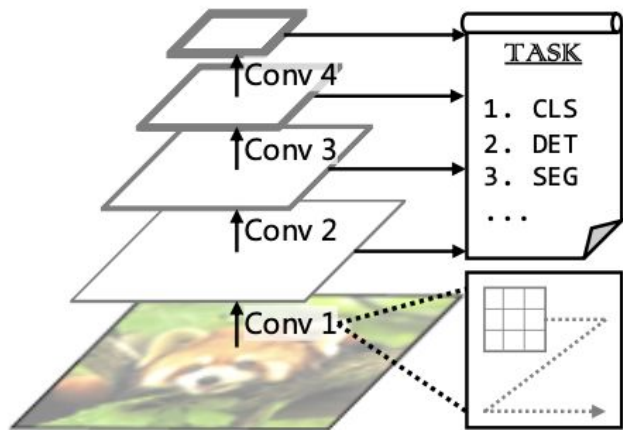# Bottleneck Transformers for Visual Recognition

## *A Self-attention model for vision*

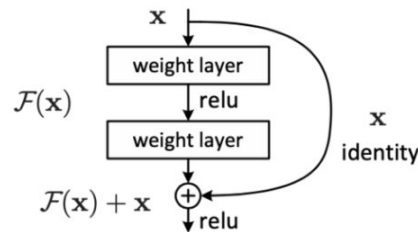Will, Tin, Tarunyaa

# Premise: Out with CNNs

**Better alternatives to pure CNNs for CV such as object detection, image classification**



(a) CNNs: VGG [53], ResNet [21], *etc.*

➜ **ResNet: Residual Neural Network**

Intermediate input added to the output of a aries of convolutional blocks to enable scaling.



They capture local patterns information but they **fail to understand long-term dependencies + require many layers.**

# Premise: Out with CNNs

➜ Longer-term memory than RNNs & LSTMs

➜ NLP to Image processing

➜ GPT, BERT

**They capture long term dependencies and don't require as many layers.**

The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
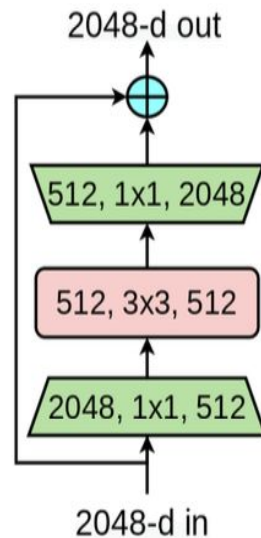The FBI is chasing a criminal on the run .
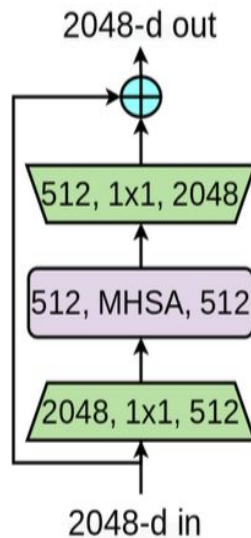
# Proposal: Using Self-Attention in Vision

**HYBRID SOLUTION: Replace spatial convolutional layers w/ multi-head self attention layer**

➔ Use convolutions to deal with large images efficiently; do spatial downsampling

➔ Letting global self-attention work on small resolutions

**Avoids processing large images w/ self-attention since its memory & computation required scales quadratically w/ spatial dimensions**

