

Lab 6: Object Detection

YOLO, FCOS, and COCO Dataset Implementation

Artificial Vision (VIAR25/26)

UVigo

October 8, 2025

Section 1

Lab Overview and Objectives

Lab 6: What You'll Implement

Core Goal

Build a complete object detection system from scratch using PyTorch, implementing both YOLO-style and FCOS-style detectors, trained and evaluated on the COCO dataset.

Three Main Components

1. YOLO-Style Detector

- Grid-based detection with cell responsibility
- Multi-component loss function (localization + confidence + classification)
- Bounding box regression and class prediction

2. FCOS-Style Anchor-Free Detector

- Per-pixel prediction of (l, t, r, b) distances
- Center-ness score for quality estimation
- Feature Pyramid Network for multi-scale detection

3. Complete Training Pipeline

- COCO dataset loader with augmentation
- Focal loss for class imbalance
- Evaluation metrics (mAP, AP@IoU thresholds)
- TensorBoard logging and visualization

Learning Objectives

By the End of This Lab, You Will:

Understand:

- Difference between anchor-based (YOLO) and anchor-free (FCOS) detection
- How loss functions balance localization, confidence, and classification
- Impact of class imbalance and why focal loss helps
- Multi-scale detection via Feature Pyramid Networks

Implement:

- Complete detection pipeline from data loading to inference
- Non-maximum suppression for duplicate removal
- COCO evaluation metrics (AP, AP50, AP75, APS, APM, APL)
- Training loop with proper batch handling and gradient updates

Analyze:

- Performance trade-offs between architectures
- Impact of hyperparameters (grid size, focal loss γ , NMS threshold)
- Model behavior on different object scales

Section 2

Setup and Dataset Preparation

Environment Setup

Required Packages

```
1 # Create virtual environment
2 python -m venv venv_lab6
3 source venv_lab6/bin/activate # Linux/Mac
4 # venv_lab6\Scripts\activate # Windows
5
6 # Install dependencies
7 pip install torch torchvision
8 pip install pycocotools
9 pip install tensorboard
10 pip install matplotlib opencv-python
11 pip install tqdm numpy scipy
```

Download COCO Dataset

```
1 # Create data directory
2 mkdir -p data/coco
3
4 # Download 2017 Train/Val images (18GB train, 1GB val)
5 wget http://images.cocodataset.org/zips/train2017.zip
6 wget http://images.cocodataset.org/zips/val2017.zip
7
8 # Download annotations (241MB)
9 wget http://images.cocodataset.org/annotations/annotations_trainval2017.zip
10
11 # Extract
12 unzip train2017.zip -d data/coco/
13 unzip val2017.zip -d data/coco/
14 unzip annotations_trainval2017.zip -d data/coco/
```

Section 3

Part 1: COCO Dataset Implementation

Understanding COCO Annotations

COCO JSON Structure

```
1 {  
2   "images": [  
3     {  
4       "id": 139,  
5       "file_name": "000000000139.jpg",  
6       "height": 426,  
7       "width": 640  
8     }  
9   ],  
10  "annotations": [  
11    {  
12      "id": 1,  
13      "image_id": 139,  
14      "category_id": 1, # 1-indexed  
15      "bbox": [x, y, width, height], # XYWH format  
16      "area": 1234.5,  
17      "iscrowd": 0  
18    }  
19  ],  
20  "categories": [  
21    {  
22      "id": 1,  
23      "name": "person",  
24      "supercategory": "person"  
25    }  
26  ]  
27 }
```

Key Points

- 80 object categories (IDs 1–90 with gaps)
- Boxes are **XYWH**; convert to **XYXY** for IoU
- Multiple objects per image; some images have none
- `iscrowd=1` marks crowd regions (special handling)

Dataset Loader Implementation Tasks

COCODataset Class Requirements

Initialization:

- Load JSON annotations using `pycocotools`
- Create mapping from COCO category IDs (1-90) to contiguous IDs (0-79)
- Filter images without annotations (optional)

`__getitem__` method:

- Load image by ID
- Retrieve all annotations for that image
- Convert bounding boxes from XYWH to XYXY format
- Apply data augmentation (resize, flip, color jitter)
- Return: image tensor, boxes tensor, labels tensor

Collate function:

- Handle variable number of objects per image
- Stack images into batch
- Keep boxes and labels as lists (different lengths)

Expected Output Format

image: [3, H, W] tensor

boxes: [N, 4] tensor in XYXY format, normalized to [0, 1]

labels: [N] tensor with class IDs in range [0, 79]

