

ReFine: Re-randomization before Fine-tuning for Cross-domain Few-shot Learning

Jaehoon Oh*
KAIST DS
Daejeon, Republic of Korea
jhoon.oh@kaist.ac.kr

Sungnyun Kim*
KAIST AI
Seoul, Republic of Korea
ksn4397@kaist.ac.kr

Namgyu Ho*
KAIST AI
Seoul, Republic of Korea
itsnamgyu@kaist.ac.kr

Jin-Hwa Kim
NAVER AI Lab
Sungnam, Republic of Korea
j1nhwa.kim@navercorp.com

Hwanjun Song[†]
NAVER AI Lab
Sungnam, Republic of Korea
hwanjun.song@navercorp.com

Se-Young Yun[†]
KAIST AI
Seoul, Republic of Korea
yunseyoung@kaist.ac.kr

ABSTRACT

Cross-domain few-shot learning (CD-FSL), where there are few target samples under extreme differences between source and target domains, has recently attracted huge attention. Recent studies on CD-FSL generally focus on transfer learning based approaches, where a neural network is pre-trained on popular labeled source domain datasets and then transferred to target domain data. Although the labeled datasets may provide suitable initial parameters for the target data, the domain difference between the source and target might hinder fine-tuning on the target domain. This paper proposes a simple yet powerful method that re-randomizes the parameters fitted on the source domain before adapting to the target data. The re-randomization resets source-specific parameters of the source pre-trained model and thus facilitates fine-tuning on the target domain, improving few-shot performance.

CCS CONCEPTS

• Computing methodologies → Machine learning.

KEYWORDS

cross-domain, few-shot, transfer learning, re-randomization

ACM Reference Format:

Jaehoon Oh, Sungnyun Kim, Namgyu Ho, Jin-Hwa Kim, Hwanjun Song, and Se-Young Yun. 2022. ReFine: Re-randomization before Fine-tuning for Cross-domain Few-shot Learning. In *Proceedings of the 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*, October 17–21, 2022, Atlanta, GA, USA. ACM, Atlanta, GA, USA, 5 pages. <https://doi.org/10.1145/3511808.3557681>

*The authors contributed equally to this research.

[†]Corresponding authors

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
CIKM '22, October 17–21, 2022, Atlanta, GA, USA

© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9236-5/22/10...\$15.00
<https://doi.org/10.1145/3511808.3557681>

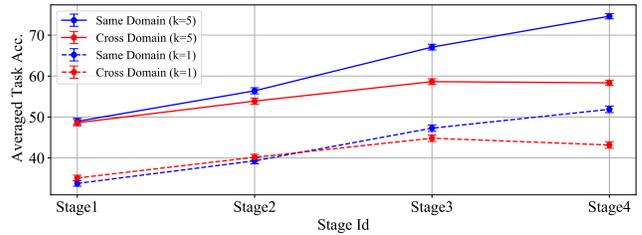


Figure 1: FSL accuracy (5-way k -shot) using the intermediate representation from each stage in ResNet10 (refer to Figure 2 for the ResNet10 structure). An average pooling layer and an auxiliary classifier are attached at the end of each stage. After pre-training on miniImageNet (source domain), the model is transferred to the same domain (blue lines) or four different target domains in the BSCD-FSL benchmark (red lines), where only the attached classifier is fine-tuned. The cross-domain accuracy is averaged on the four domains.

1 INTRODUCTION

Few-shot learning (FSL) has become an attractive field of deep learning research to tackle problems with a small number of training samples [38]. In this setting, a model is typically pre-trained on a large source dataset comprised of *base* classes from the source domain and then transferred into the target dataset comprised of few samples from unseen *novel* classes. Studies on FSL have typically assumed that the base and novel classes share the same domain, and these have followed two research directions: meta-learning [9, 23, 25, 31] and fine-tuning [4, 7, 34].

However, the source dataset and the target dataset come from considerably different domains in many real-world scenarios [13, 24]. To tackle this problem, *cross-domain few-shot learning* (CD-FSL) has recently gained significant attention, exemplified by the introduction of the BSCD-FSL benchmark dataset [13]. This benchmark considers large-scale natural image datasets as source data and four different target datasets for evaluation, each with varying levels of similarity to the source data domain. It is shown that transfer learning approaches, where a pre-trained model on the source domain is fine-tuned on the target domain, overwhelm meta-learning approaches on BSCD-FSL [13].

