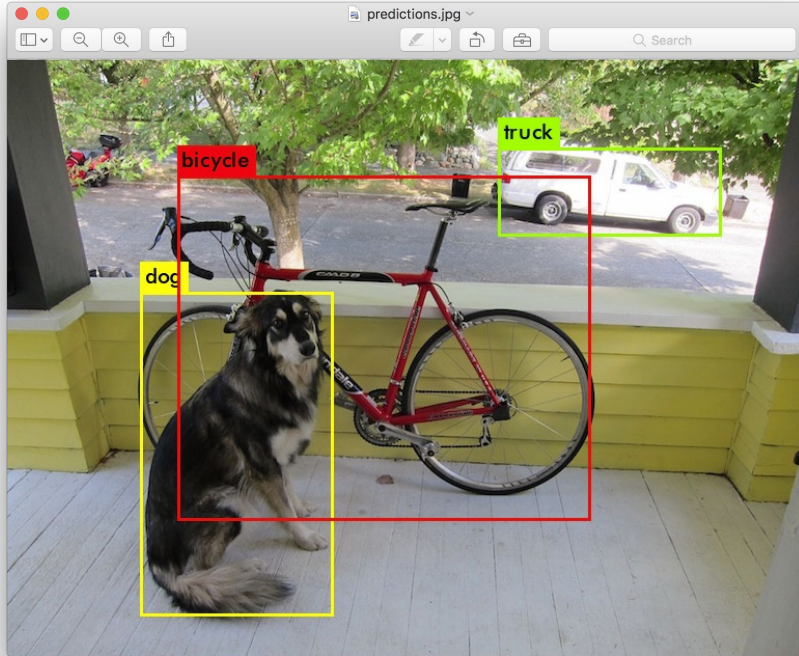


# Faster R-CNN

code approach

Evgeny Ponomarev

# Object detection

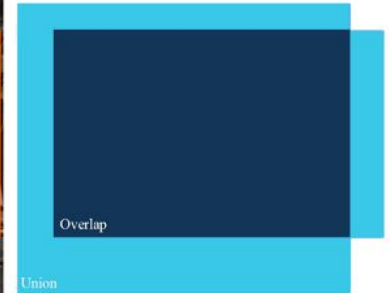


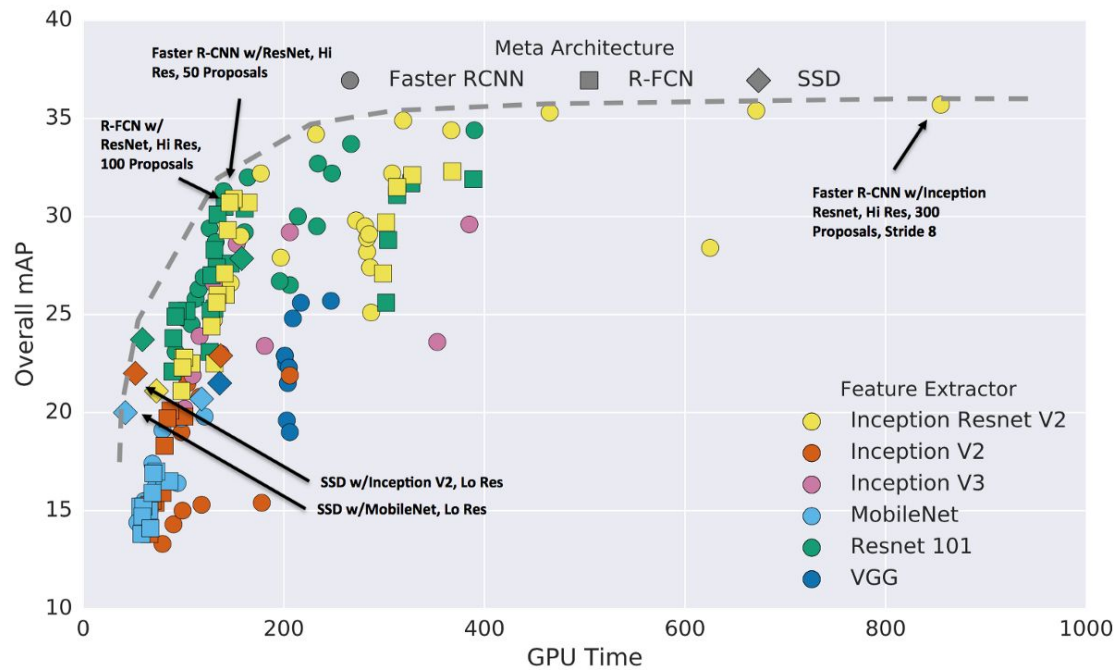
$$\text{mAP} = \frac{1}{|Q_R|} \sum_{q \in Q_R} \text{AP}(q)$$



