



ch14: Computer Vision

Image Segmentation

naive idea及问题 (P5-P7)

- Sliding Window
- Fully Convolutional

Upsampling (P9-12)

- Unpooling

Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

...
Rest of the
network

Max Unpooling

Use positions from pooling layer

填到相应的位置, 其他位置补0

1	2
3	4

Input: 2 x 2

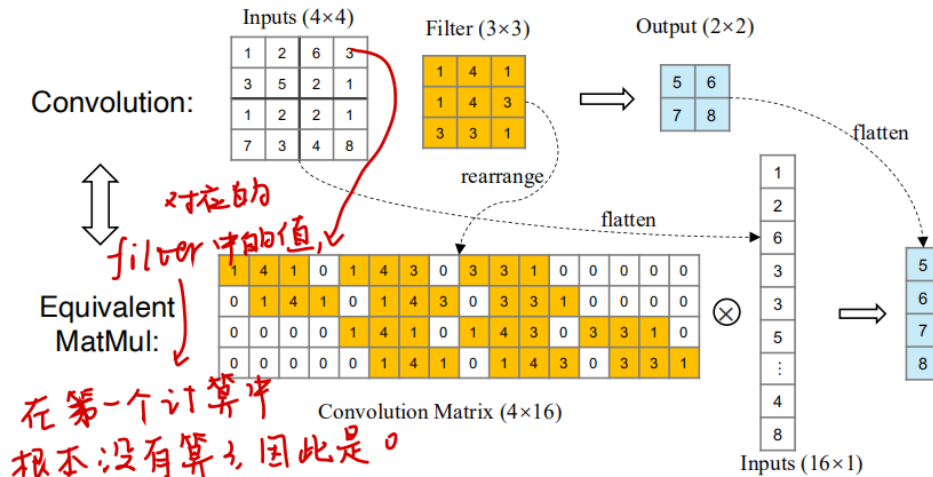
0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

Output: 4 x 4

- Transposed Convolution

• Convolution Matrix

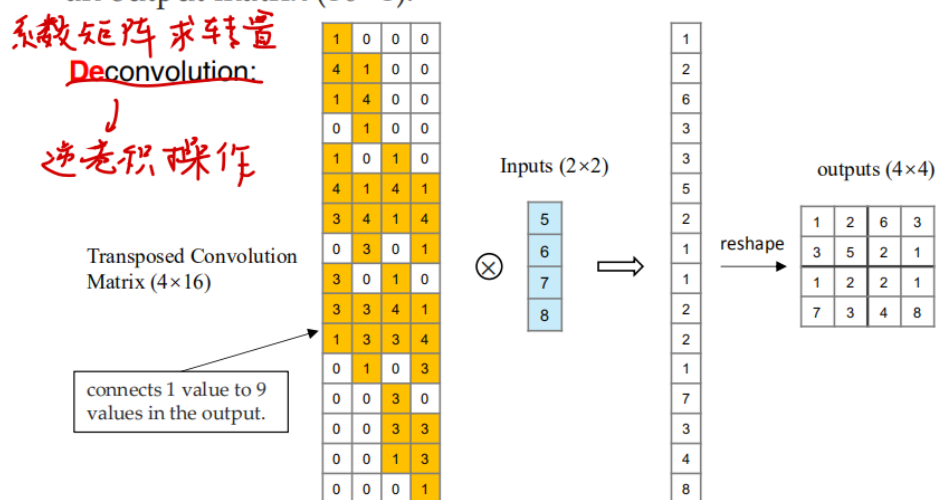
We can express a convolution operation using a **matrix multiplication** by rearranging the kernel matrix into a general matrix.



