
- Self-attention is permutation-invariant (w.r.t. the tokens):
 - the query index does not matter; we will always fetch the same value for it.
 - Is it always useful?

Transformers

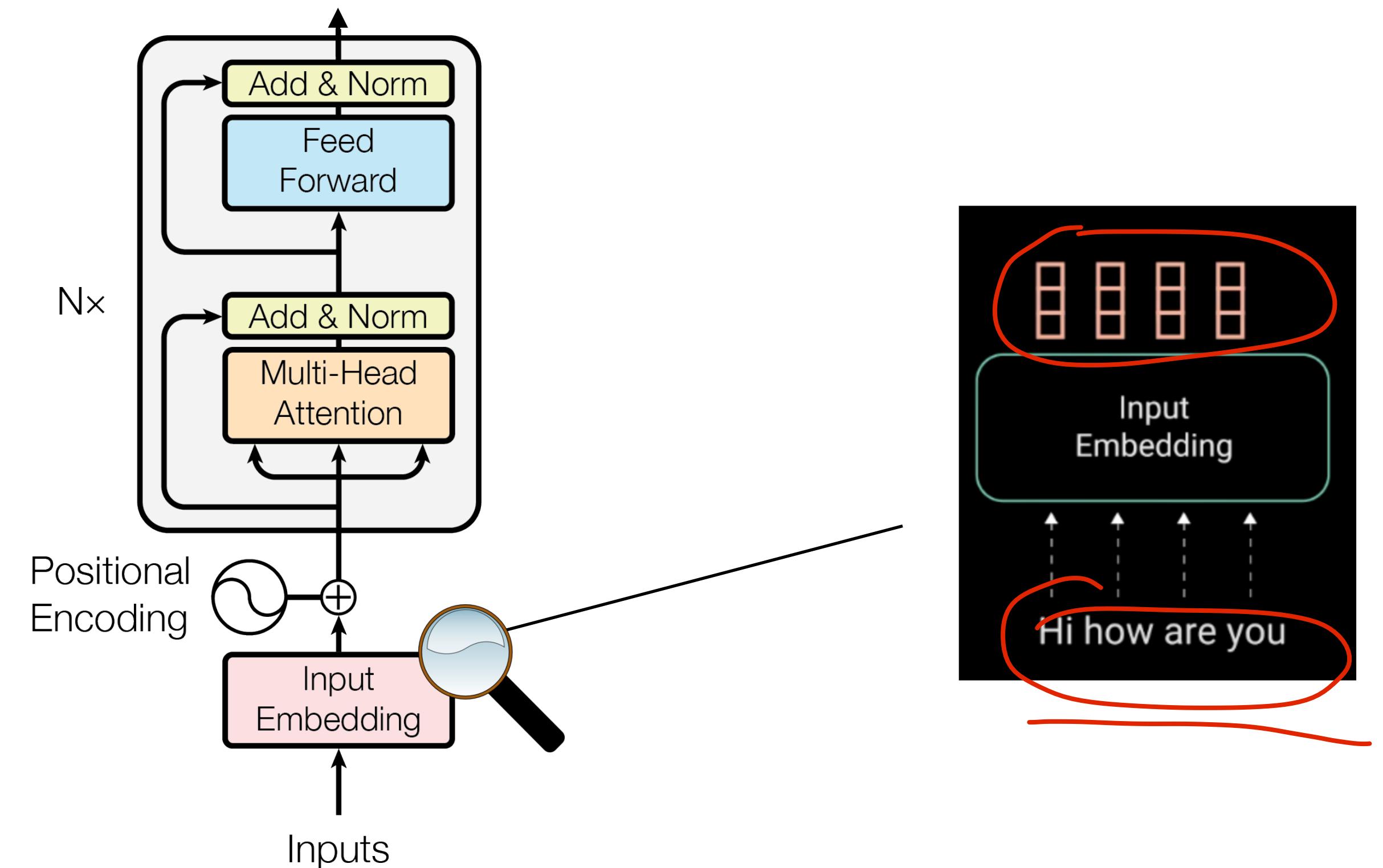
- Transformer-based encoder:



[Vaswani et al., 2017]

Transformers

- Transformer-based encoder:

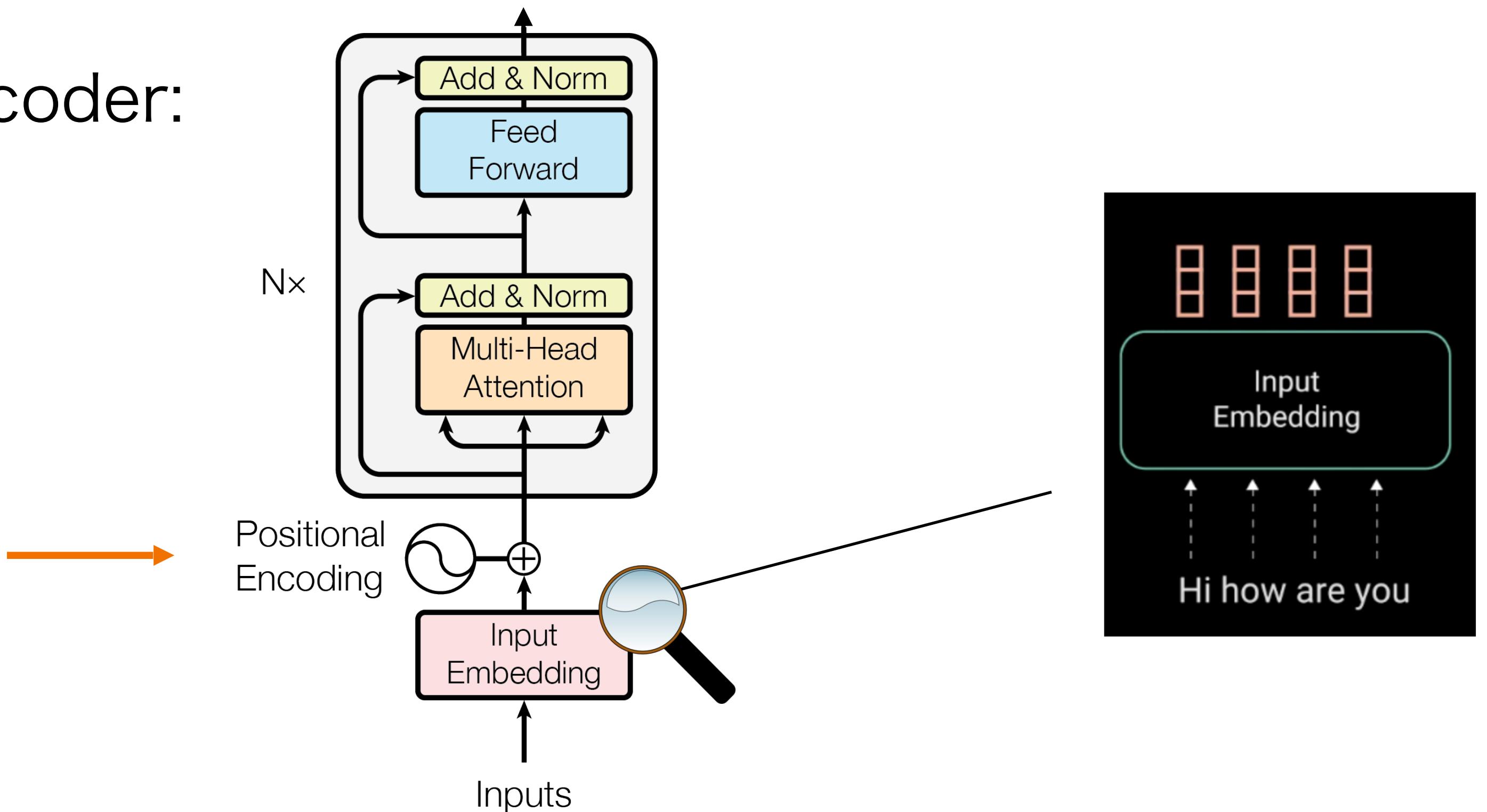


[Vaswani et al., 2017]

Transformers

- Transformer-based encoder:

Breaks
permutation
invariance.

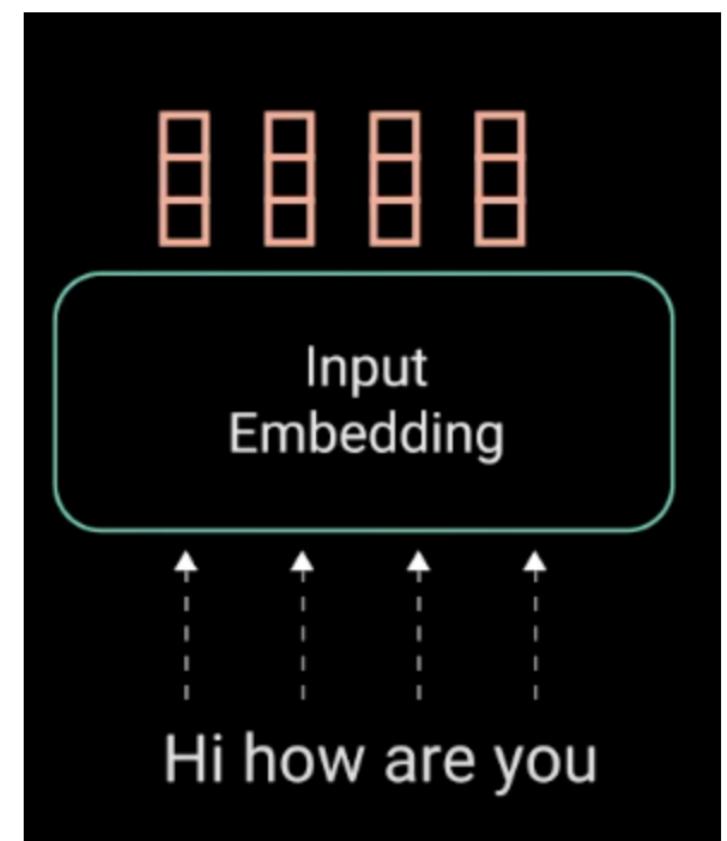
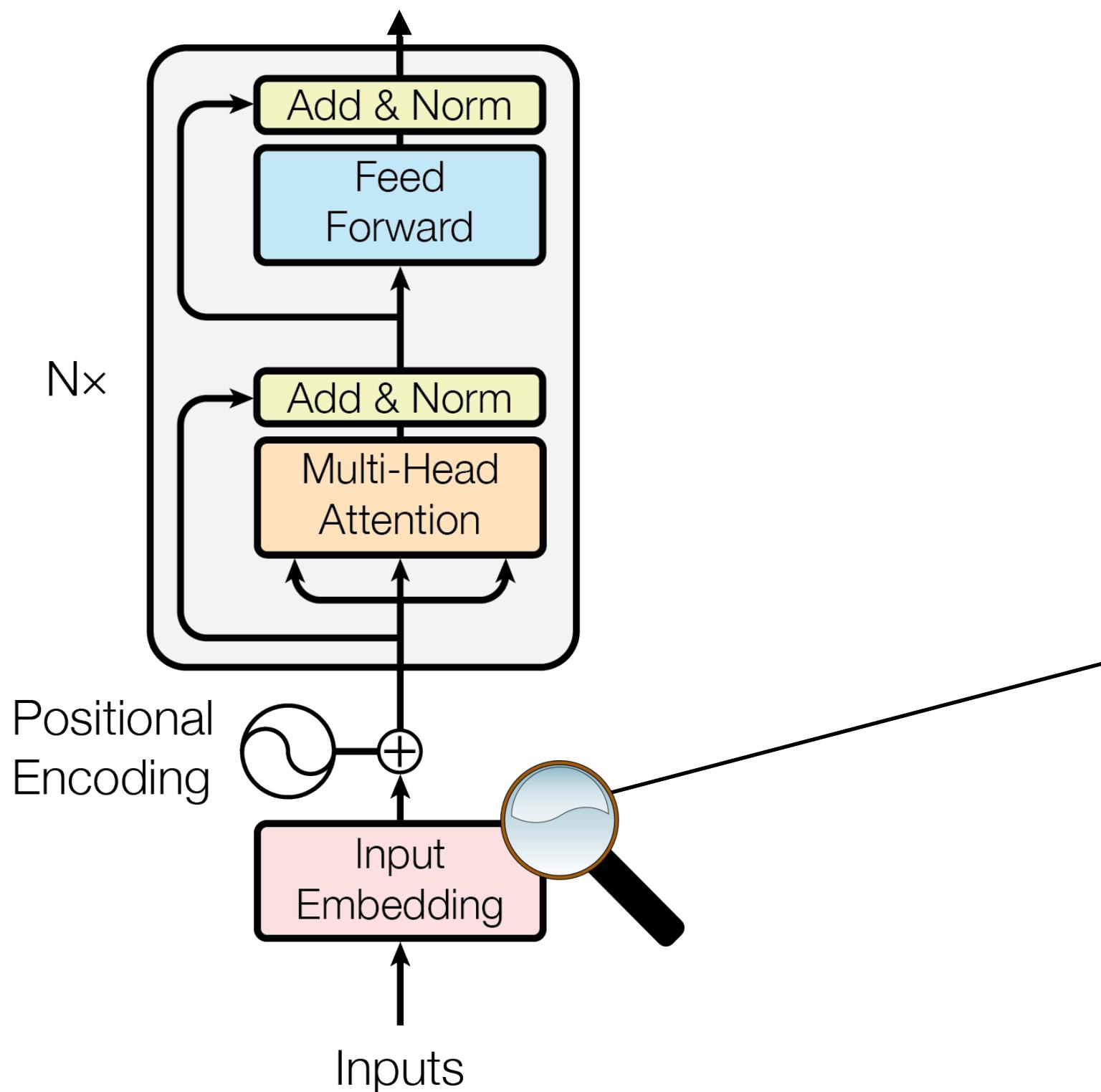


[Vaswani et al., 2017]

Transformers

- Transformer-based encoder:

Breaks
permutation
invariance.



