

YOLOv3: An Incremental Improvement全文翻译

**Amusi**

微信公众号: CVer

[关注他](#)

185 人赞同了该文章

原标题: YOLOv3: An Incremental Improvement

原作者: Joseph Redmon Ali Farhadi

翻译者: Amusi

YOLO官网: [YOLO: Real-Time Object Detection](https://pjreddie.com/detection/yolo/)论文链接: pjreddie.com/media/fileYoutube: youtube.com/watch?知乎话题: [如何评价YOLOv3: An Incremental Improvement?](#)

Amusi是一名CV初学者, 论文翻译中用到了Google, 并自己逐句检查过, 但还是会有显得晦涩的地方, 如有语法/专业名词翻译错误, 还请见谅, 并欢迎及时指出。

Abstract

我们给YOLO提供一些更新! 我们做了一些小的设计更改以使其更好。我们也训练了这个非常好的新网络。它比上次(YOLOv2)稍大一些, 但更准确。它仍然很快, 所以不用担心。在 320×320 YOLOv3运行22.2ms, 28.2 mAP, 像SSD一样准确, 但速度快三倍。当我们看看以老的0.5 IOU mAP检测指标时, YOLOv3是相当不错的。在Titan X上, 它在51 ms内实现了57.9的AP50, 与RetinaNet在198 ms内的57.5 AP50相当, 性能相似但速度快3.8倍。与往常一样, 所有代码均在pjreddie.com/yolo/。