-  Ground truth
-  Prediction

$$\text{IoU} = \frac{\text{area of overlap}}{\text{area of union}}$$



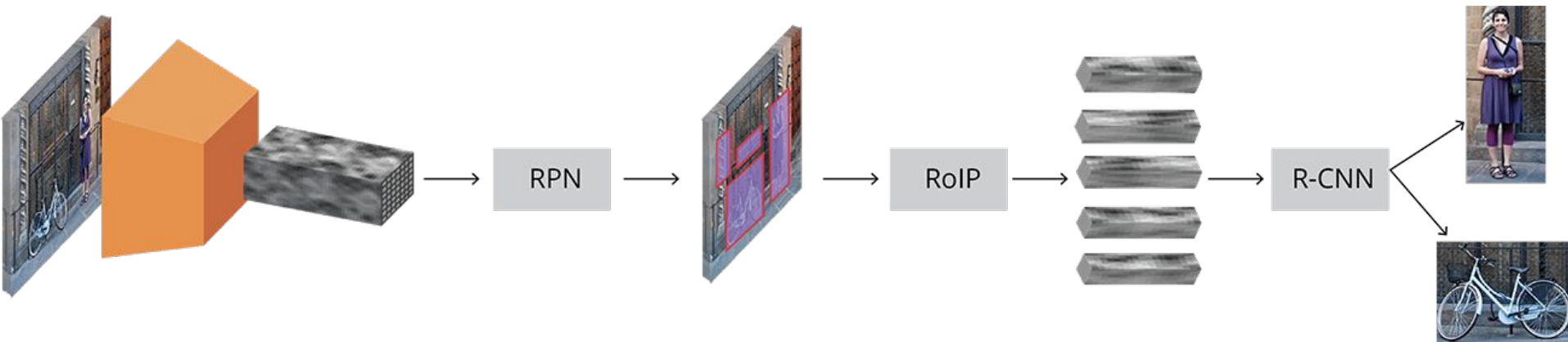


Speed vs accuracy trade-off. A higher mAP and a lower GPU Time is optimal. COCO 🐔 dataset

# Architecture

It all starts with an image, from which we want to obtain:

- a list of bounding boxes.
- a label assigned to each bounding box.
- a probability for each label and bounding box.



```
2  (backbone): Sequential(
3    (body): ResNet(
4      ... SOME RESNET LAYERS ...
5    )
6  )
7  (rpn): RPNModule(
8    (anchor_generator): AnchorGenerator(
9      (cell_anchors): BufferList()
10   )
11   (head): RPNHead(
12     (conv): Conv2d(1024, 1024, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
13     (cls_logits): Conv2d(1024, 9, kernel_size=(1, 1), stride=(1, 1))
14     (bbox_pred): Conv2d(1024, 36, kernel_size=(1, 1), stride=(1, 1))
15   )
16   (box_selector_train): RPNPostProcessor()
17   (box_selector_test): RPNPostProcessor()
18 )
19 (roi_heads): CombinedROIHeads(
20   (box): ROIBoxHead(
21     (feature_extractor): ResNet50Conv5ROIFeatureExtractor(
22       (pooler): Pooler(
23         (poolers): ModuleList(
24           (0): ROIAlign(output_size=(14, 14), spatial_scale=0.0625, sampling_ratio=0)
25         )
26       )
27       (head): ResNetHead(
28         ... SOME RESNET LAYERS ...
29       )
30     )
31     (predictor): FastRCNNPredictor(
32       (avgpool): AvgPool2d(kernel_size=7, stride=7, padding=0)
33       (cls_score): Linear(in_features=2048, out_features=21, bias=True)
34       (bbox_pred): Linear(in_features=2048, out_features=84, bias=True)
35     )
36     (post_processor): PostProcessor()
37   )
38 )
39 )
```

```
1  (backbone): Sequential(
2  (body): ResNet(
3  ... SOME RESNET LAYERS ...
4  )
5  )
6  )
7  (rpn): RPNModule(
8  (anchor_generator): AnchorGenerator(
9  (cell_anchors): BufferList()
10 )
11 (head): RPNHead(
12 (conv): Conv2d(1024, 1024, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
13 (cls_logits): Conv2d(1024, 9, kernel_size=(1, 1), stride=(1, 1))
14 (bbox_pred): Conv2d(1024, 36, kernel_size=(1, 1), stride=(1, 1))
15 )
16 (box_selector_train): RPNPostProcessor()
17 (box_selector_test): RPNPostProcessor()
18 )
19 (roi_heads): CombinedROIHeads(
20 (box): ROIBoxHead(
21 (feature_extractor): ResNet50Conv5ROIFeatureExtractor(
22 (pooler): Pooler(
23 (poolers): ModuleList(
24 (0): ROIAlign(output_size=(14, 14), spatial_scale=0.0625, sampling_ratio=0)
25 )
26 )
27 (head): ResNetHead(
28 ... SOME RESNET LAYERS ...
29 )
30 )
31 (predictor): FastRCNNPredictor(
32 (avgpool): AvgPool2d(kernel_size=7, stride=7, padding=0)
33 (cls_score): Linear(in_features=2048, out_features=21, bias=True)
34 (bbox_pred): Linear(in_features=2048, out_features=84, bias=True)
35 )
36 (post_processor): PostProcessor()
37 )
38 )
39 )
```

# Whole model forward

```
def forward(self, images, targets=None):
    """
    Arguments:
        images (list[Tensor] or ImageList): images to be processed
        targets (list[BoxList]): ground-truth boxes present in the image (optional)

    Returns:
        result (list[BoxList] or dict[Tensor]): the output from the model.
        During training, it returns a dict[Tensor] which contains the losses.
        During testing, it returns list[BoxList] contains additional fields
        like 'scores', 'labels' and 'mask' (for Mask R-CNN models).

    """
    if self.training and targets is None:
        raise ValueError("In training mode, targets should be passed")
    images = to_image_list(images)
    features = self.backbone(images.tensors)
    proposals, proposal_losses = self.rpn(images, features, targets)
    if self.roi_heads:
        x, result, detector_losses = self.roi_heads(features, proposals, targets)
    else:
        # RPN-only models don't have roi_heads
        x = features
        result = proposals
        detector_losses = {}

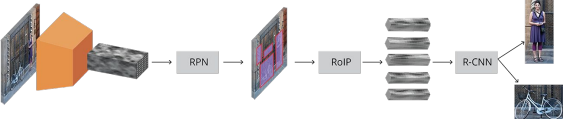
    if self.training:
        losses = {}
        losses.update(detector_losses)
        losses.update(proposal_losses)
        return losses

    return result
```

