# Transform in Computer Vision

分享人：梁瑛平          时间：2021-04-11

BEIJING INSTITUTE OF TECHNOLOGY

北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

# 目录

## CONTENTS

北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

《Attention Is All You Need》

《Attention Is All You Need》

## Encoder：

$$LayerNorm(x + Sublayer(x))$$

## Attention：

$$Attention(Q, K, V) = softmax(\frac{Q^T K}{\sqrt{d_k}})V$$

$$Q \in R^{m \times d_k} \, , \, K \in R^{m \times d_k} \, , \, V \in R^{m \times d_v}$$

《Attention Is All You Need》

$$Attention(Q, K, V) = softmax(\frac{Q^T K}{\sqrt{d_k}})V$$

$$Q \in R^{m \times d_k} \ , \ K \in R^{m \times d_k} \ , \ V \in R^{m \times d_v}$$

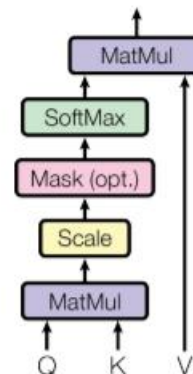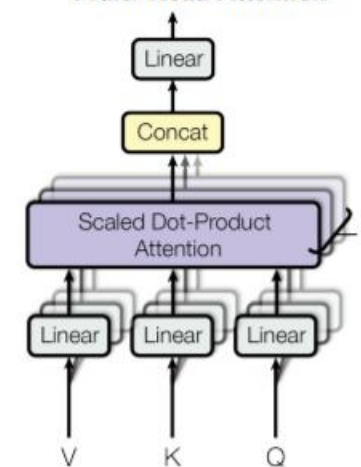**Scaled Dot-Product Attention**

MatMul

SoftMax

Mask (opt.)

