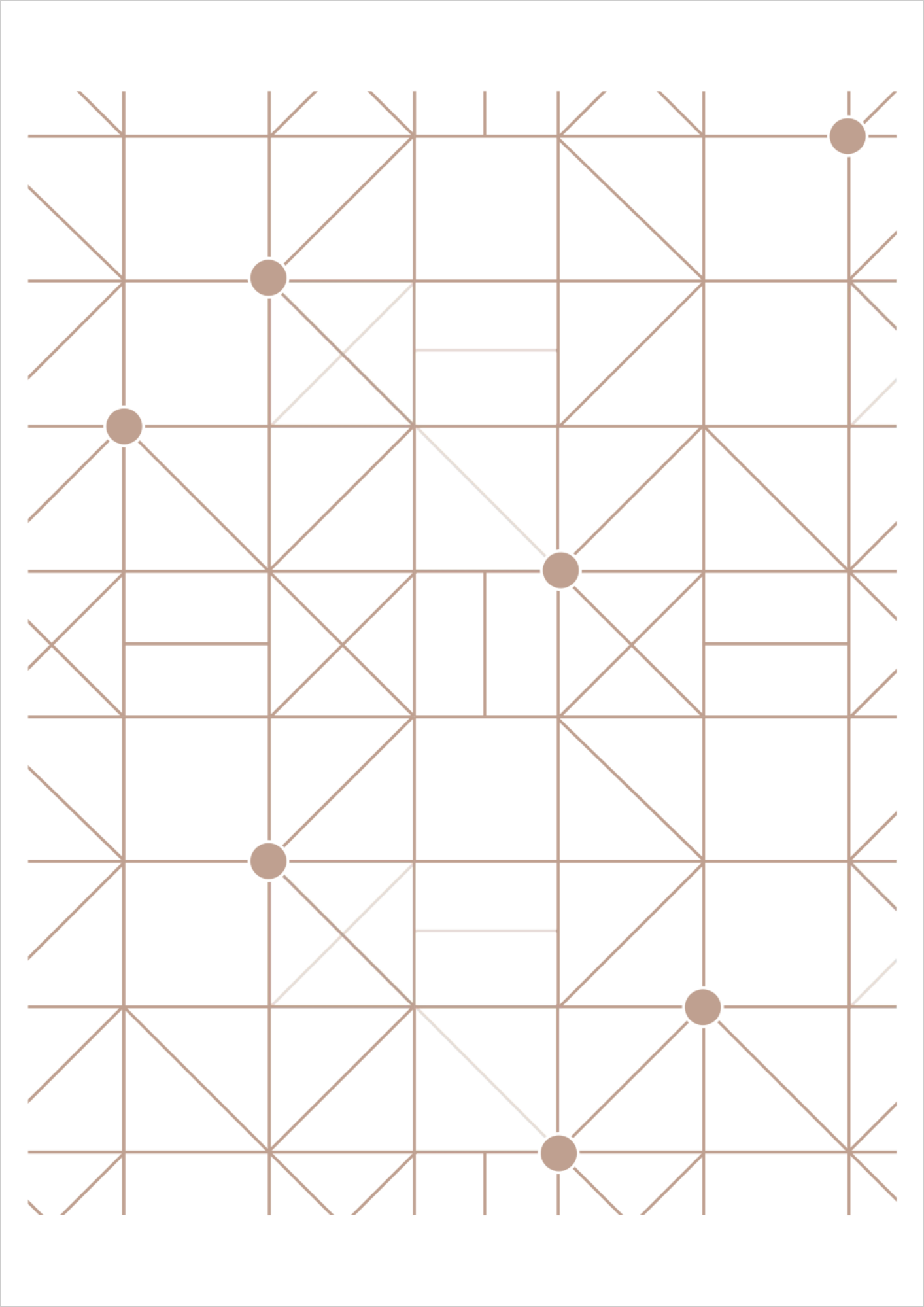


# 目标检测论文笔记

2022年5月17日 10:37



## NETNet

问题

① 物体尺度大小与标注多寡

② 大物体和新加显著区域会占据更多算力

原因

浅层特征中信息量更丰富，但浅层中低信噪比在大网络，对小物体检测不利 (FP)  
深层特征中，不同层的显著区域产生干扰，导致不同尺度的重复检测 (FP)

