
You Only Look Once: Unified, Real-Time Object Detection

Anonymous Author(s)

Affiliation

Address

email

Abstract

We present a unified system for object detection. Our system processes images in real-time at 45 frames per second, and achieves 53.5% mean average precision on the PASCAL VOC object detection benchmark. We frame object detection as a regression problem to spatially separated bounding boxes and associated class probabilities. A single neural network predicts bounding boxes and class probabilities directly from full images. Notably, our system does not rely on pre-processing, region proposals, sliding window techniques, or post-processing of bounding boxes. Due to this simplicity, our system runs hundreds to thousands of times faster than comparable detection systems like DPM or R-CNN.

1 Introduction

Humans glance at an image and instantly know what objects are in the image, where they are, and how they interact. Our fast, accurate visual system allows us to perform complex tasks like driving or grocery shopping with little conscious thought. Fast, accurate, computational object detection would allow computers to drive cars in any weather without specialized sensors, enable assistive devices to convey real-time scene information to human users, and unlock the potential for general purpose, responsive robotic systems.

Convolutional neural networks (CNNs) achieve impressive performance on classification tasks at real-time speeds. Yet the best object detection systems stare at images for seconds and still make mistakes. These shortcomings directly result from how these systems approach object detection.

Current systems repurpose classifiers to perform detection. To detect an object these systems take a classifier for that object and evaluate it at various locations and scales in a test image. Systems like deformable parts models (DPM) use a sliding window approach where the classifier is run at evenly spaced locations over the entire image [3]. More recent approaches like R-CNN use region proposal methods to first generate potential bounding boxes in an image and then run the classifier on these proposed boxes. After classification, post-processing is used to refine the bounding box, eliminate duplicate or overlapping detections, and rescore the box based on other objects in the scene [4].

These region proposal techniques typically generate between a few hundred and a few thousand potential boxes per image. Selective Search, the most common region proposal method, takes 1-2 seconds per image to even generate the potential boxes [10]. The classifier then takes additional time to evaluate each of these potentially thousands of proposals. The best performing systems require 10-20 seconds per image while even those optimized for speed do not achieve real-time performance.

Additionally, even a highly accurate classifier will produce false positives when faced with so many proposals. Due to these inherent limitations, even state-of-the-art detection systems make obvious mistakes on test images.

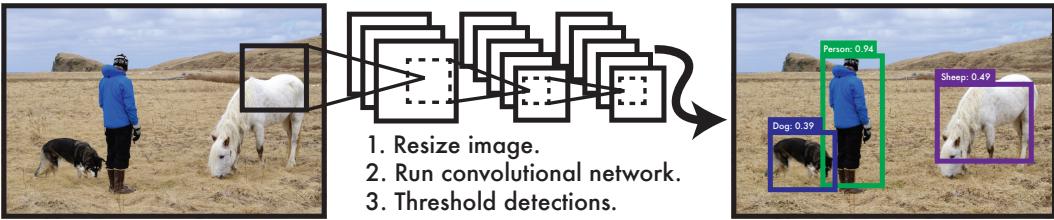


Figure 1: The YOLO Detection System. Our system (1) resizes the input image to 448×448 , (2) runs a single convolutional network on the image, and (3) thresholds the resulting detections by the model’s confidence. For all examples in this paper we use a confidence threshold of 0.2 unless otherwise noted.

Finally, these detection pipelines rely on independent techniques at every stage that cannot be optimized jointly. A typical pipeline uses Selective Search for region proposals, a convolutional network for feature extraction, a collection of one-versus-all SVMs for classification, non-maximal suppression to reduce duplicates, and a linear model to adjust the final bounding box coordinates. Selective Search tries to maximize recall while the SVMs optimize for single class accuracy and the linear model learns from localization error.

Our system is refreshingly simple. A single convolutional network simultaneously predicts multiple bounding boxes and class probabilities for those boxes. Our system enforces spatial diversity in the proposed bounding boxes without relying on non-maximal suppression to eliminate duplicates. We train our network on full images to directly optimize detection performance. At test time, a single network evaluation produces complete detections of multiple objects in multiple categories without any pre or post-processing.