Scale

MatMul

Q  K  V

**Multi-Head Attention**

Linear

Concat

Scaled Dot-Product Attention

Linear  Linear  Linear

V  K  Q

北京理工大学
BEIJING INSTITUTE OF TECHNOLOGY

《Attention Is All You Need》

## Input：



## Self-Attention：

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V$$

$$Z =$$

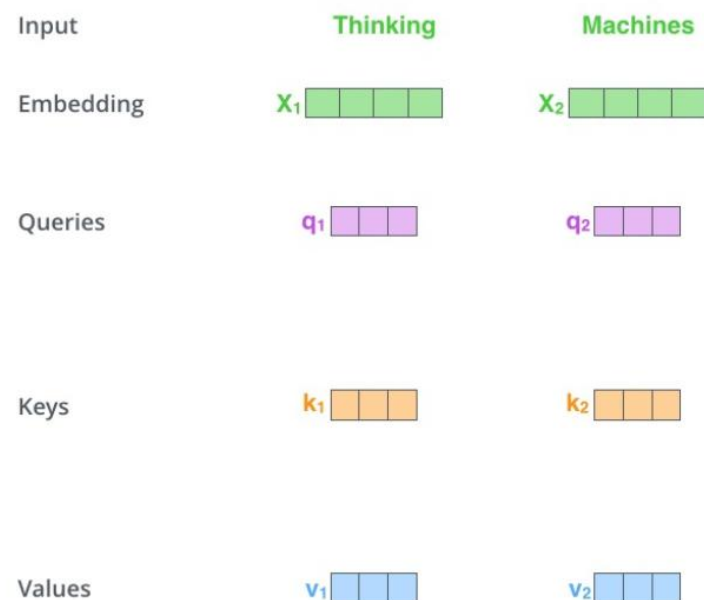## Encoders：

《Attention Is All You Need》

## Self-Attention：
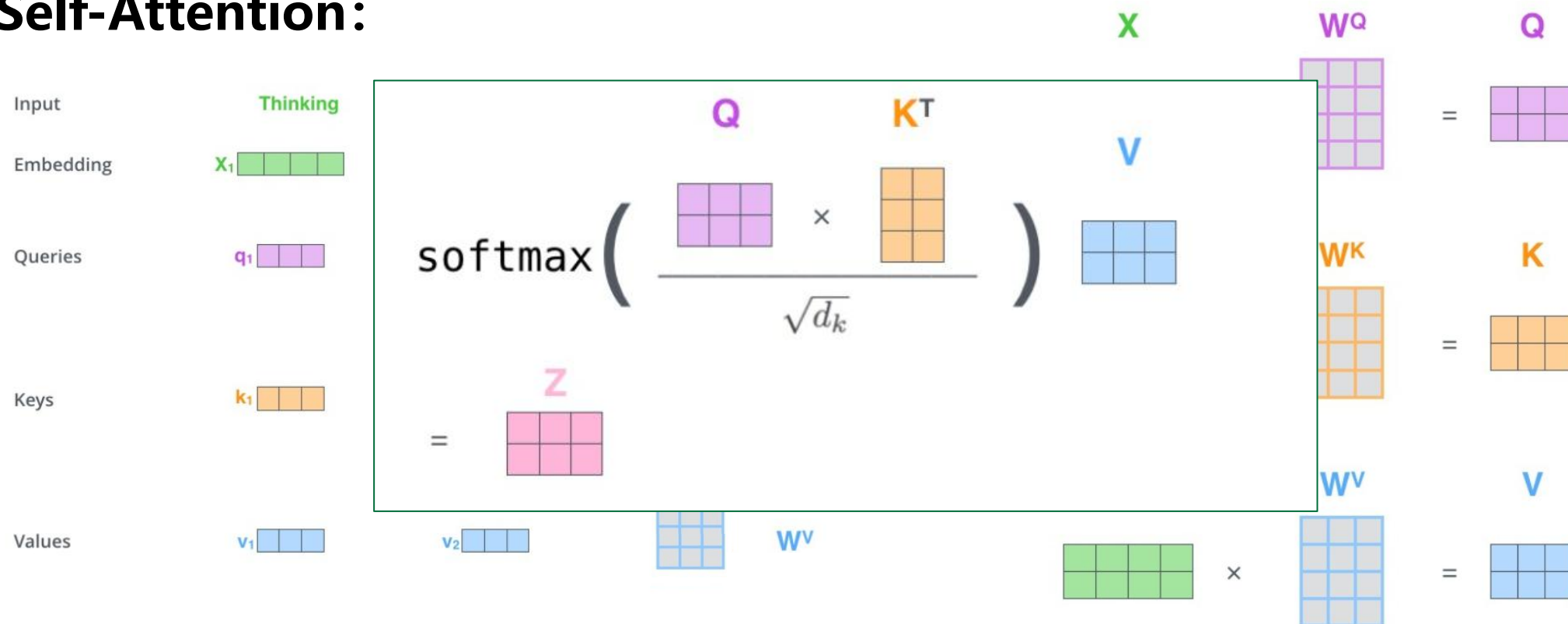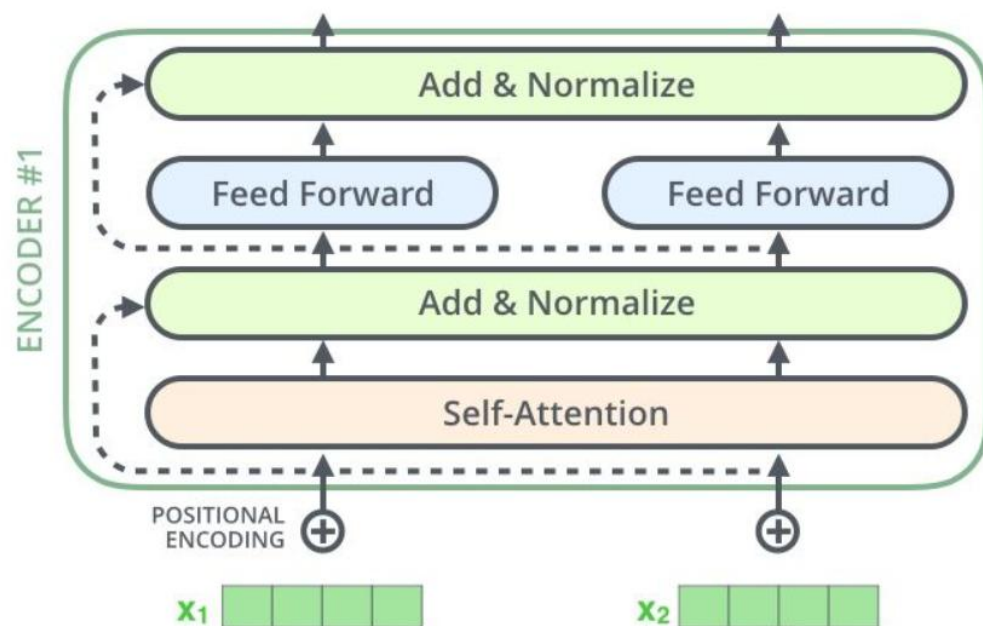
《Attention Is All You Need》