Scale-aware

ist. As shown in Fig. 1, in the shallow features (b) used for detecting small objects, the large-object features dominate the main saliency, weakening the small-object features and thus preventing the detection of small objects (e.g., the sports ball from (a) is not detected in the final result). Additionally, some parts of large objects have strong response regions on shallow features. For example, the head region in Fig. 1(e) is highlighted in (f), which leads to the wrongly detection of the head region. Thus, the features are scale-

## 解决思路

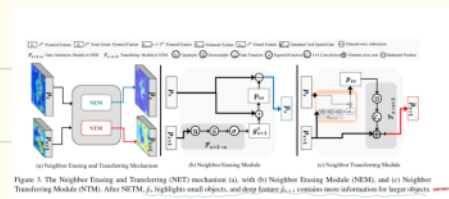
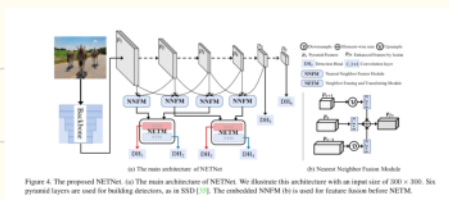
With this observation, we propose a new **Neighbor Erasing and Transferring (NET)** mechanism to rectify the pyramid features and explore scale-aware features. In **NET**, a **Neighbor Erasing Module (NEM)** is designed to **erase** the **saliency** features of large objects and **emphasize** the features of small objects in shallow layers. A **Neighbor Transferring Module (NTM)** is introduced to transfer the **erased** features and **highlight** large objects in deep layers. With this mechanism, a single-shot network called **NETNet** is constructed for scale-aware object detection. In addition, we propose to aggregate nearest neighboring pyramid features to enhance our **NET**. **NETNet** achieves 38.5% AP at a speed of 27 FPS and 32.0% AP at a speed of 55 FPS on MS COCO dataset. As a result, **NETNet** achieves a better trade-off for real-time and accurate object detection.

① 建立 NET 机制，解决 Scale-aware features

② NTM：聚合 feature feature

③ NEM：detect and erase 大物体，以调整 focus

④ NTM：将浅层中 erase and erase 的大物体特征，聚合到深层特征网络中



问题

NEM和NTM：原理？如何从海量数据中筛选出关键信息？

① anchor-free detector 网络，在检测时更灵活

NTM原理？

## DSSD

**相关问题** 小图卷积 (SSD)

缺少上下文信息组合。

**方法** 针对 SSD, 利用 ResNet 代替 VGG16

利用位置敏感卷积 (Spatial Pyramid Pooling) 定位特征图。

在特征图提取后, 使用 ResNet 模块对特征图进行处理。

**注意** 数据集在输入时, 使用的是 `convnet_v1c_products`

train, policy the same as SSD

通过调用 `mean` 函数在 `poscal_v1c_2011` 上训练, 最终 `anchors` 比例设置为  $(1, 2, 3, 4, 1)$

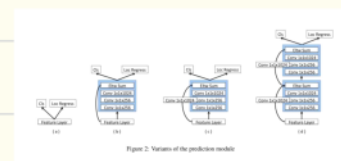
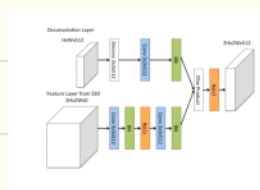
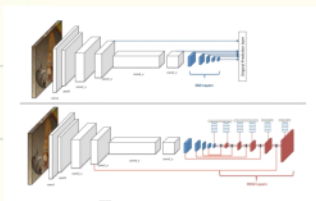
**结论** 精度相比 SSD 更高

