

Occlusion-Robustness of Convolutional Neural Networks via Inverted Cutout - Supplementary Material

Matthias Körschens

Computer Vision Group

Friedrich Schiller University Jena

07737 Jena, Germany

Email: matthias.koerschens@uni-jena.de

Paul Bodesheim

Computer Vision Group

Friedrich Schiller University Jena

07737 Jena, Germany

Email: paul.bodesheim@uni-jena.de

Joachim Denzler

Computer Vision Group

Friedrich Schiller University Jena

07737 Jena, Germany

Email: joachim.denzler@uni-jena.de

S-I. INTRODUCTION

In this supplementary document we provide several example images of both datasets used to provide a better insight into the different occlusion levels and kinds of occlusion used. Moreover, we also present the numerical results of our experiments again with the addition of the respective standard deviations to give a better indication of the significance of our results.

S-II. DATASET EXAMPLES

As mentioned in the main paper, the Occluded-Vehicles dataset, which is based on the Pascal3D+ dataset, has not only four different occlusion levels, but also four different kinds of occlusion, examples of which can be found in Figure S1. The boxes for white box occlusion, texture occlusion and noise occlusion are placed on random locations, however they are consistent over the three mentioned kinds of occlusion. The fourth kind is occlusion by segmented objects whose classes are not contained in the classes to categorize. We can see that, in contrast to the three other occlusion types, the object occlusion comprises objects in strongly differing shapes instead of rectangles. This resembles realistic occlusion scenarios in the wild more closely and gives an indication of the performance of the method in real applications.

In Figure S2 we can see several real occlusion scenarios from the images of the Occluded-COCO-Vehicles dataset. While at level L0 the object is completely visible, with rising occlusion levels the view is being blocked by poles, humans or even bigger objects to the point at which the object in question is only barely visible anymore.

S-III. EXPERIMENTAL RESULTS WITH STANDARD DEVIATION

To provide an indication of the significance of our results, we provide the experimental results from the main paper with inclusion of the standard deviation for each value in Table S-I and Table S-II. On the Occluded-Vehicles dataset, we note that for the experimental setups using A_{GAP} and A_{FC} with BCE and IC in, respectively, 10 and 11 of the 14 different

occlusion setups, the state-of-the-art results of CompNet-Multi lie outside the 2σ range our results. In general, we can also see that, regarding the mean over all occlusion setups, the results of CompNet-Multi consistently lie outside the 2σ range of our top results on both datasets. This suggests a significant improvement of our methods over CompNet-Multi of Kortylewski *et al.*[1].

S-IV. ABLATION STUDY WITH CUTOUT

To show that, while our Inverted Cutout method is similar to Cutout, it nevertheless yields significantly different results, in this part we present an ablation study with the standard Cutout augmentation [5] in the different setups. All hyperparameters used here are the same as in the main paper and the sampled cutouts also have sizes ranging from 16×16 pixels to 128×128 pixels. The results of the experiments on the Occluded-Images dataset are shown in Table S-III and the respective results for the Occluded-COCO-Vehicles dataset are shown in Table S-IV. We can also see here that, the results with Binary Cross Entropy loss are better than the typical classification Categorical Cross Entropy loss. However, in comparison with the results from our method in Table S-I and Table S-II, using the Cutout augmentation yields significantly worse results than Inverted Cutout. This is especially notable in the higher occlusion levels can be up to around 15% on both datasets when comparing the same base setup, showing a significant advantage of our method. As mentioned in the main paper, this difference is likely caused by the correlations of object parts learned when using Cutout. While the latter only removes a comparably small part from the image, it leaves most of the image intact, resulting in the network learning from co-occurring objects parts. This makes the network unable to identify the object correctly when several of the usually co-occurring object parts are missing due to occlusion. In contrast, with Inverted Cutout the network learns to identify the object from only the small unoccluded image patches, making the network more robust against partial occlusion.

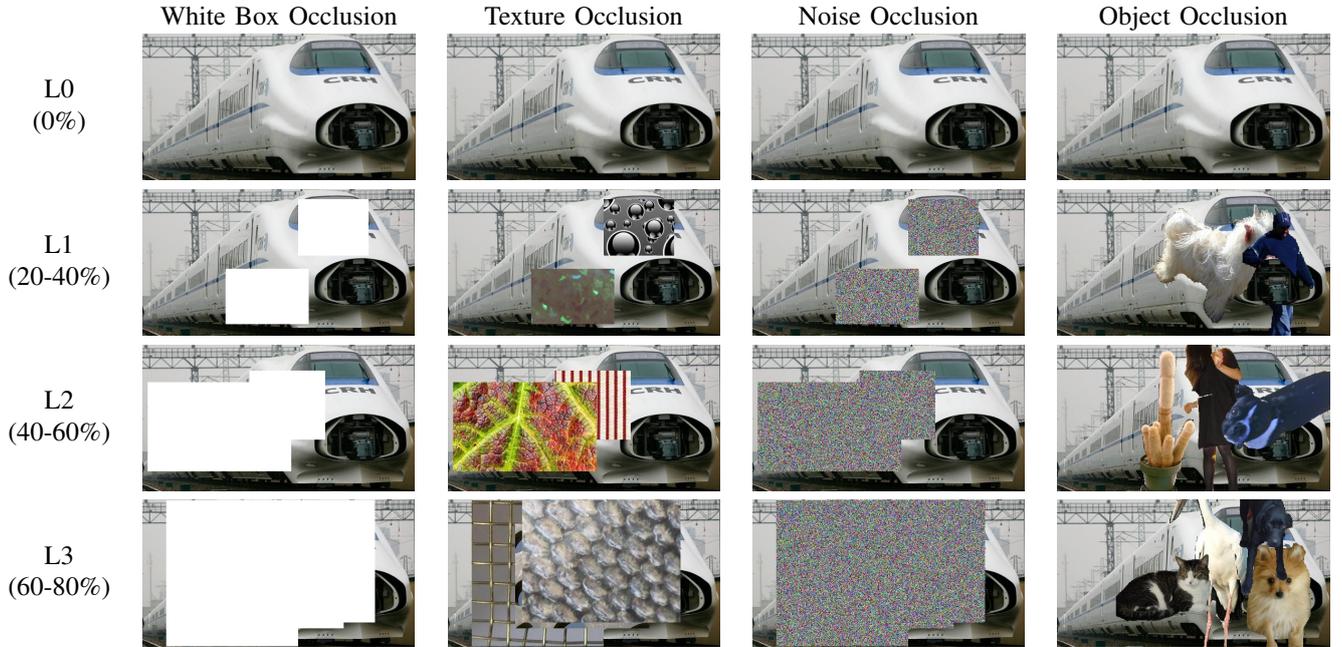


Fig. S1: Example images from the Occluded-Images Pascal3D+ dataset with several different levels and kinds of occlusion.

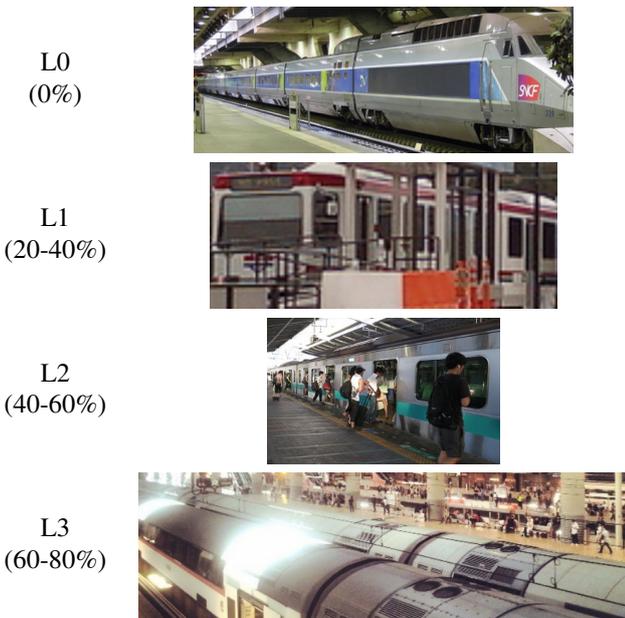


Fig. S2: Example images from the Occluded-COCO-Images dataset with several levels of occlusion.

S-V. INVERTED CUTOUT MASK SIZE ABLATION

To quantify the effect of the mask size when utilizing Inverted Cutout, we also present an ablation study using different mask sizes on both datasets. Similar to the previous ablation, the hyperparameters used are the same as in the main paper. The sizes ablated include 16×16 , 32×32 , 64×64 and 128×128 pixels. The results for the Occluded-Vehicles dataset are presented in Table S-V and the ones for Occluded-

COCO-Vehicles in Table S-VI. Generally, we can see that the larger cutout sizes usually outperform the smaller ones on both datasets. However, we also note that none of the single cutout sizes outperforms the randomly sized cutouts used in the main paper and in Table S-I and Table S-II. While on the Occluded-Vehicles dataset the difference in average occlusion accuracy in comparison to the top results from the main paper for some setups is rather small, the gaps are bigger on the Occluded-COCO-Vehicles dataset. The reason for this is likely that the more regular artificial occlusions in the Occluded-Vehicles are easier to classify despite occlusion in contrast to the much more irregular natural occlusions contained in Occluded-COCO-Vehicles, which, hence, require more sophisticated training to become more robust against more complex occlusions.