2 Unified Detection

We unify the separate components of object detection into a single neural network. Using our system, you only look once (YOLO) at an image to predict what objects are present and where they are. Our network uses features from the entire image while predicting each bounding box. It also predicts every bounding box for an image simultaneously. This means our network reasons globally about the full image and all the objects in the image. The YOLO design enables end-to-end training and real-time speeds while maintaining high average precision.

2.1 Design

Our system divides the image into a 7×7 grid. If the center of an object falls into a grid cell, that grid cell is responsible for detecting that object. Each grid cell predicts a bounding box and class probabilities associated with that bounding box.

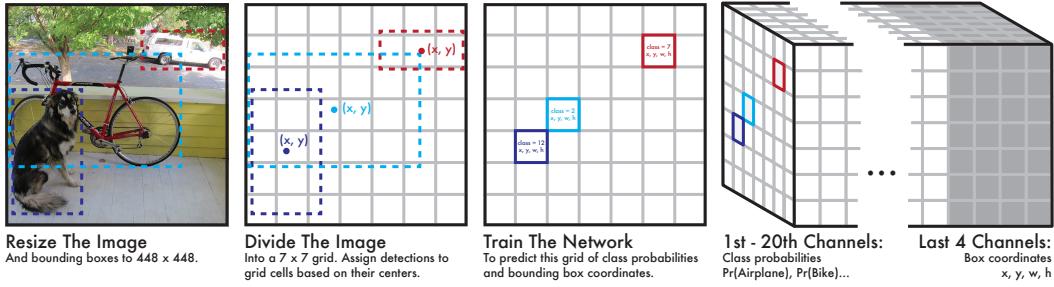


Figure 2: The Model. Our system models detection as a regression problem to a $7 \times 7 \times 24$ tensor. This tensor encodes bounding boxes and class probabilities for objects in the image.

We implement this model as a convolutional neural network and evaluate it on the PASCAL VOC detection dataset [2]. The initial convolutional layers of the network extract features from the image while the fully connected layers predict the output probabilities and coordinates.

Our network architecture is inspired by the GoogLeNet model for image classification [9]. The network has 24 convolutional layers followed by 2 fully connected layers. However, instead of the inception modules used by GoogLeNet we simply use 1×1 reduction layers followed by 3×3 convolutional layers. We also replace maxpooling layers with strided convolutions.

The final output of our network is a 7×7 grid of predictions. Each grid cell predicts 20 conditional class probabilities, and 4 bounding box coordinates.

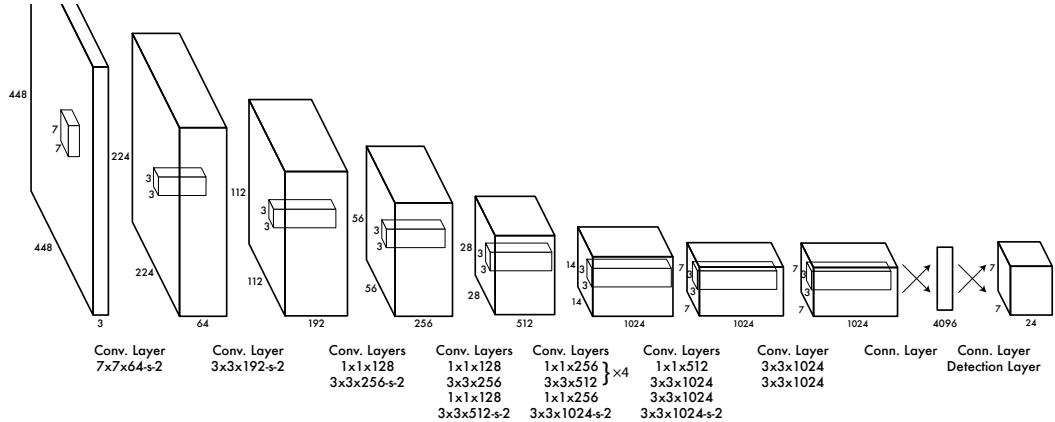


Figure 3: The Architecture. Our detection network has 24 convolutional layers followed by 2 fully connected layers. The network uses strided convolutional layers to downsample the feature space instead of maxpooling layers. Alternating 1×1 convolutional layers reduce the features space from preceding layers. We pretrain the convolutional layers on the ImageNet classification task at half the resolution (224×224 input image) and then double the resolution for detection.

