

# Automated Model Evaluation for Object Detection via Prediction Consistency and Reliability

## Supplementary Material

### A. Experimental Details

#### A.1. Detectors

We evaluate both one-stage and two-stage detection frameworks employing diverse backbones and detection heads, including RetinaNet+R50 [17, 31], RetinaNet+Swin [31, 32], Faster R-CNN+R50 [17, 39], and Faster R-CNN+Swin [32, 39]. The backbones are initialized with pre-trained classification weights, and the models are trained on a  $3\times$  schedule (36 epochs). Model training was performed using the MMDetection framework [4], with training conducted via the stochastic gradient descent (SGD) optimizer using an initial learning rate of  $10^{-3}$ , weight decay of  $10^{-4}$ , and momentum of 0.9 [45]. Synchronized SGD was employed over 4 GPUs with a total batch size of 8 images (2 images per GPU). For training, the input images are resized such that the shorter side is at most 800 pixels and the longer side does not exceed 1333 pixels, with horizontal image flipping as the only data augmentation technique. During testing, the input images are resized to a fixed dimension of  $800 \times 1333$ .

#### A.2. Computational Resources

For all experiments, an Intel(R) Xeon(R) Gold 6226R CPU, clocked at 2.90 GHz, was used. Regarding GPUs, detector training was executed on four NVIDIA RTX A5000 GPUs, whereas testing (including AutoEval experiments) was carried out on a single NVIDIA RTX A5000 GPU with PyTorch 1.13.1 and CUDA 11.6.

#### A.3. Replication of BoS Results

All BoS [45] results presented in our paper are based on our replication, with modifications made to their publicly released code. We found that the experimental settings described in the paper slightly differ from those implemented in the code, and that certain parts of the experiments are not directly reproducible. Therefore, we modified the code to match the reported performance as closely as possible.

### B. Additional Results

#### B.1. When Test Set is Extremely Difficult

Table B.1 shows the performance on extremely challenging test sets, where mAP is as low as 1–2%. To construct such a test set, we selected the bottom 10% of images (in terms of mAP) from the original test set consisting of 250 images. We evaluate the trained linear regression model not only on

the real-world dataset but also on its bottom 10% subset as a test set. This experiment demonstrates that our PCR method enables AutoEval to perform reliably, outperforming other methods even under such extreme conditions.

Meta-dataset	Corruption-based (Ours)				Avg.
Model	RetinaNet		Faster R-CNN		RMSE
	ResNet-50	Swin-T	ResNet-50	Swin-T	
BoS [45]	15.17	13.55	11.24	13.20	13.29
PCR (Ours)	<b>9.87</b>	<b>10.65</b>	<b>10.88</b>	<b>10.00</b>	<b>10.35</b>

Table B.1. Performance on the real-world dataset and its bottom 10% subset as a test set in vehicle detection with our corruption-based meta-dataset.

#### B.2. Impact of Set Size

Figure B.1 reports RMSE under varying sizes of the meta-dataset, sample set, and test set. Increasing the meta-dataset and sample set sizes provides richer information, generally resulting in lower RMSE. For smaller test set size, it remains robust with  $> 150$  samples.

#### B.3. Hyperparameter Tuning

**Confidence Threshold.** Figure B.2 (a) analyzes how changing the confidence threshold impacts mAP estimation accuracy. The results presented correspond to the mAP estimation for vehicle detection and pedestrian detection. Based on the confidence threshold, the evaluation determines for each final prediction box whether to assess the consistency between that box and its corresponding pre-NMS predictions or to focus on reliability.

**Reliability Measurement.** Figure B.2 (b) shows the RMSE for different reliability measurement methods. Reliability measures the trustworthiness of classification and localization by considering all pre-NMS predictions; therefore, relying solely on predictions with high confidence can distort this metric when numerous predictions with low confidence are omitted. Additionally, a standard sigmoid function forces low confidence to produce near-zero values, effectively rendering them negligible in summation (in the case of  $\alpha = 0$ ). To address these issues, our proposed method incorporates all pre-NMS predictions using a function that preserves the influence of those with low confidence, resulting in better performance than the method with  $\alpha = 0$ . Further, we apply a lower bound of 0.2 to the sigmoid function to prevent the impact of low-confidence score from being lost during the summation.

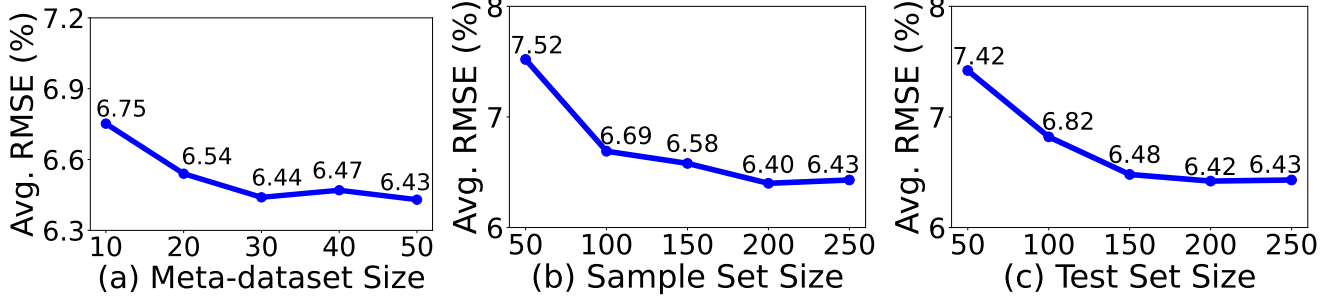


Figure B.1. Effect of set size on vehicle detection with our corruption-based meta-dataset: (a) meta-dataset, (b) sample set, and (c) test set.