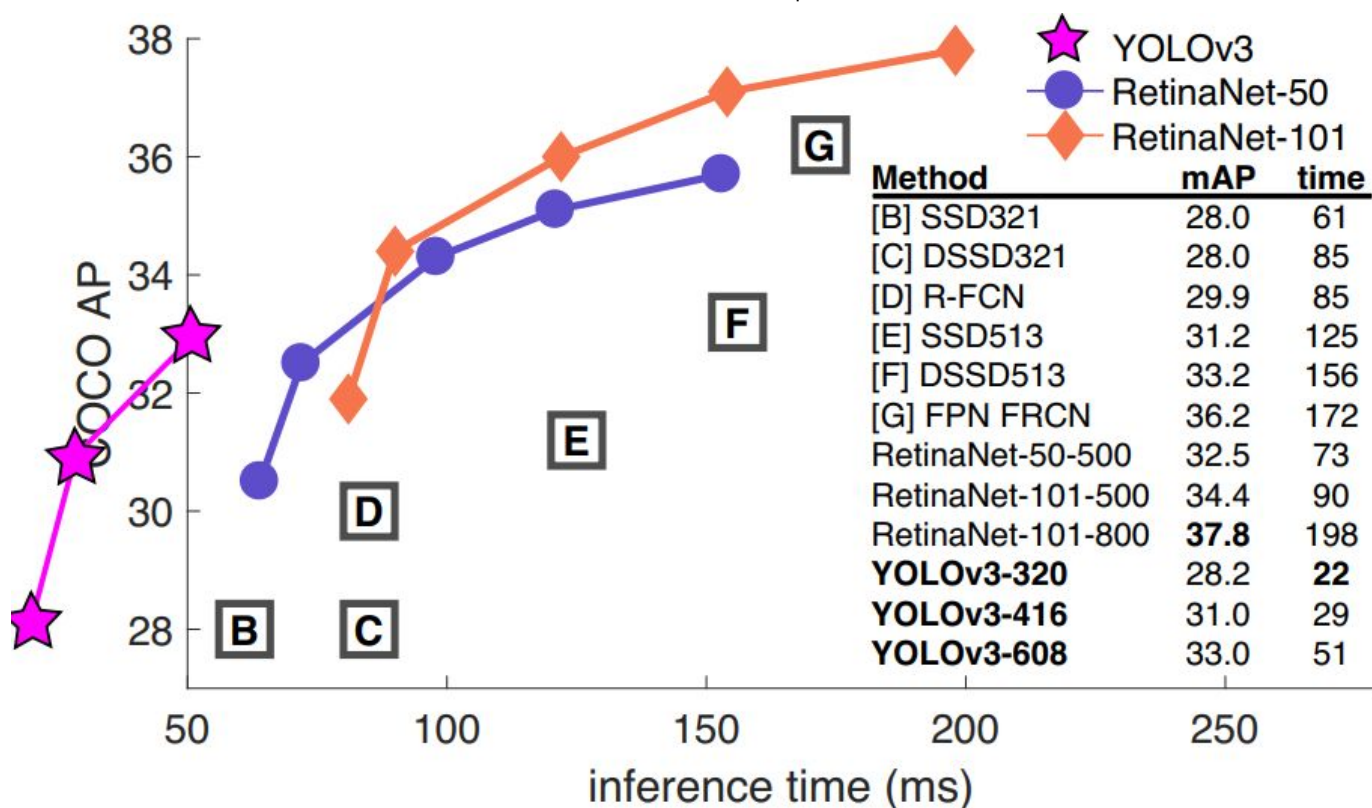


Figure 1. We adapt this figure from the Focal Loss paper [7]. YOLOv3 runs significantly faster than other detection methods with comparable performance. Times from either a K40 or Titan X, they are basically the same GPU.

1 Introduction

有时候，一年你主要只是在打电话，你知道吗？今年我没有做很多研究。我在Twitter上花了很多时间。玩了一下GAN。去年我留下了一点点的动力[10] [1]；我设法对YOLO进行了一些改进。但是诚然，没有什么比这超级有趣的了，只是一小堆 (bunch) 改变使它变得更好。我也帮助了其他人的做一些研究。

其实，这就是今天带给我们的。我们有一个camera-ready deadline，we need to cite some of the random updates I made to YOLO but we don't have a source。所以为技术报告做准备！

关于技术报告的好处是他们不需要介绍，你们都知道我们为什么来到这里。因此，这篇介绍性文章的结尾将为本文的其余部分提供signpost。首先我们会告诉你YOLOv3的详细内容。然后我们会告诉你我们是怎么做的。我们还会告诉你我们尝试过的一些没有奏效的事情。最后，我们将考虑这一切意味着什么。

2 The Deal

这里是YOLOv3的详细内容：我们主要从其他人那里获得好点子。我们也训练了一个比其他人更好的新分类器网络。我们将从头开始介绍整个系统，以便您能够理解这一切。

2.1 Bounding Box Prediction

在YOLO9000之后，我们的系统使用维度聚类（dimension clusters）作为anchor boxes来预测边界框[13]。网络为每个边界框预测4个坐标， t_x , t_y , t_w , t_h 。如果单元格从图像的左上角偏移（ c_x ; c_y ），并且之前的边界框具有宽度和高度 p_w , p_h ，则预测对应于：

$$b_x = \sigma(t_x) + c_x$$

$$b_y = \sigma(t_y) + c_y$$

$$b_w = p_w e^{t_w}$$

$$b_h = p_h e^{t_h}$$

在训练期间，我们使用平方误差损失的总和。如果对于一些坐标预测的ground truth是 \tilde{t}_* （原文中是使用^符号，但某乎自带的代码工具打不出来，所以这里我用~符号替代），我们的梯度是ground truth（由ground box计算得到）减去我们的预测： $\tilde{t}_* - t_*$ 。通过反转上面的方程可以很容易地计算出这个ground truth。

YOLOv3使用逻辑回归预测每个边界框（bounding box）的对象分数。如果先前的边界框比之前的任何其他边界框重叠ground truth对象，则该值应该为1。如果以前的边界框不是最好的，但是确实将ground truth对象重叠了一定的阈值以上，我们会忽略这个预测，按照[15]进行。我们使用阈值0.5。与[15]不同，我们的系统只为每个ground truth对象分配一个边界框。如果先前的边界框未分配给grounding box对象，则不会对坐标或类别预测造成损失。

2.2 Class Prediction

每个框使用多标签分类来预测边界框可能包含的类。我们不使用softmax，因为我们发现它对于高性能没有必要，相反，我们只是使用独立的逻辑分类器。在训练过程中，我们使用二元交叉熵损失来进行类别预测。

这个公式有助于我们转向更复杂的领域，如Open Image Dataset[5]。在这个数据集有许多重叠的标签（如女性和人物）。使用softmax会强加了一个假设，即每个框中只有一个类别，但通常情况并非如此。多标签方法更好地模拟数据。

