# Paper Review: You Only Look Once (YOLO)
## Unified, Real-Time Object Detection

Review of Redmon et al., 2016

January 6, 2026

**Abstract**

This document reviews the seminal work "You Only Look Once: Unified, Real-Time Object Detection" by Joseph Redmon et al. (2016). YOLO revolutionized object detection by reframing it as a single regression problem, processing images at 45 frames per second while achieving 63.4% mAP on PASCAL VOC. Unlike prior methods that repurpose classifiers through complex pipelines (R-CNN, DPM), YOLO uses a unified convolutional neural network architecture that simultaneously predicts bounding boxes and class probabilities directly from full images in one evaluation. This review examines the abstract, detailed architecture, and key innovations that distinguish YOLO from older detection models.

## 1 Introduction

Object detection has been a fundamental challenge in computer vision, with applications ranging from autonomous driving to assistive technologies. Prior to YOLO, state-of-the-art detection systems relied on multi-stage pipelines that were slow and difficult to optimize. YOLO introduced a paradigm shift by treating detection as a unified regression problem.

## 2 Paper Abstract Summary

The original YOLO paper presents three major contributions:

> **Key Contributions**
>
> 1. **Unified Architecture**: Reframes object detection as a single regression problem to spatially separated bounding boxes and class probabilities
>
> 2. **Real-Time Performance**: Base YOLO processes images at 45 FPS, Fast YOLO at 155 FPS while achieving double the mAP of other real-time detectors
>
> 3. **Generalization**: Learns general object representations that outperform DPM and R-CNN when applied to new domains like artwork

### 2.1 Performance Characteristics

- **Speed**: 45 FPS (YOLO) vs 0.5 FPS (Fast R-CNN)

- **Accuracy**: 63.4% mAP on PASCAL VOC 2007+2012

- **Error Profile**: More localization errors but significantly fewer background false positives (13.6% for Fast R-CNN vs 4.75% for YOLO)

# 3 YOLO Architecture

## 3.1 Overall Design Philosophy

YOLO divides the input image into an $S \times S$ grid. If the center of an object falls into a grid cell, that grid cell is responsible for detecting that object. For PASCAL VOC evaluation, YOLO uses $S = 7$.
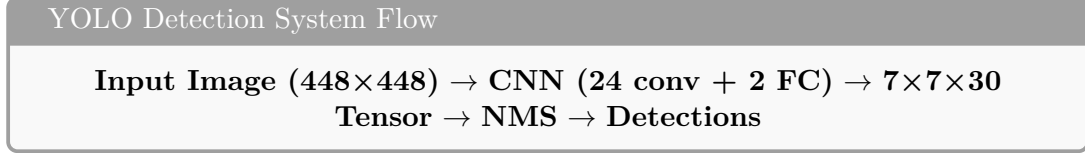
> **YOLO Detection System Flow**
>
> **Input Image (448×448) → CNN (24 conv + 2 FC) → 7×7×30 Tensor → NMS → Detections**

Figure 1: YOLO processes an image through a single neural network to predict all detections simultaneously

## 3.2 Network Architecture Details

The YOLO network consists of:

- **24 Convolutional Layers**: Extract features from the image

- **2 Fully Connected Layers**: Predict output probabilities and coordinates

- **Inspired by GoogLeNet**: Uses $1 \times 1$ reduction layers followed by $3 \times 3$ convolutional layers

Table 1: YOLO Network Layer Specifications

| Layer Type | Filter/Stride | Input Size | Output Size |
|---|---|---|---|
| Convolutional | $7 \times 7 \times 64$, s=2 | $448 \times 448 \times 3$ | $224 \times 224 \times 64$ |
| Max Pooling | $2 \times 2$, s=2 | $224 \times 224 \times 64$ | $112 \times 112 \times 64$ |
| Convolutional | $3 \times 3 \times 192$ | $112 \times 112 \times 64$ | $112 \times 112 \times 192$ |
| Max Pooling | $2 \times 2$, s=2 | $112 \times 112 \times 192$ | $56 \times 56 \times 192$ |
| Conv Layers ($\times 8$) | $1 \times 1$ and $3 \times 3$ | $56 \times 56 \times 256$ | $28 \times 28 \times 512$ |
| Conv Layers ($\times 4$) | $1 \times 1$ and $3 \times 3$ | $28 \times 28 \times 512$ | $14 \times 14 \times 1024$ |
| Conv Layers ($\times 2$) | $3 \times 3 \times 1024$ | $14 \times 14 \times 1024$ | $7 \times 7 \times 1024$ |
| Fully Connected | - | $7 \times 7 \times 1024$ | 4096 |
| Fully Connected | - | 4096 | $7 \times 7 \times 30$ |

## 3.3 Output Tensor Structure

Each grid cell predicts:

- $B = 2$ **bounding boxes**, each with 5 predictions: $(x, y, w, h, \text{confidence})$

