

# Object Detection Part I

## Exercise

Minh-Duc Bui  
Phuc-Thinh Nguyen



# Agenda

➤ Introduction

➤ Classification

➤ Classification + Localization

➤ Basic YOLOv1

➤ Question



# Introduction

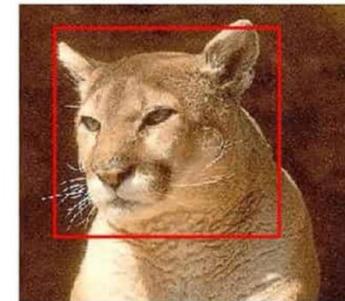
## ❖ Classification, Localization, and Object Detection

### Classification



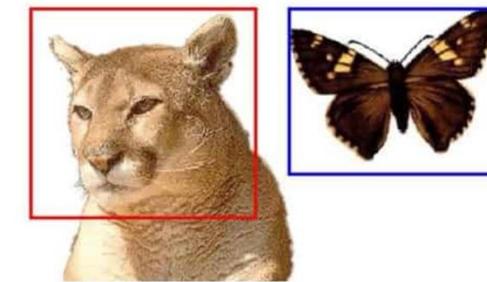
Cougar

### Classification + Localization



Cougar

### Object Detection



Cougar, Butterfly



# Classification

## ❖ Introduction



## ❖ Dataset



LARXEL · UPDATED 5 YEARS AGO

▲ 94

New Notebook

Download



### Dog and Cat Detection

Where is the good boy? (3686 images of cat and dogs with bounding boxes)

[Data Card](#)   [Code \(21\)](#)   [Discussion \(0\)](#)   [Suggestions \(0\)](#)

#### About Dataset

#### Usability ⓘ

8.75

#### About this Data

#### License

CC BY-SA 4.0

In reality, the rise of AI is due to us looking for a way to hone our capacity to show love to our pets.

The computer itself evolved so that it could house more pictures of cattos and doggos and big data is largely a result of us posting more and more pictures of our pets in the interwebs.

With this dataset, containing 3686 images with bounding boxes for 2 classes (doggos, cattos), you can finally apply your skills to a relevant and timeless problem of the human race: find the good boyes so that you can give them treats and pet them.

#### Expected update frequency

Never

#### Tags

[Arts and Entertainment](#)[Image](#)[Computer Vision](#)

## ❖ Dataset

- ▼  annotations
  - Cats\_Test0.xml
  - Cats\_Test1.xml
  - Cats\_Test10.xml
  - Cats\_Test100.xml
  - Cats\_Test1000.xml
  - Cats\_Test1001.xml
- images
  - Cats\_Test0.png
  - Cats\_Test1.png
  - Cats\_Test10.png
  - Cats\_Test100.png
  - Cats\_Test1000.png
  - Cats\_Test1001.png

Cats\_Test0.xml (559 B)

---

**About this file**

This file does not have a description yet.

**Add Suggestion**

---

```
<annotation>
  <folder>images</folder>
  <filename>Cats_Test0.png</filename>
  <size>
    <width>233</width>
    <height>350</height>
    <depth>3</depth>
  </size>
  <segmented>0</segmented>
  <object>
    <name>cat</name>
    <pose>Unspecified</pose>
    <truncated>0</truncated>
    <occluded>0</occluded>
    <difficult>0</difficult>
    <bndbox>
      <xmin>83</xmin>
      <ymin>29</ymin>
      <xmax>197</xmax>
      <ymax>142</ymax>
    </bndbox>
  </object>
</annotation>
```

Data Explorer

Version 1 (1.1 GB)

▼  annotations

- Cats\_Test0.xml
- Cats\_Test1.xml
- Cats\_Test10.xml
- Cats\_Test100.xml
- Cats\_Test1000.xml
- Cats\_Test1001.xml
- Cats\_Test1002.xml
- Cats\_Test1003.xml
- Cats\_Test1004.xml
- Cats\_Test1005.xml
- Cats\_Test1006.xml
- Cats\_Test1007.xml
- Cats\_Test1008.xml
- Cats\_Test1009.xml
- Cats\_Test101.xml
- Cats\_Test1010.xml
- Cats\_Test1011.xml
- Cats\_Test1012.xml
- Cats\_Test1013.xml
- Cats\_Test1014.xml
- Cats\_Test1015.xml
- Cats\_Test1016.xml
- Cats\_Test1017.xml
- Cats\_Test1018.xml
- Cats\_Test1019.xml
- Cats\_Test102.xml
- Cats\_Test1020.xml

## ❖ Dataset

```
1 # Data directory
2 annotations_dir = os.path.join(data_dir, 'annotations')
3 image_dir = os.path.join(data_dir, 'images')
4
5 # Get list of image files and create a dummy dataframe to split the data
6 image_files = [f for f in os.listdir(image_dir) if os.path.isfile(os.path.join(image_dir, f))]
7 df = pd.DataFrame({'image_name': image_files})
8
9 df.head()
10
✓ 0.0s
```

	image_name
0	Cats_Test3363.png
1	Cats_Test478.png
2	Cats_Test1574.png
3	Cats_Test1212.png
4	Cats_Test3405.png

## ❖ Dataset

```
1 # Split data
2 train_df, val_df = train_test_split(df, test_size=0.2, random_state=42)
3
4 # Transforms
5 transform = transforms.Compose([
6     transforms.Resize((224, 224)),
7     transforms.ToTensor(),
8     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225])
9 ])
10
11 # Datasets
12 train_dataset = ImageDataset(annotations_dir, image_dir, transform=transform)
13 val_dataset = ImageDataset(annotations_dir, image_dir, transform=transform)
14
15 # Filter datasets based on train_df and val_df
16 train_dataset.image_files = [f for f in train_dataset.image_files if f in train_df['image_name'].values]
17 val_dataset.image_files = [f for f in val_dataset.image_files if f in val_df['image_name'].values]
18
19 # DataLoaders
20 train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
21 val_loader = DataLoader(val_dataset, batch_size=32, shuffle=False)
22
✓ 1.1s
```

## ❖ Dataset

```
1 # Dataset Class
2 class ImageDataset(Dataset):
3     def __init__(self, annotations_dir, image_dir, transform=None):
4         self.annotations_dir = annotations_dir
5         self.image_dir = image_dir
6         self.transform = transform
7         self.image_files = self.filter_images_with_multiple_objects()
8
9     def filter_images_with_multiple_objects(self):
10        valid_image_files = []
11        for f in os.listdir(self.image_dir):
12            if os.path.isfile(os.path.join(self.image_dir, f)):
13                img_name = f
14                annotation_name = os.path.splitext(img_name)[0] + ".xml"
15                annotation_path = os.path.join(self.annotations_dir, annotation_name)
16
17                if self.count_objects_in_annotation(annotation_path) <= 1:
18                    valid_image_files.append(img_name)
19                else:
20                    print(
21                        f"Image {img_name} has multiple objects and will be excluded from the dataset"
22                    )
23        return valid_image_files
24
25    def count_objects_in_annotation(self, annotation_path):
26        try:
27            tree = ET.parse(annotation_path)
28            root = tree.getroot()
29            count = 0
30            for obj in root.findall("object"):
31                count += 1
32            return count
33        except FileNotFoundError:
34            return 0
35
36    def __len__(self):
37        return len(self.image_files)
```

## ❖ Dataset

```
39     def __getitem__(self, idx):
40         # Image path
41         img_name = self.image_files[idx]
42         img_path = os.path.join(self.image_dir, img_name)
43
44         # Load image
45         image = Image.open(img_path).convert("RGB")
46
47         # Annotation path
48         annotation_name = os.path.splitext(img_name)[0] + ".xml"
49         annotation_path = os.path.join(self.annotations_dir, annotation_name)
50
51         # Parse annotation
52         label = self.parse_annotation(annotation_path)
53
54         if self.transform:
55             image = self.transform(image)
56
57         return image, label
58
59     def parse_annotation(self, annotation_path):
60         tree = ET.parse(annotation_path)
61         root = tree.getroot()
62
63         label = None
64         for obj in root.findall("object"):
65             name = obj.find("name").text
66             if (
67                 label is None
68             ): # Take the first label for now. We are working with 1 label per image
69                 label = name
70
71         # Convert label to numerical representation (0 for cat, 1 for dog)
72         label_num = 0 if label == "cat" else 1 if label == "dog" else -1
73
74         return label_num
```

## ❖ Model

```
1 # Model
2 model = models.resnet18(weights=ResNet18_Weights.DEFAULT)
3 num_ftrs = model.fc.in_features
4 model.fc = nn.Linear(num_ftrs, 2) # 2 classes: cat and dog
5
6 # Device
7 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
8 model.to(device)
9
10 # Loss and Optimizer
11 criterion = nn.CrossEntropyLoss()
12 optimizer = optim.Adam(model.parameters(), lr=0.001)
```

## ❖ Training

```
1 # Training Loop
2 num_epochs = 10
3 for epoch in tqdm.tqdm(range(num_epochs), desc="Epochs"):
4     model.train()
5     for batch_idx, (data, targets) in enumerate(tqdm.tqdm(train_loader, desc="Training", leave=False)):
6         data = data.to(device)
7         targets = targets.to(device)
8
9         scores = model(data)
10        loss = criterion(scores, targets)
11
12        optimizer.zero_grad()
13        loss.backward()
14        optimizer.step()
15
16    # Validation
17    model.eval()
18    with torch.no_grad():
19        correct = 0
20        total = 0
21        for data, targets in tqdm.tqdm(val_loader, desc="Validation", leave=False):
22            data = data.to(device)
23            targets = targets.to(device)
24            scores = model(data)
25            _, predictions = scores.max(1)
26            correct += (predictions == targets).sum()
27            total += targets.size(0)
28
29    print(f'Epoch {epoch+1}/{num_epochs}, Validation Accuracy: {float(correct)/float(total)*100:.2f}%')
30
```