Section 4

Part 2: YOLO Detector Implementation

YOLO Architecture Overview

Network Structure

Backbone: ResNet18/34/50

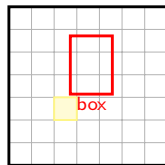
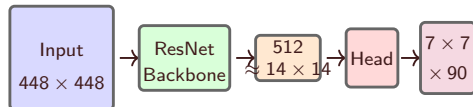
- Extract features: $\mathbb{R}^{3 \times H \times W} \rightarrow \mathbb{R}^{512 \times H/32 \times W/32}$

Detection Head:

- Conv layers to output grid: $S \times S \times (B \times 5 + C)$
- $S = 7$ (grid size), $B = 2$ (boxes per cell), $C = 80$ (classes)
- Output shape: $[7, 7, 90]$ for COCO

Output Interpretation:

- Each cell predicts B boxes
- Box: $(x, y, w, h, \text{conf}) + \text{class probs}$
- x, y : offsets relative to cell (0-1)
- w, h : relative to image (0-1)



$S \times S$ Grid ($S=7$)

YOLO Loss Function Components

Multi-Part Loss: $\mathcal{L} = \lambda_{coord}\mathcal{L}_{coord} + \mathcal{L}_{conf}^{obj} + \lambda_{noobj}\mathcal{L}_{conf}^{noobj} + \mathcal{L}_{class}$

1. Localization Loss (\mathcal{L}_{coord}):

$$\mathcal{L}_{coord} = \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{obj} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 + (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right]$$

2. Confidence Loss (Object):

$$\mathcal{L}_{conf}^{obj} = \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{obj} (C_i - \hat{C}_i)^2$$

3. Confidence Loss (No Object):

$$\mathcal{L}_{conf}^{noobj} = \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{I}_{ij}^{noobj} (C_i - \hat{C}_i)^2$$

4. Classification Loss:

$$\mathcal{L}_{class} = \sum_{i=0}^{S^2} \mathbb{I}_i^{obj} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2$$

Typical Hyperparameters

$\lambda_{coord} = 5.0$, $\lambda_{noobj} = 0.5$, $S = 7$, $B = 2$

YOLO Implementation Tasks

Step 1: Encode Ground Truth to Grid

For each bounding box:

- Find which grid cell contains the center
- Compute offsets (x, y) relative to cell $(0-1)$
- Set target confidence = 1 for responsible box
- Set class probabilities for that cell

Step 2: Forward Pass

```
1 # Input: [B, 3, 448, 448]
2 # Backbone output: [B, 512, 14, 14]
3 # Detection head output: [B, S, S, B*5+C]
4 # Reshape to: [B, S, S, B, 5+C]
5 # where last dim = [x, y, w, h, conf, class_probs...]
```

Step 3: Compute Loss

- Determine responsible boxes (highest IoU with GT)
- Apply $\mathbb{1}_{ij}^{obj}$ and $\mathbb{1}_{ij}^{noobj}$ masks
- Calculate each loss component
- Weight and sum: \mathcal{L}_{total}

Section 5

Part 3: FCOS Anchor-Free Detector

FCOS Architecture and Predictions

Per-Pixel Predictions

At each feature map location (i, j) :

Classification: C scores

$$\mathbf{p}_{ij} \in \mathbb{R}^C$$

Regression: (l, t, r, b) distances

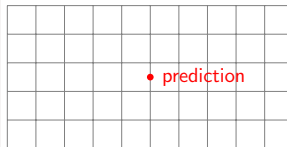
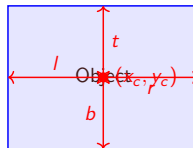
$$\mathbf{r}_{ij} = (l_{ij}, t_{ij}, r_{ij}, b_{ij}) \in \mathbb{R}^4$$

Center-ness: Quality score

$$s_{ij} = \sqrt{\frac{\min(l, r)}{\max(l, r)} \times \frac{\min(t, b)}{\max(t, b)}}$$

Key Difference from YOLO

- No grid-based assignment
- Direct per-pixel prediction
- Multi-scale via FPN levels



Feature Map

Center-ness Intuition:

- High near object center
- Low at boundaries
- Suppresses low-quality boxes

Feature Pyramid Network (FPN) for Multi-Scale

FPN Construction

Bottom-up pathway: ResNet stages

C_2, C_3, C_4, C_5 with strides 4, 8, 16, 32

Top-down pathway: Lateral connections

$$P_5 = \text{Conv}_{1 \times 1}(C_5)$$

$$P_4 = \text{Conv}_{3 \times 3}(\text{Upsample}(P_5) + \text{Conv}_{1 \times 1}(C_4))$$

$$P_3 = \text{Conv}_{3 \times 3}(\text{Upsample}(P_4) + \text{Conv}_{1 \times 1}(C_3))$$

Scale assignment:

- P_3 : objects with $\max(l, t, r, b) \in [0, 64]$ pixels
- P_4 : objects with $\max(l, t, r, b) \in [64, 128]$ pixels
- P_5 : objects with $\max(l, t, r, b) > 128$ pixels

Implementation Note

Each FPN level has independent classification and regression heads (shared weights across spatial locations).

FCOS Loss Function

Three Components

1. Classification Loss (Focal Loss):

$$\mathcal{L}_{cls} = -\frac{1}{N_{pos}} \sum_{x,y} \alpha_t (1 - p_t)^\gamma \log(p_t)$$

where $\gamma = 2.0$, $\alpha = 0.25$

2. Regression Loss (GloU):

$$\mathcal{L}_{reg} = \frac{1}{N_{pos}} \sum_{x,y} \mathbb{I}_{\{c_{x,y}^* > 0\}} \mathcal{L}_{GloU}(b_{x,y}, b_{x,y}^*)$$

3. Center-ness Loss (BCE):

$$\mathcal{L}_{centerness} = \frac{1}{N_{pos}} \sum_{x,y} \mathbb{I}_{\{c_{x,y}^* > 0\}} \text{BCE}(s_{x,y}, s_{x,y}^*)$$

Total Loss:

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda_1 \mathcal{L}_{reg} + \lambda_2 \mathcal{L}_{centerness}$$

Typical: $\lambda_1 = 1.0$, $\lambda_2 = 1.0$

FCOS Implementation Tasks

Step 1: Assign Ground Truth to Locations

```
1 for each ground truth box:
2     for each FPN level P_k:
3         for each location (x, y) in P_k:
4             # Compute (l, t, r, b) from location to box
5             l = x - bbox.x1
6             t = y - bbox.y1
7             r = bbox.x2 - x
8             b = bbox.y2 - y
9
10            # Check if inside box and within scale range
11            if l > 0 and t > 0 and r > 0 and b > 0:
12                if min_scale < max(l,t,r,b) < max_scale:
13                    # Assign as positive sample
14                    targets[k][y, x] = (l, t, r, b, class)
```

Step 2: Compute Center-ness Targets

```
1 centerness = sqrt(
2     (min(l, r) / max(l, r)) *
3     (min(t, b) / max(t, b))
4 )
```

Section 6

Part 4: Loss Functions and Training

Focal Loss Implementation

Purpose: Address Class Imbalance

In object detection, most locations are background (negative samples). Focal loss down-weights easy examples and focuses on hard ones.

Mathematical Formulation

Standard Cross-Entropy:

$$CE(p_t) = -\log(p_t)$$

Focal Loss:

$$FL(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

where:

- $p_t = p$ if $y = 1$ (positive), else $p_t = 1 - p$
- $\gamma \geq 0$ is focusing parameter (typical: 2)
- α_t is class weighting (typical: 0.25 for positive)

Behavior

- When $p_t \rightarrow 1$ (easy example): $(1 - p_t)^\gamma \rightarrow 0$, loss $\rightarrow 0$
- When $p_t \rightarrow 0$ (hard example): $(1 - p_t)^\gamma \rightarrow 1$, loss remains high
- γ controls the rate of down-weighting

Non-Maximum Suppression (NMS)

Purpose

Remove duplicate detections for the same object by keeping only the highest confidence box among overlapping predictions.

Algorithm

```
1 def nms(bboxes, scores, iou_threshold=0.5):
2     # Sort by confidence descending
3     indices = scores.argsort(descending=True)
4     keep = []
5
6     while len(indices) > 0:
7         # Keep highest confidence box
8         current = indices[0]
9         keep.append(current)
10
11        # Compute IoU with remaining boxes
12        ious = compute_iou(bboxes[current], bboxes[indices[1:]])
13
14        # Remove overlapping boxes
15        mask = ious <= iou_threshold
16        indices = indices[1:][mask]
17
18    return keep
```

Key Parameters

`iou_threshold`: Higher \Rightarrow more boxes kept (default: 0.5)