- $C = 20$ **conditional class probabilities** (for PASCAL VOC)

**Final prediction tensor**: $7 \times 7 \times 30$

$$\text{Tensor size} = S \times S \times (B \times 5 + C) = 7 \times 7 \times (2 \times 5 + 20) = 7 \times 7 \times 30 \tag{1}$$

### 3.4 Prediction Components

#### 3.4.1 Bounding Box Predictions

Each bounding box consists of 5 predictions:

- $(x, y)$: Center coordinates relative to grid cell bounds

- $(w, h)$: Width and height relative to whole image

- Confidence: $\Pr(\text{Object}) \times \text{IOU}_{\text{pred}}^{\text{truth}}$

#### 3.4.2 Class Probability

Each grid cell predicts $C$ conditional class probabilities:

$$\Pr(\text{Class}_i | \text{Object}) \tag{2}$$

At test time, these are combined:

$$\boxed{\Pr(\text{Class}_i | \text{Object}) \times \Pr(\text{Object}) \times \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) \times \text{IOU}_{\text{pred}}^{\text{truth}}} \tag{3}$$

## 4 Key Differentiating Features from Older Models

### 4.1 1. Unified Detection vs. Multi-Stage Pipelines

> **YOLO's Single-Stage Approach**
>
> **Older Methods (R-CNN, DPM):**
>
> - Region proposal generation (Selective Search)
>
> - Feature extraction for each proposal
>
> - SVM classification
>
> - Bounding box regression
>
> - Non-max suppression
>
> **YOLO:**
>
> - Single convolutional network pass
>
> - Simultaneous prediction of all boxes and classes
>
> - Non-max suppression

**Impact**: This enables end-to-end optimization and eliminates the need to train separate components.

## 4.2  2. Global Context Reasoning

Table 2: Comparison of Detection Approaches

| Aspect | Older Methods | YOLO |
|---|---|---|
| Image View | Local patches/regions | Entire image |
| Context | Limited to proposal region | Global scene context |
| Background Errors | 13.6% (Fast R-CNN) | 4.75% |
| Speed | 0.5-6 FPS | 45-155 FPS |
| Proposals | 2000 (Selective Search) | 98 (grid-based) |

**Key Insight**: By seeing the entire image during training and test time, YOLO implicitly encodes contextual information about classes and their appearance, reducing false positives on background regions.

## 4.3  3. Direct Detection as Regression

> **Detection as Regression Problem**
>
> **Traditional Approach:**
>
> - Repurpose image classifier for detection
> - Sliding window or region proposals
> - Multiple evaluations per image
>
> **YOLO Approach:**
>
> - Direct regression from image pixels to bounding box coordinates
> - Single network evaluation
> - Optimized directly for detection performance

## 4.4  4. Loss Function: Multi-Part Optimization

YOLO uses a sum-squared error loss with weighted components:

$$
\begin{aligned}
\mathcal{L} =& \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
& + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
& + \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\
& + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^{B} \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\
& + \sum_{i=0}^{S^2} \mathbb{1}_{i}^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2
\end{aligned}
\tag{4}
$$

Where:

- $\lambda_{\text{coord}} = 5$: Increases loss from bounding box coordinates
- $\lambda_{\text{noobj}} = 0.5$: Decreases loss from confidence predictions for boxes without objects
- $\mathbb{1}_{ij}^{\text{obj}}$: Indicator if object appears in cell $i$ and box $j$ is responsible
- Square root for $(w, h)$: Makes small deviations in large boxes matter less

## 4.5  5. Speed-Accuracy Trade-off

**Performance Comparison on PASCAL VOC 2007:**

- **Fast YOLO**: 155 FPS, 52.7% mAP
- **YOLO**: 45 FPS, 63.4% mAP
- **Fast R-CNN**: 0.5 FPS, 70.0% mAP
- **Faster R-CNN (VGG-16)**: 7 FPS, 73.2% mAP
- **DPM**: 30 FPS, 26.1% mAP

Figure 2: YOLO achieves real-time performance with competitive accuracy

**Analysis**: YOLO is the only method operating at true real-time speeds (>30 FPS) while maintaining high accuracy. It trades some localization precision for speed and reduced background errors.

## 4.6  6. Activation Function

YOLO uses a **leaky rectified linear activation** for all layers except the final layer:

$$\phi(x) = \begin{cases} x, & \text{if } x > 0 \\ 0.1x, & \text{otherwise} \end{cases} \tag{5}$$

The final layer uses a **linear activation function** to predict bounding box coordinates and probabilities.