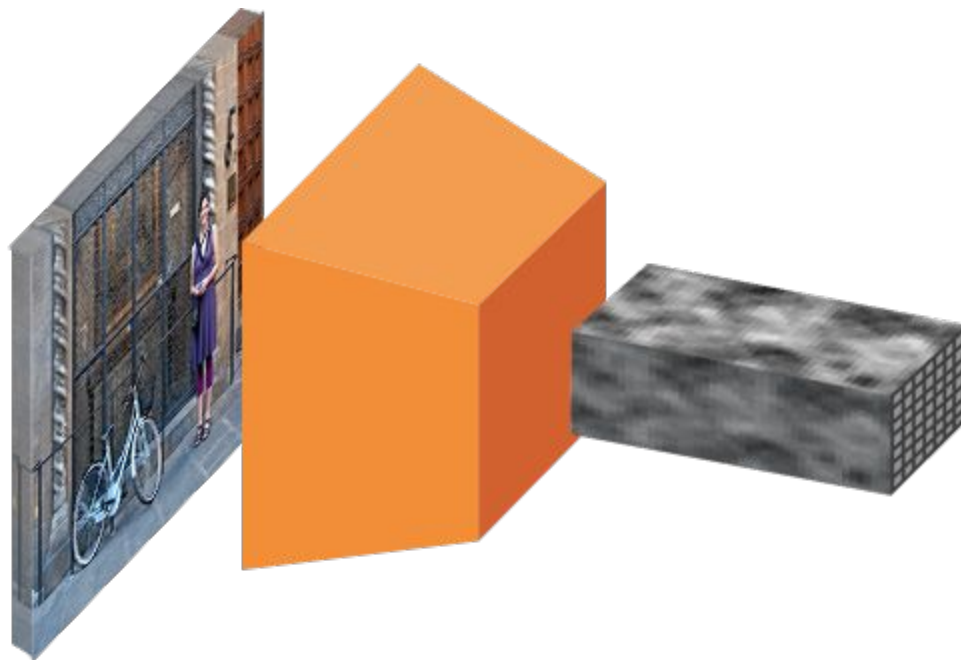
```
2 (backbone): Sequential(
3   (body): ResNet(
4     ... SOME RESNET LAYERS ...
5   )
6 )
7 (rpn): RPNModule(
8   (anchor_generator): AnchorGenerator(
9     (cell_anchors): BufferList()
10  )
11  (head): RPNHead(
12    (conv): Conv2d(1024, 1024, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
13    (cls_logits): Conv2d(1024, 9, kernel_size=(1, 1), stride=(1, 1))
14    (bbox_pred): Conv2d(1024, 36, kernel_size=(1, 1), stride=(1, 1))
15  )
16  (box_selector_train): RPNPostProcessor()
17  (box_selector_test): RPNPostProcessor()
18 )
19 (roi_heads): CombinedROIHeads(
20   (box): ROIBoxHead(
21     (feature_extractor): ResNet50Conv5ROIFeatureExtractor(
22       (pooler): Pooler(
23         (poolers): ModuleList(
24           (0): ROIAlign(output_size=(14, 14), spatial_scale=0.0625, sampling_ratio=0)
25         )
26       )
27       (head): ResNetHead(
28         ... SOME RESNET LAYERS ...
29       )
30     )
31     (predictor): FastRCNNPredictor(
32       (avgpool): AvgPool2d(kernel_size=7, stride=7, padding=0)
33       (cls_score): Linear(in_features=2048, out_features=21, bias=True)
34       (bbox_pred): Linear(in_features=2048, out_features=84, bias=True)
35     )
36     (post_processor): PostProcessor()
37   )
38 )
39 )
```



```
1 generalizedRCNN(  
2 (backbone): Sequential(  
3 (body): ResNet(  
4 (stem): StemWithFixedBatchNorm(  
5 (conv1): Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3), bias=False)  
6 (bn1): FrozenBatchNorm2d()  
7 )  
8 (layer1): Sequential(  
9 (0): BottleneckWithFixedBatchNorm(  
10 (downsample): Sequential(  
11 (0): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)  
12 (1): FrozenBatchNorm2d()  
13 )  
14 (conv1): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)  
15 (bn1): FrozenBatchNorm2d()  
16 (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)  
17 (bn2): FrozenBatchNorm2d()  
18 (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)  
19 (bn3): FrozenBatchNorm2d()  
20 )  
21 (1): BottleneckWithFixedBatchNorm(...)  
22 )  
23 (layer2): Sequential(  
24 (0): BottleneckWithFixedBatchNorm(...)  
25 (1): BottleneckWithFixedBatchNorm(...)  
26 (2): BottleneckWithFixedBatchNorm(...)  
27 (3): BottleneckWithFixedBatchNorm(...)  
28 )  
29 (layer3): Sequential(  
30 (0): BottleneckWithFixedBatchNorm(...)  
31 (1): BottleneckWithFixedBatchNorm(...)  
32 (2): BottleneckWithFixedBatchNorm(...)  
33 (3): BottleneckWithFixedBatchNorm(...)  
34 (4): BottleneckWithFixedBatchNorm(...)  
35 (5): BottleneckWithFixedBatchNorm(...)  
36 )  
37 )  
38 )  
39 (rpn): RPNModule(  
40 (anchor_generator): AnchorGenerator(  
41 (cell_anchors): BufferList()  
42 )  
43 )  
44 )
```