2048-d out

512, 1x1, 2048

512, 3x3, 512

2048, 1x1, 512

2048-d in

**ResNet Bottleneck**

2048-d out

512, 1x1, 2048

512, MHSA, 512

2048, 1x1, 512

2048-d in
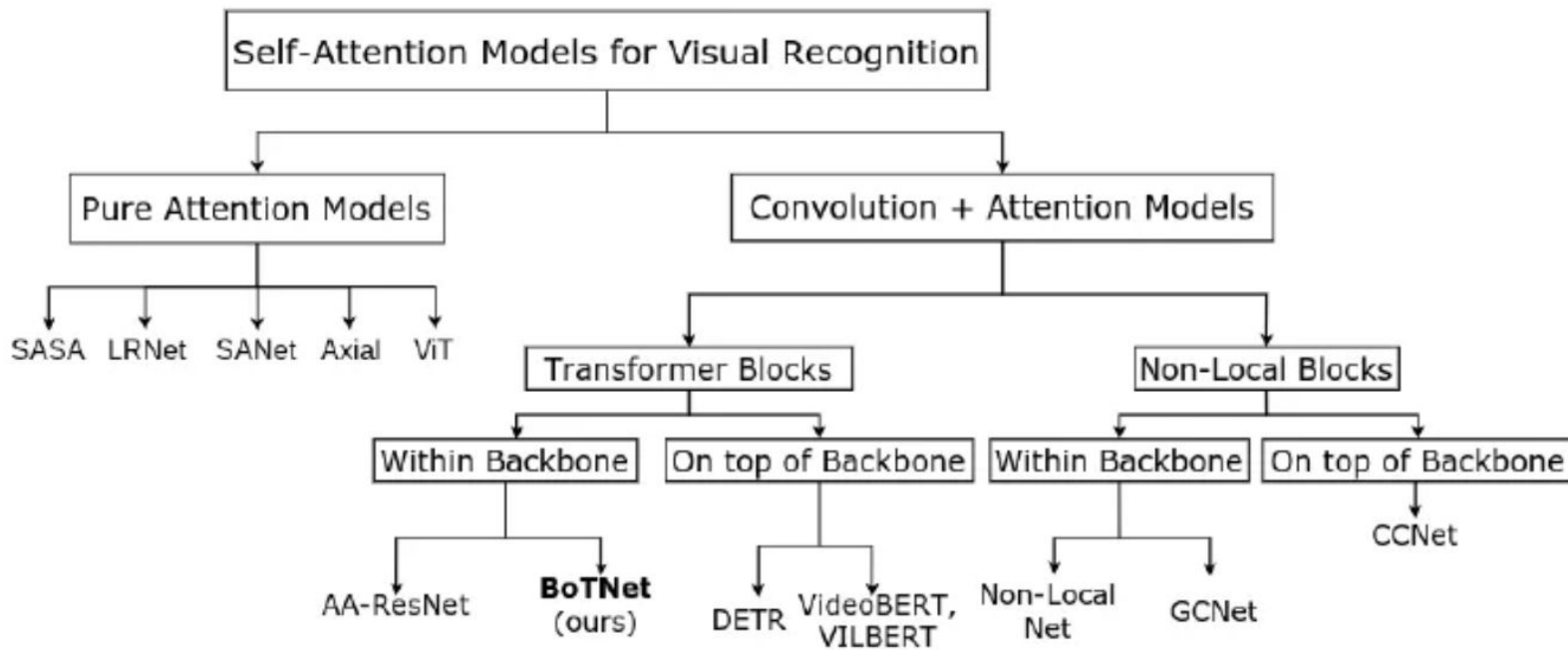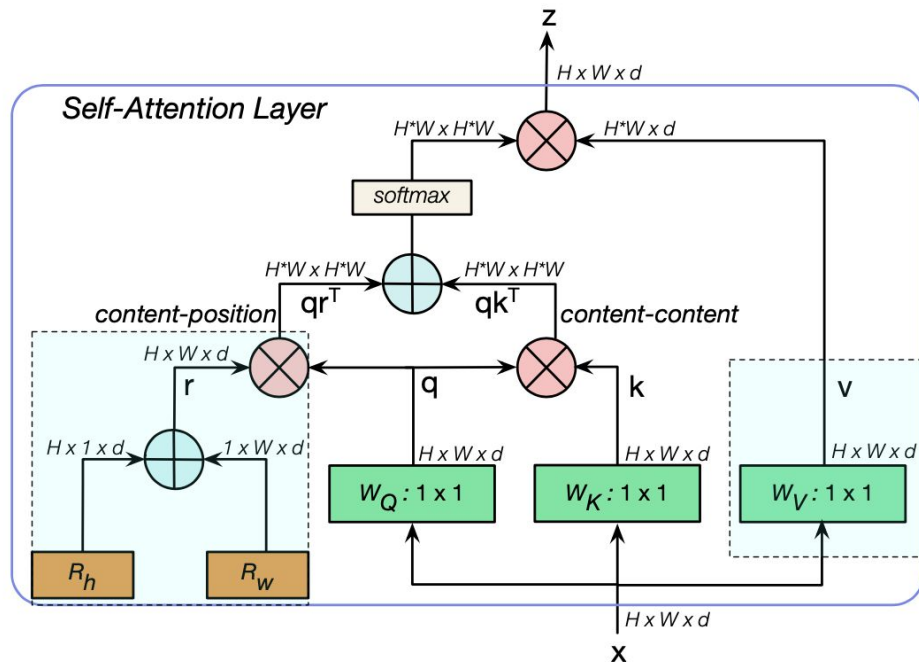
**Bottleneck Transformer**