2.2 Training

We pretrain our convolutional layers on the ImageNet 1000-class competition dataset [7]. For pre-training we use the first 20 convolutional layers from 3 followed by a maxpooling layer and two fully connected layers. We train this network for approximately a week and achieve top-5 accuracy of 86% on the ImageNet 2012 validation set.

We then adapt the model to perform detection. We add four convolutional layers and two fully connected layers with randomly initialized weights. Detection often requires fine-grained visual information so we increase the input resolution of the network from 224×224 to 448×448 .

Our final layer predicts both class probabilities and bounding box coordinates. We normalize the bounding box width and height by the image width and height so that they fall between 0 and 1. We parameterize the bounding box x and y coordinates to be offsets of a particular grid cell location so they are also bounded between 0 and 1. We use a logistic activation function to reflect these constraints on the final layer. All other layers use the following leaky rectified linear activation:

$$\phi(x) = \begin{cases} 1.1x, & \text{if } x > 0 \\ .1x, & \text{otherwise} \end{cases} \quad (1)$$

We optimize for sum-squared error in the output of our model. We chose sum-squared error because it is easy to optimize however it does not perfectly align with our actual goal of maximizing average precision. It weights localization error equally with classification error which may not be ideal. To remedy this, we use a scaling factor to adjust weight given to error from coordinate predictions versus error from class probabilities. In our final model we use a scaling factor of 4.

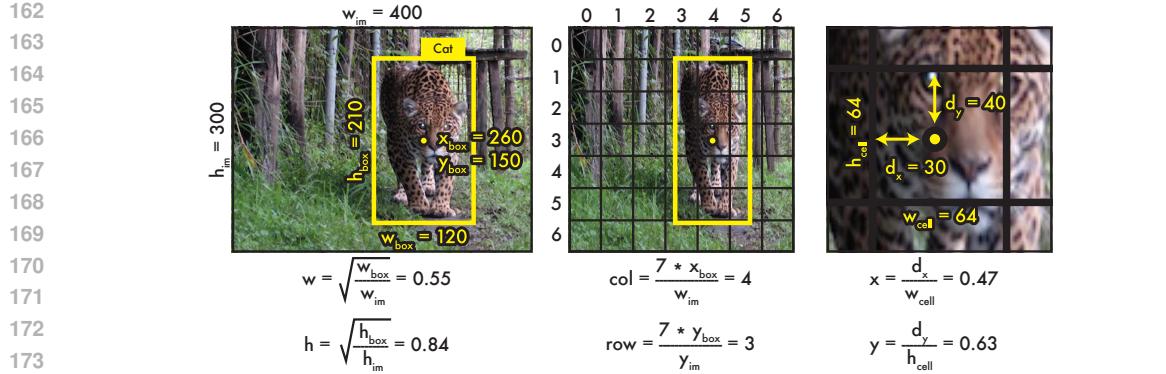


Figure 4: The Coordinate System. This image shows an example transformation from bounding box coordinates to the coordinates used by our output layer.

Sum-squared error also equally weights errors in large boxes and small boxes. Our error metric should reflect that small deviations in large boxes matter less than in small boxes. To partially address this we predict the square root of the bounding box width and height instead of the width and height directly.

If cell i predicts class probabilities $\hat{p}_i(\text{aeroplane}), \hat{p}_i(\text{bicycle})\dots$ and the bounding box $\hat{x}_i, \hat{y}_i, \hat{w}_i, \hat{h}_i$ then our full loss function for an example is:

$$\sum_{i=0}^{49} \left(\mathbb{1}_{\text{obj}}((x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 + (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2) + \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \right) \quad (2)$$

Note that if there is no object in a cell we do not consider any loss from the bounding box coordinates predicted by that cell. In this case, there is no ground truth bounding box so we only penalize the associated probabilities with that region.

We train the network for about 120 epochs on the training and validation data sets from PASCAL VOC 2007 and 2012 as well as the test set from 2007, a total of 21k images. Throughout training we use a batch size of 64, a momentum of 0.9 and a decay of 0.0005. We use two learning rates during training: 10^{-2} and 10^{-3} . Training diverges if we use the higher learning rate, 10^{-2} , from the start. We use the lower rate, 10^{-3} , for one epoch so that the randomly initialized weights in the final layers can settle to a reasonable value. Then we train with the following learning rate schedule: 10^{-2} for 80 epochs, and 10^{-3} for 40 epochs.