- Positional encoding is a matrix of size distinct values
corresponding to each row

[Vaswani et al., 2017]

Positional encoding

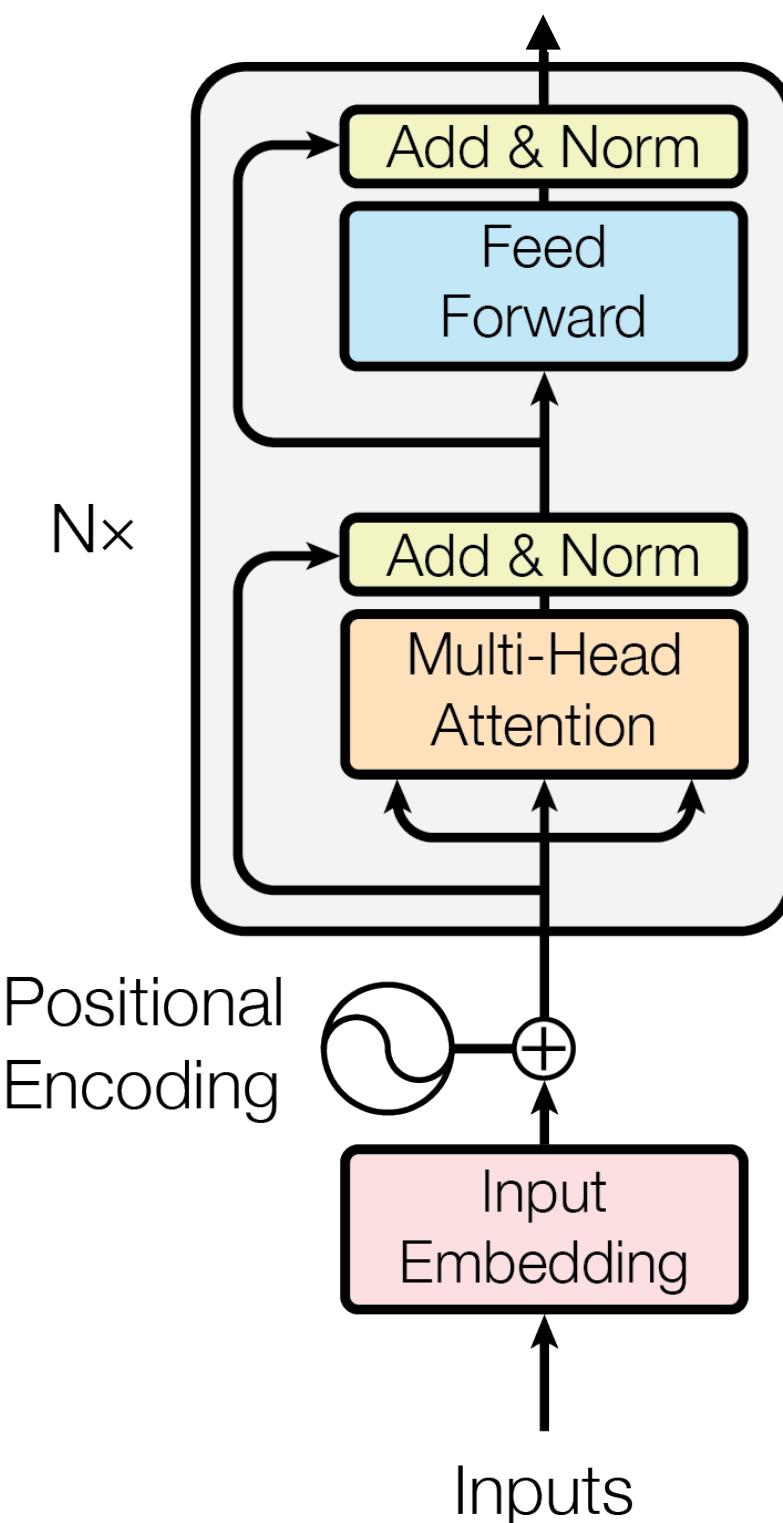
- Positional encoding stores a unique value (vector) corresponding to a token index.

- It can be a learnable model parameter or fixed, for example:

$$PE_{(pos,2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$

$$PE_{(pos,2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$

- It introduces the notion of spatial affinity (e.g. distance between two words in sentence; two patches in an image, etc.)
- It breaks permutation invariance.



[Vaswani et al., 2017]

Transformers in computer vision

Vision Transformer (ViT)

- An all-round transformer architecture
 - Competitive with CNNs (classification accuracy)*
 - Main idea: Train on sequences of image patches; only minor changes to the original (NLP) Transformer model.
 - Can be complemented with a CNN model.

(Dosovitskiy et al., 2020)

Vision Transformer

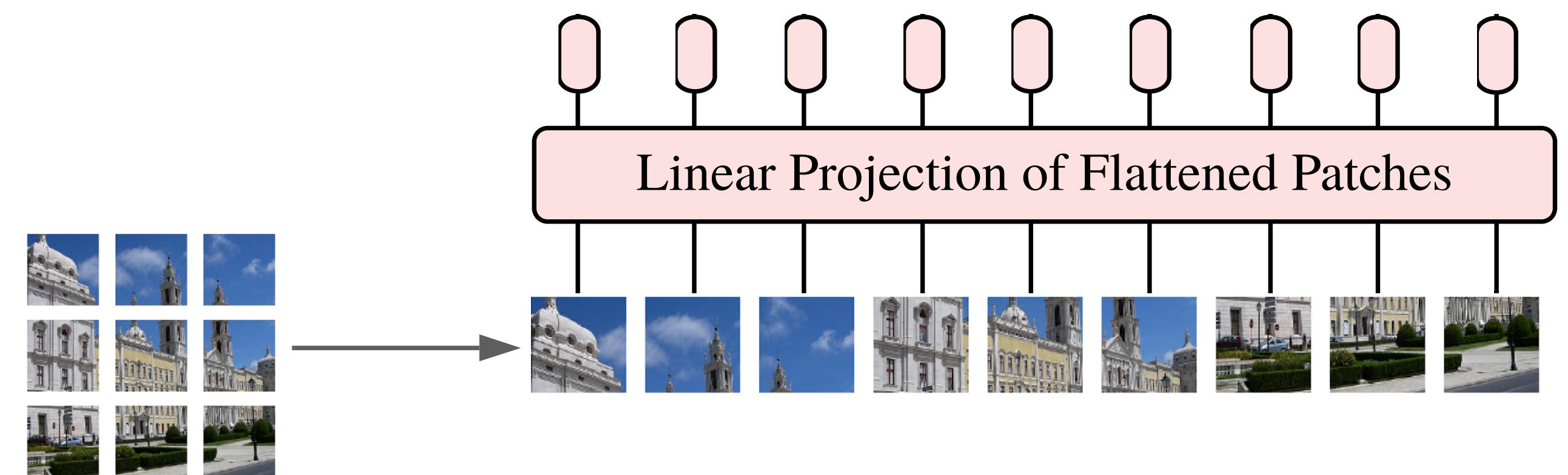
1. Split an image into fixed-sized patches.



(Dosovitskiy et al., 2020)

Vision Transformer

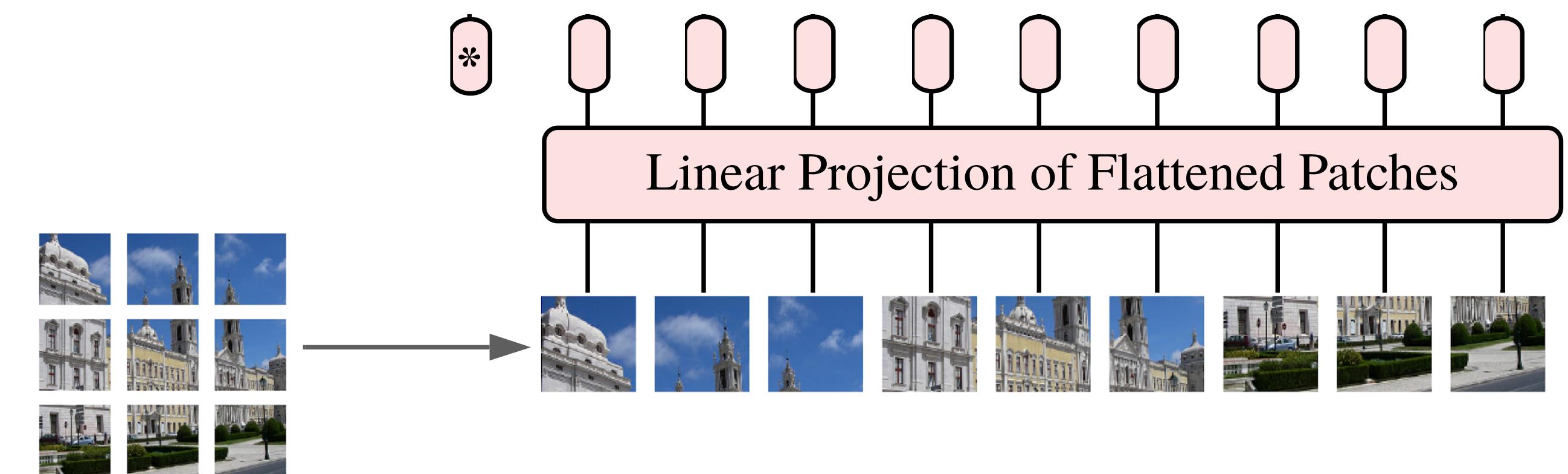
1. Split an image into fixed-sized patches.
2. Linearly embed each patch (a fully connected layer).



(Dosovitskiy et al., 2020)

Vision Transformer

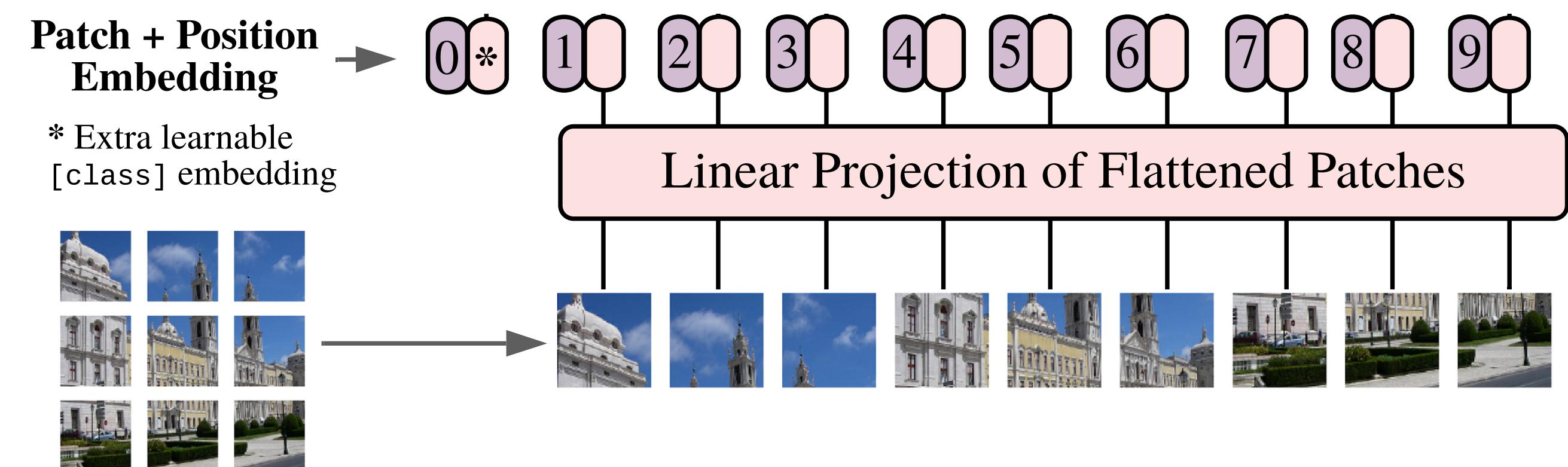
1. Split an image into fixed-sized patches.
2. Linearly embed each patch (a fully connected layer).
3. Attach an extra “[class]” embedding (learnable).



(Dosovitskiy et al., 2020)

Vision Transformer

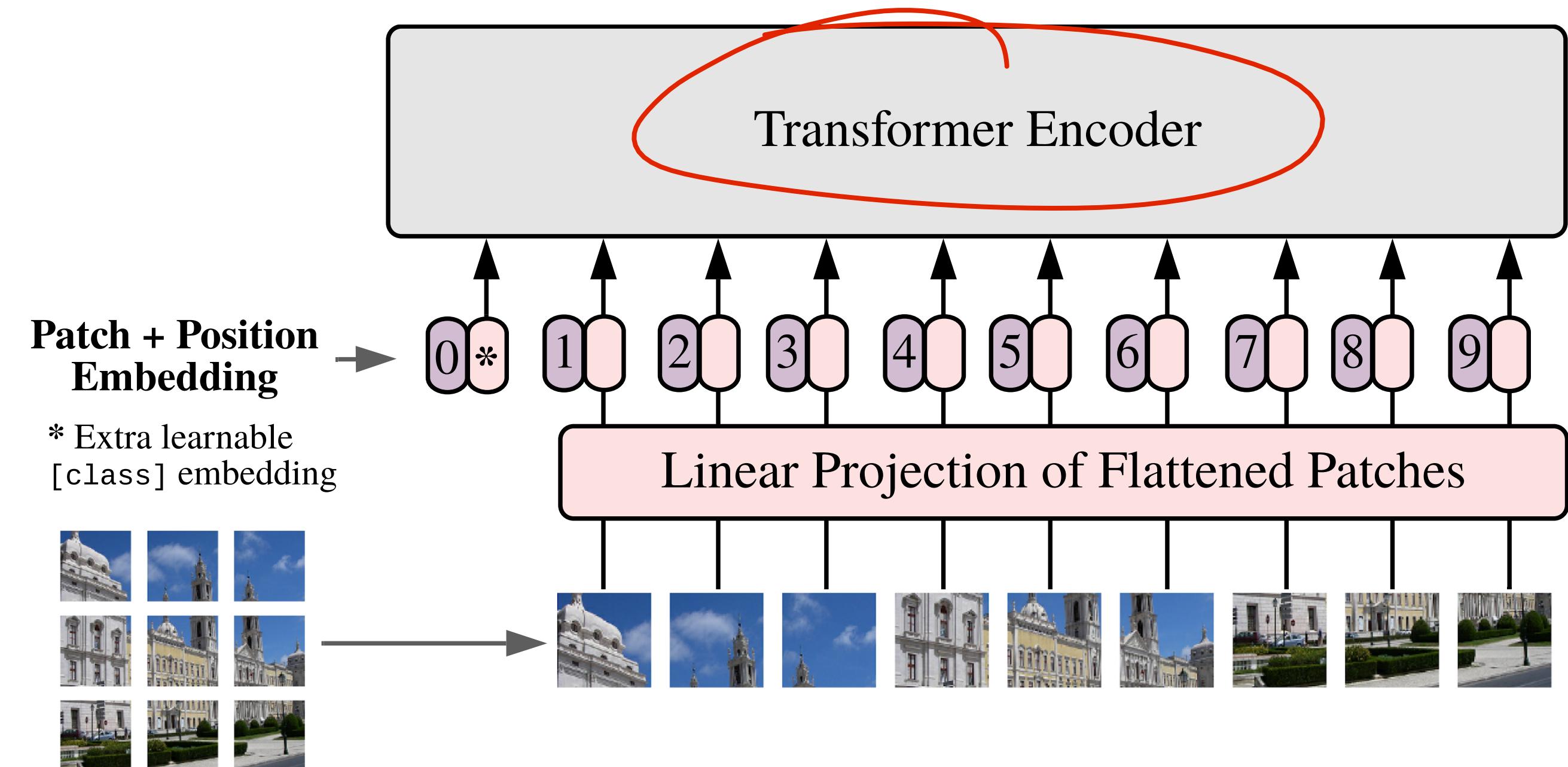
1. Split an image into fixed-sized patches.
2. Linearly embed each patch (a fully connected layer).
3. Attach an extra “[class]” embedding (learnable).
4. Add positional embeddings.