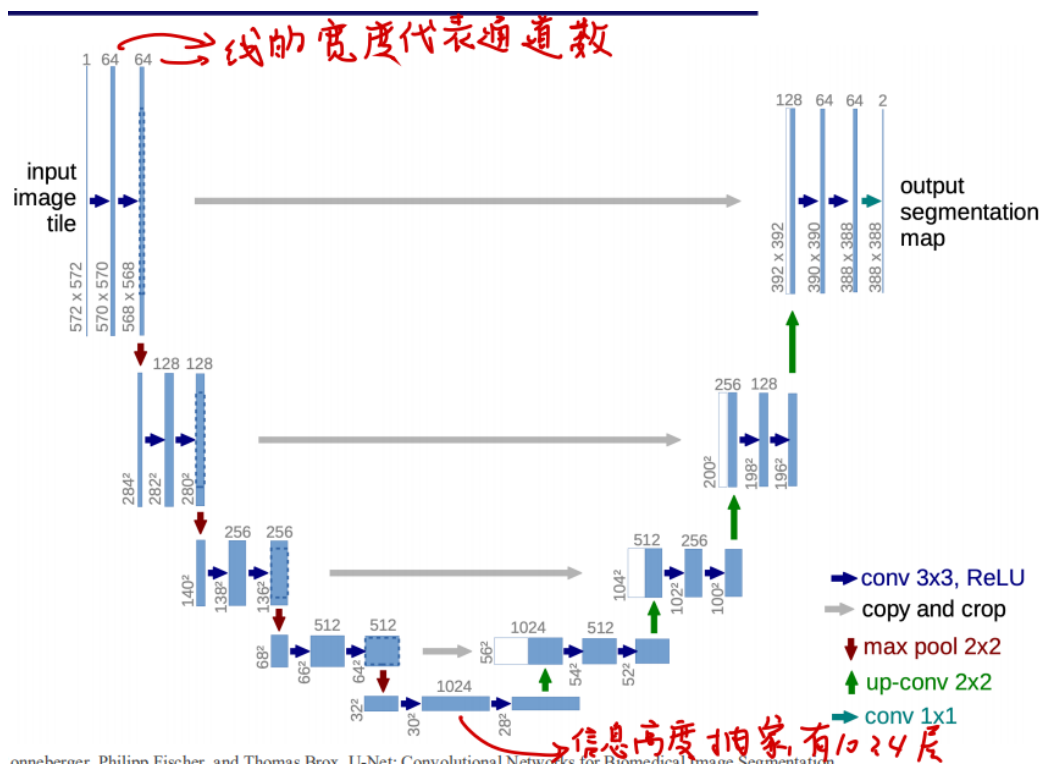
• Transposed Convolution Matrix

Transpose the convolution matrix C (4×16) to C^T (16×4).

We can matrix-multiply C^T (16×4) with a column vector (4×1) to generate an output matrix (16×1).



U-Net (P15)