## Self-Attention：

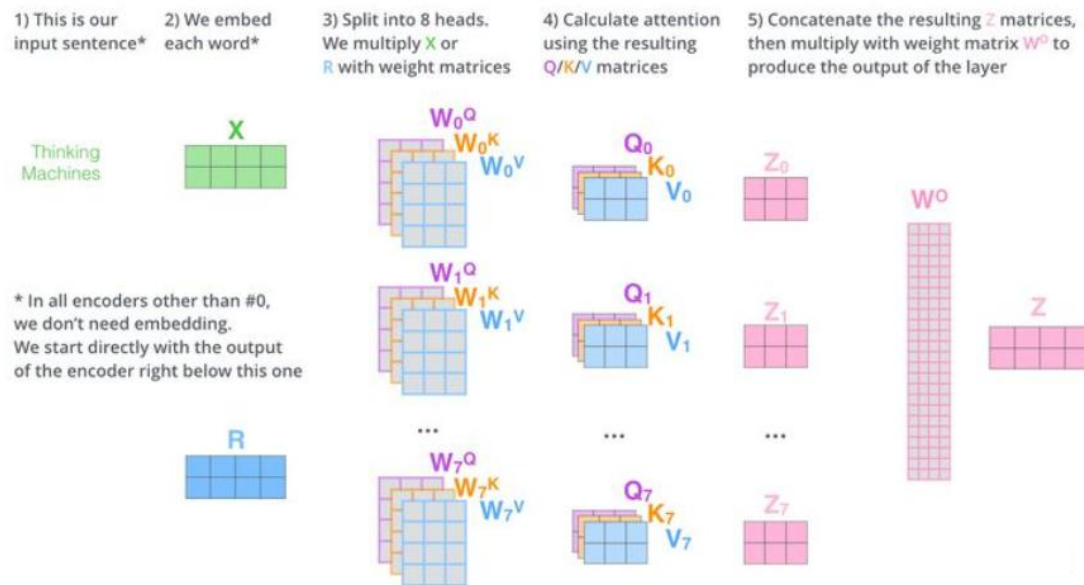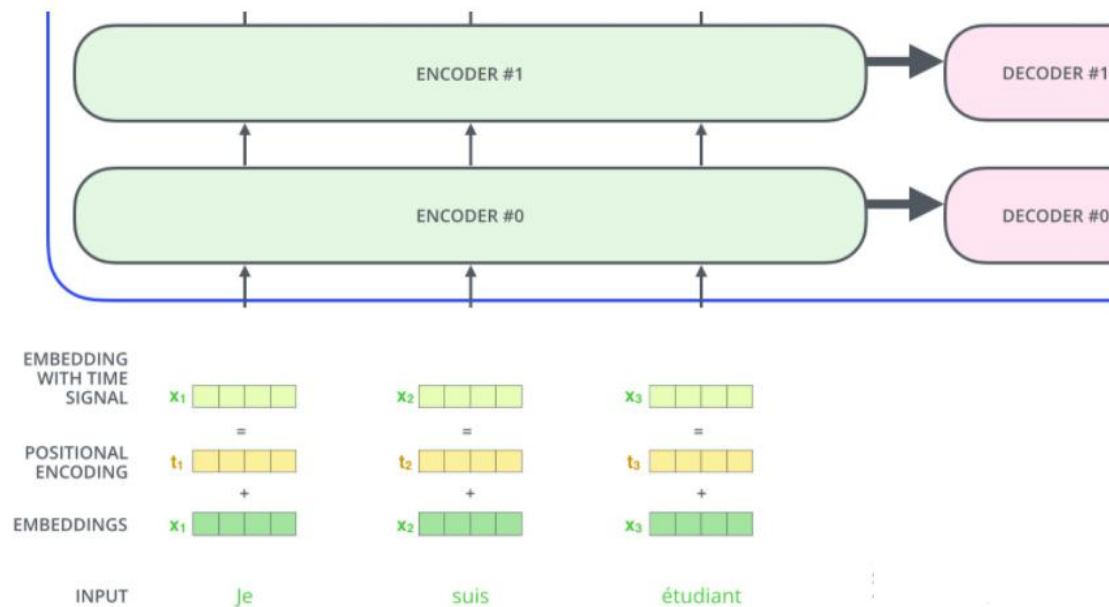《Attention Is All You Need》

## Skip Connection：



## Multi-Head：



图14：Multi-Head Attention

《Attention Is All You Need》

## Position Embedding：



$$PE(pos, 2i) = sin(\frac{pos}{10000^{\frac{2i}{d_{model}}}})$$

$$PE(pos, 2i+1) = cos(\frac{pos}{10000^{\frac{2i}{d_{model}}}})$$
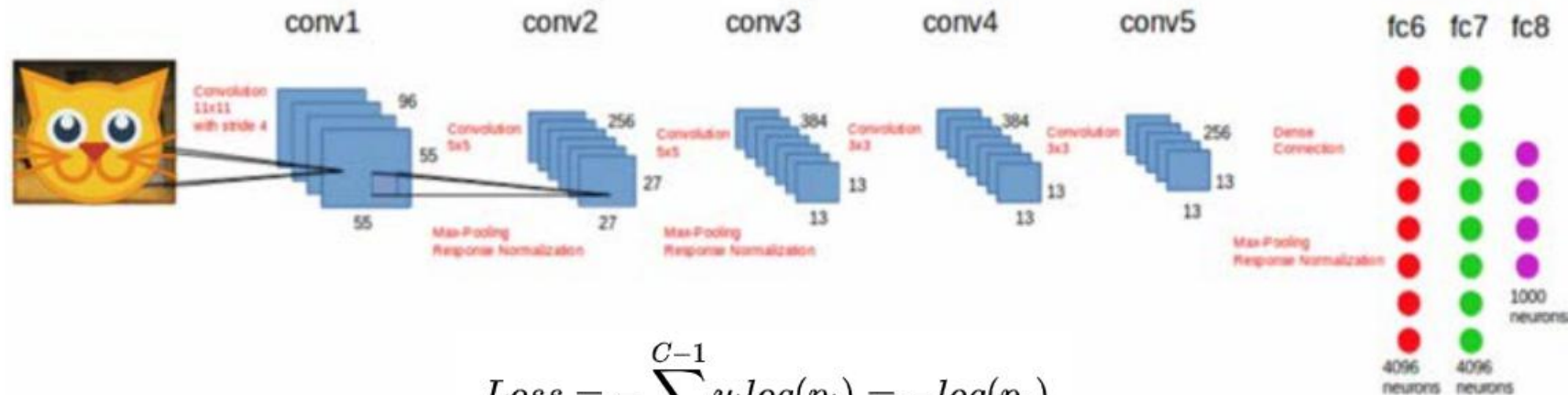
General image classification