(Dosovitskiy et al., 2020)

Vision Transformer

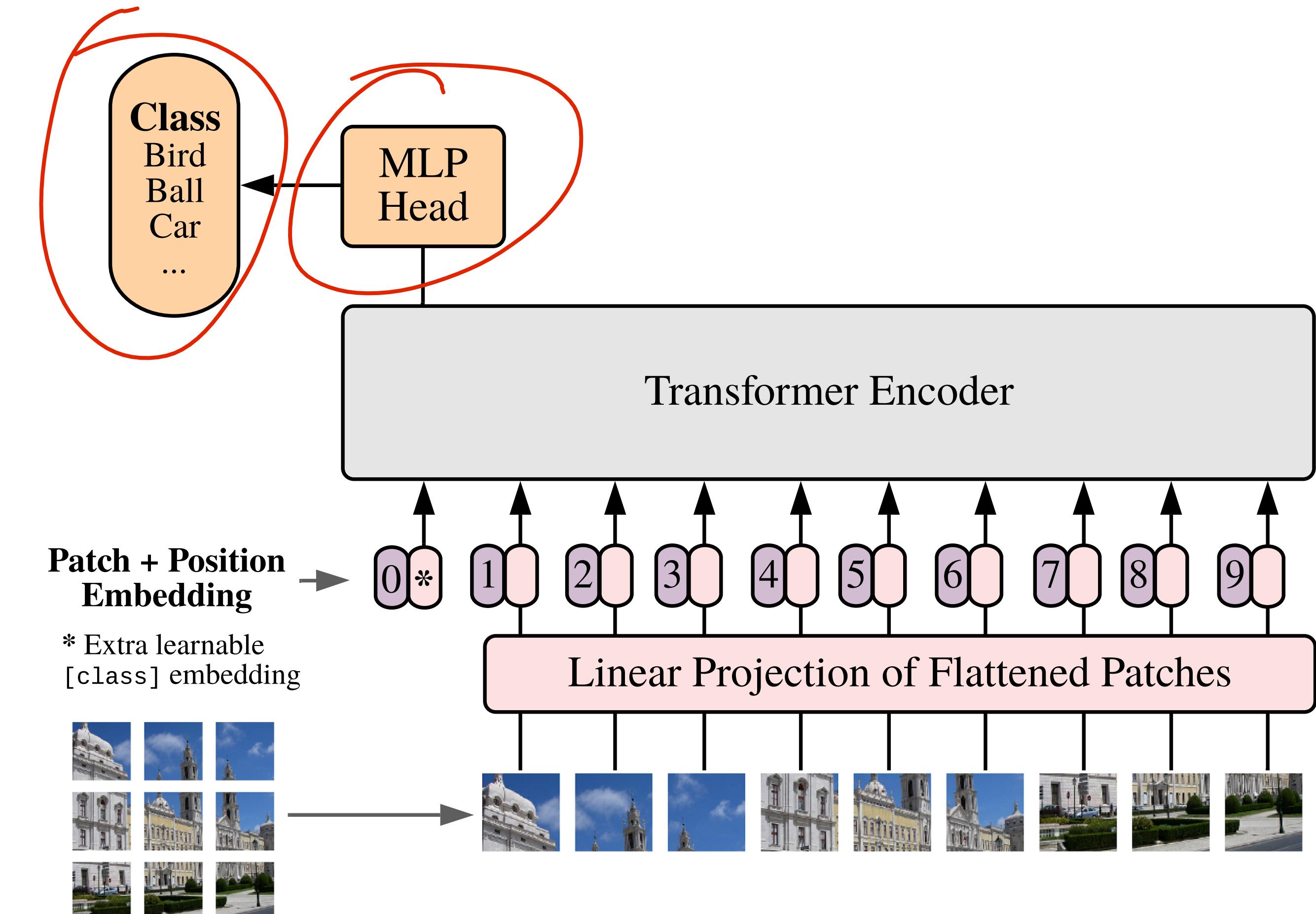
1. Split an image into fixed-sized patches.
2. Linearly embed each patch (a fully connected layer).
3. Attach an extra “[class]” embedding (learnable).
4. Add positional embeddings.
5. Feed the sequence to the standard Transformer encoder.



(Dosovitskiy et al., 2020)

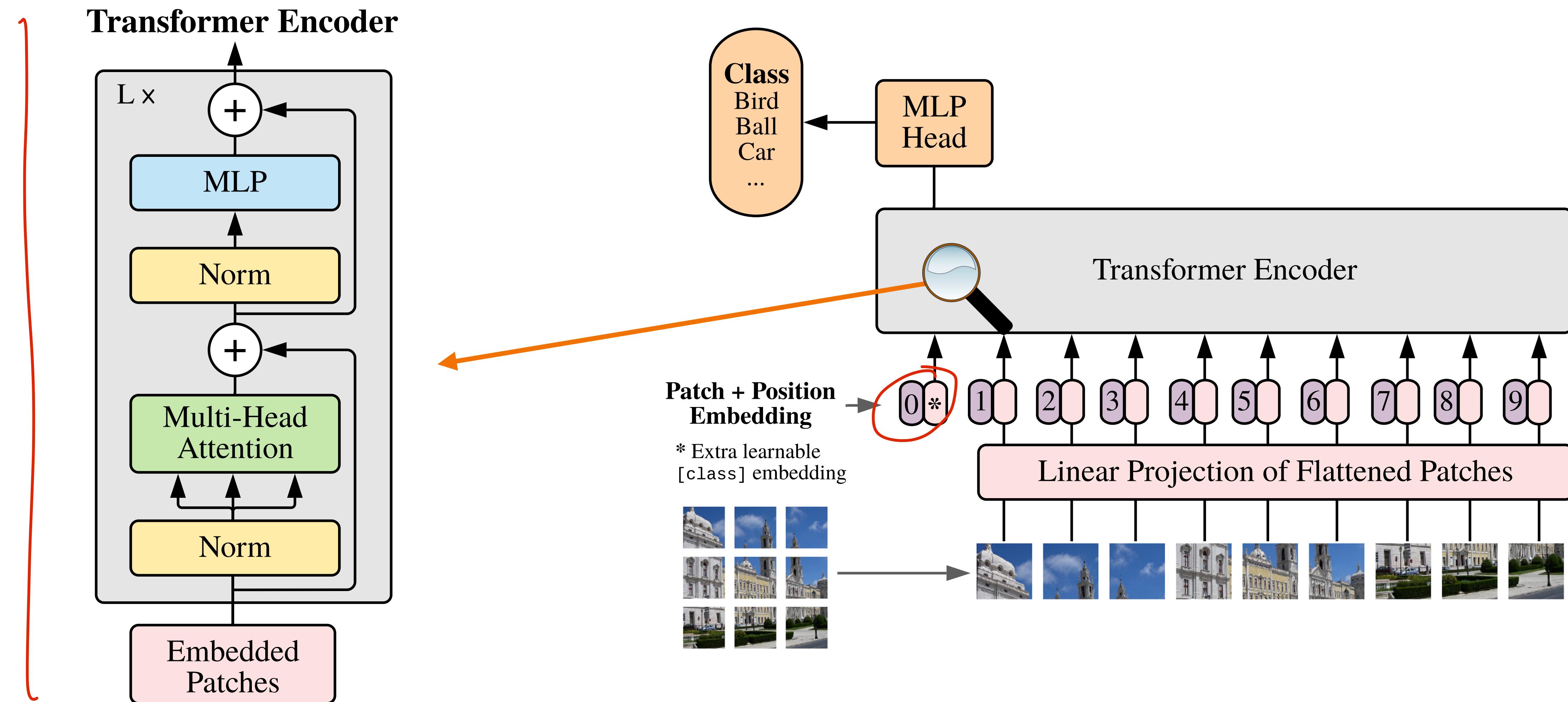
Vision Transformer

1. Split an image into fixed-sized patches.
2. Linearly embed each patch (a fully connected layer).
3. Attach an extra “[class]” embedding (learnable).
4. Add positional embeddings.
5. Feed the sequence to the standard Transformer encoder.
6. Predict the class with an MLP using the output [class] token.



(Dosovitskiy et al., 2020)

ViT Encoder



(Dosovitskiy et al., 2020)

Experiments with ViT

- ViT performs well only when pre-trained on large JFT dataset (300M images). (QUIZ: Why?)

Experiments with ViT

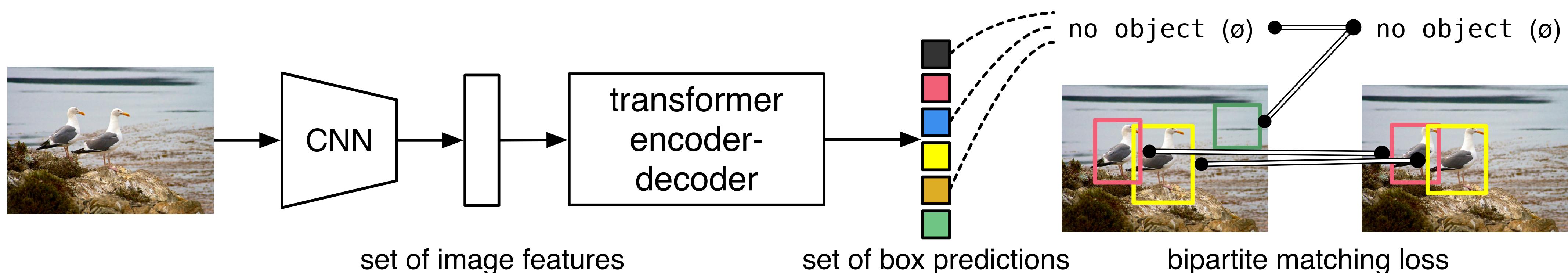
- ViT performs well only when pre-trained on large JFT dataset (300M images). (QUIZ: Why?)
- ViT does not have the inductive biases of CNNs:
 - locality (self-attention is global);
 - 2D neighbourhood structure (positional embedding is inferred from data);
 - translation invariance.

Experiments with ViT

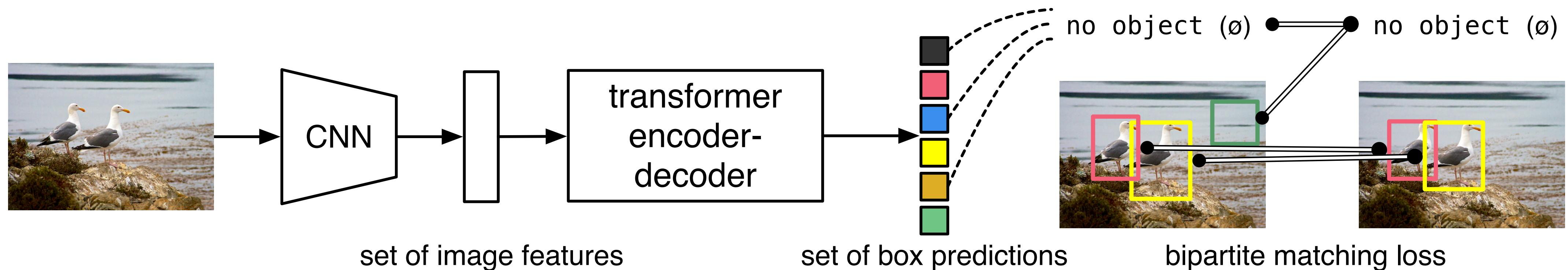
- ViT performs well only when pre-trained on large JFT dataset (300M images). (QUIZ: Why?)
- ViT does not have the inductive biases of CNNs:
 - locality (self-attention is global);
 - 2D neighbourhood structure (positional embedding is inferred from data);
 - translation invariance.
- Nevertheless: now we use the same computational framework for two distinct modalities: language and vision!

Transformer for detection?

- Would it make sense to adapt the Transformer to object detection?
- Recall that object detection is a set prediction problem
 - i.e. we do not care about the ordering of the bounding boxes.
- Transformers are well-suited for processing sets.
- Directly formulate the task as set prediction!

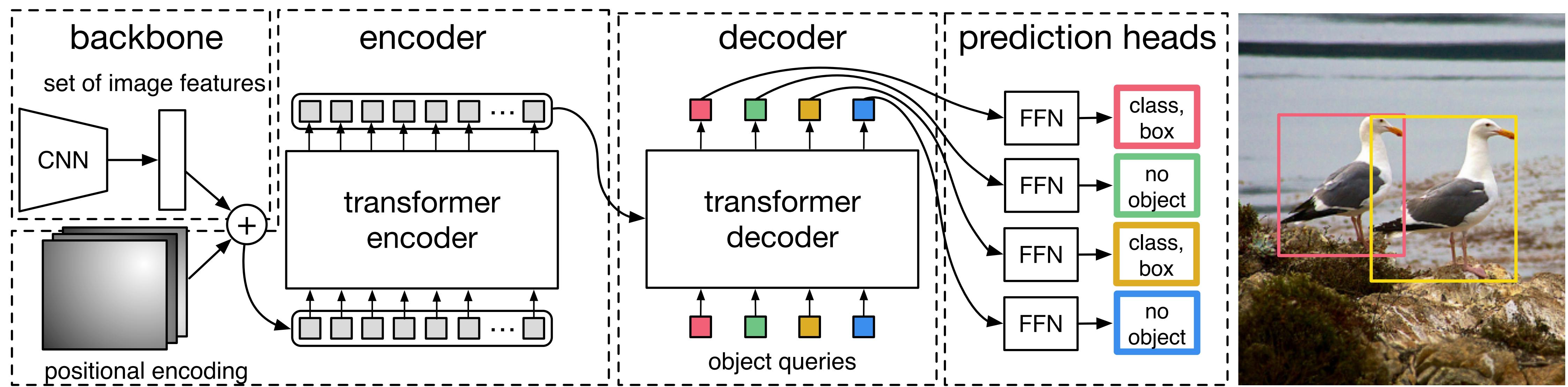


Detection: DETR

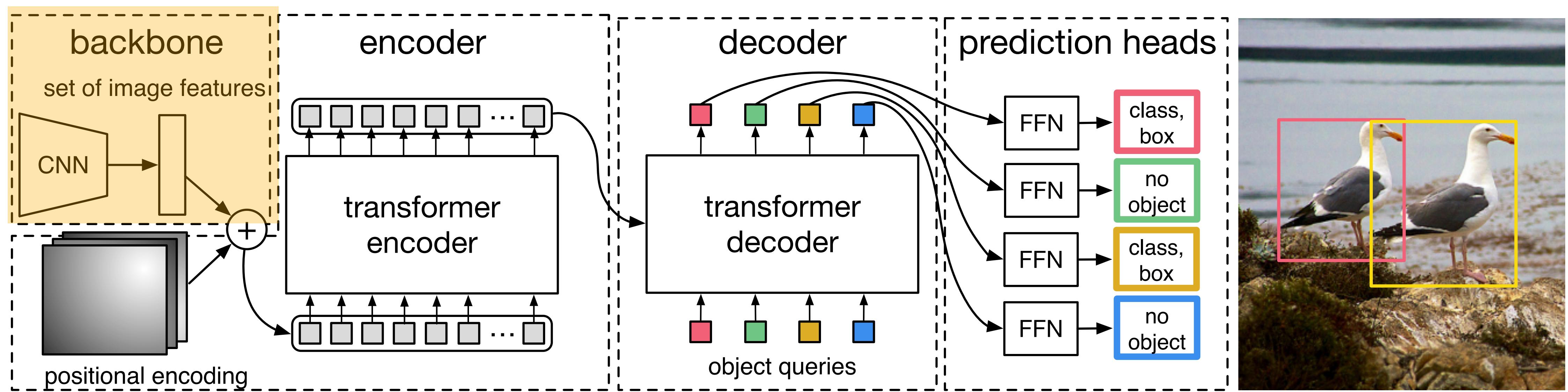


- The CNN predicts local feature embeddings.
- The Transformer predicts the bounding boxes in parallel.
- During training, we uniquely assigns predictions to ground truth boxes with Hungarian matching.
- No need for non-maximum suppression.

