

# YOLO: 统一的实时目标检测

Joseph Redmon\*, Santosh Divvala\*†, Ross Girshick¶, Ali Farhadi\*†  
University of Washington\*, Allen Institute for AI†, Facebook AI Research¶  
<http://pjreddie.com/yolo/>

## 摘要

我们提出了 YOLO，一种新的目标检测方法。以前的目标检测工作重复利用分类器来完成检测任务。相反，我们将目标检测框架看作回归问题，从空间上分割边界框和相关的类别概率。单个神经网络在一次评估中直接从整个图像上预测边界框和类别概率。由于整个检测流水线是单一网络，因此可以直接对检测性能进行端到端的优化。

我们的统一架构非常快。我们的基础 YOLO 模型以 45 帧/秒的速度实时处理图像。Fast YOLO 是 YOLO 的一个较小版本，每秒能处理惊人的 155 帧图像，同时实现其它实时检测器两倍的 mAP。与最先进的检测系统相比，YOLO 虽然存在较多的定位错误，但很少将背景预测成假阳性（译者注：其它先进的目标检测算法将背景预测成目标的概率较大）。最后，YOLO 能学习到目标非常通用的表示。当从自然图像到艺术品等其它领域泛化时，它都优于其它检测方法，包括 DPM 和 R-CNN。

## 1. 引言

人们瞥一眼图像，立即知道图像中的物体是什么，它们在哪里以及它们如何相互作用。人类的视觉系统是快速和准确的，使我们能够执行复杂的任务，例如如驾驶车辆时不会刻意地进行思考或思想。快速、准确的目标检测算法可以让计算机在没有专用传感器的情况下驾

驶汽车，使辅助设备能够向人类用户传达实时的场景信息，并具有解锁通用目的和响应机器人系统的潜力。

目前的检测系统重复利用分类器来执行检测。为了检测目标，这些系统为该目标提供一个分类器，并在不同的位置对其进行评估，并在测试图像中进行缩放。像**可变形部件模型（DPM）**这样的系统使用滑动窗口方法，其分类器在整个图像的均匀间隔的位置上运行[10]。

最近的许多方法，如 R-CNN 使用 region proposal 方法首先在图像中生成潜在的边界框，然后在这些提出的框上运行分类器。在分类之后，通过后处理对边界框进行修正，消除重复的检测，并根据场景中的其它目标重新定位边界框[13]。这些复杂的流程很慢，很难优化，因为每个单独的组件都必须单独进行训练。

我们将目标检测重构并看作为单一的回归问题，直接从图像像素到边界框坐标和类别概率。使用我们的系统，您只需要在图像上看一次（you only look once, YOLO），以预测出现的目标和位置。

YOLO 新奇又很简单：如图 1 所示。单个卷积网络同时预测这些框的多个边界框和类别概率值。YOLO 在全图像上训练并直接优化检测性能。这种统一的模型比传统的目标检测方法有一些好处。

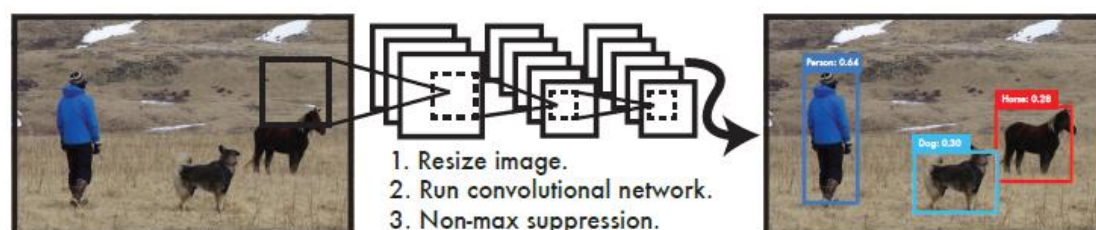


图 1: YOLO 检测系统。用 YOLO 处理图像简单直接。我们的系统（1）将输入图像调整为  $448 \times 448$ ，（2）在图像上运行单个卷积网络，以及（3）由模型的置信度对所得到的检测进行阈值处理。

首先，YOLO 速度非常快。由于我们将检测视为回归问题，所以我们不需要复杂的流程。测试时我们在一张新图像上简单的运行我们的神经网络来预测检测。我们的基础网络以每秒 45 帧的速度运行，在 Titan X GPU 上没有批处理，快速版本运行速度超过 150fps。这意味着我们可以在不到 25 毫秒的延迟内实时处理流媒体视频。此外，YOLO 实现了其它实时系统两倍以上 mAP。关于我们的系统在网络摄像头实时运行的演示，请参阅我们的项目网页：<http://pjreddie.com/yolo/>。

其次，YOLO 在进行预测时，会对图像进行全局地推理。与基于滑动窗口和 region proposal 的技术不同，YOLO 在训练期间和测试时会看到整个图像，所以它隐式地编码了关于类的上下文信息以及它们的外形。Fast R-CNN 是一种顶级的检测方法[14]，但因为它看不到更大的上下文，所以在图像中会将背景块误检为目标。与 Fast R-CNN 相比，YOLO 的背景误检数量少了一半。

第三，YOLO 学习目标的泛化表示。当在自然的图像上进行训练并对艺术作品进行测试时，YOLO 大幅优于 DPM 和 R-CNN 等顶级检测方法。由于 YOLO 具有高度泛化能力，因此在应用于新领域或碰到非正常输入时很少出故障。

YOLO 在准确度上仍然落后于最先进的检测系统。虽然它可以快速识别图像中的目标，但它仍在努力精确定位一些目标，尤其是一些小目标。我们在实验中会进一步检查这些权衡。

我们所有的训练和测试代码都是开源的。各种预训练模型也都可以下载。

## 2. 统一的检测

我们将目标检测的单独组件集成到单个神经网络中。我们的网络使用整个图像的特征来预测每个边界框。它还可以同时预测一张图像中的所有类别的所有边界框。这意味着我们的网络全面地推理整张图像和图像中的所有目标。YOLO 设计可实现端到端训练和实时的速度，同时保持较高的平均精度。

