

---

# Multi-view Feature Augmentation with Adaptive Class Activation Mapping

---

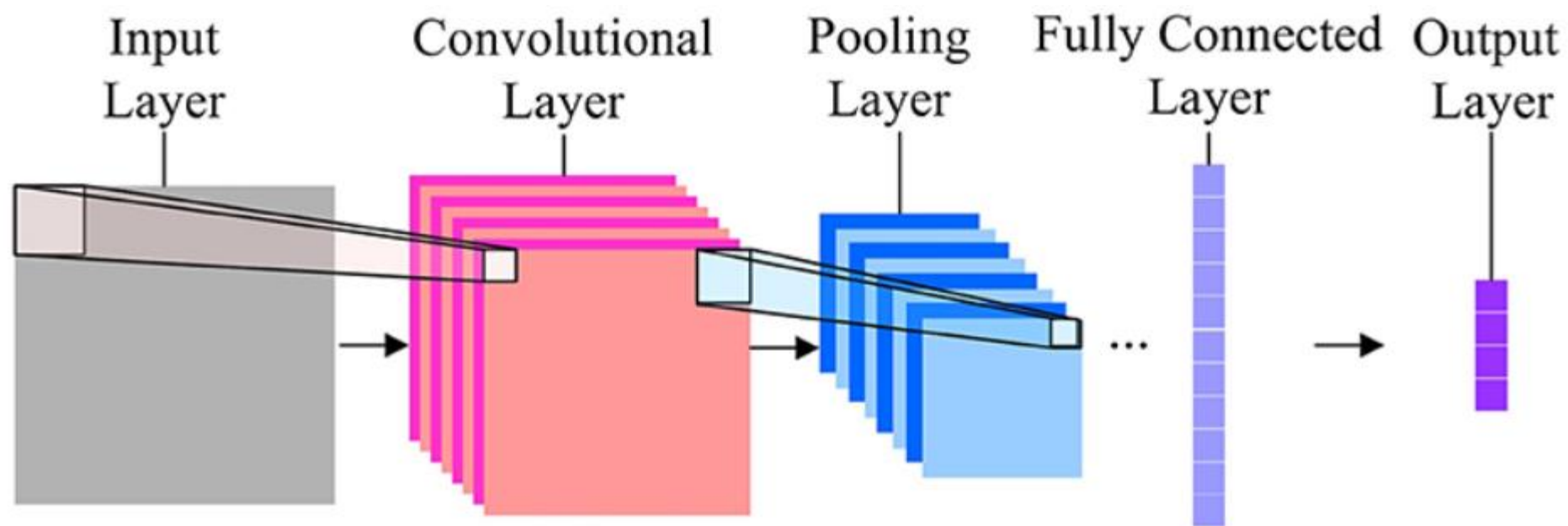
Xiang Gao

Yingjie Tian

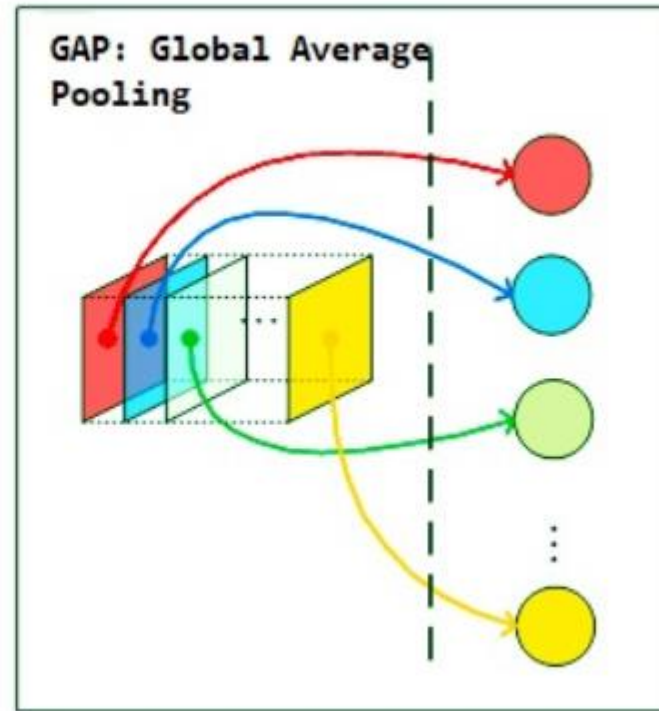
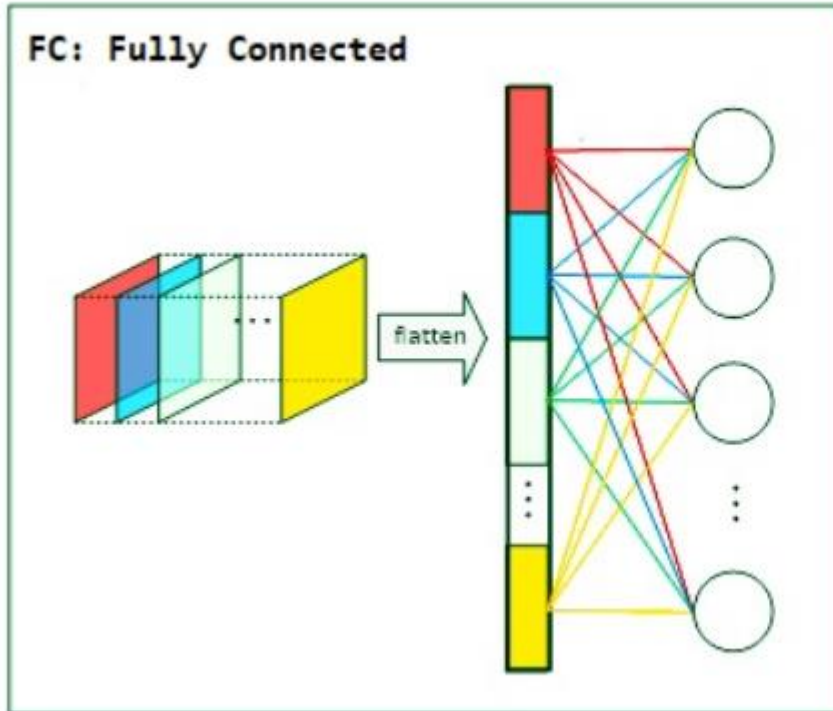
Zhiquan Qi

IJCAI 2021

# CNN



CNN  $\longrightarrow$  GAP(Global Average Pooling)



