

CD-FSOD: A BENCHMARK FOR CROSS-DOMAIN FEW-SHOT OBJECT DETECTION

Wuti Xiong

Center for Machine Vision and Signal Analysis, University of Oulu, Finland
wuti.xiong@oulu.fi

ABSTRACT

In this paper, we propose a study of the cross-domain few-shot object detection (CD-FSOD) benchmark, consisting of image data from a diverse data domain. On the proposed benchmark, we evaluate state-of-art FSOD approaches, including meta-learning FSOD approaches and fine-tuning FSOD approaches. The results show that these methods tend to fail, and even underperform the naive fine-tuning model. We analyze the reasons for their failure and introduce a strong baseline that uses a mutually-beneficial manner to alleviate the overfitting problem. Our approach is remarkably superior to existing approaches by significant margins (2.0% on average) on the proposed benchmark. Our code is available at <https://github.com/FSOD/CD-FSOD>.

Index Terms— Few-shot Object Detection, Cross-domain.

1. INTRODUCTION

Few-shot object detection (FSOD) aims to detect novel classes of objects with a few annotated instances. In the previous FSOD setting [1, 2], a detector is pre-training on the source dataset consisting of base classes and then transferred into the target dataset consisting of novel classes with few instances, where base classes and novel classes are disjoint but share similar data domains. However, this underlying assumption does not apply to some real-world scenarios because it is difficult or impossible to collect a sufficient amount of data in these domains. This leads to a new FSOD problem, where the detector must resort to pre-training in the base classes from a different domain. In these cases, even humans have trouble recognizing new categories that vary too greatly between examples or differ from prior experience [3, 4]. Thus, finding new approaches to tackle the problem remains a challenging but desirable goal.

Although conventional FSOD benchmarks [1, 2] are well established, no works study FSOD across different domains. To fill this gap, In this paper, we introduce the study of *Cross-Domain Few-Shot Object Detection* (CD-FSOD) benchmark (As shown in Figure 1), which covers three target datasets: ArTaxOr [6], UODD [7] and DIOR [8]. On the proposed benchmark, we conduct extensive experiments to evaluate existing FSOD approaches (including meta-learning approaches

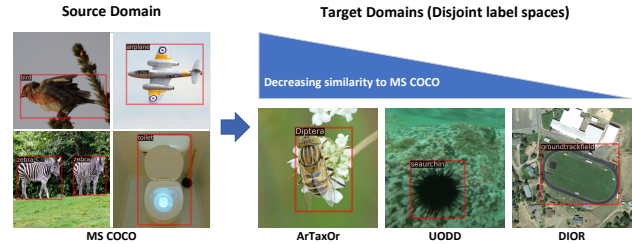


Fig. 1: The CD-FSOD benchmark. MS COCO [5] is used for source training, and domains of varying dissimilarity from natural images are used for target evaluation.

[2, 9, 10] and fine-tuning approaches [11, 12, 13]). The results show that existing FSOD approaches can not achieve satisfactory performance and even underperform the naive fine-tuning model due to freezing parameters. Even without freezing parameters, fine-tuning methods struggle to outperform the naive transfer model while meta-learning methods still fail. This finding shows that existing FSOD methods cannot work for CD-FSOD, and there is an urgent need to develop new methods.

Besides, we introduce a novel distillation-based baseline, which enable a “flywheel effect” that the student and teacher can mutually reinforce each other so that both get better and better as the training goes on. Specifically, EMA (Exponential Moving Average) enables the teacher model to ensemble the student models in different time steps. The student’s weights are optimized by the distillation loss between the pseudo-labels generated by the teacher and the predictions by the students on the same image. Our approach outperforms existing FSOD approaches by a large margin on the proposed benchmark. In summary, our main contributions are as follows: (1) we established the CD-FSOD benchmark, where there is a very large domain difference between the base and target datasets; (2) on the proposed benchmark, we evaluate existing FSOD approaches, and analyze the reasons for their failure; (3) we introduce a strong baseline that achieves state-of-the-art performance on the proposed benchmark.