Fine-grained image classification

The most common form of a ConvNet architecture stacks a few CONV-RELU layers, follows them with POOL layers, and repeats this pattern until the image has been merged spatially to a small size.
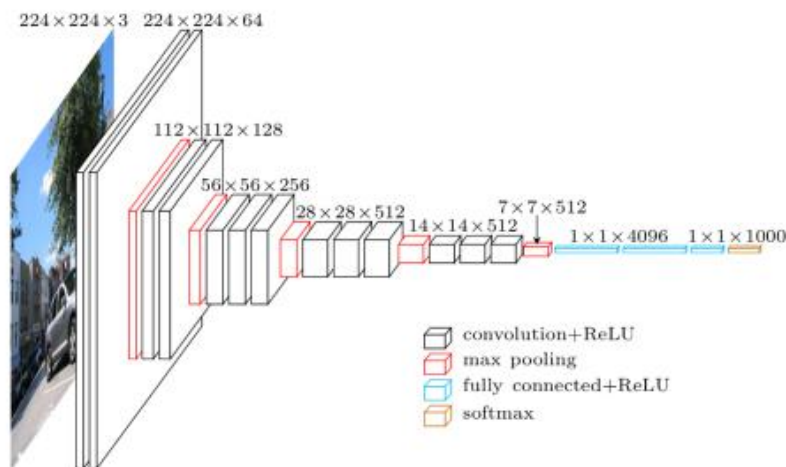
## ConvNet Architecture

$INPUT \rightarrow [[CONV \rightarrow RELU]*N \rightarrow POOL?]*M \rightarrow [FC \rightarrow RELU]*K \rightarrow FC$

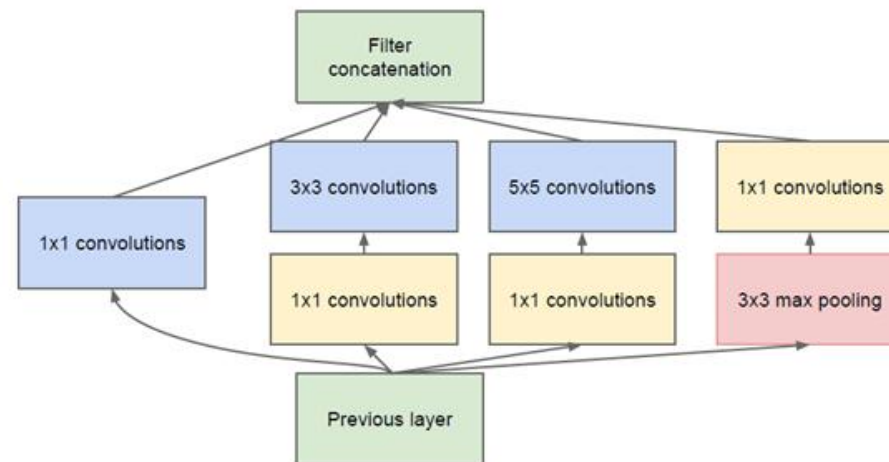where the *indicates repetition, and the $POOL?$ indicates an optional pooling layer.



$$Loss = -\sum_{i=0}^{C-1} y_i log(p_i) = -log(p_c)$$

## VGG:



224×224×3  224×224×64
112×112×128
56×56×256
28×28×512
14×14×512
7×7×512
1×1×4096  1×1×1000