Epoch 1/10, Validation Accuracy: 93.50%

Epoch 2/10, Validation Accuracy: 95.12%

Epoch 3/10, Validation Accuracy: 89.43%

Epoch 4/10, Validation Accuracy: 81.84%

Epoch 5/10, Validation Accuracy: 94.58%

Epoch 6/10, Validation Accuracy: 92.01%

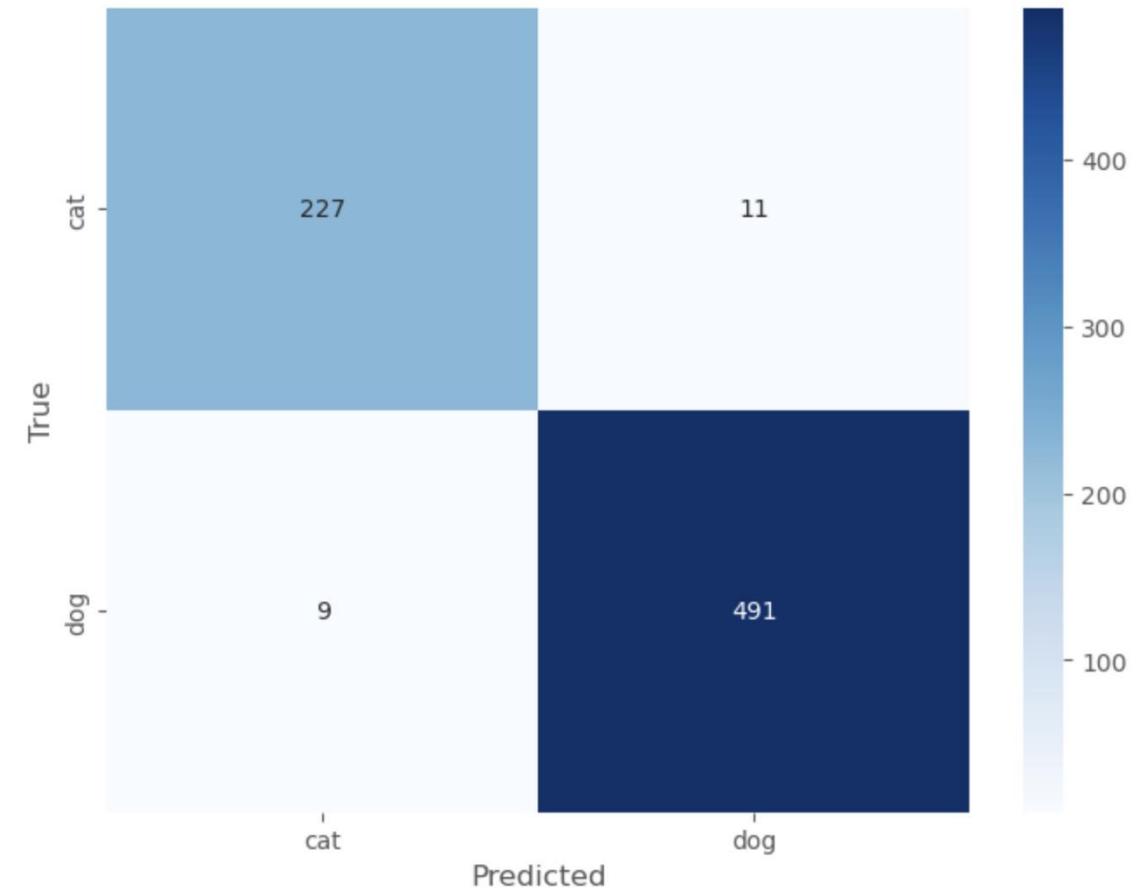
Epoch 7/10, Validation Accuracy: 92.14%

Epoch 8/10, Validation Accuracy: 93.77%

Epoch 9/10, Validation Accuracy: 89.97%

Epoch 10/10, Validation Accuracy: 87.53%

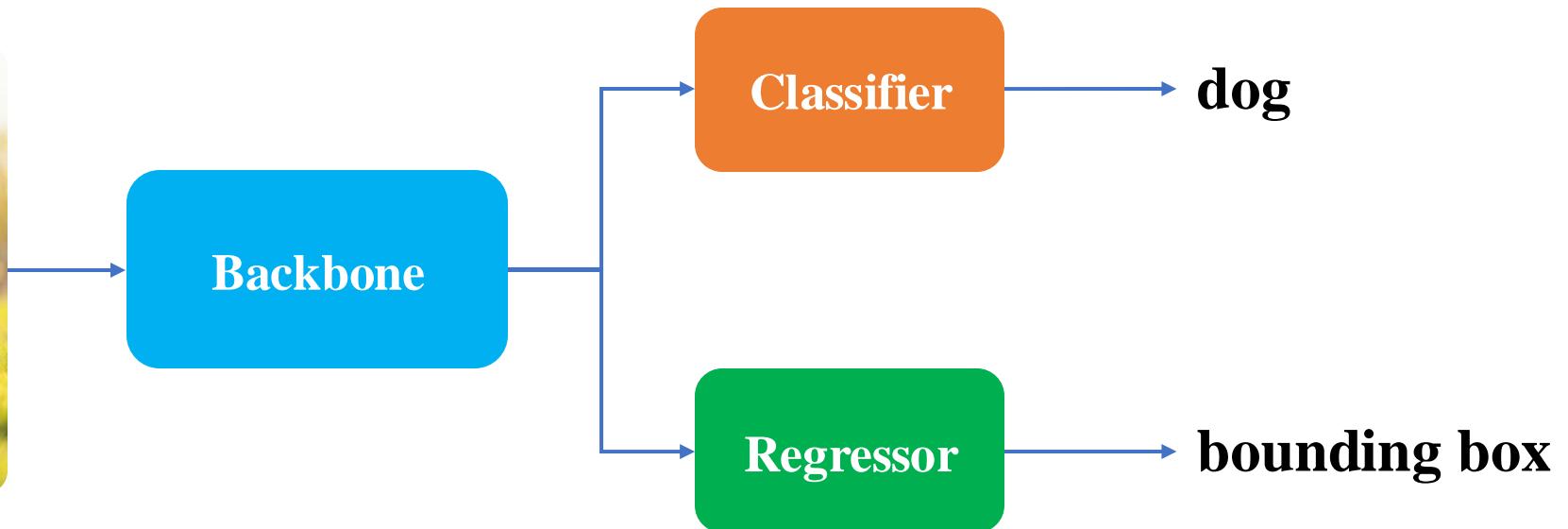
## ❖ Inference





# Classification + Localization

## ❖ Introduction



## ❖ Dataset

```
35     def __getitem__(self, idx):
36         # Image path
37         img_name = self.image_files[idx]
38         img_path = os.path.join(self.image_dir, img_name)
39
40         # Load image
41         image = Image.open(img_path).convert("RGB")
42
43         # Annotation path
44         annotation_name = os.path.splitext(img_name)[0] + ".xml"
45         annotation_path = os.path.join(self.annotations_dir, annotation_name)
46
47         # Parse annotation
48         label, bbox = self.parse_annotation(annotation_path) # Get both label and bbox
49
50         if self.transform:
51             image = self.transform(image)
52
53     return image, label, bbox
```

## ❖ Dataset

```
1 <annotation>
2   <folder>images</folder>
3   <filename>Cats_Test0.png</filename>
4   <size>
5     <width>233</width>
6     <height>350</height>
7     <depth>3</depth>
8   </size>
9   <segmented>0</segmented>
10  <object>
11    <name>cat</name>
12    <pose>Unspecified</pose>
13    <truncated>0</truncated>
14    <occluded>0</occluded>
15    <difficult>0</difficult>
16    <bndbox>
17      <xmin>83</xmin>
18      <ymin>29</ymin>
19      <xmax>197</xmax>
20      <ymax>142</ymax>
21    </bndbox>
22  </object>
23 </annotation>
```

```
55 def parse_annotation(self, annotation_path):
56   tree = ET.parse(annotation_path)
57   root = tree.getroot()
58
59   # Get image size for normalization
60   image_width = int(root.find('size/width').text)
61   image_height = int(root.find('size/height').text)
62
63   label = None
64   bbox = None
65   for obj in root.findall('object'):
66     name = obj.find('name').text
67     if label is None: # Take the first label
68       label = name
69     # Get bounding box coordinates
70     xmin = int(obj.find('bndbox/xmin').text)
71     ymin = int(obj.find('bndbox/ymin').text)
72     xmax = int(obj.find('bndbox/xmax').text)
73     ymax = int(obj.find('bndbox/ymax').text)
74
75     # Normalize bbox coordinates to [0, 1]
76     bbox = [
77       xmin / image_width,
78       ymin / image_height,
79       xmax / image_width,
80       ymax / image_height,
81     ]
82
83   # Convert label to numerical representation (0 for cat, 1 for dog)
84   label_num = 0 if label == 'cat' else 1 if label == 'dog' else -1
85
86   return label_num, torch.tensor(bbox, dtype=torch.float32)
```