2. PROPOSED BENCHMARK.

Rethinking FSOD Benchmarks. In previous FSOD work [1, 2, 10, 9, 11, 12, 13], two benchmarks have been widely adopted: MS COCO [5] and PASCAL VOC [14]. As for PASCAL VOC, there are three random split groups, and each of them covers 20 categories, which are randomly divided into 15 base classes and 5 novel classes. Each novel category has $K = 1, 2, 3, 5, 10$ objects sampled from the combination of VOC07 and VOC12 train/val set for few-shot detection training. As for MS COCO, the 60 categories disjoint with VOC are denoted as base classes while the remaining 20 classes are used as novel classes with $K = 1, 2, 3, 5, 10, 30$ shots. 5k images from the validation set are used for evaluation and the rest are used for training. While these benchmarks contributed to the research progress in FSOD, they have a limitation. As we discussed in Section 1, these benchmarks sample base classes and novel classes from a single dataset. There is a major issue that occurs commonly in practice: by the nature of the problem, collecting data from the same domain for many FSOD tasks is difficult. Under these circumstances, useful knowledge may still be effectively transferred across different domains, implying that approaches designed in the FSOD setting may not continue to perform well when applied to different domains, such as biological natural images and satellite images. Currently, no works study this scenario.

CD-FSOD Benchmark. To explore FSOD across a wide range of domains, we propose to build the benchmark using datasets from a wide range of domains rather than just a subset of natural image datasets. Our proposed benchmark include a base dataset (MS COCO [5]) and 3 target datasets from diverse domains: ArTaxOr [6], UODD [7] and DIOR [8]. The selected datasets reflect well-curated real-world use cases for few-shot object detection. In addition, collecting enough examples from the above domains is often difficult, expensive, or in some cases not possible. The similarity of these datasets to the MS COCO dataset, from high to low, is as follows: 1) ArTaxOr images are natural but are fine-grained (specific to biology); 2) UODD images are less similar as the poor visibility and low color contrast, but are still color images of natural scenes; 3) DIOR images are the most dissimilar as they have lost perspective distortion. The statistics of the target dataset are shown in Table 1. Similar to the previous FSOD setting [1, 2], the model is trained from a base dataset where each class has abundantly annotated instances, then is adapted to the target dataset where each class only has K ($K = 1, 5, 10$) instances.

Domain	Dataset	Classs	Train images	Test images
Biology	ArTaxOr	7	13,991	1,383
Underwater	UODD	3	3,194	506
Aerial	DIOR	20	18,463	5,000

Table 1: The statistics of the target datasets.

3. PROPOSED METHOD

Overview. As shown in Figure 2 (a), our approach consists of two stages: the training stage and the testing stage. At the training stage, we simply train the detector using the base data. At the beginning of the fine-tuning, we duplicate the initialized detector into the student model. The student first goes through a burn-in step, i.e. training the student with the standard detection supervision losses [15] on K -shot target instances. Then its weights are copied into the student and the teacher to initiate the distillation step. As shown in Figure 2 (b), in the distillation step, the teacher and student are trained in a mutually-beneficial manner, where the teacher promotes the student by the distillation loss, and its weights are updated by the student model via exponential mean average (EMA). The proposed method consists of two branches: the supervised branch and the distillation branch. The final loss L is the sum of supervised loss L_S and distillation loss L_D .

$$L = L_S + \lambda L_D \quad (1)$$

where the λ is a hyper-parameter. As pointed out by prior works [16], a key factor in improving the teacher is the diversity of student; thus, we use strongly augmented images as input for the student, but we use weakly augmented images as input to the teacher to provide reliable pseudo-labels.