DETR: A closer look

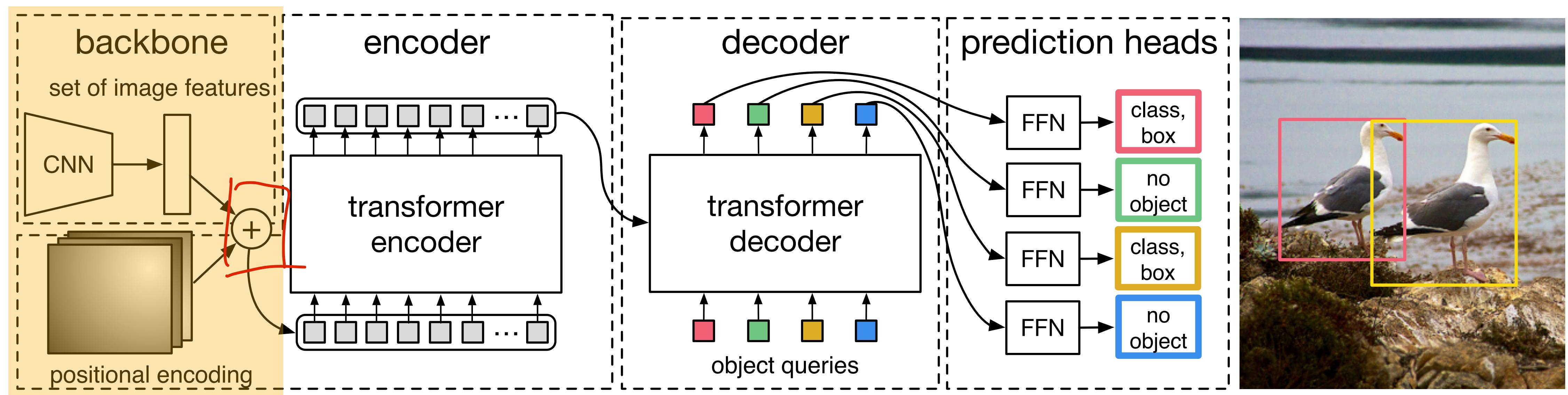


DETR: A closer look



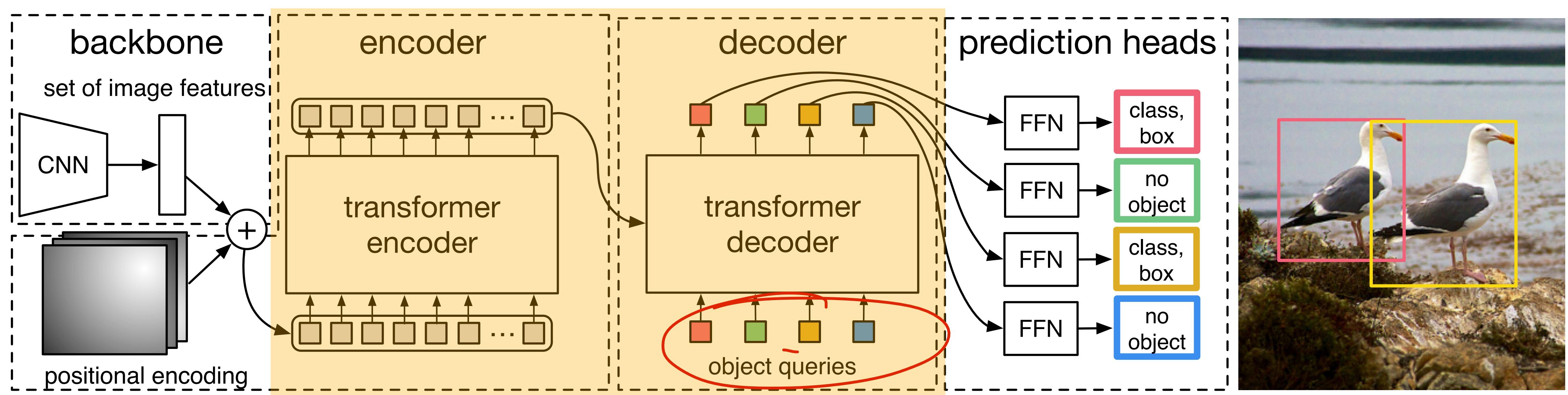
- DETR uses a conventional CNN backbone to learn a 2D representation of an input image.

DETR: A closer look



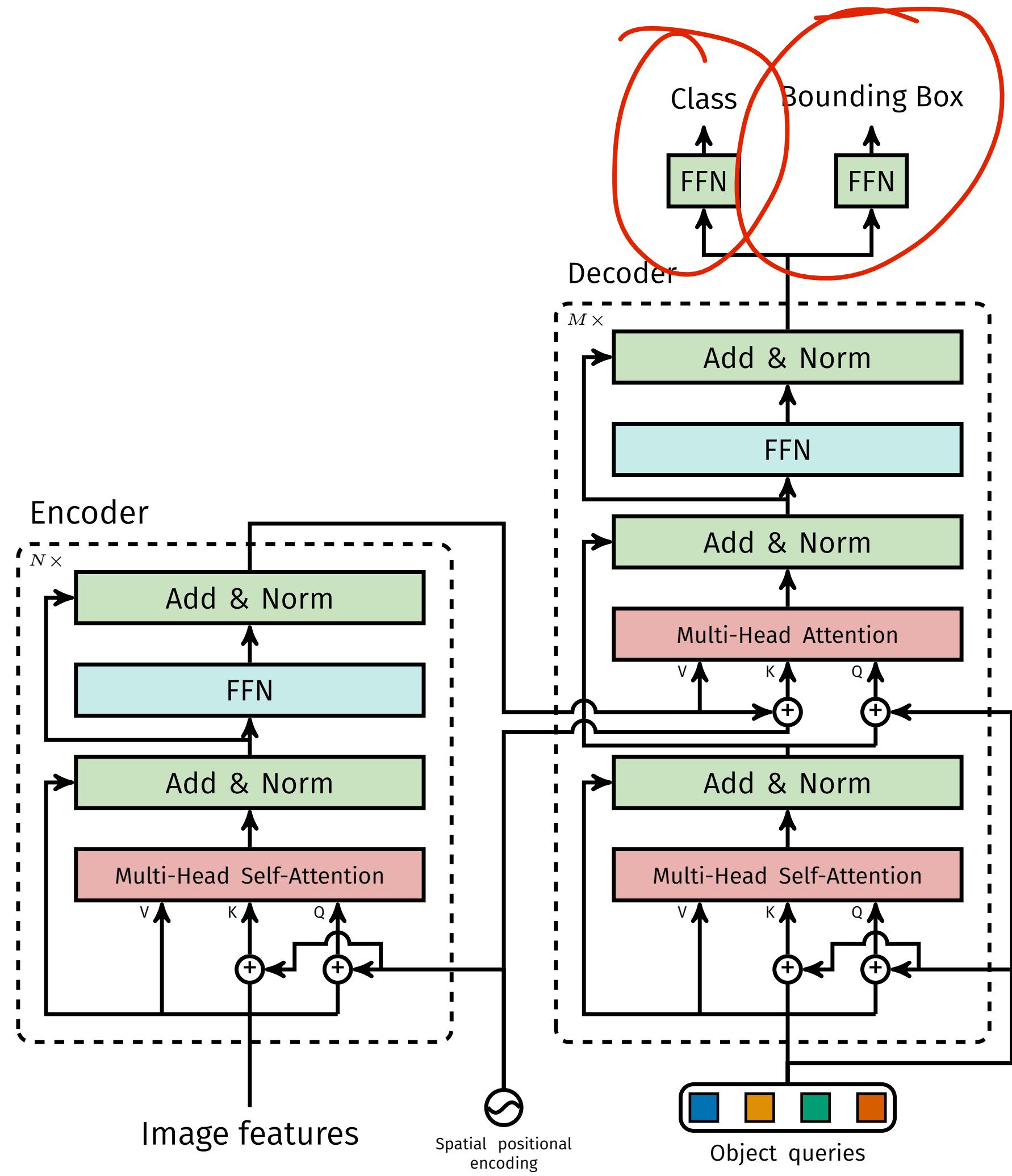
- The model flattens it and supplements it with a positional encoding before passing it to the Transformer.

DETR: A closer look



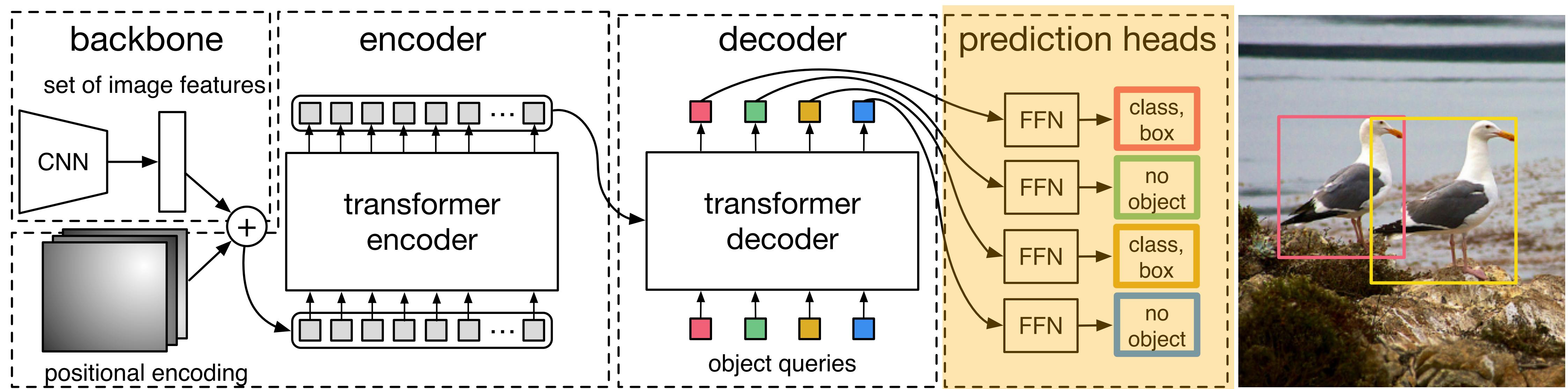
- The Transformer encodes the input sequence.
- We supply **object queries** (learnable positional encoding) to the decoder, which jointly attends to the encoder output.

DETR: Transformer architecture



- The Transformer is very similar to [Vaswani et al., 2017].
- Output (for each object query):
 - the class logits;
 - the bounding box (normalised coordinates).

DETR: A closer look



- Each output embedding of the decoder is passed to a shared feed forward network (FFN) that predicts either a detection (class and bounding box) or a “no object” class.

DETR: The loss function

- Hungarian matching between the predictions (queries) and the ground truth.

DETR: The loss function

- Hungarian matching between the predictions (queries) and the ground truth.



$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

Optimal assignment computed
in the first step

DETR: The loss function

- Hungarian matching between the predictions (queries) and the ground truth.

It also contain empty
ground-truth

$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

Loss for the class

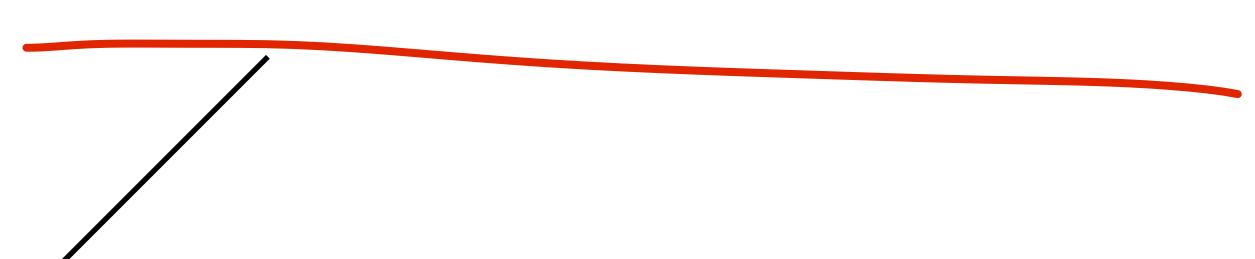


DETR: The loss function

- Hungarian matching between the predictions (queries) and the ground truth.