# Motivation

tench



English  
springer



golf ball



# Framework

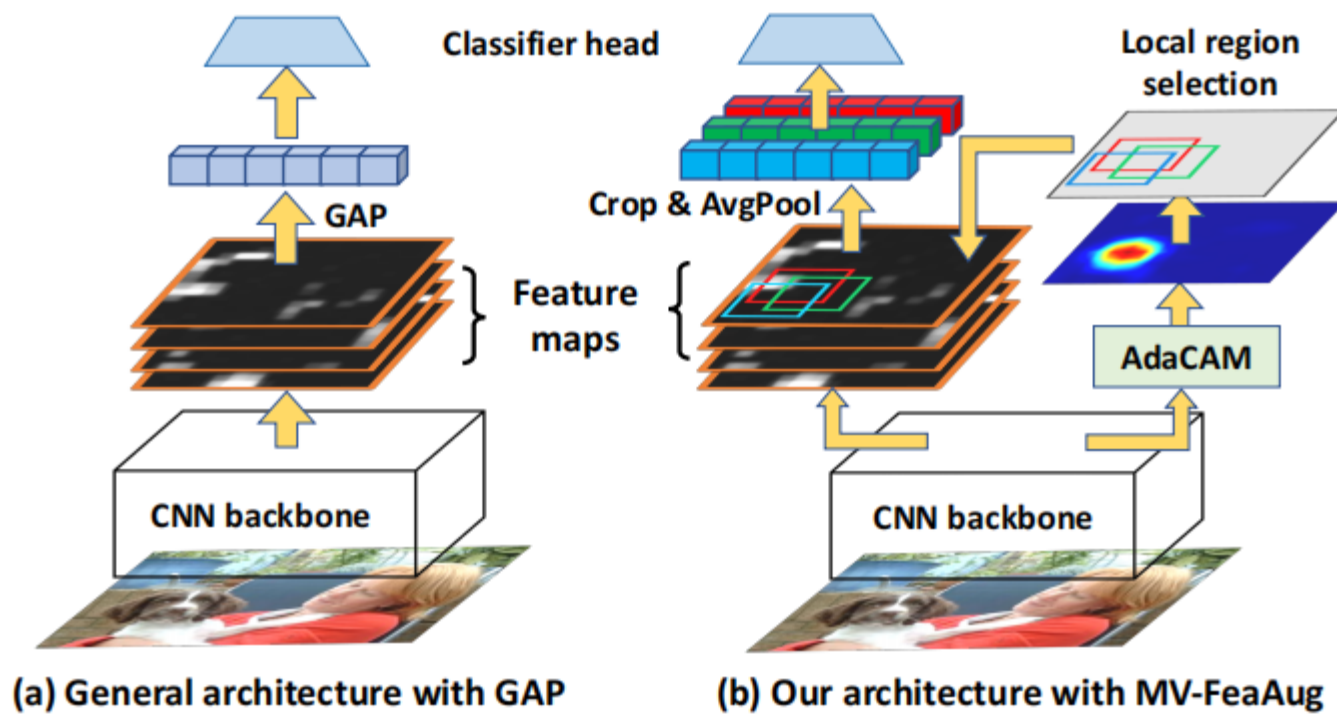


Figure 2: Comparison between the image classification architecture with the general GAP (left) and our MV-FeaAug module (right). We sample diverse local features around the class-discriminative region of the final convolutional feature maps as multi-view local image representations for ensemble classification, as compared with GAP that extracts only global-view image representation.