2.3 Prediction Across Scales

YOLOv3预测3种不同尺度的框 (boxes)。我们的系统使用类似的概念来提取这些尺度的特征，以形成金字塔网络[6]。从我们的基本特征提取器中，我们添加了几个卷积层。其中最后一个预测了3-d张量编码边界框，对象和类别预测。在我们的COCO实验[8]中，我们预测每个尺度的3个框，所以对于4个边界框偏移量，1个目标性预测和80个类别预测，张量为 $N \times N \times [3 * (4 + 1 + 80)]$ 。

接下来，我们从之前的两层中取得特征图 (feature map)，并将其上采样2倍。我们还从网络中的较早版本获取特征图，并使用element-wise addition将其与我们的上采样特征进行合并。这种方法使我们能够从早期特征映射中的上采样特征和更细粒度的信息中获得更有意义的语义信息。然后，我们再添加几个卷积层来处理这个组合的特征图，并最终预测出一个相似的张量，虽然现在是两倍的大小。

我们再次执行相同的设计来预测最终尺度的方框。因此，我们对第三种尺度的预测将从所有先前的计算中获益，并从早期的网络中获得细粒度的特征。

我们仍然使用k-means聚类来确定我们的边界框的先验。我们只是选择了9个聚类 (clusters) 和3个尺度 (scales)，然后在整个尺度上均匀分割聚类。在COCO数据集上，9个聚类是：

(10×13) ; (16×30) ; (33×23) ; (30×61) ; (62×45) ; (59×119) ; (116×90) ;
(156×198) ; (373×326) 。

2.4 Feature Extractor

我们使用新的网络来实现特征提取。我们的新网络是用于YOLOv2，Darknet-19中的网络和那些新颖的残差网络的混合方法。我们的网络使用连续的 3×3 和 1×1 卷积层，但现在也有一些shortcut连接，该网络明显更大。它有53个卷积层，所以我们称之为..... Darknet-53！（手动滑稽...）

	Type	Filters	Size	Output
	Convolutional	32	3×3	256×256
	Convolutional	64	$3 \times 3 / 2$	128×128
1x	Convolutional	32	1×1	
	Convolutional	64	3×3	
	Residual			128×128
	Convolutional	128	$3 \times 3 / 2$	64×64
2x	Convolutional	64	1×1	
	Convolutional	128	3×3	
	Residual			64×64
	Convolutional	256	$3 \times 3 / 2$	32×32
8x	Convolutional	128	1×1	
	Convolutional	256	3×3	
	Residual			32×32
	Convolutional	512	$3 \times 3 / 2$	16×16
8x	Convolutional	256	1×1	
	Convolutional	512	3×3	
	Residual			16×16
	Convolutional	1024	$3 \times 3 / 2$	8×8
4x	Convolutional	512	1×1	
	Convolutional	1024	3×3	
	Residual			8×8
	Avgpool		Global	
	Connected		1000	
	Softmax			

Table 1. **Darknet-53.**

这个新网络比Darknet-19功能强大得多，而且比ResNet-101或ResNet-152更有效。以下是一些ImageNet结果：

Backbone	Top-1	Top-5	Bn Ops	BFLOP/s	FPS
Darknet-19 [13]	74.1	91.8	7.29	1246	171
ResNet-101[3]	77.1	93.7	19.7	1039	53
ResNet-152 [3]	77.6	93.8	29.4	1090	37
Darknet-53	77.2	93.8	18.7	1457	78

Table 2. Comparison of backbones. Accuracy, billions of operations, billion floating point operations per second, and FPS for various networks.