## ❖ Model

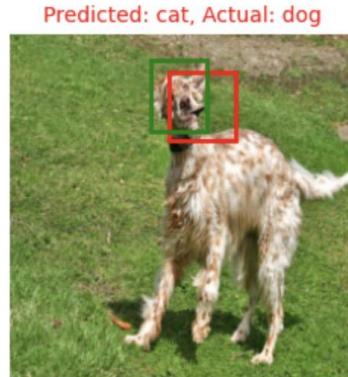
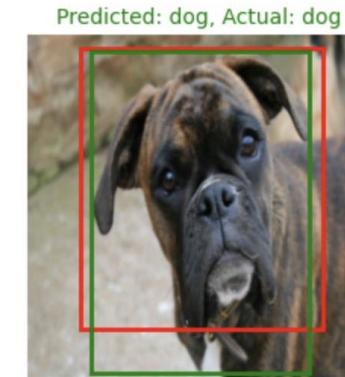
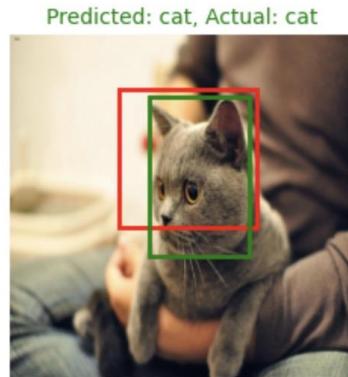
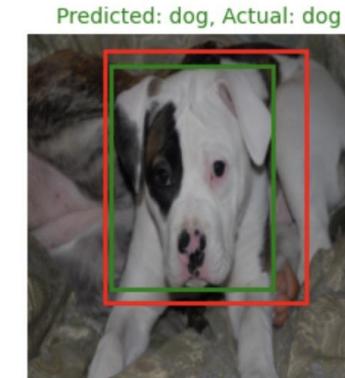
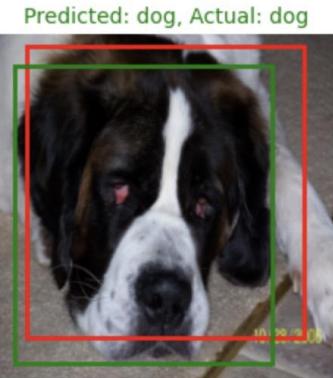
```
1 # Model with Two Heads
2 class TwoHeadedModel(nn.Module):
3     def __init__(self, num_classes=2):
4         super(TwoHeadedModel, self).__init__()
5         self.base_model = models.resnet18(weights=ResNet18_Weights.DEFAULT)
6         self.num_ftrs = self.base_model.fc.in_features
7
8         # Remove the original fully connected layer
9         self.base_model.fc = nn.Identity()
10
11     # Classification head
12     self.classifier = nn.Linear(self.num_ftrs, num_classes)
13
14     # Bounding box regression head
15     self.regressor = nn.Linear(self.num_ftrs, 4)
16
17     def forward(self, x):
18         x = self.base_model(x)
19         class_logits = self.classifier(x)
20         bbox_coords = torch.sigmoid(self.regressor(x))
21         return class_logits, bbox_coords
22
```

```
1 # Model
2 model = TwoHeadedModel()
3
4 # Device
5 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
6 model.to(device)
7
8 # Loss and Optimizer
9 criterion_class = nn.CrossEntropyLoss()
10 criterion_bbox = nn.MSELoss()
11 optimizer = optim.Adam(model.parameters(), lr=0.001)
12
```

## ❖ Training

```
1 # Training Loop
2 num_epochs = 10
3 for epoch in tqdm.tqdm(range(num_epochs), desc="Epochs"):
4     model.train()
5     for batch_idx, (data, targets, bboxes) in enumerate(tqdm.tqdm(train_loader, desc="Training", leave=False)):
6         data = data.to(device)
7         targets = targets.to(device)
8         bboxes = bboxes.to(device)
9
10        scores, pred_bboxes = model(data)
11        loss_class = criterion_class(scores, targets)
12        loss_bbox = criterion_bbox(pred_bboxes, bboxes)
13        loss = loss_class + loss_bbox # Combine losses
14
15        optimizer.zero_grad()
16        loss.backward()
17        optimizer.step()
18
19    # Validation
20    model.eval()
21    with torch.no_grad():
22        correct = 0
23        total = 0
24        total_loss_bbox = 0
25        total_samples = 0
26        for data, targets, bboxes in tqdm.tqdm(val_loader, desc="Validation", leave=False):
27            data = data.to(device)
28            targets = targets.to(device)
29            bboxes = bboxes.to(device)
30
31            scores, pred_bboxes = model(data)
32            _, predictions = scores.max(1)
33            correct += (predictions == targets).sum()
34            total += targets.size(0)
35
36            # Calculate bbox loss for monitoring (optional)
37            total_loss_bbox += criterion_bbox(pred_bboxes, bboxes).item() * data.size(0)
38            total_samples += data.size(0)
39
40            avg_loss_bbox = total_loss_bbox / total_samples
41
42            print(f'Epoch {epoch+1}/{num_epochs}, Validation Accuracy: {float(correct)/float(total)*100:.2f}%, '
43                  f'Avg. Bbox Loss: {avg_loss_bbox:.4f}')
44
```

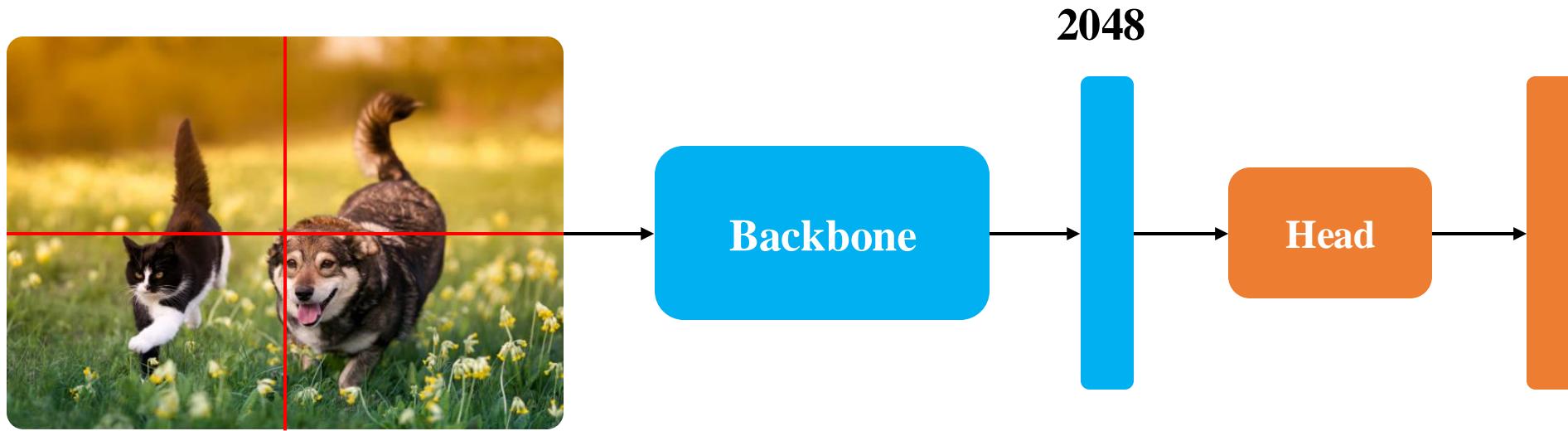
## ❖ Inference





# Classification (>1 object) + Localization

## ❖ Introduction



$$\begin{aligned} & n\_cell * (n\_class + bbox) \\ & 4 * (2 + 4) \end{aligned}$$

## ❖ Dataset

```
1 class MyDataset(Dataset):
2     def __init__(self, annotations_dir, image_dir, transform=None):
3         self.annotations_dir = annotations_dir
4         self.image_dir = image_dir
5         self.transform = transform
6         self.image_files = self.filter_images_with_multiple_objects()
7
8     def filter_images_with_multiple_objects(self):
9         valid_image_files = []
10        for f in os.listdir(self.image_dir):
11            if os.path.isfile(os.path.join(self.image_dir, f)):
12                img_name = f
13                annotation_name = os.path.splitext(img_name)[0] + ".xml"
14                annotation_path = os.path.join(self.annotations_dir, annotation_name)
15
16                if self.count_objects_in_annotation(annotation_path) == 1:
17                    valid_image_files.append(img_name)
18        return valid_image_files
19
20    def count_objects_in_annotation(self, annotation_path):
21        try:
22            tree = ET.parse(annotation_path)
23            root = tree.getroot()
24            count = 0
25            for obj in root.findall("object"):
26                count += 1
27            return count
28        except FileNotFoundError:
29            return 0
```

```
31    def parse_annotation(self, annotation_path):
32        tree = ET.parse(annotation_path)
33        root = tree.getroot()
34
35        # Get image size for normalization
36        image_width = int(root.find("size/width").text)
37        image_height = int(root.find("size/height").text)
38
39        label = None
40        bbox = None
41        for obj in root.findall("object"):
42            name = obj.find("name").text
43            if label is None: # Take the first label
44                label = name
45            # Get bounding box coordinates
46            xmin = int(obj.find("bndbox/xmin").text)
47            ymin = int(obj.find("bndbox/ymin").text)
48            xmax = int(obj.find("bndbox/xmax").text)
49            ymax = int(obj.find("bndbox/ymax").text)
50
51            # Normalize bbox coordinates to [0, 1]
52            bbox = [
53                xmin / image_width,
54                ymin / image_height,
55                xmax / image_width,
56                ymax / image_height,
57            ]
58
59            # Convert label to numerical representation (0 for cat, 1 for dog)
60            label_num = 0 if label == "cat" else 1 if label == "dog" else -1
61
62            return label_num, torch.tensor(bbox, dtype=torch.float32)
63
64    def __len__(self):
65        return len(self.image_files)
```

```
67 def __getitem__(self, idx):
68     img1_file = self.image_files[idx]
69     img1_path = os.path.join(self.image_dir, img1_file)
70
71     annotation_name = os.path.splitext(img1_file)[0] + ".xml"
72     img1_annotations = self.parse_annotation(
73         os.path.join(self.annotations_dir, annotation_name)
74     )
75
76     idx2 = random.randint(0, len(self.image_files) - 1)
77     img2_file = self.image_files[idx2]
78     img2_path = os.path.join(self.image_dir, img2_file)
79
80     annotation_name = os.path.splitext(img2_file)[0] + ".xml"
81     img2_annotations = self.parse_annotation(
82         os.path.join(self.annotations_dir, annotation_name)
83     )
84
85     img1 = Image.open(img1_path).convert("RGB")
86     img2 = Image.open(img2_path).convert("RGB")
87
88     # Horizontal merge
89     merged_image = Image.new(
90         "RGB", (img1.width + img2.width, max(img1.height, img2.height))
91     )
92     merged_image.paste(img1, (0, 0))
93     merged_image.paste(img2, (img1.width, 0))
94     merged_w = img1.width + img2.width
95     merged_h = max(img1.height, img2.height)
96
97     merged_annotations = []
98
99     # No change for objects from img1, already normalized
100    merged_annotations.append(
101        {"bbox": img1_annotations[1].tolist(), "label": img1_annotations[0]}
102    )
103
104    # Adjust bbox coordinates for objects from img2 AND normalize
105    new_bbox = [
106        (img2_annotations[1][0] * img2.width + img1.width) / merged_w, # Normalize xmin
107        img2_annotations[1][1] * img2.height / merged_h, # Normalize ymin
108        (img2_annotations[1][2] * img2.width + img1.width) / merged_w, # Normalize xmax
109        img2_annotations[1][3] * img2.height / merged_h, # Normalize ymax
110    ]
111
112    merged_annotations.append({"bbox": new_bbox, "label": img2_annotations[0]})
113
114    # Convert merged image to tensor
115    if self.transform:
116        merged_image = self.transform(merged_image)
117    else:
118        merged_image = transforms.ToTensor()(merged_image)
119
120    # Convert annotations to 1D tensors, with shape (4,) for bbox and (1,) for label
121    annotations = torch.zeros((len(merged_annotations), 5))
122    for i, ann in enumerate(merged_annotations):
123        annotations[i] = torch.cat((torch.tensor(ann["bbox"])), torch.tensor([ann["label"]])))
124
125    return merged_image, annotations
126
```