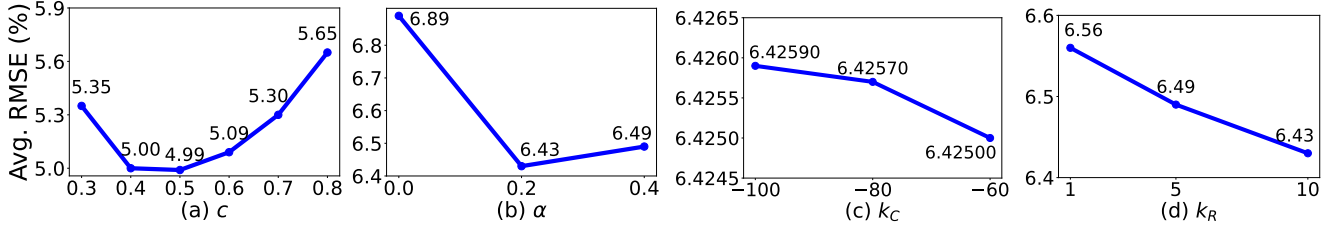


Figure B.2. (a) Effect of confidence threshold  $c$  in Eq. (4) and Eq. (5) (b) Effect of  $\alpha$  in  $\sigma_R$ . (c–d) Effect of slope parameters  $k_C$  and  $k_R$  in  $\sigma_C$  and  $\sigma_R$ , respectively.

**Slope parameters.** We identify the optimal values of  $k_C$  and  $k_R$ , the slope parameters for the sigmoid functions used in the definitions of consistency and reliability scores, respectively. Figure B.2 (middle and right) shows the experimental results for consistency and reliability, respectively. The role of this function varies between the two scores. For consistency, a large negative scaling factor of  $k$  is used to effectively filter out final predictions with low confidence. In contrast, reliability employs a smaller positive scaling factor because it considers all confidence scores of pre-NMS predictions rather than focusing only on high-confidence ones.

#### B.4. Multi-Class Detection

We conducted experiments not only on single-class detection but also on multi-class detection. The detector is trained on a COCO training set. Also, we evaluate our method on five datasets: Cityscapes, COCO, ExDark, KITTI, and Self-driving. Multi-class detection requires considering the relationships between classes, which necessitates more detailed evaluation metrics for AutoEval. PCR demonstrates stable performance not only for the car and person classes but also for multi-class detection, exceeding BoS [45], the top-performing existing method.

#### B.5. Vehicle Detection

Tables B.3 to B.10 present the mAP estimation results for vehicle detectors, categorized by the dataset used for testing the mAP estimation model and the detector architecture.

Meta-dataset	Corruption-based (Ours)				Avg. RMSE
Model	RetinaNet		Faster R-CNN		
	ResNet-50	Swin-T	ResNet-50	Swin-T	
BoS [45]	10.04	<b>5.07</b>	9.38	<b>7.34</b>	7.96
PCR (Ours)	<b>7.09</b>	7.06	<b>8.35</b>	7.80	<b>7.58</b>

Table B.2. Performance for multi-class detection with corruption-based meta-dataset (Ours).

ture. Tables B.3 to B.6 present the results of a linear regression model trained on meta-datasets generated using the transformations from BoS [45]. Our proposed PCR method, when evaluated under this setting, achieves performance estimation comparable to existing methods. And Tables B.7 to B.10 display results on meta-datasets constructed using the ImageNet-C [19] transformations. Specifically, when considering the meta-dataset derived from the ImageNet-C transformation, our approach achieved the lowest RMSE—averaged over all datasets—across every evaluated detector architecture. This observation suggests that our method generally exhibits robust estimation performance, but it performs particularly well under realistic transformation conditions, further emphasizing its potential effectiveness in a wide range of practical applications.

#### B.6. Pedestrian Detection

Tables B.11 to B.18 report the mAP estimation results for pedestrian detectors, categorized by the dataset used for

Method	COCO	BDD	Cityscapes	DETRAC	ExDark	Kitti	Self-driving	Roboflow	Udacity	Traffic	Avg. RMSE
	34.6	32.8	40.5	40.4	28.4	42.4	32.0	33.2	29.4	26.2	
PS [20]	<u>1.70</u>	7.67	3.84	12.23	6.47	11.60	<u>2.58</u>	<u>2.18</u>	<b>0.08</b>	<u>0.19</u>	4.85
ES [40]	6.39	3.83	15.28	4.25	<b>0.77</b>	18.85	4.85	5.30	2.02	3.59	6.51
AC [15]	6.88	5.04	10.55	44.09	2.89	15.37	4.23	2.31	2.02	1.57	9.50
ATC [10]	4.53	7.19	11.48	<b>0.84</b>	4.67	15.79	6.06	<b>0.86</b>	3.71	<b>0.05</b>	5.52
BoS [45]	<b>1.20</b>	<b>0.47</b>	<b>2.06</b>	4.53	4.57	<u>7.71</u>	2.80	2.85	4.06	0.82	<u>3.11</u>
PCR	6.07	<u>1.93</u>	<u>2.85</u>	<u>1.39</u>	<u>2.08</u>	<b>7.28</b>	<b>1.47</b>	3.00	<u>1.99</u>	2.08	<b>3.01</b>