特征提取 ResNet-101 特征图, 规模大

增加 3 个卷积核提取特征, 信息冗余

使用 3 个卷积核提取。

## 网络结构



**讨论** 可能不同? (只是同一种类型的卷积, 但卷积核的大小和步长不同)

对于不同的卷积核, 使用不同的卷积核, 使用不同的卷积核。

## DSOD

limited receptive depth 问题  
使用 pre-convolutional module, 限制网络的有效深度;

learning loss  
多层网络与浅层网络的 loss function 权重设置有所不同, 避免互相干扰

Dimensional reduction  
网络中层的特征图尺寸过大 (因为图像, depth image)

解决方法  
就是, 将 110 个 40x40 的特征图按照顺序分成 5 组

在每层提取特征并, 利用 DenseBlock 代替 VGG (Deep Superpixels)

200 个特征图, 通过 DenseBlock 提取, 一半特征, 一半卷积

将 DenseBlock 层中的特征图按顺序, 每隔 10 个 DenseBlock 提取一次 [7x7 conv 64x2 + 3x3 maxpooling 64x2]  $\Rightarrow$  1x64 conv 64x2 + 3x3 conv 64x1 + 5x64 conv 64x1 + 2x2 maxpooling 64x2

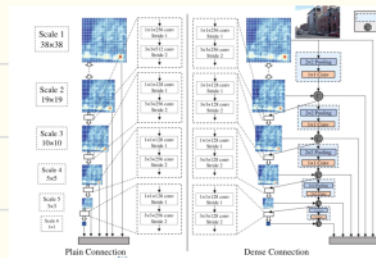
总结  
DSOD 模型与数量无关

DSOD 从 1 开始训练, 效果与使用 pre-convolutional 模型相近, 甚至更好。

## 网络结构:

Layers		Output Size (Input 3 × 300 × 300)	DSOD
Stem	Convolution	64 × 150 × 150	3 × 3 conv, stride 2
	Convolution	64 × 150 × 150	3 × 3 conv, stride 1
	Convolution	128 × 150 × 150	3 × 3 conv, stride 1
	Pooling	128 × 75 × 75	2 × 2 max pool, stride 2
Dense Block (1)		416 × 75 × 75	1 × 1 conv 3 × 3 conv × 6
	Transition Layer (1)	416 × 75 × 75	1 × 1 conv
Dense Block (2)		800 × 38 × 38	1 × 1 conv 3 × 3 conv × 8
	Transition Layer (2)	800 × 38 × 38	1 × 1 conv
Dense Block (3)		800 × 19 × 19	2 × 2 max pool, stride 2 1 × 1 conv × 8
	Transition Layer (3)	1184 × 19 × 19	1 × 1 conv
Dense Block (4)		1568 × 19 × 19	1 × 1 conv 3 × 3 conv × 8
	Transition Layer (4)	1568 × 19 × 19	1 × 1 conv
DSOD Prediction Layers		-	Plain/Dense

Table 1: DSOD architecture (growth rate  $k = 40$  in each dense block).





## RSSD

问题: SSD 是反特征图级联, 逐步不同特征图间建立

小目标, 重要规则

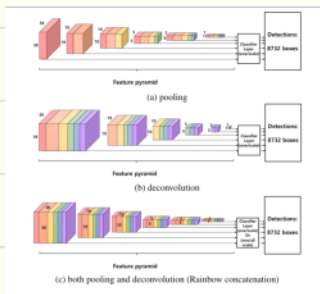
解决为: 通过 pooling 和 deconvolution 方法抽取每一层特征图的特征 (channels) (映射层关系图)

后在用于检测的每个特征图, 因此可以不需要检测 (可以投入检测的泛化能力, 避免重复计算降低运算效率, 避免检测不理想的性能表现)

结论: 通过建立不同层之间特征图的关系, 避免重复计算的问题, 自然提高检测精度

使用多个特征图检测, 精度提高, 但增加了模型的复杂度

网络结构:



问题: channel 个, 精度, why?