我们的系统将输入图像分成  $S \times S$  的网格。如果一个目标的中心落入一个网格单元中，该网格单元负责检测该目标。

每个网格单元预测这些盒子的  $B$  个边界框和置信度分数。这些置信度分数反映了该模型对盒子是否包含目标的置信度，以及它预测盒子的准确程度。在形式上，我们将置信度定义为  $\text{Pr}(\text{Object}) * \text{IOU}_{\text{truthpred}}$ 。如果该单元格中不存在目标，则置信度分数应为零。否则，我们希望置信度分数等于预测框与真实值之间联合部分的交集（IOU）。

每个边界框包含 5 个预测： $x$ 、 $y$ 、 $w$ 、 $h$  和置信度。 $(x, y)$  坐标表示边界框相对于网格单元边界框的中心。宽度和高度是相对于整张图像预测的。最后，置信度预测表示预测框与实际边界框之间的 IOU。

每个网格单元还预测  $C$  个条件类别概率  $\text{Pr}(\text{Class} | \text{Object})$ 。这些概率以包含目标的网格单元为条件。每个网格单元我们只预测的一组类别概率，而不管边界框的数量  $B$  是多少。

在测试时，我们乘以条件类概率和单个盒子的置信度预测值：

$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}} \quad (1)$$

它为我们提供了每个框特定类别的置信度分数。这些分数编码了该类出现在框中的概率以及预测框拟合目标的程度。

为了在 Pascal VOC 上评估 YOLO，我们使用  $S=7$ ， $B=2$ 。Pascal VOC 有 20 个标注类，所以  $C=20$ 。我们最终的预测是  $7 \times 7 \times 30$  的张量。

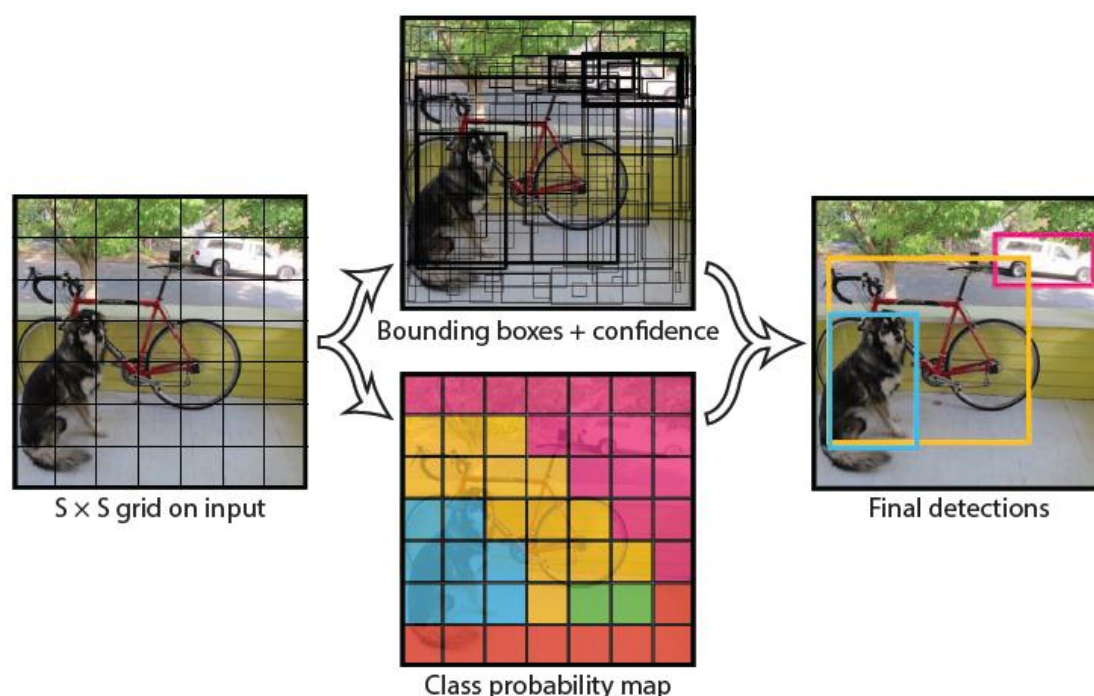


图 2：模型。我们的系统将检测任务构建为回归问题。它将图像分成  $S \times S$  的网格，并且对于每个网格单元都预测  $B$  个边界框、这些边界框的置信度以及  $C$  个类别概率。这些预测结果被编码为  $S \times S \times (B \times 5 + C)$  的张量。

## 2.1. 网络设计

我们将此模型实现为卷积神经网络，并在 Pascal VOC 检测数据集[9]上进行评估。网络的初始卷积层从图像中提取特征，而全连接层预测输出概率和坐标。

我们的网络架构受到 GoogLeNet 图像分类模型的启发[34]。我们的网络有 24 个卷积层，后面是 2 个全连接层。我们只使用  $1 \times 1$  降维层，后面是  $3 \times 3$  卷积层，这与 Lin 等人[22]的模型结构类似，而不是 GoogLeNet 使用的 Inception 模块。完整的网络如图 3 所示。

