



Group-wise Inhibition based Feature Regularization for Robust Classification

Haozhe Liu[†], Haoqian Wu[†], Weicheng Xie^{*}, Feng Liu^{*}, Linlin Shen

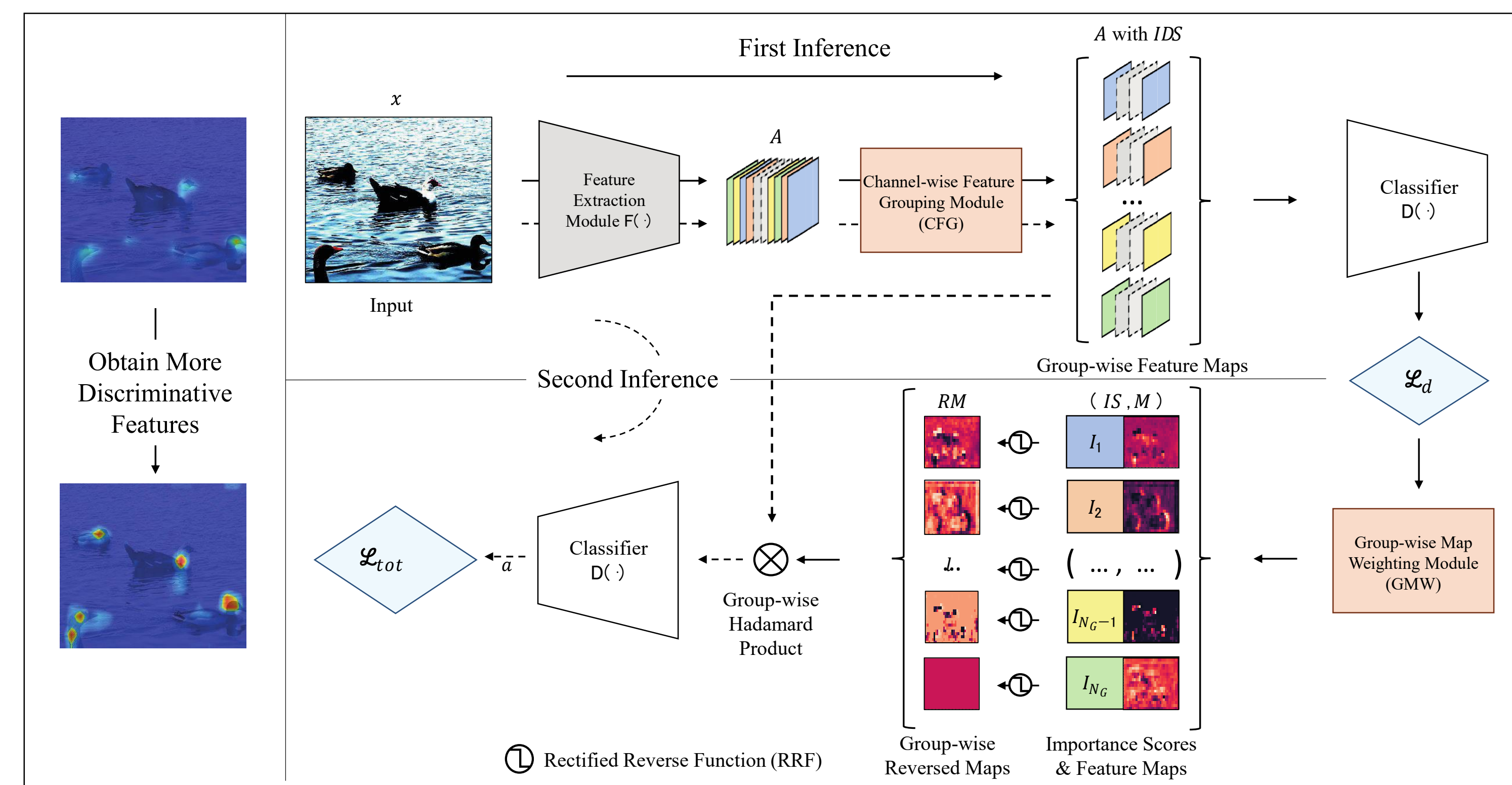
Institute of Computer Vision, College of Computer Science and Software Engineering
Shenzhen University, Shenzhen 518060, China