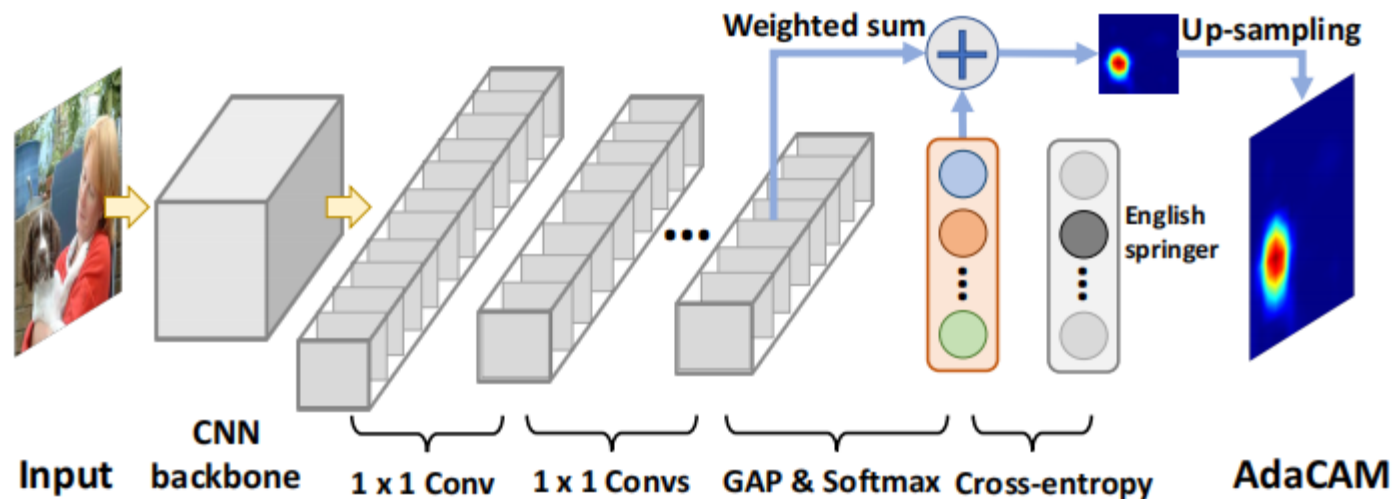
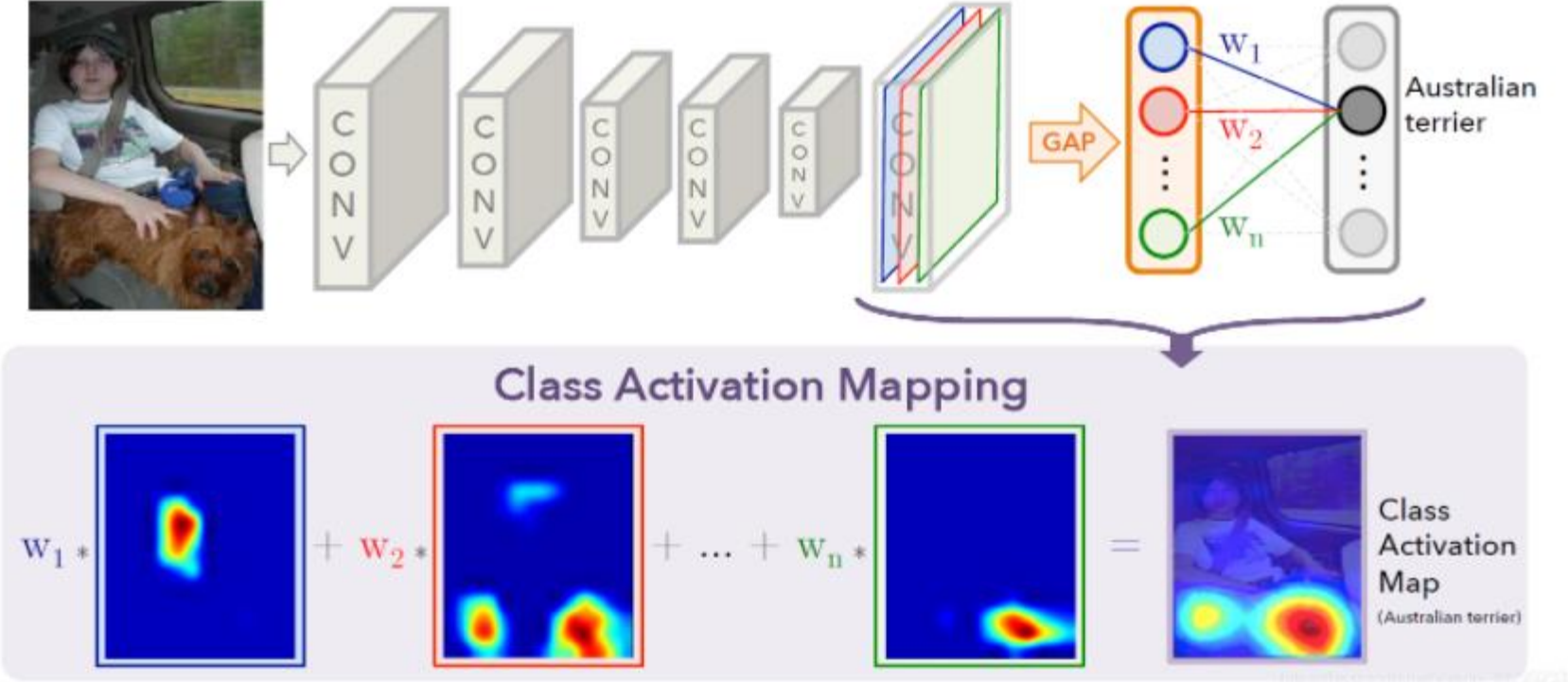


Figure 3: Adaptive class activation mapping (AdaCAM). We replace the traditional classifier head made up of [GAP→MLP→Softmax] with [MLPConv→GAP→Softmax] (MLPConv comprises consecutive  $\text{Conv}_{1 \times 1}$  layers joined by non-linear activations) to maintain spatial resolution of feature maps. The AdaCAM is obtained by performing channel-wise weighted sum of the last convolutional feature maps with respect to the softmax logit vector.

