

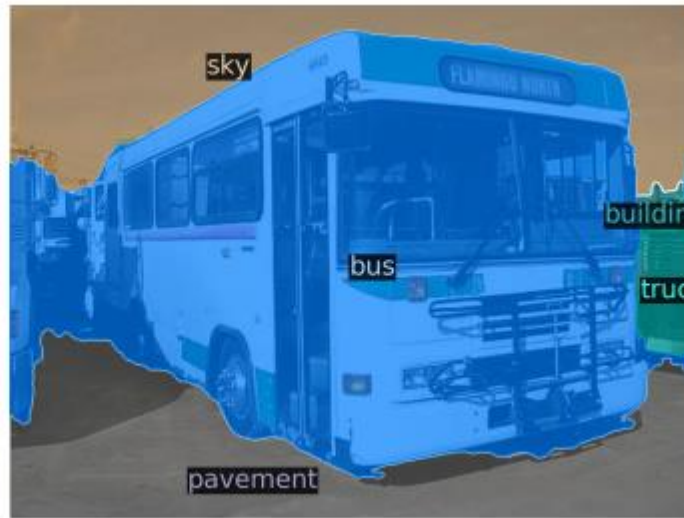
DEPT

# Paper Details

- Paper Title: End-to-End Object Detection with Transformers
- Publication Date: 28 May 2020
- Publisher: Nicolas Carion , Francisco Massa , ...
- Affiliation: Facebook AI

# Objectives

- end-to-end object detection

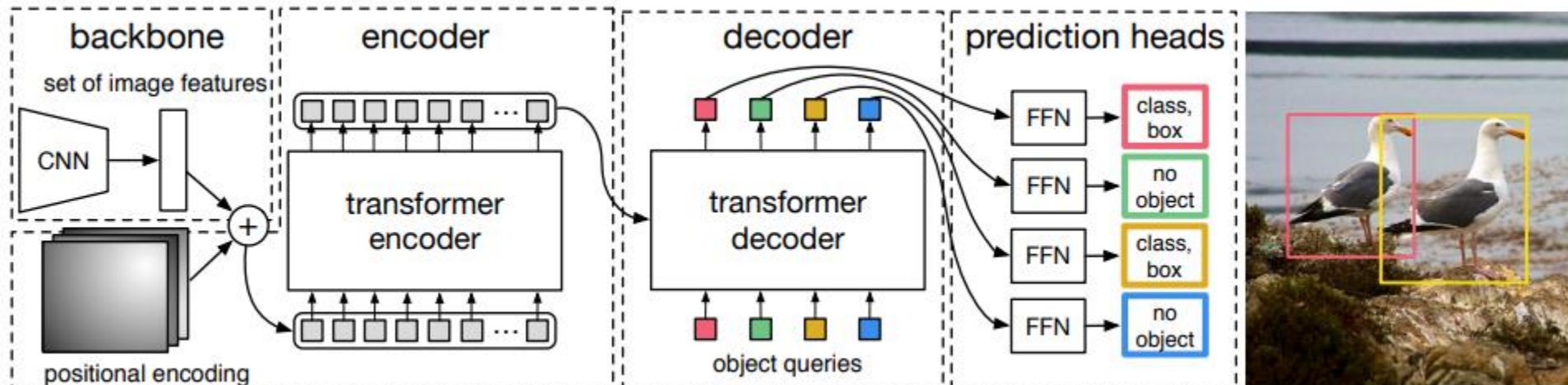


# Motivation

- detectors do set prediction task in an indirect way
  - simplify the pipeline
- Transformer in other fields
  - Bridge the gap

# Overall Architecture

- Backbone
- Transformer Encoder
- Transformer Decoder
- Prediction feed-forward network