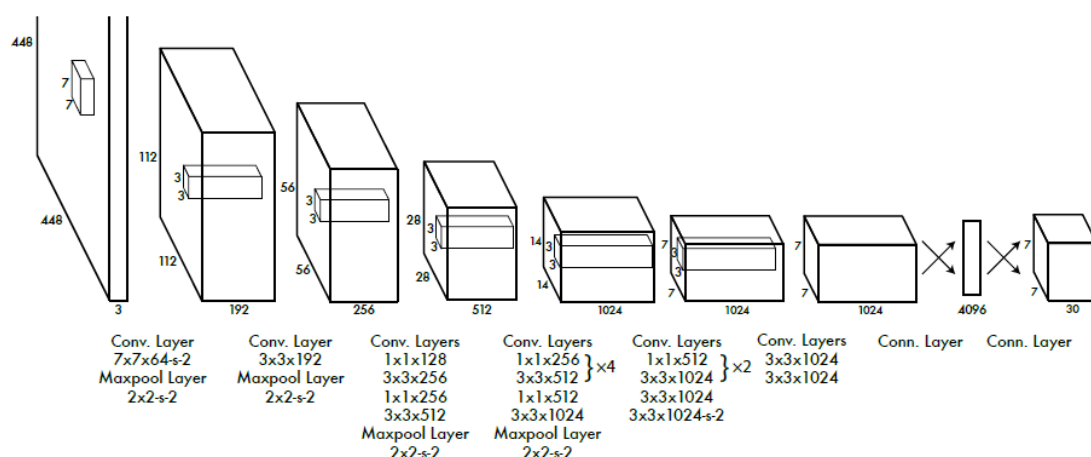


图 3：模型架构。我们的检测网络有 24 个卷积层，其次是 2 个全连接层。交替使用  $1 \times 1$  卷积层减少了前面层的特征空间。我们在 ImageNet 分类任务上以一半的分辨率( $224 \times 224$  的输入图像)预训练卷积层，然后将分辨率加倍来进行检测。

我们还训练了快速版本的 YOLO，旨在挑战快速目标检测的界限。快速 YOLO 使用具有较少卷积层（9 层而非 24 层）的神经网络，在这些层中使用较少的卷积核。除了网络规模之外，YOLO 和快速 YOLO 的所有训练和测试参数都是相同的。

我们网络的最终输出是  $7 \times 7 \times 30$  的预测张量。

## 2.2. 训练

我们在 ImageNet 1000 类竞赛数据集[30]上预训练我们的卷积层。对于预训练，我们使用图 3 中的前 20 个卷积层，接着是平均池化层和全连接层。我们对这个网络进行了大约一周的训练，并且在 ImageNet 2012 验证集上获得了单一裁剪图像 88% 的 top-5 准确率，



与 Caffe 模型池中的 GoogLeNet 模型相当。我们使用 Darknet 框架进行所有的训练和推断[26]。

然后我们转换模型来进行检测任务。Ren 等人表明，预训练网络中增加卷积层和连接层可以提高性能[29]。按照他们的例子，我们添加了四个卷积层和两个全连接层，并且对权重进行随机初始化。检测通常需要细粒度的视觉信息，因此我们将网络的输入分辨率从  $224 \times 224$  变为  $448 \times 448$ 。

我们的最后一层预测类别概率和边界框坐标。我们通过图像宽度和高度来归一化边界框的宽度和高度，使它们落在 0 和 1 之间。我们将边界框 x 和 y 坐标参数化为特定网格单元位置的偏移量，所以它们边界也在 0 和 1 之间。

我们对最后一层使用线性激活函数，所有其它层使用下面的 leaky ReLU 激活函数：

$$\phi(x) = \begin{cases} x, & \text{if } x > 0 \\ 0.1x, & \text{otherwise} \end{cases} \quad (2)$$

我们优化了模型输出的平方和误差。我们使用平方和误差是因为它很容易进行优化，但是它并不完全符合我们最大化平均精度的目标。分类误差与定位误差的权重是一样的，这可能并不理想。另外，在每张图像中，许多网格单元不包含任何对象。这将导致这些单元格的“置信度”分数为零，通常压倒了包含目标的单元格的梯度。这可能导致模型不稳定，从而导致训练过早发散。

为了改善这一点，我们增加了边界框坐标预测的损失，并减少了不包含目标边界框的置信度预测的损失。我们使用两个参数  $\lambda_{\text{coord}}$  和  $\lambda_{\text{noobj}}$  来完成这个调整工作。我们设置  $\lambda_{\text{coord}}=5$  和  $\lambda_{\text{noobj}}=0.5$ 。

平方和误差也可以在大盒子和小盒子中同样加权误差。我们的误差指标应该反映出，大盒子中小偏差的重要性不如小盒子中小偏差的重要性。为了部分解决这个问题，我们直接预测边界框宽度和高度的平方根，而不是宽度和高度。

YOLO 每个网格单元预测多个边界框。在训练时，每个目标我们只需要一个边界框预测器来负责。我们根据哪个预测器的预测值与真实值之间具有当前最高的 IOU 来指定哪个预测器“负责”预测该目标。这导致边界框预测器之间的专一化。每个预测器可以更好地预测特定大小、长宽比或目标的类别，从而改善整体召回率。

在训练期间，我们优化以下由多部分组成的损失函数：

$$\begin{aligned}
& \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
& + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
& + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\
& + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\
& + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad (3)
\end{aligned}$$



其中  $1obj_i$  表示目标是否出现在网格单元  $i$  中， $1obj_{ij}$  表示网格单元  $i$  中的第  $j$  个边界框预测器“负责”该预测。

注意，如果目标存在于该网格单元中（前面讨论的条件类别概率），则损失函数仅惩罚分类误差。如果预测器“负责”真实边界框（即该网格单元中具有最高 IOU 的预测器），则它也仅惩罚边界框坐标误差。

我们在 Pascal VOC 2007 和 2012 的训练和验证数据集上进行了大约 135 个迭代周期的网络训练。在 Pascal VOC 2012 上进行测试时，我们的训练包含了 Pascal VOC 2007 的测试数据。在整个训练过程中，我们使用了 64 的批大小，0.9 的动量和 0.0005 的衰减率。

我们的学习率方案如下：对于第一个迭代周期，我们慢慢地将学习率从  $10^{-3}$  提高到  $10^{-2}$ 。如果我们从高学习率开始，我们的模型往往会由于梯度不稳定而发散。我们继续以  $10^{-2}$  的学习率训练 75 个迭代周期，然后用  $10^{-3}$  的学习率训练 30 个迭代周期，最后用  $10^{-4}$  的学习率训练 30 个迭代周期。

为了避免过度拟合，我们使用 dropout 和大量的数据增强。在第一个连接层之后，dropout 层使用  $rate=0.5$  的比例，防止层之间的互相适应[18]。对于数据增强，我们引入高达原始图像 20% 大小的随机缩放和转换。我们还在 HSV 色彩空间中使用高达 1.5 的因子来随机调整图像的曝光和饱和度。