In summary, bigger cutouts were found to be more effective, however, the combination of different cutout sizes by randomly sampling the size during the training still outperforms using constant sizes.

S-VI. FEATURE EXTRACTION LAYER ABLATION

While in the main paper we used, similar to the top results from [1], the combined features from the two last convolutional blocks of VGG-16, there is also the option to use only the output of each of these blocks without the respective other. In this section, we will discuss this possibility. Experimental results for this setup are shown in Table S-VII and Table S-VIII for the trials on the Occluded-Vehicles and Occluded-COCO-Vehicles dataset, respectively, and the hyperparameters are the same as before.

In both tables, one can see that, when using the layers separately, the network also performs comparably well against

TABLE S-I: The accuracies (in %) and their respective standard deviations (denoted by \pm) of our method in comparison with previously introduced methods on the Pascal3D+ Occluded-Vehicles dataset. The values of methods used for comparison have been taken from [1]. * marks the results received after fine-tuning the whole network. The occlusion types are: w - white box occlusion, n - noise box occlusion, t - texture occlusion, o - occlusion by segmented objects. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy loss, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and the one utilizing a convolutional layer, respectively.

Occ. Area	L0: 0%	L1: 20-40%				L2: 40-60%				L3: 60-80%				Mean
Occ. Type	-	w	n	t	o	w	n	t	o	w	n	t	o	
VGG [2]	99.2	96.9	97.0	96.5	93.8	92.0	90.3	89.9	79.6	67.9	62.1	59.5	62.2	83.6
CoD [3]	92.1	92.7	92.3	91.7	92.3	87.4	89.5	88.7	90.6	70.2	80.3	76.9	87.1	87.1
VGG+CoD [3]	98.3	96.8	95.9	96.2	94.4	91.2	91.8	91.3	91.4	71.6	80.7	77.3	87.2	89.5
TDAPNet [4]	99.3	98.4	98.6	98.5	97.4	96.1	97.5	96.6	91.6	82.1	88.1	82.7	79.8	92.8
CompNet-Multi [1]	99.3	98.6	98.6	98.8	97.9	98.4	98.4	97.8	94.6	91.7	90.7	86.7	88.4	95.4
$A_{GAP}+CCE$	99.5 ± 0.0	97.5 ± 0.5	97.4 ± 0.2	97.1 ± 0.4	92.9 ± 1.2	90.3 ± 3.0	88.2 ± 1.7	88.3 ± 1.9	68.7 ± 4.1	58.2 ± 7.4	47.2 ± 4.3	46.6 ± 3.9	47.1 ± 4.5	78.4 ± 2.0
$A_{GAP}+CCE^*$	99.5 ± 0.1	97.6 ± 0.5	97.5 ± 0.2	97.3 ± 0.3	93.1 ± 1.0	90.5 ± 2.8	88.9 ± 1.5	88.9 ± 1.7	69.2 ± 4.0	59.3 ± 7.6	48.5 ± 4.3	47.3 ± 3.9	47.4 ± 4.4	78.9 ± 2.0
$A_{FC}+CCE$	99.7 ± 0.0	97.9 ± 0.3	98.0 ± 0.4	98.0 ± 0.4	95.2 ± 0.4	93.2 ± 1.2	92.7 ± 1.1	91.4 ± 0.9	80.4 ± 1.7	61.8 ± 2.7	56.2 ± 2.8	51.9 ± 1.8	60.8 ± 2.8	82.9 ± 1.1
$A_{FC}+CCE^*$	99.7 ± 0.1	98.1 ± 0.3	98.2 ± 0.3	98.1 ± 0.4	95.3 ± 0.5	93.3 ± 1.0	93.1 ± 1.2	91.6 ± 0.8	80.4 ± 2.0	62.0 ± 2.4	56.8 ± 2.9	52.4 ± 1.4	61.1 ± 2.9	83.1 ± 1.0
$A_{GAP}+BCE$	98.5 ± 0.3	93.7 ± 0.9	92.7 ± 1.7	92.6 ± 1.9	84.9 ± 4.6	83.2 ± 2.2	80.4 ± 3.2	79.0 ± 4.4	59.9 ± 7.6	49.3 ± 5.0	43.4 ± 6.2	41.2 ± 4.5	41.8 ± 7.1	72.4 ± 3.3
$A_{GAP}+BCE^*$	99.2 ± 0.3	96.8 ± 0.9	97.1 ± 0.8	97.0 ± 0.9	94.5 ± 0.6	90.0 ± 1.6	91.0 ± 2.4	90.6 ± 1.9	83.2 ± 3.5	62.9 ± 3.0	64.6 ± 6.9	61.4 ± 6.6	68.4 ± 5.6	84.3 ± 2.3
$A_{FC}+BCE$	99.8 ± 0.1	98.7 ± 0.1	98.5 ± 0.1	98.7 ± 0.2	97.0 ± 0.3	95.4 ± 0.3	94.7 ± 0.5	94.1 ± 0.5	83.6 ± 2.3	66.4 ± 0.7	60.9 ± 1.9	59.5 ± 1.3	61.4 ± 2.9	85.3 ± 0.6
$A_{FC}+BCE^*$	99.7 ± 0.1	98.7 ± 0.1	98.6 ± 0.1	98.7 ± 0.2	97.0 ± 0.4	95.3 ± 0.2	94.8 ± 0.4	94.1 ± 0.5	83.8 ± 2.3	66.3 ± 0.6	61.1 ± 2.0	59.3 ± 1.6	61.8 ± 3.2	85.3 ± 0.6
$A_{GAP}+CCE+IC$	99.5 ± 0.1	99.2 ± 0.1	99.1 ± 0.1	98.9 ± 0.3	96.8 ± 0.5	99.0 ± 0.1	97.7 ± 0.2	96.3 ± 0.4	90.8 ± 2.5	95.2 ± 0.3	86.6 ± 1.8	74.4 ± 2.7	78.0 ± 5.4	93.2 ± 0.9
$A_{GAP}+CCE+IC^*$	99.7 ± 0.0	99.4 ± 0.1	99.3 ± 0.1	99.0 ± 0.0	97.8 ± 0.2	99.2 ± 0.2	98.3 ± 0.2	96.9 ± 0.4	92.8 ± 1.4	96.4 ± 0.3	87.0 ± 1.9	75.4 ± 2.3	82.3 ± 3.4	94.1 ± 0.6
$A_{FC}+CCE+IC$	99.5 ± 0.0	99.5 ± 0.0	99.2 ± 0.1	99.3 ± 0.1	97.9 ± 0.1	99.2 ± 0.1	98.2 ± 0.1	97.5 ± 0.2	93.7 ± 0.5	96.8 ± 0.4	87.8 ± 1.9	78.9 ± 1.6	85.1 ± 0.7	94.8 ± 0.3
$A_{FC}+CCE+IC^*$	99.6 ± 0.1	99.6 ± 0.1	99.2 ± 0.2	99.4 ± 0.1	98.6 ± 0.1	99.3 ± 0.2	97.9 ± 0.5	97.5 ± 0.2	95.7 ± 0.6	95.9 ± 0.8	83.0 ± 4.8	80.4 ± 1.3	88.9 ± 1.1	95.0 ± 0.6
$A_{GAP}+BCE+IC$	99.1 ± 0.1	98.9 ± 0.1	98.3 ± 0.1	98.3 ± 0.1	97.4 ± 0.2	98.0 ± 0.3	96.4 ± 0.3	96.7 ± 0.1	93.1 ± 0.6	92.8 ± 0.6	83.8 ± 1.7	80.5 ± 1.4	84.8 ± 1.0	93.7 ± 0.3
$A_{GAP}+BCE+IC^*$	99.7 ± 0.1	99.6 ± 0.1	99.4 ± 0.1	99.4 ± 0.1	98.9 ± 0.1	99.2 ± 0.1	98.4 ± 0.2	98.1 ± 0.3	96.9 ± 0.2	95.8 ± 0.3	87.4 ± 2.4	85.4 ± 1.1	92.1 ± 1.0	96.2 ± 0.2
$A_{FC}+BCE+IC$	99.6 ± 0.1	99.6 ± 0.1	99.3 ± 0.1	99.3 ± 0.1	98.7 ± 0.0	99.4 ± 0.1	98.3 ± 0.1	98.1 ± 0.2	95.4 ± 0.6	97.0 ± 0.1	90.5 ± 1.3	82.8 ± 1.2	87.4 ± 1.0	95.8 ± 0.3
$A_{FC}+BCE+IC^*$	99.7 ± 0.1	99.7 ± 0.1	99.4 ± 0.1	99.4 ± 0.1	98.9 ± 0.1	99.6 ± 0.1	98.8 ± 0.2	98.5 ± 0.2	96.2 ± 0.3	97.9 ± 0.1	90.8 ± 1.6	83.9 ± 1.0	89.4 ± 1.0	96.3 ± 0.2

occlusion. We can also see that using the P5 layer outperforms the P4 layer in most instances, which is expected due to the more sophisticated features generated by the later network layers. While the overall performance is similar to the combined one, we can, however, see some differences. Most notably, the accuracy of the networks with separate blocks is often much worse when training only the last layer instead of fine-tuning the complete network. Hence, the combination of both produces much more valuable features that can be used when training the complete network is not desired. When fine-tuning the complete network, however, the performance using the separate layers comes closer to the combined approach, with the P5 layer even surpassing the latter in one instance. Nevertheless, the combined approach overall outperforms the two separate layers, which was also found in the investigations by Kortylewski *et al.* in [1].