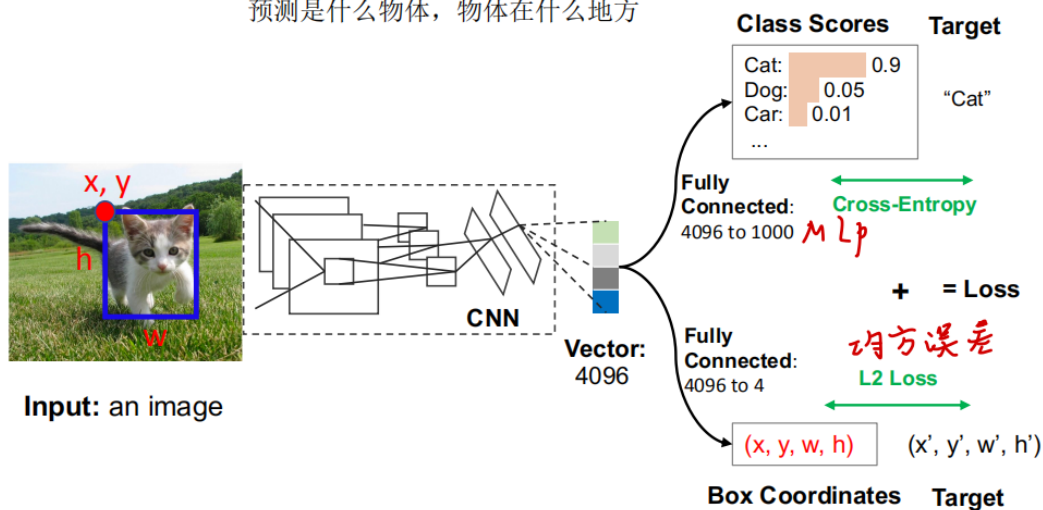


Object Detection

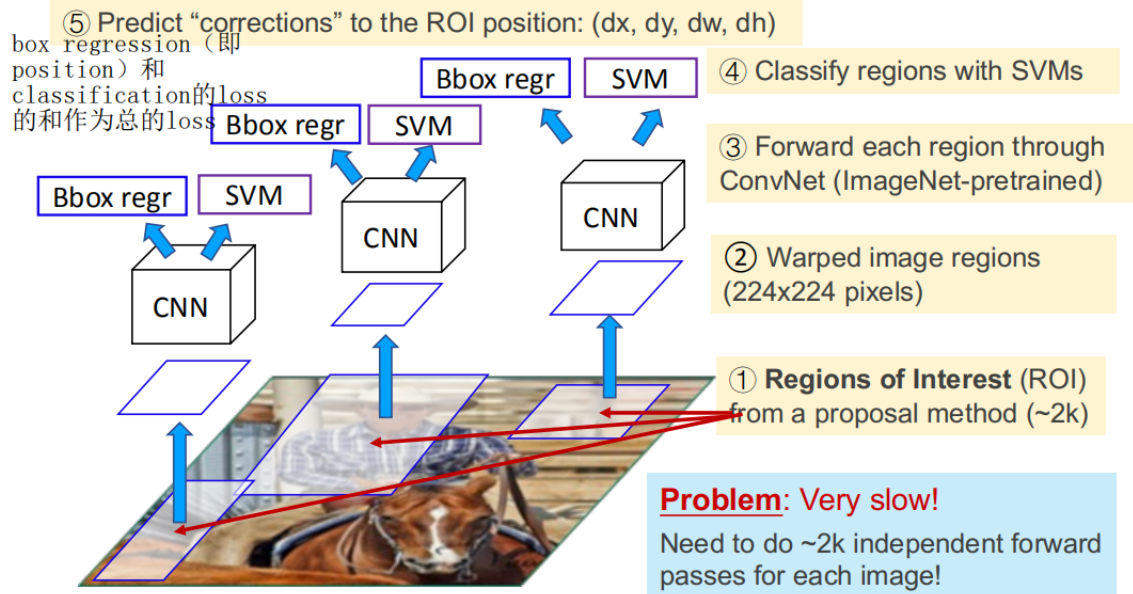
Single Object (P19)