Abstract

The convolutional neural network (CNN) is vulnerable to degraded images with even very small variations (e.g. corrupted and adversarial samples). One of the possible reasons is that CNN pays more attention to the most discriminative regions, but ignores the auxiliary features when learning, leading to the lack of feature diversity for final judgment. In our method, we propose to dynamically suppress significant activation values of CNN by group-wise inhibition, but not fixedly or randomly handle them when training. The feature maps with different activation distribution are then processed separately to take the feature independence into account. CNN is finally guided to learn richer discriminative features hierarchically for robust classification according to the proposed regularization. Our method is comprehensively evaluated under multiple settings, including classification against corruptions, adversarial attacks and low data regime. Extensive experimental results show that the proposed method can achieve significant improvements in terms of both robustness and generalization performances, when compared with the state-of-the-art methods. *Code is available at https://github.com/LinusWu/TENET_Training.*

Method



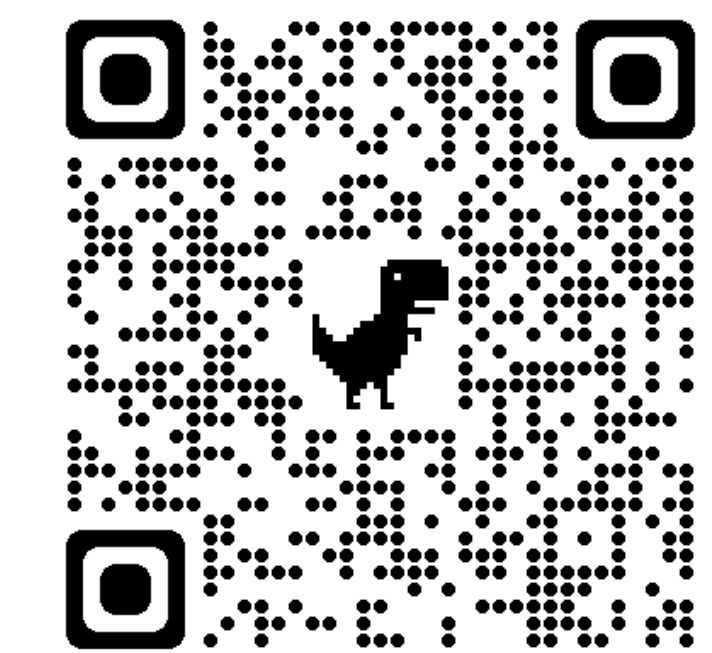
PIPELINE:

CNNs consist of the feature extraction module $F(\cdot)$ and the classifier $D(\cdot)$. In the first inference, feature maps A encoded with $F(\cdot)$ are divided into multiple groups by the CFG module, and loss \mathcal{L}_d is calculated based on $D(\cdot)$. Reversed maps RM are then derived using GMW module and RRF. In the second inference, the Hadamard Product of A (with IDS) and RM is fed to $D(\cdot)$ to calculate the loss \mathcal{L}_{total} .

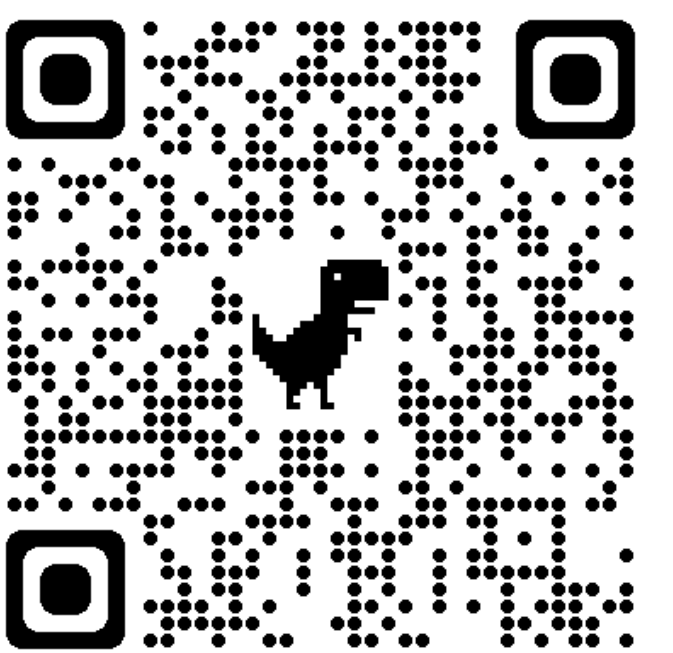
Conclusion

- No limitation for network architecture
- Flexible to complement with other solutions
- Improving both robustness and generalization