S-VII. RESNET50-BACKBONE ABLATION

While in the main paper we utilized only the VGG-16 network architecture for the comparability to previous works, in this section we will also take a look into the application of Inverted Cutout in conjunction with another popular backbone architecture: ResNet50 [6]. During the training process with the ResNets, we utilized slightly different hyperparameters, which are shown in Table S-IX. Hyperparameters not mentioned are the same as in the main paper.

The experimental results for this ablation study are shown in Table S-X and Table S-XI for the Occluded-Vehicles and Occluded-COCO-Vehicles datasets, respectively.

Generally, it is visible that for ResNet50, Inverted Cutout is also strongly beneficial, especially when utilizing the A_{FC} aggregation module in comparison to the setups without IC. Moreover, we notice that, when training the model without IC,

TABLE S-II: Accuracies (in %) and their respective standard deviations (denoted by \pm) of our experiments on the Occluded-COCO-Vehicles dataset. We compare our approaches with previously introduced methods based on classification accuracy. The values of methods used for comparison have been taken from [1]. * marks the results received after fine-tuning the network. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy losses, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and with a convolutional layer, respectively.

Train Data	MS-COCO				
Occ. Area	L0	L1	L2	L3	Avg
VGG [2]	99.1	88.7	78.8	63.0	82.4
VGG [2] + Cutout	99.3	90.9	87.5	75.3	88.3
TDAPNet [4]	99.4	88.8	87.9	69.9	86.5
CompNet-Multi [1]	99.4	95.3	90.9	86.3	93.0
$A_{GAP}+CCE$	99.2 ± 0.2	85.3 ± 2.0	80.0 ± 1.6	68.8 ± 4.0	83.3 ± 1.8
$A_{GAP}+CCE^*$	99.3 ± 0.1	85.3 ± 1.6	80.2 ± 1.2	68.8 ± 4.0	83.4 ± 1.7
$A_{FC}+CCE$	99.6 ± 0.0	89.9 ± 0.5	83.7 ± 0.7	70.2 ± 2.0	85.9 ± 0.4
$A_{FC}+CCE^*$	99.6 ± 0.0	90.2 ± 0.5	84.2 ± 0.9	72.3 ± 2.6	86.6 ± 0.7
$A_{GAP}+BCE$	99.3 ± 0.0	86.3 ± 1.9	78.3 ± 8.5	69.5 ± 11.5	83.4 ± 5.4
$A_{GAP}+BCE^*$	99.3 ± 0.0	86.4 ± 1.6	78.3 ± 8.8	69.5 ± 11.6	83.4 ± 5.4
$A_{FC}+BCE$	99.7 ± 0.1	90.6 ± 0.5	84.9 ± 0.7	74.0 ± 2.2	87.3 ± 0.7
$A_{FC}+BCE^*$	99.7 ± 0.1	90.7 ± 0.5	85.5 ± 0.9	74.3 ± 1.8	87.5 ± 0.6
$A_{GAP}+CCE+IC$	99.1 ± 0.1	90.6 ± 0.5	87.2 ± 1.1	87.0 ± 1.5	91.0 ± 0.6
$A_{GAP}+CCE+IC^*$	99.5 ± 0.1	94.5 ± 0.2	90.8 ± 1.3	89.4 ± 3.0	93.5 ± 1.0
$A_{FC}+CCE+IC$	99.4 ± 0.1	94.2 ± 0.8	90.0 ± 0.8	82.9 ± 3.0	91.6 ± 1.0
$A_{FC}+CCE+IC^*$	99.4 ± 0.1	94.7 ± 0.6	90.8 ± 1.0	84.2 ± 0.7	92.3 ± 0.4
$A_{GAP}+BCE+IC$	99.1 ± 0.1	89.4 ± 0.7	88.0 ± 0.5	86.0 ± 2.4	90.6 ± 0.8
$A_{GAP}+BCE+IC^*$	99.4 ± 0.0	94.6 ± 0.8	91.6 ± 0.5	92.8 ± 1.8	94.6 ± 0.4
$A_{FC}+BCE+IC$	99.4 ± 0.1	94.1 ± 0.6	92.0 ± 0.4	86.3 ± 1.7	93.0 ± 0.3
$A_{FC}+BCE+IC^*$	99.5 ± 0.0	95.4 ± 0.4	92.1 ± 0.4	88.7 ± 2.4	93.9 ± 0.7

in general training only the aggregation layer often yields better results than fine-tuning the complete network. This effect can likely be attributed to the batch normalization layers [7] contained in the ResNet50 architecture. Batch normalization appears to make the network unable to deal with occlusion if it has not seen any occlusion examples during training. This strong detrimental effect, however, is mitigated when using IC, where the results surpass the ones after training the aggregation layer alone.

When comparing the results on both datasets, we also notice that the network trained without IC performs better on the more natural occlusions seen in the Occluded-COCO-Vehicles dataset in contrast to the artificial occlusions from the Pascal3D+-based Occluded-Vehicles dataset. The reason

for this is likely that in the former bigger connected object parts are visible in contrast to more randomly spread small object parts like in the images from the Occluded-Vehicles dataset. Such bigger object parts might make the class easier to identify under occlusion for the ResNet.

Nevertheless, on both datasets, the usage of Inverted Cutout improves the results in comparison to the ones generated after training without IC. Therefore, our novel augmentation method is also beneficial for network architecture beyond VGG-16, yielding results comparable with the latter.

REFERENCES

- [1] A. Kortylewski, J. He, Q. Liu, and A. L. Yuille, "Compositional convolutional neural networks: A deep architecture with innate robustness to partial occlusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8940–8949.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [3] A. Kortylewski, Q. Liu, H. Wang, Z. Zhang, and A. Yuille, "Combining compositional models and deep networks for robust object classification under occlusion," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2020, pp. 1333–1341.
- [4] M. Xiao, A. Kortylewski, R. Wu, S. Qiao, W. Shen, and A. Yuille, "Tdapnet: Prototype network with recurrent top-down attention for robust object classification under partial occlusion," *arXiv preprint arXiv:1909.03879*, 2019.
- [5] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *arXiv preprint arXiv:1708.04552*, 2017.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [7] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.

TABLE S-III: **Cutout Ablations** - The accuracies (in %) and their respective standard deviations (denoted by \pm) of the Cutout ablations on the Pascal3D+ Occluded-Vehicles dataset. * marks the results received after fine-tuning the whole network. The occlusion types are: w - white box occlusion, n - noise box occlusion, t - texture occlusion, o - occlusion by segmented objects. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy loss, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and the one utilizing a convolutional layer, respectively.

Occ. Area	L0: 0%	L1: 20-40%				L2: 40-60%				L3: 60-80%				Mean
Occ. Type	-	w	n	t	o	w	n	t	o	w	n	t	o	
$A_{GAP}+CCE+Cutout$	99.7	99.2	99.1	98.8	95.8	98.5	97.3	95.3	79.9	87.2	75.5	65.1	57.0	88.3
	± 0.1	± 0.1	± 0.0	± 0.1	± 0.4	± 0.2	± 0.6	± 0.5	± 1.4	± 2.6	± 4.1	± 2.6	± 2.6	± 0.5
$A_{GAP}+CCE+Cutout^*$	99.7	99.5	99.3	99.0	96.8	98.8	97.7	96.0	86.1	92.1	82.5	70.7	67.2	91.2
	± 0.1	± 0.0	± 0.1	± 0.1	± 0.2	± 0.1	± 0.2	± 0.4	± 1.5	± 0.9	± 2.6	± 1.7	± 1.5	± 0.4
$A_{FC}+CCE+Cutout$	99.7	99.5	99.3	99.2	97.4	98.7	97.8	96.6	85.7	90.7	77.6	71.7	68.1	90.9
	± 0.0	± 0.1	± 0.1	± 0.2	± 0.2	± 0.3	± 0.6	± 0.7	± 2.1	± 1.0	± 3.9	± 2.4	± 1.8	± 0.7
$A_{FC}+CCE+Cutout^*$	99.5	98.9	97.7	98.5	97.7	97.7	90.9	95.9	92.6	85.6	54.2	75.6	82.7	89.8
	± 0.0	± 0.1	± 1.0	± 0.2	± 0.4	± 0.5	± 6.2	± 0.4	± 0.7	± 2.6	± 22.1	± 2.6	± 1.9	± 2.3
$A_{GAP}+BCE+Cutout$	99.2	97.5	97.1	96.6	92.2	94.8	92.5	90.5	76.6	73.4	66.3	59.8	58.8	84.2
	± 0.1	± 0.3	± 0.4	± 0.3	± 1.5	± 0.9	± 0.7	± 0.9	± 4.8	± 4.4	± 3.2	± 3.3	± 6.4	± 1.6
$A_{GAP}+BCE+Cutout^*$	99.7	99.2	99.0	98.9	98.0	98.2	97.0	96.4	91.5	88.2	79.5	77.6	80.3	92.6
	± 0.1	± 0.2	± 0.3	± 0.1	± 0.4	± 0.2	± 0.3	± 0.3	± 1.3	± 1.8	± 3.5	± 1.7	± 2.7	± 0.6
$A_{FC}+BCE+Cutout$	99.8	99.5	99.5	99.5	98.5	99.2	98.4	97.7	90.3	93.1	82.5	75.6	73.0	92.8
	± 0.0	± 0.1	± 0.1	± 0.0	± 0.2	± 0.1	± 0.3	± 0.2	± 0.9	± 0.6	± 2.8	± 1.8	± 1.5	± 0.3
$A_{FC}+BCE+Cutout^*$	99.8	99.6	99.5	99.5	98.6	99.4	98.4	97.9	91.1	93.6	82.6	76.3	75.0	93.2
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.1	± 0.5	± 0.3	± 1.0	± 0.8	± 3.7	± 2.1	± 2.1	± 0.4

TABLE S-IV: **Cutout Ablations** Accuracies (in %) and their respective standard deviations (denoted by \pm) of our Cutout ablation study on the Occluded-COCO-Vehicles dataset. * marks the results received after fine-tuning the network. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy losses, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and with a convolutional layer, respectively.