In this regard, recent works have attempted to extract better representations during the pre-training phase by exploiting unlabeled data from the target domain [16, 22, 24] or reconstructing the images with an autoencoder to enhance the generalization of a model [18]. While these works focus on developing better pre-training methods, we suppose the fine-tuning phase is also a crucial research direction. Das et al. [6] were aware of the importance of fine-tuning for CD-FSL, however, their framework using a mask generator is highly complicated to use.

In this paper, we present a new perspective to tackle the domain gap issue in CD-FSL: *not all the pre-trained parameters from the source domain are desirable on the target domain*. We posit that parameters in deeper layers of a pre-trained feature extractor may be detrimental for target domain adaptation, as they contain domain-specific information belonging to the source domain. This is demonstrated in Figure 1, where we use fixed image features from different stages of a pre-trained backbone and analyze the change in few-shot performance. We observe different trends for same-domain and cross-domain scenarios. While accuracy increases consistently with feature depth in the same-domain case (the blue lines), the accuracy decreases when using features from the last stage in the cross-domain case (the red lines).

Motivated by these findings, we propose a novel method, **ReFine** (**Re**-randomization before **Fine**-tuning), where we re-randomize the top layers of the feature extractor after supervised training on the source domain, before fine-tuning on the target domain. This is effective for CD-FSL because it helps reduce the learning bias towards the source domain by simply re-randomizing the domain-specific layer. It can also be implemented by adding a few lines of code and can be easily combined with other recent CD-FSL methods. This simplicity and flexibility allows it to be easily adapted in practical uses for CD-FSL. Contrary to the prior works that have focused on improving universal representations during the pre-training phase [22, 24], our method focuses on removing source-specific features obtained during pre-training to aid the fine-tuning.

Our contributions are summarized as follows:

- We propose a simple yet effective algorithm called ReFine, which re-randomizes the fitted parameters on the source domain and then fine-tunes the partially re-randomized model. This puts forward a new perspective for adapting to novel classes of the target domain in CD-FSL.
- We demonstrate improved performance for CD-FSL when our re-randomization technique is used, and provide an in-depth analysis on where and how to re-randomize.

2 RELATED WORKS

Few-shot learning (FSL) has been studied in two research directions, meta-learning and fine-tuning. Regarding the meta-learning approach, a meta-trained model is evaluated after fast adaptation on a few train sets. The meta-training procedure resembles the episodic evaluating procedure. Meta-learning approaches include learning good initialized parameters [9, 10, 23, 34], a metric space [3, 31, 32, 36], and update rule or optimization [2, 11, 27]. Regarding the fine-tuning approach [4, 7, 34], a pre-trained model is typically evaluated after fine-tuning. During the pre-training procedure, the model is trained in a mini-batch manner.

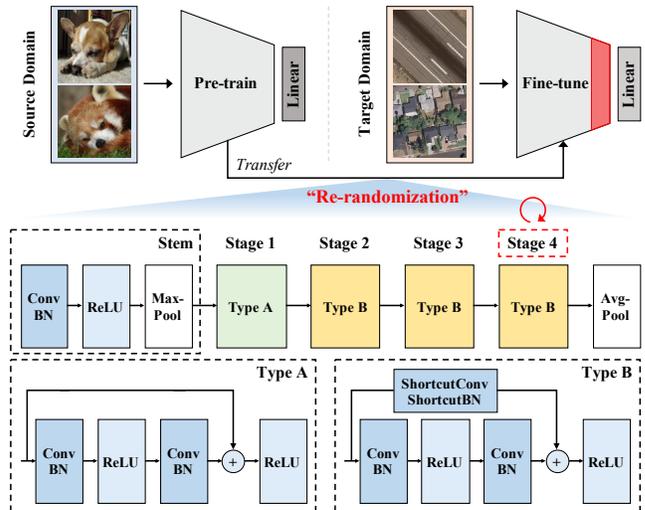


Figure 2: Overview of our proposed algorithm, ReFine, with the structure of ResNet10 backbone network.