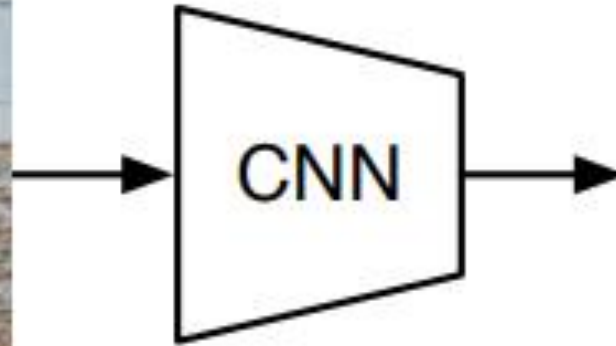


# Overall Architecture

- Backbone



$$x_{\text{img}} \in \mathbb{R}^{3 \times H_0 \times W_0}$$



$$\mathbb{R}^{C \times H \times W}$$

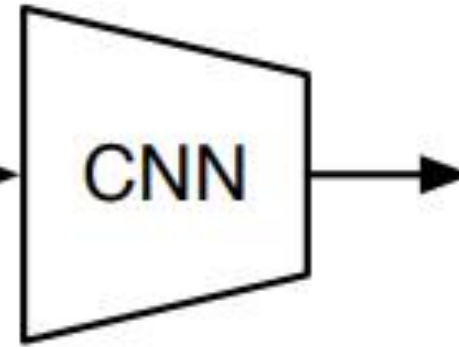
$$C = 2048 \text{ and } H, W = \frac{H_0}{32}, \frac{W_0}{32}.$$