Train Data	MS-COCO				
Occ. Area	L0	L1	L2	L3	Avg
$A_{GAP}+CCE+Cutout$	98.7	86.1	80.9	71.2	84.2
	± 0.4	± 0.9	± 2.0	± 5.1	± 1.7
$A_{GAP}+CCE+Cutout^*$	99.0	88.2	82.3	66.8	84.1
	± 0.2	± 2.9	± 5.7	± 4.3	± 3.2
$A_{FC}+CCE+Cutout$	99.6	91.6	84.4	71.2	86.7
	± 0.1	± 0.6	± 3.2	± 3.1	± 1.6
$A_{FC}+CCE+Cutout^*$	99.2	90.4	79.5	63.7	83.2
	± 0.1	± 0.6	± 2.4	± 5.9	± 1.8
$A_{GAP}+BCE+Cutout$	99.2	87.9	83.1	76.7	86.7
	± 0.1	± 0.7	± 1.5	± 1.4	± 0.8
$A_{GAP}+BCE+Cutout^*$	99.2	91.2	87.2	75.3	88.2
	± 0.2	± 1.0	± 2.2	± 4.8	± 1.9
$A_{FC}+BCE+Cutout$	99.6	92.1	86.3	75.0	88.2
	± 0.1	± 0.8	± 1.2	± 2.2	± 0.8
$A_{FC}+BCE+Cutout^*$	99.6	94.1	87.6	73.3	88.6
	± 0.1	± 0.8	± 1.2	± 3.1	± 1.0

TABLE S-V: **Inverted Cutout Mask Size Ablation** The accuracies (in %) and their respective standard deviations (denoted by \pm) of the cutout size ablation study of our Inverted Cutout method. * marks the results received after fine-tuning the whole network. The occlusion types are: w - white box occlusion, n - noise box occlusion, t - texture occlusion, o - occlusion by segmented objects. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy loss, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and the one utilizing a convolutional layer, respectively.

