

Task-4

YOLO(You only look once)

Its distinctive approach, examining the entire image just once to identify objects and their positions.

In conventional methods employing two stage detection processes But yolo treats object detection as a regression problem

Like in conventional object detection models used to work in two stages like find the object locations like drawing bounding boxed and classification with refinement like to classify what's in each region so this method was slow but

Regression means predicting continuous numerical values like every real value is possible.

What actually CNN is ?

It is like a robot that looks at an image piece by piece understands patterns like edges , corners , eyes ,faces and uses that to recognize what's in the image ..

Layer 1:

"Edge Detector" :- Looks for simple things like straight lines , curves ...etc.

Layer 2

"Shape detector" use those edges to recognize basic shapes.

Layer 3:

part detector: recognizes the part based from above

Mixes all the above ...

How YOLO divides the image into grids

It expects a fixed - size input

Then the convolutional network processes the image the image is turned into feature map

Like if we take an image the kernel/filter first takes some part of image and squeezes it to get the output and similarly it operates patch wise to give the complete output

How YOLO creates a bounding box

Yolo divides the image into an $S \times S$ grid and each grid cell is responsible for detecting objects whose center falls inside it.

Each grid cell doesn't just predict one box ,it predicts multiple boxes using anchor boxes

For each anchor box YOLO predicts offset for x-center , offset for y-center , Log space scale for width and log-space scale for height ,confidence and probabilities for each class

In the YOLO paradigm, a single convolutional neural network is employed to predict bounding boxes and class probabilities for an entire image

The YOLO architecture introduced the end-to-end, differentiable approach to object detection by **unifying the tasks of bounding box regression and object classification into a single neural network**

Core components	Description
Backbone	CNN is responsible for encoding the image information into feature maps at varying scales
Neck	A series of layers designed to integrate and refine feature representations
Head	Generates predictions for object bounding boxes and class labels based on processed

Core components	Description
	features

The two primary training techniques employed in YOLOv5 :

Data Augmentation:

Data augmentation is the process of **creating new, slightly modified copies of your existing data** to help a model **see more variety** and generalize better

Some common augmentations like flip, rotation , zoom/crop , Color jitter , translate , Resize/stretch ,Adding Noise like that so on....

Loss function:

It is a composite metric

A

loss function measures **how bad or wrong** the model's predictions are compared to the actual correct answers.

In regression we have seen MSE and RMSE are loss functions in regression (continuous)...

GIoU(Generalized Intersection over Union):

One of the component of loss function

It measures how well the predicted bounding box overlaps with the true bounding box (the ground truth)

Like in normal IoU it measures overlap but is zero if boxes don't overlap.

GIoU extends it to give meaningful feedback even when boxes don't overlap. so the model can still learn to improve box placement

Minimizing the difference guarantees better box localization

Objectness Loss

This tells how confident the model is that an object exists inside the predicted bounding box

Classification Loss

It measures how well the model predicts the correct class inside the bounding box.

$$Totalloss = GIoULoss + ObjectnessLoss + ClassificationLoss$$

Data Augmentation

scaling - Adjustments to image size

Color space manipulation - Modifications to color channels

Mosaic augmentation - A novel technique that combines four images into four randomly sized tiles

Mosaic data augmentation has proven particularly effective in addressing the challenge of small object detection

Mosaic augmentation, initially introduced in the YOLOv3 PyTorch repository and later integrated into YOLOv5, has shown significant effectiveness in improving the detection of small objects. By combining multiple images into one, this technique increases the diversity of object scales and spatial configurations, enhancing the model's ability to generalize across object sizes and improve accuracy on smaller targets.

Bounding Box Anchors

YOLOv3 introduced a new approach for generating anchor boxes using **K-means clustering** combined with **genetic algorithms**. This technique derives anchor box dimensions from the actual distribution of bounding boxes within the training dataset, rather than relying on generic settings such as those from COCO. This is especially beneficial for custom object detection tasks where object sizes and aspect ratios may differ significantly.

YOLOv5 uses anchor boxes to initialize bounding box predictions, computing final coordinates as offsets from these anchors.