# Overall Architecture

- Backbone

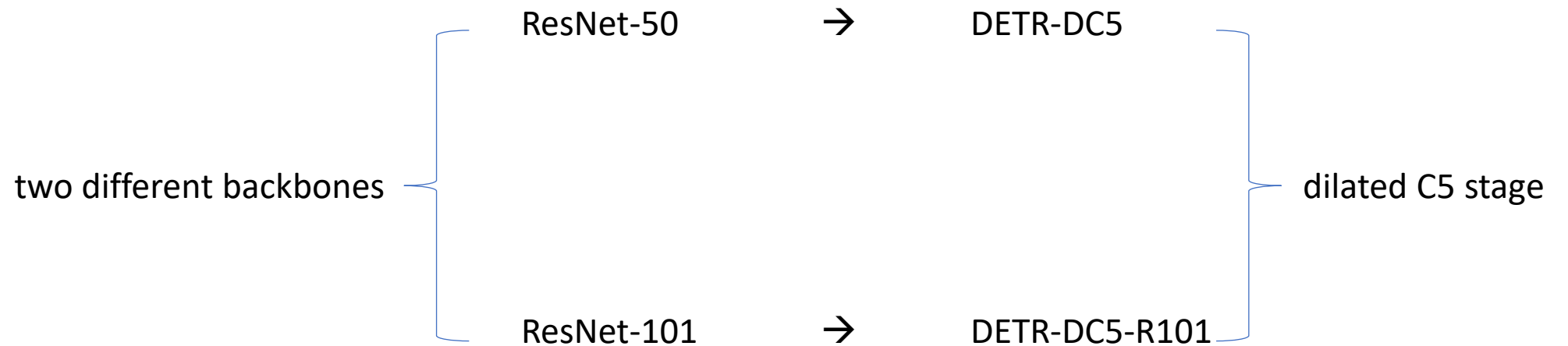


$$x_{\text{img}} \in \mathbb{R}^{3 \times H_0 \times W_0}$$



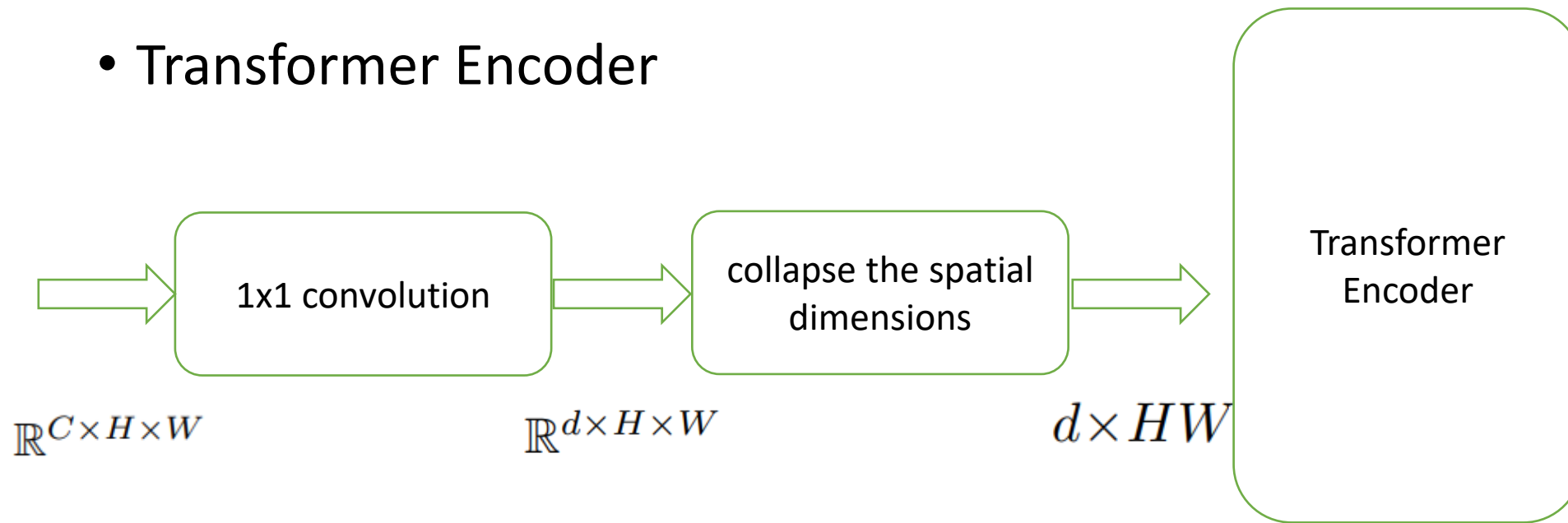
$$\mathbb{R}^{C \times H \times W}$$

ImageNet-pretrained ResNet model



# Overall Architecture

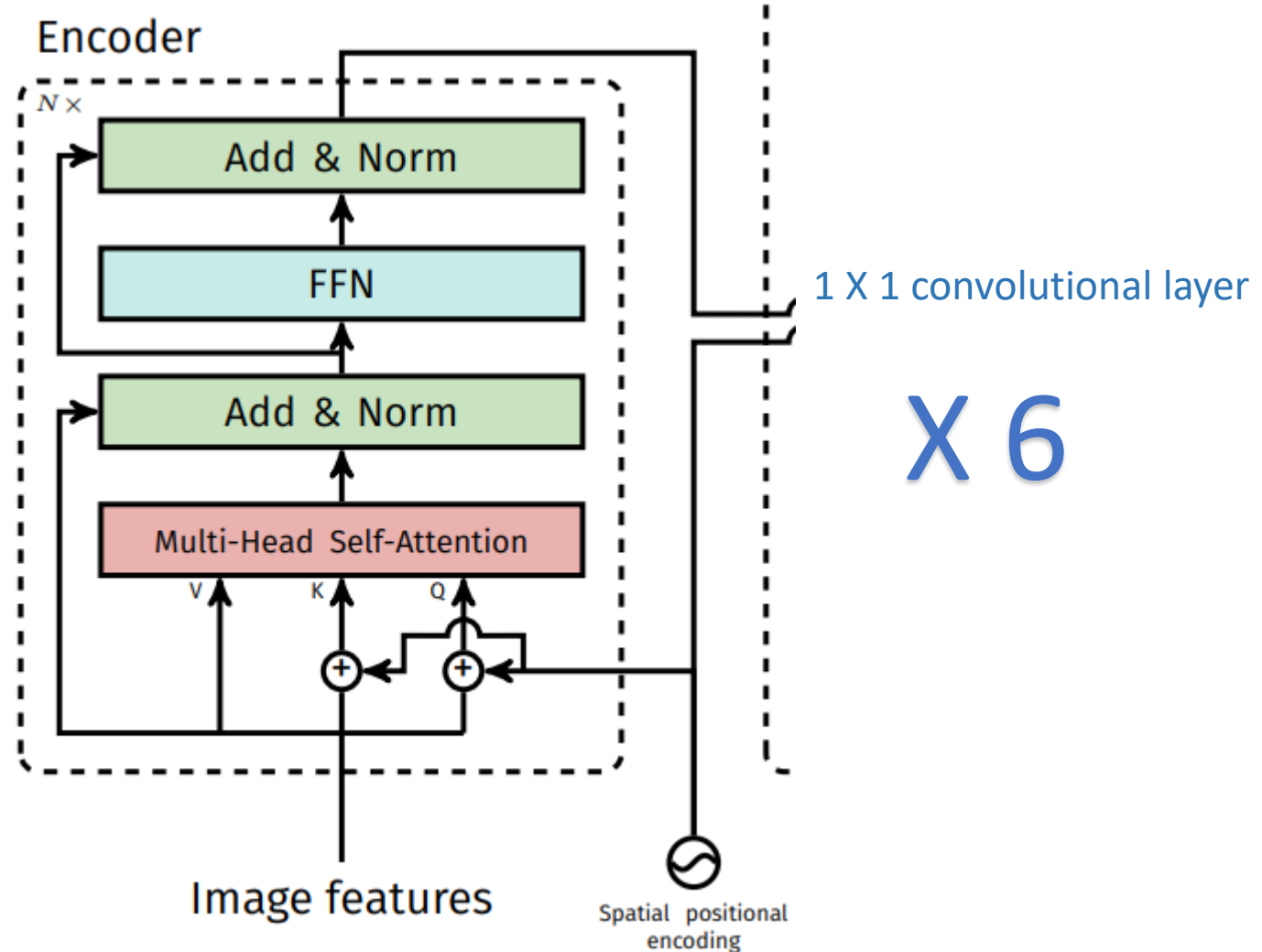
- Backbone
- Transformer Encoder





# Overall Architecture

- Backbone
- Transformer Encoder



# Overall Architecture

- Backbone

$$\text{Self-Attention}(\mathbf{X})_{t,:} := \text{softmax}(\mathbf{A}_{t,:}) \mathbf{X} \mathbf{W}_{val},$$

- Transformer Encoder
  - positional encodings

$$\mathbf{A} := \mathbf{X} \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} \mathbf{X}^{\top}$$

$$\mathbf{A} := (\mathbf{X} + \mathbf{P}) \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} (\mathbf{X} + \mathbf{P})^{\top}$$

# Overall Architecture

- Backbone