Occ. Area	L0: 0%	L1: 20-40%				L2: 40-60%				L3: 60-80%				Mean	
Occ. Type	-	w	n	t	o	w	n	t	o	w	n	t	o		
Mask size 16	$A_{GAP}+CCE+IC$	56.4	55.1	53.4	53.8	52.6	52.7	51.7	51.6	49.8	46.4	45.7	43.7	48.0	50.8
		± 4.1	± 3.1	± 3.1	± 2.9	± 2.1	± 3.0	± 2.1	± 1.8	± 2.3	± 2.2	± 3.0	± 2.1	± 4.1	± 1.4
	$A_{GAP}+CCE+IC^*$	62.6	58.9	59.8	59.8	48.0	54.5	54.6	54.8	41.5	49.1	41.2	42.3	35.0	50.9
		± 7.2	± 8.0	± 5.6	± 6.3	± 6.4	± 8.2	± 6.2	± 4.8	± 6.3	± 7.9	± 10.8	± 2.3	± 6.2	± 5.7
	$A_{FC}+CCE+IC$	48.5	50.9	48.2	48.4	47.3	52.1	48.4	48.4	47.2	55.8	47.1	46.9	47.1	48.9
		± 1.2	± 2.0	± 1.0	± 1.1	± 0.3	± 2.3	± 1.0	± 0.9	± 0.2	± 2.8	± 0.6	± 0.7	± 0.1	± 1.1
	$A_{FC}+CCE+IC^*$	89.7	88.3	77.0	85.5	80.5	86.6	56.9	80.8	76.6	80.8	12.6	66.0	72.2	73.3
		± 1.3	± 0.9	± 2.4	± 1.5	± 2.1	± 1.1	± 3.8	± 1.6	± 2.3	± 0.8	± 3.8	± 2.1	± 2.5	± 1.1
	$A_{GAP}+BCE+IC$	49.4	50.8	49.4	47.2	42.0	49.3	48.1	44.2	37.8	41.7	42.2	34.9	35.6	44.1
	± 9.9	± 9.1	± 8.2	± 11.2	± 11.9	± 7.3	± 8.3	± 11.1	± 10.6	± 6.3	± 9.1	± 10.7	± 9.9	± 9.1	
$A_{GAP}+BCE+IC^*$	69.4	69.8	65.3	68.1	61.9	68.3	55.4	63.3	58.1	60.7	32.3	45.3	52.0	59.2	
	± 1.0	± 2.7	± 4.0	± 1.5	± 3.4	± 3.7	± 10.9	± 1.2	± 4.1	± 6.5	± 15.3	± 3.8	± 6.0	± 1.8	
$A_{FC}+BCE+IC$	55.2	59.0	54.5	54.6	50.4	61.3	54.7	54.3	49.6	65.2	54.1	51.6	49.2	54.9	
	± 3.0	± 3.0	± 2.7	± 2.6	± 1.5	± 2.8	± 2.0	± 1.9	± 1.1	± 2.6	± 1.7	± 1.5	± 1.2	± 2.1	
$A_{FC}+BCE+IC^*$	73.7	75.5	70.4	71.6	67.5	75.6	67.3	68.6	64.4	73.5	54.6	59.5	62.8	68.1	
	± 4.6	± 4.2	± 4.0	± 4.4	± 5.2	± 4.0	± 3.7	± 3.9	± 4.9	± 3.5	± 1.7	± 2.8	± 4.5	± 3.8	
Mask size 32	$A_{GAP}+CCE+IC$	77.0	80.0	74.2	73.9	69.4	80.3	73.5	69.8	62.3	77.5	66.9	56.2	53.1	70.3
		± 2.3	± 2.1	± 1.6	± 1.9	± 2.4	± 1.9	± 1.8	± 1.8	± 2.4	± 2.2	± 1.1	± 1.1	± 3.2	± 1.4
	$A_{GAP}+CCE+IC^*$	90.9	92.9	88.7	89.4	81.8	93.1	86.5	85.0	73.8	90.2	75.8	67.9	63.5	83.1
		± 1.0	± 0.7	± 1.1	± 1.1	± 1.6	± 0.9	± 1.4	± 1.3	± 2.1	± 0.8	± 2.4	± 1.8	± 2.5	± 1.1
	$A_{FC}+CCE+IC$	88.3	88.4	84.0	84.5	82.3	87.7	82.0	80.4	78.9	85.1	73.6	69.1	74.8	81.5
		± 1.3	± 0.9	± 1.1	± 0.9	± 1.0	± 1.2	± 1.2	± 0.9	± 0.9	± 0.6	± 1.1	± 1.2	± 1.5	± 0.9
	$A_{FC}+CCE+IC^*$	97.9	97.9	96.4	97.2	95.2	97.1	91.6	94.0	91.8	94.0	59.0	78.7	85.6	90.5
		± 0.1	± 0.1	± 0.2	± 0.1	± 0.4	± 0.2	± 2.7	± 0.3	± 0.3	± 0.5	± 12.8	± 1.0	± 1.1	± 1.4
	$A_{GAP}+BCE+IC$	75.4	77.7	72.3	74.3	72.6	78.6	71.3	72.8	70.1	74.4	65.4	63.9	64.5	71.8
	± 3.3	± 3.0	± 2.2	± 2.4	± 1.7	± 3.4	± 1.5	± 1.6	± 3.3	± 6.3	± 1.6	± 2.2	± 4.6	± 2.7	
$A_{GAP}+BCE+IC^*$	96.4	96.3	95.5	95.8	93.1	95.9	92.4	93.6	90.3	91.7	67.9	80.4	85.1	90.3	
	± 0.6	± 0.6	± 0.8	± 0.6	± 0.8	± 0.7	± 2.1	± 0.7	± 0.9	± 1.9	± 10.3	± 0.7	± 1.5	± 1.4	
$A_{FC}+BCE+IC$	91.2	90.7	87.9	88.4	88.5	90.1	86.6	85.9	85.4	87.8	79.6	75.5	81.2	86.1	
	± 0.7	± 0.9	± 1.1	± 1.1	± 1.1	± 0.7	± 0.8	± 1.3	± 1.1	± 1.3	± 1.0	± 1.3	± 1.3	± 1.0	
$A_{FC}+BCE+IC^*$	97.0	97.0	95.6	96.0	94.1	96.4	93.7	93.5	90.1	93.8	84.2	82.9	85.7	92.3	
	± 0.4	± 0.6	± 1.1	± 0.8	± 1.2	± 0.7	± 1.3	± 1.4	± 1.4	± 1.0	± 1.9	± 1.3	± 1.0	± 1.0	
Mask size 64	$A_{GAP}+CCE+IC$	98.3	98.3	97.4	97.4	94.9	98.0	96.0	94.6	89.6	95.4	87.9	77.5	80.4	92.7
		± 0.1	± 0.1	± 0.1	± 0.2	± 0.5	± 0.1	± 0.2	± 0.2	± 1.5	± 0.2	± 0.7	± 1.2	± 2.8	± 0.4
	$A_{GAP}+CCE+IC^*$	98.9	99.1	98.8	98.5	97.1	98.8	97.7	96.6	92.7	97.2	89.4	80.8	84.5	94.6
		± 0.1	± 0.1	± 0.1	± 0.2	± 0.3	± 0.2	± 0.2	± 0.5	± 1.2	± 0.2	± 2.2	± 0.9	± 3.1	± 0.5
	$A_{FC}+CCE+IC$	99.1	99.1	98.6	98.5	96.8	98.5	97.4	96.2	91.2	95.6	87.3	80.9	83.6	94.1
		± 0.1	± 0.1	± 0.1	± 0.2	± 0.1	± 0.1	± 0.1	± 0.3	± 0.3	± 0.4	± 1.1	± 0.6	± 0.8	± 0.1
	$A_{FC}+CCE+IC^*$	99.5	99.4	98.7	99.1	98.4	99.3	95.9	97.6	96.4	96.3	68.9	80.8	90.9	93.9
		± 0.1	± 0.1	± 0.4	± 0.1	± 0.1	± 0.1	± 2.4	± 0.2	± 0.4	± 0.2	± 13.2	± 0.9	± 0.9	± 1.4
	$A_{GAP}+BCE+IC$	97.8	97.7	96.7	97.1	94.6	96.9	94.9	95.4	89.9	93.2	83.2	81.5	81.7	92.4
	± 0.4	± 0.2	± 0.3	± 0.4	± 0.7	± 0.3	± 0.5	± 0.5	± 1.2	± 0.5	± 3.7	± 3.2	± 3.5	± 1.0	
$A_{GAP}+BCE+IC^*$	99.4	99.5	99.0	99.4	98.6	99.2	97.7	98.3	96.5	96.4	77.7	85.1	91.7	95.3	
	± 0.1	± 0.0	± 0.1	± 0.2	± 0.1	± 0.0	± 0.4	± 0.2	± 0.2	± 0.4	± 6.4	± 2.3	± 0.3	± 0.6	
$A_{FC}+BCE+IC$	99.3	99.3	99.0	98.9	97.9	98.9	98.1	97.6	94.9	96.3	91.5	85.0	87.9	95.7	
	± 0.1	± 0.1	± 0.1	± 0.2	± 0.3	± 0.1	± 0.3	± 0.2	± 0.4	± 0.2	± 0.7	± 0.7	± 0.7	± 0.1	
$A_{FC}+BCE+IC^*$	99.6	99.7	99.4	99.3	98.5	99.5	98.4	98.5	96.6	97.6	87.0	84.9	91.0	96.2	
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.1	± 0.1	± 0.1	± 0.4	± 0.2	± 2.1	± 0.3	± 0.6	± 0.2	
Mask size 128	$A_{GAP}+CCE+IC$	99.7	99.4	99.2	99.0	97.7	99.2	97.8	96.8	91.2	95.9	84.7	74.5	78.2	93.3
		± 0.0	± 0.0	± 0.1	± 0.1	± 0.2	± 0.1	± 0.4	± 0.3	± 1.1	± 0.4	± 3.0	± 1.9	± 2.8	± 0.7
	$A_{GAP}+CCE+IC^*$	99.7	99.6	99.4	99.1	97.6	99.3	98.0	96.9	91.2	96.4	81.8	73.2	79.3	93.2
		± 0.0	± 0.1	± 0.1	± 0.0	± 0.3	± 0.1	± 0.3	± 0.3	± 1.1	± 0.3	± 3.3	± 1.3	± 2.1	± 0.5
	$A_{FC}+CCE+IC$	99.7	99.5	99.3	99.3	98.4	99.4	98.2	97.6	93.8	95.3	80.9	75.6	83.8	93.9
		± 0.0	± 0.0	± 0.1	± 0.1	± 0.1	± 0.1	± 0.3	± 0.1	± 1.0	± 0.3	± 2.9	± 1.2	± 2.5	± 0.5
	$A_{FC}+CCE+IC^*$	99.7	99.5	99.3	99.3	98.8	99.1	97.7	97.4	95.6	93.1	78.4	77.1	88.6	94.1
		± 0.1	± 0.0	± 0.1	± 0.1	± 0.2	± 0.1	± 0.5	± 0.4	± 0.4	± 0.4	± 6.8	± 1.9	± 0.6	± 0.5
	$A_{GAP}+BCE+IC$	99.3	98.9	98.6	98.8	97.5	98.1	96.4	96.9	92.2	91.2	78.4	77.5	81.3	92.7
	± 0.0	± 0.1	± 0.1	± 0.2	± 0.4	± 0.3	± 0.6	± 0.3	± 1.3	± 0.8	± 2.7	± 1.8	± 2.9	± 0.8	
$A_{GAP}+BCE+IC^*$	99.8	99.7	99.5	99.5	99.1	99.4	98.6	98.6	96.7	94.2	80.7	82.5	89.7	95.2	
	± 0.0	± 0.0	± 0.1	± 0.1	± 0.1	± 0.1	± 0.3	± 0.3	± 0.3	± 1.5	± 3.5	± 1.6	± 1.0	± 0.3	
$A_{FC}+BCE+IC$	99.8	99.6	99.4	99.4	98.9	99.5	98.5	98.4	95.4	96.2	84.1	79.7	85.2	94.9	
	± 0.0	± 0.1	± 0.1	± 0.1	± 0.1										

TABLE S-VI: **Inverted Cutout Mask Size Ablation** Accuracies (in %) and their respective standard deviations (denoted by \pm) of our cutout size ablations on the Occluded-COCO-Vehicles dataset. * marks the results received after fine-tuning the network. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy losses, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and with a convolutional layer, respectively.