# CAM



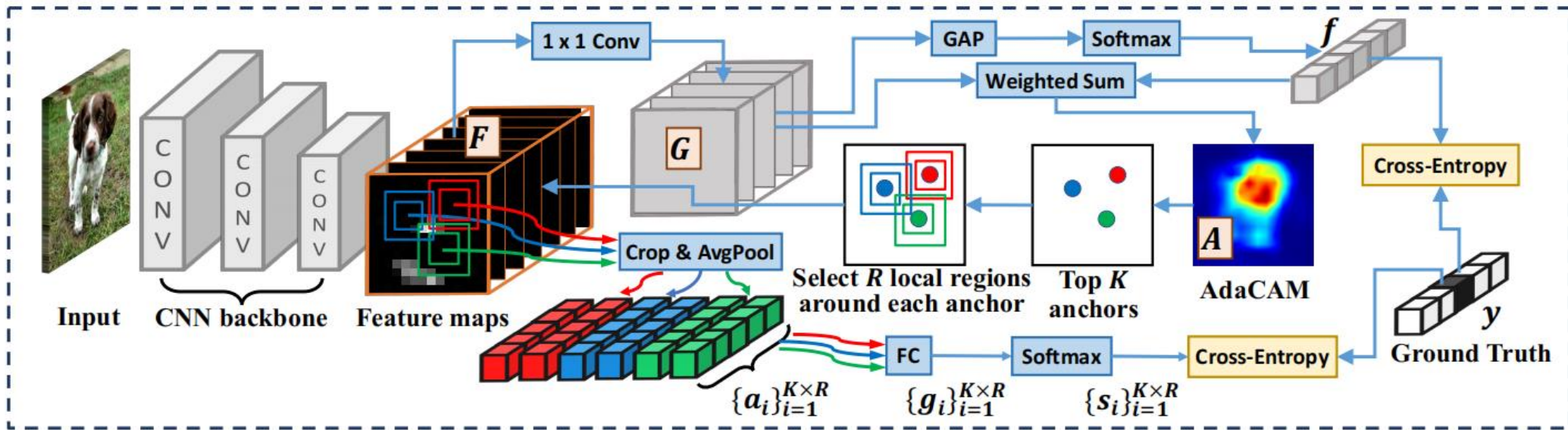


Figure 4: Overview of MV-FeaAug. We concurrently train an auxiliary classifier head comprised of only one  $1 \times 1$  convolutional layer for dynamic generation of AdaCAM, based on which we sample multiple local representations on the final convolutional feature maps as multi-view inputs to the main classifier head. The main classifier head comprises (but not restricted to) a single fully-connected layer.

$$\begin{aligned}
 L_{total} = & -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log(f_c^{(n)}) \\
 & -\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^{K \times R} \sum_{c=1}^C y_c^{(n)} \log(s_{i,c}^{(n)}), \quad (9)
 \end{aligned}$$