- convolution+ReLU
- max pooling
- fully connected+ReLU
- softmax

## GoogleNet:



Filter concatenation

3x3 convolutions | 5x5 convolutions | 1x1 convolutions

1x1 convolutions | 1x1 convolutions | 1x1 convolutions | 3x3 max pooling

Previous layer

## ResNet:



64-d
3x3, 64
relu
3x3, 64
relu

256-d
1x1, 64
relu
3x3, 64
relu
1x1, 256
relu

《An Image is Worth 16x16 Words:Transformers for Image Recognition at Scale》

《An Image is Worth 16x16 Words:Transformers for Image Recognition at Scale》

**Problems to be solved:**

1. Positions

2. Classifications

3. Unordered set

**Faster-RCNN：**



Figure 2: Faster R-CNN is a single, unified network for object detection. The RPN module serves as the 'attention' of this unified network.
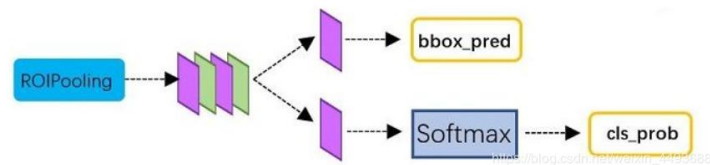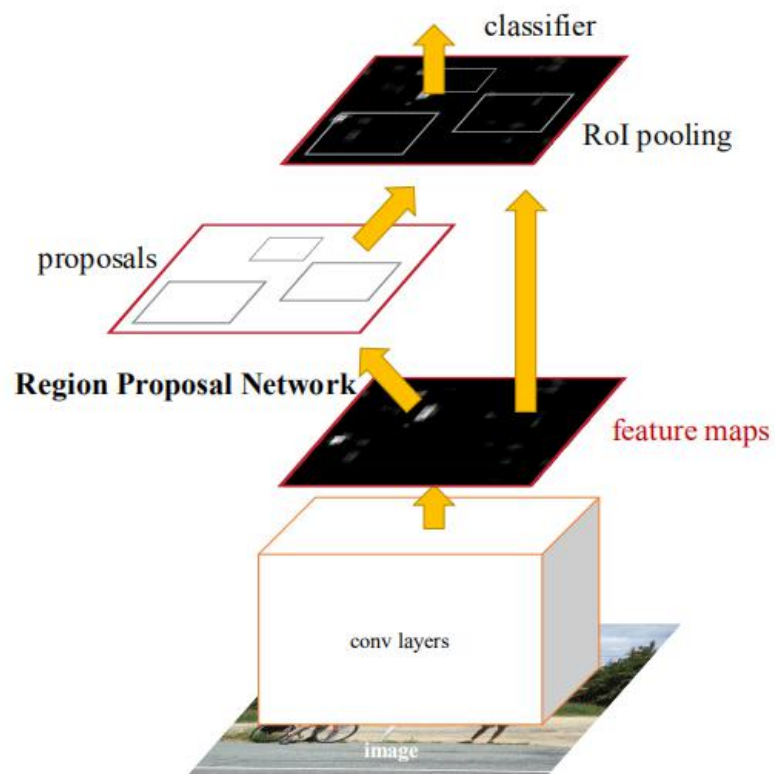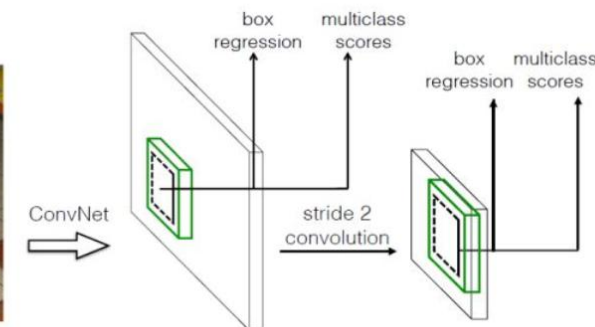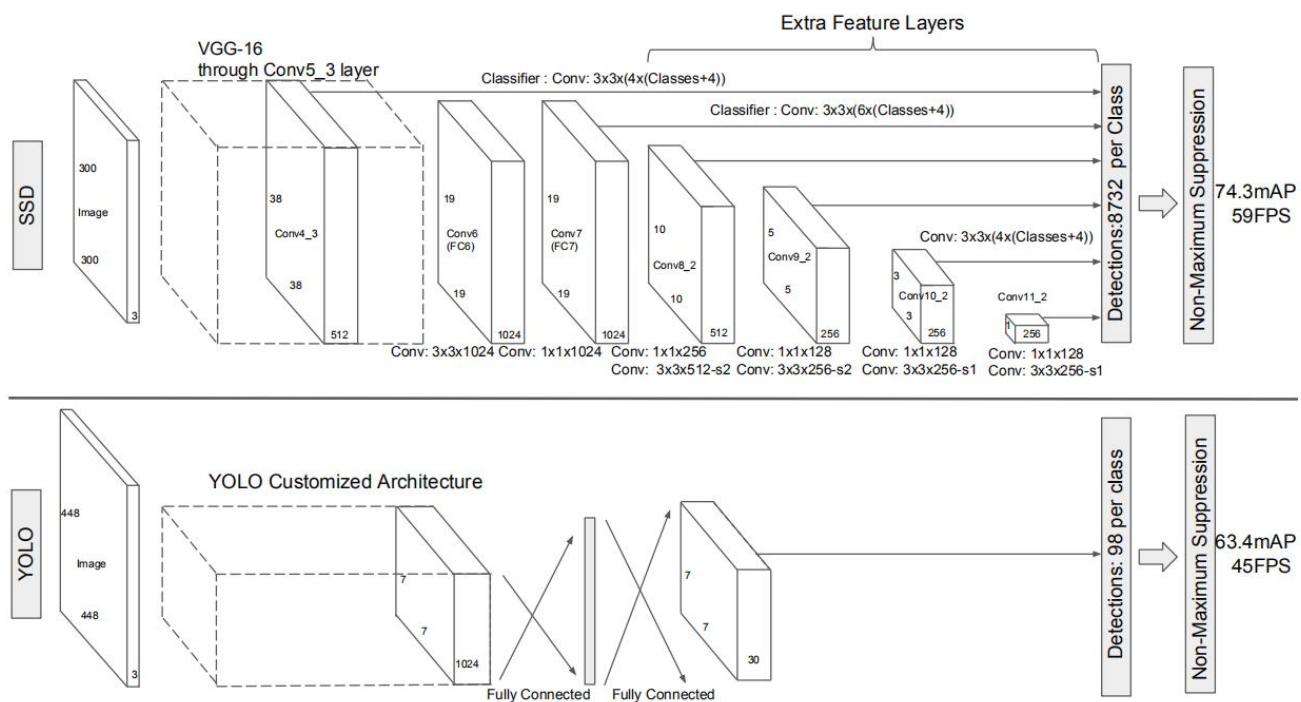
## SSD & YOLO：