# Proposal: Using Self-Attention in Vision

**HYBRID SOLUTION: Replace spatial convolutional layers w/ multi-head self attention layer**

# Method: Positional Encoding

**Making the attention-operation *position aware***



- ➤ Global attention is performed on a 2D feature map

- ➤ Split relative position encodings, *Rh* and *Rw*, for height and width respectively

# Method: Model Architecture

**Low relative overhead**

| stage | output | ResNet-50 | BoTNet-50 |
|---|---|---|---|
| c1 | $512 \times 512$ | $7 \times 7$, 64, stride 2 | $7 \times 7$, 64, stride 2 |
| c2 | $256 \times 256$ | $3 \times 3$ max pool, stride 2 <br> $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$ | $3 \times 3$ max pool, stride 2 <br> $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$ |
| c3 | $128 \times 128$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$ |
| c4 | $64 \times 64$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$ |
| c5 | $32 \times 32$ | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 512 \\ \text{MHSA}, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$ |
| # params. | | $25.5 \times 10^6$ | $20.8 \times 10^6$ |
| M.Adds | | $85.4 \times 10^9$ | $102.98 \times 10^9$ |
| TPU steptime | | 786.5 ms | 1032.66 ms |

➔ Only difference is the use MHSA layer in c5

➔ BoT50 has only 1.2 x multiple-adds and 1.3 x training overheads with 1.2x fewer parameters.

# Results: Comparison w/ ResNet on Coco

*BoT50 is better than R50 and R101, competitive with R152*

*Relative positional encoding boosts performance*

| Backbone | $AP^{bb}$ | $AP^{mk}$ |
|---|---|---|
| R50 | 42.1 | 37.7 |
| BoT50 | 43.6 (+ **1.5**) | 38.9 (+ **1.2**) |
| R101 | 43.3 | 38.4 |
| BoT101 | 45.5 (+ **2.2**) | 40.4 (+ **2.0**) |
| R152 | 44.2 | 39.1 |
| BoT152 | 46.0 (+ **1.8**) | 40.6 (+ **1.5**) |

| Backbone | Att. Type | $AP^{bb}$ | $AP^{mk}$ |
|---|---|---|---|
| R50 | - | 42.1 | 37.7 |
| BoT50 | $qk^T$ | 42.7 (+ **0.6**) | 38.3 (+ **0.6**) |
| BoT50 | $qr_{relative}^T$ | 43.1 (+ **1.0**) | 38.4 (+ **0.7**) |
| BoT50 | $qk^T + qr_{relative}^T$ | 43.6 (+ **1.5**) | 38.9 (+ **1.2**) |
| BoT50 | $qk^T + qr_{abs}^T$ | 42.5 (+ **0.4**) | 38.1 (+ **0.4**) |