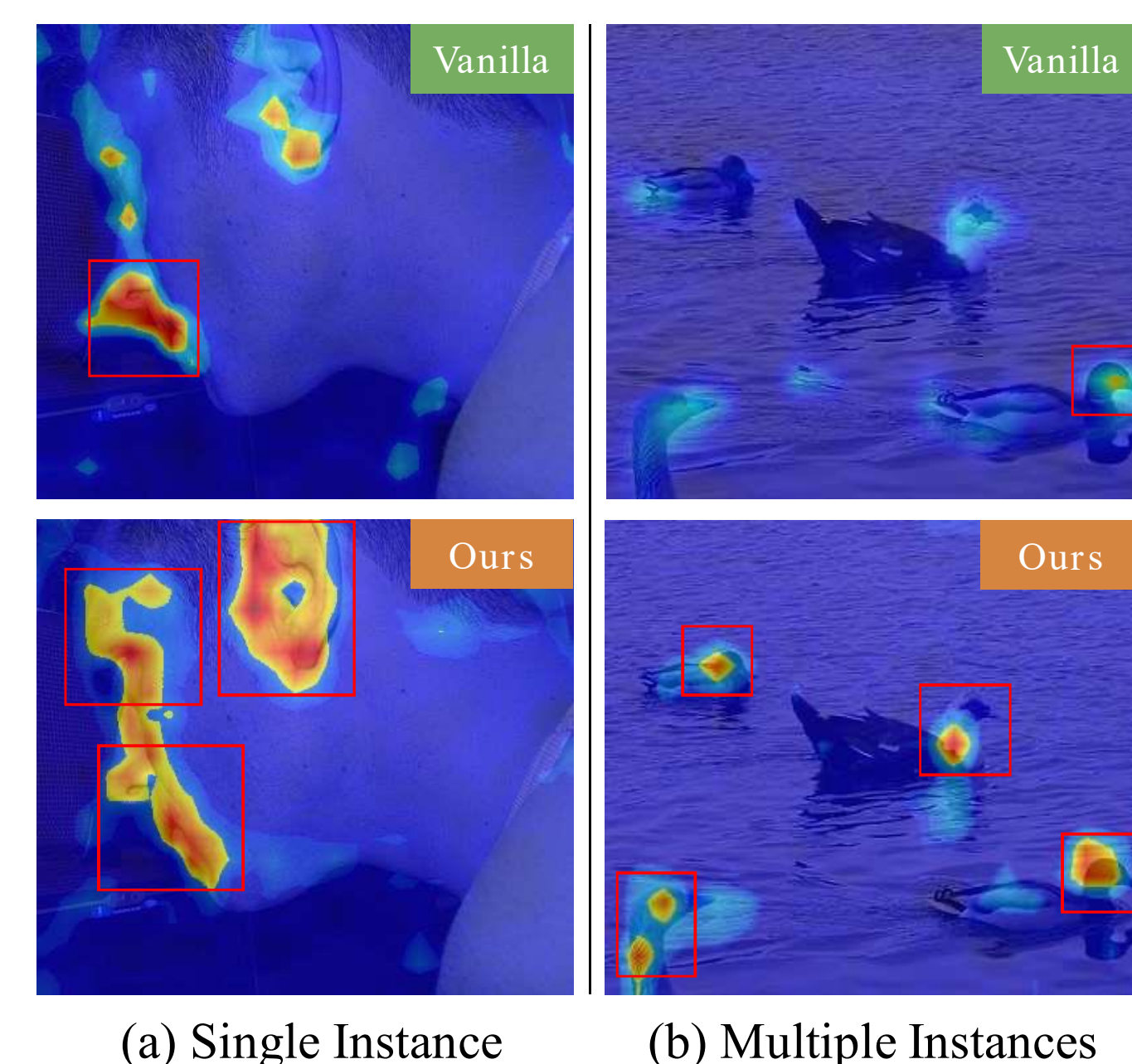
Paper(Arxiv):



Code(Github):