## Loss Function:

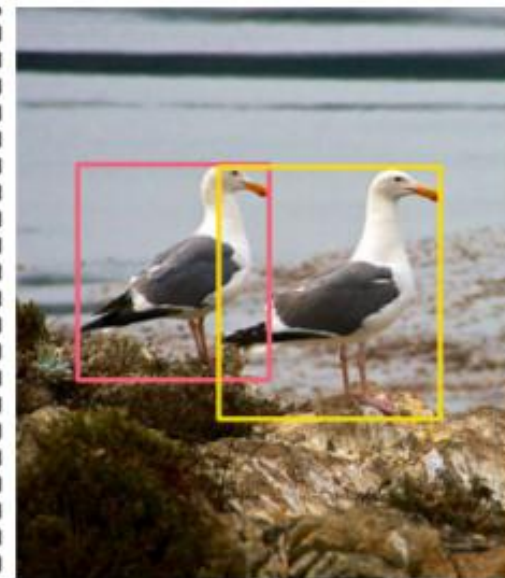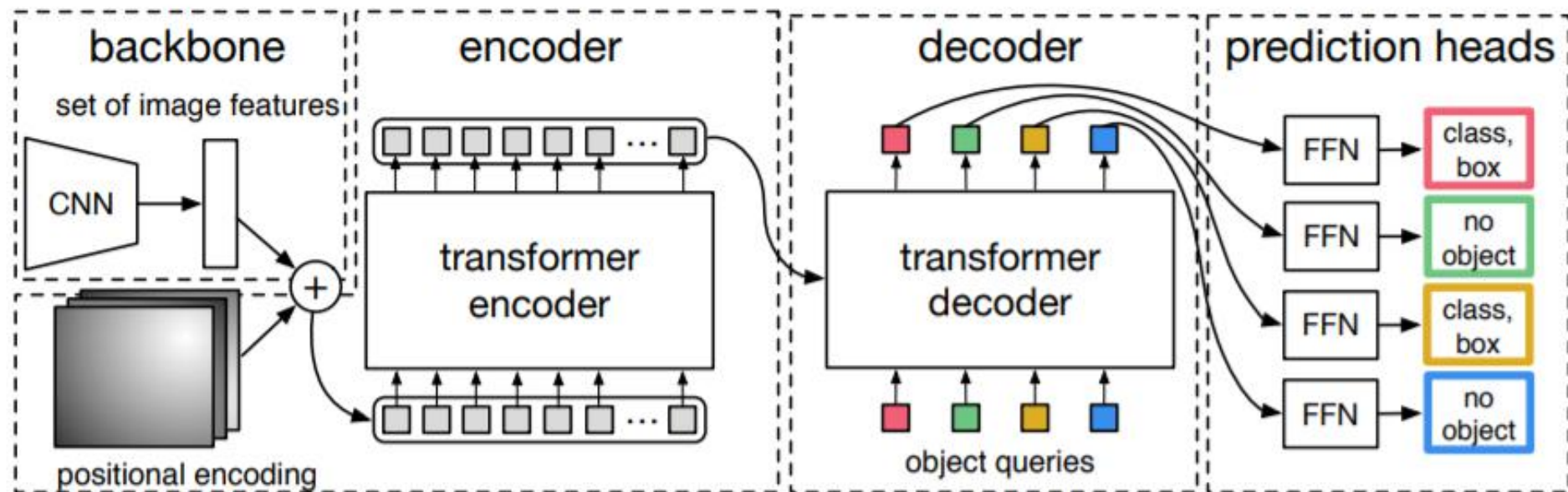$$t_\text{x} = (x - x_\text{a})/w_\text{a}, \quad t_\text{y} = (y - y_\text{a})/h_\text{a},$$
$$t_\text{w} = \log(w/w_\text{a}), \quad t_\text{h} = \log(h/h_\text{a}),$$
$$t_\text{x}^* = (x^* - x_\text{a})/w_\text{a}, \quad t_\text{y}^* = (y^* - y_\text{a})/h_\text{a},$$
$$t_\text{w}^* = \log(w^*/w_\text{a}), \quad t_\text{h}^* = \log(h^*/h_\text{a}),$$

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*)$$
$$+ \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*).$$