$$\text{Self-Attention}(\mathbf{X})_{t,:} := \text{softmax}(\mathbf{A}_{t,:}) \mathbf{X} \mathbf{W}_{val},$$

- Transformer Encoder
  - positional encodings

$$\mathbf{A} := \mathbf{X} \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} \mathbf{X}^{\top}$$

$$\mathbf{A} := (\mathbf{X} + \mathbf{P}) \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} (\mathbf{X} + \mathbf{P})^{\top}$$

Absolute  
Encoding

$$\mathbf{A}_{q,k}^{\text{abs}} = (\mathbf{X}_{q,:} + \mathbf{P}_{q,:}) \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} (\mathbf{X}_{k,:} + \mathbf{P}_{k,:})^{\top}$$

Relative  
Encoding

$$\mathbf{A}_{q,k}^{\text{rel}} := \mathbf{X}_{q,:}^{\top} \mathbf{W}_{qry}^{\top} \mathbf{W}_{key} \mathbf{X}_{k,:} + \mathbf{X}_{q,:}^{\top} \mathbf{W}_{qry}^{\top} \widehat{\mathbf{W}}_{key} \mathbf{r}_{\delta} + \mathbf{u}^{\top} \mathbf{W}_{key} \mathbf{X}_{k,:} + \mathbf{v}^{\top} \widehat{\mathbf{W}}_{key} \mathbf{r}_{\delta}$$

