

# 1 Appendix

## 2 1.1 Proof

3 **Proposition 1.1.** *For any probability functions  $p$  and  $q$  of training distribution and testing distribution,*  
 4 *diversity shift  $D_{\text{diversity}}$  and attribute shift  $D_{\text{attribute}}$  are inclusively bounded between 0 and 1.*

5 *Proof.* Obviously, both  $D_{\text{diversity}}$  and  $D_{\text{attribute}}$  are positive. Then, we prove the upper bound by the  
 6 triangle inequality as followed:

$$D_{\text{diversity}} = \frac{1}{2} \int_{\mathbf{S}} |p(z) - q(z)| dz \leq \frac{1}{2} \int_{\mathbf{S}} [p(z) + q(z)] dz \leq 1 \quad (1)$$

7 Similarly, we have the following inequality:

$$\begin{aligned} D_{\text{attribute}} &= \frac{1}{2} \int_{\mathbf{T}} \sqrt{p(z) \cdot q(z)} \sum_{y_1, y_2} |p(y_1, y_2|z) - q(y_1, y_2|z)| dz \\ &\leq \frac{1}{2} \int_{\mathbf{T}} \sqrt{p(z) \cdot q(z)} \sum_{y_1, y_2} [p(y_1, y_2|z) + q(y_1, y_2|z)] dz \\ &= \frac{1}{2} \int_{\mathbf{T}} 2\sqrt{p(z) \cdot q(z)} dz = \int_{\mathbf{T}} \sqrt{p(z) \cdot q(z)} dz \leq 1 \end{aligned} \quad (2)$$

8 The second inequality is due to triangle inequality.  $\square$

## 9 1.2 Datasets

10 **BDD100K** The original BDD100K [26] contains 80,000 labeled images (70,000 for training and  
 11 10,000 for validation) with ten annotated object categories, including bike, bus, car, motor, person,  
 12 rider, traffic light, traffic sign, train and truck. Each image has three attribute labels which indicate  
 13 the condition, including the weather, scene and time for data collection. We remove the images with  
 14 an undefined attribute label and separate the rest into three OOD environments.

15 **Sim10K** Sim10k [10] is a synthetic dataset containing 10,000 images (8,000 for training, 1,000 for  
 16 validating and 1,000 for testing) with bounding box annotations for cars, which is rendered with the  
 17 Grand Theft Auto V (GTA5) game engine. On DetectBench, we use the training and validating data  
 18 to construct Sim2real benchmark.

19 **Cityscapes** Cityscapes [5] is a large-scale database which focuses on urban street scenes. The dataset  
 20 consists of around 5000 fine annotated images (2975 for training, 500 for validating and the rest for  
 21 testing) with eight annotated instance categories. On DetectBench, we consider the car recognition  
 22 task to construct Sim2real benchmark for simplicity and without loss of generality.

23 **CtrlShift** CtrlShift is a synthetic dataset to analyze the two-dimension shift on OOD object detection.  
 24 It provides an API to generate the training set and the testing set with specific choices of  $\rho_{\text{diversity}}$   
 25 and  $\rho_{\text{attribute}}$ .

## 26 1.3 Implementation Details

27 To evaluate the object detection algorithms, we use the models and the pre-trained weights provided  
 28 by mmdetection [4].

29 For domain generalization algorithms on OOD object detection, we derive the implementations using  
 30 Faster R-CNN [19] with ResNet-50 FPN backbone [8] from torchvision. The whole network is  
 31 optimized by Stochastic Gradient Descent with learning rate 0.02, momentum 0.9 and weight decay  
 32 0.0005.