Cross-domain few-shot learning (CD-FSL) addresses a problem when the source and target domains are extremely different, which is a more real-world scenario for FSL [13, 35]. Initially, Tseng et al. [35] proposed feature-wise transformation (FWT) that learns scale-and-shift meta-parameters using pseudo-unseen target data during meta-training. However, it showed poor performance on the recently released BSCD-FSL benchmark [13], consisting of four target datasets collected from different domains. In general, fine-tuning based approaches have been shown to outperform meta-learning based approaches such as FWT [13]. Therefore, recent CD-FSL studies have proposed their algorithms under a pre-training and fine-tuning scheme. These works have mainly focused on improving the pre-training phase, so that the pre-trained model is more suitable for adaptation to the target domain.

Re-randomization¹ has been widely studied in the field of language tasks [33, 41], in particular related to BERT, which is one of the most popular fine-tuning based language models. An interesting observation from Zhang et al. [41] is that re-randomizing the topmost block in BERT increases the performance for downstream tasks by reducing the fine-tuning workload. Concurrently to this observation, Tamkin et al. [33] examined the relations between the partial re-randomization of BERT and transferability of the layers. Meanwhile, in a visual task, Alabdulmohsin et al. [1] showed that placing more emphasis on the early layers of a convolutional neural network helps improve generalization. There have been similar attempts in meta-learning based FSL, e.g., zeroing the context vector for adaptation in each new task [42] and setting the classifier weight to have the same row vector (for any-shot problem) [8]. However, to the best of our knowledge, our work is the first to investigate the impact of re-randomization in fine-tuning based approaches for better CD-FSL.

¹Although some literature use the term *re-initialization*, we distinguish it from *re-randomization* because *re-initialization* reverts the values to the previously initialized ones. Refer to [40] for a formal definition. For a more concrete comparison, we have also dealt with *re-initialization* in Section 4.4.

3 REFINe: RE-RANDOMIZATION BEFORE FINE-TUNING

The objective of fine-tuning based CD-FSL algorithms is to learn a backbone f on the source data D_B with base classes C_B , extracting meaningful representations on the target data D_N with novel classes C_N , where $C_B \cap C_N = \emptyset$. However, because there is no access to target data, the pre-trained model is biased towards the source domain, especially in the upper layers that pertain to classification of base classes. To mitigate this, we re-randomize the upper layers of the pre-trained backbone f to reset source-fitted parameters, depicted in Figure 2. Specifically, the weights of convolutional layers are re-randomized to uniform distributions [14]. The scaling and shifting parameters of batch normalization layers are reset to ones and zeros, respectively.

The reason why *upper* layers of the backbone f are re-randomized is that more domain-specific representations are extracted as the depth increases in convolutional neural networks [1, 19, 26, 39]. Re-randomization of upper layers helps the training loss to escape from local minima attributed to D_B and allows bottom-level layers to be sufficiently updated, alleviating the gradient vanishing problem [17, 29]. This is in line with previous works which show that representation change is beneficial for CD-FSL [23, 35].

Finally, fine-tuning and evaluation are performed with episodes, each representing distinct tasks, sampled from the labeled target data D_N . Each episode consists of a support set D_s , used to fine-tune the partially re-randomized pre-trained model, and a query set D_q , used to evaluate after the fine-tuning. To sample an episode (D_s, D_q), n classes are first selected from C_N , and subsequently, k and k_q samples are selected per class for support and query sets, respectively, where $n = 5$ and $k \in \{1, 5\}$ in general.

4 EXPERIMENTS

We introduce the experimental setup in Section 4.1 and compare **ReFine** (ours) with two baselines in Section 4.2: (1) **Linear** is a linear probing method to fine-tune only the classifier layer; (2) **Transfer** is a transfer learning method to fine-tune the entire network without using re-randomization². We further investigate where and how to re-randomize in Section 4.3 and Section 4.4, respectively.