Training Pipeline Structure

Main Training Loop

```
1 for epoch in range(num_epochs):
2     model.train()
3
4     for images, boxes, labels in train_loader:
5         # Forward pass
6         predictions = model(images)
7
8         # Compute loss
9         loss = compute_loss(predictions, boxes, labels)
10
11        # Backward pass
12        optimizer.zero_grad()
13        loss.backward()
14        optimizer.step()
15
16        # Logging
17        if step % log_interval == 0:
18            writer.add_scalar('train/loss', loss.item(), step)
19
20    # Validation
21    if epoch % val_interval == 0:
22        metrics = evaluate(model, val_loader)
23        writer.add_scalar('val/mAP', metrics['mAP'], epoch)
24
25    # Save checkpoint
26    save_checkpoint(model, optimizer, epoch)
```

Section 7

Part 5: Evaluation Metrics

COCO Evaluation Metrics

Primary Metrics

AP (Average Precision): Mean AP over IoU thresholds [0.5:0.05:0.95]

$$\text{mAP} = \frac{1}{10} \sum_{t=0.5}^{0.95} \text{AP}_t$$

AP50: AP at IoU threshold 0.5 (PASCAL VOC metric)

AP75: AP at IoU threshold 0.75 (stricter localization)

Scale-Specific Metrics

- **APS:** Small objects (area < 32² pixels)
- **APM:** Medium objects (32² ≤ area < 96²)
- **APL:** Large objects (area ≥ 96²)

Implementation

Use `pycocotools.cocoeval.COCOeval` for official evaluation:

- Matches predictions to ground truth via IoU
- Computes precision-recall curves
- Aggregates across categories and IoU thresholds

Evaluation Code Structure

Evaluation Function

```
1 def evaluate(model, data_loader, device):
2     model.eval()
3     coco_results = []
4
5     with torch.no_grad():
6         for images, image_ids in data_loader:
7             # Forward pass
8             predictions = model(images.to(device))
9
10            # Post-process: NMS
11            detections = post_process(predictions)
12
13            # Convert to COCO format
14            for det, img_id in zip(detections, image_ids):
15                for box, score, label in det:
16                    coco_results.append({
17                        'image_id': img_id,
18                        'category_id': label,
19                        'bbox': [x, y, w, h], # XYWH
20                        'score': score
21                    })
22
23    # Evaluate using COCO API
24    coco_eval = COCOeval(coco_gt, coco_dt, 'bbox')
25    coco_eval.evaluate()
26    coco_eval.accumulate()
27    coco_eval.summarize()
28
29    return {
30        'mAP': coco_eval.stats[0],
31        'AP50': coco_eval.stats[1],
32        'AP75': coco_eval.stats[2]
33    }
```

Section 8

Implementation Roadmap

Step-by-Step Implementation Guide

Phase 1: Data Pipeline (2-3 hours)

- 1 Implement COCODataset class with annotation parsing
- 2 Add data augmentation (resize, flip, normalize)
- 3 Create custom collate function for variable-length boxes
- 4 Test: Visualize loaded images with ground truth boxes

Phase 2: YOLO Detector (3-4 hours)

- 1 Build YOLO model with ResNet backbone
- 2 Implement grid-based target encoding
- 3 Write YOLO loss function with all components
- 4 Test: Forward pass and loss computation

Phase 3: FCOS Detector (3-4 hours)

- 1 Implement FPN for multi-scale features
- 2 Build FCOS heads (classification, regression, centerness)
- 3 Write focal loss and GloU loss
- 4 Implement target assignment logic

Step-by-Step Implementation Guide (cont.)

Phase 4: Training Pipeline (2-3 hours)

- 1 Set up training loop with optimizer and scheduler
- 2 Add TensorBoard logging for losses and metrics
- 3 Implement checkpoint saving and loading
- 4 Start training and monitor convergence

Phase 5: Evaluation and Post-Processing (2-3 hours)

- 1 Implement NMS for inference
- 2 Write evaluation function using COCO API
- 3 Compute mAP, AP50, AP75, and scale-specific metrics
- 4 Visualize predictions on validation set

Phase 6: Analysis and Report (2 hours)

- 1 Compare YOLO vs FCOS performance
- 2 Analyze scale-specific performance (APS, APM, APL)
- 3 Experiment with hyperparameters (learning rate, loss weights)
- 4 Write report with findings and visualizations

Section 9

Tips and Common Pitfalls

Implementation Tips

Data Loading

- **Memory:** Use `num_workers > 0` in `DataLoader` for parallel loading
- **Normalization:** Apply ImageNet mean/std for pretrained backbones
- **Augmentation:** Strong augmentation helps (flip, color jitter, crop)