# Experiments

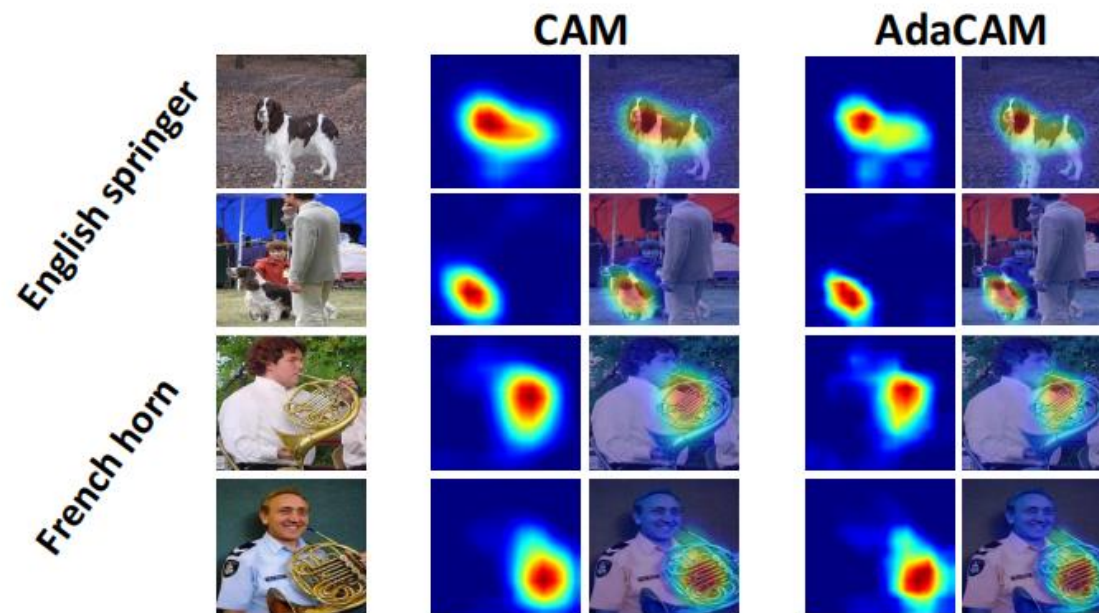
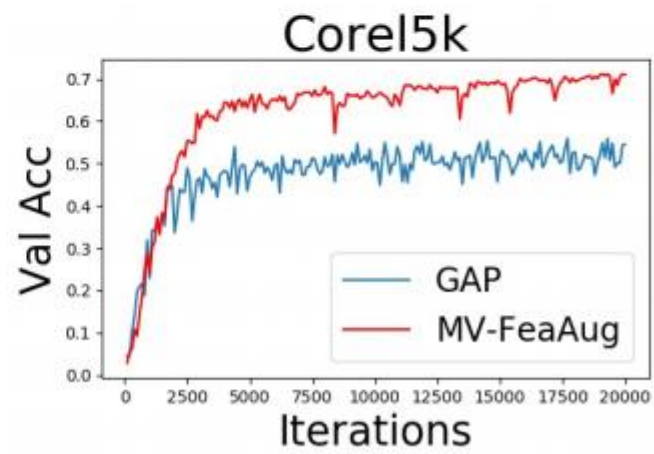
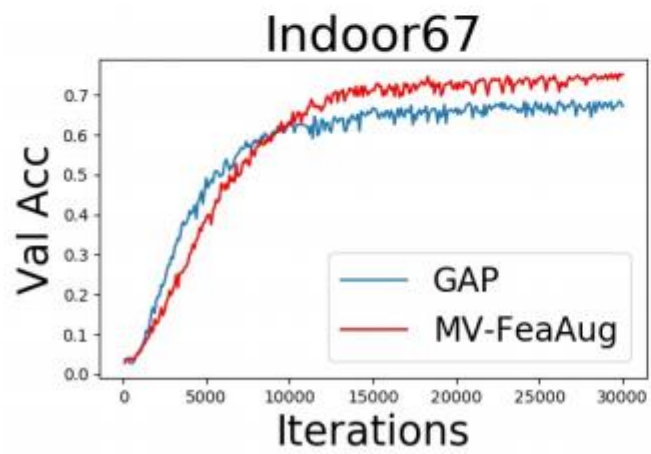
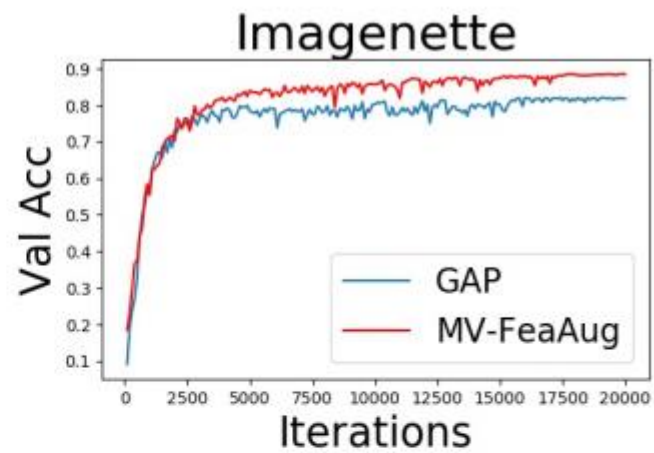
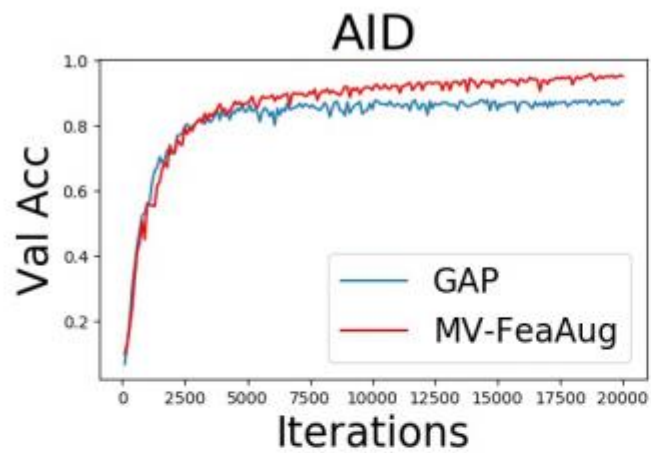


Figure 5: Visual comparison between CAM and our AdaCAM (evaluated on Imagenette validation set) in object localization. Refer to supplementary materials for more results.