## CenterNet (Object as point)

**问题** 传统的目标检测器大多都依靠目标潜在框(先验框),费时,效率不高,还需要做额外的后处理(NMS)

传统检测器下采样因子为16,目标细节信息丢失严重

**思想** 将目标  $\Rightarrow$  point (Center point of bounding box)

利用 keypoint 预测 center point, 然后回归目标其他属性 (size, orientation...)

## CenterNet 与 anchor-based 区别:

Our approach is closely related to anchor-based one-stage approaches [33, 36, 43]. A center point can be seen as a single shape-agnostic anchor (see Figure 3). However, there are a few important differences. First, our CenterNet assigns the "anchor" based solely on location, not box overlap [18]. We have no manual thresholds [18] for foreground and background classification. Second, we only have one positive "anchor" per object, and hence do not need Non-Maximum Suppression (NMS) [2]. We simply extract local peaks in the keypoint heatmap [4, 39]. Third, CenterNet uses a larger output resolution (output stride of 4) compared to traditional object detectors [21, 22] (output stride of 16). This eliminates the need for multiple anchors [47].

**方法** 输入:  $I \in \mathbb{R}^{W \times H \times 3}$

keypoint-heatmap  $\begin{cases} \hat{Y} \in \mathbb{R}^{\frac{W}{K} \times \frac{H}{K} \times C}, \hat{Y} \in [0, 1] \\ \hat{Y}_{x,y,C} = 1 \text{ (keypoint)} \quad \hat{Y}_{x,y,C} = 0 \text{ (background)} \end{cases}$

ground-truth  $\begin{cases} p = (x, y), \text{类别为 } C \Rightarrow \hat{p} = \lfloor \frac{x}{K}, \frac{y}{K} \rfloor \\ \text{ground-truth heatmap: } Y \in \mathbb{R}^{\frac{H}{K} \times \frac{W}{K} \times C} \in [0, 1] \text{ 通过高斯核 } Y_{x,y,C} = \exp\left(-\frac{(x-\tilde{p}_x)^2 + (y-\tilde{p}_y)^2}{2\sigma^2}\right) \end{cases}$

预测偏移:  $\hat{O} \in \mathbb{R}^{\frac{W}{K} \times \frac{H}{K} \times 2}$

预测 size:  $\hat{S} \in \mathbb{R}^{\frac{W}{K} \times \frac{H}{K} \times 2}$

损失函数: keypoint: Focal Loss offset:  $L_1$  loss size:  $L_1$  Loss

对于每个  $(x, y)$  预测  $1C+4$  个值



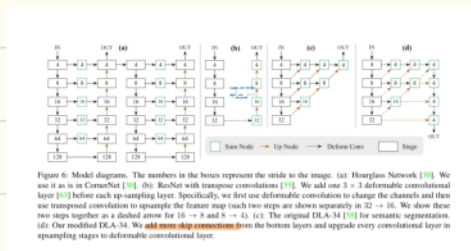
预测:

- ① 对每一类  $C$ , 设其最大期望所求点 (peak)  $(Y_{\max}, a_c)$   
 ② 对每一类  $C$ , 设其有  $n_c$  个最大期望点 (peak)  
 ③ 对每一类  $C$ , 设其有  $n_c$  个最大期望点 (peak)  
 ④ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑤ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑥ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑦ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑧ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑨ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑩ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑪ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑫ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑬ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑭ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑮ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑯ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑰ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑱ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑲ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ⑳ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉑ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉒ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉓ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉔ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉕ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉖ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉗ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉘ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉙ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉚ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉛ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉜ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉝ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉞ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㉟ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊱ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊲ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊳ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊴ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊵ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊶ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊷ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊸ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊹ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊺ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊻ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊼ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊽ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊾ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值  
 ㊿ 在某一类  $C$  的峰值  $\rightarrow$  其邻域内, 设其有  $n_c$  个峰值

问题是: 如果两个目标的关联点重合, 那么只能识别其中一个.

基础网络: Hourglass, ResNet, DLA.

