

Appendix of BagCAMs

Lei Zhu¹, Qian Chen², Lujia Jin³, Yunfei You³, and Yanye Lu^{3*}

¹ Institute of Medical Technology, Peking University

² Department of Biomedical Engineering, Peking University

³ Institute of Biomedical Engineering, Peking University Shenzhen Graduate School
<https://github.com/zh460045050/BagCAMs>
 zhulei@stu.pku.edu.cn, yanye.lu@pku.edu.cn

1 Key symbols List

Here we give the pivotal symbols list in Table. 1 to enhance the readability of our main paper. In general, we use ***italic bold*** uppercase characters to denote the matrices. Vectors are denoted with lowercase. Sets are noted by blackboard bold, for example \mathcal{F} .

Table 1. The Meaning of Some Pivotal Symbols

$\mathbf{X} \in \mathbb{R}^{3 \times N}$	RGB feature of the input image.
$\mathbf{Z} \in \mathbb{R}^{C \times N}$	Pixel-level features of \mathbf{X}
$\hat{\mathbf{P}} \in \mathbb{R}^{K \times N}$	Initial Pixel-level localization scores of BagCAMs.
$\mathbf{P}^* \in \mathbb{R}^{K \times N}$	Output Pixel-level localization scores of BagCAMs.
$\mathbf{\Lambda} \in \mathbb{R}^{N \times N}$	Co-efficient matrix of regional localizers
$\mathbf{W} \in \mathbb{R}^{K \times C}$	Weight matrix of the classifier.
$\mathbf{z} \in \mathbb{R}^{C \times 1}$	Image-level features of \mathbf{X} .
$\mathbf{p} \in \mathbb{R}^{K \times 1}$	Localization score of a spatial position (column vector of $\hat{\mathbf{P}}$).
$\mathbf{s} \in \mathbb{R}^{K \times 1}$	Image-level classification scores.
$\bar{\mathbf{s}} \in \mathbb{R}^{K \times 1}$	Logarithmic Image-level classification scores.
$\mathbf{y} \in \mathbb{R}^{K \times 1}$	Image-level ground truth mask.
N^I	The spatial resolution of input \mathbf{X} .
N	The spatial resolution of feature \mathbf{Z} .
C	The number of channel for feature.
K	The number of object class.
\mathcal{F}^*	The base localizer set of BagCAMs that contains N^2 localizers.
\mathcal{P}^*	The set of localization scores generate by localizers in \mathcal{F}^* .
$e(\cdot) : \mathbb{R}^{3 \times N} \rightarrow \mathbb{R}^{C \times N}$	The feature extractor, implemented by ResNet/Inception.
$c(\cdot) : \mathbb{R}^{C \times 1} \rightarrow \mathbb{R}^{K \times 1}$	The image classifier, implemented by fully-connected layer.
$f(\cdot) : \mathbb{R}^{C \times N} \rightarrow \mathbb{R}^{K \times N}$	The object localizer, derived by CAM-based methods.

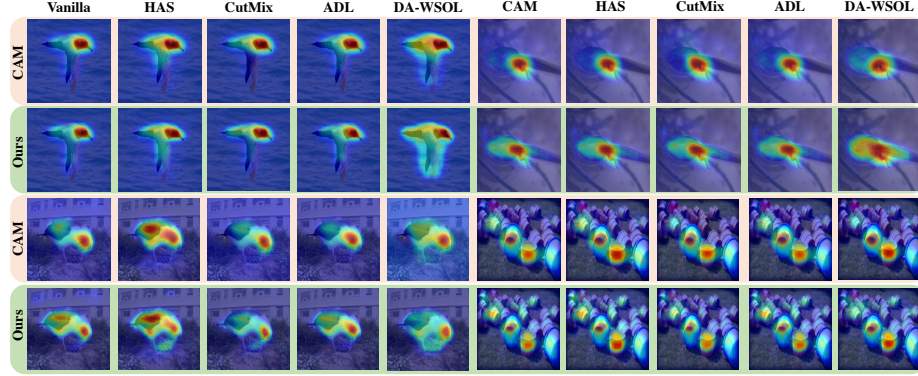


Fig. 1. Visualization on replacing CAM into BagCAMs for different WSOL methods.