- Object detection = Classification + Localization

预测是什么物体，物体在什么地方

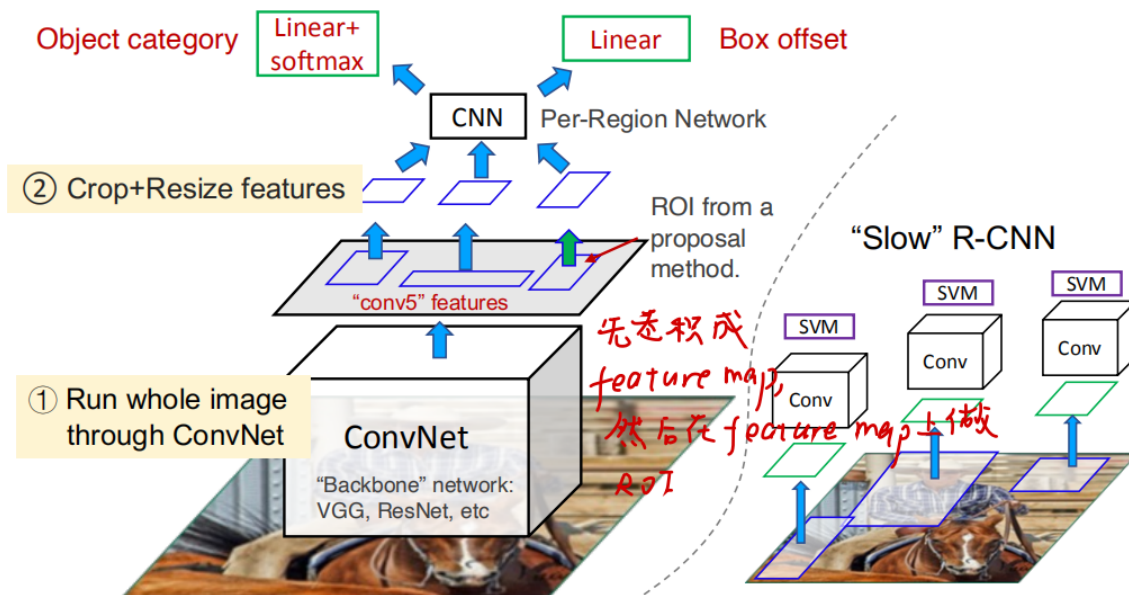


R-CNN (P22)



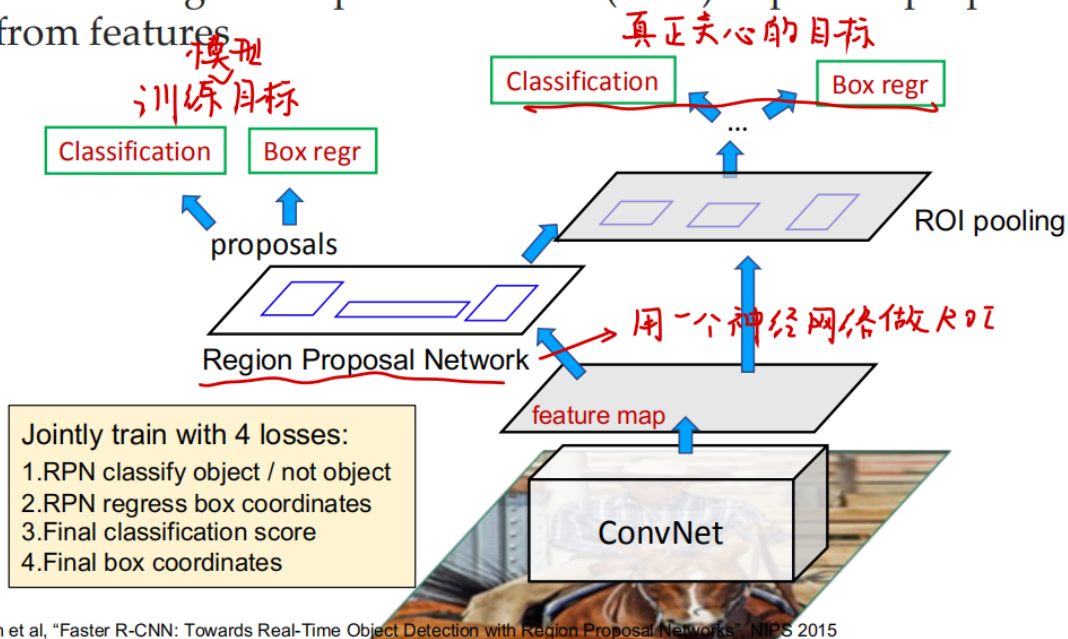
Fast R-CNN (P23)

- Idea:** Pass the image through convnet before cropping!



