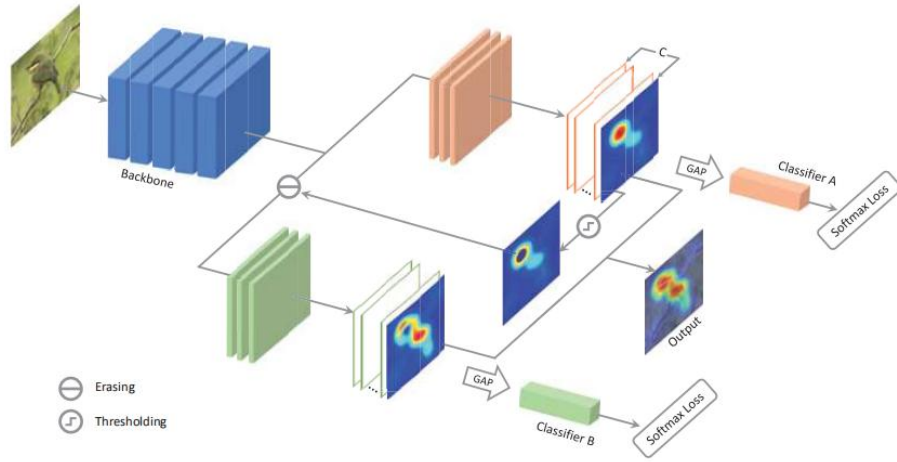


论文学习 1: Adversarial Complementary Learning for weakly supervised object localization

这篇能应用到当前的框架比较多

主要思想：对于一张给定输入，由于 CAM 会出现背景激活过度或者前景激活不足的情况，构造两个独立网络（记为 A、B）对其进行特征提取，当 A 提取完特征后，在 B 提取特征时把 A 的特征进行擦除，以强迫 B 网络对未激活的特征部分进行学习，最后将两个网络的特征图输出融合在一起。

缺点：需要训练多个独立网络来获得对象区域，势必会花费更多的训练时间和计算资源；在没有别的监督信息指导下，网络并不一定总能发现新的对象区域。



基本架构

简要说明：Backbone 是一个全卷积网络，用于特征提取，然后接两个分类器，其输入特征不同，B 的特征输入会在原输入特征上擦除 A 的定位映射特征。具体是，先对 A 的特征定位映射加以阈值判断，将大于阈值（认为特征明显）的部分设为 0，以对抗的方式擦除 B 的输入部分，从而激励 B 学习到其他的 A 未关注的特征，最后将两个分支结合起来获得完整定位区域。

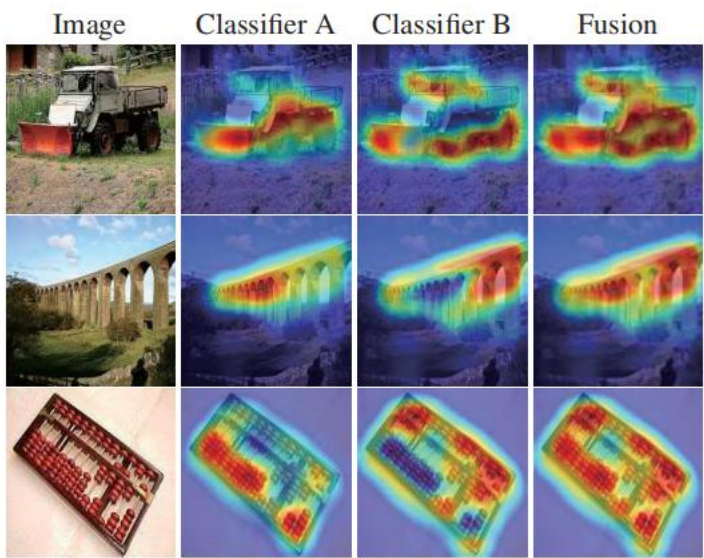
Algorithm 1 Training algorithm for ACoL

Input: Training data $I = \{(I_i, y_i)\}_{i=1}^N$, threshold δ

- 1: **while** training is not convergent **do**
- 2: Update feature maps $S \leftarrow f(\theta_0, I_i)$
- 3: Extract localization map $M^A \leftarrow f(\theta_A, S, y_i)$
- 4: Discover the discriminative region $R = \bar{M}^A > \delta$
- 5: Obtain erased feature maps $\tilde{S} \leftarrow \text{erase}(S, R)$
- 6: Extract localization map $M^B \leftarrow f(\theta_B, \tilde{S}, y_i)$
- 7: Obtain fused map $\bar{M}_{i,j}^{fuse} = \max(\bar{M}_{i,j}^A, \bar{M}_{i,j}^B)$
- 8: Update θ_0, θ_A and θ_B
- 9: **end while**