2 Proof of the formulation of base localizer

Here we detailed the proof of the formulation of the base localizer of our BagCAMs defined in our Sec.3.2:

$$f_n^m(\mathbf{x})_k = \sum_{c_1} s_k \left(1 + \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_1, m}} \mathbf{Z}_{c_1, m} \right) \sum_{c_2} \left(\frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_2, n}} \mathbf{x}_{c_2} \right) . \quad (1)$$

Proof. When adopting PCS [4] to initialize the coarse localization score $\hat{P}_{k, m}$ of our BagCAMs, the formulation of f_n^m obtained by RLG becomes:

$$f_n^m(\mathbf{x})_k = \sum_{c_2} \frac{\partial (\sum_{c_1} \frac{\partial s_k}{\partial \mathbf{Z}_{c_1, m}} \mathbf{Z}_{c_1, m})}{\partial \mathbf{Z}_{c_2, n}} \mathbf{x}_{c_2} = \sum_{c_1, c_2} \left(\frac{\partial s_k}{\partial \mathbf{Z}_{c_2, n}} + \frac{\partial^2 s_k}{\partial \mathbf{Z}_{c_1, m} \partial \mathbf{Z}_{c_2, n}} \mathbf{Z}_{c_1, m} \right) \mathbf{x}_{c_2} . \quad (2)$$

Considering that ReLU is used as the activation function of the extractor, based on the prior work [1], the partial derivative $\frac{\partial^2 s_k}{\partial \mathbf{Z}_{c_1, m} \partial \mathbf{Z}_{c_2, n}}$ can be reformulated as:

$$\frac{\partial s_k}{\partial \mathbf{Z}_{c_2, n}} = s_k \frac{\bar{s}_k}{\partial \mathbf{Z}_{c_2, n}} , \quad \frac{\partial^2 s_k}{\partial \mathbf{Z}_{c_1, m} \partial \mathbf{Z}_{c_2, n}} = s_k \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_1, m}} \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_2, n}} . \quad (3)$$

Finally, taking Eq. 3 into Eq. 2, the formulation of Eq. 1 can be obtained:

$$\begin{aligned} f_n^m(\mathbf{x})_k &= \sum_{c_1, c_2} \left(s_k \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_2, n}} + s_k \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_1, m}} \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_2, n}} \mathbf{Z}_{c_1, m} \right) \mathbf{x}_{c_2} \\ &= \sum_{c_1} s_k \left(1 + \frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_1, m}} \mathbf{Z}_{c_1, m} \right) \sum_{c_2} \left(\frac{\partial \bar{s}_k}{\partial \mathbf{Z}_{c_2, n}} \mathbf{x}_{c_2} \right) . \end{aligned} \quad (4)$$

□

3 Visualizations of InceptionV3 extractor

In our main paper, we only included the visualization of our method with the ResNet50 extractor. Here we give more visualization when using InceptionV3 [3] as the backbone extractor. In detail, Fig. 1 shows that replacing CAM with BagCAMs for different WSOL methods can also enhance the performance of the baseline methods under the InceptionV3. Moreover, our BagCAMs still achieve better performance when localizing objects based on features of intermediate layers than other methods [1,2,4] as indicated in Fig. 2. These are in accordance with the results with the ResNet50 backbone as discussed in our main paper.

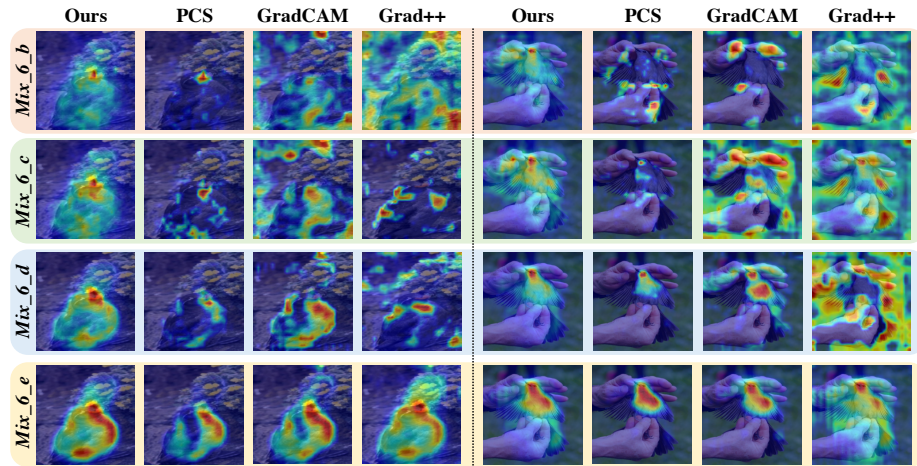


Fig. 2. Localization map generated based on different InceptionV3 layer.

References

1. Chattopadhyay, A., Sarkar, A., Howlader, P., Balasubramanian, V.N.: Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 839–847. IEEE (2018)
2. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: Proceedings of the IEEE/CVF International Conference on Computer Bision (ICCV). pp. 618–626 (2017)
3. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2818–2826 (2016)
4. Tan, C., Gu, G., Ruan, T., Wei, S., Zhao, Y.: Dual-gradients localization framework for weakly supervised object localization. In: Proceedings of the 28th ACM International Conference on Multimedia. pp. 1976–1984 (2020)