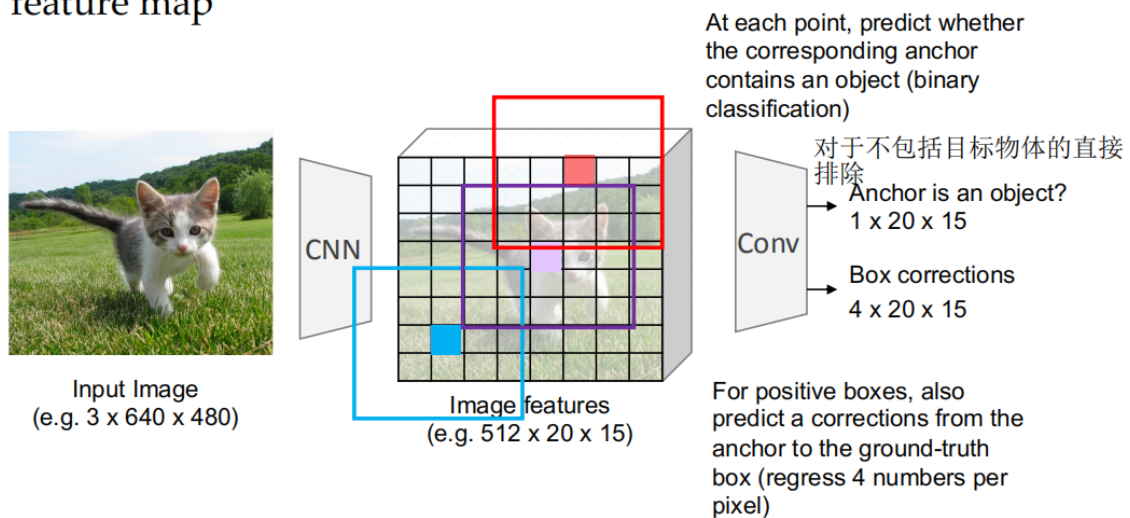
Faster R-CNN (P24)

- Insert a Region Proposal Network (RPN) to predict proposals from features



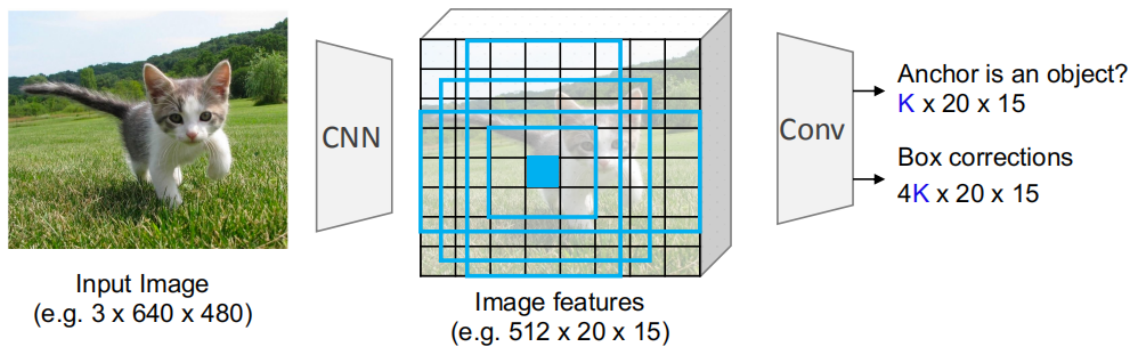
Region Proposal Network (P25-26)

- Imagine an **anchor box** of fixed size at each point in the feature map



- In practice use **K** different anchor boxes of different size / scale at each point 每个关注像素点周围有K个大小不同的box 选择不同的范围