4.1 Experimental Setup

Datasets. For the source domain dataset, we use miniImageNet [36] and tieredImageNet (tieredIN) [28]. For the target domain, we use the BSCD-FSL benchmark [13], which consists of four different datasets: CropDisease [21], EuroSAT [15], ISIC [5], and ChestX [37], in order of similarity to miniIN.

Backbone and Training Setup. We use ResNet10 for miniIN and ResNet18 for tieredIN. Figure 2 describes the ResNet10 backbone. A family of ResNet consists of one stem module and four stages. The stem module consists of Conv-BN-ReLU-MaxPool layers. Each stage includes one or more convolution blocks, where resolution is halved and the number of channels is doubled in the first block, and they are maintained in the following blocks. For the pre-training and fine-tuning setups, we follow Guo et al. [13]

²Many meta-learning based approaches such as MAML, ProtoNet, ProtoNet+FWT, and MetaOptNet have worse performance than Transfer, which is shown in [13].

Table 1: 5-way k -shot accuracy over 600 tasks on {miniIN, tieredIN} \rightarrow {BSCD-FSL}. For ReFine, topmost layers in the last stage are re-randomized (see Section 4.3). Mean and 95% confidence interval are reported.

Source dataset	Methods	Target dataset			
		$k = 1$		$k = 5$	
		CropDisease	EuroSAT	ISIC	ChestX
miniImageNet	Linear	65.73 \pm .87	88.68 \pm .53	54.35 \pm .92	75.96 \pm .67
	Transfer	57.57 \pm .92	88.04 \pm .57	51.54 \pm .86	79.33 \pm .66
	ReFine	68.93\pm.84	90.75\pm.49	64.14\pm.82	82.36\pm.57
tieredImageNet	Linear	70.88\pm.90	90.04 \pm .49	50.84 \pm .93	69.36 \pm .73
	Transfer	63.93 \pm .85	85.73 \pm .60	50.62 \pm .86	72.24 \pm .65
	ReFine	67.39 \pm .89	90.96\pm.50	51.21\pm.82	74.39\pm.72
miniImageNet	Linear	30.42 \pm .54	42.97 \pm .56	22.17 \pm .37	25.80 \pm .43
	Transfer	32.31 \pm .63	49.67 \pm .62	21.82 \pm .40	26.10 \pm .44
	ReFine	35.30\pm.59	51.68\pm.63	22.48\pm.41	26.76\pm.42
tieredImageNet	Linear	28.14 \pm .55	37.20 \pm .53	22.33 \pm .40	25.03 \pm .41
	Transfer	32.31\pm.60	46.36\pm.65	22.49\pm.41	25.76\pm.41
	ReFine	28.24 \pm .48	38.83 \pm .54	21.68 \pm .36	24.83 \pm .37

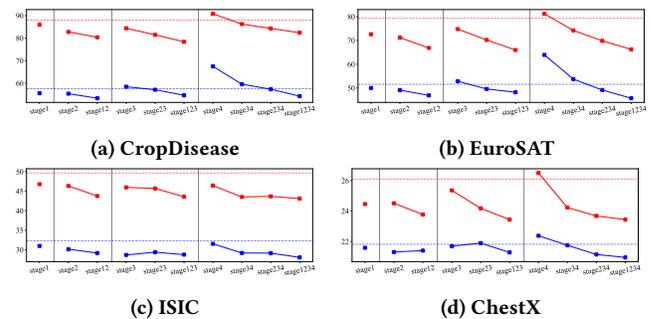


Figure 3: Accuracy trends according to the re-randomized stage(s). x -axis indicates re-randomized stage(s), and the blue and red lines indicate 1-/5-shot performance (%), respectively. The dashed lines are the performances of Transfer.

4.2 Performance Comparison

Table 1 describes the 5-way k -shot performance of Linear, Transfer, and ReFine in which a model is pre-trained on miniIN or tieredIN and then fine-tuned on BSCD-FSL. In most cases, ReFine outperforms Linear and Transfer. This implies that random parameters are generally better than the source-fitted parameters, especially of the topmost layers, for fine-tuning initialization. Meanwhile, in the ISIC and ChestX data, we observed that it might be advantageous to transfer the source information to the target without re-randomization when the source data becomes larger.