$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

Loss for the bounding-box



DETR: Bounding Box loss

- The bounding box loss $\mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)})$:

$$\lambda_{\text{iou}} \mathcal{L}_{\text{iou}}(b_i, \hat{b}_{\sigma(i)}) + \lambda_{\text{L1}} \|b_i - \hat{b}_{\sigma(i)}\|_1$$

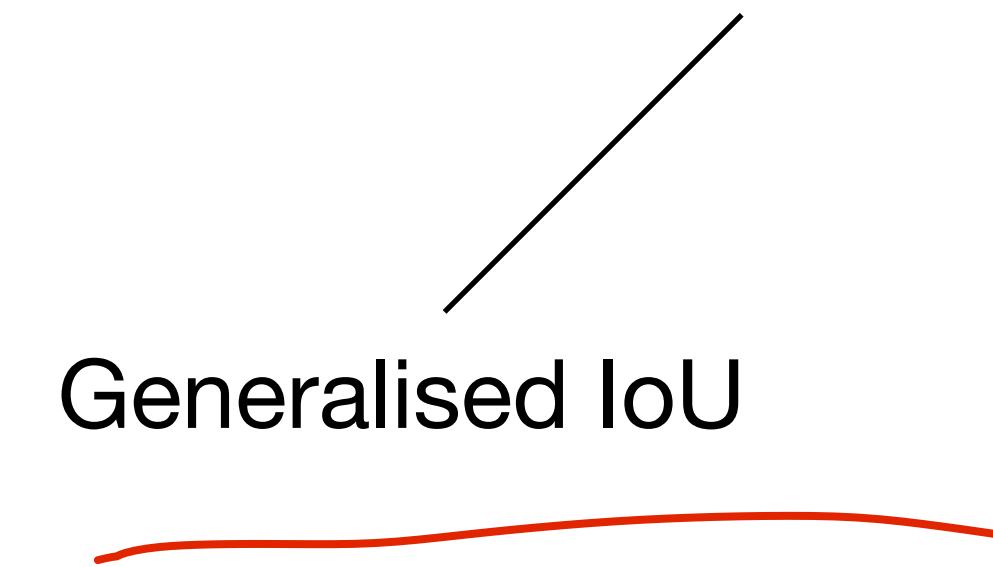
hyperparameters



DETR: Bounding Box loss

- The bounding box loss $\mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)})$:

$$\lambda_{\text{iou}} \mathcal{L}_{\text{iou}}(b_i, \hat{b}_{\sigma(i)}) + \lambda_{\text{L1}} \|b_i - \hat{b}_{\sigma(i)}\|_1$$

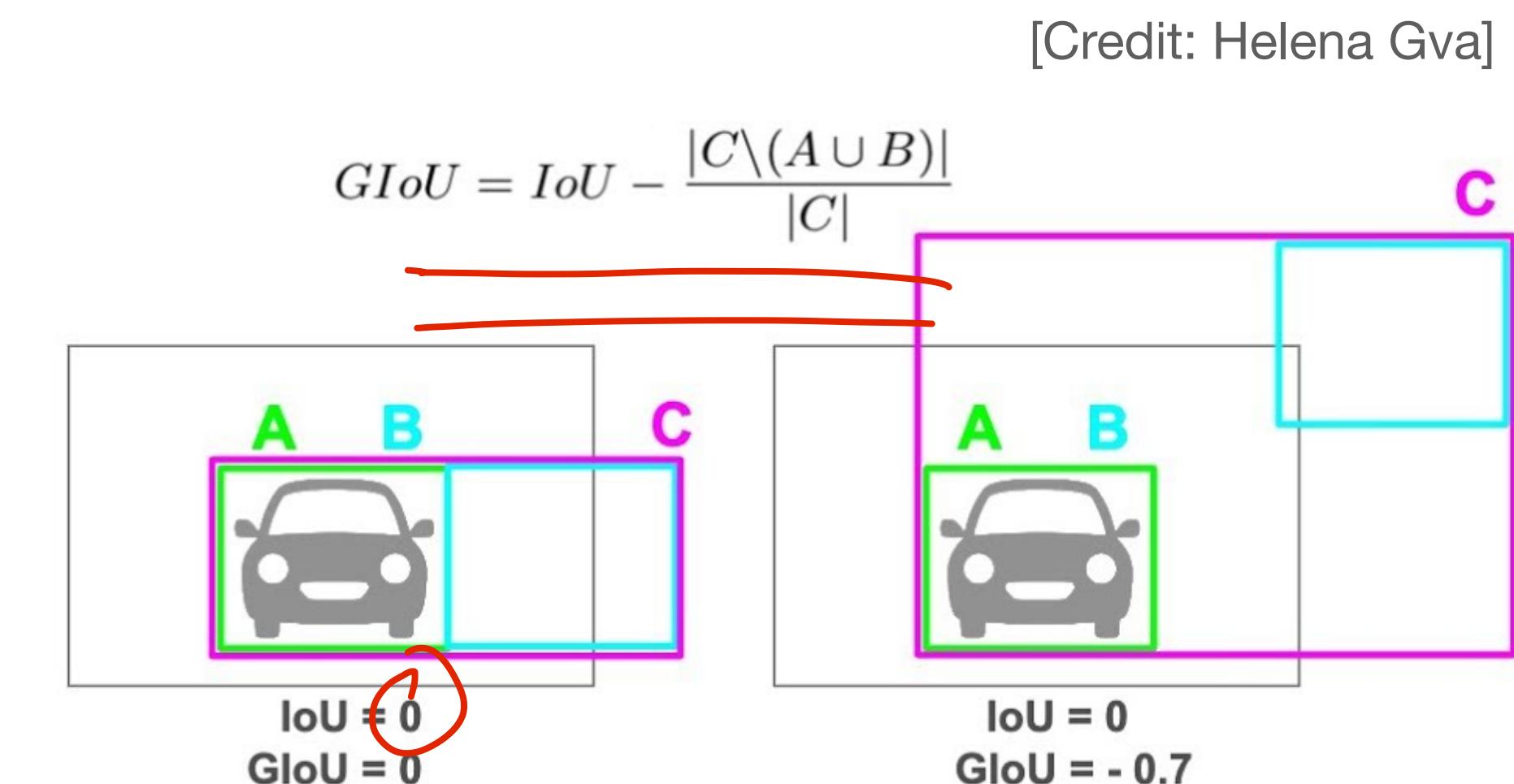


DETR: Bounding Box loss

- The bounding box loss $\mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)})$:

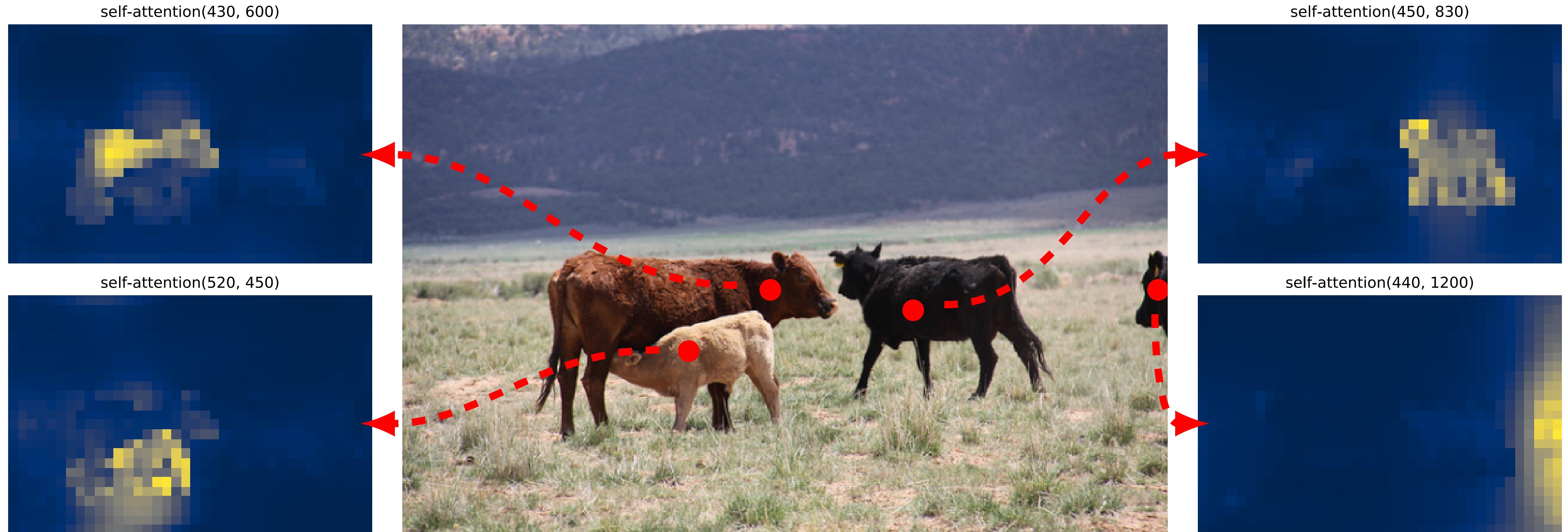
$$\lambda_{\text{iou}} \mathcal{L}_{\text{iou}}(b_i, \hat{b}_{\sigma(i)}) + \lambda_{\text{L1}} \|b_i - \hat{b}_{\sigma(i)}\|_1$$

Generalised IoU



Idea: prevent vanishing gradients when IoU is zero.

DETR: Qualitative results



The encoder is able to separate individual instances. Predictions are made with baseline DETR model on a validation set image.

DETR: Summary

- Accurate detection with a (relatively) simple architecture.
- No need for non-maximum suppression
 - We can simply disregard bounding boxes with “empty” class, or with low confidence.
- Accurate panoptic segmentation with a minor extension.

Issues:

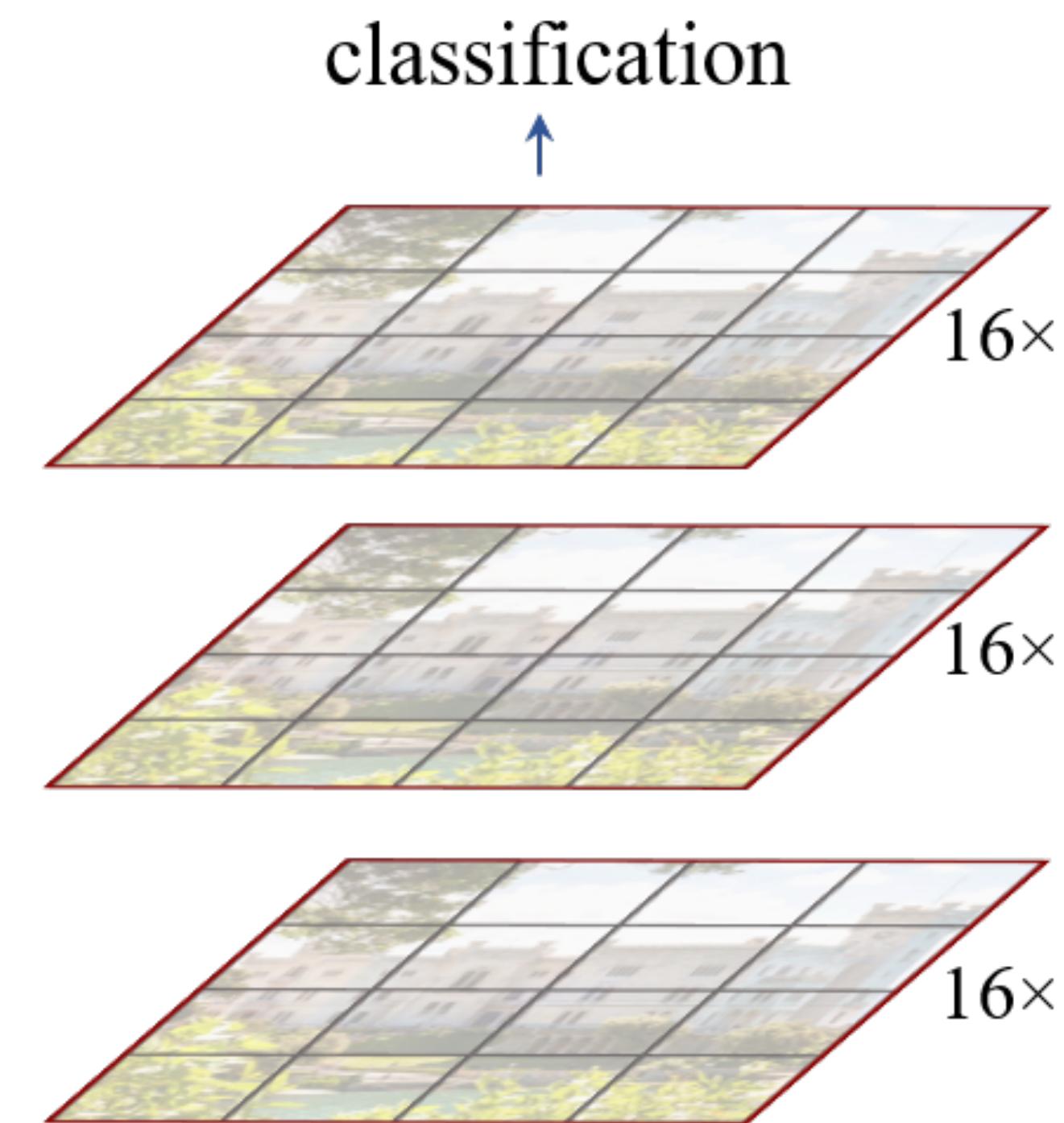
- High computational and memory requirements (especially the memory).
- Slow convergence / long training.

See: Zhu et al., “Deformable DETR: Deformable Transformers for End-to-End Object Detection” (ICLR 2021).

Swin Transformer

Recall ViT:

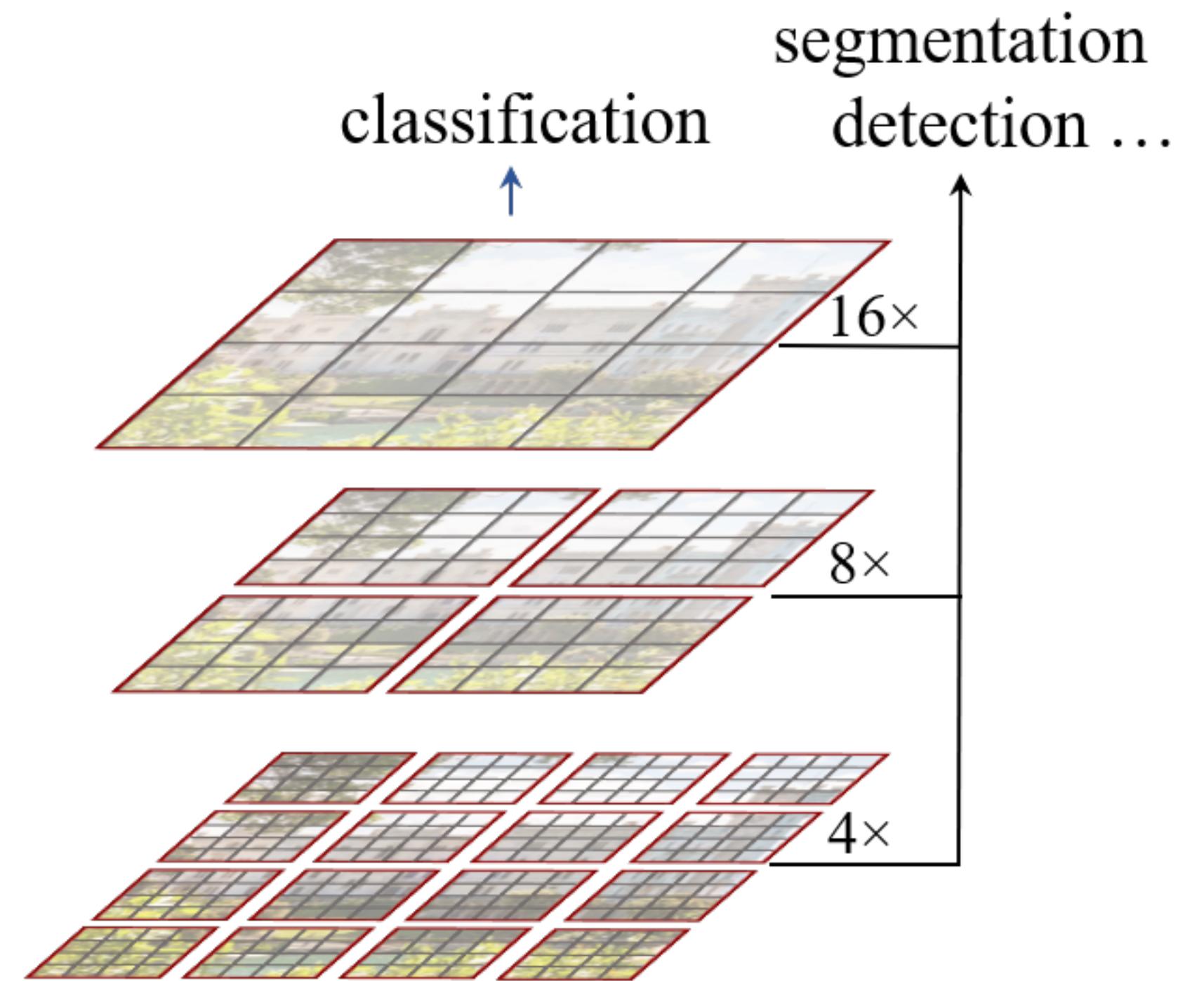
- We process patch embeddings with a series of self-attention layers;
 - Quadratic computational complexity. *= patch size*
- The number of tokens remains the same in all layers (QUIZ: How many / what does it depend on?)
- In CNNs, we gradually decrease the resolution, which has computation benefits (and increases the receptive field size).
- Do ViTs benefit from the same strategy?