# Overall Architecture

- Backbone
- Transformer Encoder
  - positional encodings

$$\text{Self-Attention}(\mathbf{X})_{t,:} := \text{softmax}(\mathbf{A}_{t,:}) \mathbf{X} \mathbf{W}_{val},$$

$$\mathbf{A} := \mathbf{X} \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} \mathbf{X}^{\top}$$

$$\mathbf{A} := (\mathbf{X} + \mathbf{P}) \mathbf{W}_{qry} \mathbf{W}_{key}^{\top} (\mathbf{X} + \mathbf{P})^{\top}$$

Absolute  
Encoding

| spatial pos. enc. |                  | output pos. enc.<br>decoder | AP          |          | AP <sub>50</sub> |          |
|-------------------|------------------|-----------------------------|-------------|----------|------------------|----------|
| encoder           | decoder          |                             |             | $\Delta$ |                  | $\Delta$ |
| none              | none             | learned at input            | 32.8        | -7.8     | 55.2             | -6.5     |
| sine at input     | sine at input    | learned at input            | 39.2        | -1.4     | 60.0             | -1.6     |
| learned at attn.  | learned at attn. | learned at attn.            | 39.6        | -1.0     | 60.7             | -0.9     |
| none              | sine at attn.    | learned at attn.            | 39.3        | -1.3     | 60.3             | -1.4     |
| sine at attn.     | sine at attn.    | learned at attn.            | <b>40.6</b> | -        | <b>61.6</b>      | -        |