Model	Datasets									
	Imagenette	Caltech101	Corel5k	Scene15	Indoor67	Action40	Event8	UCMLU	RSSCN7	AID
VGG16_GAP	82.04%	70.42%	54.31%	80.86%	68.25%	52.19%	80.00%	90.95%	91.25%	87.61%
VGG16_MFA	<b>88.45%</b>	<b>76.05%</b>	<b>70.74%</b>	<b>85.53%</b>	<b>75.27%</b>	<b>60.36%</b>	<b>85.75%</b>	<b>96.67%</b>	<b>95.29%</b>	<b>94.55%</b>
ResNet50_GAP	83.66%	71.84%	54.55%	81.41%	70.14%	53.43%	81.08%	91.67%	91.43%	89.48%
ResNet41_GAP	82.47%	70.75%	52.72%	80.17%	68.90%	52.56%	79.88%	90.58%	90.54%	88.31%
ResNet41_MFA	<b>88.60%</b>	<b>76.19%</b>	<b>69.84%</b>	<b>85.00%</b>	<b>75.59%</b>	<b>60.61%</b>	<b>85.25%</b>	<b>96.43%</b>	<b>94.71%</b>	<b>94.89%</b>
ResNeXt50_GAP	84.10%	72.65%	56.20%	83.18%	71.26%	54.71%	82.95%	92.65%	92.32%	90.24%
ResNeXt41_GAP	82.79%	71.54%	55.66%	81.84%	69.95%	53.50%	81.89%	91.79%	91.18%	88.96%
ResNeXt41_MFA	<b>88.86%</b>	<b>76.78%</b>	<b>71.48%</b>	<b>86.29%</b>	<b>76.38%</b>	<b>61.77%</b>	<b>87.14%</b>	<b>97.31%</b>	<b>95.14%</b>	<b>95.29%</b>
MobileNet_GAP	81.26%	70.26%	53.86%	80.24%	67.79%	52.48%	80.50%	90.24%	90.54%	86.87%
MobileNet24_GAP	80.60%	69.78%	51.77%	79.34%	66.96%	51.73%	79.25%	89.37%	89.61%	86.04%
MobileNet24_MFA	<b>87.55%</b>	<b>75.74%</b>	<b>69.31%</b>	<b>84.97%</b>	<b>74.53%</b>	<b>59.95%</b>	<b>84.75%</b>	<b>95.48%</b>	<b>93.82%</b>	<b>93.52%</b>

Table 1: Comparison between GAP and our MV-FeaAug (MFA) module with  $K = 50$  in validation accuracy based on different CNN backbones. Our MV-FeaAug impressively boosts model performance with simply one more  $Conv_{1 \times 1}$  layer than GAP-based counterpart.

Model	Datasets			Backbone Parameters
	Caltech101	AID	Indoor67	
VGG16_GAP	70.42%	87.61%	68.25%	14.03M
VGG16_SE_GAP	70.96%	88.44%	68.75%	14.18M
VGG16_CBAM_GAP	71.53%	89.37%	69.24%	14.61M
VGG16_SK_GAP	71.25%	89.63%	69.08%	19.55M
VGG16_GC_GAP	71.87%	90.04%	69.51%	14.76M
VGG16_mixup_GAP	71.55%	89.50%	69.67%	14.03M
VGG16_cutout_GAP	70.90%	88.25%	69.14%	14.03M
VGG16_cutmix_GAP	72.10%	90.16%	70.03%	14.03M
VGG16_MFA (ours)	<b>76.05%</b>	<b>94.55%</b>	<b>75.27%</b>	14.03M

Table 2: Comparison of our MV-FeaAug (MFA) with related visual attention modules and data augmentation methods in validation accuracy. We use VGG16 backbone for all the related methods, and set  $K = 50$ ,  $L_R = [3, 5, 7, 9]$  in our MV-FeaAug module.

Datasets	GAP	MV-FeaAug				
		K=20	K=30	K=40	K=50	K=60
Imagenette	82.04%	86.27%	87.62%	88.13%	88.45%	88.72%
Indoor67	68.25%	73.51%	74.34%	74.89%	75.27%	75.55%
Action40	53.21%	58.24%	59.36%	59.98%	60.36%	60.68%
UCMLU	90.95%	94.52%	95.33%	96.02%	96.67%	97.05%

Table 3: Study of the impact of  $K$  to validation accuracy on different datasets. The backbone network is VGG16.

