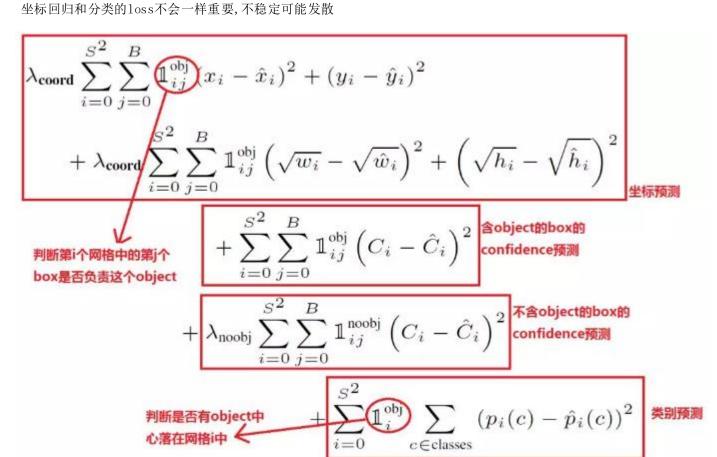
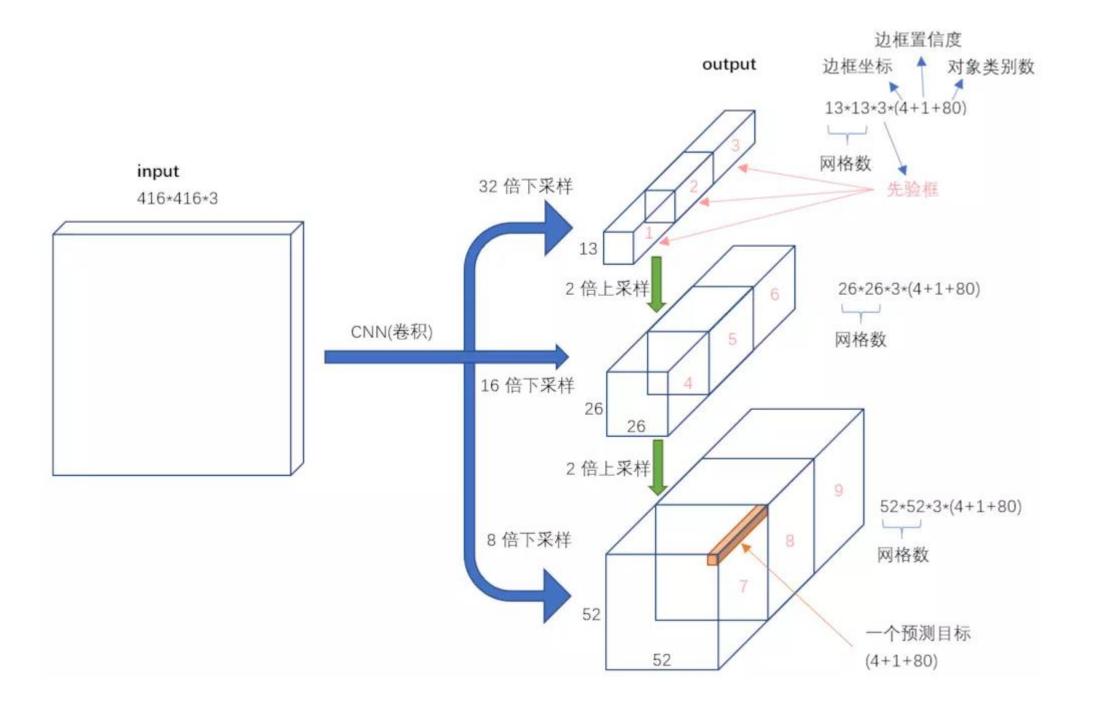
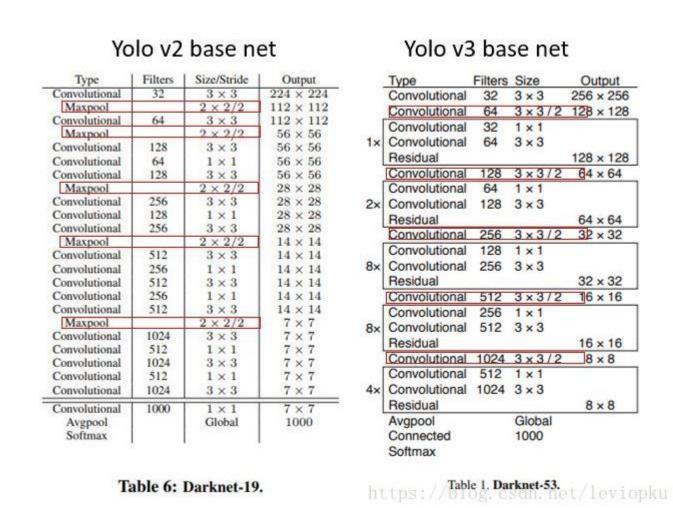