Training Stability

- **Gradient clipping:** Use `torch.nn.utils.clip_grad_norm_()` (max norm = 10)
- **Warmup:** Start with low learning rate for first few epochs
- **Loss weighting:** Balance localization vs classification losses
- **Batch size:** Larger is better (8-16 on GPUs), use gradient accumulation if needed

Debugging

- **Overfitting check:** Train on small subset (10 images) first
- **Visualization:** Plot predictions vs ground truth frequently
- **NaN losses:** Check for division by zero, use `torch.clamp()`
- **IoU computation:** Ensure boxes are in correct format (XYXY vs XYWH)

Common Pitfalls and Solutions

Problem: Low mAP ($\downarrow 10\%$)

Possible Causes:

- Incorrect target encoding (grid assignment, box format)
- Wrong loss weighting (confidence loss too dominant)
- NMS threshold too aggressive

Solutions: Visualize targets vs predictions, check loss components individually

Problem: Training Doesn't Converge

Possible Causes:

- Learning rate too high (try $1e-4$ instead of $1e-3$)
- Missing gradient clipping
- Imbalanced loss components

Solutions: Reduce LR, add warmup, check loss curves in TensorBoard

Problem: Good Training Loss, Poor Validation

Cause: Overfitting or domain shift

Solutions: More augmentation, dropout, smaller model, regularization

Expected Results and Benchmarks

Target Performance (after 20 epochs)

YOLO (Grid 7×7 , ResNet18 backbone):

- mAP: 15-20% (lightweight, fast)
- AP50: 30-35%
- Inference: 30-40 FPS on GPU

FCOS (FPN, ResNet50 backbone):

- mAP: 25-30% (better than YOLO)
- AP50: 42-48%
- APS: 10-15% (improved small object detection)
- Inference: 15-20 FPS on GPU

Note

These are expected results with limited training. State-of-the-art detectors achieve 40-50+ mAP with longer training, larger backbones, and advanced techniques.

GPU Requirements

- Minimum: 6GB VRAM (batch size 4-8)
- Recommended: 11GB+ VRAM (batch size 16)
- CPU training possible but 10-20 \times slower

Section 10

Deliverables and Grading

Required Submissions

1. Code (60%)

- Complete implementation of YOLO and FCOS detectors
- COCO dataset loader with augmentation
- Training and evaluation scripts
- All code must run without errors

2. Trained Models (20%)

- Checkpoint files for both YOLO and FCOS
- Training logs (TensorBoard events)
- Minimum 15 epochs of training

3. Report (20%)

- Architecture descriptions with diagrams
- Training curves (loss, mAP over epochs)
- Evaluation results table (mAP, AP50, AP75, APS, APM, APL)
- Comparison analysis: YOLO vs FCOS
- Example predictions (5-10 images with visualizations)
- Discussion of failures and potential improvements

Grading Rubric

Detailed Breakdown

Code Quality (60 points):

- Dataset implementation: 10 pts
- YOLO model and loss: 15 pts
- FCOS model and loss: 15 pts
- Training pipeline: 10 pts
- Evaluation and NMS: 10 pts

Model Performance (20 points):

- YOLO mAP > 15%: 8 pts (partial credit for 10-15%)
- FCOS mAP > 25%: 12 pts (partial credit for 20-25%)

Report Quality (20 points):

- Completeness and clarity: 8 pts
- Results presentation: 6 pts
- Analysis depth: 6 pts

Bonus Points (up to 10%)

Additional features: Soft-NMS, DIoU loss, data augmentation experiments, etc.

Section 11

Resources and Support

Helpful Resources

Documentation

- **PyTorch**: <https://pytorch.org/docs/>
- **COCO API**: <https://github.com/cocodataset/cocoapi>
- **Torchvision**: <https://pytorch.org/vision/stable/index.html>

Papers

- **YOLO**: Redmon et al., "You Only Look Once" (CVPR 2016)
- **FCOS**: Tian et al., "FCOS: Fully Convolutional One-Stage" (ICCV 2019)
- **FPN**: Lin et al., "Feature Pyramid Networks" (CVPR 2017)
- **Focal Loss**: Lin et al., "Focal Loss for Dense Object Detection" (ICCV 2017)

Reference Implementations

- Ultralytics YOLO: <https://github.com/ultralytics/yolov5>
- MMDetection: <https://github.com/open-mmlab/mmdetection>
- Detectron2: <https://github.com/facebookresearch/detectron2>