### 2.3. 推断

与训练时一样，预测测试图像的检测只需要一次网络评估。在 Pascal VOC 上，每张图像上网络预测 98 个边界框（译者注：每张图

像被划分成  $7 \times 7$  的格子，每个格子预测两个边界框，总共 98 个边界框）和每个框的类别概率。YOLO 在测试时非常快，因为它只需要运行一次网络评估，不像基于分类器的方法。

网格设计强化了边界框预测中的空间多样性。通常很明显一个目标落在哪一个网格单元中，而网络只能为每个目标预测一个边界框。然而，一些大的目标或靠近多个网格单元边界的目标可以被多个网格单元很好地定位。非极大值抑制（译者注：NMS）可以用来修正这些多重检测。对于 R-CNN 或 DPM 而言，虽然非极大值抑制对性能影响不大，但会增加 2-3% 的 mAP。

## 2.4. YOLO 的缺点

YOLO 对边界框预测强加空间约束，因为每个网格单元只预测两个框，只能有一个类别。这个空间约束限制了我们的模型可以预测的邻近目标的数量。我们的模型对成群出现的小物体（比如鸟群）预测效果较差。

由于我们的模型学习从数据中预测边界框，因此它很难泛化到新的、不常见的长宽比或配置中的目标。我们的模型也使用相对较粗糙的特征来预测边界框，因为我们的架构具有来自输入图像的多个下采样层。

最后，当我们训练一个近似检测性能的损失函数时，我们的损失函数会同样的对待小边界框与大边界框的误差。大边界框的小误差通常是良性的，但小边界框的小误差对 IOU 的影响要大得多。我们的主要误差来源是定位误差。

### 3. 与其它检测系统的比较

目标检测是计算机视觉中的核心问题。检测流程通常从输入图像上提取一组鲁棒特征（Haar [25]，SIFT [23]，HOG [4]，卷积特征[6]）开始。然后，分类器[36,21,13,10]或定位器[1,32]被用来识别特征空间中的目标。这些分类器或定位器在整个图像上或在图像中的一些子区域上以滑动窗口的方式运行[35,15,39]。我们将 YOLO 检测系统与几种顶级检测框架进行比较，以突出了主要的相似性和差异性。

**可变形部件模型（DPM）。**可变形零件模型（DPM）使用滑动窗口方法进行目标检测[10]。DPM 使用不相交的流程来提取静态特征，对区域进行分类，预测高分区域的边界框等。我们的系统用单个卷积神经网络替换所有这些不同的部分。网络同时进行特征提取、边界框预测、非极大值抑制和上下文推理。网络内嵌训练特征而不是静态特征，并优化它们完成检测任务。我们的统一架构获得了比 DPM 更快、更准确的模型。

**R-CNN。**R-CNN 及其变种使用 region proposals 而不是滑动窗口来查找图像中的目标。Selective Search [35]产生潜在的边界框、卷积网络提取特征、SVM 对边界框进行评分、线性模型调整边界框、非极大值抑制消除重复检测。这个复杂流程的每个阶段都必须独立地进行精确调整，所得到的系统非常慢，测试时每张图像需要超过 40 秒 [14]。

YOLO 与 R-CNN 有一些相似之处。每个网格单元提出潜在的边界框并使用卷积特征对这些框进行评分。但是，我们的系统对网格单

元提出进行了空间限制，这有助于缓解对同一目标的多次检测。我们的系统还提出了更少的边界框，每张图像只有 98 个，而 **Selective Search** 则需要 2000 个左右。最后，我们的系统将这些单独的组件组合成一个单一的、共同优化的模型。

**其它快速检测器。**Fast 和 Faster R-CNN 通过共享计算和使用神经网络替代 **Selective Search** 来提出区域加速 R-CNN 框架[14] [28]。虽然它们提供了比 R-CNN 更快的速度和更高的准确度，但两者仍然不能达到实时性能。

许多研究工作集中在加快 DPM 流程上[31] [38] [5]。它们加速 HOG 计算，使用级联，并将计算推动到 GPU 上。但是，实际上 DPM [31]实时运行只达到 30Hz。

**YOLO** 不是试图优化大型检测流程的单个组件，而是完全抛弃流程，为提升检测速度而重新设计。

像人脸或行人等单类别的检测器可以被高度优化，因为他们只需处理更少的多样性[37]。**YOLO** 是一种通用的检测器，可以学习同时检测不同的多个目标。

**Deep MultiBox。**与 R-CNN 不同，Szegedy 等人训练了一个卷积神经网络来预测感兴趣区域(ROI)[8]，而不是使用 **Selective Search**。**MultiBox** 还可以通过用单类别预测替换置信度预测来执行单目标检测。然而，**MultiBox** 无法执行通用的目标检测，并且仍然只是一个较大的检测流程中的一部分，需要进一步的对图像块进行分类。**YOLO**

和 MultiBox 都使用卷积网络来预测图像中的边界框，但是 YOLO 是一个完整的检测系统。

**OverFeat.** Sermanet 等人训练了一个卷积神经网络来完成定位工作，并使该定位器进行检测[32]。OverFeat 高效地执行滑动窗口检测，但它仍然是一个不相交的系统。OverFeat 优化了定位，而不是检测性能。像 DPM 一样，定位器在进行预测时只能看到局部信息。OverFeat 不能推断全局上下文，因此需要大量的后处理来完成连贯的检测。

**MultiGrasp.** 我们的工作在设计上类似于 Redmon 等[27]的 grasp 检测。我们对边界框预测的网格方法是基于 MultiGrasp 系统 grasp 的回归分析。然而，grasp 检测比目标检测任务要简单得多。MultiGrasp 只需要为包含一个目标的图像预测一个可以 grasp 的区域。不必估计目标的大小、位置或目标边界或预测目标的类别，只需找到适合 grasp 的区域。YOLO 可以预测图像中多个类别的多个目标的边界框和类别概率。