每个网络都使用相同的设置进行训练，并以 256×256 的单精度测试进行测试。运行时间是在Titan X上以 256×256 进行测量的。因此，Darknet-53可与state-of-the-art的分类器相媲美，但浮点运算更少，速度更快。Darknet-53比ResNet-101更好，速度更快1：5倍。Darknet-53与ResNet-152具有相似的性能，速度提高2倍。

Darknet-53也可以实现每秒最高的测量浮点运算。这意味着网络结构可以更好地利用GPU，从而使其评估效率更高，速度更快。这主要是因为ResNets的层数太多，效率不高。

2.5 Training

我们仍然训练完整的图像，没有hard negative mining or any of that stuff。我们使用多尺度训练，大量的data augmentation，batch normalization，以及所有标准的东西。我们使用Darknet神经网络框架进行训练和测试[12]。

3 How We Do

YOLOv3非常好！请参见表3。就COCO的mAP指标而言，它与SSD variants相当，但速度提高了3倍。尽管如此，它仍然比像RetinaNet这样的其他模型落后很多。

	backbone	AP	AP ₅₀	AP ₇₅	AP _S	AP _M	AP _L
<i>Two-stage methods</i>							
Faster R-CNN+++ [3]	ResNet-101-C4	34.9	55.7	37.4	15.6	38.7	50.9
Faster R-CNN w FPN [6]	ResNet-101-FPN	36.2	59.1	39.0	18.2	39.0	48.2
Faster R-CNN by G-RMI [4]	Inception-ResNet-v2 [19]	34.7	55.5	36.7	13.5	38.1	52.0
Faster R-CNN w TDM [18]	Inception-ResNet-v2-TDM	36.8	57.7	39.2	16.2	39.8	52.1
<i>One-stage methods</i>							
YOLOv2 [13]	DarkNet-19 [13]	21.6	44.0	19.2	5.0	22.4	35.5
SSD513 [9, 2]	ResNet-101-SSD	31.2	50.4	33.3	10.2	34.5	49.8
DSSD513 [2]	ResNet-101-DSSD	33.2	53.3	35.2	13.0	35.4	51.1
RetinaNet [7]	ResNet-101-FPN	39.1	59.1	42.3	21.8	42.7	50.2
RetinaNet [7]	ResNeXt-101-FPN	40.8	61.1	44.1	24.1	44.2	51.2
YOLOv3 608 × 608	Darknet-53	33.0	57.9	34.4	18.3	35.4	41.9

Table 3. I'm seriously just stealing all these tables from [7] they take soooo long to make from scratch. Ok, YOLOv3 is doing alright. Keep in mind that RetinaNet has like $3.8 \times$ longer to process an image. YOLOv3 is much better than SSD variants and comparable to state-of-the-art models on the AP₅₀ metric.

然而，当我们在 $\text{IOU} = 0.5$ （或者图表中的AP50）看到mAP的“旧”检测度量时，YOLOv3非常强大。它几乎与RetinaNet相当，并且远高于SSD variants。这表明YOLOv3是一个非常强大的检测器，擅长为目标生成像样的框（boxes）。However, performance drops significantly as the IOU threshold increases indicating YOLOv3 struggles to get the boxes perfectly aligned with the object.

在过去，YOLO在小目标的检测上表现一直不好。然而，现在我们看到了这种趋势的逆转。随着新的多尺度预测，我们看到YOLOv3具有相对较高的APS性能。但是，它在中等和更大尺寸的物体上的表现相对较差。需要更多的研究来达到这个目的。当我们在AP50指标上绘制精确度和速度时（见图3），我们看到YOLOv3与其他检测系统相比具有显著的优势。也就是说，速度越来越快。