Table B.3. Comparison of AutoEval methods for vehicle detection with augmentation-based meta-dataset from [45]. RetinaNet with ResNet-50 trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO	BDD	Cityscapes	DETRAC	ExDark	Kitti	Self-driving	Roboflow	Udacity	Traffic	Avg. RMSE
	35.8	33.5	40.2	39.9	25.8	42.1	30.6	27.3	28.7	25.5	
PS [20]	9.29	4.84	5.93	22.51	2.61	17.66	5.05	<u>0.42</u>	2.16	3.38	7.38
ES [40]	10.85	<b>1.21</b>	14.04	<b>4.25</b>	<b>0.91</b>	16.91	<u>1.59</u>	1.19	1.60	<u>3.27</u>	<u>5.58</u>
AC [15]	9.55	<u>3.49</u>	11.59	31.93	<u>1.84</u>	16.28	2.27	<u>0.42</u>	<b>0.87</b>	3.49	8.17
ATC [10]	<u>8.85</u>	4.70	11.20	62.18	2.22	15.54	3.06	<b>0.06</b>	<u>1.04</u>	<b>3.13</b>	11.20
BoS [45]	<b>2.78</b>	5.57	<b>2.45</b>	<u>8.88</u>	10.29	<b>4.23</b>	10.30	13.87	10.59	7.90	7.69
PCR	9.50	3.97	<u>5.70</u>	10.50	2.87	<u>7.80</u>	<b>1.20</b>	8.34	1.34	3.47	<b>5.47</b>

Table B.4. Comparison of AutoEval methods for vehicle detection with augmentation-based meta-dataset from [45]. RetinaNet with Swin Transformer trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

testing the mAP estimation model and the detector architecture. Tables B.11 to B.14 resent the results of a linear regression model trained on meta-datasets generated using the transformations from BoS [45]. Our proposed PCR appears to exhibit outstanding performance estimation relative to existing methods in most cases. And Tables B.15 to B.18 display results on meta-datasets constructed using the ImageNet-C [19] transformations. In this setting, our proposed PCR shows superior performance estimation relative to existing methods. Consistent with the vehicle detection case, it recorded the lowest RMSE across all detector architectures.

## B.7. Correlation Results

In addition to RMSE, correlation is also an important evaluation metric. The related experimental results are shown in Figure B.3 and Figure B.4, based on meta-dataset built using the ImageNet-C [19] transformations. The correlation is computed between the mAP of each sample set in the meta-dataset and its corresponding AutoEval score. PCR shows a higher correlation than BoS [45] indicating that the applied transformations effectively differentiate characteristics of different sample sets. Therefore, PCR not only achieves the best performance on the test sets but also shows

superior consistency across the training meta-dataset.

Method	COCO 36.0	BDD 32.2	Cityscapes 40.5	DETRAC 36.2	ExDark 25.9	Kitti 38.9	Self-driving 29.4	Roboflow 29.4	Udacity 27.8	Traffic 27.7	Avg. RMSE
PS [20]	8.85	5.96	8.48	2.36	2.99	4.14	2.63	8.02	5.01	1.10	4.95
ES [40]	13.57	<u>4.28</u>	12.64	<b>0.52</b>	<b>0.26</b>	8.67	<b>0.16</b>	12.44	<u>3.32</u>	1.58	5.74
AC [15]	11.79	5.63	11.15	7.83	<u>0.50</u>	8.53	1.40	6.47	<b>2.06</b>	1.42	5.68
ATC [10]	8.80	7.66	10.60	4.15	2.66	5.63	2.16	3.30	3.64	<b>0.28</b>	4.89
BoS [45]	<b>6.93</b>	7.23	<b>5.01</b>	2.32	4.11	<b>0.66</b>	<u>0.95</u>	<u>2.31</u>	3.50	<u>0.30</u>	<b>3.33</b>
PCR	<u>7.15</u>	<b>3.98</b>	<u>8.36</u>	<u>0.90</u>	3.33	6.40	3.41	<b>1.00</b>	4.53	0.97	<u>4.00</u>

Table B.5. Comparison of AutoEval methods for vehicle detection with augmentation-based meta-dataset from [45]. Faster R-CNN with ResNet-50 trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 38.7	BDD 34.4	Cityscapes 42.6	DETRAC 36.4	ExDark 31.6	Kitti 44.4	Self-driving 31.4	Roboflow 31.9	Udacity 29.2	Traffic 27.3	Avg. RMSE
PS [20]	7.16	4.39	6.06	18.93	2.15	7.42	2.76	4.64	6.95	<u>2.33</u>	6.28
ES [40]	12.75	7.17	9.34	<u>7.27</u>	<u>1.32</u>	8.97	<b>0.39</b>	7.58	7.90	<b>1.88</b>	6.46
AC [15]	8.34	3.79	8.56	19.68	2.68	11.26	<u>2.45</u>	<u>4.40</u>	<b>4.49</b>	4.05	6.97
ATC [10]	7.03	3.63	7.62	23.44	1.61	9.86	3.30	<b>3.18</b>	5.10	3.78	6.86
BoS [45]	<b>0.42</b>	<b>1.04</b>	<u>5.55</u>	<b>3.65</b>	2.32	<u>7.40</u>	6.66	13.30	9.19	6.19	<u>5.57</u>
PCR	<u>6.78</u>	<u>1.68</u>	<b>4.56</b>	14.31	<b>0.27</b>	<b>7.05</b>	2.67	5.24	<u>5.03</u>	3.17	<b>5.08</b>