Train Data	MS-COCO						
Occ. Area	L0	L1	L2	L3	Avg		
$A_{GAP}+CCE+IC$	43.2	19.1	14.8	23.3	25.1		
	± 7.5	± 2.9	± 3.1	± 9.2	± 4.9		
$A_{GAP}+CCE+IC^*$	54.7	28.9	21.4	20.5	31.4		
	± 11.0	± 7.5	± 7.8	± 6.5	± 7.8		
$A_{FC}+CCE+IC$	72.3	39.6	30.0	29.8	42.9		
	± 3.4	± 3.6	± 1.6	± 2.6	± 2.1		
$A_{FC}+CCE+IC^*$	87.9	69.3	61.6	46.6	66.3		
	± 1.5	± 2.2	± 2.8	± 3.1	± 1.6		
$A_{GAP}+BCE+IC$	50.1	25.7	21.4	34.9	33.0		
	± 2.9	± 1.8	± 1.3	± 2.5	± 1.2		
$A_{GAP}+BCE+IC^*$	64.3	34.5	25.2	21.6	36.4		
	± 6.1	± 8.7	± 8.1	± 8.1	± 7.6		
$A_{FC}+BCE+IC$	75.5	41.5	32.1	32.5	45.4		
	± 1.2	± 0.9	± 1.0	± 2.0	± 0.9		
$A_{FC}+BCE+IC^*$	83.8	62.8	53.2	40.8	60.1		
	± 3.8	± 3.1	± 4.0	± 3.4	± 3.1		
Mask size 16	$A_{GAP}+CCE+IC$	34.1	21.3	15.8	16.8	22.0	
		± 1.3	± 1.0	± 1.6	± 8.6	± 2.6	
	$A_{GAP}+CCE+IC^*$	74.6	36.9	25.5	30.8	42.0	
		± 5.4	± 7.9	± 4.8	± 5.2	± 3.9	
	$A_{FC}+CCE+IC$	95.8	81.6	69.3	52.4	74.8	
		± 0.5	± 2.0	± 1.9	± 3.7	± 1.8	
	$A_{FC}+CCE+IC^*$	97.9	88.2	84.8	76.0	86.7	
		± 0.2	± 0.8	± 1.2	± 2.8	± 1.1	
	$A_{GAP}+BCE+IC$	82.5	71.4	58.3	59.9	68.0	
		± 6.0	± 2.9	± 4.2	± 8.8	± 4.9	
Mask size 32	$A_{GAP}+BCE+IC^*$	95.0	82.1	77.8	73.6	82.1	
		± 0.6	± 0.7	± 1.8	± 3.0	± 1.3	
	$A_{FC}+BCE+IC$	95.5	81.0	67.0	50.7	73.5	
		± 0.5	± 1.8	± 1.6	± 4.8	± 2.0	
	$A_{FC}+BCE+IC^*$	98.4	88.2	85.5	73.3	86.3	
		± 0.2	± 0.6	± 0.7	± 2.3	± 0.9	
	Mask size 64	$A_{GAP}+CCE+IC$	97.7	86.4	83.7	76.0	86.0
			± 0.4	± 1.3	± 1.5	± 3.0	± 0.9
		$A_{GAP}+CCE+IC^*$	99.0	91.8	89.6	84.2	91.1
			± 0.1	± 1.1	± 1.1	± 4.4	± 1.6
$A_{FC}+CCE+IC$		99.1	92.2	88.0	76.0	88.8	
		± 0.1	± 0.3	± 2.0	± 1.5	± 0.9	
$A_{FC}+CCE+IC^*$		99.3	93.4	89.4	84.2	91.6	
		± 0.1	± 0.7	± 0.7	± 2.8	± 1.0	
$A_{GAP}+BCE+IC$		97.9	87.3	85.3	82.2	88.2	
		± 0.3	± 1.0	± 1.3	± 1.0	± 0.5	
Mask size 128	$A_{GAP}+BCE+IC^*$	99.3	92.7	88.6	86.3	91.7	
		± 0.1	± 0.4	± 0.4	± 2.6	± 0.8	
	$A_{FC}+BCE+IC$	99.1	92.7	90.1	79.1	90.3	
		± 0.1	± 0.4	± 0.6	± 2.2	± 0.6	
	$A_{FC}+BCE+IC^*$	99.4	94.3	91.0	87.0	92.9	
		± 0.0	± 0.6	± 0.6	± 2.1	± 0.5	
	Mask size 16	$A_{GAP}+CCE+IC$	99.0	90.0	86.0	76.0	87.8
			± 0.1	± 0.7	± 1.0	± 2.3	± 0.9
		$A_{GAP}+CCE+IC^*$	99.5	94.2	88.4	86.6	92.2
			± 0.1	± 0.4	± 0.5	± 2.6	± 0.8
$A_{FC}+CCE+IC$		99.5	93.8	88.9	79.1	90.3	
		± 0.1	± 0.6	± 1.3	± 2.6	± 0.6	
$A_{FC}+CCE+IC^*$		99.2	93.8	87.2	76.7	89.2	
		± 0.1	± 0.5	± 0.1	± 4.4	± 1.3	
$A_{GAP}+BCE+IC$		99.1	90.0	86.8	78.8	88.7	
		± 0.1	± 0.4	± 1.2	± 2.1	± 0.6	
Mask size 128	$A_{GAP}+BCE+IC^*$	99.5	93.9	89.6	84.9	92.0	
		± 0.1	± 1.0	± 0.4	± 3.2	± 1.1	
	$A_{FC}+BCE+IC$	99.6	94.7	90.4	83.2	92.0	
		± 0.1	± 0.3	± 0.9	± 1.1	± 0.2	
	$A_{FC}+BCE+IC^*$	99.6	94.8	90.8	87.7	93.2	
		± 0.1	± 0.5	± 1.1	± 1.4	± 0.6	

TABLE S-VII: **Feature Extraction Layer Ablations** The accuracies (in %) and their respective standard deviations (denoted by \pm) of the feature extraction layer ablation study of our Inverted Cutout method. * marks the results received after fine-tuning the whole network. The occlusion types are: w - white box occlusion, n - noise box occlusion, t - texture occlusion, o - occlusion by segmented objects. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy loss, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and the one utilizing a convolutional layer, respectively.

Occ. Area	L0: 0%	L1: 20-40%				L2: 40-60%				L3: 60-80%				Mean	
Occ. Type	-	w	n	t	o	w	n	t	o	w	n	t	o		
Extraction Block P4	$A_{GAP}+CCE+IC$	99.4	99.0	98.6	98.4	96.5	98.7	97.6	95.6	90.3	94.1	85.1	71.4	77.6	92.5
		± 0.1	± 0.2	± 0.1	± 0.1	± 0.2	± 0.1	± 0.3	± 0.2	± 1.0	± 0.4	± 1.5	± 0.7	± 1.5	± 0.3
	$A_{GAP}+CCE+IC^*$	99.5	99.3	99.1	98.8	97.4	98.9	98.3	95.7	92.3	95.6	87.0	71.8	81.8	93.5
		± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.1	± 0.1	± 0.3	± 0.7	± 0.2	± 1.4	± 1.3	± 1.7	± 0.2
	$A_{FC}+CCE+IC$	99.4	99.4	98.9	99.0	97.1	99.0	97.5	96.5	91.0	95.7	83.7	75.2	80.0	93.3
		± 0.1	± 0.1	± 0.1	± 0.1	± 0.4	± 0.0	± 0.2	± 0.1	± 1.1	± 0.2	± 0.6	± 0.6	± 2.2	± 0.3
	$A_{FC}+CCE+IC^*$	99.6	99.5	99.0	99.4	98.3	99.3	97.4	97.0	94.9	95.1	79.6	77.2	86.7	94.1
		± 0.0	± 0.0	± 0.2	± 0.1	± 0.0	± 0.1	± 0.2	± 0.2	± 0.5	± 1.1	± 2.5	± 0.9	± 1.1	± 0.3
	$A_{GAP}+BCE+IC$	98.7	98.2	98.0	97.7	96.2	97.2	95.6	94.5	88.9	90.6	84.2	73.7	76.8	91.6
		± 0.2	± 0.1	± 0.2	± 0.3	± 0.5	± 0.5	± 0.6	± 0.5	± 1.1	± 1.6	± 2.0	± 1.2	± 1.8	± 0.7
$A_{GAP}+BCE+IC^*$	99.5	99.4	98.9	99.2	98.6	99.0	96.7	97.5	96.0	95.5	78.0	82.0	89.7	94.6	
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.8	± 0.2	± 0.4	± 0.5	± 2.6	± 1.2	± 0.8	± 0.2	
$A_{FC}+BCE+IC$	99.5	99.4	99.2	99.1	98.2	99.2	98.2	97.5	93.0	95.8	90.4	80.0	82.8	94.8	
	± 0.1	± 0.1	± 0.1	± 0.0	± 0.2	± 0.1	± 0.2	± 0.1	± 0.5	± 0.3	± 0.1	± 0.5	± 1.4	± 0.2	
$A_{FC}+BCE+IC^*$	99.6	99.6	99.4	99.3	98.5	99.3	98.6	97.6	94.4	97.2	91.8	80.8	85.5	95.5	
	± 0.1	± 0.1	± 0.1	± 0.0	± 0.1	± 0.1	± 0.2	± 0.2	± 0.2	± 0.2	± 0.8	± 0.6	± 0.7	± 0.1	
Extraction Block P5	$A_{GAP}+CCE+IC$	99.3	99.1	98.5	98.5	94.6	97.9	95.9	95.6	79.6	90.6	74.2	70.5	58.6	88.7
		± 0.1	± 0.1	± 0.1	± 0.2	± 0.5	± 0.1	± 0.5	± 0.3	± 2.4	± 0.7	± 3.5	± 2.1	± 3.7	± 0.9
	$A_{GAP}+CCE+IC^*$	99.8	99.6	99.4	99.3	97.7	99.2	97.8	97.9	91.1	95.3	79.3	78.8	77.1	93.2
		± 0.0	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.5	± 0.1	± 1.3	± 0.1	± 3.4	± 0.6	± 3.2	± 0.6
	$A_{FC}+CCE+IC$	99.2	99.2	99.0	98.8	96.8	98.5	96.7	96.4	86.9	92.4	72.9	70.8	69.7	90.6
		± 0.1	± 0.1	± 0.2	± 0.1	± 0.2	± 0.1	± 0.6	± 0.4	± 2.1	± 0.2	± 2.9	± 1.2	± 3.9	± 0.3
	$A_{FC}+CCE+IC^*$	99.6	99.5	99.4	99.4	98.7	99.2	98.8	98.2	96.3	96.1	88.3	83.4	89.6	95.9
		± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.1	± 0.2	± 0.1	± 0.2	± 0.8	± 1.6	± 0.3	± 1.3	± 0.1
	$A_{GAP}+BCE+IC$	99.0	98.7	98.2	98.3	95.6	97.6	95.8	95.9	86.6	90.1	78.2	75.2	69.3	90.7
		± 0.1	± 0.2	± 0.1	± 0.2	± 0.5	± 0.3	± 0.4	± 0.2	± 1.6	± 0.5	± 2.5	± 1.5	± 2.9	± 0.6
$A_{GAP}+BCE+IC^*$	99.6	99.6	99.4	99.4	98.9	99.2	98.9	98.5	96.8	96.9	91.8	86.2	91.0	96.6	
	± 0.0	± 0.1	± 0.0	± 0.0	± 0.1	± 0.1	± 0.2	± 0.1	± 0.3	± 0.4	± 1.6	± 0.6	± 0.7	± 0.1	
$A_{FC}+BCE+IC$	99.4	99.3	98.9	99.1	97.5	98.6	97.1	97.1	88.9	93.5	78.2	74.2	71.8	91.8	
	± 0.1	± 0.0	± 0.1	± 0.1	± 0.3	± 0.2	± 0.2	± 0.3	± 1.8	± 0.4	± 1.6	± 0.2	± 3.9	± 0.4	
$A_{FC}+BCE+IC^*$	99.7	99.6	99.4	99.4	98.9	99.5	98.6	98.6	95.0	97.0	86.6	84.7	84.3	95.5	
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 1.1	± 0.2	± 1.7	± 1.0	± 2.3	± 0.2	
Extraction Block P4 & P5	$A_{GAP}+CCE+IC$	99.5	99.2	99.1	98.9	96.8	99.0	97.7	96.3	90.8	95.2	86.6	74.4	78.0	93.2
		± 0.1	± 0.1	± 0.1	± 0.3	± 0.5	± 0.1	± 0.2	± 0.4	± 2.5	± 0.3	± 1.8	± 2.7	± 5.4	± 0.9
	$A_{GAP}+CCE+IC^*$	99.7	99.4	99.3	99.0	97.8	99.2	98.3	96.9	92.8	96.4	87.0	75.4	82.3	94.1
		± 0.0	± 0.1	± 0.1	± 0.0	± 0.2	± 0.2	± 0.2	± 0.4	± 1.4	± 0.3	± 1.9	± 2.3	± 3.4	± 0.6
	$A_{FC}+CCE+IC$	99.5	99.5	99.2	99.3	97.9	99.2	98.2	97.5	93.7	96.8	87.8	78.9	85.1	94.8
		± 0.0	± 0.0	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.5	± 0.4	± 1.9	± 1.6	± 0.7	± 0.3
	$A_{FC}+CCE+IC^*$	99.6	99.6	99.2	99.4	98.6	99.3	97.9	97.5	95.7	95.9	83.0	80.4	88.9	95.0
		± 0.1	± 0.1	± 0.2	± 0.1	± 0.1	± 0.2	± 0.5	± 0.2	± 0.6	± 0.8	± 4.8	± 1.3	± 1.1	± 0.6
	$A_{GAP}+BCE+IC$	99.1	98.9	98.3	98.3	97.4	98.0	96.4	96.7	93.1	92.8	83.8	80.5	84.8	93.7
		± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.3	± 0.3	± 0.1	± 0.6	± 0.6	± 1.7	± 1.4	± 1.0	± 0.3
$A_{GAP}+BCE+IC^*$	99.7	99.6	99.4	99.4	98.9	99.2	98.4	98.1	96.9	95.8	87.4	85.4	92.1	96.2	
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.3	± 0.2	± 0.3	± 2.4	± 1.1	± 1.0	± 0.2	
$A_{FC}+BCE+IC$	99.6	99.6	99.3	99.3	98.7	99.4	98.3	98.1	95.4	97.0	90.5	82.8	87.4	95.8	
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.0	± 0.1	± 0.1	± 0.2	± 0.6	± 0.1	± 1.3	± 1.2	± 1.0	± 0.3	
$A_{FC}+BCE+IC^*$	99.7	99.7	99.4	99.4	98.9	99.6	98.8	98.5	96.2	97.9	90.8	83.9	89.4	96.3	
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.2	± 0.3	± 0.1	± 1.6	± 1.0	± 1.0	± 0.2	