Swin Transformer

Swin Transformer:

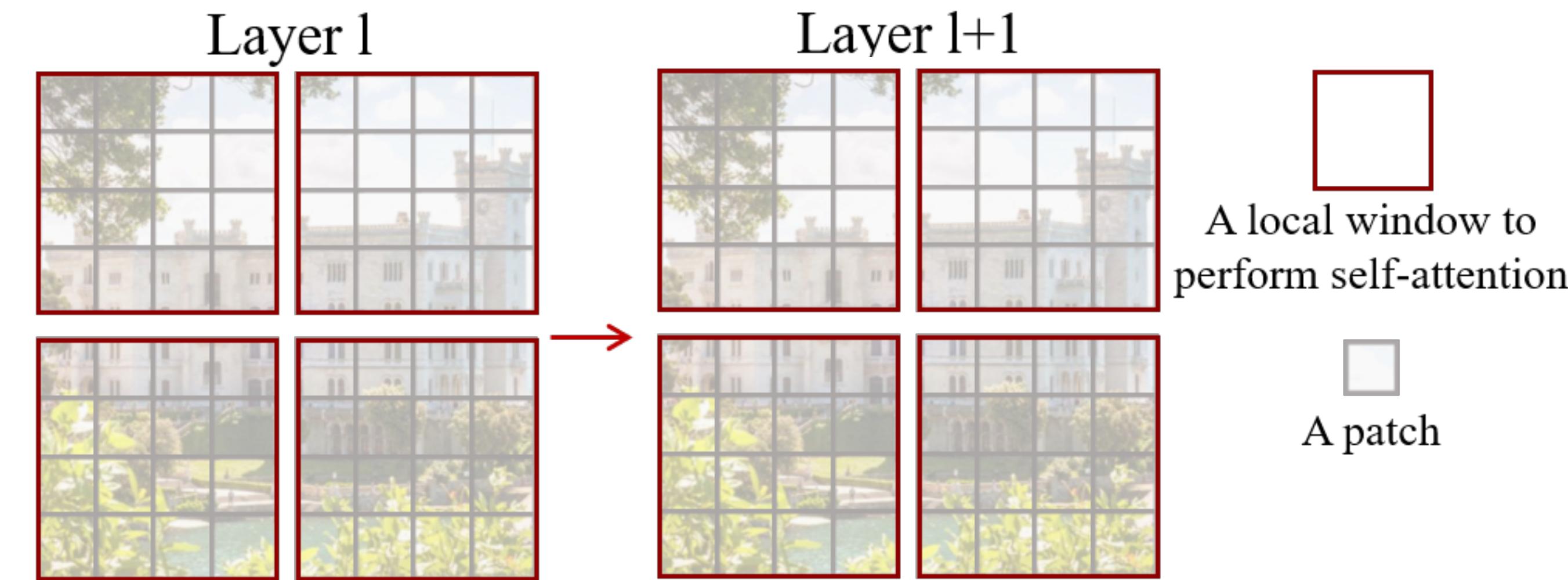
- Construct a hierarchy of image patches.
- Attention windows have a fixed size.
- Linear computational complexity.
- Output representation is compatible with standard backbones (e.g. ResNet)
 - We can test many downstream tasks (e.g. object detection)



Swin Transformer

A naive implementation will hurt context reasoning:

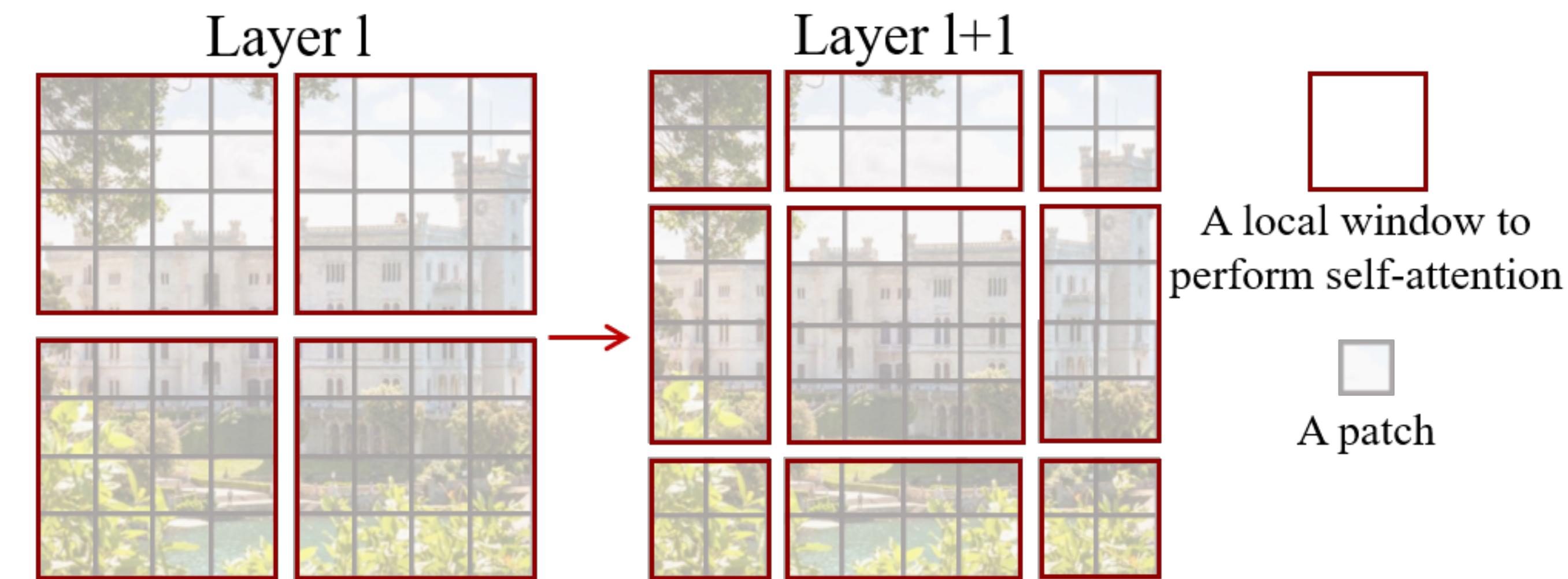
- At any given level (except the last one), our context is not global anymore:



Swin Transformer

A naive implementation will hurt context reasoning:

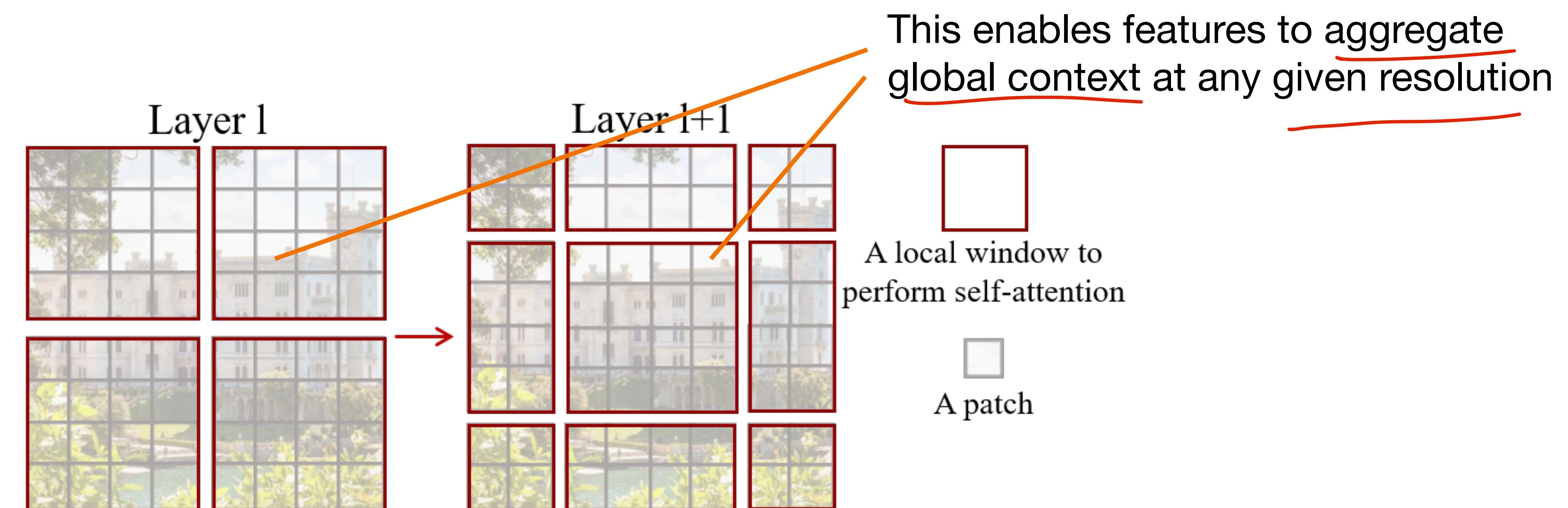
- Solution: alternate the layout of local windows:



Swin Transformer

A naive implementation will hurt context reasoning:

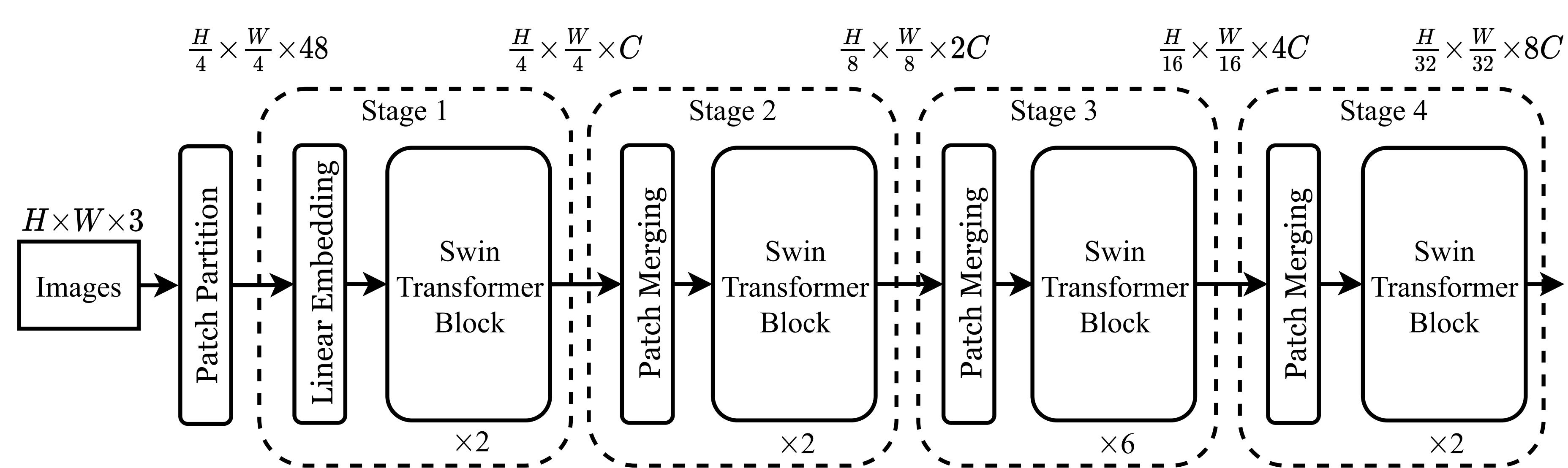
- Solution: alternate the layout of local windows:



Swin Transformer

Successively decrease the number of tokens using patch merging:

- concatenate 2×2 C-dim patches into a feature vector ($4 \times C$ -dim);
- linearly transform to a $2 \times C$ dimensional vector.