## 4. 实验

首先，我们在 PASCAL VOC 2007 上比较了 YOLO 和其它的实时检测系统。为了理解 YOLO 和 R-CNN 变种之间的差异，我们探索了 YOLO 和 R-CNN 性能最高的版本之一 Fast R-CNN[14]在 VOC 2007 上错误率。根据不同的误差曲线，我们的研究显示 YOLO 可以用来重新评估 Fast R-CNN 检测，并减少背景假阳性带来的误差，从而显著提升性能。我们还展示了在 VOC 2012 上的结果，并与目前最

先进的方法比较了 mAP。最后，在两个艺术品数据集上我们显示了 YOLO 可以比其它检测器更好地泛化到新领域。

#### 4.1. 与其它实时系统的比较

目标检测方面的许多研究工作都集中在对标准检测流程[5], [38], [31], [14], [17], [28]提升速度上。然而，只有 Sadeghi 等真正研究出了一个实时运行的检测系统（每秒 30 帧或更好）[31]。我们将 YOLO 与他们 DPM 的 GPU 实现进行了比较，其在 30Hz 或 100Hz 下运行。虽然其它的研究工作没有达到实时检测的标准，我们也比较了它们的相对 mAP 和速度来检查目标检测系统中精度—性能之间的权衡。

快速 YOLO 是 PASCAL 上最快的目标检测方法；据我们所知，它是现有的最快的目标检测器。具有 52.7% 的 mAP，其精度是前人实时检测精度的两倍以上。YOLO 将 mAP 提升到 63.4% 的同时保持了实时检测的性能。

我们还使用 VGG-16 训练了 YOLO。（译者注：YOLO 使用了作者自己开发的 DarkNet 模型为 baseline）这个模型比 YOLO 更准确，但速度慢得多。这个模型可以用来与依赖于 VGG-16 的其它检测系统作比较，但由于它比实时的 YOLO 更慢，本文的剩余部分将重点放在我们更快的模型上。

Fastest DPM 可以在不牺牲太多 mAP 的情况下有效地加速 DPM，但仍然会将实时性能降低 2 倍[38]。与神经网络方法相比，DPM 相对较低的检测精度也是其限制。



减去 R 的 R-CNN 用静态边界框 proposals 取代 Selective Search [20]。虽然速度比 R-CNN 更快，但仍然达不到实时，并且由于没有好的边界框 proposals，准确性受到了严重影响。

Fast R-CNN 加快了 R-CNN 的分类阶段，但是仍然依赖 selective search，每张图像需要花费大约 2 秒来生成边界框 proposals。因此，它具有很高的 mAP，但是 0.5 fps 的速度仍离实时性很远。

最近 Faster R-CNN 用神经网络替代了 selective search 来提出边界框，类似于 Szegedy 等[8]。在我们的测试中，他们最精确的模型达到了 7fps，而较小的、不太精确的模型运行速度达到 18fps。VGG-16 版本的 Faster R-CNN 要高出 10mAP，但速度比 YOLO 慢 6 倍。Zeiler-Fergus 的 Faster R-CNN 只比 YOLO 慢了 2.5 倍，但也不太准确。

Real-Time Detectors	Train	mAP	FPS
100Hz DPM [30]	2007	16.0	100
30Hz DPM [30]	2007	26.1	30
Fast YOLO	2007+2012	52.7	155
YOLO	2007+2012	63.4	45
Less Than Real-Time			
Fastest DPM [37]	2007	30.4	15
R-CNN Minus R [20]	2007	53.5	6
Fast R-CNN [14]	2007+2012	70.0	0.5
Faster R-CNN VGG-16[27]	2007+2012	73.2	7
Faster R-CNN ZF [27]	2007+2012	62.1	18
YOLO VGG-16	2007+2012	66.4	21

表 1: Pascal VOC 2007 上的实时检测系统。比较了快速检测器的性能和速度。Fast YOLO 是 Pascal VOC 检测记录中速度最快的检测器，其精度仍然是其它实时检测器的两倍。YOLO 比快速版本更精确 10mAP，同时在速度上仍保持了实时性。

## 4.2. VOC 2007 检测误差分析

为了进一步检查 YOLO 和最先进的检测器之间的差异，我们详细分析了 VOC 2007 的结果。我们将 YOLO 与 Fast R-CNN 进行比较，因为 Fast R-CNN 是 PASCAL 上性能最高的检测器之一并且它的检测代码是可公开得到的。

我们使用 Hoiem 等人[19]的方法和工具。对于测试时的每个类别，我们只关注这个类别的前 N 个预测。每个预测要么归为正确，要么根据错误类型进行归类：

- **Correct:** 分类正确且  $\text{IOU} > 0.5$ 。
- **Localization:** 分类正确但  $0.1 < \text{IOU} < 0.5$ 。
- **Similar:** 分类的类别相似且  $\text{IOU} > 0.1$ 。
- **Other:** 类别错误， $\text{IOU} > 0.1$ 。
- **Background:** 分类为其它任何目标， $\text{IOU} < 0.1$ 。

图 4 显示了在所有的 20 个类别上每种错误类型平均值的分解图。

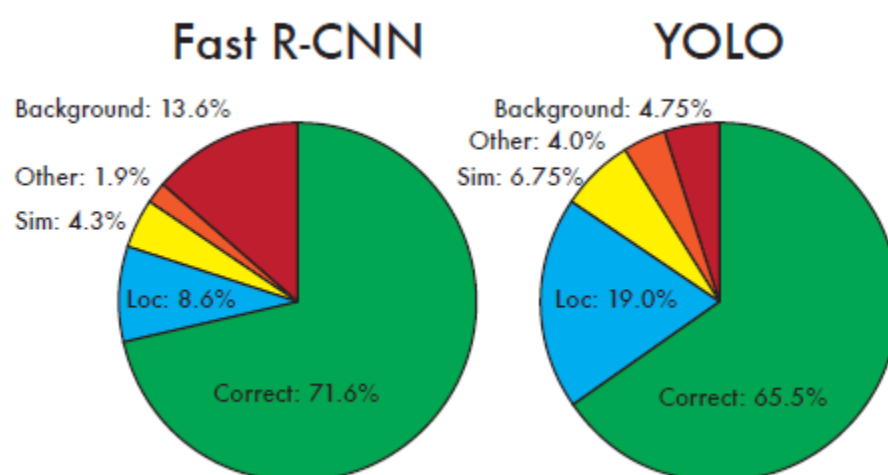


图 4，误差分析：Fast R-CNN 对比 YOLO。这些图显示了各种类别的前 N 个预测中定位错误和背景错误的百分比（N = #表示目标在那个类别中）。

YOLO 需要改善正确定位目标。定位误差占 YOLO 模型误差的大多数，比其它误差错误来源总合都多。Fast R-CNN 定位误差少很

多，但背景误差更多。它的检测结果中 13.6% 是不包含任何目标的假阳性。Fast R-CNN 与 YOLO 相比，将背景预测成目标的可能性高出近 3 倍。（译者注：根据图 4， $13.6/4.75=2.86$ ）