TABLE S-VIII: **Feature Extraction Layer Ablations** Accuracies (in %) and their respective standard deviations (denoted by \pm) of our feature extraction layer ablations on the Occluded-COCO-Vehicles dataset. * marks the results received after fine-tuning the network. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy losses, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and with a convolutional layer, respectively.

Train Data	MS-COCO				
Occ. Area	L0	L1	L2	L3	Avg
$A_{GAP}+CCE+IC$	98.5	87.6	82.9	71.2	85.1
	± 0.1	± 1.2	± 2.4	± 4.1	± 1.3
$A_{GAP}+CCE+IC^*$	99.2	92.9	89.7	84.6	91.6
	± 0.1	± 0.1	± 0.5	± 2.6	± 0.8
$A_{FC}+CCE+IC$	99.4	92.7	88.5	79.1	89.9
	± 0.0	± 0.3	± 1.4	± 3.5	± 1.1
$A_{FC}+CCE+IC^*$	99.5	94.1	90.4	79.5	90.9
	± 0.1	± 0.0	± 0.7	± 2.6	± 0.8
$A_{GAP}+BCE+IC$	98.8	88.1	85.4	79.8	88.0
	± 0.1	± 0.6	± 1.3	± 1.8	± 0.7
$A_{GAP}+BCE+IC^*$	99.1	93.2	89.8	89.4	92.9
	± 0.2	± 0.5	± 0.6	± 2.8	± 0.7
$A_{FC}+BCE+IC$	99.3	92.8	88.7	80.1	90.2
	± 0.0	± 0.5	± 0.6	± 1.5	± 0.4
$A_{FC}+BCE+IC^*$	99.5	94.7	90.8	87.0	93.0
	± 0.1	± 0.3	± 1.5	± 3.1	± 1.2
$A_{GAP}+CCE+IC$	98.9	88.2	84.6	73.6	86.3
	± 0.1	± 0.8	± 1.2	± 0.6	± 0.5
$A_{GAP}+CCE+IC^*$	99.6	94.2	91.7	85.3	92.7
	± 0.1	± 0.6	± 1.0	± 0.6	± 0.5
$A_{FC}+CCE+IC$	99.4	92.7	86.5	78.1	89.2
	± 0.1	± 0.3	± 1.6	± 3.5	± 1.3
$A_{FC}+CCE+IC^*$	99.6	95.1	92.2	84.6	92.9
	± 0.1	± 0.2	± 0.4	± 0.6	± 0.1
$A_{GAP}+BCE+IC$	99.1	89.8	85.7	76.0	87.7
	± 0.1	± 0.2	± 2.1	± 1.5	± 0.6
$A_{GAP}+BCE+IC^*$	99.4	94.9	91.1	88.4	93.4
	± 0.0	± 0.2	± 1.2	± 3.4	± 1.0
$A_{FC}+BCE+IC$	99.5	93.1	88.2	81.8	90.6
	± 0.0	± 0.3	± 0.9	± 2.0	± 0.7
$A_{FC}+BCE+IC^*$	99.5	95.2	94.0	90.4	94.8
	± 0.1	± 0.6	± 0.7	± 0.0	± 0.2
$A_{GAP}+CCE+IC$	99.1	90.6	87.2	87.0	91.0
	± 0.1	± 0.5	± 1.1	± 1.5	± 0.6
$A_{GAP}+CCE+IC^*$	99.5	94.5	90.8	89.4	93.5
	± 0.1	± 0.2	± 1.3	± 3.0	± 1.0
$A_{FC}+CCE+IC$	99.4	94.2	90.0	82.9	91.6
	± 0.1	± 0.8	± 0.8	± 3.0	± 1.0
$A_{FC}+CCE+IC^*$	99.4	94.7	90.8	84.2	92.3
	± 0.1	± 0.6	± 1.0	± 0.7	± 0.4
$A_{GAP}+BCE+IC$	99.1	89.4	88.0	86.0	90.6
	± 0.1	± 0.7	± 0.5	± 2.4	± 0.8
$A_{GAP}+BCE+IC^*$	99.4	94.6	91.6	92.8	94.6
	± 0.0	± 0.8	± 0.5	± 1.8	± 0.4
$A_{FC}+BCE+IC$	99.4	94.1	92.0	86.3	93.0
	± 0.1	± 0.6	± 0.4	± 1.7	± 0.3
$A_{FC}+BCE+IC^*$	99.5	95.4	92.1	88.7	93.9
	± 0.0	± 0.4	± 0.4	± 2.4	± 0.7

TABLE S-IX: The hyperparameters used in the training process of the ResNet50-networks. * denotes hyperparameters used during fine-tuning of the network, A_{GAP} denotes the aggregation module with global average pooling, and A_{FC} the module with kernel size equal to the input feature map. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy losses, respectively.

Dataset	Agg.	Loss	LR	EP	LR*	EP*
Pascal3D+	A_{GAP}	CCE	1e-3	90	1e-5	90
		BCE	1e-3	90	1e-5	90
	A_{FC}	CCE	1e-4	90	1e-5	90
		BCE	1e-4	90	1e-5	90
MSCOCO	A_{GAP}	CCE	1e-3	180	1e-4	90
		BCE	1e-4	180	1e-4	90
	A_{FC}	CCE	1e-4	180	1e-5	90
		BCE	1e-4	180	1e-4	90

TABLE S-X: **ResNet50 Ablations** The accuracies (in %) and their respective standard deviations (denoted by \pm) of the ResNet50 ablation study of our Inverted Cutout method. * marks the results received after fine-tuning the whole network. The occlusion types are: w - white box occlusion, n - noise box occlusion, t - texture occlusion, o - occlusion by segmented objects. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy loss, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and the one utilizing a convolutional layer, respectively.