## ❖ Model

```
1 class SimpleYOLO(nn.Module):
2     def __init__(self, num_classes):
3         super(SimpleYOLO, self).__init__()
4         self.backbone = models.resnet50(weights=ResNet50_Weights.DEFAULT)
5         self.num_classes = num_classes
6
7         # Remove the final classification layer of ResNet
8         self.backbone = nn.Sequential(*list(self.backbone.children())[:-2])
9
10        # Add the YOLO head
11        self.fcs = nn.Linear(
12            2048, 2 * 2 * (4 + self.num_classes)
13        ) # 2 is for the number of grid cell
14
15    def forward(self, x):
16        # x shape: (batch_size, C, H, W)
17        features = self.backbone(x)
18        features = F.adaptive_avg_pool2d(features, (1, 1)) # shape: (batch_size, 2048, 1, 1)
19        features = features.view(features.size(0), -1) # shape: (batch_size, 2048)
20        features = self.fcs(features)
21
22        return features
23
```

## ❖ LOSS

```
1 def calculate_loss(output, targets, device, num_classes):
2     mse_loss = nn.MSELoss()
3     ce_loss = nn.CrossEntropyLoss()
4
5     batch_size = output.shape[0]
6     total_loss = 0
7
8     output = output.view(batch_size, 2, 2, 4 + num_classes) # Reshape to (batch_size, grid_y, grid_x, 4 + num_classes)
9
10    for i in range(batch_size): # Iterate through each image in the batch
11        for j in range(len(targets[i])): # Iterate through objects in the image
12
13            # Determine which grid cell the object's center falls into
14            # Assuming bbox coordinates are normalized to [0, 1]
15            bbox_center_x = (targets[i][j][0] + targets[i][j][2]) / 2
16            bbox_center_y = (targets[i][j][1] + targets[i][j][3]) / 2
17
18            grid_x = int(bbox_center_x * 2) # Multiply by number of grid cells (2 in this case)
19            grid_y = int(bbox_center_y * 2)
20
21            # 1. Classification Loss for the responsible grid cell
22            # Convert label to one-hot encoding only for this example
23            # One hot encoding
24            label_one_hot = torch.zeros(num_classes, device=device)
25            label_one_hot[int(targets[i][j][4])] = 1
26
27            # Classification loss (using CrossEntropyLoss)
28            classification_loss = ce_loss(output[i, grid_y, grid_x, 4:], label_one_hot)
29
30            # 2. Regression Loss for the responsible grid cell
31            bbox_target = targets[i][j][:4].to(device)
32            regression_loss = mse_loss(output[i, grid_y, grid_x, :4], bbox_target)
33
34            # 3. No Object Loss (for other grid cells)
35            no_obj_loss = 0
36            for other_grid_y in range(2):
37                for other_grid_x in range(2):
38                    if other_grid_y != grid_y or other_grid_x != grid_x:
39                        # MSE loss for predicting no object (all zeros)
40                        no_obj_loss += mse_loss(output[i, other_grid_y, other_grid_x, :4], torch.zeros(4, device=device))
41
42            total_loss += classification_loss + regression_loss + no_obj_loss
43
44    return total_loss / batch_size # Average loss over the batch
45
```

```

1 def train_model(
2     model, train_loader, val_loader, optimizer, num_epochs, device, num_classes
3 ):
4     best_val_accuracy = 0.0
5     train_losses = []
6     val_losses = []
7     train_accuracies = []
8     val_accuracies = []
9
10    for epoch in tqdm.tqdm(range(num_epochs), desc="Epochs"):
11        model.train()
12        running_loss = 0.0
13
14        for images, targets in tqdm.tqdm(train_loader, desc="Batches", leave=False):
15            images = images.to(device)
16
17            optimizer.zero_grad()
18            output = model(images)
19
20            total_loss = calculate_loss(output, targets, device, num_classes)
21
22            total_loss.backward()
23            optimizer.step()
24            running_loss += total_loss.item()
25
26        epoch_loss = running_loss / len(train_loader)
27        train_losses.append(epoch_loss)
28
29    # Validation
30    val_loss, val_accuracy = evaluate_model(model, val_loader, device, num_classes)
31    val_losses.append(val_loss)
32    val_accuracies.append(val_accuracy)
33
34    # Save the best model
35    if val_accuracy > best_val_accuracy:
36        best_val_accuracy = val_accuracy
37        torch.save(model.state_dict(), "best_model.pth")
38
39    return train_losses, val_losses, train_accuracies, val_accuracies
40

```

```

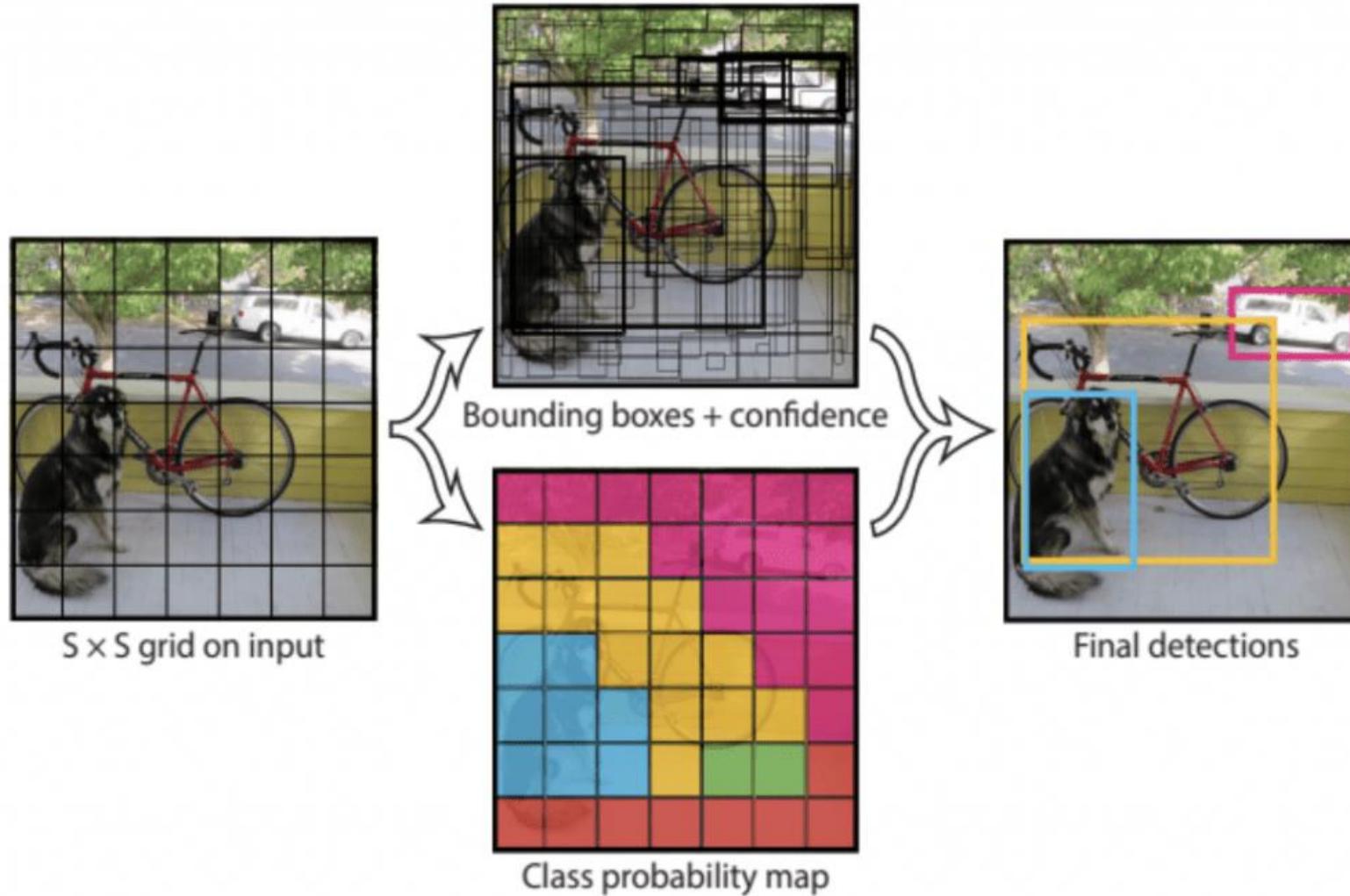
1 def evaluate_model(model, data_loader, device, num_classes):
2     model.eval()
3     running_loss = 0.0
4     all_predictions = []
5     all_targets = []
6
7     with torch.no_grad():
8         for images, targets in tqdm.tqdm(data_loader, desc="Validation", leave=False):
9             images = images.to(device)
10
11             output = model(images)
12
13             total_loss = calculate_loss(output, targets, device, num_classes)
14             running_loss += total_loss.item()
15
16             # Reshape output to (batch_size, grid_y, grid_x, 4 + num_classes)
17             output = output.view(images.shape[0], 2, 2, 4 + num_classes)
18
19             # Collect predictions and targets for mAP calculation
20             for batch_idx in range(images.shape[0]):
21                 for target in targets[batch_idx]:
22                     # Determine responsible grid cell
23                     bbox_center_x = (target[0] + target[2]) / 2
24                     bbox_center_y = (target[1] + target[3]) / 2
25                     grid_x = int(bbox_center_x * 2)
26                     grid_y = int(bbox_center_y * 2)
27
28                     # Class prediction (index of max probability)
29                     prediction = output[batch_idx, grid_y, grid_x, 4:].argmax().item()
30                     all_predictions.append(prediction)
31
32                     all_targets.append(target[4].item())
33
34             val_loss = running_loss / len(data_loader)
35
36             # Convert lists to tensors for PyTorch's metric functions
37             all_predictions = torch.tensor(all_predictions, device=device)
38             all_targets = torch.tensor(all_targets, device=device)
39
40             # Calculate accuracy
41             val_accuracy = (all_predictions == all_targets).float().mean()
42
43     return val_loss, val_accuracy.item()
44