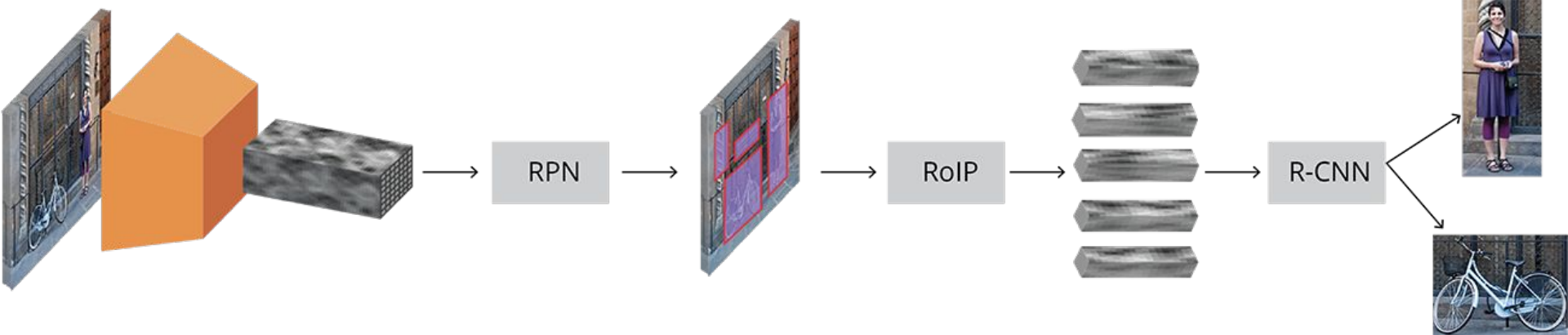


# Backbone



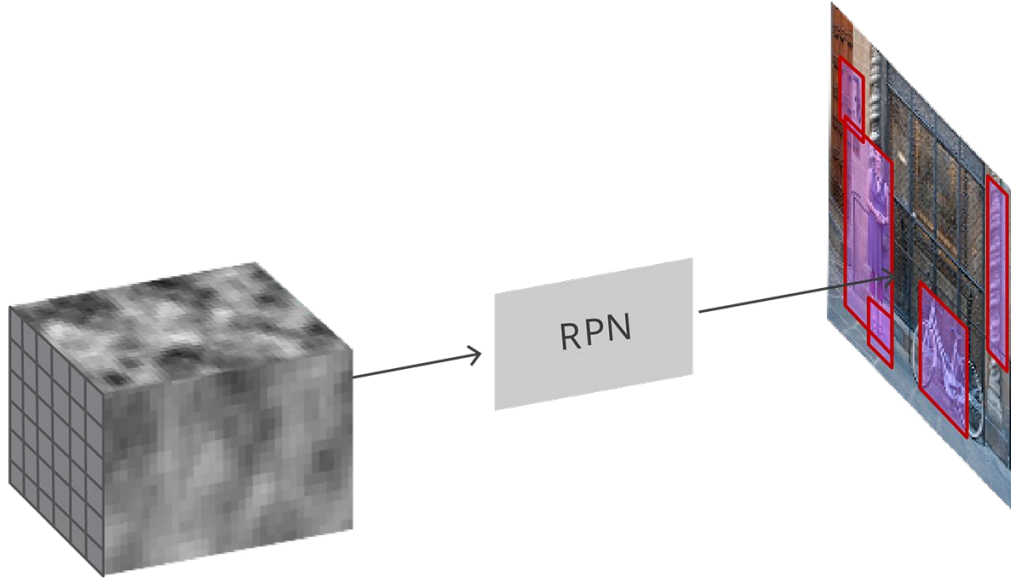
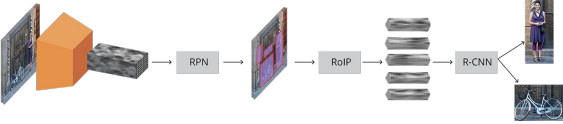
Input:  
image batch  
typ: [1, 3, 1066, 800]

Output:  
feature map  
typ: [1, 1024, 67, 50]



```
2  (backbone): Sequential(
3  (body): ResNet(
4  ... SOME RESNET LAYERS ...
5  )
6  )
7  (rpn): RPNModule(
8  (anchor_generator): AnchorGenerator(
9  (cell_anchors): BufferList()
10 )
11 (head): RPNHead(
12 (conv): Conv2d(1024, 1024, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
13 (cls_logits): Conv2d(1024, 9, kernel_size=(1, 1), stride=(1, 1))
14 (bbox_pred): Conv2d(1024, 36, kernel_size=(1, 1), stride=(1, 1))
15 )
16 (box_selector_train): RPNPostProcessor()
17 (box_selector_test): RPNPostProcessor()
18 )
19 (roi_heads): CombinedROIHeads(
20 (box): ROIBoxHead(
21 (feature_extractor): ResNet50Conv5ROIFeatureExtractor(
22 (pooler): Pooler(
23 (poolers): ModuleList(
24 (0): ROIALign(output_size=(14, 14), spatial_scale=0.0625, sampling_ratio=0)
25 )
26 )
27 (head): ResNetHead(
28 ... SOME RESNET LAYERS ...
29 )
30 )
31 (predictor): FastRCNNPredictor(
32 (avgpool): AvgPool2d(kernel_size=7, stride=7, padding=0)
33 (cls_score): Linear(in_features=2048, out_features=21, bias=True)
34 (bbox_pred): Linear(in_features=2048, out_features=84, bias=True)
35 )
36 (post_processor): PostProcessor()
37 )
38 )
39 )
```