**Surpasses previous best published model on ResNet**

# Results: Comparison w/ ResNet on Coco

*BoTNet benefits from training on larger images*

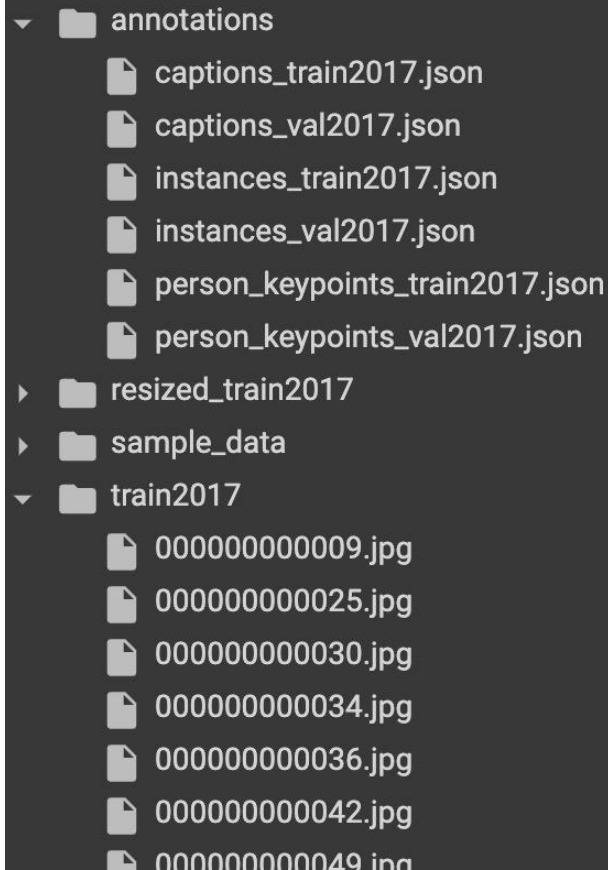| Backbone | res | $AP^{bb}$ | $AP^{mk}$ |
|----------|-----|-----------|-----------|
| R50 | 1280 | 44.0 | 39.5 |
| BoT50 | 1024 | 45.9 (**+ 1.9**) | 40.7 (**+ 1.2**) |
| BoT50 | 1280 | 46.1 (**+ 2.1**) | 41.2 (**+ 1.8**) |
| R101 | 1280 | 46.4 | 41.2 |
| BoT101 | 1024 | 47.4 (**+ 1.0**) | 42.0 (**+ 0.8**) |
| BoT101 | 1280 | 47.9 (**+ 1.5**) | 42.4 (**+ 1.2**) |