Supervised branch. In the supervised branch, we compute the supervised detection losses (classification loss L_{cls} and localization loss L_{loc}) for the student model. With the K -shot target data $D_s = \{x_i^s, y_i^s\}_{i=1}^{N_s}$, the supervised detection loss L_S is written as:

$$L_S = \sum_i L_{cls}(x_i^s, y_i^s) + L_{loc}(x_i^s, y_i^s) \quad (2)$$

Distillation branch. As shown in Figure 2 (b), the teacher and student share the same architecture and are initialized with the same weights after the burn-in step. An image is processed independently by both the student and the teacher. The teacher is used to generate thousands of box candidates for the weakly augmented version. After NMS [15]) is performed, only candidates with the foreground score higher than a threshold δ are retained as the pseudo boxes p_i^d . Then the distillation loss is obtained by calculating the detection loss between the student predictions x_i^d and the pseudo-labels.

$$L_D = \sum_i L_{cls}(x_i^d, p_i^d) + L_{loc}(x_i^d, p_i^d) \quad (3)$$

Model update. EMA update has been shown to be successful in many prior works [17, 18, 19]. Thus, we use it to alleviate the overfitting problem in the CD-FSOD setting. Specifically, we detach the student and the teacher. After obtaining the pseudo-labels from the teacher, only the learnable weights of

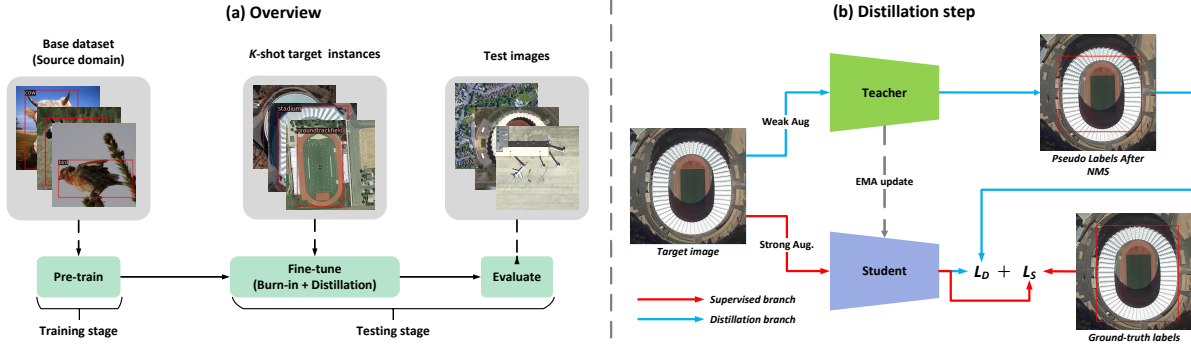


Fig. 2: (a) The overview of our proposed approach. (b) Distillation step.

the student W_s is updated via back-propagation

$$W_s \leftarrow W_s + \gamma \frac{\partial L}{\partial W_s} \quad (4)$$

where the γ denotes the learning rate. Then, the teacher model weights W_t are updated from the student model weights W_s by exponential moving average (EMA) [20]. At each iteration, we update the teacher weights by:

$$W_t \leftarrow \alpha W_t + (1 - \alpha) W_s \quad (5)$$

where the α is a hyper-parameter.

4. EXPERIMENT

4.1. Implementation details

For a fair comparison, we follow previous work [2, 10, 9, 11, 12, 13] to use Faster-RCNN [15] with FPN [21] and ResNet-50 backbone [22] to build the student and teacher. For generating the pseudo boxes, we use confidence threshold $\delta = 0.7$. For the data augmentation, we apply random horizontal flips for weak augmentation and randomly add color jittering, grayscale, Gaussian blur, and cutout patches for strong augmentations. Our implementation builds upon the Detectron2 framework. For the baselines [2, 10, 9, 11, 12, 13], we use the official implementations: A-RPN¹, H-GCN², Meta-RCNN³, TFA⁴, FSCE⁵, DeFRCN⁶. FRCN-ft is a Faster R-CNN [15] detector which is simply trained on the base dataset, then fine-tuned on the K -shot target instances. The teacher is used for the inference and evaluation of test images.