```

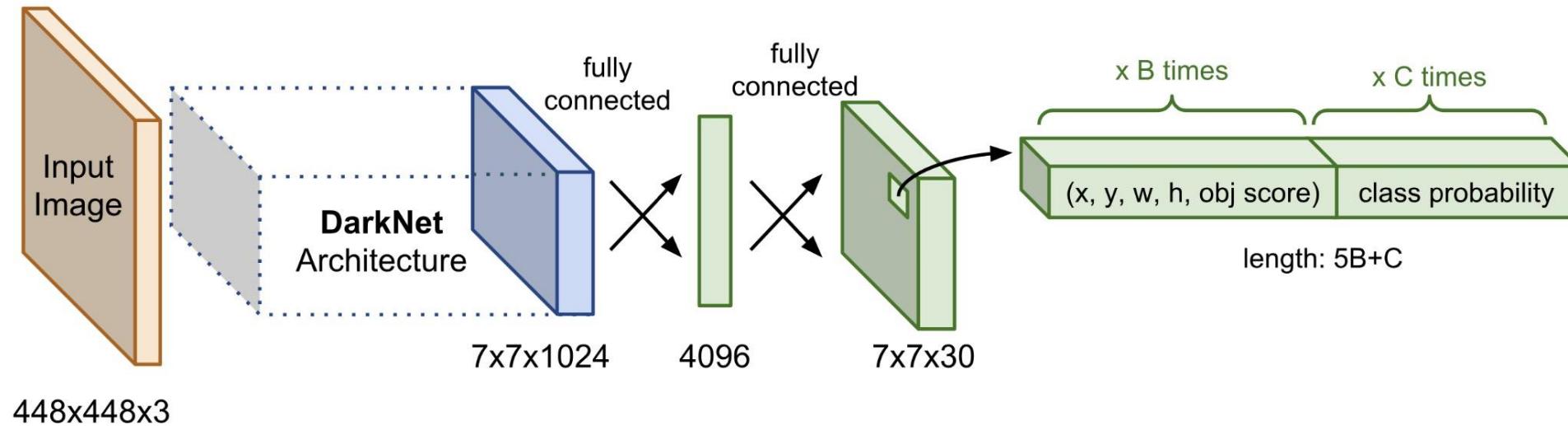


# Basic YOLOv1

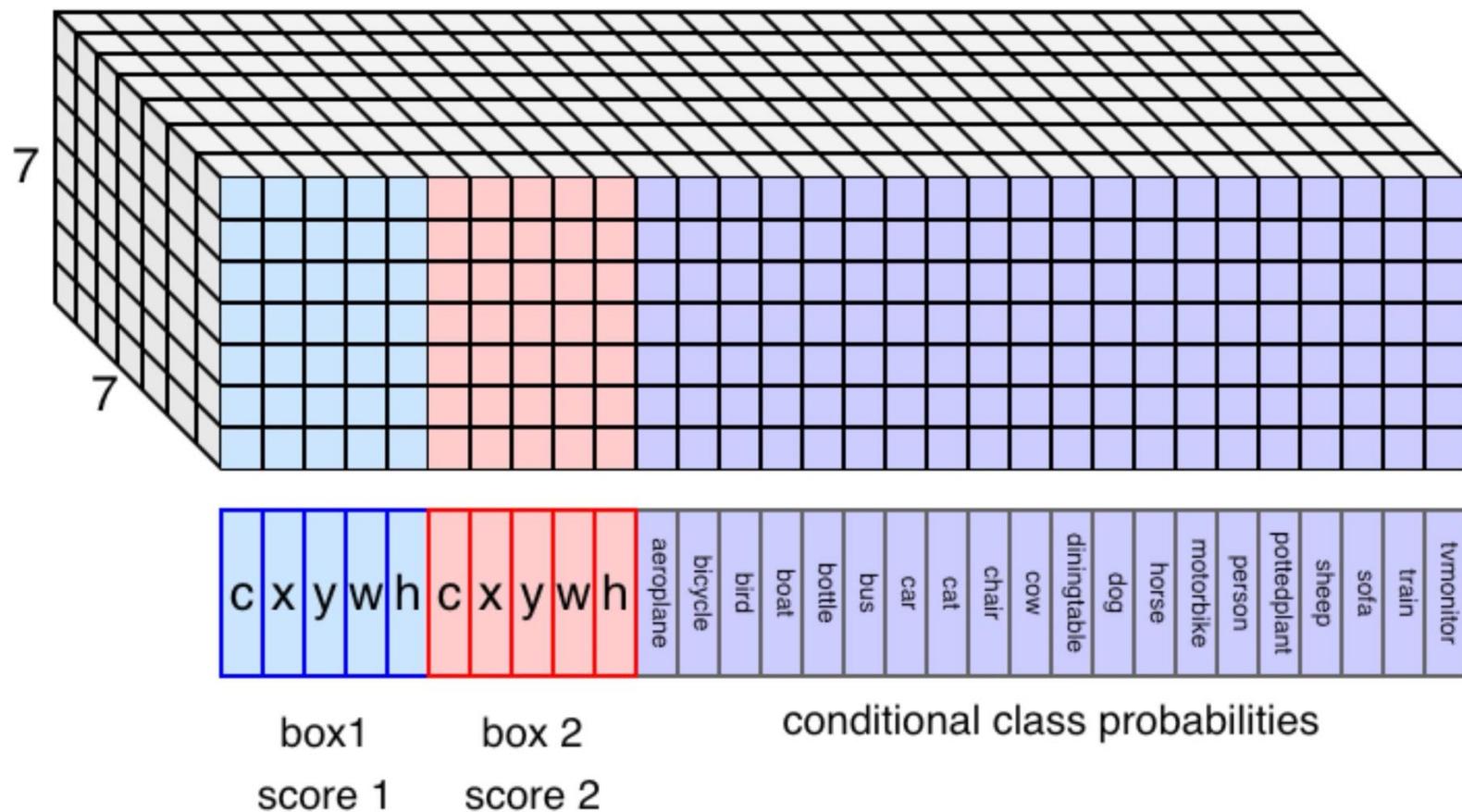
## ❖ Motivation



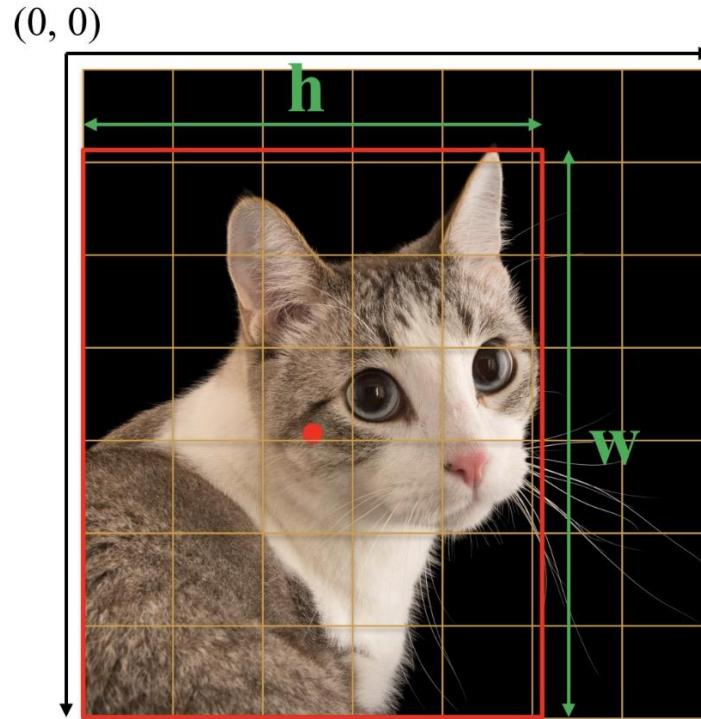
## ❖ Model



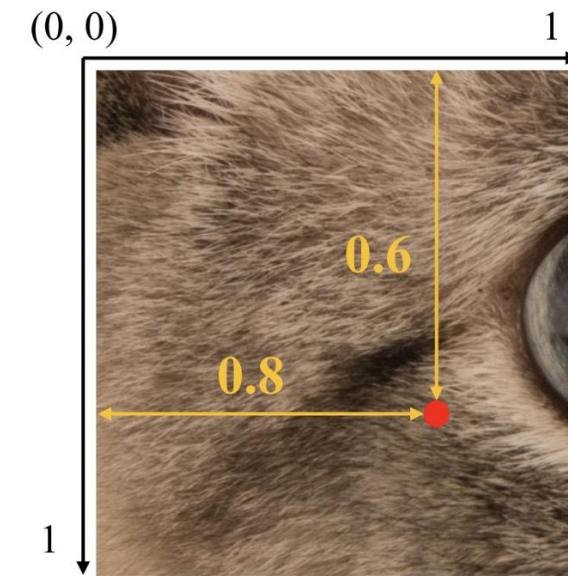
## ❖ Output



## ❖ Bbox



rescale



x	y	w	h
0.8	0.6	5.1	6.2

## ❖ Dataset

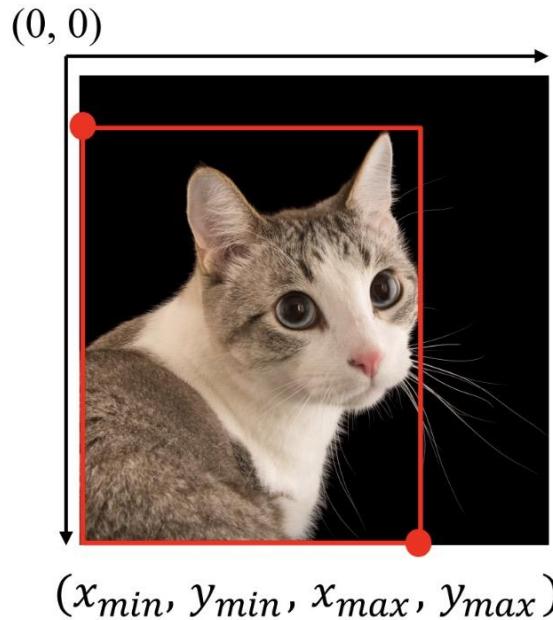
```
1 class CustomVOCDataset(torchvision.datasets.VOCDetection):
2     def __init__(self, class_mapping, S=7, B=2, C=20, custom_transforms=None):
3         # Initialize YOLO-specific configuration parameters.
4         self.S = S  # Grid size S x S
5         self.B = B  # Number of bounding boxes
6         self.C = C  # Number of classes
7         self.class_mapping = class_mapping  # Mapping of class names to class indices
8         self.custom_transforms = custom_transforms
9 
```

## ❖ Dataset

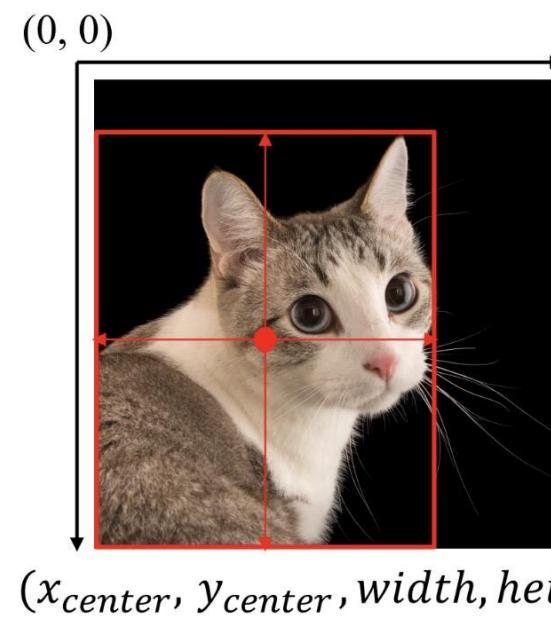
```
10     def __getitem__(self, index):
11         # Get an image and its target (annotations) from the VOC dataset.
12         image, target = super(CustomVOCdataset, self).__getitem__(index)
13         img_width, img_height = image.size
14
15         # Convert target annotations to YOLO format bounding boxes.