每个关注像素点周围有K个大小不同的box, 选择不同的范围

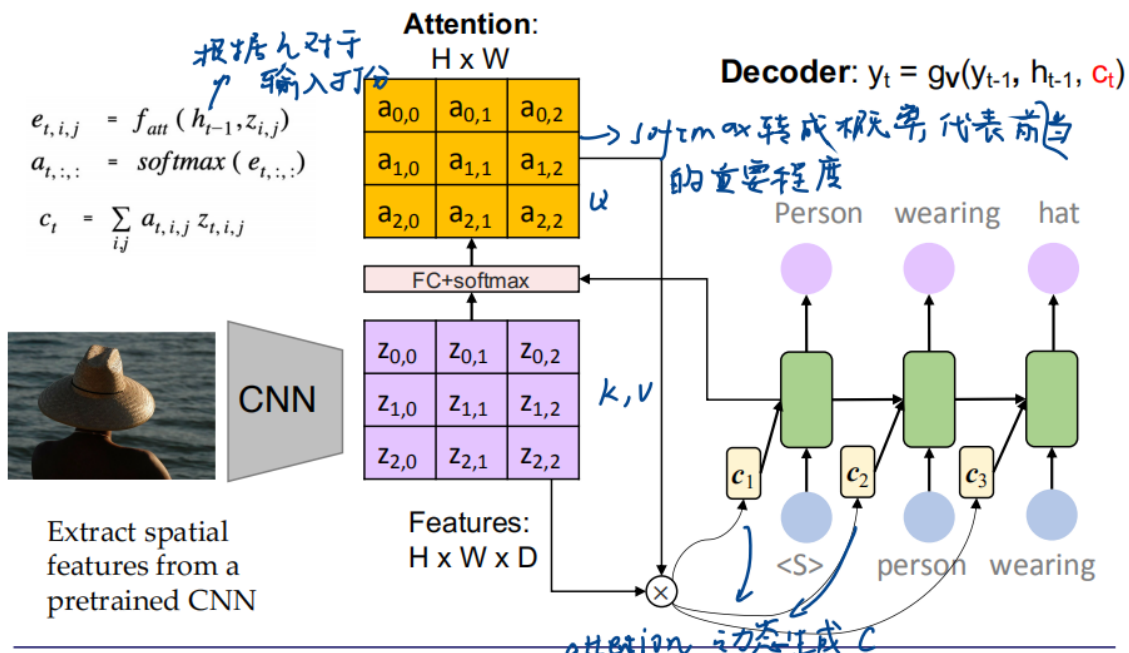


- Sort the $K \times 20 \times 15$ boxes by their “objectness” score, take top ~ 300 as our proposals