To avoid overfitting we use dropout and extensive data augmentation. A dropout layer with a probability $\text{Pr} = .5$ after the first connected layer prevents co-adaptation between layers. For data augmentation we introduce random scaling and translations of up to 10% of the original image size. We also randomly adjust the exposure and saturation of the image by up to a factor of 2.

2.2.1 Nuisance Variables

Each grid cell predicts class probabilities for that area of the image. There are 49 cells with a possible 20 classes each yielding 980 predicted probabilities per image. Most of these will be probabilities will be zero since only a few objects appear in any given image. Left unchecked, this imbalance pushes all of the probabilities to zero, leading to divergence during training.

To overcome this, we add an extra variable to each grid location, the probability that any object exists in that location regardless of class. Thus instead of 20 class probabilities we have 1 "objectness" probability, $\text{Pr}(\text{Object})$, and 20 conditional probabilities: $\text{Pr}(\text{Airplane}|\text{Object})$, $\text{Pr}(\text{Bicycle}|\text{Object})$, etc.

216 To get the unconditional probability for an object class at a given location we simply multiply the
217 "objectness" probability by the conditional class probability:
218

$$\Pr(\text{Cat}) = \Pr(\text{Object}) * \Pr(\text{Cat}|\text{Object}) \quad (3)$$

221 We can optimize these probabilities independently or jointly using a novel "detection layer" in our
222 convolutional network. During the initial stages of training we optimize them independently to
223 improve model stability. We update the "objectness" probabilities at every location however we
224 only update the conditional probabilities at locations that actually contain an object. This means
225 there are far fewer probabilities getting pushed towards zero.

226 During later stages of training we optimize the unconditioned probabilities by performing the re-
227 quired multiplications in the network and calculating error based on the result.
228

229 2.2.2 Predicting IOU 230

231 Like most detection systems, our network has trouble precisely localizing small objects. While it
232 may correctly predict that an object is present in an area of the image, if it does not predict a precise
233 enough bounding box the detection is counted as a false positive.

234 We want YOLO to have some notion of uncertainty in its probability predictions. Instead of predict-
235 ing 1-0 probabilities we can scale the target class probabilities by the IOU of the predicted bounding
236 box with the ground truth box for a region. When YOLO predicts good bounding boxes it is also
237 encouraged to predict high class probabilities. For poor bounding box predictions it learns to predict
238 lower confidence probabilities.

239 We do not train to predict IOU from the beginning, only during the second stage of training. It is not
240 necessary for good performance but it does boost our mean average precision by 3-4%.
241

242 2.3 Inference 243

244 Just like in training, predicting detections for a test image only requires one network evaluation.
245 The network predicts 49 bounding boxes per image and class probabilities for each box. YOLO is
246 extremely fast at test time since it only requires a single network evaluation unlike classifier-based
247 methods.

248 Non-Maximal Suppression 249

250 The grid design enforces spatial diversity in the bounding box predictions. Often it is clear which
251 grid cell an object falls in to and the network only predicts one box for each object. However, some
252 large objects or objects near the border of multiple cells can be well localized by multiple cells. Non-
253 maximal suppression can be used to fix these multiple detections. While not critical to performance
254 as it is for R-CNN or DPM, non-maximal suppression adds 2-3% in mAP.

255 2.4 Benefits and Limitations 256

257 This model and architecture form a unified pipeline for object detection. It is simple to construct
258 and can be trained directly on full images. Unlike classifier-based approaches, YOLO is trained on a
259 loss function that directly corresponds to detection performance and every piece of the pipeline can
260 be trained jointly.

261 The network reasons about the entire image during inference which makes it less likely to predict
262 background false positives than sliding window or proposal region detectors. Moreover, it predicts
263 detections with only a single network evaluation, making it extremely fast.

264 A cell is only responsible for objects centered in that cell so we restrict the coordinates to fall within
265 that cell. This imposes a strong spatial constraint on the predictions which has both benefits and
266 drawbacks. Cells rarely predict multiple bounding boxes for the same object.
267

268 Conversely, this spatial constraint limits the number of nearby objects that our model can predict.
269 If two objects fall into the same cell our model can only predict one of them. Our model struggles
with small objects that appear in groups, such as flocks of birds.