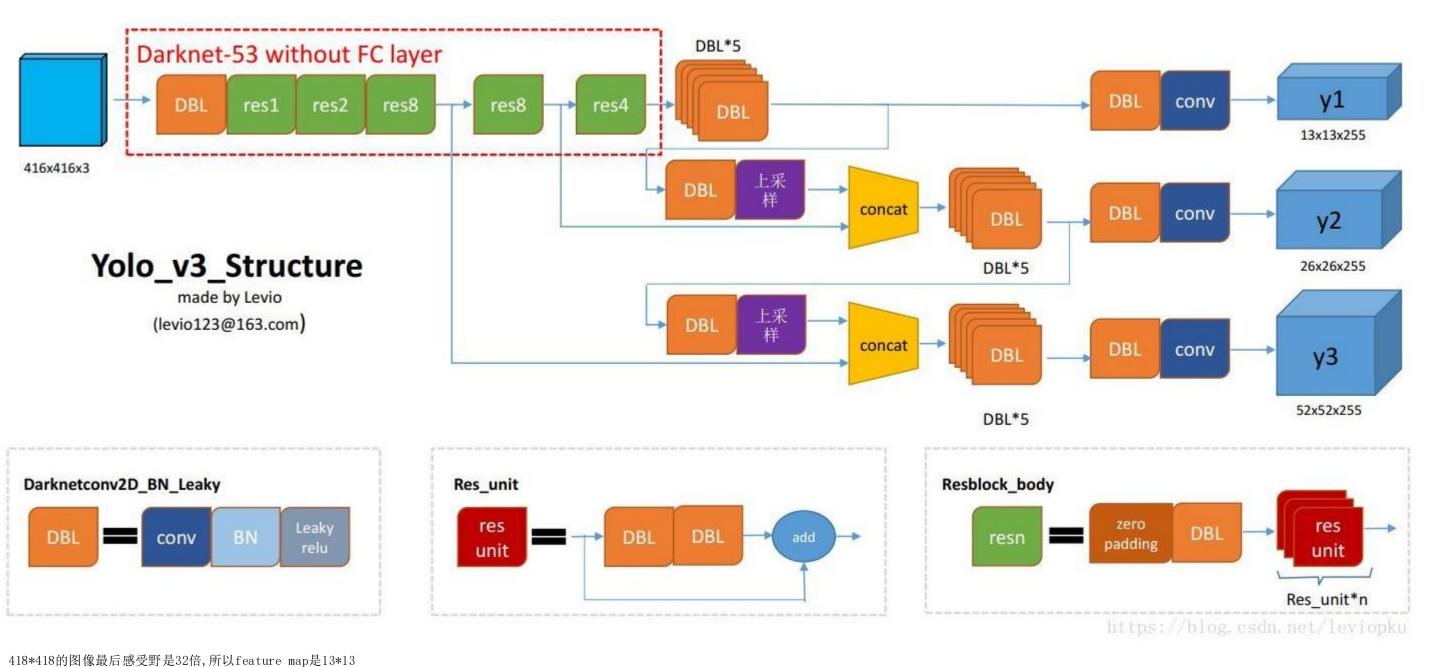
448*448->7*7 每个pixel都要回归预测B=2(2个框)以及两个框分别的conf,以及分类的prob(voc中cls=20)7*7 * (2*(4+1) + 20) = 7*7 * 30







张量的尺寸变换是通过改变卷积核的步长来实现的



上采样的方法来实现这种多尺度的feature map=>13*13的yolo 26*26的yolo 52*52的yolo 992*992的图像最后感受野是32倍,所以feature map是31*13 上采样的方法来实现这种多尺度的feature map=>31*31的yolo 62*62的yolo 124*124的yolo 若每个yolo层预测的anchor(3)个数*(cls_prob(coco:80)+bbox(4)+conf(1))=3*85=255 13*13 yolo层的每个pixe1*255 = 13*13 * 255 若是31*31的yolo就是31*31 * 255 26*26 yolo层的每个pixe1*255 = 26*26 * 255 52*52 yolo层的每个pixel*255 = 52*52 * 255 若每个yolo层预测的anchor(2)个数*(cls_prob(2)+bbox(4)+conf(1))=2*7=14 13*13 yolo层的每个pixe1*14= 13*13 * 14 若是31*31的yolo就是31*31 * 14 26*26 yolo层的每个pixel*14= 26*26 * 14 52*52 yolo层的每个pixel*14= 52*52 * 14 其中每一层yolo的anchor的size是自己取的(通过聚类),采用多尺度来对不同size的目标进行检测,越精细的grid cell就可以检测出越精细的物体 就是最后越小的feature map适用最大的anchor,因为最小的feature map感受野最大 可以自己换backbone(densenet等) 或是根据实际需要 修正darknet-53 anchor (3) 个数就意味着418*418的每张图会生成13*13*3+26*26*3+52*52*3 = 10647个bbox, 992*992的每张图会生成31*31*3+62*62*3+124*124*3 = 60543个bbox 一张图中真正的object是少数, 所以真正能预测object的bbox只占极少数, 大部分bbox都是BG, 造成了数据极大的不平衡 one-stage不如two-stage(faster rcnn)的部分就是two-stage有个rpn层会提前过滤大多数无用的BG部分,并保证正例和负例数据的平衡 所有就有后来的RetinaNet提出的Focal Loss