Occ. Area	L0: 0%	L1: 20-40%				L2: 40-60%				L3: 60-80%				Mean
Occ. Type	-	w	n	t	o	w	n	t	o	w	n	t	o	
$A_{GAP}+CCE$	99.7	98.7	98.5	98.5	93.6	94.4	92.9	93.1	72.5	69.1	61.3	61.8	52.8	83.6
	± 0.0	± 0.2	± 0.2	± 0.2	± 0.5	± 0.4	± 0.3	± 0.5	± 1.2	± 1.7	± 0.7	± 2.0	± 1.8	± 0.5
$A_{GAP}+CCE^*$	99.7	98.3	98.3	97.9	92.2	91.9	90.0	89.9	70.8	58.4	51.3	51.7	52.2	80.2
	± 0.1	± 0.5	± 0.3	± 0.3	± 2.4	± 1.3	± 1.9	± 0.4	± 6.5	± 4.5	± 8.1	± 2.5	± 5.7	± 1.4
$A_{FC}+CCE$	99.8	98.9	98.4	98.4	94.2	94.3	93.2	92.6	76.2	64.9	59.8	57.8	57.9	83.6
	± 0.1	± 0.2	± 0.3	± 0.1	± 0.9	± 0.5	± 1.9	± 1.1	± 1.6	± 3.1	± 6.2	± 4.9	± 2.6	± 1.5
$A_{FC}+CCE^*$	99.8	97.9	97.3	97.6	91.8	91.2	88.0	89.2	72.3	53.6	46.8	47.5	53.6	79.0
	± 0.1	± 0.8	± 0.8	± 0.6	± 2.0	± 3.2	± 4.0	± 3.0	± 6.0	± 11.7	± 9.6	± 7.5	± 7.5	± 4.0
$A_{GAP}+BCE$	99.6	98.8	98.5	98.6	94.0	95.1	93.1	93.0	75.8	70.6	60.5	59.9	58.3	84.3
	± 0.0	± 0.1	± 0.2	± 0.1	± 0.5	± 0.3	± 0.4	± 0.1	± 1.0	± 1.2	± 1.2	± 1.9	± 1.2	± 0.2
$A_{GAP}+BCE^*$	99.8	98.7	97.9	97.9	91.6	92.6	87.3	89.3	70.2	59.8	45.8	48.1	51.2	79.3
	± 0.1	± 0.3	± 0.5	± 0.5	± 2.1	± 1.5	± 2.6	± 1.7	± 3.8	± 2.8	± 3.5	± 2.5	± 4.0	± 1.8
$A_{FC}+BCE$	99.8	99.1	98.7	98.8	94.3	95.1	94.1	93.8	76.6	65.5	60.8	57.7	57.5	84.0
	± 0.1	± 0.1	± 0.2	± 0.0	± 0.2	± 0.3	± 0.4	± 0.2	± 1.0	± 1.1	± 1.3	± 1.8	± 1.2	± 0.3
$A_{FC}+BCE^*$	99.8	98.9	98.6	98.5	92.2	93.5	91.0	90.8	70.9	56.3	48.8	50.3	53.5	80.2
	± 0.0	± 0.2	± 0.3	± 0.4	± 1.7	± 0.8	± 1.6	± 1.7	± 2.9	± 2.1	± 3.7	± 4.4	± 3.3	± 1.7
$A_{GAP}+CCE+IC$	99.5	98.9	96.6	97.7	90.7	96.9	90.1	91.5	71.2	87.2	61.2	57.4	53.3	84.0
	± 0.1	± 0.1	± 0.2	± 0.2	± 1.0	± 0.5	± 0.6	± 0.6	± 3.3	± 0.6	± 3.1	± 1.6	± 3.6	± 1.0
$A_{GAP}+CCE+IC^*$	99.7	99.5	99.3	99.3	98.3	99.2	97.9	97.7	92.0	95.7	79.0	78.1	78.7	93.4
	± 0.0	± 0.0	± 0.1	± 0.1	± 0.2	± 0.1	± 0.2	± 0.3	± 1.2	± 0.4	± 3.8	± 1.5	± 1.8	± 0.6
$A_{FC}+CCE+IC$	99.5	99.4	98.3	98.8	96.8	98.7	95.2	96.0	89.0	92.9	72.8	70.9	78.2	91.3
	± 0.1	± 0.1	± 0.2	± 0.1	± 0.3	± 0.1	± 1.0	± 0.3	± 0.6	± 0.8	± 3.9	± 2.6	± 1.6	± 0.6
$A_{FC}+CCE+IC^*$	99.8	99.5	99.3	99.4	98.7	99.4	98.3	98.2	95.7	96.9	83.1	79.2	87.2	95.0
	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.1	± 0.3	± 0.2	± 1.0	± 0.4	± 1.8	± 1.1	± 1.8	± 0.4
$A_{GAP}+BCE+IC$	99.3	99.1	97.7	98.3	92.9	97.7	93.0	94.4	76.5	89.9	66.4	68.8	60.8	87.3
	± 0.1	± 0.1	± 0.4	± 0.2	± 0.6	± 0.2	± 1.1	± 0.3	± 0.8	± 0.8	± 1.2	± 0.6	± 1.8	± 0.1
$A_{GAP}+BCE+IC^*$	99.8	99.6	99.4	99.5	98.7	99.2	98.2	98.5	95.0	95.6	85.1	84.5	83.8	95.1
	± 0.0	± 0.1	± 0.1	± 0.1	± 0.1	± 0.1	± 0.2	± 0.2	± 0.5	± 0.4	± 2.7	± 1.5	± 0.9	± 0.4
$A_{FC}+BCE+IC$	99.5	99.4	98.8	99.0	96.4	98.9	96.8	97.1	88.8	93.4	80.1	77.5	77.3	92.6
	± 0.1	± 0.0	± 0.1	± 0.1	± 0.1	± 0.1	± 0.4	± 0.3	± 1.0	± 0.3	± 2.0	± 1.6	± 1.9	± 0.5
$A_{FC}+BCE+IC^*$	99.8	99.6	99.4	99.5	99.0	99.5	98.8	98.7	95.8	96.9	87.5	85.5	86.1	95.9
	± 0.1	± 0.0	± 0.1	± 0.1	± 0.2	± 0.1	± 0.2	± 0.3	± 0.7	± 0.4	± 3.6	± 1.2	± 1.3	± 0.5

TABLE S-XI: **ResNet50 Ablations** Accuracies (in %) and their respective standard deviations (denoted by \pm) of our ResNet50 ablations on the Occluded-COCO-Vehicles dataset. * marks the results received after fine-tuning the network. CCE and BCE denote the Categorical (Softmax) and Binary Cross Entropy losses, respectively. A_{GAP} and A_{FC} denote the aggregation module with global average pooling and with a convolutional layer, respectively.

Train Data	MS-COCO				
Occ. Area	L0	L1	L2	L3	Avg
$A_{GAP}+CCE$	99.6 ± 0.1	92.4 ± 0.5	89.0 ± 0.7	79.8 ± 1.1	90.2 ± 0.3
$A_{GAP}+CCE^*$	99.5 ± 0.1	89.3 ± 1.4	79.8 ± 2.0	67.5 ± 6.2	84.0 ± 2.3
$A_{FC}+CCE$	99.6 ± 0.1	93.1 ± 0.9	87.8 ± 1.8	75.0 ± 3.5	88.9 ± 1.5
$A_{FC}+CCE^*$	99.5 ± 0.1	90.0 ± 1.0	82.7 ± 1.6	63.4 ± 3.3	83.9 ± 1.4
$A_{GAP}+BCE$	99.5 ± 0.0	92.1 ± 0.2	88.5 ± 0.3	72.9 ± 0.6	88.3 ± 0.1
$A_{GAP}+BCE^*$	99.7 ± 0.1	91.1 ± 0.8	84.6 ± 2.9	71.6 ± 4.2	86.8 ± 1.9
$A_{FC}+BCE$	99.6 ± 0.1	93.5 ± 0.2	88.7 ± 0.4	75.7 ± 1.1	89.4 ± 0.3
$A_{FC}+BCE^*$	99.5 ± 0.1	91.3 ± 0.4	84.1 ± 2.7	74.7 ± 4.1	87.4 ± 0.9
$A_{GAP}+CCE+IC$	99.2 ± 0.1	91.4 ± 0.5	85.4 ± 2.8	72.6 ± 1.9	87.2 ± 1.2
$A_{GAP}+CCE+IC^*$	99.4 ± 0.1	93.1 ± 0.1	90.0 ± 1.6	85.6 ± 3.1	92.0 ± 1.2
$A_{FC}+CCE+IC$	99.6 ± 0.0	94.8 ± 0.4	91.3 ± 1.1	87.0 ± 1.2	93.2 ± 0.5
$A_{FC}+CCE+IC^*$	99.4 ± 0.1	94.2 ± 0.8	91.0 ± 1.3	88.4 ± 2.5	93.3 ± 1.1
$A_{GAP}+BCE+IC$	99.4 ± 0.0	90.6 ± 0.3	84.2 ± 0.6	76.7 ± 1.7	87.7 ± 0.6
$A_{GAP}+BCE+IC^*$	99.5 ± 0.0	93.7 ± 0.4	92.2 ± 1.1	88.4 ± 2.1	93.4 ± 0.7
$A_{FC}+BCE+IC$	99.5 ± 0.1	95.2 ± 0.3	90.6 ± 1.0	85.3 ± 0.6	92.6 ± 0.1
$A_{FC}+BCE+IC^*$	99.6 ± 0.0	94.7 ± 0.4	92.6 ± 0.5	90.8 ± 2.0	94.4 ± 0.7