### 4.3. 结合 Fast R-CNN 和 YOLO

YOLO 比 Fast R-CNN 的背景误检要少得多。通过使用 YOLO 消除 Fast R-CNN 的背景检测，我们获得了显著的性能提升。对于 R-CNN 预测的每个边界框，我们检查 YOLO 是否预测一个类似的框。如果是这样，我们根据 YOLO 预测的概率和两个盒子之间的重叠来对这个预测进行改进。

最好的 Fast R-CNN 模型在 VOC 2007 测试集上达到了 71.8% 的 mAP。当与 YOLO 结合时，其 mAP 增加了 3.2% 达到了 75.0%。我们也尝试将最好的 Fast R-CNN 模型与其它几个版本的 Fast R-CNN 结合起来。这些模型组合产生了 0.3-0.6% 的小幅增加，详见表 2。

	mAP	Combined	Gain
Fast R-CNN	71.8	-	-
Fast R-CNN (2007 data)	<b>66.9</b>	72.4	.6
Fast R-CNN (VGG-M)	59.2	72.4	.6
Fast R-CNN (CaffeNet)	57.1	72.1	.3
YOLO	63.4	<b>75.0</b>	<b>3.2</b>

表 2: VOC 2007 模型组合实验。我们检验了各种模型与 Fast R-CNN 最佳版本结合的效果。Fast R-CNN 的其它版本只产生了很小的改进，而 YOLO 则提供了显著的性能提升。

来自 YOLO 的提升不仅仅是模型集成的副产品，因为组合不同版本的 Fast R-CNN 几乎没有什么改进。相反，正是因为 YOLO 在测试时出现了各种各样的误差，所以在提高 Fast R-CNN 的性能方面非常有效。

遗憾的是，这个组合并没有从 YOLO 的速度中受益，因为我们分别运行每个模型，然后将结果组合起来。但是，由于 YOLO 速度如此之快，与 Fast R-CNN 相比，不会增加任何显著的计算时间。

#### 4.4. VOC 2012 上的结果

在 VOC 2012 测试集上，YOLO 获得了 57.9% 的 mAP。这低于现有的最好技术，如表 3 所示其接近于使用 VGG-16 的原始 R-CNN。我们的系统与其最接近的竞争对手相比，需要改善在小目标上的检测。在水瓶、绵羊和电视/显示器等类别上，YOLO 的得分比 R-CNN 或 Feature Edit 低 8–10%。然而，在猫和火车等其它类别上 YOLO 实现了更高的性能。

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_MORE_DATA [11]	73.9	85.5	82.9	76.6	57.8	62.7	79.4	77.2	86.6	55.0	79.1	62.2	87.0	83.4	84.7	78.9	45.3	73.4	65.8	80.3	74.0
HyperNet_VGG	71.4	84.2	78.5	73.6	55.6	53.7	78.7	79.8	87.7	49.6	74.9	52.1	86.0	81.7	83.3	81.8	48.6	73.5	59.4	79.9	65.7
HyperNet_LSP	71.3	84.1	78.3	73.3	55.5	53.6	78.6	79.6	87.5	49.5	74.9	52.1	85.6	81.6	83.2	81.6	48.4	73.2	59.3	79.7	65.6
Fast R-CNN + YOLO	70.7	83.4	78.5	73.5	55.8	43.4	79.1	73.1	89.4	49.4	75.5	57.0	87.5	80.9	81.0	74.7	41.8	71.5	68.5	82.1	67.2
MR_CNN_S_CNN [11]	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
Faster R-CNN [27]	70.4	84.9	79.8	74.3	53.9	49.8	77.5	75.9	88.5	45.6	77.1	55.3	86.9	81.7	80.9	79.6	40.1	72.6	60.9	81.2	61.5
DEEP_ENS_COYO	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
NoC [28]	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 [14]	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
UMICH_FGS_STRUCT	66.4	82.9	76.1	64.1	44.6	49.4	70.3	71.2	84.6	42.7	68.6	55.8	82.7	77.1	79.9	68.7	41.4	69.0	60.0	72.0	66.2
NUS_NIN_C2000 [7]	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 [7]	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
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 BB [13]	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
R-CNN VGG [13]	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
YOLO	57.9	77.0	67.2	57.7	38.3	22.7	68.3	55.9	81.4	36.2	60.8	48.5	77.2	72.3	71.3	63.5	28.9	52.2	54.8	73.9	50.8
Feature Edit [32]	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
R-CNN BB [13]	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 [16]	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 [13]	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

表 3: PASCAL VOC 2012 排行榜。截至 2015 年 11 月 6 日，YOLO 与完整 comp4（允许外部数据）公开排行榜进行了比较。显示了各种检测方法的 mAP 和每类的平均精度。YOLO 是唯一达到实时的检测器。Fast R-CNN + YOLO 是评分第四高的方法，比 Fast R-CNN 提升了 2.3%。

我们联合的 Fast R-CNN + YOLO 模型是性能最高的检测方法之一。Fast R-CNN 从与 YOLO 的组合中获得了 2.3% 的提高，在公开排行榜上提升了 5 位。