网络结构:



SCRDet (一种适用于小, 杂, 旋转目标的3类别检测器)

针对小目标: 通过特征融合和特征采样, 设计了一种特征融合结构。

针对密集布景问题: 设计了一种结合有监督像素级网络和多尺度网络的结构, 抑制背景, 突出目标特征。

针对旋转问题: 通过添加旋转因子设计了一种改进的 Smooth L1 损失, 专门用于解决旋转边界回归中的边界问题。

网络框架:

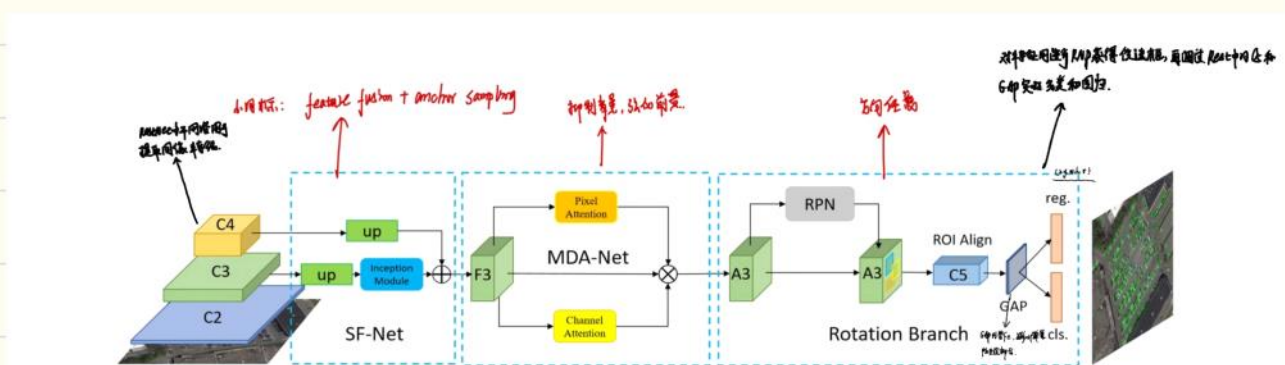


Figure 1: SCRDet includes SF-Net, MDA-Net against small and cluttered objects and rotation branch for rotated objects.

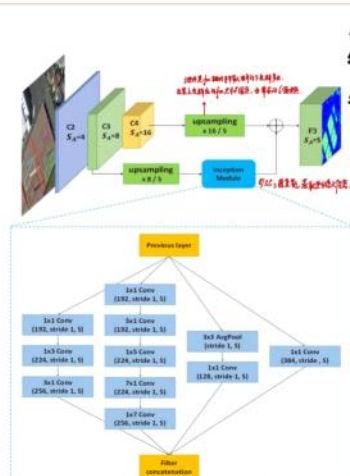


Figure 3: SF-Net. F3 has a small  $S_A$ , while fully considering the feature fusion and adaptability to different scales.

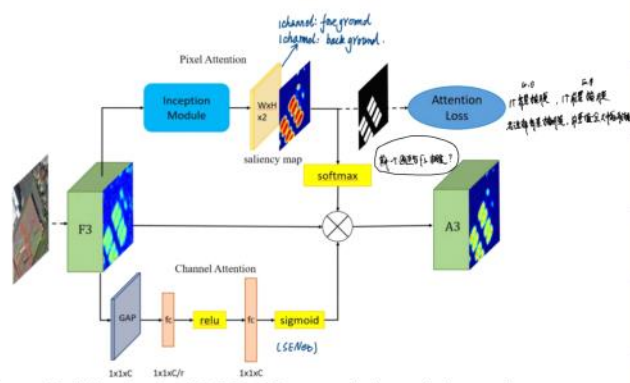


Figure 5: The devised MDA-Net consisting of channel attention network and pixel attention network.