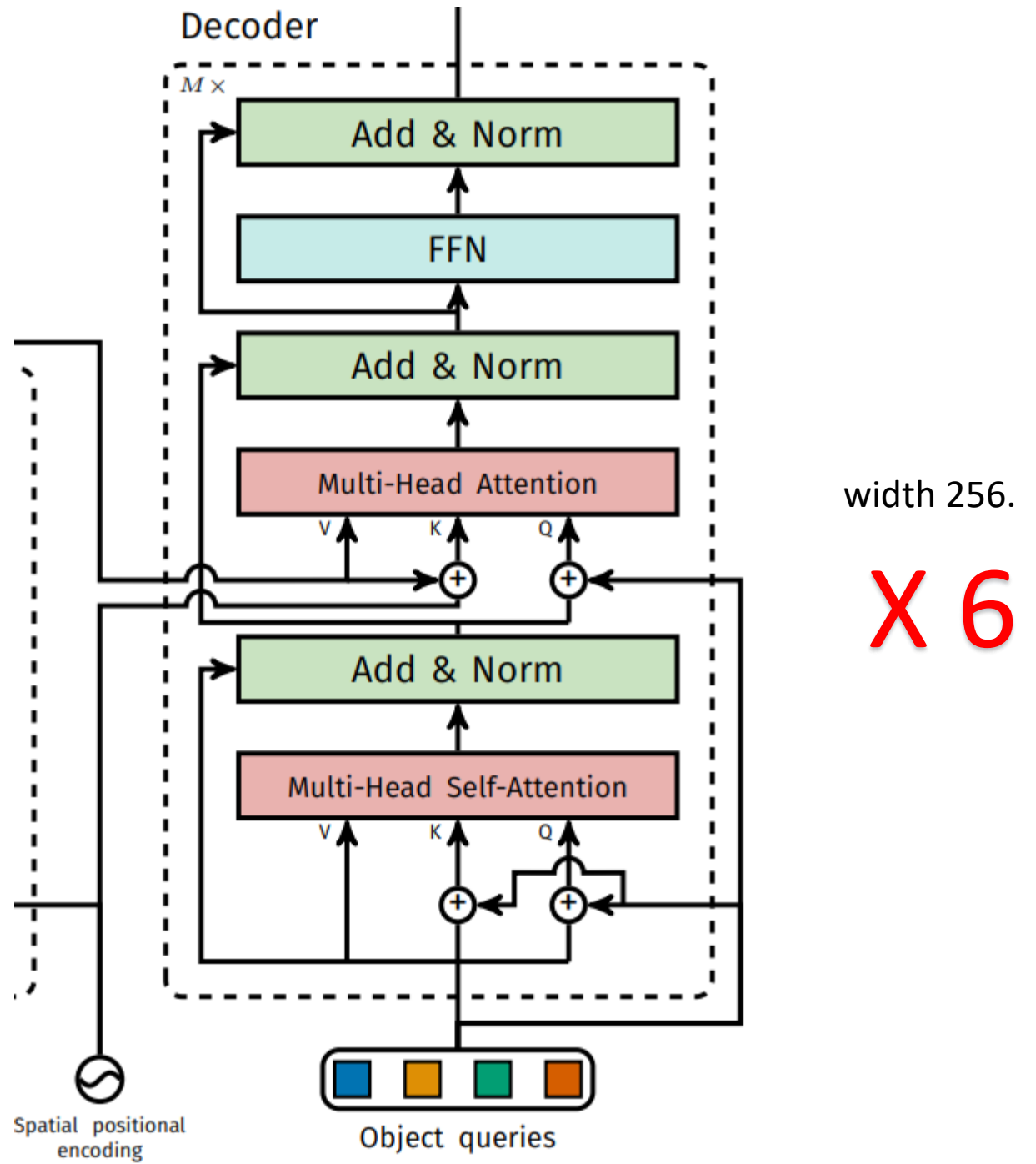
Relative  
Encoding

## Attention Augmented Convolutional Networks

| Position Encodings | mAP <sub>COCO</sub> | mAP <sub>50</sub> | mAP <sub>75</sub> |
|--------------------|---------------------|-------------------|-------------------|
| None               | 37.7                | 56.0              | 40.2              |
| CoordConv [29]     | 37.4                | 55.5              | 40.1              |
| Relative (ours)    | 38.2                | 56.5              | 40.7              |

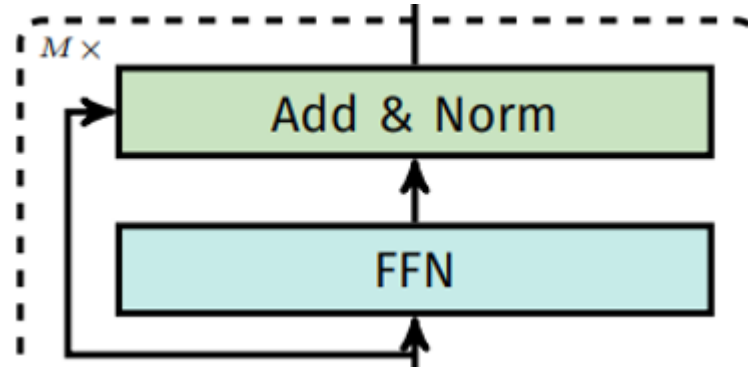
# Overall Architecture

- Backbone
- Transformer Encoder
- Transformer Decoder



# Overall Architecture

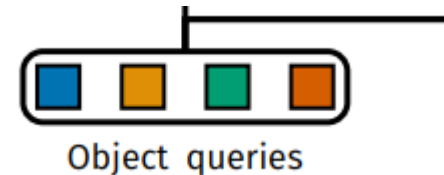
- Backbone
- Transformer Encoder
- Transformer Decoder