Swin Transformer: Results

- More efficient and accurate than ViT and some of CNNs (despite using lower resolution input).
- Note: No pre-training on large datasets!
- Improved scalability on large datasets (compare ImageNet-1K and ImageNet-22K)

ImageNet-22K

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
R-101x3 [38]	384^2	388M	204.6G	-	84.4
R-152x4 [38]	480^2	937M	840.5G	-	85.4
ViT-B/16 [20]	384^2	86M	55.4G	85.9	84.0
ViT-L/16 [20]	384^2	307M	190.7G	27.3	85.2
Swin-B	224^2	88M	15.4G	278.1	85.2
Swin-B	384^2	88M	47.0G	84.7	86.4
Swin-L	384^2	197M	103.9G	42.1	87.3

ImageNet-1K

method	image size	#param.	FLOPs	throughput (image / s)	ImageNet top-1 acc.
RegNetY-4G [48]	224^2	21M	4.0G	1156.7	80.0
RegNetY-8G [48]	224^2	39M	8.0G	591.6	81.7
RegNetY-16G [48]	224^2	84M	16.0G	334.7	82.9
EffNet-B3 [58]	300^2	12M	1.8G	732.1	81.6
EffNet-B4 [58]	380^2	19M	4.2G	349.4	82.9
EffNet-B5 [58]	456^2	30M	9.9G	169.1	83.6
EffNet-B6 [58]	528^2	43M	19.0G	96.9	84.0
EffNet-B7 [58]	600^2	66M	37.0G	55.1	84.3
ViT-B/16 [20]	384^2	86M	55.4G	85.9	77.9
ViT-L/16 [20]	384^2	307M	190.7G	27.3	76.5
DeiT-S [63]	224^2	22M	4.6G	940.4	79.8
DeiT-B [63]	224^2	86M	17.5G	292.3	81.8
DeiT-B [63]	384^2	86M	55.4G	85.9	83.1
Swin-T	224^2	29M	4.5G	755.2	81.3
Swin-S	224^2	50M	8.7G	436.9	83.0
Swin-B	224^2	88M	15.4G	278.1	83.5
Swin-B	384^2	88M	47.0G	84.7	84.5

Swin Transformer: Summary

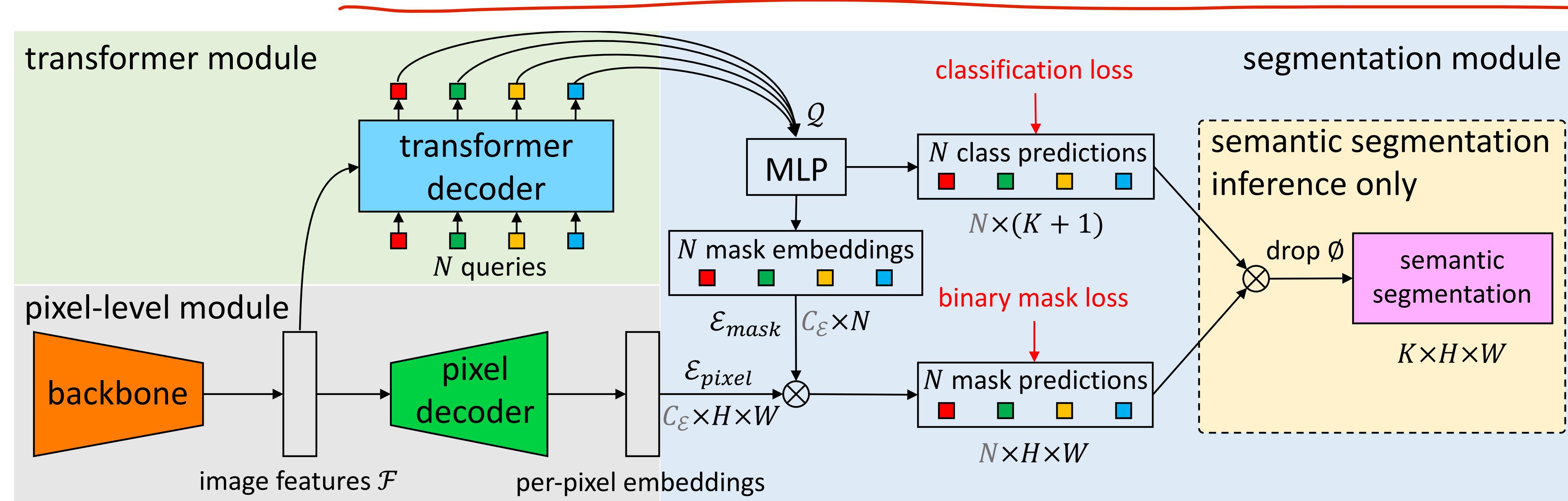
- Reconciles the inductive biases from CNNs with Transformers;
- State-of-the-art on image classification, object detection and semantic segmentation;
- Linear computational complexity.

Swin Transformer: Summary

- Reconciles the inductive biases from CNNs with Transformers;
- State-of-the-art on image classification, object detection and semantic segmentation;
- Linear computational complexity.
- Demonstrates that Transformers are strong vision models across a range of classic downstream applications.
- Best paper award at ICCV '21.

MaskFormer

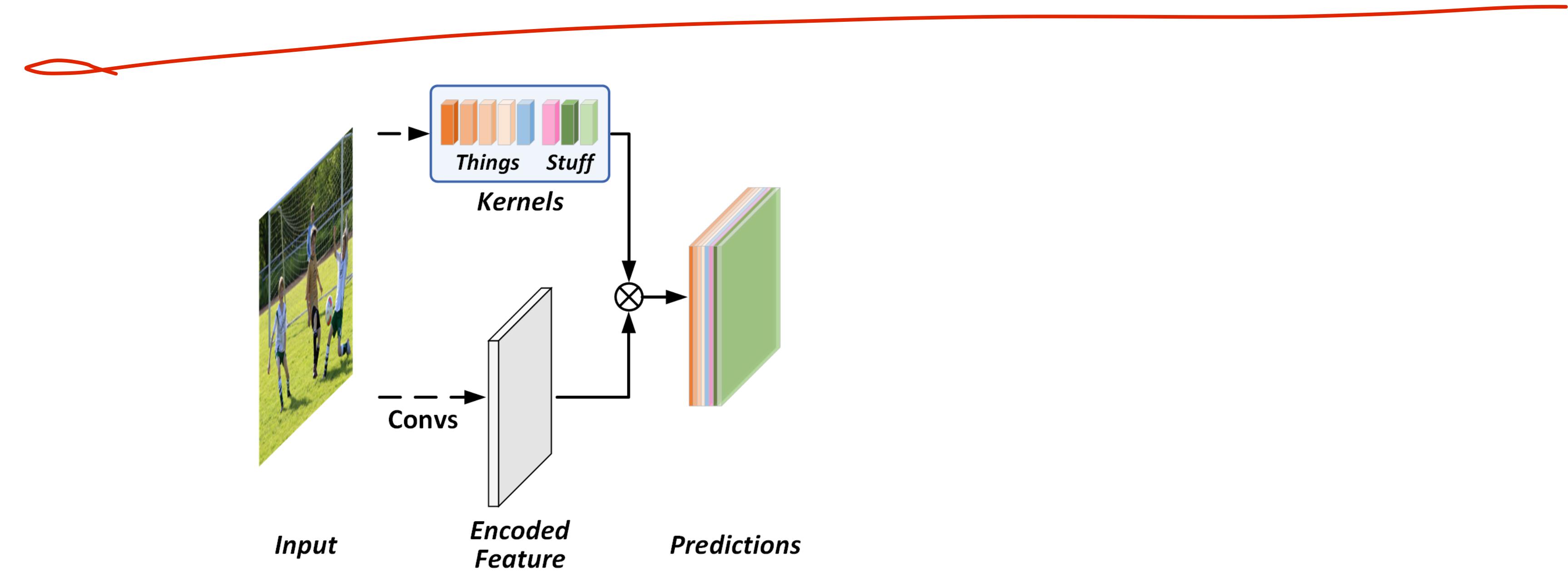
MaskFormer: a unified model for semantic and panoptic segmentation



[Cheng et al., 2021]

Recall: Panoptic FCN

Panoptic FCN: a unified model for semantic and panoptic segmentation

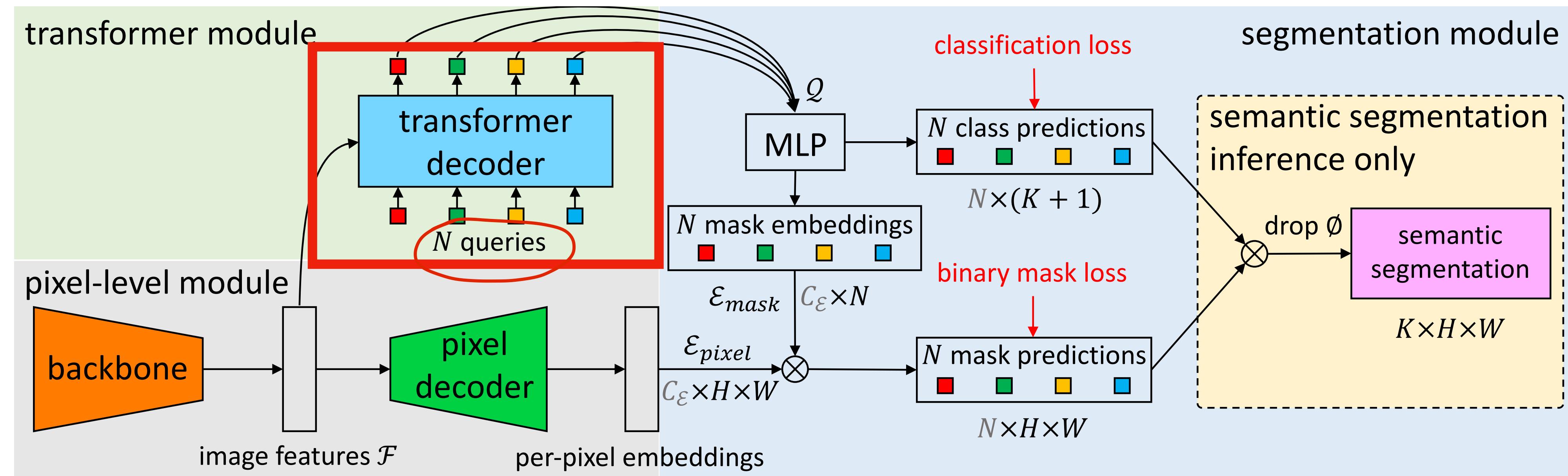


MaskFormer's idea: Compute the kernels using learnable queries with a Transformer.

[Cheng et al., 2021]

MaskFormer

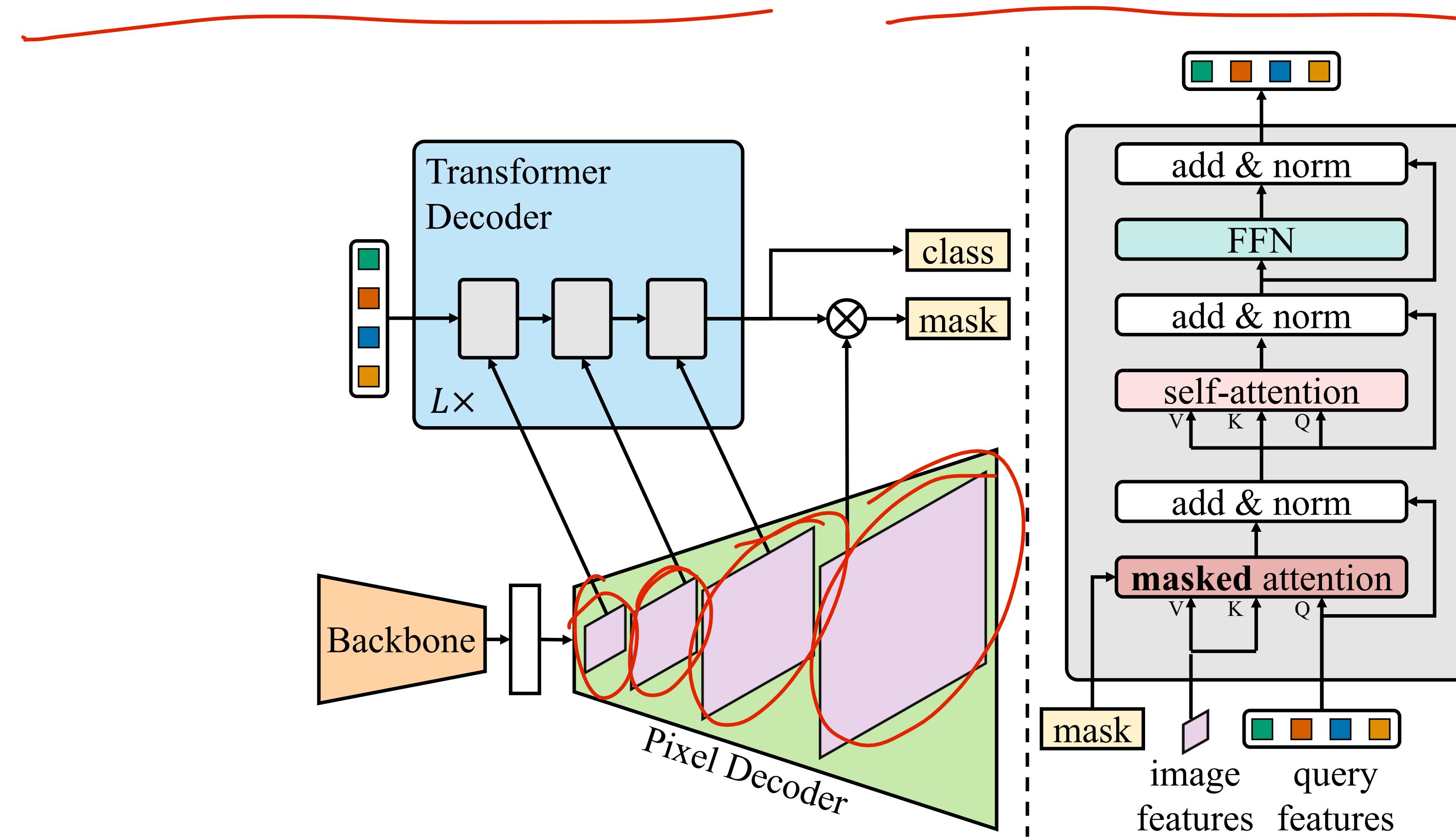
MaskFormer: a unified model for semantic and panoptic segmentation



Works well, but has troubles with small objects/segments.

Mask2Former

Idea: Attend with self-attention at multiple scales of the feature hierarchy:



[Cheng et al., 2021]

Masked attention

Idea: constrain attention only to the foreground area
 corresponding to the query.