16         boxes = convert_to_yolo_format(target, img_width, img_height, self.class_mapping)
17
18         just_boxes = boxes[:,1:]
19         labels = boxes[:,0]
20
21         # Transform the image and bounding boxes
22         if self.custom_transforms:
23             sample = {
24                 'image': np.array(image),
25                 'bboxes': just_boxes,
26                 'labels': labels
27             }
28
29             sample = self.custom_transforms(**sample)
30             image = sample['image']
31             boxes = sample['bboxes']
32             labels = sample['labels']
```

## ❖ Dataset

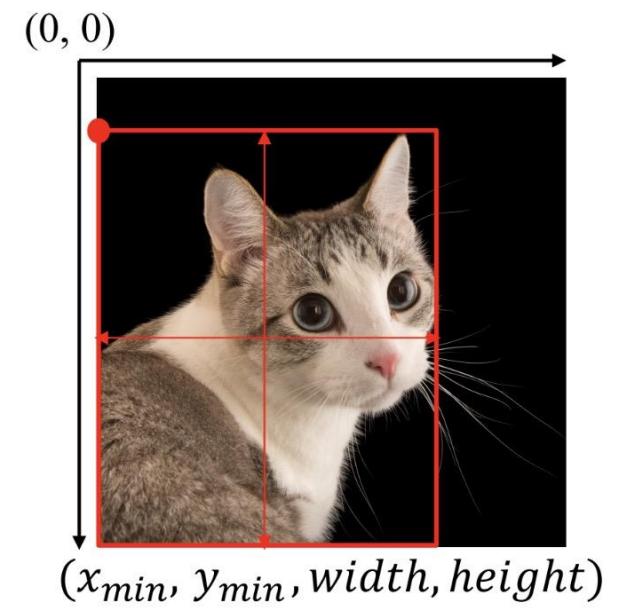
VOC



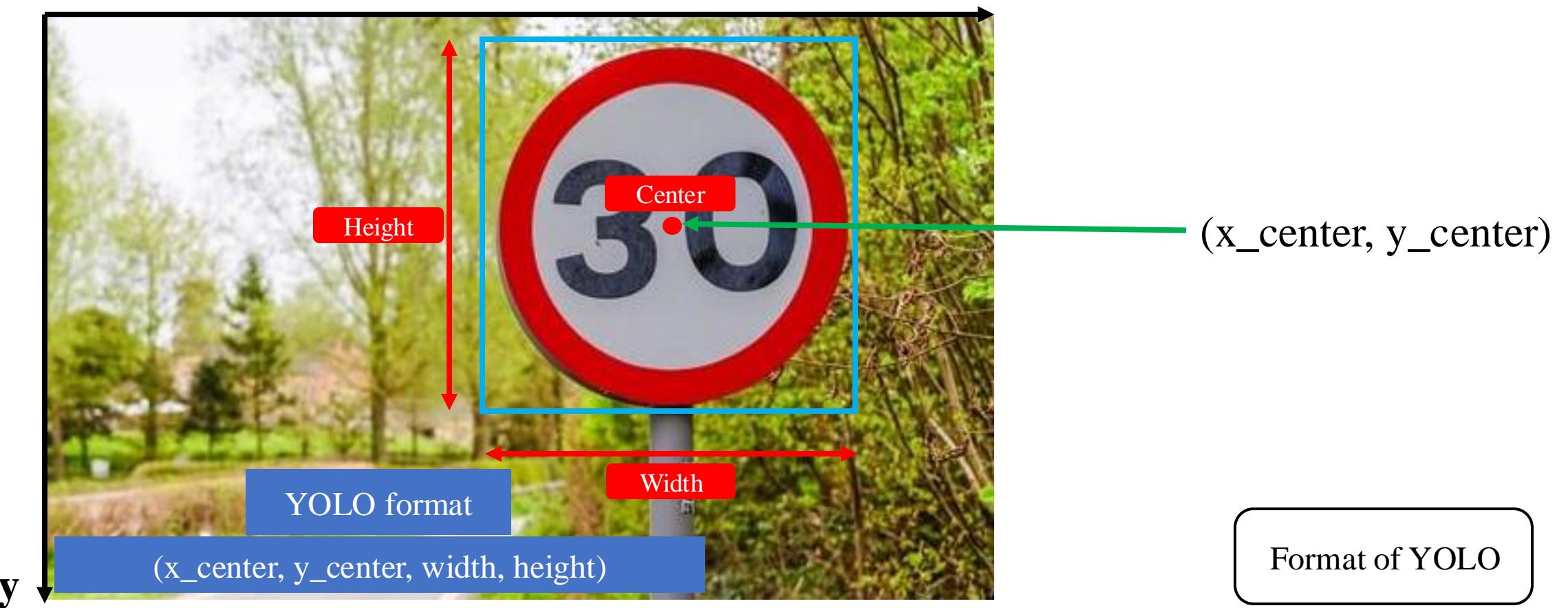
YOLO



COCO



## ❖ Dataset



# ❖ Dataset

## Detailed Explanation

### 1. Input Values:

```
x, y = box coordinates # normalized coordinates (0-1) representing center of object
self.S = 7 # grid size (7x7 grid)
```

### 2. Grid Cell Calculation:

```
i, j = int(self.S * y), int(self.S * x)
```

- `self.S * x`: Scales the normalized x-coordinate (0-1) to grid coordinates (0-7)

- `self.S * y`: Scales the normalized y-coordinate (0-1) to grid coordinates (0-7)
- `int()`: Floors the values to get the grid cell indices