# Region Proposal Network



Input:

image;

targets - list(bbox)

$\text{bbox} \sim (c, x, y, h, w)$

Output:

predicted bboxes

losses

# Region Proposal Network forward

```
def forward(self, images, features, targets=None):

    objectness, rpn_box_regression = self.head(features)
    anchors = self.anchor_generator(images, features)

    if self.training:
        return self._forward_train(anchors, objectness, rpn_box_regression, targets)
    else:
        return self._forward_test(anchors, objectness, rpn_box_regression)

def _forward_train(self, anchors, objectness, rpn_box_regression, targets):
    if self.cfg.MODEL.RPN_ONLY:
        # When training an RPN-only model, the loss is determined by the
        # predicted objectness and rpn_box_regression values and there is
        # no need to transform the anchors into predicted boxes; this is an
        # optimization that avoids the unnecessary transformation.
        boxes = anchors
    else:
        # For end-to-end models, anchors must be transformed into boxes and
        # sampled into a training batch.
        with torch.no_grad():
            boxes = self.box_selector_train(
                anchors, objectness, rpn_box_regression, targets
            )
    loss_objectness, loss_rpn_box_reg = self.loss_evaluator(
        anchors, objectness, rpn_box_regression, targets
    )
    losses = {
        "loss_objectness": loss_objectness,
        "loss_rpn_box_reg": loss_rpn_box_reg,
    }
    return boxes, losses
```

Input:

image;

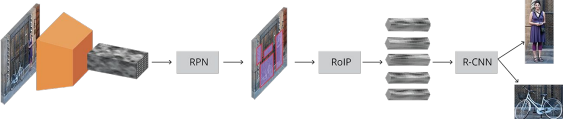
targets - list(bbox)

$\text{bbox} \sim (c, x, y, h, w)$

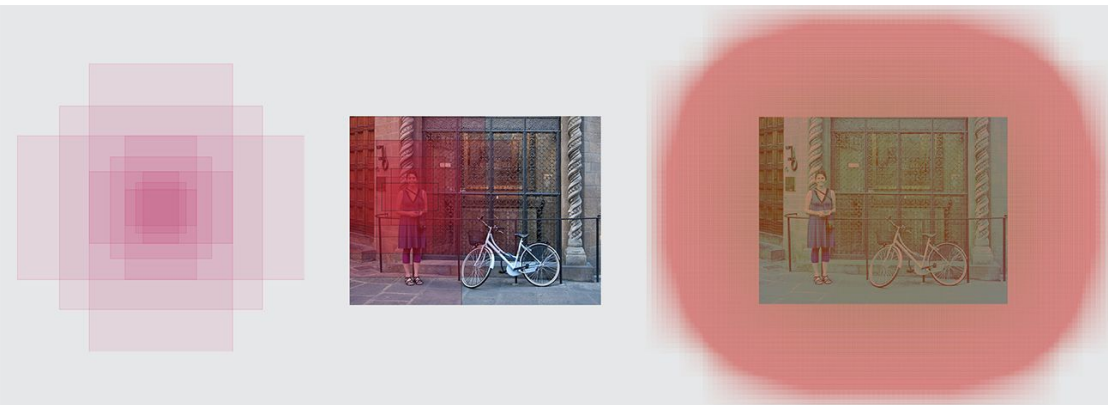
Output:

predicted bboxes

losses

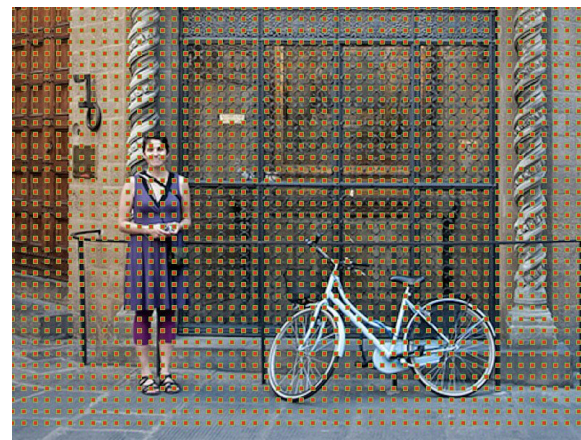


# Region Proposal Network. Anchors



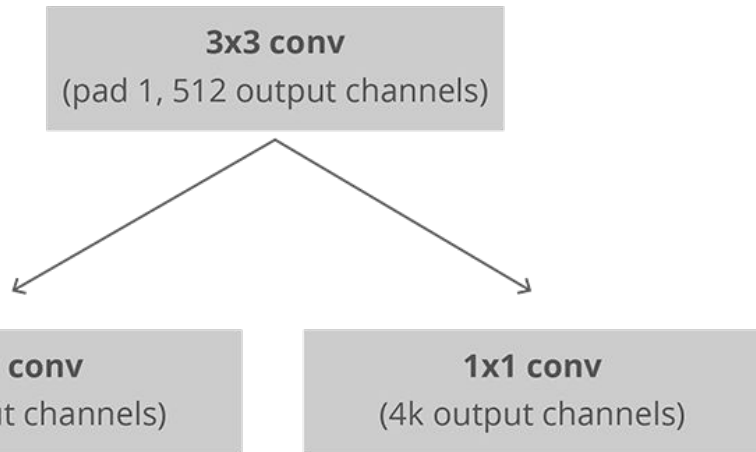
Get set of fixed-size boxes on image that corresponded to each point on feature map

```
def forward(self, image_list, feature_maps):
    grid_sizes = [feature_map.shape[-2:] for feature_map in feature_maps]
    anchors_over_all_feature_maps = self.grid_anchors(grid_sizes)
    anchors = []
    for i, (image_height, image_width) in enumerate(image_list.image_sizes):
        anchors_in_image = []
        for anchors_per_feature_map in anchors_over_all_feature_maps:
            boxlist = BoxList(
                anchors_per_feature_map, (image_width, image_height), mode="xyxy"
            )
            self.add_visibility_to(boxlist)
            anchors_in_image.append(boxlist)
        anchors.append(anchors_in_image)
    return anchors
```



# Region Proposal Network. Regressor and objectness

```
(rpn): RPNModule(  
  (anchor_generator): AnchorGenerator(  
    (cell_anchors): BufferList(  
    )  
  )  
  (head): RPNHead(  
    (conv): Conv2d(1024, 1024, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))  
    (cls_logits): Conv2d(1024, 9, kernel_size=(1, 1), stride=(1, 1))  
    (bbox_pred): Conv2d(1024, 36, kernel_size=(1, 1), stride=(1, 1))  
  )  
  (box_selector_train): RPNPostProcessor()  
  (box_selector_test): RPNPostProcessor()  
)
```



Regressor computes tilts for given anchors

$$t_x = (x - x_a) / w_a$$

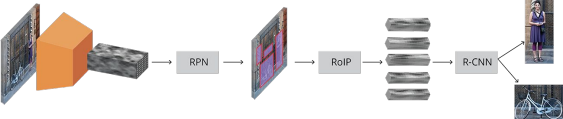
$$t_y = (y - y_a) / h_a$$

$$t_w = \log(w / w_a)$$

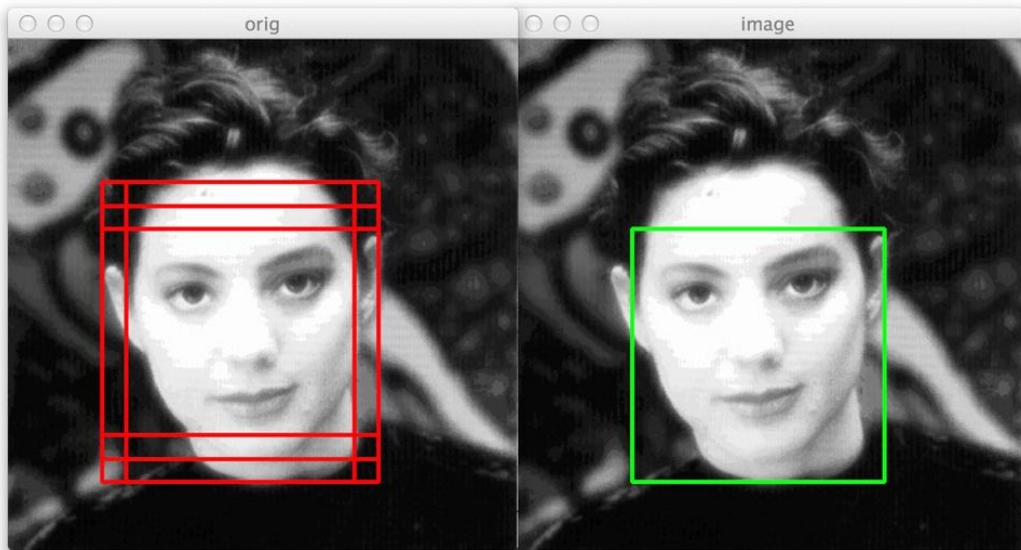
$$t_h = \log(h / h_a)$$

Objectness predict probability of not-background





# Non-maximum suppression



- Start from  $bbox(MAX)$  with highest IoU (intersection over union) score with its *ground-truth bbox* (predicted)
- exclude all bboxes with  $iou(THIS, MAX) \geq THRESH$
- Select the next prediction bbox with the highest IoU from the left bboxes until no bbox left.
- return top of the sorted bbox list:

`list(Bbox)[:2000]`

# Region Proposal Network loss

$$L(p, u, t, t^*) = L_{obj}(p, u) + [u > 0] \cdot L_{reg}(t, t^*)$$

$L_{obj}$  = cross-entropy loss

$L_{reg}$  = smooth L1 loss

$$t_x = (x - x_a) / w_a$$

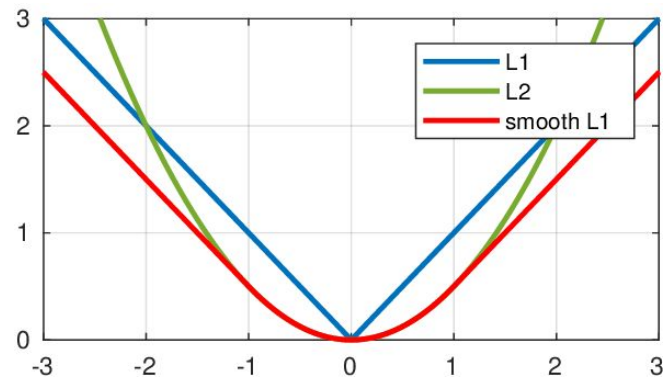
$$t_y = (y - y_a) / h_a$$

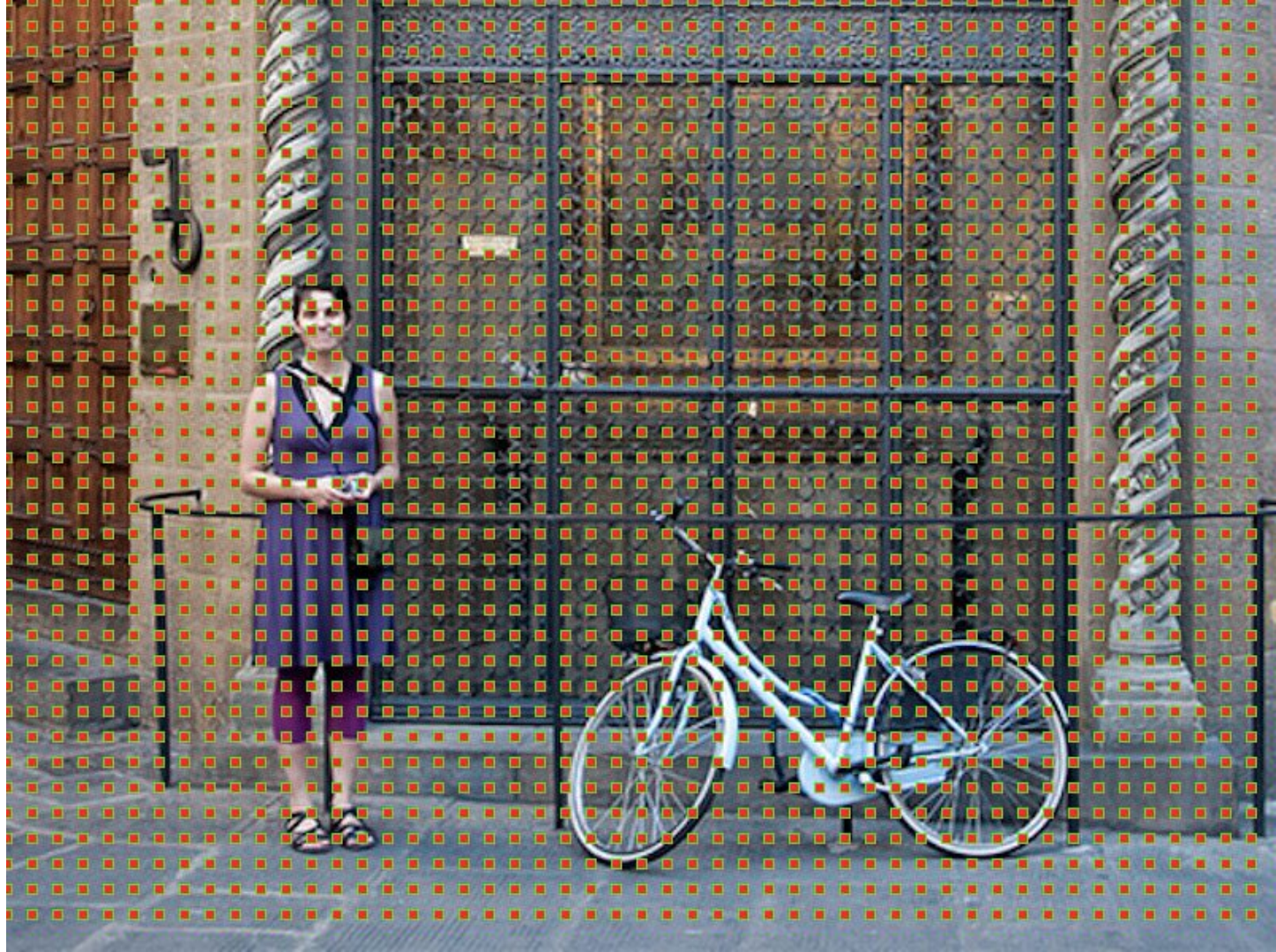
$$t_w = \log(w / w_a)$$

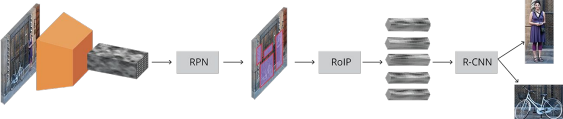
$$t_h = \log(h / h_a)$$

$(x_a, y_a, w_a, h_a)$  – anchor bbox

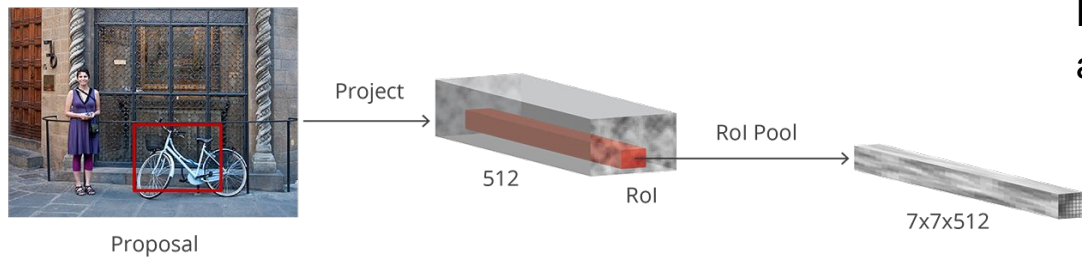
$t^*$  – tilts for ground-truth







# Regions of Interest(RoI) Pool/Align

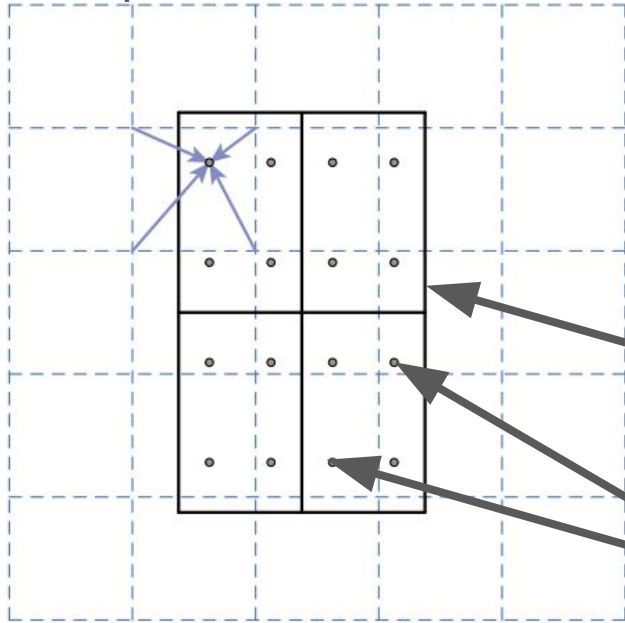


Project bounding box into feature space and then align or pool to fixed size

```
(roi_heads): CombinedROIHeads(
  (box): ROIBoxHead(
    (feature_extractor): ResNet50Conv5ROIFeatureExtractor(
      (pooler): Pooler(
        (poolers): ModuleList(
          (0): ROIAlign(output_size=(14, 14), spatial_scale=0.0625, sampling_ratio=0)
        )
      )