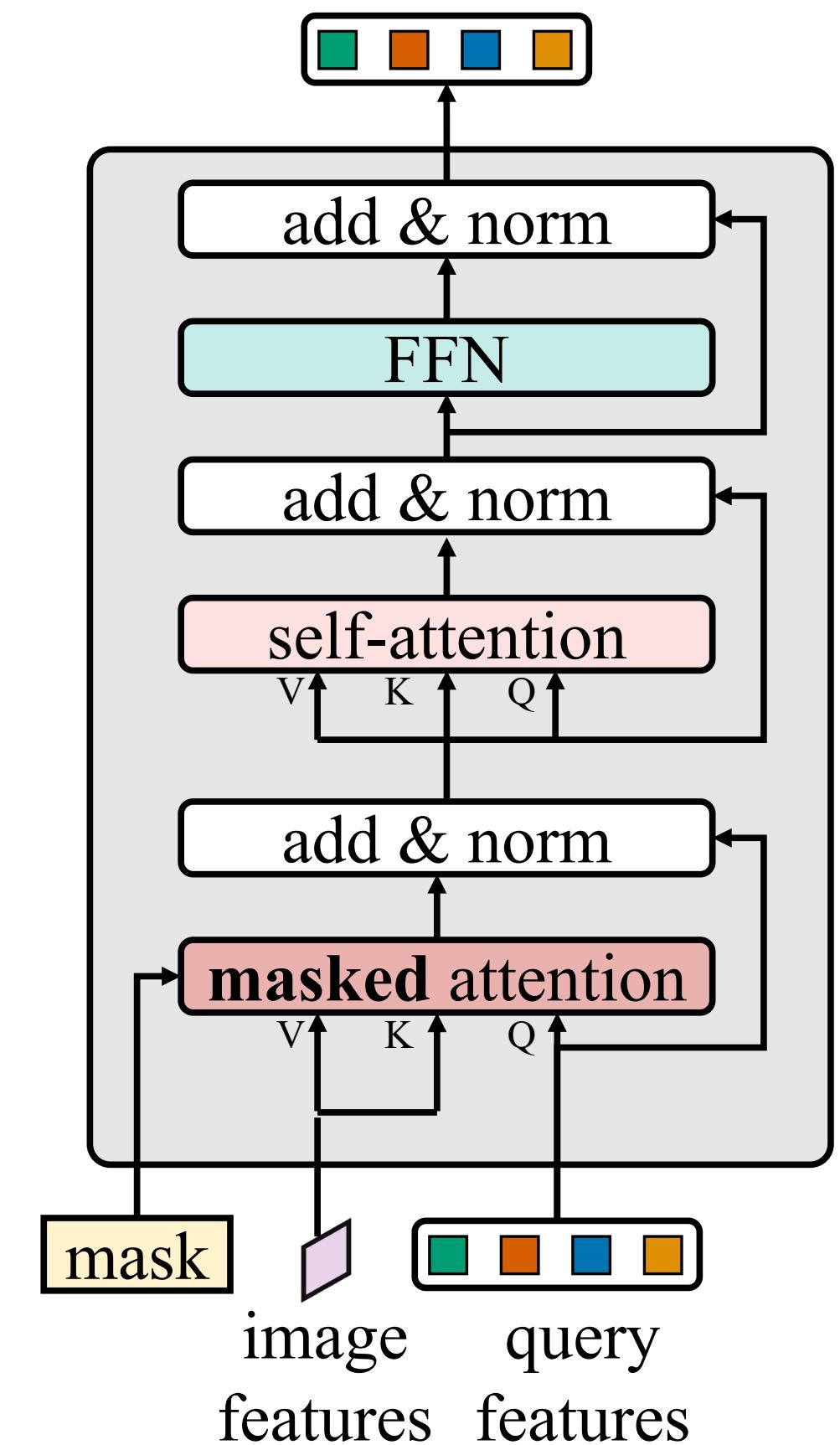
Standard self-attention: $\mathbf{X}_l = \text{softmax}(\mathbf{Q}_l \mathbf{K}_l^T) \mathbf{V}_l + \mathbf{X}_{l-1}$

Masked attention: $\mathbf{X}_l = \text{softmax}(\mathcal{M}_{l-1} + \mathbf{Q}_l \mathbf{K}_l^T) \mathbf{V}_l + \mathbf{X}_{l-1}$

$$\mathcal{M}_{l-1}(x, y) = \begin{cases} 0 & \text{if } \mathbf{M}_{l-1}(x, y) = 1 \\ -\infty & \text{otherwise} \end{cases}$$

Where does this come from?

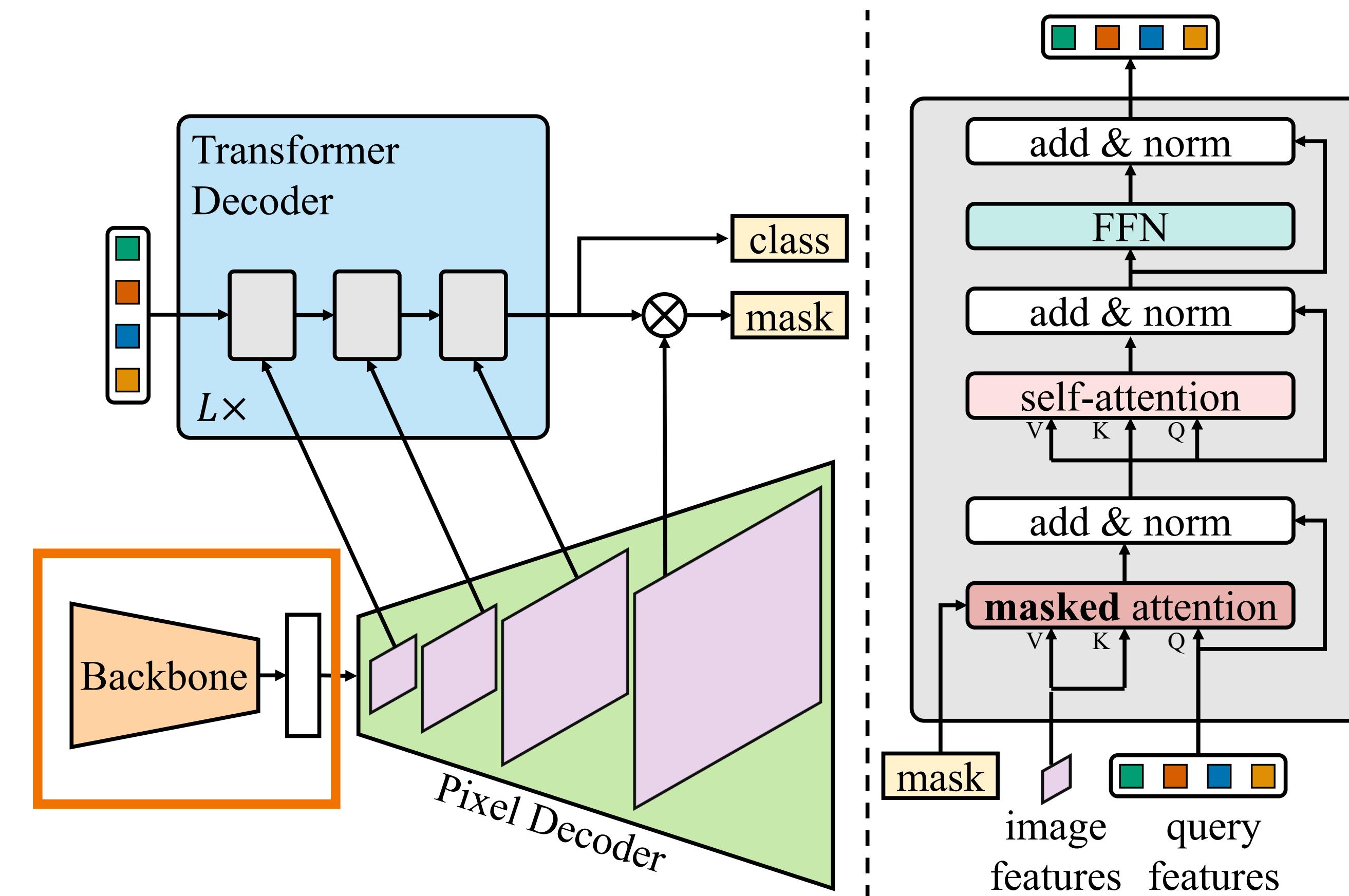
mask binarisation



[Cheng et al., 2021]

Mask2Former

Idea: Attend with self-attention at multiple scales of the feature hierarchy:

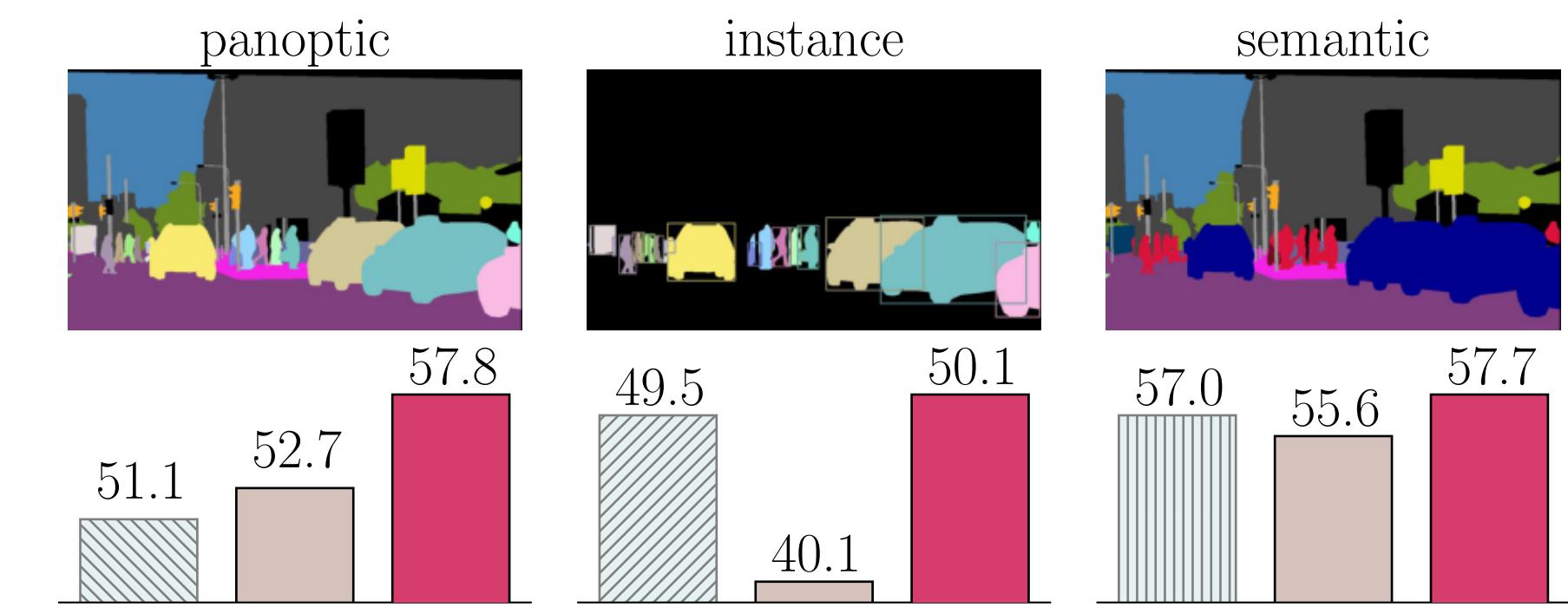


[Cheng et al., 2021]

Mask2Former: Summary

A unified model:

- note that we do not talk about “stuff” and “things” anymore
- queries abstract those notions away.
- Achieves state-of-the-art accuracy across segmentation tasks:



Universal architectures:

■ Mask2Former (ours) ■ MaskFormer

SOTA specialized architectures:

■ Max-DeepLab ■ Swin-HTC++ ■ BEiT

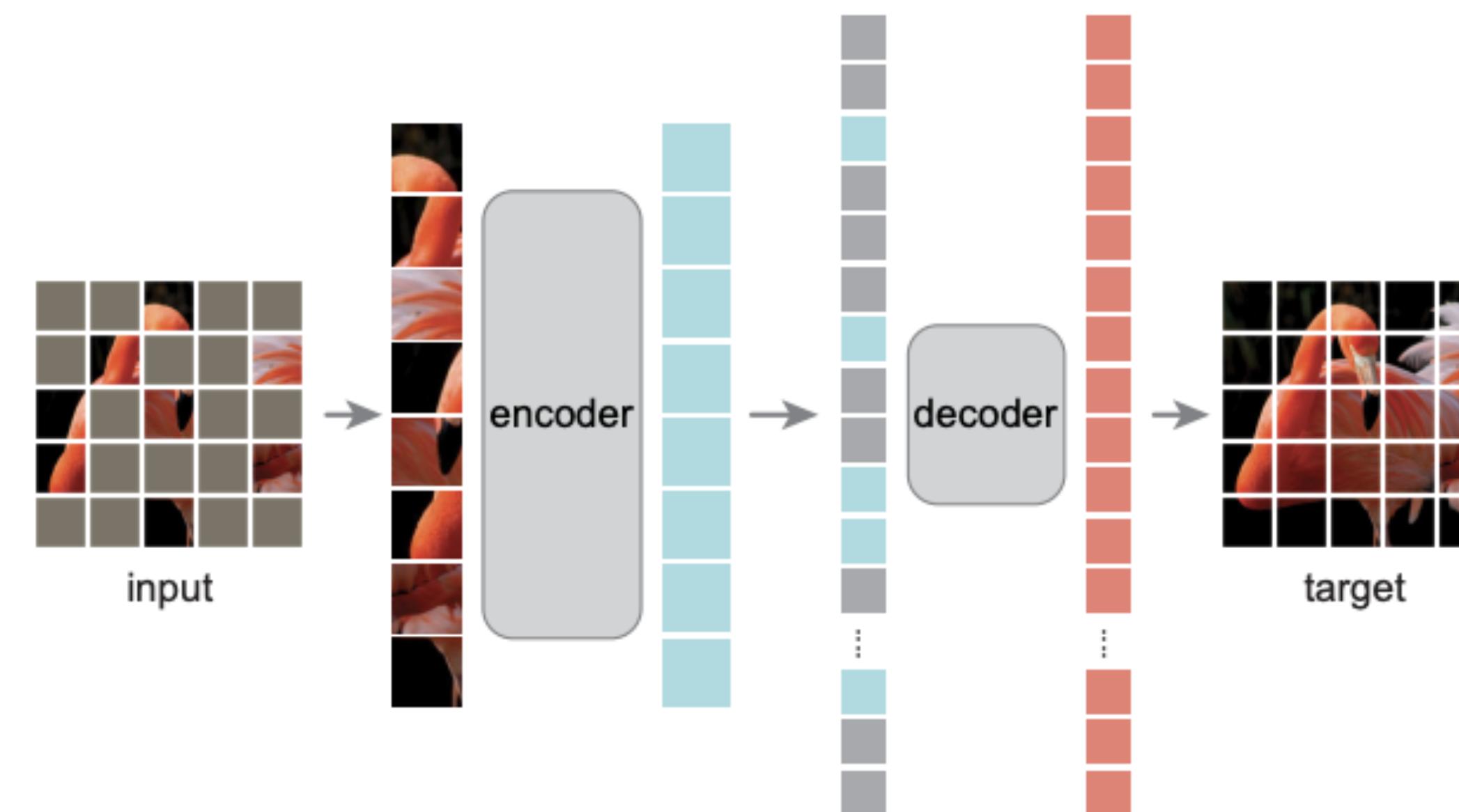
Conclusions

- Transformers have revolutionised the field of NLP, achieving incredible results.
- We observe massive impact on computer vision (DETR, ViT).
- Complementing CNNs, Transformers have reached state-of-the-art in object classification, detection, tracking and image generation.

Conclusions

- Transformers have revolutionised the field of NLP, achieving incredible results.
- We observe massive impact on computer vision (DETR, ViT).
- Complementing CNNs, Transformers have reached state-of-the-art in object classification, detection, tracking and image generation.
- A grain of salt: This often comes at an increased computational budget (larger GPUs, longer training).

Applications: Unsupervised learning



[He et al., 2021]

Round-up

Reminders:

- Ongoing course evaluation (deadline: this Friday, 22.12).
- Exam registration (deadline: 15.01).
- Next lecture: January 9, 2023.

Happy holidays!



Computer Vision III: Transformers

Nikita Araslanov
19.12.2023

Adapted from:
Prof. Laura Leal-Taixé
<https://dvl.in.tum.de>