## 33 1.4 Further Results

34 **Task complexity.** To analyse the IID condition on CtrlShift, which indicates both  $D_{\text{attribute}}$  and  
 35  $D_{\text{diversity}}$  equal zero, we propose a hyper-parameter task complexity  $\alpha$  to measure the difficulty of the

Table 1: Illustration of the two setting on Sim2real.  $\text{sim}_{\text{train}}$  and  $\text{sim}_{\text{val}}$  indicate the training set and validating set from original Sim10K [10] while  $\text{city}_{\text{train}}$  and  $\text{city}_{\text{val}}$  are from original Cityscapes [5]. **Quantity** indicates the number of image. **Total** counts the total number of training and testing domains respectively.

Setting	Split	Train	Test	Quantity	Total
part-sim-part-real	$\text{sim}_{\text{train}}$	✓		8000	8500
	$\text{sim}_{\text{val}}$		✓	1000	
	$\text{city}_{\text{train}}$		✓	2975	3975
	$\text{city}_{\text{val}}$	✓		500	
all-sim-all-real	$\text{sim}_{\text{train}}$	✓		8000	9000
	$\text{sim}_{\text{val}}$	✓		1000	
	$\text{city}_{\text{train}}$		✓	2975	3475
	$\text{city}_{\text{val}}$		✓	500	

Table 2: The experimental results of object detection algorithms on the all-sim-all-real of Sim2real. Mem (GB)<sup>†</sup> and Inf time (fps)<sup>†</sup> are from mmdetection [4].

Detector	Backbone	Mem <sup>†</sup>	fps <sup>†</sup>	AP
Faster R-CNN [19]	X-101	10.3	9.4	35.6
RetinaNet [15]	X-101	10.0	8.7	38.0
Mask R-CNN [7]	X-101	10.7	8.0	36.7
CornetNet [13]	Hourglass104	13.9	4.2	21.6
YOLOv3 [18]	DarkNet-53	7.4	48.1	28.2
FCOS [23]	X-101	10.0	9.7	37.9
Cascade R-CNN [3]	X-101	10.7	-	40.5
MS R-CNN [9]	R-X101	11.0	8.0	35.7
Libra R-CNN [16]	X-101	10.8	8.5	35.3
DH R-CNN [25]	R-50	6.8	9.5	33.8
VarifocalNet [27]	X-101	-	-	42.3
Sparse R-CNN [22]	R-101	-	-	40.3
Deformable [28]	R-50	-	-	37.4
YOLOX [6]	YOLOX-x	28.1	-	36.4

task. The difficulty is adjusted by using  $1 - \alpha$  percent novel data in the testing set in addition to the original training data. The experimental results are shown in Figure 1. The generalization ability of each algorithm drops with the increase of task complexity.

**Sim2real benchmark.** The training set of the Sim2real results reported in the main manuscript comprises of the training data from Sim10K [10] and the validating data from Cityscapes [5], while the testing set comprises of the training data from Cityscapes [5] and the validating data from Sim10K [10] (noted by part-sim-part-real, more details can be found in Table 1.4). We reported the experimental results on all-sim-all-real in Table 2 and Table 3.

Table 3: The experimental results of domain generalization algorithms on the all-sim-all-real of Sim2real.

Algorithm	hyper-parameters	AP
ERM [24]	-	32.8
IB-ERM [1]	$\lambda_{ib} = 100$	18.3
IRM [2]	$\lambda_{irm} = 1$	32.7
MMD [14]	$\gamma_{mmd} = 1$	33.2
CORAL [21]	$\gamma_{mmd} = 1$	32.5
VREx [12]	$\lambda_{vrex} = 1$	32.4
GS [17]	$\lambda_{reg} = 0.1$	31.4
IGA [11]	$\lambda_{penalty} = 1000$	33.4
GroupDRO [20]	$\eta_{groupdro} = 0.01$	31.9