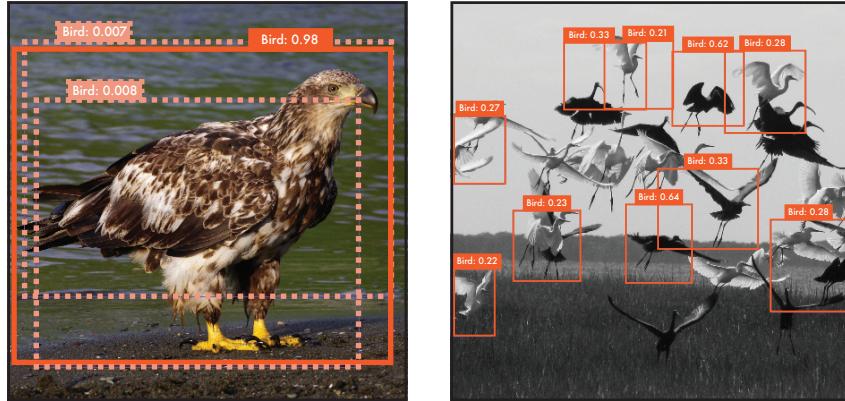


Figure 5: Global reasoning. Our model makes local predictions but it reasons about the entire image. When multiple grid cells can predict good bounding boxes for an object, the model picks the best cell and downweights surrounding predictions. **Problems with small objects in groups.** The spatial constraint we impose restricts the model’s ability to predict boxes for multiple small objects that are next to each other in an image. Our model cannot predict correct bounding boxes for multiple objects if they fall into the same grid cell.

Finally, our loss function treats errors the same in small bounding boxes versus large bounding boxes. A small error in a large box is generally benign but a small error in a small box has a much greater effect on IOU. Our main source of error is incorrect localizations.

3 Comparison to Other Detection Systems

Deformable parts models. Deformable parts models (DPM) use a sliding window approach to object detection [3]. DPM uses a disjoint pipeline to extract static features, classify regions, predict bounding boxes for high scoring regions, etc. Our system replaces all of these disparate parts with a single convolutional neural network. The network performs feature extraction, bounding box prediction, non-maximal suppression, and contextual reasoning all concurrently. Instead of static features, the network trains the features in-line and optimizes them for the detection task. Our unified architecture leads to a faster, more accurate model than DPM.

R-CNN. Another class of detection algorithms uses region proposals instead of sliding windows to find objects in images. These systems use region proposal methods like Selective Search [10] to generate potential bounding boxes in an image. Instead of scanning through every region in the window, now the classifier only has to score a small subset of potential regions in an image. Good region proposal methods maintain high recall despite greatly limiting the search space. This performance comes at a cost. Selective Search, even in “fast mode” takes about 2 seconds to propose regions for an image.

R-CNN shares many design aspects with DPM. After region proposal, R-CNN uses the same multi-stage pipeline of feature extraction (using CNNs instead of HOG), SVM scoring, non-maximal suppression, and bounding box prediction using a linear model [4].

YOLO shares some similarities with R-CNN. Each grid cell proposes a potential bounding box and then scores that bounding box using convolutional features. However, our system puts spatial constraints on the grid cell proposals which helps mitigate multiple detections of the same object. Our system also proposes far fewer bounding boxes, only 49 per image compared to about 2000. Finally, our system combines these individual components into a single, jointly optimized model.

Deep Multibox. Most R-CNN variants use Selective Search to generate region proposals. Szegedy et al. instead train a convolutional neural network to predict regions of interest [1]. MultiBox can also perform single object detection by replacing the confidence prediction with a single class prediction. However, MultiBox cannot perform general object detection and is still just a piece in a larger detection pipeline.