### Example:

- If  $x = 0.3$ ,  $y = 0.6$  and  $S = 7$ :
  - $j = \text{int}(7 \cdot 0.3) = \text{int}(2.1) = 2$  (column index)
  - $i = \text{int}(7 \cdot 0.6) = \text{int}(4.2) = 4$  (row index)

```

# Create an empty label matrix for YOLO ground truth.
label_matrix = torch.zeros((self.S, self.S, self.C + 5 * self.B))

boxes = torch.tensor(boxes, dtype=torch.float32)
labels = torch.tensor(labels, dtype=torch.float32)
image = torch.as_tensor(image, dtype=torch.float32)

# Iterate through each bounding box in YOLO format.
for box, class_label in zip(boxes, labels):
    x, y, width, height = box.tolist()
    class_label = int(class_label)

    # Calculate the grid cell (i, j) that this box belongs to.
    i, j = int(self.S * y), int(self.S * x)
    x_cell, y_cell = self.S * x - j, self.S * y - i

    # Calculate the width and height of the box relative to the grid cell.
    width_cell, height_cell = (
        width * self.S,
        height * self.S,
    )

    # If no object has been found in this specific cell (i, j) before:
    if label_matrix[i, j, 20] == 0:
        # Mark that an object exists in this cell.
        label_matrix[i, j, 20] = 1

        # Store the box coordinates as an offset from the cell boundaries.
        box_coordinates = torch.tensor(
            [x_cell, y_cell, width_cell, height_cell]
        )

        # Set the box coordinates in the label matrix.
        label_matrix[i, j, 21:25] = box_coordinates

        # Set the one-hot encoding for the class label.
        label_matrix[i, j, class_label] = 1

    return image, label_matrix

```

## ❖ Dataset

### 3. Cell-Relative Coordinates:

```
x_cell, y_cell = self.S * x - j, self.S * y - i
```

This calculates the object's position relative to the grid cell's top-left corner:

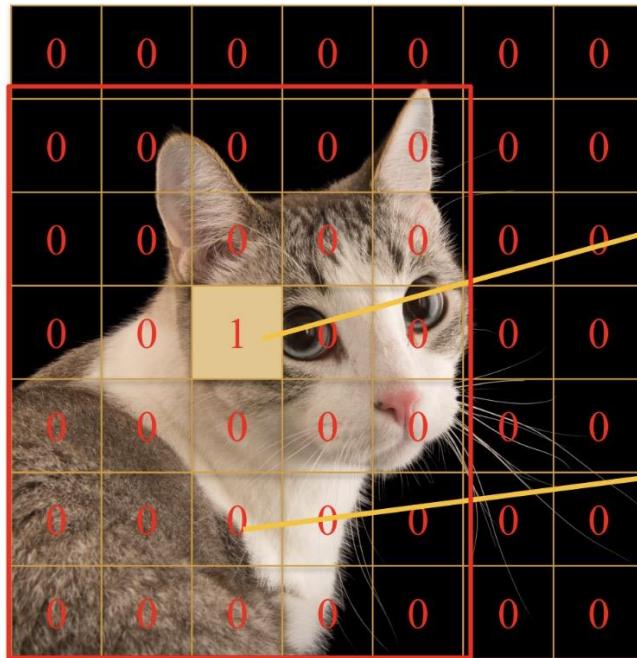
- `self.S * x - j`: Horizontal offset within the cell (0-1)
- `self.S * y - i`: Vertical offset within the cell (0-1)

Example continued:

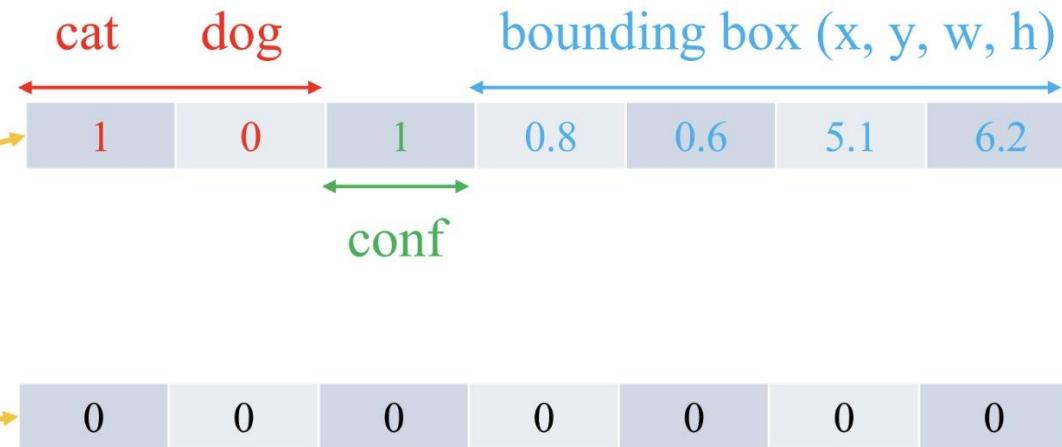
- $x_{\text{cell}} = 7 \cdot 0.3 - 2 = 2.1 - 2 = 0.1$
- $y_{\text{cell}} = 7 \cdot 0.6 - 4 = 4.2 - 4 = 0.2$

```
34 # Create an empty label matrix for YOLO ground truth.
35 label_matrix = torch.zeros((self.S, self.S, self.C + 5 * self.B))
36
37 boxes = torch.tensor(boxes, dtype=torch.float32)
38 labels = torch.tensor(labels, dtype=torch.float32)
39 image = torch.as_tensor(image, dtype=torch.float32)
40
41 # Iterate through each bounding box in YOLO format.
42 for box, class_label in zip(boxes, labels):
43     x, y, width, height = box.tolist()
44     class_label = int(class_label)
45
46     # Calculate the grid cell (i, j) that this box belongs to.
47     i, j = int(self.S * y), int(self.S * x)
48     x_cell, y_cell = self.S * x - j, self.S * y - i
49
50     # Calculate the width and height of the box relative to the grid cell.
51     width_cell, height_cell = (
52         width * self.S,
53         height * self.S,
54     )
55
56     # If no object has been found in this specific cell (i, j) before:
57     if label_matrix[i, j, 20] == 0:
58         # Mark that an object exists in this cell.
59         label_matrix[i, j, 20] = 1
60
61         # Store the box coordinates as an offset from the cell boundaries.
62         box_coordinates = torch.tensor(
63             [x_cell, y_cell, width_cell, height_cell]
64         )
65
66         # Set the box coordinates in the label matrix.
67         label_matrix[i, j, 21:25] = box_coordinates
68
69         # Set the one-hot encoding for the class label.
70         label_matrix[i, j, class_label] = 1
71
72 return image, label_matrix
```

## ❖ Target for training



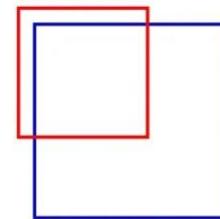
For example, we have 2 classes: *cat* and *dog*.



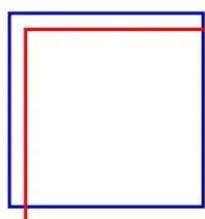
For each object, **only 1 grid cell** is responsible for the prediction.

## ❖ IoU

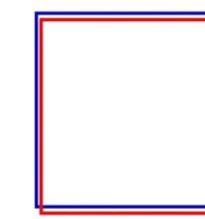
$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



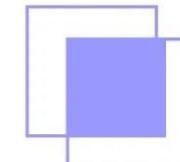
Poor



Good

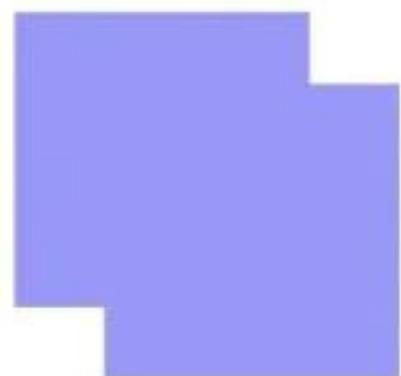
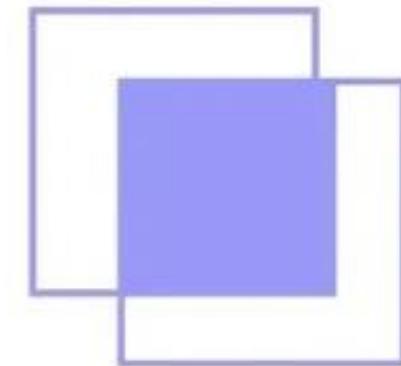


Excellent

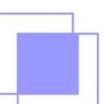


## ❖ IoU

```
1 def intersection_over_union(boxes_preds, boxes_labels, box_format="midpoint"):
2 >     """
3
4     # Check if the box format is "midpoint"
5     if box_format == "midpoint":...
6
7
8     # Check if the box format is "corners"
9     if box_format == "corners":...
10
11
12     # Calculate coordinates of the intersection rectangle
13     x1 = torch.max(box1_x1, box2_x1)
14     y1 = torch.max(box1_y1, box2_y1)
15     x2 = torch.min(box1_x2, box2_x2)
16     y2 = torch.min(box1_y2, box2_y2)
17
18     # Compute the area of the intersection rectangle, clamp(0) to handle cases where they do not overlap
19     intersection = (x2 - x1).clamp(0) * (y2 - y1).clamp(0)
20
21     # Calculate the areas of the predicted and ground truth boxes
22     box1_area = abs((box1_x2 - box1_x1) * (box1_y2 - box1_y1))
23     box2_area = abs((box2_x2 - box2_x1) * (box2_y2 - box2_y1))
24
25     # Calculate the Intersection over Union, adding a small epsilon to avoid division by zero
26     return intersection / (box1_area + box2_area - intersection + 1e-6)
```