Motivation



GENERAL IDEA

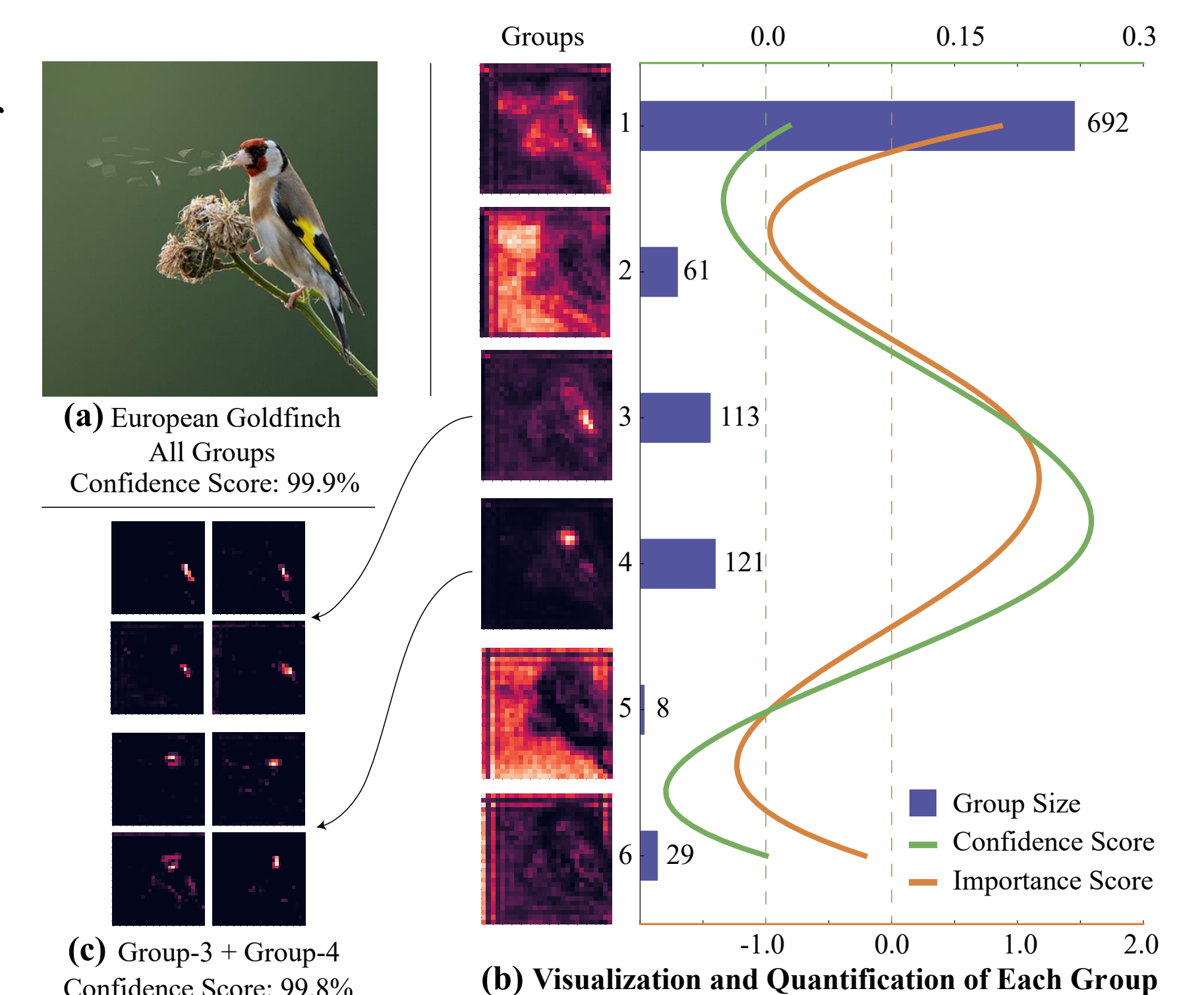
CNNs can locate the most discriminative regions for both single-instance and multi-instance samples with the regularization method, while neglecting other auxiliary features that are critical for the recognition. The lack of auxiliary features may lead to insufficient feature diversity, which consequently results in a feature space with low-dimension for classification and limits the robustness.

Experimental Results

Task-[protocol]	Dataset	Previous SOTA	Gain
Standard Classification -[4]	PASCAL VOC 2012[9]	Group Orthogonal Training [4]	2.9%
Robustness against Adversarial Attack-[8, 24]	CIFAR-10/100 [18]	A. T. [24]	5.75%
		Augmix[14]	15.56%*
Robustness against Common Corruption-[13, 14, 21]	CIFAR-10/100-C [13]	Augmix[14]	1.77%
	ImageNet-C [13]		2.8%†
Generalization -[2]	CUB-200 [26]	GLICO [2]	2.75%

To evaluate the performance of the proposed method, extensive experiments are carried on publicly-available data sets, including PASCAL VOC 2012, CIFAR-10/100, ImageNet-C and CUB-200.

The visualization and quantification of the feature maps extracted by the 3rd residual block of ResNet-50 using TENET Training. (a) An input image with the label of European Goldfinch. (b) The activation distribution, the corresponding importance and confidence scores of each group clustered by CFG module. (c) The example feature maps selected from the 3rd and 4th groups.



The heatmap visualization of feature maps encoded with ResNet-50, based on Grad-CAM with or without the proposed method. Our method locates more diverse discriminative regions (in red boxes) for both single-instance (a) and multiple-instance (b) samples.