4.2. Main Results

As shown in Table 2, our proposed approach outperforms existing FSOD approaches in all settings. Overall, our ap-

¹<https://github.com/fanql15/FewX>

²<https://github.com/GuangxingHan/QA-FewDet>

³<https://github.com/guangxinghan/meta-faster-r-cnn>

⁴<https://github.com/ucbdrive/few-shot-object-detection>

⁵<https://github.com/megvii-research/FSCE>

⁶<https://github.com/er-muyue/DeFRCN>

proach produces an average 2.0% improvement over the second-best approach on the three datasets. We further observe that all approaches obtain performance gains without freezing parameters. This confirms that freezing some parameters [23, 24, 25] can not alleviate the overfitting problem in the CD-FSOD. Moreover, these approaches still do not show satisfactory performance. The meta-learning approaches still fail to outperform naive fine-tuned models in all settings. This suggests that meta-learning approaches use supervision for pre-training and cannot mimic distant domain datasets, which leads them to overfit the source data and generalize poorly to distant target domains. The fine-tuning approaches DeFRCN and FSCE have only a slight performance improvement over FRCN-ft. This suggests that these approaches tailored for FSOD do not work with CD-FSOD. There is a desirable need to develop approaches that work under both FSOD and CD-FSOD.

4.3. Ablation Studies

In this section, we show the ablation experiments on the DIOR dataset.

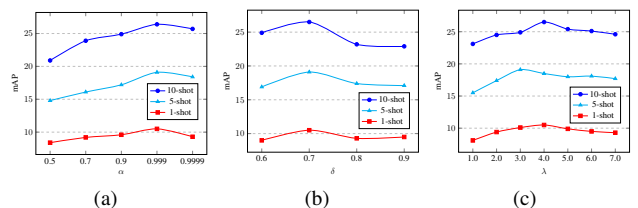


Fig. 3: Ablation studies for (a) EMA rate α , (b) pseudo-labeling threshold δ , and (c) distillation loss weights λ .

EMA and Distillation. As shown in Table 3, the EMA and distillation both improve the performance but EMA achieves better performance than distillation. The combination of distillation and EMA, leads to even better performance. We also observe that the performance of the students increases with the performance of the teacher. This means that our proposed approach enables students and teachers to progress together in a mutually beneficial manner. Extensive research has shown

Method/Shot	ArTaxOr			UODD			DIOR			Avg
	1	5	10	1	5	10	1	5	10	
A-RPN [11]	2.5 _{-1.1}	8.1 _{-3.1}	13.9 _{-4.2}	3.3 _{-1.0}	8.4 _{-2.3}	10.8 _{-1.6}	7.5 _{-1.2}	17.1 _{-2.7}	20.3 _{-2.4}	10.2 _{-2.1}
Meta-RCNN [9]	2.8 _{-0.9}	8.5 _{-2.4}	14.0 _{-3.7}	3.6 _{-0.8}	8.8 _{-2.1}	11.2 _{-1.3}	7.8 _{-2.3}	17.7 _{-2.5}	20.6 _{-1.8}	10.6 _{-1.9}
H-GCN[10]	2.6 _{-0.6}	8.2 _{-1.9}	14.2 _{-3.3}	3.8 _{-0.7}	7.7 _{-1.5}	11.0 _{-1.6}	7.9 _{-1.9}	18.0 _{-2.6}	20.9 _{-2.2}	10.5 _{-1.8}
TFA w/cos [2]	3.1 _{-2.3}	8.8 _{-5.0}	14.8 _{-7.7}	4.4 _{-1.7}	8.7 _{-2.2}	11.8 _{-4.6}	8.0 _{-4.1}	18.1 _{-7.8}	20.5 _{-7.1}	10.9 _{-4.7}
FSCE [12]	3.7 _{-1.9}	10.2 _{-4.3}	15.9 _{-5.1}	3.9 _{-1.1}	9.6 _{-2.9}	12.0 _{-3.6}	8.6 _{-3.0}	18.7 _{-3.8}	21.9 _{-3.6}	11.6 _{-3.2}
DeFRCN [13]	3.6 _{-0.7}	9.9 _{-1.1}	15.5 _{-1.0}	4.5 _{-0.8}	9.9 _{-1.0}	12.1 _{-1.4}	9.3 _{-1.3}	18.9 _{-1.2}	22.9 _{-2.2}	11.8 _{-1.2}
FRCN-ft	3.4	9.3	15.2	4.1	9.2	12.3	8.4	18.3	21.2	11.2
Ours	5.1	12.5	18.1	5.9	12.2	14.5	10.5	19.1	26.5	13.8