**Scales well w/ larger images**

# COCO dataset

```
annotation{
    "id"            : int,
    "image_id"      : int,
    "category_id"   : int,
    "segmentation"  : RLE or [polygon],
    "area"          : float,
    "bbox"          : [x,y,width,height],
    "iscrowd"       : 0 or 1,
}

categories[{
    "id"            : int,
    "name"          : str,
    "supercategory" : str,
}]
```
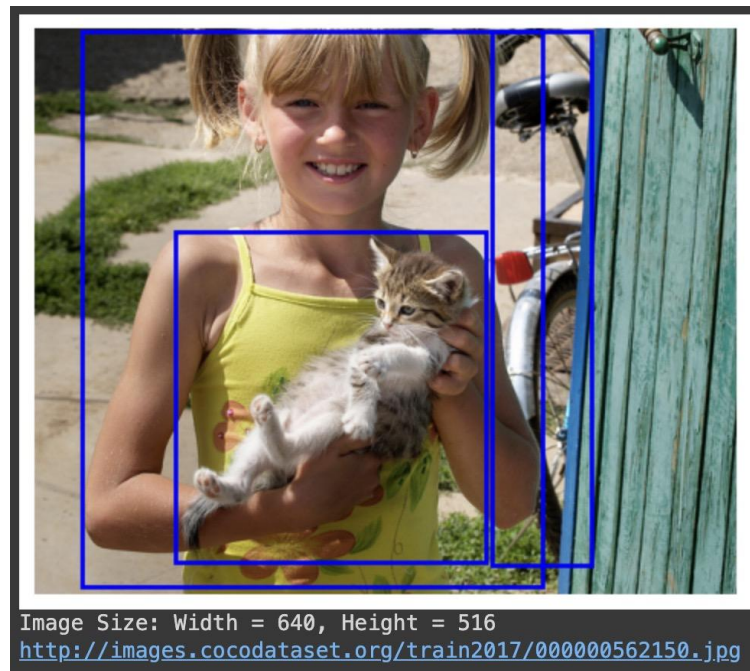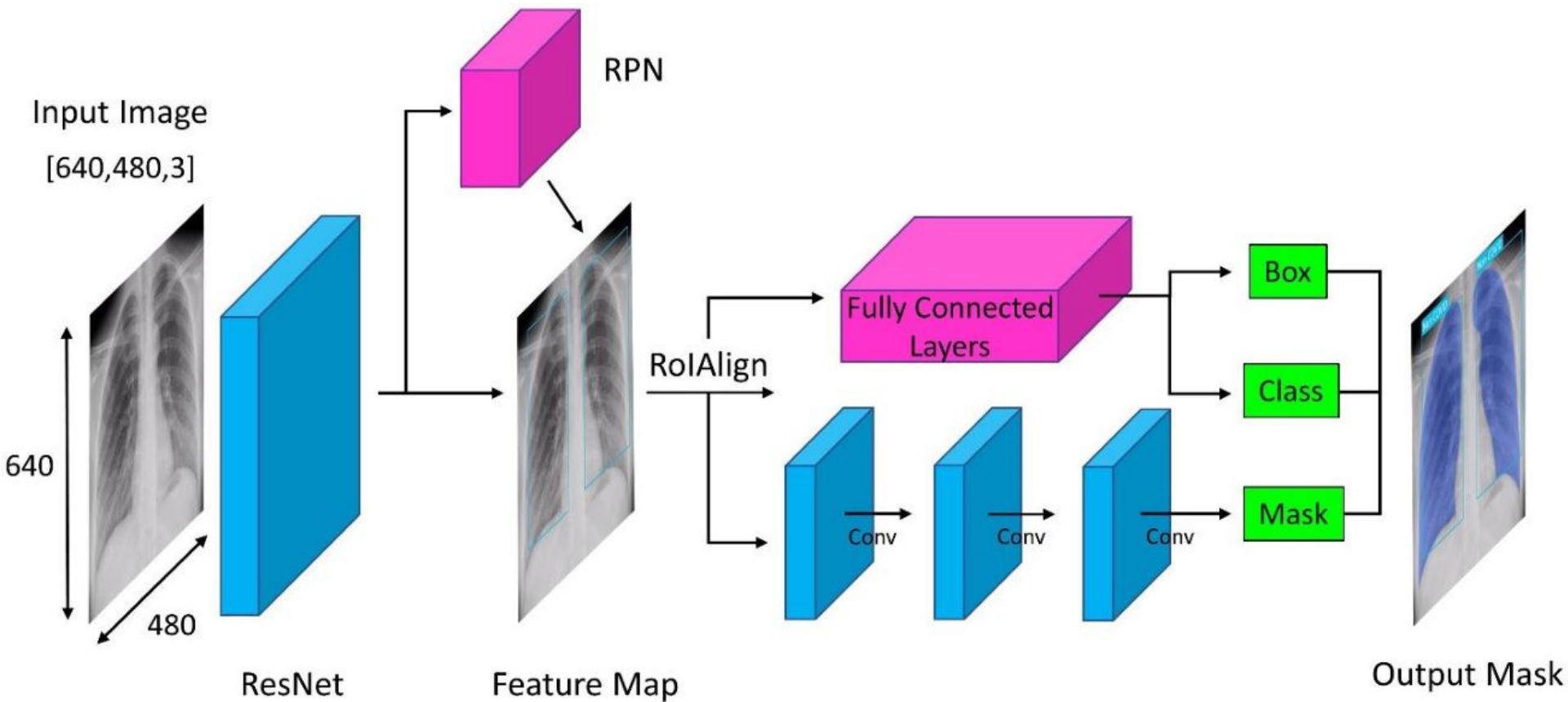
- ▼ 📁 annotations
  - 📄 captions_train2017.json
  - 📄 captions_val2017.json
  - 📄 instances_train2017.json
  - 📄 instances_val2017.json
  - 📄 person_keypoints_train2017.json
  - 📄 person_keypoints_val2017.json
- ▶ 📁 resized_train2017
- ▶ 📁 sample_data
- ▼ 📁 train2017
  - 📄 000000000009.jpg
  - 📄 000000000025.jpg
  - 📄 000000000030.jpg
  - 📄 000000000034.jpg
  - 📄 000000000036.jpg
  - 📄 000000000042.jpg
  - 📄 000000000049.jpg