| VOC 2012 test | mAP | aero | bike | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | mbike | person | plant | sheep | sofa | train | tv |
|--------------------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| MR.CNN.S.CNN | 70.7 | 85.0 | 79.6 | 71.5 | 55.3 | 57.7 | 76.0 | 73.9 | 84.6 | 50.5 | 74.3 | 61.7 | 85.5 | 79.9 | 81.7 | 76.4 | 41.0 | 69.0 | 61.2 | 77.7 | 72.1 |
| DEEP.LENS.COCO | 70.1 | 84.0 | 79.4 | 71.6 | 51.9 | 51.1 | 74.1 | 72.1 | 88.6 | 48.3 | 73.4 | 57.8 | 86.1 | 80.0 | 80.7 | 70.4 | 46.6 | 69.6 | 68.8 | 75.9 | 71.4 |
| Fast R-CNN + YOLO | 70.0 | 83.0 | 78.4 | 73.4 | 55.7 | 42.5 | 78.2 | 72.7 | 89.5 | 48.2 | 74.0 | 56.4 | 87.2 | 80.8 | 80.7 | 74.4 | 41.1 | 70.0 | 67.1 | 81.2 | 66.0 |
| NoC | 68.8 | 82.8 | 79.0 | 71.6 | 52.3 | 53.7 | 74.1 | 69.0 | 84.9 | 46.9 | 74.3 | 53.1 | 85.0 | 81.3 | 79.5 | 72.2 | 38.9 | 72.4 | 59.5 | 76.7 | 68.1 |
| Fast R-CNN | 68.4 | 82.3 | 78.4 | 70.8 | 52.3 | 38.7 | 77.8 | 71.6 | 89.3 | 44.2 | 73.0 | 55.0 | 87.5 | 80.5 | 80.8 | 72.0 | 35.1 | 68.3 | 65.7 | 80.4 | 64.2 |
| NUS.NIN.C2000 | 63.8 | 80.2 | 73.8 | 61.9 | 43.7 | 43.0 | 70.3 | 67.6 | 80.7 | 41.9 | 69.7 | 51.7 | 78.2 | 75.2 | 76.9 | 65.1 | 38.6 | 68.3 | 58.0 | 68.7 | 63.3 |
| BabyLearning | 63.2 | 78.0 | 74.2 | 61.3 | 45.7 | 42.7 | 68.2 | 66.8 | 80.2 | 40.6 | 70.0 | 49.8 | 79.0 | 74.5 | 77.9 | 64.0 | 35.3 | 67.9 | 55.7 | 68.7 | 62.6 |
| R-CNN VGG BB | 62.4 | 79.6 | 72.7 | 61.9 | 41.2 | 41.9 | 65.9 | 66.4 | 84.6 | 38.5 | 67.2 | 46.7 | 82.0 | 74.8 | 76.0 | 65.2 | 35.6 | 65.4 | 54.2 | 67.4 | 60.3 |
| NUS.NIN | 62.4 | 77.9 | 73.1 | 62.6 | 39.5 | 43.3 | 69.1 | 66.4 | 78.9 | 39.1 | 68.1 | 50.0 | 77.2 | 71.3 | 76.1 | 64.7 | 38.4 | 66.9 | 56.2 | 66.9 | 62.7 |
| R-CNN VGG | 59.2 | 76.8 | 70.9 | 56.6 | 37.5 | 36.9 | 62.9 | 63.6 | 81.1 | 35.7 | 64.3 | 43.9 | 80.4 | 71.6 | 74.0 | 60.0 | 30.8 | 63.4 | 52.0 | 63.5 | 58.7 |
| Feature Edit | 56.3 | 74.6 | 69.1 | 54.4 | 39.1 | 33.1 | 65.2 | 62.7 | 69.7 | 30.8 | 56.0 | 44.6 | 70.0 | 64.4 | 71.1 | 60.2 | 33.3 | 61.3 | 46.4 | 61.7 | 57.8 |
| YOLO | 53.6 | 71.6 | 62.2 | 55.2 | 35.9 | 23.2 | 62.4 | 53.7 | 78.1 | 34.0 | 52.9 | 38.7 | 72.2 | 67.3 | 66.3 | 62.1 | 25.5 | 50.0 | 46.9 | 67.4 | 46.3 |
| R-CNN BB | 53.3 | 71.8 | 65.8 | 52.0 | 34.1 | 32.6 | 59.6 | 60.0 | 69.8 | 27.6 | 52.0 | 41.7 | 69.6 | 61.3 | 68.3 | 57.8 | 29.6 | 57.8 | 40.9 | 59.3 | 54.1 |
| SDS | 50.7 | 69.7 | 58.4 | 48.5 | 28.3 | 28.8 | 61.3 | 57.5 | 70.8 | 24.1 | 50.7 | 35.9 | 64.9 | 59.1 | 65.8 | 57.1 | 26.0 | 58.8 | 38.6 | 58.9 | 50.7 |
| R-CNN | 49.6 | 68.1 | 63.8 | 46.1 | 29.4 | 27.9 | 56.6 | 57.0 | 65.9 | 26.5 | 48.7 | 39.5 | 66.2 | 57.3 | 65.4 | 53.2 | 26.2 | 54.5 | 38.1 | 50.6 | 51.6 |