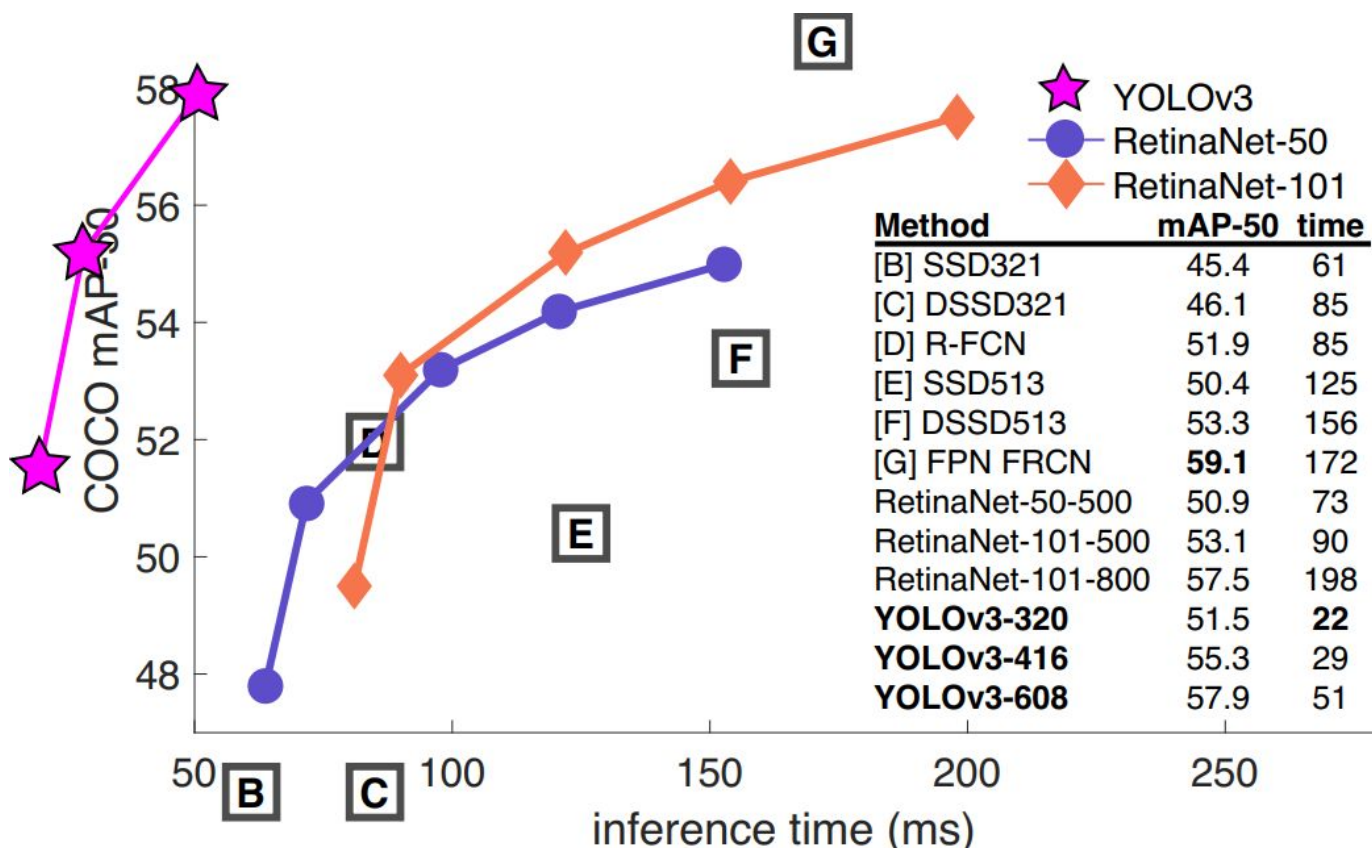


Figure 3. Again adapted from the [7], this time displaying speed/accuracy tradeoff on the mAP at .5 IOU metric. You can tell YOLOv3 is good because it's very high and far to the left. Can you cite your own paper? Guess who's going to try, this guy → [14].

注：这个图只能用Amazing来概括！！！！

4 Things We Tried That Didn't Work

我们在研究YOLOv3时尝试了很多东西。很多都不起作用。这是我们可以记住的东西。

Anchor box x , y offset predictions. 我们尝试使用正常anchor box预测机制，这里你使用线性激活来预测 x , y offset作为box的宽度或高度的倍数。我们发现这种方法降低了模型的稳定性，并且效果不佳。

Linear x, y predictions instead of logistic. 我们尝试使用线性激活来直接预测 x, y offset而不是逻辑激活。这导致mAP下降了几个点。

Focal loss. 我们尝试使用focal loss。它使得mAp降低了2个点。YOLOv3对focal loss解决的问题可能已经很强大，因为它具有单独的对象预测和条件类别预测。因此，对于大多数例子来说，类别预测没有损失？或者其他的东西？我们并不完全确定。

Dual IOU thresholds and truth assignment。Faster R-CNN在训练期间使用两个IOU阈值。如果一个预测与ground truth重叠达到0.7，它就像是一个正样本，如果达到0.3-0.7，它被忽略，如果小于0.3，这是一个负样本的例子。我们尝试了类似的策略，但无法取得好成绩。