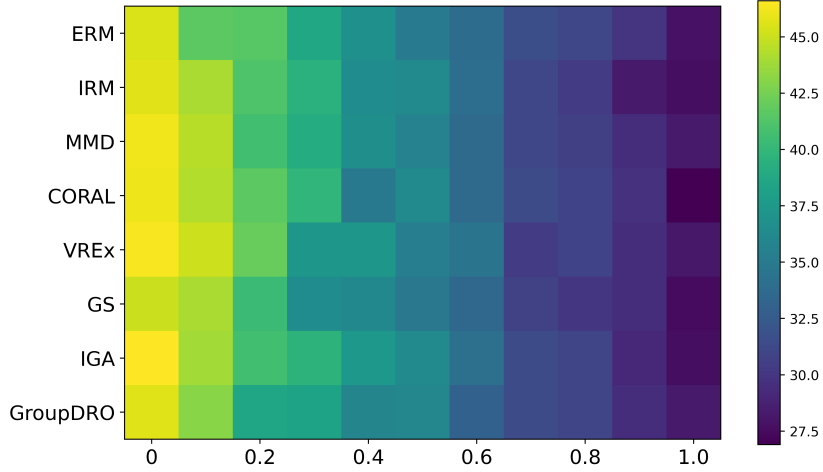


Figure 1: X-axis is task complexity  $\alpha$ . Each block indicates the AP(%).

## References

- [1] Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-distribution generalization. *NeurIPS*, 34, 2021.
- [2] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv:1907.02893*, 2019.
- [3] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: High quality object detection and instance segmentation. *TPAMI*, page 1–1, 2019.
- [4] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv:1906.07155*, 2019.
- [5] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, pages 3213–3223, 2016.
- [6] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. YOLOX: Exceeding yolo series in 2021. *arXiv:2107.08430*, 2021.
- [7] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. *ICCV*, Oct 2017.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

- [9] Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. Mask scoring r-cnn. In *CVPR*, 2019.
- [10] Matthew Johnson-Roberson, Charles Barto, Rounak Mehta, Sharath Nittur Sridhar, Karl Rosaen, and Ram Vasudevan. Driving in the matrix: Can virtual worlds replace human-generated annotations for real world tasks? In *ICRA*, 2016.
- [11] Masanori Koyama and Shoichiro Yamaguchi. Out-of-distribution generalization with maximal invariant predictor. 2021.
- [12] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In *ICML*, pages 5815–5826. PMLR, 2021.
- [13] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, pages 765–781. Springer Verlag, 2018.
- [14] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *CVPR*, pages 5400–5409, 2018.
- [15] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017.
- [16] Jiangmiao Pang, Kai Chen, Jianping Shi, Huajun Feng, Wanli Ouyang, and Dahua Lin. Libra r-cnn: Towards balanced learning for object detection. In *CVPR*, 2019.
- [17] Mohammad Pezeshki, Sékou-Oumar Kaba, Yoshua Bengio, Aaron Courville, Doina Precup, and Guillaume Lajoie. Gradient starvation: A learning proclivity in neural networks. *arXiv: Learning*, 2020.
- [18] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv:1804.02767*, 2018.
- [19] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS*, 28, 2015.
- [20] Shiori Sagawa, Pang Wei Koh, Tatsunori Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv: Learning*, 2019.
- [21] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *ECCV*, pages 443–450. Springer, 2016.
- [22] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, and Ping Luo. SparseR-CNN: End-to-end object detection with learnable proposals. *arXiv:2011.12450*, 2020.
- [23] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. *arXiv:1904.01355*, 2019.
- [24] Vladimir Vapnik. *Statistical Learning Theory*. Wiley, 1998.
- [25] Yue Wu, Yinpeng Chen, Lu Yuan, Zicheng Liu, Lijuan Wang, Hongzhi Li, and Yun Fu. Rethinking classification and localization for object detection. *arXiv:1904.06493*, 2019.
- [26] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and Trevor Darrell. Bdd100k: A diverse driving video database with scalable annotation tooling. *arXiv:1805.04687*, 2(5):6, 2018.
- [27] Haoyang Zhang, Ying Wang, Feras Dayoub, and Niko Sünderhauf. Varifocalnet: An iou-aware dense object detector. *arXiv:2008.13367*, 2020.
- [28] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2021.