Table 1: PASCAL VOC 2012 Leaderboard. YOLO compared with the full comp4 (outside data allowed) public leaderboard as of June 5th, 2015. Mean average precision and per-class average precision are shown for a variety of detection methods. YOLO is the top detection method that is not based on the R-CNN detection framework. Fast R-CNN + YOLO is one of the top methods overall, with a 1.6% boost over Fast R-CNN and the highest average precision in 6 out of 20 classes.

Both YOLO and MultiBox use a convolutional network to predict bounding boxes and confidence scores. However, YOLO additionally predicts class-specific probabilities making it able to perform general object detection. Instead of predicting bounding boxes based on clusters, YOLO uses the grid approach to enforce spatial diversity in predicted clusters.

Overfeat. Sermanet et al. train a convolutional neural network to perform localization and adapt that localizer to perform detection [8]. Overfeat efficiently performs sliding window detection but it is still a disjoint system. Overfeat optimizes for localization, not detection performance. Like DPM, the localizer only sees local information when making a prediction. Overfeat cannot reason about global context and thus requires significant post-processing to produce coherent detections.

4 Experiments

We present detection results for the PASCAL VOC 2012 dataset and compare our mean average precision (mAP) and runtime to other top detection methods. We also perform error analysis on the VOC 2007 dataset. We compare our results to Fast R-CNN, one of the highest performing versions of R-CNN [5]. We use publicly available runs of Fast R-CNN available on GitHub. Finally we show that a combination of our method with Fast R-CNN gives a significant performance boost.

On the VOC 2012 test set we achieve 53.6 mAP. This is lower than the current state of the art, closer to R-CNN based methods that use AlexNet. Our system struggles with small objects compared to its closest competitors. On categories like `bottle`, `sheep`, and `tv/monitor` YOLO scores 8-10 percentage points lower than R-CNN or Feature Edit. However, on other categories like `cat` and `horse` YOLO achieves significantly higher performance. We further investigate the source of these performance disparities in section 4.3.

4.2 Speed

At test time YOLO processes images at 45 frames per second. This is considerably faster than classifier-based methods with similar mAP. Normal R-CNN using AlexNet or the small VGG network take 400-500x longer to process images. The recently proposed Fast R-CNN shares convolutional features between the bounding boxes but still relies on selective search for bounding box proposals which accounts for the bulk of their processing time. YOLO is still around 100x faster than Fast R-CNN.

| | mAP | Prediction Time | Images/Sec | Run Time |
|------------------------|------|-----------------|------------|----------|
| R-CNN (VGG-16) | 66.9 | 48.2 hr | 0.02 fps | 1500x |
| Fast R-CNN (VGG-16) | 59.2 | 3.1 hr | 0.45 fps | 100x |
| R-CNN (Small VGG) | 60.2 | 14.4 hr | 0.09 fps | 500x |
| Fast R-CNN (Small VGG) | 59.2 | 2.9 hr | 0.48 fps | 93x |
| R-CNN (CaffeNet) | 58.5 | 12.2 hr | 0.11 fps | 409x |
| Fast R-CNN (CaffeNet) | 57.1 | 2.8 hr | 0.48 fps | 93x |
| YOLO | 58.8 | 110 sec | 45 fps | - |

Table 2: Prediction Timing. mAP and timing information for R-CNN, Fast R-CNN, and YOLO on the VOC 2007 test set. Timing information is given both as frames per second and the time each method takes to process the full 4952 image set. The final column shows the relative speed of YOLO compared to that method.

4.3 VOC 2007 Error Analysis

Using the methodology and tools of Hoiem et al. [6] we analyze our performance on the VOC 2007 test set. We compare YOLO to Fast R-CNN using VGG-16, one of the highest performing object detectors.

Figure 6 compares frequency of localization and background errors between Fast R-CNN and YOLO. A detection is considered a localization error if it overlaps a true positive in the image but by less than the required 50% IOU. A detection is a background error if the box does not overlap any objects of any class in the image.