我们非常喜欢我们目前的表述，似乎至少在局部最佳状态。有些技术可能最终会产生好的结果，也许他们只是需要一些调整来稳定训练。

5 What This All Means

YOLOv3是一个很好的检测器。速度很快，很准确。COCO平均AP介于0.5和0.95 IOU指标之间的情况并不如此。但是，对于检测度量0.5 IOU来说非常好。

为什么我们要改变指标？最初的COCO论文只是含有这个神秘的句子：“一旦评估服务器完成，就会添加完整的评估指标的讨论”。Russakovsky等人报告说，人类很难区分IOU为0.3还是0.5。

“训练人们目视检查一个IOU值为0.3的边界框，并将它与IOU 0.5区分开来是一件非常困难的事情。” [16]如果人类很难区分这种差异，那么它有多重要？

但是也许更好的问题是：“现在我们有了这些检测器 (detectors)，我们要做什么？”很多做这项研究的人都在Google和Facebook上。我想至少我们知道这项技术是非常好的，绝对不会被用来收集您的个人信息，并将其出售给.....等等，您是说这就是它的用途？

那么其他大量资助视觉研究的人都是军人，他们从来没有做过任何可怕的事情，例如用新技术杀死很多人哦等等.....

我有很多希望，大多数使用计算机视觉的人都是做的快乐，研究了很多好的应用，比如计算一个国家公园内的斑马数量[11]，或者追踪它们在它们周围徘徊时的猫[17]。但是计算机视觉已经被用于可疑的应用，作为研究人员，我们有责任至少考虑我们的工作可能会造成的伤害，并考虑如何减轻它的影响。我们非常珍惜这个世界。（作者走心了.....）

最后，不要@我。（因为我终于退出了Twitter）。

创新点

使用金字塔网络

用逻辑回归替代softmax作为分类器

Darknet-53

不足

速度确实快了，但mAP没有明显提升，特别是IOU > 0.5时。

参考文献

[1] Analogy. *Wikipedia*, Mar 2018. 1

[2] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg. Dssd: Deconvolutional single shot detector. *arXiv preprint arXiv:1701.06659*, 2017. 3

[3] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3

[4] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. 3

[5] I. Krasin, T. Duerig, N. Alldrin, V. Ferrari, S. Abu-El-Haija,

A. Kuznetsova, H. Rom, J. Uijlings, S. Popov, A. Veit,

S. Belongie, V. Gomes, A. Gupta, C. Sun, G. Chechik, D. Cai, Z. Feng, D. Narayanan, and K. Murphy. Openimages: A public dataset for large-scale multi-label and multi-class image classification. *Dataset available from github.com/openimages*, 2017. 2

[6] T.-Y. Lin, P. Dollar, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017. 2, 3

[7] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar. Focal loss for dense object detection. *arXiv preprint arXiv:1708.02002*, 2017. 1, 3, 4

- [8] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 2
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 3
- [10] I. Newton. *Philosophiae naturalis principia mathematica*. William Dawson & Sons Ltd., London, 1687. 1
- [11] J. Parham, J. Crall, C. Stewart, T. Berger-Wolf, and D. Rubenstein. Animal population censusing at scale with citizen science and photographic identification. 2017. 4
- [12] J. Redmon. Darknet: Open source neural networks in c. pjreddie.com/darknet/, 2013–2016. 3
- [13] J. Redmon and A. Farhadi. Yolo9000: Better, faster, stronger. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 6517–6525. IEEE, 2017. 1, 2, 3
- [14] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *arXiv*, 2018. 4
- [15] 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. 2
- [16] O. Russakovsky, L.-J. Li, and L. Fei-Fei. Best of both worlds: human-machine collaboration for object annotation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2121–2131, 2015. 4
- [17] M. Scott. Smart camera gimbal bot scanlime:027, Dec 2017.
- 4
- [18] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta. Beyond skip connections: Top-down modulation for object detection. *arXiv preprint arXiv:1612.06851*, 2016. 3
- [19] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. 2017. 3