Image Captioning

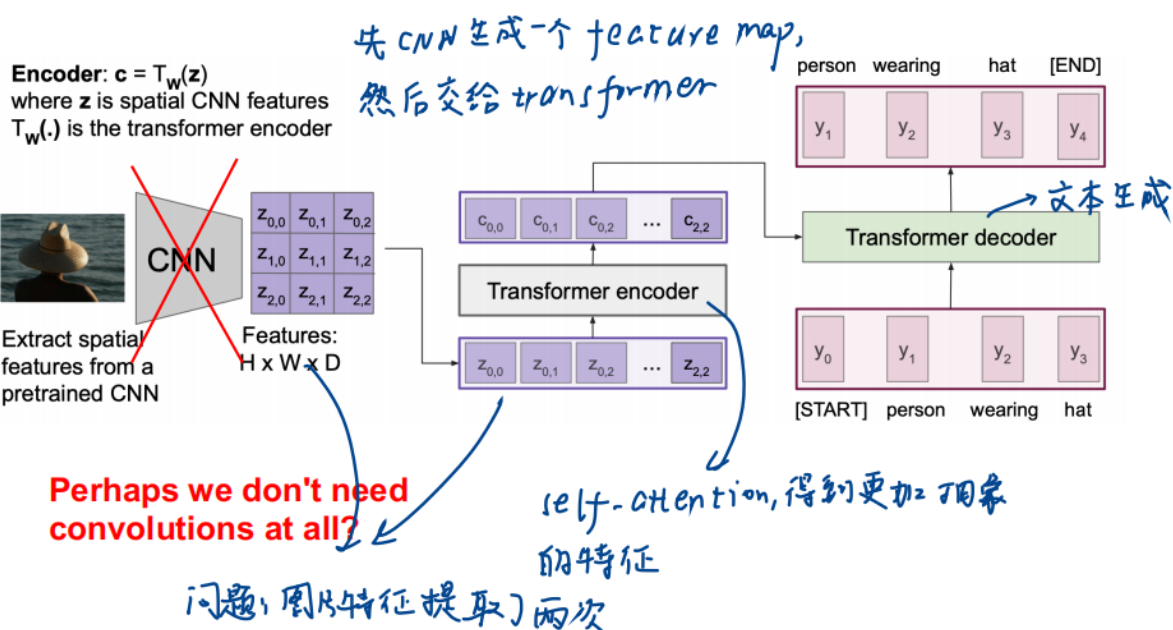
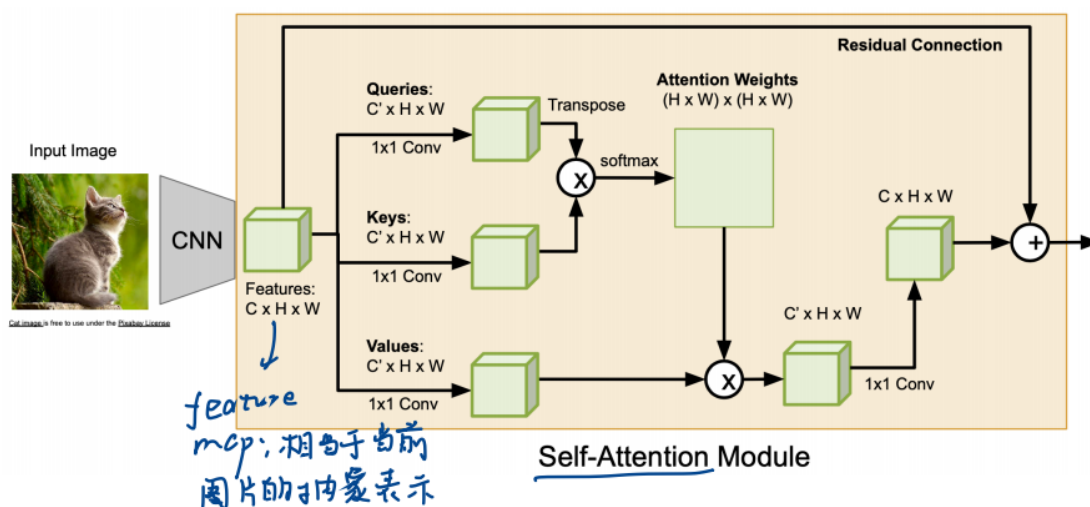
Encoder-Decoder with Attention (P31)