#### 4.5. 泛化能力：艺术品中的行人检测

用于目标检测的学术数据集以相同分布获取训练和测试数据。在现实世界的应用中，很难预测所有可能的样本，而且测试数据可能与系统之前看到的不同[3]。（译者注：也就是说测试数据与模型训练的数据可能在风格、模式、目标表现形式等方面有很大的区别，例如模型训练时的数据是用相机拍摄的图像，而测试时用油画作品图像，此时由于油画作品比较抽象，模型就可能无法很好地识别其中的目标）我们在 Picasso 数据集上[12]和 People-Art 数据集[3]上将 YOLO 与其它的检测系统进行比较，这两个数据集常用于测试艺术品中的行人检测。

图 5 所示为 YOLO 和其它检测方法之间性能比较的结果。作为参考，我们在 person 上提供 VOC 2007 的检测 AP，其中所有模型仅在 VOC 2007 数据上训练。在 Picasso 数据集上测试的模型使用 VOC 2012 进行了训练，而 People-Art 数据集测试的模型使用 VOC 2010 进行了训练。

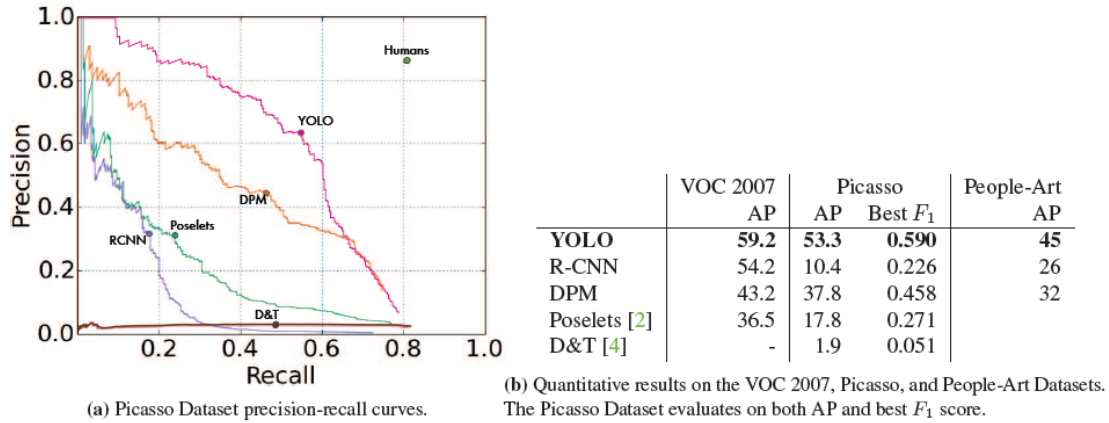


图 5: Picasso 和 People-Art 数据集上的泛化结果。

R-CNN 在 VOC 2007 上的 AP 较高。然而，当应用于艺术品时，R-CNN 明显性能下降。R-CNN 使用 Selective Search 来调整自然图像



的边界框 proposals。R-CNN 中的分类器步骤只能看到小区域，并且需要很好的边界框 proposals。

DPM 在应用于艺术品时保持了其 AP。之前的工作认为 DPM 表现良好，因为它具有目标形状和布局的强大空间模型。虽然 DPM 不会像 R-CNN 那样 AP 下降很厉害，但它的 AP 一开始就较低（译者注：DPM AP 本来就低，已经没有下降空间了，哈哈）。

YOLO 在 VOC 2007 上有很好的性能，在应用于艺术品时其 AP 下降幅度低于其它方法。像 DPM 一样，YOLO 建模了目标的大小和形状以及目标和目标通常出现的位置之间的关系。艺术品和自然图像在像素级别上有很大不同，但是它们在目标的大小和形状方面是相似的，因此 YOLO 仍然可以预测好的边界框和检测结果。

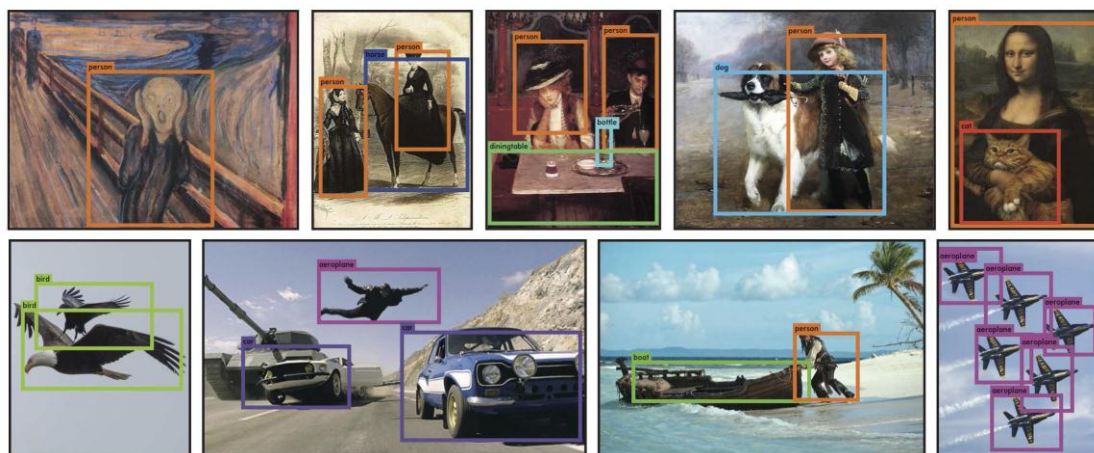


图 6：定性结果。YOLO 在源于网络的艺术品和自然图像上的运行结果。虽然它将人误检成了飞机，但它大部分上是准确的。

## 5. 现实环境下的实时检测

YOLO 是一种快速、精确的目标检测器，非常适合计算机视觉应用。我们将 YOLO 连接到网络摄像头，并验证它是否能保持实时性能，包括从摄像头获取图像并显示检测结果的时间。



最终的系统是交互式的并且是参与式的。虽然 YOLO 单独地处理图像，但当连接到网络摄像头时，其功能类似于跟踪系统，可在目标移动和外观变化时检测目标。系统演示和源代码可以在我们的项目网站上找到：<http://pjreddie.com/yolo/>。

## 6. 结论

我们介绍了 YOLO，一种统一的目标检测模型。我们的模型构建简单，可以直接在整张图像上进行训练。与基于分类器的方法不同，YOLO 直接在对应检测性能的损失函数上训练，并且整个模型统一训练。