$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



```
1 def non_max_suppression(bboxes, iou_threshold, confidence_threshold, box_format="corners"):
2     """
3         Perform Non-Maximum Suppression (NMS) to filter overlapping bounding boxes.
4
5     Args:
6         bboxes (list): List of bounding boxes, each containing [class_id, confidence, x1, y1, x2, y2]
7         iou_threshold (float): IoU threshold to determine box overlap
8         confidence_threshold (float): Minimum confidence score to keep a box
9         box_format (str): Format of box coordinates - "corners" or "midpoint"
10
11    Returns:
12        list: Filtered list of bounding boxes after NMS
13    """
14
15    if not isinstance(bboxes, list):
16        raise TypeError("bboxes must be a list")
17
18    # Keep only high confidence predictions
19    confident_boxes = [box for box in bboxes if box[1] > confidence_threshold]
20
21    # Sort by confidence score (highest first)
22    confident_boxes.sort(key=lambda x: x[1], reverse=True)
23
24    kept_boxes = []
25
26    while confident_boxes:
27        # Take box with highest confidence
28        best_box = confident_boxes.pop(0)
29
30        # Keep remaining boxes that either:
31        # 1. Are from a different class OR
32        # 2. Don't overlap significantly with the best box
33        confident_boxes = [
34            box for box in confident_boxes
35            if (box[0] != best_box[0]) or # Different class
36            (intersection_over_union(
37                torch.tensor(best_box[2:]),
38                torch.tensor(box[2:]),
39                box_format=box_format
40            ) < iou_threshold) # Low overlap
41        ]
42
43        kept_boxes.append(best_box)
44
45    return kept_boxes
```

## ❖ Non-Maximum Suppression (NMS)

### NMS Process:

#### 1. Initial Filtering (confidence > 0.5):

- All boxes pass (0.9, 0.8, 0.7, 0.85 > 0.5)

#### 2. Sorting by confidence:

```
bboxes = [  
    [0, 0.9, 0.2, 0.3, 0.4, 0.5],    # Highest confidence  
    [1, 0.85, 0.7, 0.7, 0.9, 0.9],  
    [0, 0.8, 0.21, 0.31, 0.41, 0.51],  
    [0, 0.7, 0.19, 0.29, 0.39, 0.49]  
]
```

#### 3. First Iteration:

- Choose first box (0.9 confidence)
- Compare with remaining boxes:
  - Box 2 (0.85): Different class (1 vs 0) → Keep
  - Box 3 (0.8): Same class, high IoU → Remove
  - Box 4 (0.7): Same class, high IoU → Remove

#### 4. Second Iteration:

- Only box with class 1 remains
- Added to final list

### Final Output:

```
bboxes_after_nms = [  
    [0, 0.9, 0.2, 0.3, 0.4, 0.5],    # Best box for class 0  
    [1, 0.85, 0.7, 0.7, 0.9, 0.9]    # Best box for class 1  
]
```

### NMS Benefits:

- Removes duplicate detections
- Keeps highest confidence predictions
- Handles multiple classes independently
- Improves precision by reducing false positives

## ❖ Loss

Localization loss

$$\begin{aligned}
 & \left[ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \right. \\
 & \quad \left. + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \right] \\
 & \quad \text{1 when there is object, } \text{0 when there is no object (in the box j of the cell i)} \\
 & \quad \text{Bounding box location (x, y) when there is object} \\
 & \quad \text{Bounding box size (w, h) when there is object} \\
 & \quad \left[ + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \right. \\
 & \quad \left. + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \right] \\
 & \quad \text{Confidence when there is object} \\
 & \quad \text{Confidence when there is no object (in the box j of the cell i)} \\
 & \quad \left[ + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \right] \\
 & \quad \text{Class probabilities when there is object (in the cell i)} \\
 & \quad \text{1 when there is no object, } \text{0 when there is object}
 \end{aligned}$$

huytranvan2010@gmail.com

```

1 class YoloLoss(nn.Module):
2 >     def __init__(self, S=7, B=2, C=20):...
19
20     def forward(self, predictions, target):
21         # Reshape the predictions to the shape (BATCH_SIZE, S*S(C+B*5))
22         predictions = predictions.reshape(-1, self.S, self.S, self.C + self.B * 5)
23
24         # Calculate Intersection over Union (IoU) for the two predicted bounding boxes
25         # with the target bounding box.
26         iou_b1 = intersection_over_union(predictions[..., 21:25], target[..., 21:25])
27         iou_b2 = intersection_over_union(predictions[..., 26:30], target[..., 21:25])
28         ious = torch.cat([iou_b1.unsqueeze(0), iou_b2.unsqueeze(0)], dim=0)
29
30         # Get the box with the highest IoU among the two predictions.
31         # Note that bestbox will have an index of 0 or 1, indicating which box is better.
32         iou_maxes, bestbox = torch.max(ious, dim=0)
33         exists_box = target[..., 20].unsqueeze(3) # This represents Iobj_i in the paper
34
35         # ===== #
36         # FOR BOX COORDINATES #
37         # ===== #
38
39         # Set the boxes with no objects to zero. Choose one of the two predictions
40         # based on the bestbox index calculated earlier.
41         box_predictions = exists_box * (
42             (
43                 bestbox * predictions[..., 26:30]
44                 + (1 - bestbox) * predictions[..., 21:25]
45             )
46         )
47         box_targets = exists_box * target[..., 21:25]
48
49         # Take the square root of width and height to ensure positive values.
50         box_predictions[..., 2:4] = torch.sign(box_predictions[..., 2:4]) * torch.sqrt(
51             torch.abs(box_predictions[..., 2:4] + 1e-6)
52         )
53         box_targets[..., 2:4] = torch.sqrt(box_targets[..., 2:4])
54
55         box_loss = self.mse(
56             torch.flatten(box_predictions, end_dim=-2),
57             torch.flatten(box_targets, end_dim=-2),
58         )

```

```

61         # ===== #
62         # FOR OBJECT LOSS #
63         # ===== #
64         # pred_box represents the confidence score of the box with the highest IoU.
65         pred_box = (
66             bestbox * predictions[..., 25:26] + (1 - bestbox) * predictions[..., 20:21]
67         )
68         object_loss = self.mse(
69             torch.flatten(exists_box * pred_box),
70             torch.flatten(exists_box * target[..., 20:21]),
71         )
72
73         # ===== #
74         # FOR NO OBJECT LOSS #
75         # ===== #
76         no_object_loss = self.mse(
77             torch.flatten((1 - exists_box) * predictions[..., 20:21], start_dim=1),
78             torch.flatten((1 - exists_box) * target[..., 20:21], start_dim=1),
79         )
80         no_object_loss += self.mse(
81             torch.flatten((1 - exists_box) * predictions[..., 25:26], start_dim=1),
82             torch.flatten((1 - exists_box) * target[..., 20:21], start_dim=1)
83         )
84
85         # ===== #
86         # FOR CLASS LOSS #
87         # ===== #
88         class_loss = self.mse(
89             torch.flatten(exists_box * predictions[..., :20], end_dim=-2),
90             torch.flatten(exists_box * target[..., :20], end_dim=-2),
91         )
92
93         # Calculate the final loss by combining the above components.
94         loss = (
95             self.lambda_coord * box_loss # First term
96             + object_loss # Second term
97             + self.lambda_noobj * no_object_loss # Third term
98             + class_loss # Fourth term
99         )
100
101         return loss

```

```

4  class Yolov1(nn.Module):
5      def __init__(self, in_channels=3, **kwargs):
6          """
7              Initializes the YOLOv1 model with ResNet34 as the backbone.
8
9              Args:
10                  in_channels (int): Number of input channels. Default is 3 for RGB images.
11                  **kwargs: Additional keyword arguments such as split_size, num_boxes, num_classes.
12
13      super(Yolov1, self).__init__()
14      self.split_size = kwargs.get("split_size", 7) # S
15      self.num_boxes = kwargs.get("num_boxes", 2) # B
16      self.num_classes = kwargs.get("num_classes", 20) # C
17
18      # Initialize ResNet34 backbone from timm
19      # `features_only=True` returns a list of feature maps from specified stages
20      # `out_indices=[4]` corresponds to the last layer's output
21      self.backbone = timm.create_model(
22          "resnet34",
23          pretrained=True,
24          features_only=True,
25          out_indices=[4],
26          in_chans=in_channels,
27      )
28
29      # Adaptive pooling to ensure the spatial dimensions match (S x S)
30      self.pool = nn.AdaptiveAvgPool2d((self.split_size, self.split_size))
31
32      # Fully connected layers
33      self.fcs = self._create_fcs()
34

```

```

35 @torch.autocast(device_type="cuda", dtype=torch.float16)
36 def forward(self, x):
37     """
38         Forward pass of the YOLOv1 model.
39
40         Args:
41             x (torch.Tensor): Input tensor of shape (N, C, H, W).
42
43         Returns:
44             torch.Tensor: Output tensor containing bounding box predictions.
45
46         """
47         # Extract features using ResNet34 backbone
48         features = self.backbone(x)[
49             0
50         ] # timm returns a list; take the first (and only) element
51
52         # Apply adaptive pooling to match spatial dimensions
53         pooled_features = self.pool(features)
54
55         # Pass through fully connected layers
56         output = self.fcs(pooled_features)
57
58         return output
59
60 def _create_fcs(self):
61     """
62         Creates the fully connected layers for YOLOv1.
63
64         Returns:
65             nn.Sequential: Sequential container of fully connected layers.
66
67         """
68         S, B, C = self.split_size, self.num_boxes, self.num_classes
69
70         # Calculate the input features for the first linear layer
71         # ResNet34's last feature map has 512 channels
72         # After pooling to (S x S), the total features are 512 * S * S
73         input_features = 512 * S * S
74
75         return nn.Sequential(
76             nn.Flatten(),
77             nn.Linear(input_features, 4096),
78             nn.Dropout(0.0),
79             nn.LeakyReLU(0.1),
80             nn.Linear(4096, S * S * (C + B * 5)),
81         )
82

```