SCRDet是一个 two-stage 目标检测器。

SCRDet由 SF-Net, MDA-Net 以及 Rotation Branch 三部分组成, SF-Net 以及 MDA-Net 通过不断的对目标进行提取特征, 最后利用 Rotation Branch 进行分类回归。

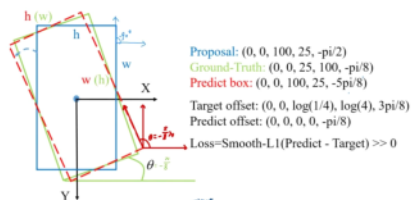


Figure 6: Boundary discontinuity of the rotation angle.

respectively (likewise for  $y, w, n, v$ ).

### 3.4. Loss Function

The multi-task loss is used which is defined as follows:

$$L = \frac{\lambda_1}{N} \sum_{n=1}^N t_n' \sum_{j \in \{x, y, w, h, \theta\}} \frac{L_{reg}(v_{nj}', v_{nj})}{|L_{reg}(v_{nj}', v_{nj})|} \left| -\log(IoU) \right| + \frac{\lambda_2}{h \times w} \sum_i^h \sum_j^w L_{att}(u_{ij}', u_{ij}) + \frac{\lambda_3}{N} \sum_{n=1}^N L_{cls}(p_n, t_n) \quad (3)$$

where  $N$  indicates the number of proposals.  $t_n$  represents

## AugFPN

问题: FPN存在的问题:

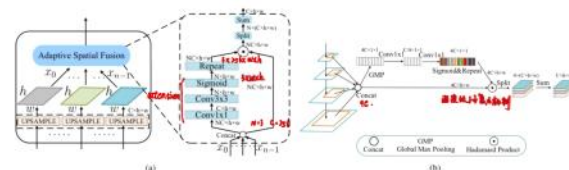
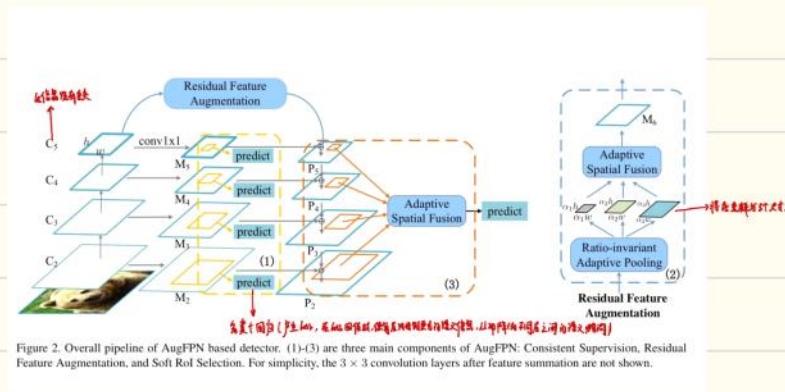
- ① **before fusion**: 在 fusion 之前, 需对用  $1 \times 1$  conv 对图片进行降维, 使 channel 相同, 才能相加。然而不同 stage 的 feature 不同, 包含的语义信息不同, 把两个语义信息相差较大的 feature 直接相加, 将削弱自底向上特征的表达能力。
- ② **top-down feature fusion**: 高层的 feature 经过  $1 \times 1$  conv 后才与底层 feature 相加, 导致信息丢失。
- ③ **after feature fusion**: 特征融合后直接对特征进行预测, 没有预测层也包含冗余的信息, 在检测时会对检测结果有影响。

思路: Consistent supervision (一致性监督): 为了特征融合前, 不同尺度特征之间在语义上差距。

Residual Feature Augmentation (剩余特征增强): 在特征融合中, 通过残差特征增强提取不变的下层信息, 以减少特征映射在最高层特征金字塔中的信息丢失。

Soft RoI selection: 使用软 ROI 选择, 在特征融合后自适应的筛选更好的 ROI 特征。

结构:



$$\begin{aligned}
L_{rcnn} = & \lambda \sum_{M=2}^5 (L_{cls,M}(p_M, t^*) + \beta[t^* > 0]L_{loc,M}(d_M, b^*)) \\
& + \sum_{P=2}^5 (L_{cls,P}(p, t^*) + \beta[t^* > 0]L_{loc,P}(d, b^*)).
\end{aligned}
\tag{1}$$