Output: \bar{M}^{fuse}

结果展示：

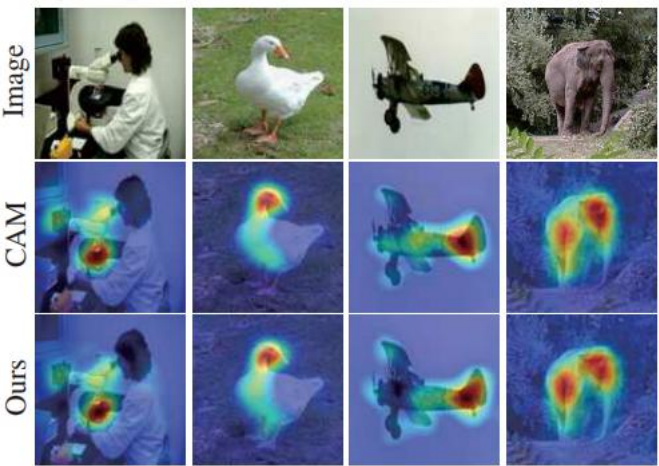


一些新的方法：

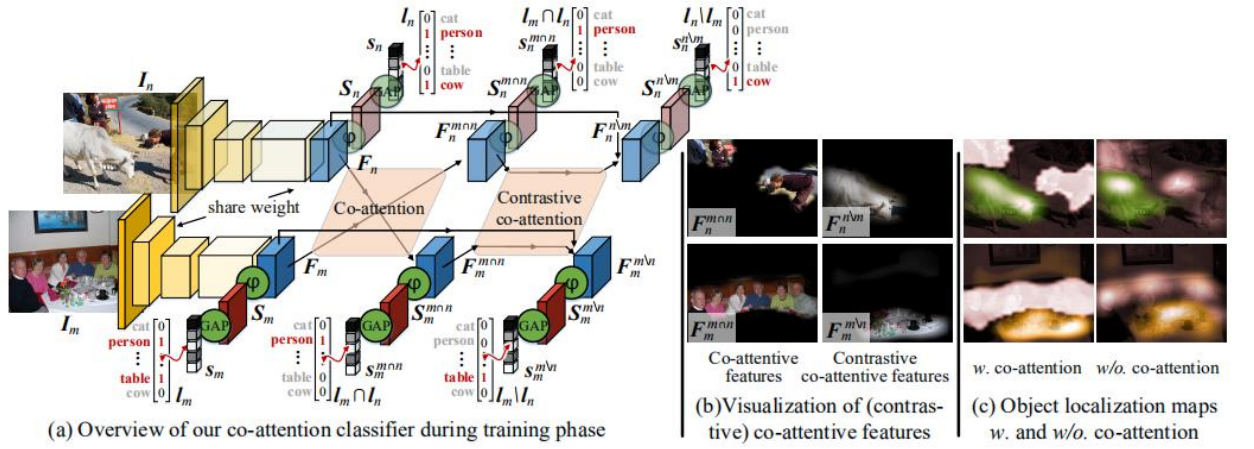
CAM 新生成思路：之前的 CAM 是通过某层特征图和网络输出之间求梯度等一些运算后得到，且不能直接做端到端的训练（图片输入麻烦）。这篇文章从数学角度证明了可以通过一步步骤得到与 CAM 效果极为相似的图。

Given the output feature maps S of an FCN, we add a convolutional layer of C channels with the kernel size of 1×1 , stride 1 on top of the feature maps S . Then, the output is fed into a GAP layer followed by a softmax layer for classification

通过这种修正方法，对象定位的映射就可以直接在前向传递中获得，而不需要 CAM 的后处理步骤。



作者处理结果对比图



基本框架

思想概述：第一部分依然是特征提取，输入位两张图片 m, n ，输出特征图记为 F_m 和 F_n ，随后引入两个注意力机制，其中一个是同类别的注意力机制，另一个是异类别的注意力机制，在这篇文章里，输入的两张图拥有同样的类别（比如 **person**）和不同的类别（比如 m 中有 **cow**，而 n 中有 **table**），两个注意力机制是用于寻找相同的部分和不同的部分。

首先是同类别注意力部分：

their correlations. We first compute the affinity matrix P between F_m and F_n :

$$P = F_m^T W_P F_n \in \mathbb{R}^{HW \times HW}, \quad (2)$$

where $F_m \in \mathbb{R}^{C \times HW}$ and $F_n \in \mathbb{R}^{C \times HW}$ are flattened into matrix formats, and $W_P \in \mathbb{R}^{C \times C}$ is a learnable matrix. The affinity matrix P stores similarity scores corresponding to all pairs of positions in F_m and F_n , i.e., the $(i, j)^{th}$ element of P gives the similarity between i^{th} location in F_m and j^{th} location in F_n .

Then P is normalized column-wise to derive attention maps across F_m for each position in F_n , and row-wise to derive attention maps across F_n for each position in F_m :

$$A_m = \text{softmax}(P) \in [0, 1]^{HW \times HW}, \quad A_n = \text{softmax}(P^T) \in [0, 1]^{HW \times HW}, \quad (3)$$

$$F_m^{m \cap n} = F_n A_n \in \mathbb{R}^{C \times H \times W}, \quad F_n^{m \cap n} = F_m A_m \in \mathbb{R}^{C \times H \times W}, \quad (4)$$

原文如上，大概是首先将 F_m 和 F_n 展平，通过一个可学习的矩阵 W_P 与 F_m 和 F_n 作乘操作得到相似性矩阵 P ， P 中的第 (i, j) 元素表示了 F_m 和 F_n 在第 (i, j) 位置的相似度大小，然后在一些归一化等操作后，将得到的相似度矩阵分别乘以原来的特征图，这样做就把相似的部分强化了而把不相似的部分弱化了。



可视化后的样子，只有人的部分显示了。

异类别的操作原理也差不多。

想法：这些是基于一张图中有多个实物标签做的（像是有 **person**、**table**、**cow**），因此它的分割难度可能更大。对于肿瘤 MRI 图，一般只具有背景、肿瘤和头部（骨骼那种框架），而在之前观察健康样本的 CAM 时，也确实骨骼的权值较高，疾病样本出问题的也有在骨骼位置权值高的。有没有必要引入多标签描述一份样本（头颅+肿瘤，头颅+无肿瘤），去弱化掉同样标签的头颅部分？这么做和直接判别两张图是否属于同一类会不会有区别？还有没有其他可以延伸的地方呢？