Table B.6. Comparison of AutoEval methods for vehicle detection with augmentation-based meta-dataset from [45]. Faster R-CNN with Swin Transformer trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 34.6	BDD 32.8	Cityscapes 40.5	DETRAC 40.4	ExDark 28.4	Kitti 42.4	Self-driving 32.0	Roboflow 33.2	Udacity 29.4	Traffic 26.2	Avg. RMSE
PS [20]	<u>11.71</u>	14.90	<u>14.48</u>	<u>9.31</u>	<b>3.97</b>	<u>21.11</u>	10.68	11.59	<u>8.03</u>	7.18	<u>11.30</u>
ES [40]	14.23	12.59	21.40	31.70	7.24	23.48	11.13	12.83	8.74	5.39	14.87
AC [15]	15.11	13.45	20.87	20.82	7.22	24.02	11.88	12.72	9.78	6.44	14.23
ATC [10]	14.05	14.59	20.65	18.41	6.51	24.19	12.75	<u>11.56</u>	10.78	7.46	14.10
BoS [45]	13.77	<u>12.16</u>	21.94	20.59	7.01	23.06	<u>10.49</u>	12.46	8.16	<u>5.32</u>	13.50
PCR	<b>10.37</b>	<b>6.18</b>	<b>8.14</b>	<b>6.30</b>	<u>4.30</u>	<b>13.20</b>	<b>4.98</b>	<b>7.96</b>	<b>1.48</b>	<b>1.34</b>	<b>6.43</b>

Table B.7. Comparison of AutoEval methods for vehicle detection with corruption-based meta-dataset (Ours). RetinaNet with ResNet-50 trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 35.8	BDD 33.5	Cityscapes 40.2	DETRAC 39.9	ExDark 25.8	Kitti 42.1	Self-driving 30.6	Roboflow 27.3	Udacity 28.7	Traffic 25.5	Avg. RMSE
PS [20]	12.79	7.87	<u>4.23</u>	38.03	<b>1.11</b>	21.21	8.71	<b>4.32</b>	5.96	<u>0.12</u>	10.44
ES [40]	13.94	<b>2.91</b>	16.53	<b>1.85</b>	3.54	19.70	3.52	<u>4.51</u>	4.57	0.93	7.20
AC [15]	13.68	6.21	13.87	34.55	2.92	20.12	5.49	4.92	4.68	<b>0.08</b>	10.65
ATC [10]	12.76	7.92	13.00	53.19	<u>2.42</u>	19.06	6.70	4.58	4.97	0.55	12.52
BoS [45]	<b>6.14</b>	4.02	9.78	<u>8.04</u>	3.17	<u>13.51</u>	<b>0.12</b>	4.92	<b>1.11</b>	0.97	<u>5.18</u>
PCR	<u>9.67</u>	<u>3.52</u>	<b>3.09</b>	8.68	4.62	<b>5.27</b>	<u>0.87</u>	7.42	<u>1.73</u>	3.14	<b>4.80</b>

Table B.8. Comparison of AutoEval methods for vehicle detection with corruption-based meta-dataset (Ours). RetinaNet with Swin Transformer trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 36.0	BDD 32.2	Cityscapes 40.5	DETRAC 36.2	ExDark 25.9	Kitti 38.9	Self-driving 29.4	Roboflow 29.4	Udacity 27.8	Traffic 27.7	Avg. RMSE
PS [20]	16.36	13.26	<u>16.39</u>	<b>9.57</b>	4.89	<u>13.86</u>	<u>5.17</u>	13.81	<u>2.62</u>	8.15	10.41
ES [40]	19.41	<u>12.04</u>	20.07	12.84	7.02	17.44	7.49	15.59	4.79	8.64	12.53
AC [15]	18.37	13.16	19.15	<u>9.61</u>	6.86	17.37	6.57	12.44	5.64	8.49	11.77
ATC [10]	17.14	15.01	20.02	12.20	6.01	17.21	7.23	11.05	5.55	8.34	11.98
BoS [45]	<u>14.64</u>	13.04	16.96	13.47	<u>3.93</u>	14.51	7.63	<u>6.59</u>	4.61	<u>7.81</u>	<u>10.32</u>
PCR	<b>11.65</b>	<b>5.51</b>	<b>13.47</b>	12.83	<b>1.92</b>	<b>13.20</b>	<b>4.49</b>	<b>1.94</b>	<b>1.99</b>	<b>2.00</b>	<b>6.90</b>