Table 2: The performance (mAP) on the CD-FSOD benchmark. The best results are in bold. Red numbers indicate performance degradation due to frozen parameters.

EMA	distillation	1		5		10	
		S	T	S	T	S	T
✓		8.4	9.2	18.3	18.7	21.2	24.8
	✓	8.6	8.9	18.4	18.5	22.3	23.1
✓	✓	9.5	10.5	18.8	19.1	23.6	26.5

Table 3: The effect of EMA and the distillation. “S” and “T” represent the student and the teacher respectively.

that there is an overfitting problem in FSOD. We argue that EMA and distillation can effectively alleviate the problem in FSOD. With EMA, the teacher can be seen as an average model of the student over different steps, so it is more stable and robust. And the distillation loss can be seen as a regularization, which can improve the generalization of the student.

EMA rate. We evaluate the model using various EMA rates α from 0.5 to 0.9999, and present the mAP results in Figure 3 (a). When the EMA ratio is small (e.g., $\alpha = 0.5$), the student contributes more to the teacher model in each iteration, which leads to an unstable teacher model with a lower mAP. This situation can be stabilized and improved as the EMA ratio α increases. It performs the best mAP when the EMA ratio α reaches 0.999. However, if the EMA rate α keeps increasing, the teacher model performance will degrade because the teacher model mainly derives the next model weights from the previous teacher model weights.

Pseudo-labeling thresholding. Pseudo-labeling thresholding plays an important role in the distillation loss, as it can filter the low-confidence predicted bounding boxes. As shown in Figure 3 (b), if the threshold is too low (e.g. $\alpha = 0.6$), the mAP of the model is low because the model predicts more unreliable bounding boxes. On the other hand, the performance of a model using an excessively high threshold (e.g., $\alpha = 0.9$) degrades as it cannot predict a sufficient number of bounding boxes in its generated pseudo-labels.

Distillation loss weight. To examine the effect of the distillation loss weight, we vary the distillation loss weight λ from 1.0 to 8.0. As shown in Figure 3 (c), with a lower distillation loss weight $\lambda = 1.0$, the model has a lower performance.

On the other hand, we observe that the model performs the best with the loss weight $\lambda = 4.0$ (1-shot and 10-shot) or 3.0 (5-shot).

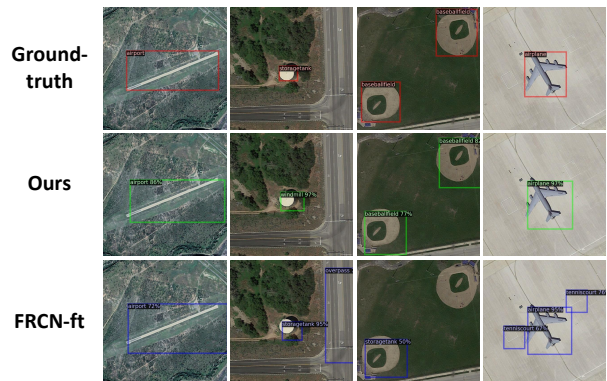


Fig. 4: Detection examples based on 10-shot setting. FRCN-ft leads to three types of false detections: missed detection of certain objects (the third example in the third row), incorrect detection of background (the second example in the third row) and inaccurate localization (the first example in the third row). For these examples, our approach reduces the occurrence of these errors.