N input embeddings

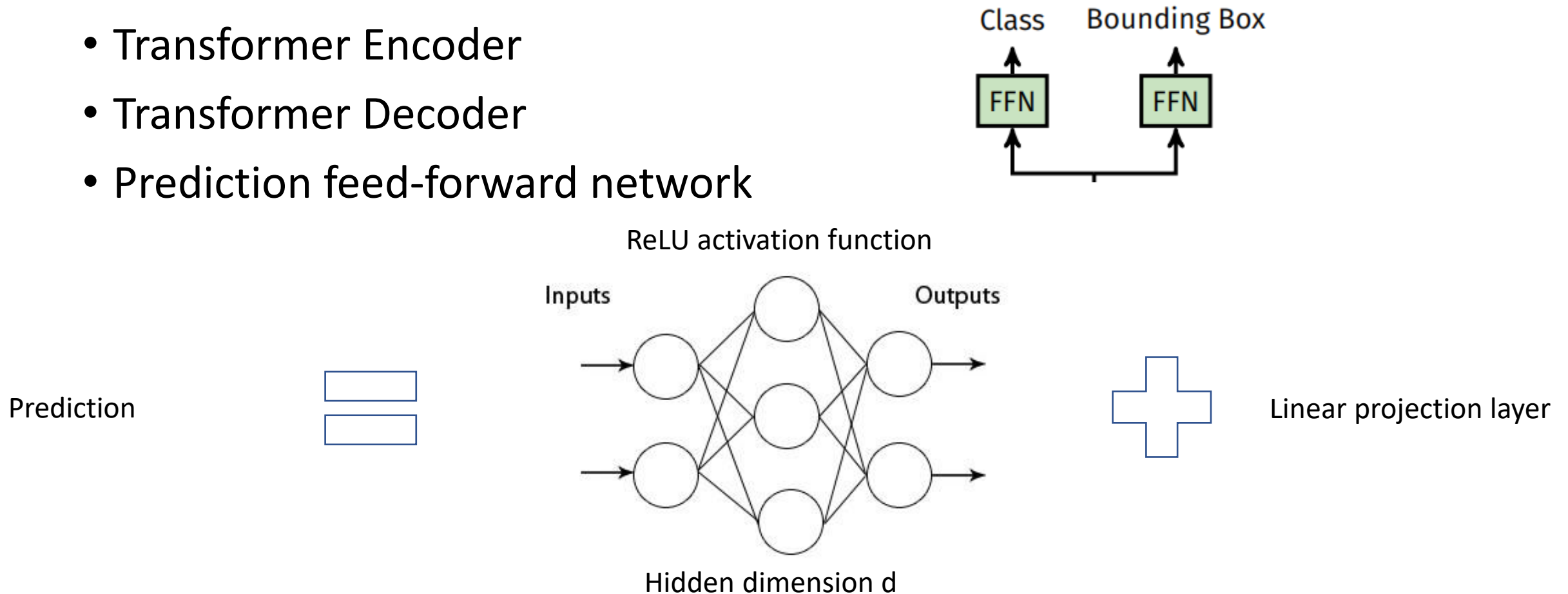
different

learnt positional encodings

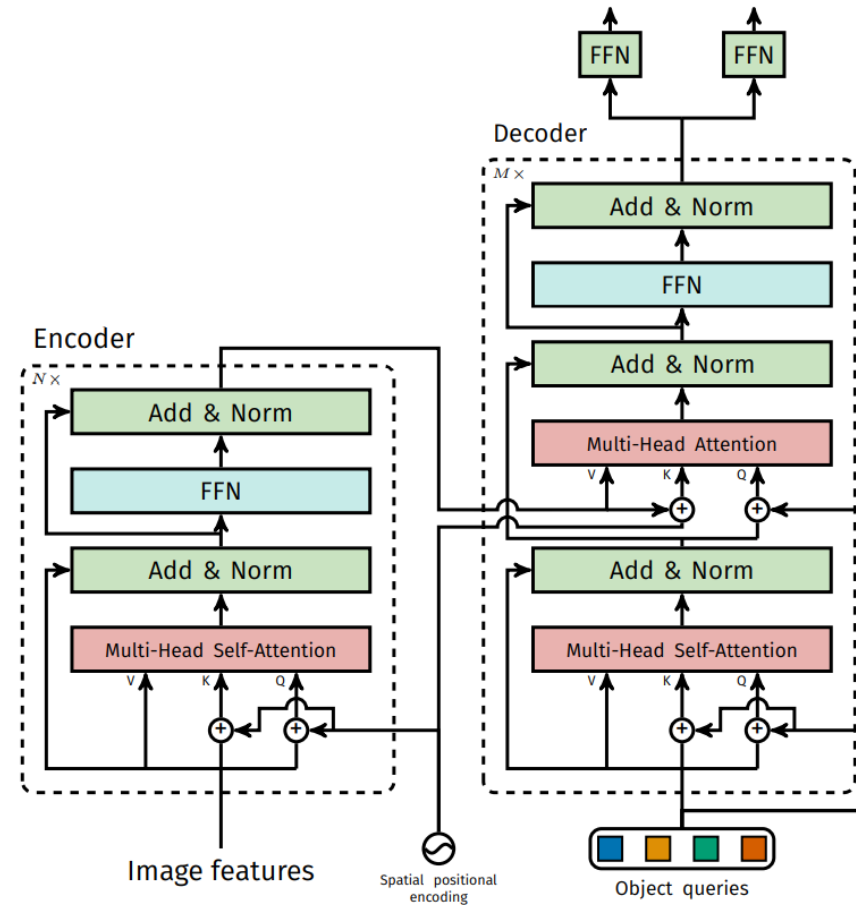


# Overall Architecture

- Backbone
- Transformer Encoder
- Transformer Decoder
- Prediction feed-forward network



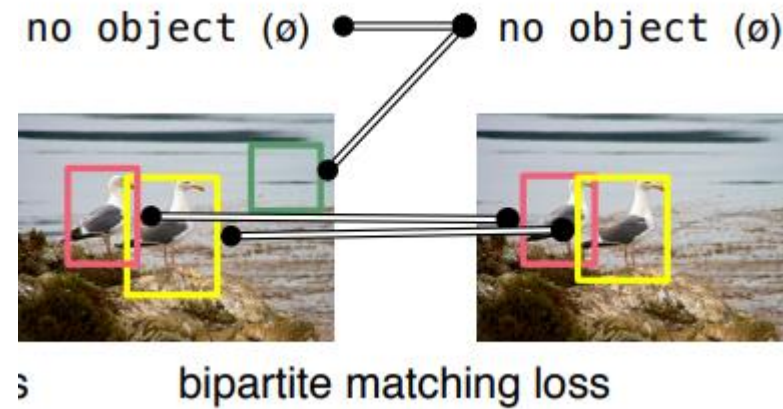
# Overall Architecture





# Prediction Loss

- Bipartite matching



$$\hat{\sigma} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_i^N \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)})$$

# Prediction Loss

- Bipartite matching
  - Hungarian Algorithm

$$\hat{\sigma} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_i^N \mathcal{L}_{\text{match}}(y_i, \hat{y}_{\sigma(i)})$$

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

Trick : we down-weight the log-probability term when  $c_i = \emptyset$  by a factor 10 to account for class imbalance

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

Classification loss

l1 bounding box distance loss

GloU loss

# Experiments

- Dataset
- Training

# Experiments

- Dataset
  - COCO 2017
    - Average 7 instance in an image (maximal 63 instances)
    - AP refers to bbox AP

# Experiments

- Dataset
- Training
  - Whole Architecture
  - Data Process
  - Backbone
  - Transformer
  - Losses

# Experiments

- Dataset
- Training
  - Whole Architecture
    - AdamW with improved weight decay  $10^{-4}$ , maximum gradient norm

# Experiments

- Dataset
- Training
  - Whole Architecture
  - Data Process
    - Scale augmentation resize images

# Experiments

- Dataset
- Training
  - Whole Architecture
  - Data Process
  - Backbone
    - Backbone batch normalization weights and statistics are frozen during training
    - Learning rate  $10^{-5}$



# Experiments

- Dataset
- Training
  - Whole Architecture
  - Data Process
  - Backbone
  - Transformer
    - Learning rate  $10^{-4}$
    - Dropout 0.1

# Experiments

- Dataset
- Training
  - Whole Architecture
  - Data Process
  - Backbone
  - Transformer
  - Losses

$$\lambda_{L1} = 5 \text{ and } \lambda_{iou} = 2$$

# Panoptic segmentations

- Predicting Box
- Mask Head

# Panoptic segmentations

- Predicting Box
- Mask Head