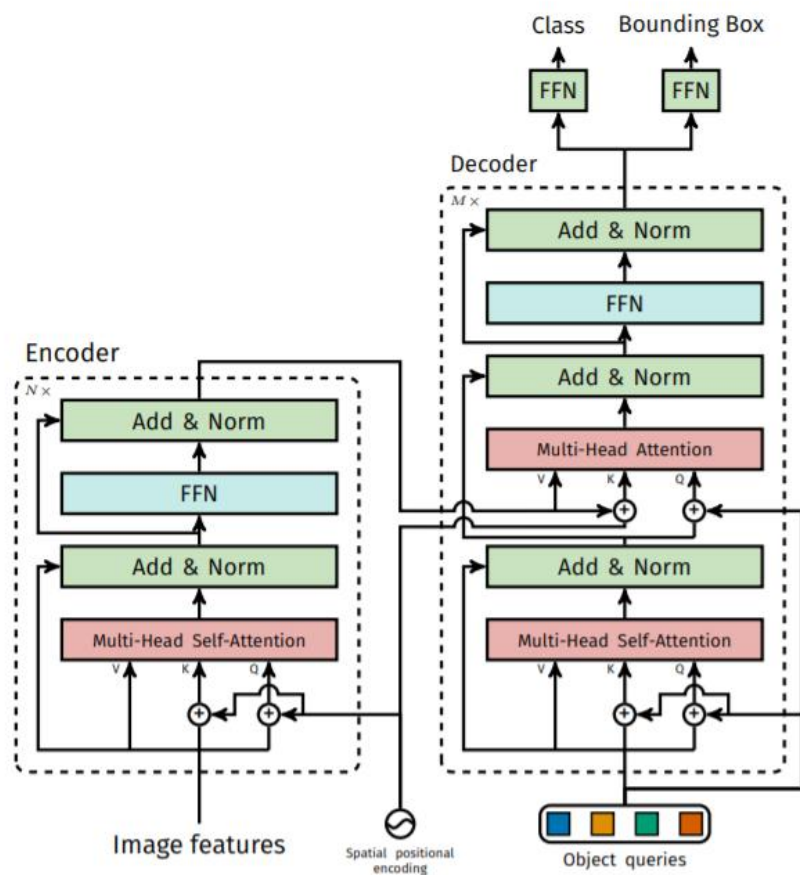# Object Detection：using transformer
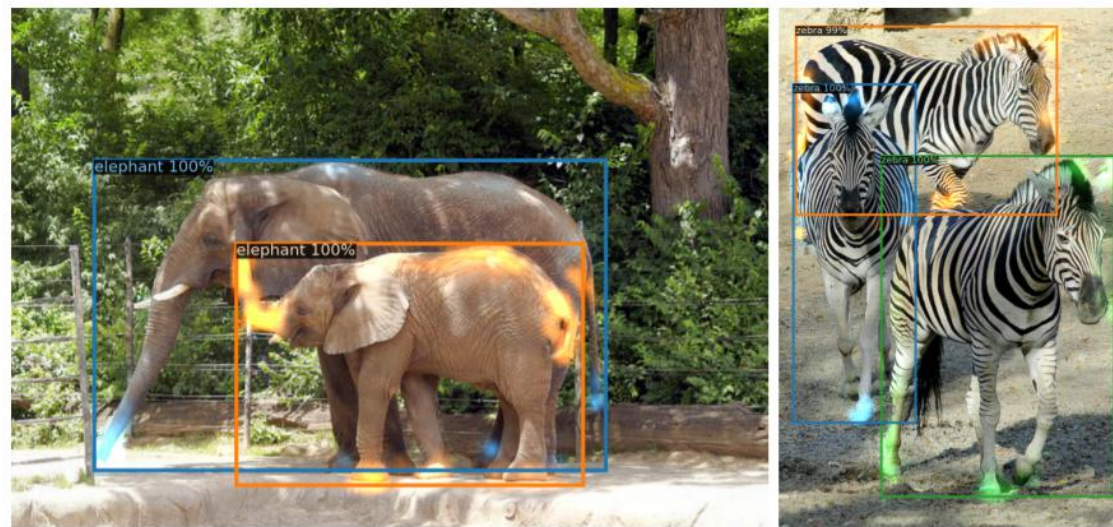
《End-to-End Object Detection with Transformers》

《End-to-End Object Detection with Transformers》



$$\hat{\sigma} = \arg\min \sum_{}^{N} \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)}),$$