4.3 Ablation Study on Where to Re-randomize

We demonstrate that re-randomizing the extractor at the topmost layers is essential. We only consider Transfer as a baseline for fair comparison because ReFine fine-tunes the entire network. Figure 3 shows the test accuracy according to the re-randomized stage(s).

Table 2: 5-way k -shot accuracy over 600 tasks on {miniIN} \rightarrow {BSCD-FSL} according to the parts of re-randomization in the last stage. The topmost layers are boldfaced.

Path	Layer	Re-randomization layer											
Original	Conv1	✓	✓	✓	✓				✓	✓	✓	✓	
	BN1	✓							✓				
	Conv2		✓	✓	✓					✓	✓	✓	
	BN2		✓	✓	✓					✓	✓	✓	
ShortCut	ShortCutConv					✓	✓	✓	✓	✓	✓	✓	
	ShortCutBN						✓	✓	✓	✓	✓	✓	
1-shot													
	CropDisease	66.85	68.93	62.34	62.43	67.74	52.60	64.32	58.22	66.41	68.83	67.74	
	EuroSAT	61.99	64.14	59.19	58.91	62.96	53.03	56.29	57.33	62.01	62.77	63.90	
	ISIC	31.91	35.30	34.32	30.96	32.88	29.42	32.45	30.22	30.55	33.21	31.17	
	ChestX	22.20	22.48	22.00	21.53	22.08	21.19	21.55	21.60	21.99	22.82	22.46	
5-shot													
	CropDisease	89.78	90.75	87.29	89.12	89.99	87.78	90.03	89.44	90.03	90.89	90.82	
	EuroSAT	81.02	82.36	79.16	80.10	81.01	78.87	80.77	81.13	81.56	82.24	81.22	
	ISIC	49.73	51.68	51.90	46.17	50.20	46.49	48.00	46.21	46.85	49.59	46.44	
	ChestX	26.30	26.76	25.41	25.51	26.00	25.72	26.39	26.29	26.07	26.60	26.50	

In Figure 3, we observe that re-randomizing only the last stage is the best. This is indicated by the performance decrease from re-randomizing more stages within each subdivision separated by vertical lines, and by the best performance in the rightmost subdivision when only one stage is re-randomized.

Furthermore, we investigate layer-wise re-randomization within the last stage for more granular analysis on where to re-randomize. Table 2 describes the results according to the re-randomized layers in the last stage. Re-randomizing only {Conv2, BN2} shows the best performance overall. We conclude that *re-randomizing the topmost layers, excluding the shortcut path, in the last stage is a good rule of thumb*. A similar trend appears when the model is pre-trained on tieredIN, as described in Table 3.

4.4 Ablation Study on How to Re-randomize

Table 4 shows that re-randomizing the parameters following uniform distribution is generally the best practice. Uniform and Normal indicate that the values are sampled from the uniform and normal distribution. Orthogonal indicates the weights are randomized as an orthogonal matrix, as described in Saxe et al. [30]. Sparse indicates the weights are randomized as a sparse matrix, where non-zero elements are sampled from the zero-mean normal distribution, as described in Martens et al. [20]. Lottery refers to re-initialization, i.e., resetting parameters to their initial state, prior to training. In the model pruning literature, the lottery ticket hypothesis [12] suggests that re-initialization can improve performance. However, we find that re-randomization is better suited in the case of CD-FSL. We believe that although re-initialization can be helpful in the original domain, this is not true under domain differences.

5 CONCLUSION

We propose **ReFine** (Re-randomization before **Fi**ne-tuning), a simple yet effective method for CD-FSL, that involves resetting the parameters fitted to the source domain in order to maximize the efficacy of few-shot adaptation to the labeled target dataset. We demonstrate that our method outperforms conventional baselines

Table 3: 5-way k -shot accuracy over 600 tasks on {tieredIN} \rightarrow {BSCD-FSL} according to the parts of re-randomization in the last stage. The topmost layers are boldfaced.