- Each context vector will attend to different image regions

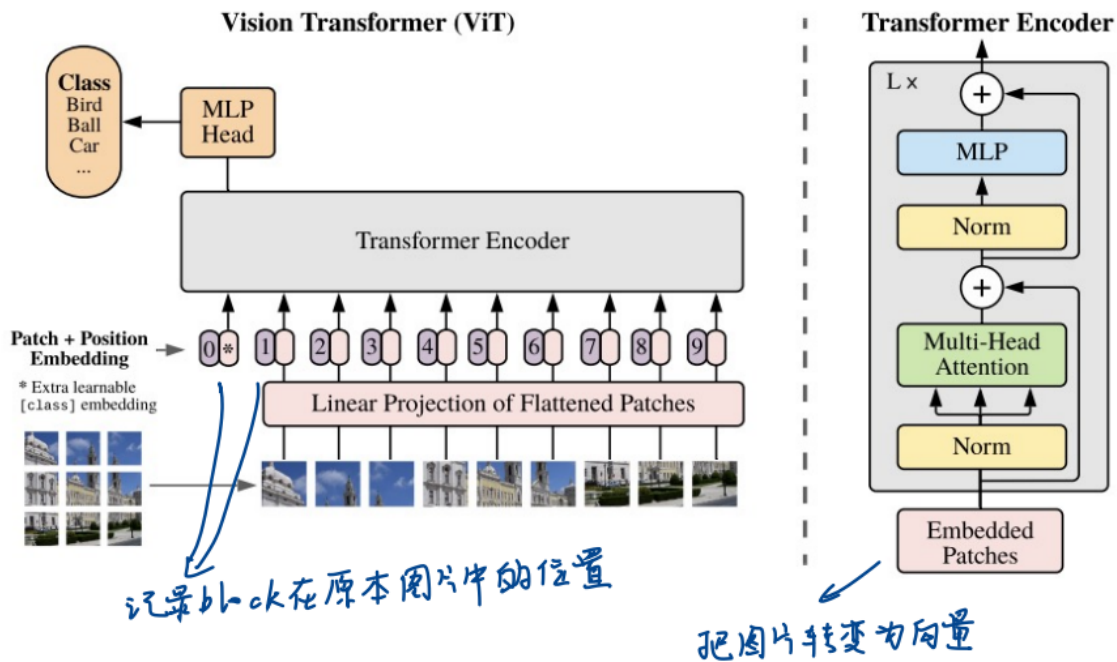


CNN with Self-Attention (P33-34)

- We can also apply **self-attention** to CNN features.

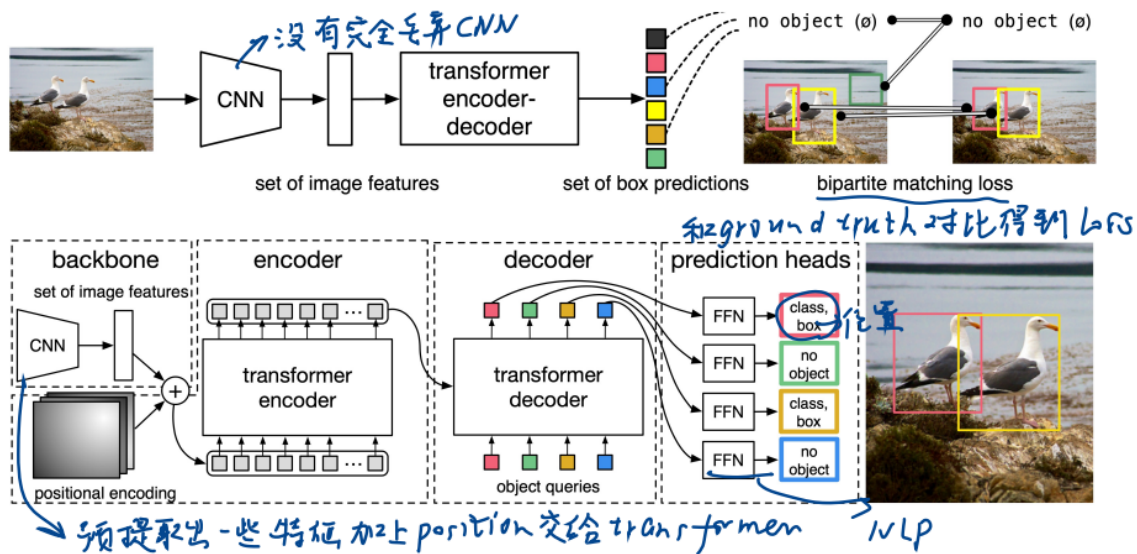


Vision Transformer (ViT) (P36)



Vision Transformers for Object Detection (P36)

• DETR



BEiT: Pre-Trained Models for Images (P39)

- Backbone Network: Standard Transformer
- Pre-Training: Masked Image Modeling