快速 YOLO 是文献中最快的通用目的的目标检测器，YOLO 推动了实时目标检测的最新技术。YOLO 还很好地泛化到新领域，使其成为要求快速、强大目标检测应用的理想选择。

**致谢：**这项工作得到了 ONR N00014-13-1-0720、NSF IIS-1338054 和 The Allen Distinguished Investigator Award 的部分支持。

## 参考文献

- [1] M. B. Blaschko and C. H. Lampert. Learning to localize objects with structured output regression. In Computer Vision–ECCV 2008, pages 2–15. Springer, 2008. 4
- [2] L. Bourdev and J. Malik. Poselets: Body part detectors trained using 3d human pose annotations. In International Conference on Computer Vision (ICCV), 2009. 8
- [3] H. Cai, Q. Wu, T. Corradi, and P. Hall. The cross-depiction problem: Computer vision algorithms for recognising objects in artwork and in photographs. arXiv preprint arXiv:1505.00110, 2015. 7
- [4] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, volume 1, pages 886–893. IEEE, 2005. 4, 8
- [5] T. Dean, M. Ruzon, M. Segal, J. Shlens, S. Vijaya-narasimhan, J. Yagnik, et al. Fast, accurate detection of 100,000 object classes on a single machine. In Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on, pages 1814–1821. IEEE, 2013. 5

- [6] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. arXiv preprint arXiv:1310.1531, 2013. 4
- [7] J. Dong, Q. Chen, S. Yan, and A. Yuille. Towards unified object detection and semantic segmentation. In *Computer Vision–ECCV 2014*, pages 299–314. Springer, 2014. 7
- [8] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 2155–2162. IEEE, 2014. 5, 6
- [9] M. Everingham, S. M. A. Eslami, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, Jan. 2015. 2
- [10] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(9):1627–1645, 2010. 1, 4
- [11] S. Gidaris and N. Komodakis. Object detection via a multi-region & semantic segmentation-aware CNN model. CoRR, abs/1505.01749, 2015. 7
- [12] S. Ginosar, D. Haas, T. Brown, and J. Malik. Detecting people in cubist art. In *Computer Vision–ECCV 2014 Workshops*, pages 101–116. Springer, 2014. 7
- [13] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 580–587. IEEE, 2014. 1, 4, 7
- [14] R. B. Girshick. Fast R-CNN. CoRR, abs/1504.08083, 2015. 2, 5, 6, 7
- [15] S. Gould, T. Gao, and D. Koller. Region-based segmentation and object detection. In *Advances in neural information processing systems*, pages 655–663, 2009. 4
- [16] B. Hariharan, P. Arbeláez, R. Girshick, and J. Malik. Simultaneous detection and segmentation. In *Computer Vision–ECCV 2014*, pages 297–312. Springer, 2014. 7
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. arXiv preprint arXiv:1406.4729, 2014. 5
- [18] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580, 2012. 4
- [19] D. Hoiem, Y. Chodpathumwan, and Q. Dai. Diagnosing error in object detectors. In *Computer Vision–ECCV 2012*, pages 340–353. Springer, 2012. 6
- [20] K. Lenc and A. Vedaldi. R-cnn minus r. arXiv preprint arXiv:1506.06981, 2015. 5, 6
- [21] R. Lienhart and J. Maydt. An extended set of haar-like features for rapid object detection. In *Image Processing. 2002. Proceedings. 2002 International Conference on*, volume 1, pages I–900. IEEE, 2002. 4
- [22] M. Lin, Q. Chen, and S. Yan. Network in network. CoRR, abs/1312.4400, 2013. 2
- [23] D. G. Lowe. Object recognition from local scale-invariant features. In *Computer vision, 1999. The proceedings of the seventh IEEE international conference on*, volume 2, pages 1150–1157. Ieee, 1999. 4

- [24] D. Mishkin. Models accuracy on imagenet 2012 val. <https://github.com/BVLC/caffe/wiki/Models-accuracy-on-ImageNet-2012-val>. Accessed: 2015-10-2. 3
- [25] C. P. Papageorgiou, M. Oren, and T. Poggio. A general framework for object detection. In *Computer vision, 1998. sixth international conference on*, pages 555–562. IEEE, 1998. 4
- [26] J. Redmon. Darknet: Open source neural networks in c. <http://pjreddie.com/darknet/>, 2013–2016. 3
- [27] J.Redmon and A.Angelova. Real-time grasp detection using convolutional neural networks. *CoRR*, abs/1412.3128, 2014. 5
- [28] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *arXiv preprint arXiv:1506.01497*, 2015. 5, 6, 7
- [29] S. Ren, K. He, R. B. Girshick, X. Zhang, and J. Sun. Object detection networks on convolutional feature maps. *CoRR*, abs/1504.06066, 2015. 3, 7
- [30] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 2015. 3
- [31] M. A. Sadeghi and D. Forsyth. 30hz object detection with dpm v5. In *Computer Vision–ECCV 2014*, pages 65–79. Springer, 2014. 5, 6
- [32] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. *CoRR*, abs/1312.6229, 2013. 4, 5
- [33] Z.Shen and X.Xue. Do more dropouts in pool5 feature maps for better object detection. *arXiv preprint arXiv:1409.6911*, 2014. 7
- [34] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *CoRR*, abs/1409.4842, 2014. 2
- [35] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *International journal of computer vision*, 104(2):154–171, 2013. 4, 5
- [36] P. Viola and M. Jones. Robust real-time object detection. *International Journal of Computer Vision*, 4:34–47, 2001. 4
- [37] P. Viola and M. J. Jones. Robust real-time face detection. *International journal of computer vision*, 57(2):137–154, 2004. 5
- [38] J. Yan, Z. Lei, L. Wen, and S. Z. Li. The fastest deformable part model for object detection. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*, pages 2497–2504. IEEE, 2014. 5, 6
- [39] C.L.Zitnick and P.Dollár.Edgeboxes:Locating object proposals from edges. In *Computer Vision–ECCV 2014*, pages 391–405. Springer, 2014. 4