Table B.9. Comparison of AutoEval methods for vehicle detection with corruption-based meta-dataset (Ours). Faster R-CNN with ResNet-50 trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 38.7	BDD 34.4	Cityscapes 42.6	DETRAC 36.4	ExDark 31.6	Kitti 44.4	Self-driving 31.4	Roboflow 31.9	Udacity 29.2	Traffic 27.3	Avg. RMSE
PS [20]	<b>9.18</b>	6.82	<b>4.37</b>	15.67	5.04	<b>6.76</b>	1.24	7.60	6.00	0.33	6.30
ES [40]	16.25	11.24	17.05	8.66	6.83	18.81	5.92	9.95	<b>0.42</b>	5.07	10.02
AC [15]	11.05	<u>6.60</u>	<u>7.00</u>	15.33	6.82	12.27	0.68	8.44	2.69	1.43	7.23
ATC [10]	10.08	6.89	8.50	10.29	<u>5.08</u>	11.86	<u>0.50</u>	6.65	<u>2.28</u>	<u>0.72</u>	<u>6.29</u>
BoS [45]	10.63	7.58	15.20	<u>7.97</u>	4.54	17.46	3.12	<b>0.54</b>	0.36	1.08	6.85
PCR	<u>9.52</u>	<b>5.38</b>	7.61	<b>0.77</b>	<b>2.64</b>	<u>10.47</u>	<b>0.28</b>	<u>6.08</u>	2.22	<b>0.15</b>	<b>4.51</b>

Table B.10. Comparison of AutoEval methods for vehicle detection with corruption-based meta-dataset (Ours). Faster R-CNN with Swin Transformer trained on the COCO dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 28.3	Caltech 12.0	citypersons 17.0	cityscapes 21.5	crowdhuman 34.2	ECP 24.1	ExDark 27.0	kitti 15.8	self-driving 17.4	Avg. RMSE
PS [20]	<u>6.61</u>	4.46	<b>0.13</b>	4.21	35.67	7.62	5.72	<u>0.60</u>	2.24	7.47
ES [40]	10.97	<u>2.53</u>	2.58	<b>0.34</b>	11.29	6.04	9.12	0.97	2.36	5.13
AC [15]	7.68	3.34	2.35	<b>0.34</b>	17.50	6.59	5.91	0.61	<u>0.99</u>	5.03
ATC [10]	8.68	3.49	0.81	1.84	28.37	7.36	4.76	<b>0.12</b>	<b>0.46</b>	6.21
BoS [45]	<b>4.29</b>	4.59	1.42	1.14	<b>9.49</b>	<b>4.17</b>	<b>0.39</b>	4.44	3.99	<u>3.77</u>
PCR	6.83	<b>1.27</b>	<u>0.46</u>	1.40	<u>10.52</u>	<u>5.41</u>	<u>3.87</u>	0.76	2.03	<b>3.62</b>

Table B.11. Comparison of AutoEval methods for pedestrian detection with augmentation-based meta-dataset from [45]. RetinaNet with ResNet-50 trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 24.6	Caltech 13.2	citypersons 17.9	cityscapes 21.6	crowdhuman 31.0	ECP 24.8	ExDark 22.0	kitti 16.2	self-driving 17.1	Avg. RMSE
PS [20]	6.15	3.52	0.90	2.63	39.55	5.97	4.80	<b>0.07</b>	2.54	7.35
ES [40]	6.40	<u>0.77</u>	<u>0.80</u>	<u>0.97</u>	<b>5.03</b>	8.45	4.26	1.92	3.39	3.55
AC [15]	5.42	1.76	0.34	1.97	17.46	8.22	3.20	<u>0.08</u>	1.03	4.39
ATC [10]	6.16	3.10	0.46	3.32	41.66	8.49	<u>2.81</u>	0.82	<b>0.30</b>	7.46
BoS [45]	<b>1.49</b>	4.03	5.29	1.98	<u>5.78</u>	<b>2.08</b>	<b>2.65</b>	3.34	1.85	<u>3.17</u>
PCR	<u>4.20</u>	<b>0.36</b>	<b>0.30</b>	<b>0.03</b>	7.92	<u>5.82</u>	4.02	2.62	<u>0.82</u>	<b>2.90</b>

Table B.12. Comparison of AutoEval methods for pedestrian detection with augmentation-based meta-dataset from [45]. RetinaNet with Swin Transformer trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 28.3	Caltech 11.7	citypersons 17.7	cityscapes 22.6	crowdhuman 36.3	ECP 26.1	ExDark 26.0	kitti 18.1	self-driving 15.9	Avg. RMSE
PS [20]	<u>4.44</u>	6.43	<u>1.16</u>	<b>2.08</b>	<u>11.06</u>	8.29	1.24	10.56	10.56	6.20
ES [40]	9.82	8.88	1.45	2.92	20.53	7.93	6.28	2.20	13.42	8.16
AC [15]	7.57	8.79	1.43	<u>2.36</u>	17.45	8.91	2.91	<b>1.22</b>	10.55	6.80
ATC [10]	8.44	6.97	<b>0.91</b>	4.21	18.23	9.76	<u>0.75</u>	6.02	7.51	6.98
BoS [45]	6.80	<u>6.38</u>	1.47	2.59	14.54	<u>7.85</u>	<b>0.62</b>	7.14	<u>5.50</u>	<u>5.88</u>
PCR	<b>4.06</b>	<b>0.08</b>	3.45	2.59	<b>3.30</b>	<b>4.46</b>	2.72	<u>1.33</u>	<b>3.08</b>	<b>2.79</b>