## 4.7  7. Training Strategy

1. **Pretraining**: First 20 convolutional layers trained on ImageNet at $224 \times 224$ resolution
2. **Detection Fine-tuning**: Add 4 convolutional + 2 fully connected layers, increase resolution to $448 \times 448$
3. **Data Augmentation**: Random scaling and translations up to 20%, exposure and saturation adjustments
4. **Regularization**: Dropout (rate = 0.5) after first connected layer
5. **Learning Rate Schedule**:
   - Epochs 1-75: $10^{-2}$
   - Epochs 76-105: $10^{-3}$
   - Epochs 106-135: $10^{-4}$

# 5  Comparison with Previous Detection Systems

## 5.1  vs. Deformable Parts Models (DPM)

| Feature | DPM | YOLO |
|---|---|---|
| Approach | Sliding window | Grid-based regression |
| Pipeline | Disjoint (feature extraction, classification, localization separate) | Unified single network |
| Features | Static (HOG) | Learned end-to-end |
| Speed | 15-30 FPS | 45-155 FPS |
| Accuracy (VOC 2007) | 30.4% mAP | 63.4% mAP |

Table 3: YOLO vs. Deformable Parts Models

## 5.2  vs. R-CNN Family

**R-CNN Pipeline:**

1. Selective Search: Generate 2000 region proposals ($\sim$2 seconds)

2. CNN: Extract features for each proposal

3. SVM: Classify each region

4. Regression: Adjust bounding boxes

5. NMS: Eliminate duplicates

**Time**: >40 seconds per image

**Fast R-CNN Improvements:**

- Shares computation across proposals

- Still relies on Selective Search

- 0.5 FPS ($\sim$2 seconds per image)

**Faster R-CNN:**

- Replaces Selective Search with Region Proposal Network

- 7-18 FPS depending on base network

- Still slower than YOLO

## 5.3 Error Analysis: YOLO vs. Fast R-CNN

Table 4: Error Type Distribution (Top N Detections)

| Error Type | Fast R-CNN | YOLO |
|---|---|---|
| Correct | 71.6% | 65.5% |
| Localization Error | 8.6% | 19.0% |
| Similar Class | 4.3% | 6.75% |
| Other Class | 1.9% | 4.0% |
| Background | **13.6%** | **4.75%** |

**Key Finding**: YOLO makes almost $3\times$ fewer background errors than Fast R-CNN because it sees the full image and reasons about global context.

## 5.4 vs. OverFeat

- **OverFeat**: Sliding window detection using CNN, optimizes for localization not detection

- **YOLO**: Direct detection optimization, reasons globally about context

- **Difference**: OverFeat is a disjoint system requiring post-processing, while YOLO is unified

# 6 Limitations and Trade-offs

> **YOLO Limitations**
>
> 1. **Spatial Constraints**: Each grid cell predicts only 2 boxes and one class, limiting detection of nearby objects (e.g., flocks of birds)
>
> 2. **Generalization**: Struggles with objects in unusual aspect ratios or configurations not seen in training
>
> 3. **Coarse Features**: Multiple downsampling layers result in relatively coarse features for small object detection
>
> 4. **Localization Errors**: Main source of error is incorrect localization, especially for small objects

# 7 Generalization Capability

A unique strength of YOLO is its ability to generalize to new domains:

Table 5: Generalization Performance on Artwork Datasets

| Method | VOC 2007 AP | Picasso AP | People-Art AP |
|---|---|---|---|
| YOLO | 59.2 | 53.3 | 45 |
| R-CNN | 54.2 | 10.4 | 26 |
| DPM | 43.2 | 37.8 | 32 |

**Insight**: YOLO's global reasoning about object size, shape, and relationships enables better transfer to artwork, where pixel-level appearance differs dramatically from natural images.

# 8   Architectural Innovations Summary

> **What Makes YOLO Different**
>
> 1. **Single Regression Problem**: Direct prediction from pixels to detections
>
> 2. **Global Context**: Entire image visible during prediction
>
> 3. **Unified Architecture**: End-to-end trainable single network
>
> 4. **Grid-Based Prediction**: $S \times S$ grid with $B$ boxes per cell
>
> 5. **Real-Time Speed**: 45-155 FPS vs. $<7$ FPS for previous methods
>
> 6. **Joint Optimization**: Simultaneous prediction of all boxes and classes
>
> 7. **Fewer Background Errors**: $3\times$ reduction compared to Fast R-CNN
>
> 8. **Strong Generalization**: Transfer learning to new domains

# 9   Conclusion

YOLO represents a fundamental paradigm shift in object detection by treating it as a unified regression problem rather than a multi-stage classification pipeline. Its key innovations—single-pass detection, global context reasoning, and direct bounding box regression—enable real-time performance while maintaining competitive accuracy. While YOLO trades some localization precision for speed, it significantly reduces background false positives and generalizes better to new domains.

The architecture's simplicity and effectiveness have made YOLO a foundational work in computer vision, spawning numerous improvements (YOLOv2-v8) that further refine the balance between speed and accuracy. YOLO demonstrated that object detection could be fast enough for real-time applications like autonomous driving while remaining accurate enough for practical deployment.

*"You only look once at an image to predict what objects are present and where they are."*
*— Redmon et al., 2016*