Image Size: Width = 640, Height = 424
http://images.cocodataset.org/train2017/000000372938.jpg

Image Size: Width = 640, Height = 516
http://images.cocodataset.org/train2017/000000562150.jpg
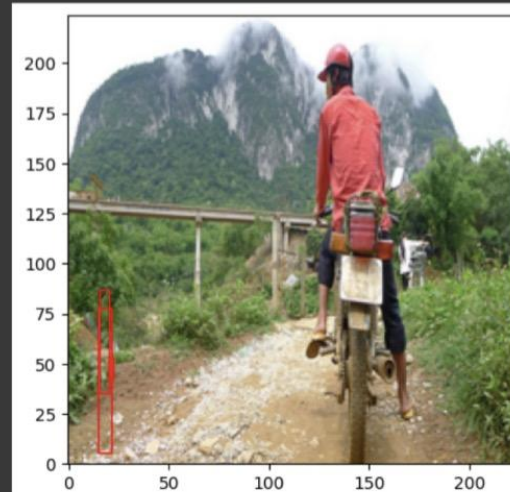
Instance segmentation

Object detection

Mask R-CNN Architecture

```
1  import torch.nn as nn
2  from torchvision.models.detection import MaskRCNN
3  from torchvision.models.detection.rpn import AnchorGenerator
4  from torchvision.ops import MultiScaleRoIAlign
5
6
7  backbone = ResNet(Bottleneck, [3, 4, 6, 3], resolution=(224, 224), heads=4)
8
9  # you are effectively informing the rest of the MaskRCNN model about the shape
10 # of the tensors it will receive from the backbone. This ensures that subsequent
11 # layers can be correctly configured to work with these tensors.
12 backbone.out_channels = 2048
13
14 anchor_generator = AnchorGenerator(sizes=((32, 64, 128, 256, 512),), aspect_ratios=((0.5, 1.0, 2.0),))
15
16 roi_pooler = MultiScaleRoIAlign(featmap_names=['0'], output_size=7, sampling_ratio=2)
17
18 mask_roi_pooler = MultiScaleRoIAlign(featmap_names=['0'], output_size=14, sampling_ratio=2)
19
20 # Define the model
21 model = MaskRCNN(backbone, num_classes=91,  # COCO has 80 classes + background
22                  rpn_anchor_generator=anchor_generator,
23                  box_roi_pool=roi_pooler,
24                  mask_roi_pool=mask_roi_pooler)
```

# Our Code

➔ **Image Classification on CIFAR**

➔ **Does slightly better than CNN and uses fewer parameters**

➔ **Exploration of the architecture**

➔ **[Here](#)**

# Thank you!