Table B.13. Comparison of AutoEval methods for pedestrian detection with augmentation-based meta-dataset from [45]. Faster R-CNN with ResNet-50 trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 27.6	Caltech 12.9	citypersons 19.7	cityscapes 24.0	crowdhuman 35.2	ECP 27.3	ExDark 27.7	kitti 20.2	self-driving 17.4	Avg. RMSE
PS [20]	3.75	4.97	<b>1.26</b>	3.63	3.73	7.03	<u>1.28</u>	13.60	2.52	4.64
ES [40]	9.76	9.11	1.45	5.04	16.96	8.43	7.53	1.91	<u>1.81</u>	6.89
AC [15]	<u>2.82</u>	5.60	<u>1.34</u>	<u>1.54</u>	9.08	6.98	1.89	<u>1.60</u>	5.80	4.07
ATC [10]	4.08	5.70	1.59	3.99	<u>0.82</u>	7.88	<b>0.90</b>	7.01	2.99	3.88
BoS [45]	<b>2.47</b>	<u>4.49</u>	2.08	<b>1.51</b>	<b>0.75</b>	<b>5.57</b>	1.64	8.71	4.88	<u>3.57</u>
PCR	3.57	<b>1.65</b>	2.34	2.07	3.14	<u>5.76</u>	2.19	<b>0.18</b>	<b>0.52</b>	<b>2.38</b>

Table B.14. Comparison of AutoEval methods for pedestrian detection with augmentation-based meta-dataset from [45]. Faster R-CNN with Swin Transformer trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 28.3	Caltech 12.0	citypersons 17.0	cityscapes 21.5	crowdhuman 34.2	ECP 24.1	ExDark 27.0	kitti 15.8	self-driving 17.4	Avg. RMSE
PS [20]	<u>9.05</u>	0.83	4.32	8.69	43.82	10.96	8.42	6.00	<b>1.27</b>	10.37
ES [40]	14.76	5.59	<u>1.98</u>	<u>3.96</u>	12.29	<u>9.50</u>	12.89	7.34	7.86	8.46
AC [15]	11.90	3.29	2.66	5.57	<u>11.94</u>	11.04	10.20	6.94	4.02	<u>7.51</u>
ATC [10]	12.30	2.39	3.54	6.33	15.24	11.03	<u>7.99</u>	5.61	3.90	7.59
BoS [45]	16.72	<b>0.10</b>	5.10	<b>9.72</b>	25.11	12.32	14.50	<u>3.60</u>	4.94	10.23
PCR	<b>7.13</b>	<u>0.27</u>	<b>0.05</b>	<u>1.19</u>	<b>6.32</b>	<b>7.09</b>	<b>5.27</b>	<b>2.59</b>	<u>2.15</u>	<b>3.56</b>

Table B.15. Comparison of AutoEval methods for pedestrian detection with corruption-based meta-dataset (Ours). RetinaNet with ResNet-50 trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 24.6	Caltech 13.2	citypersons 17.9	cityscapes 21.6	crowdhuman 31.0	ECP 24.8	ExDark 22.0	kitti 16.2	self-driving 17.1	Avg. RMSE
PS [20]	<u>8.23</u>	<u>0.57</u>	<u>1.35</u>	<u>4.52</u>	62.54	<u>7.51</u>	7.44	3.60	<b>0.65</b>	10.71
ES [40]	9.94	4.78	3.15	5.34	12.68	11.52	7.52	5.03	6.65	7.40
AC [15]	9.35	3.62	3.88	5.97	<b>0.97</b>	12.00	7.12	4.74	5.48	<u>5.90</u>
ATC [10]	10.37	2.28	4.96	7.55	<u>4.86</u>	12.55	6.95	4.09	4.96	6.51
BoS [45]	9.25	<b>0.13</b>	2.59	6.27	16.62	9.72	<b>5.51</b>	<u>2.00</u>	<u>2.68</u>	6.08
PCR	<b>4.89</b>	4.27	<b>0.12</b>	<b>0.74</b>	11.61	<b>6.90</b>	<u>5.56</u>	<b>0.55</b>	4.05	<b>4.30</b>

Table B.16. Comparison of AutoEval methods for pedestrian detection with corruption-based meta-dataset (Ours). RetinaNet with Swin Transformer trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.



Method	COCO 28.3	Caltech 11.7	citypersons 17.7	cityscapes 22.6	crowdhuman 36.3	ECP 26.1	ExDark 26.0	kitti 18.1	self-driving 15.9	Avg. RMSE
PS [20]	<u>9.66</u>	1.46	<u>2.74</u>	<u>6.24</u>	<u>15.07</u>	<u>11.96</u>	<b>4.82</b>	4.62	5.00	<u>6.84</u>
ES [40]	16.97	1.63	5.64	10.73	26.25	14.59	13.95	5.82	<u>0.37</u>	10.66
AC [15]	14.47	2.05	4.48	8.60	20.47	14.24	10.65	5.70	1.71	9.15
ATC [10]	14.30	<b>0.14</b>	6.19	9.57	20.29	14.64	7.84	<u>0.76</u>	1.24	8.33
BoS [45]	15.88	<b>0.14</b>	6.86	10.30	24.67	14.46	11.39	4.15	2.94	10.09
PCR	<b>9.56</b>	1.20	<b>1.45</b>	<b>3.73</b>	<b>14.29</b>	<b>9.18</b>	<u>6.14</u>	<b>0.07</b>	<b>0.32</b>	<b>5.10</b>

Table B.17. Comparison of AutoEval methods for pedestrian detection with corruption-based meta-dataset (Ours). Faster R-CNN with ResNet-50 trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.