5. CONCLUSION

In this paper, we formally introduce the study of the cross-domain few-shot object detection (CD-FSOD) benchmark, which covers several target domains with varying similarities to natural images. On the proposed benchmarks, we evaluate existing FSOD approaches and analyze the reasons for their failure. Then, we introduce a strong baseline that achieves state-of-the-art performance on the proposed benchmark. In the future, we will work on developing novel approaches for both FSOD and CD-FSOD.

Acknowledgement. The authors wish to acknowledge CSC IT Center for Science, Finland, for computational resources.

6. REFERENCES

- [1] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell, "Few-shot object detection via feature reweighting," in *CVPR*, 2019.
- [2] Qi Fan, Wei Zhuo, Chi-Keung Tang, and Yu-Wing Tai, "Few-shot object detection with attention-rpn and multi-relation detector," in *CVPR*, 2020.
- [3] Brenden Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua Tenenbaum, "One shot learning of simple visual concepts," in *Proceedings of the annual meeting of the cognitive science society*, 2011.
- [4] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, 2015.
- [5] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, "Microsoft coco: Common objects in context," in *ECCV*, 2014.
- [6] Geir Drange, "Arthropod taxonomy orders object detection dataset," <https://www.kaggle.com/datasets/mistag/arthropod-taxonomy-orders-object-detection-dataset>, 2019.
- [7] Lihao Jiang, Yi Wang, Qi Jia, Shengwei Xu, Yu Liu, Xin Fan, Haojie Li, Risheng Liu, Xinwei Xue, and Ruili Wang, "Underwater species detection using channel sharpening attention," in *ACM MM*, 2021.
- [8] Ke Li, Gang Wan, Gong Cheng, Liqiu Meng, and Junwei Han, "Object detection in optical remote sensing images: A survey and a new benchmark," *ISPRS J. Photogramm. Remote Sens.*, 2020.
- [9] Guangxing Han, Shiyuan Huang, Jiawei Ma, Yicheng He, and Shih-Fu Chang, "Meta faster r-cnn: Towards accurate few-shot object detection with attentive feature alignment," in *AAAI*, 2022.
- [10] Guangxing Han, Yicheng He, Shiyuan Huang, Jiawei Ma, and Shih-Fu Chang, "Query adaptive few-shot object detection with heterogeneous graph convolutional networks," in *ICCV*, 2021.
- [11] Xin Wang, Thomas Huang, Joseph Gonzalez, Trevor Darrell, and Fisher Yu, "Frustratingly simple few-shot object detection," in *ICML*, 2020.
- [12] Bo Sun, Banghuai Li, Shengcai Cai, Ye Yuan, and Chi Zhang, "FSCE: few-shot object detection via contrastive proposal encoding," in *CVPR*, 2021.
- [13] Limeng Qiao, Yuxuan Zhao, Zhiyuan Li, Xi Qiu, Jianan Wu, and Chi Zhang, "Defrcn: Decoupled faster r-cnn for few-shot object detection," in *CVPR*, 2021.
- [14] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman, "The pascal visual object classes (voc) challenge," *IJCV*, 2010.
- [15] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *TPAMI*, 2016.
- [16] Yen-Cheng Liu, Chih-Yao Ma, Zijian He, Chia-Wen Kuo, Kan Chen, Peizhao Zhang, Bichen Wu, Zsolt Kira, and Peter Vajda, "Unbiased teacher for semi-supervised object detection," in *ICLR*, 2021.
- [17] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *ICLR*, 2015.
- [18] Sergey Ioffe and Christian Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015.
- [19] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*, 2020.
- [20] Antti Tarvainen and Harri Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *NIPS*, 2017.
- [21] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [23] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris, "A broader study of cross-domain few-shot learning," in *ECCV*, 2020.
- [24] Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang, "Cross-domain few-shot classification via learned feature-wise transformation," in *ICLR*, 2020.
- [25] Ashraf Islam, Chun-Fu Richard Chen, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, and Richard Radke, "Dynamic distillation network for cross-domain few-shot recognition with unlabeled data," *NeurIPS*, 2021.