    (head): ResNetHead(
      ... SOME RESNET LAYERS ...
    )
  )
  (predictor): FastRCNNPredictor(
    (avgpool): AvgPool2d(kernel_size=7, stride=7, padding=0)
    (cls_score): Linear(in_features=2048, out_features=21, bias=True)
    (bbox_pred): Linear(in_features=2048, out_features=84, bias=True)
  )
  (post_processor): PostProcessor()
)
```

# Regions of Interest(RoI) Pool/Align

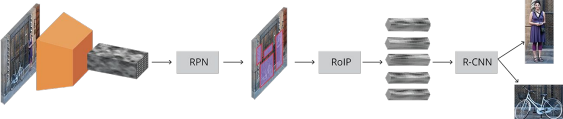
feature map



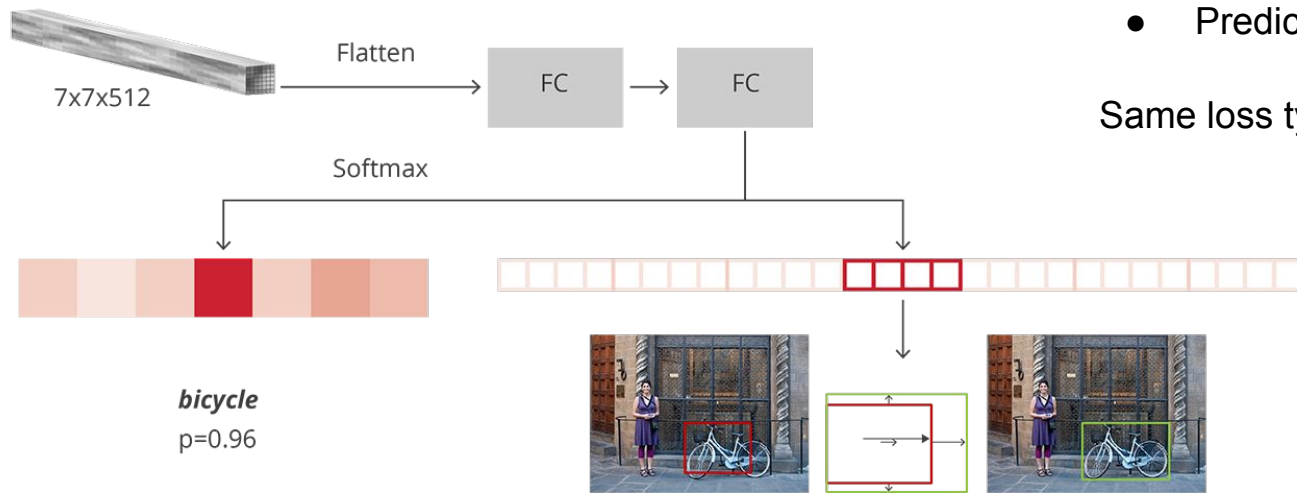
feature map

BBox projection from Image to feature map  
here 2x2, actually typ. 7x7 or 14x14

Points for bilinear interpolation



# (Regions of Interest) Head



- Finalise tilts for bounding box for each class
- Predict class probability

Same loss type as for RPN, but multiclass



person: 0.82

person: 0.83

person: 0.98

person: 0.99

person: 0.99

person: 0.76

person: 0.89

person: 0.97

cup: 0.96

person: 0.80

person: 0.99

person: 1.00

person: 1.00

cup: 0.99

cup: 1.00