Method	COCO 27.6	Caltech 12.9	citypersons 19.7	cityscapes 24.0	crowdhuman 35.2	ECP 27.3	ExDark 27.7	kitti 20.2	self-driving 17.4	Avg. RMSE
PS [20]	<u>6.80</u>	0.91	5.08	6.93	<b>6.13</b>	<u>10.11</u>	<b>4.00</b>	11.81	1.56	5.93
ES [40]	14.01	2.05	5.24	10.05	22.81	13.31	13.27	6.33	2.92	10.00
AC [15]	7.92	<u>0.26</u>	<u>3.23</u>	<u>6.08</u>	8.36	11.30	7.01	3.29	<b>0.85</b>	<u>5.37</u>
ATC [10]	8.49	<b>0.02</b>	6.04	8.23	9.24	11.95	5.24	2.63	1.71	5.95
BoS [45]	10.49	0.58	4.25	8.24	17.48	11.86	7.91	<u>1.14</u>	<u>0.91</u>	6.98
PCR	<b>6.13</b>	3.42	<b>1.80</b>	<b>2.93</b>	<u>7.11</u>	<b>8.56</b>	<u>4.21</u>	<b>0.04</b>	2.44	<b>4.07</b>

Table B.18. Comparison of AutoEval methods for pedestrian detection with corruption-based meta-dataset (Ours). Faster R-CNN with Swin Transformer trained on the CrowdHuman dataset is used, and its ground truth mAP (%) is reported for various unseen test sets (as indicated in the header). RMSE is employed to quantify the accuracy of the mAP estimation. The best result in each dataset is highlighted in **bold**, and the second-best is underlined.



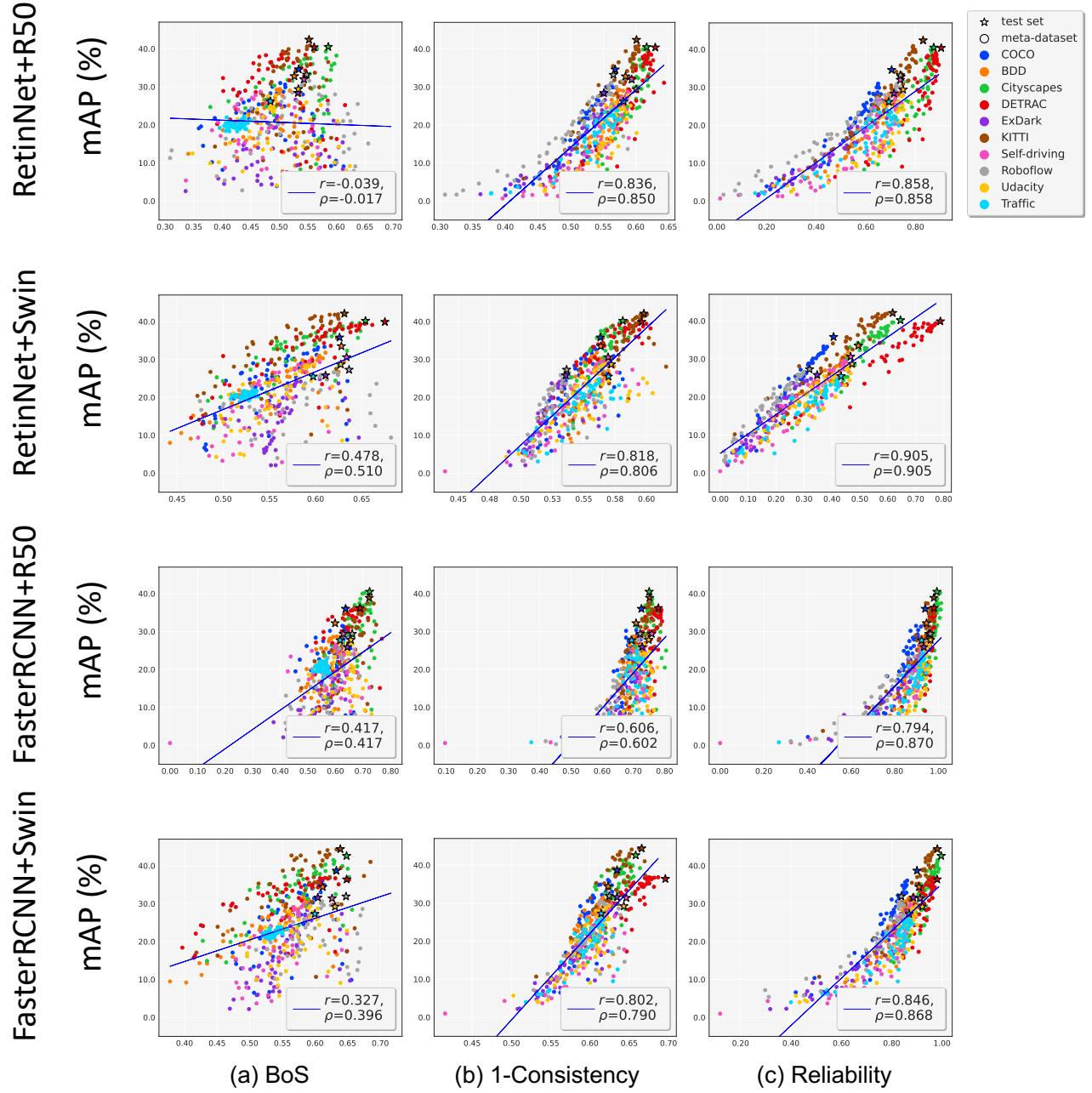
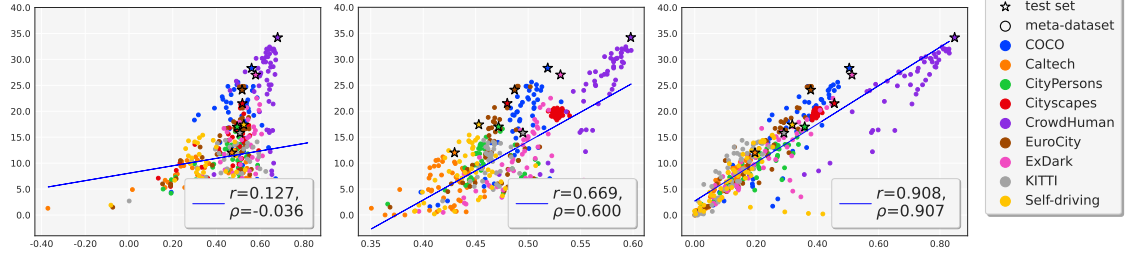