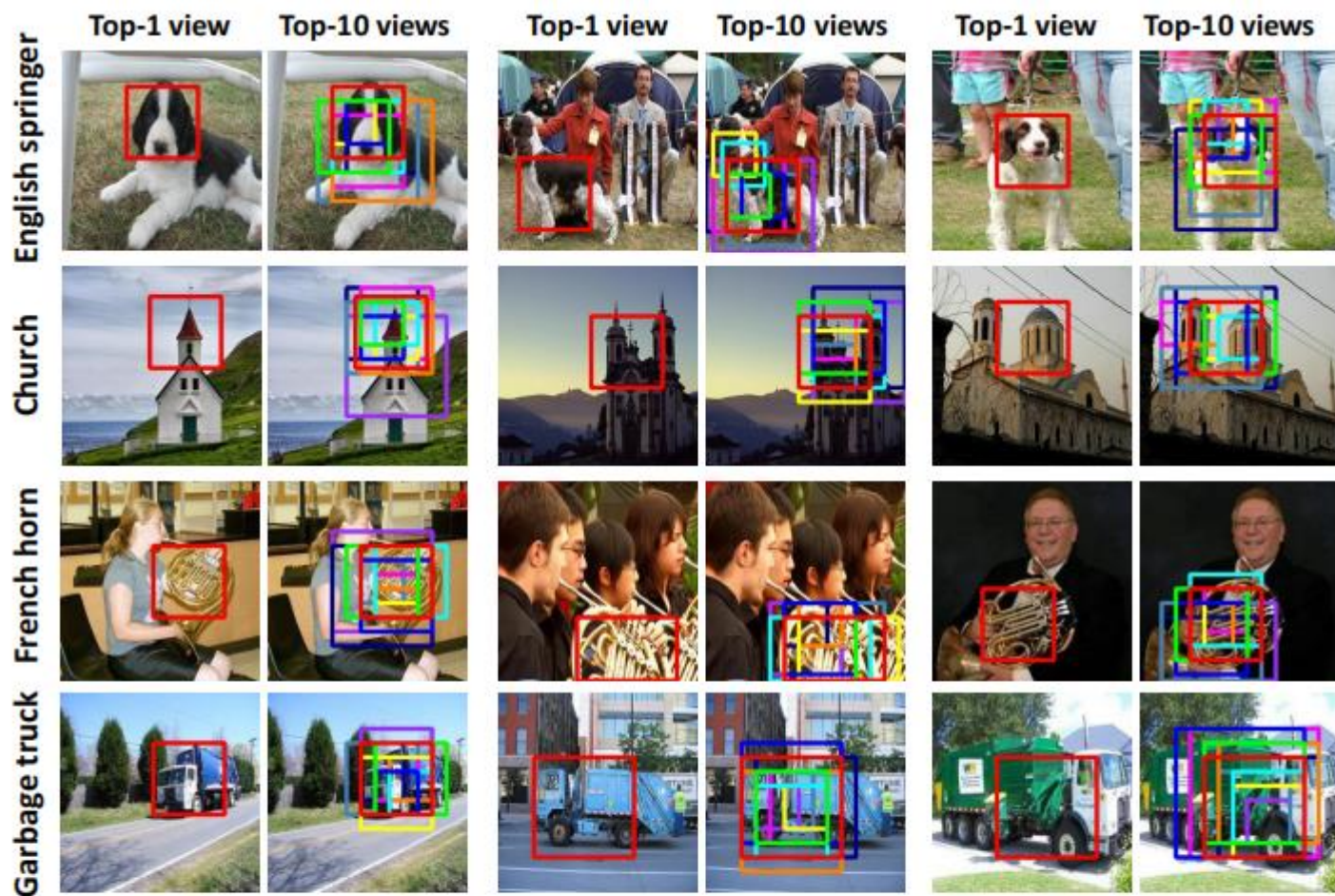


Figure 7: Visualization of the top-1 and top-10 local regions that correspond to local features with the highest prediction confidence of the corresponding image class. Please zoom in for better resolution and refer to supplementary materials for more results.