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

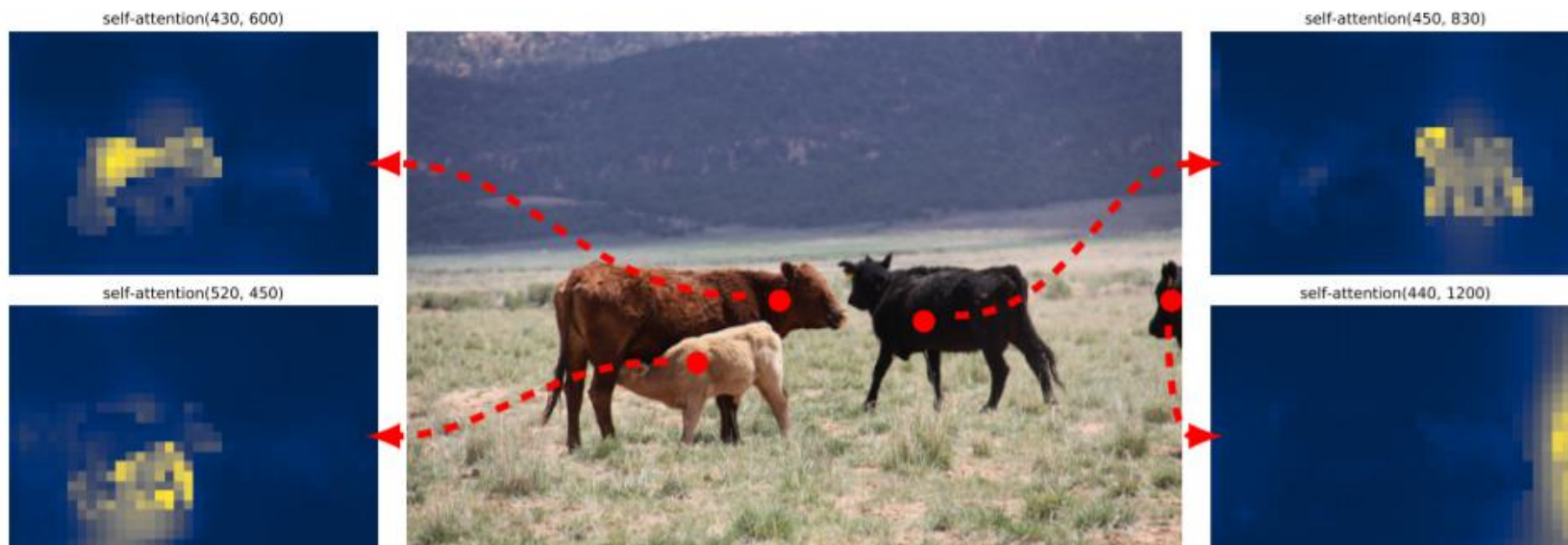《End-to-End Object Detection with Transformers》



Fig. 3: Encoder self-attention for a set of reference points. The encoder is able to separate individual instances. Predictions are made with baseline DETR model on a validation set image.

Person
Bicycle
Background

**FCN：**



Image

Convolved
Feature

《Rethinking Semantic Segmentation from a Sequence-to-Sequence Perspective with Transformers》

《Rethinking Semantic Segmentation from a Sequence-to-Sequence Perspective with Transformers》

$$\text{query} = Z^{l-1}\mathbf{W}_Q, \ \text{key} = Z^{l-1}\mathbf{W}_K, \ \text{value} = Z^{l-1}\mathbf{W}_V,$$

$$SA(Z^{l-1}) = Z^{l-1} + \text{softmax}\left(\frac{Z^{l-1}\mathbf{W}_Q(Z\mathbf{W}_K)^\top}{\sqrt{d}}\right)(Z^{l-1}\mathbf{W}_V).$$

and project their concatenated outputs: $MSA(Z^{l-1}) = [SA_1(Z^{l-1}); SA_2(Z^{l-1}); \cdots ; SA_m(Z^{l-1})]\mathbf{W}_O$, where $\mathbf{W}_O \in \mathbb{R}^{md \times C}$. $d$ is typically set to $C/m$. The output of

$$Z^l = MSA(Z^{l-1}) + MLP(MSA(Z^{l-1})) \in \mathbb{R}^{L \times C}.$$
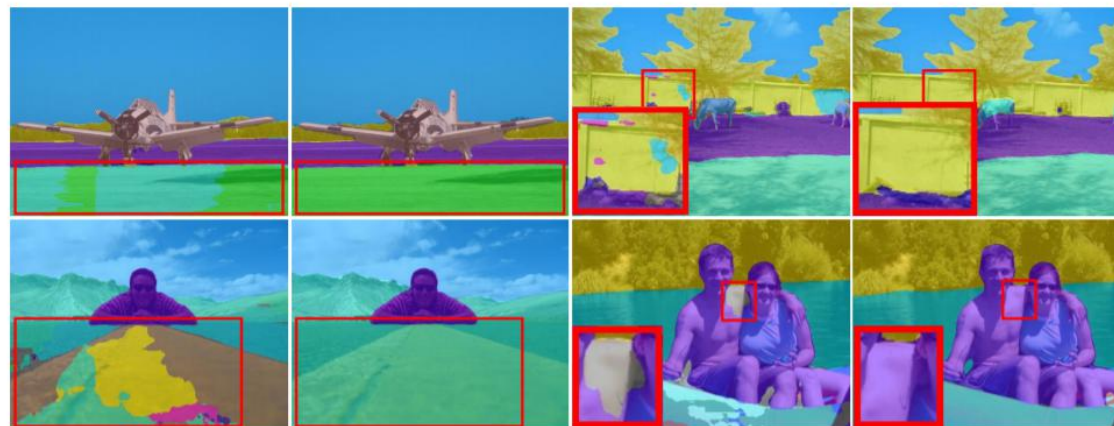
$$\{Z^1, Z^2, \cdots, Z^{L_e}\}$$



Figure 3. **Qualitative results on Pascal Context:** SETR (right column) vs. dilated FCN baseline (left column) in each pair. Best viewed in color and zoom in.