Loss Calculation

The YOLOv5 loss function consists of three key components:

1. **Class Prediction Loss (L_{cls})** – Binary Cross-Entropy (BCE)
2. **Objectness Loss (L_{obj})** – Binary Cross-Entropy (BCE)
3. **Localization Loss (L_{loc})** – Complete Intersection over Union (CIoU)

Bit Floating Point Precision

YOLOv5 supports 16-bit floating point (FP16) precision, a feature provided by the PyTorch framework. This allows for faster training and inference by reducing the memory and computational footprint compared to standard 32-bit precision.

CSP Backbone

The CSP architecture addresses redundancy in gradient flow by splitting feature maps into two paths: one processed through a dense block and another passed directly to the next stage. This approach:

- Reduces model complexity,
- Decreases computational cost,
- Maintains or improves model performance.

Inspired by **DenseNet**, CSP promotes enhanced feature propagation and reusability, while also mitigating the vanishing gradient problem.

YOLOv5 variants use **CSPResNext50** and **CSPDarknet53** as backbones, which build on DenseNet by preserving part of the feature maps in an unaltered form, enabling better learning and lower computational overhead.

PA-Net Neck

YOLOv5 employs the **Path Aggregation Network (PA-Net)** for feature aggregation, following its successful use in YOLOv4. PA-Net enhances information flow between feature levels and helps in accurate object localization at different scales.

YOLOv5 Models

The YOLOv5 architecture offers a family of five model variants, designed to balance accuracy, speed, and computational resources. These models range from lightweight implementations suitable for edge computing to high-capacity models optimized for maximum accuracy.

Model	Full Name	Parameters	Size (Approx.)	Key Features
YOLOv5n	Nano	1.9M	~2.5 MB (INT8), ~4 MB (FP32)	Optimized for edge devices and IoT; fastest and most lightweight; compatible with OpenCV DNN.
YOLOv5s	Small	7.2M	–	Baseline model; balances performance and speed; ideal for CPU-based inference.
YOLOv5m	Medium	21.2M	–	Mid-range model; offers higher accuracy with manageable computational load.
YOLOv5l	Large	46.5M	–	Designed for high-precision tasks; better at detecting small or dense objects.
YOLOv5x	Extra Large	86.7M	–	Highest mAP and accuracy; suitable for scenarios where computational resources are not constrained.

5. YOLOv5 Annotation Format

YOLOv5 uses a format similar to the traditional YOLO Darknet TXT format, but enhanced with support from a YAML configuration file. The annotations consist of .txt files where each line represents a bounding box in the format:

All coordinates are normalized between 0 and 1 relative to the image dimensions.

Additionally, YOLOv5 requires a .yaml file to define:

- Class labels
- Number of classes
- Paths to the training and validation datasets

Annotation Tool Compatibility:

To facilitate dataset preparation, various tools support export to YOLOv5 format:

Annotation tools

Tool	Main Use
Roboflow	Image labeling, preprocessing, export
Labellmg	Manual bounding box drawing
CVAT	Advanced annotation (web UI)
VOTT	Tagging images for object detection

Integration platforms

Platform	Purpose	Type
Deci	Model compilation and quantization for faster inference	Optimization
ClearML	Experiment tracking and remote training	Training/Tracking
Roboflow	Data annotation + export + dataset hosting	Annotation + Integration
Weights and Biases	Training monitoring, visualization, experiment comparison	Logging/Tracking

I used chatgpt for the above table info

These tools may require intermediate conversions to ensure full compatibility with the expected YOLOv5 format.

Efficient dataset labeling is essential for supervised learning workflows. YOLOv5 supports integration with several popular platforms that streamline annotation, training, and model deployment.

Key innovations:

1. Advances in Architecture:

Efficiency and detection accuracy are greatly increased by combining the **CSP backbone** and **PA-Net neck**. While PA-Net improves multi-scale feature fusion, CSP minimizes redundant gradients.

2. Versatility of Model:

Models in the YOLOv5 family are optimized for a variety of platforms, including microcontrollers and GPUs, allowing for smooth deployment in cloud, edge, and mobile settings.

3. Creative Training Approaches:

Methods such as **mosaic augmentation** increase the strength to data variation and the detection of small objects. Support for **FP16 precision** speeds up inference on compatible hardware and uses less memory.