	Re-randomization layer												
Block1.Conv1										✓		✓	
Block1.BN1										✓		✓	
Block1.Conv2									✓	✓		✓	
Block1.BN2									✓	✓		✓	
Block1.Conv3							✓	✓	✓		✓	✓	
Block1.BN3							✓	✓	✓		✓	✓	
Block1.ShortCutConv										✓	✓	✓	
Block1.ShortCutBN										✓	✓	✓	
Block2.Conv1							✓	✓	✓	✓	✓	✓	
Block2.BN1							✓	✓	✓	✓	✓	✓	
Block2.Conv2							✓	✓	✓	✓	✓	✓	
Block2.BN2							✓	✓	✓	✓	✓	✓	
Block2.Conv3	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
Block2.BN3	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
1-shot													
	CropDisease	67.39	68.31	60.98	52.84	48.43	42.61	51.82	52.78	51.37	49.28		
	EuroSAT	51.21	48.18	36.16	35.19	34.22	35.60	38.01	40.35	40.60	40.37		
	ISIC	28.24	28.06	27.02	26.64	26.12	26.94	26.24	26.35	26.42	26.70		
	ChestX	21.68	21.24	21.31	21.12	21.19	21.14	21.32	21.08	21.21	21.06		
5-shot													
	CropDisease	90.96	90.84	90.25	87.25	86.44	84.06	83.22	83.00	84.36	83.17		
	EuroSAT	74.39	74.03	71.54	67.58	66.26	62.66	60.22	62.07	63.17	60.40		
	ISIC	38.83	38.76	37.29	37.85	38.75	39.85	37.29	38.35	39.63	40.91		
	ChestX	24.83	24.90	24.64	24.08	23.66	23.54	23.23	22.88	23.15	22.89		

Table 4: Analysis on the initializing distribution of ReFine. Sparse distribution initializes parameters with 20% sparsity. Lottery indicates re-initialization.

Shot	Distribution	CropDisease	EuroSAT	ISIC	ChestX
1	Uniform	68.93±.84	64.14±.82	35.30±.59	22.48±.41
	Normal	69.34±.86	60.85±.82	31.35±.58	22.38±.39
	Orthogonal	67.96±.84	59.71±.83	31.05±.59	22.50±.38
	Sparse	69.07±.84	61.21±.82	31.10±.61	22.52±.39
	Lottery	61.53±.92	61.30±.88	31.27±.57	21.87±.36
5	Uniform	90.75±.49	82.36±.57	51.68±.63	26.76±.42
	Normal	91.31±.48	81.97±.58	46.92±.61	26.27±.43
	Orthogonal	91.01±.50	81.92±.58	45.73±.59	26.71±.43
	Sparse	91.33±.47	81.33±.60	45.62±.58	26.36±.43
	Lottery	89.81±.50	81.96±.58	48.10±.62	26.45±.43

under the CD-FSL setup. Furthermore, we investigate where and how to re-randomize the pre-trained models. We believe that our research will inspire CD-FSL researchers with the concept of removing information that is specific to the source domain.

ACKNOWLEDGMENTS

This work was supported by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government(MSIT) (No.2022-0-00641, XVoice: Multi-Modal Voice Meta Learning; No.2019-0-00075, Development of AI Autonomy and Knowledge Enhancement for AI Agent Collaboration).