Figure B.3. Correlation results for vehicle detection. Each color represents a different dataset. Markers “★” and “○” indicate the test sets (real-world datasets) and the train sets (meta-dataset), respectively. The values in the legend and the blue line show the correlation coefficients between the mAP of the meta-dataset and the AutoEval score ( $r$ : Pearson’s correlation,  $\rho$ : Spearman’s rank correlation) and the fitted regression line.

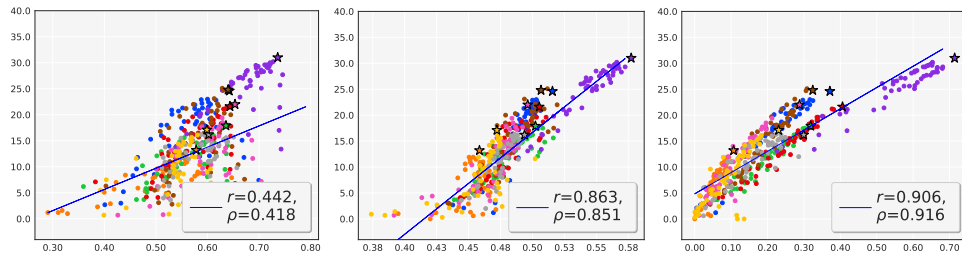
RetinNet+R50

mAP (%)



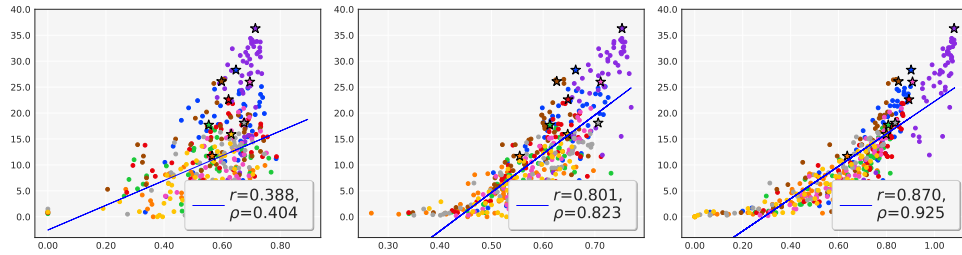
RetinNet+Swin

mAP (%)



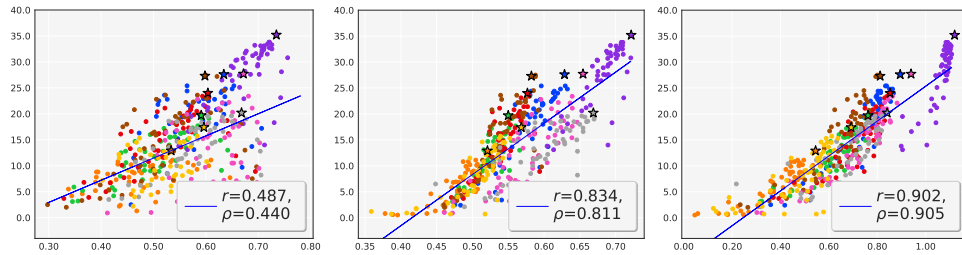
FasterRCNN+R50

mAP (%)



FasterRCNN+Swin

mAP (%)



(a) BoS

(b) 1-Consistency

(c) Reliability

Figure B.4. Correlation results for pedestrian detection. Each color represents a different dataset. Markers “ $\star$ ” and “ $\circ$ ” indicate the test sets (real-world datasets) and the train sets (meta-dataset), respectively. The values in the legend and the blue line show the correlation coefficients between the mAP of the meta-dataset and the AutoEval score ( $r$ : Pearson’s correlation,  $\rho$ : Spearman’s rank correlation) and the fitted regression line.