Fast R-CNN makes around the same number of localization errors and background errors. Over all 20 classes in the top N detections 13.6% are localization errors and 8.6% are background errors.

YOLO makes far more localization errors but relatively few background errors. Averaged across all classes, of its top N detections 24.7% are localization errors and a mere 4.3% are background errors. This is about twice the number of localization errors but half the number of background detections.

YOLO uses global context to evaluate detections while R-CNN only sees small portions of the image. Many of the background detections made by R-CNN are obviously not objects when shown in context.

4.4 Combining Fast R-CNN and YOLO

YOLO makes far fewer background mistakes than Fast R-CNN. By using YOLO to eliminate background detections from Fast R-CNN we get a significant boost in performance. For every bounding box that R-CNN predicts we check to see if YOLO predicts a similar box. If it does, we give that prediction a boost based on the probability predicted by YOLO and the overlap between the two boxes. This reorders the detections to favor those predicted by both systems. Since we still use Fast R-CNN’s bounding box predictions we do not introduce any localization error. Thus we take advantage of the best aspects of both systems.

The best Fast R-CNN model achieves a mAP of 71.8% on the VOC 2007 test set. When combined with YOLO, its mAP increases by 2.9% to 74.7%. We also tried combining the top Fast R-CNN

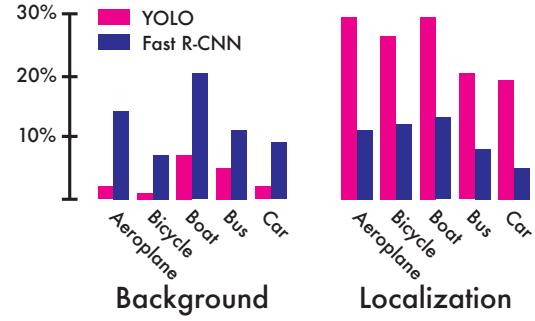


Figure 6: Error Analysis: Fast R-CNN vs. YOLO These charts show the percentage of localization and background errors in the top N detections for various categories ($N = \#$ objects in that category).

| | mAP | Combined mAP | % Pt Increase |
|-----------------------------|------|--------------|---------------|
| Fast R-CNN, No Augmentation | - | 71.8 | - |
| Fast R-CNN (2007 data) | 66.9 | 72.4 | .6 |
| Fast R-CNN (Small VGG) | 59.2 | 72.4 | .6 |
| Fast R-CNN (CaffeNet) | 57.1 | 72.1 | .3 |
| YOLO | 59.1 | 74.7 | 2.9 |

438
439 **Table 3: Model combination experiments on VOC 2007.**
440
441

442 model with several other training versions of Fast R-CNN. Those ensembles produced small in-
443 creases in mAP between .3 and .6%, see table ?? for details.

444 Using this combination strategy we achieve a significant boost on the VOC 2012 and 2010 test sets
445 as well, around 2%. The combined Fast R-CNN + YOLO model is the second highest performing
446 model on the VOC 2012 leaderboard.

448 5 Conclusion

450 We introduce YOLO, a unified detection pipeline. YOLO is extremely fast and moderately accurate.
451

452 References

- 454 [1] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks.
455 In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 2155–2162. IEEE,
456 2014.
- 457 [2] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal
458 visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–
136, Jan. 2015.
- 459 [3] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discrimi-
460 natively trained part based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
461 32(9):1627–1645, 2010.
- 462 [4] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection
463 and semantic segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference
on*, pages 580–587. IEEE, 2014.
- 464 [5] R. B. Girshick. Fast R-CNN. *CoRR*, abs/1504.08083, 2015.
- 465 [6] D. Hoiem, Y. Chodpathumwan, and Q. Dai. Diagnosing error in object detectors. In *Computer Vision–
ECCV 2012*, pages 340–353. Springer, 2012.
- 466 [7] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla,
467 M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *Internat-
ional Journal of Computer Vision (IJCV)*, 2015.
- 468 [8] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition,
469 localization and detection using convolutional networks. *CoRR*, abs/1312.6229, 2013.
- 470 [9] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabi-
471 novich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014.
- 472 [10] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders. Selective search for object recogni-
473 tion. *International journal of computer vision*, 104(2):154–171, 2013.