REFERENCES

- [1] Ibrahim Alabdulmohsin, Hartmut Maennel, and Daniel Keysers. 2021. The Impact of Reinitialization on Generalization in Convolutional Neural Networks. *arXiv preprint arXiv:2109.00267* (2021).
- [2] Marcin Andrychowicz, Misha Denil, Sergio Gómez Colmenarejo, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando de Freitas. 2016. Learning to learn by gradient descent by gradient descent. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*. 3988–3996.
- [3] Haoxing Chen, Huaxiong Li, Yaohui Li, and Chunlin Chen. 2021. Multi-level metric learning for few-shot image recognition. *arXiv preprint arXiv:2103.11383* (2021).
- [4] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. 2019. A Closer Look at Few-shot Classification. In *ICLR*.
- [5] Noel Codella, Veronica Rotemberg, Philipp Tschandl, M Emre Celebi, Stephen Dusza, David Gutman, Brian Helba, Aadi Kalloo, Konstantinos Liopyris, Michael Marchetti, et al. 2019. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368* (2019).
- [6] Debansmit Das, Sungrack Yun, and Fatih Porikli. 2021. ConFeSS: A Framework for Single Source Cross-Domain Few-Shot Learning. In *International Conference on Learning Representations*.
- [7] Guneet Singh Dhillon, Pratik Chaudhari, Avinash Ravichandran, and Stefano Soatto. 2020. A Baseline for Few-Shot Image Classification. In *ICLR*.
- [8] Rafael Rego Drummond, Lukas Brinkmeyer, Josif Grabocka, and Lars Schmidt-Thieme. 2020. HIDRA: Head initialization across dynamic targets for robust architectures. In *ICDM*. 397–405.
- [9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*. 1126–1135.
- [10] Chelsea Finn, Kelvin Xu, and Sergey Levine. 2018. Probabilistic model-agnostic meta-learning. *arXiv preprint arXiv:1806.02817* (2018).
- [11] Sebastian Flennerhag, Andrei A. Rusu, Razvan Pascanu, Francesco Visin, Hujun Yin, and Raia Hadsell. 2020. Meta-Learning with Warped Gradient Descent. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rkeiQIBFPB>
- [12] Jonathan Frankle and Michael Carbin. 2018. The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks. In *International Conference on Learning Representations*.
- [13] Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. 2020. A broader study of cross-domain few-shot learning. In *ECCV*. 124–141.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*. 1026–1034.
- [15] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. 2019. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 12, 7 (2019), 2217–2226.
- [16] Ashraf Islam, Chun-Fu Chen, Rameswar Panda, Leonid Karlinsky, Rogerio Feris, and Richard J Radke. 2021. Dynamic Distillation Network for Cross-Domain Few-Shot Recognition with Unlabeled Data. *arXiv preprint arXiv:2106.07807* (2021).
- [17] Xingjian Li, Haoyi Xiong, Haozhe An, Cheng-Zhong Xu, and Dejing Dou. 2020. RIFLE: Backpropagation in Depth for Deep Transfer Learning through Re-Initializing the Fully-connected Layer. In *ICML*. 6010–6019.
- [18] Hanwen Liang, Qiong Zhang, Peng Dai, and Juwei Lu. 2021. Boosting the Generalization Capability in Cross-Domain Few-shot Learning via Noise-enhanced Supervised Autoencoder. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9424–9434.
- [19] Hartmut Maennel, Ibrahim Alabdulmohsin, Ilya Tolstikhin, Robert JN Baldock, Olivier Bousquet, Sylvain Gelly, and Daniel Keysers. 2020. What Do Neural Networks Learn When Trained With Random Labels? *arXiv preprint arXiv:2006.10455* (2020).
- [20] James Martens et al. 2010. Deep learning via hessian-free optimization.. In *ICML*, Vol. 27. 735–742.
- [21] Sharada P Mohanty, David P Hughes, and Marcel Salathé. 2016. Using deep learning for image-based plant disease detection. *Frontiers in plant science* 7 (2016), 1419.
- [22] Jaehoon Oh, Sungnyun Kim, Namgyu Ho, Jin-Hwa Kim, Hwanjun Song, and Se-Young Yun. 2022. Understanding Cross-Domain Few-Shot Learning: An Experimental Study. *arXiv preprint arXiv:2202.01339* (2022).
- [23] Jaehoon Oh, Hyungjun Yoo, ChangHwan Kim, and Se-Young Yun. 2021. BOIL: Towards Representation Change for Few-shot Learning. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=umldUL8rMH>
- [24] Cheng Perng Phoo and Bharath Hariharan. 2021. Self-training For Few-shot Transfer Across Extreme Task Differences. In *ICLR*.
- [25] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. 2019. Rapid learning or feature reuse? towards understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157* (2019).
- [26] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. 2019. Transfusion: Understanding transfer learning for medical imaging. *arXiv preprint arXiv:1902.07208* (2019).
- [27] Sachin Ravi and Hugo Larochelle. 2016. Optimization as a model for few-shot learning. In *ICLR*.
- [28] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. 2018. Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676* (2018).
- [29] Youngmin Ro, Jongwon Choi, Dae Ung Jo, Byeongho Heo, Jongin Lim, and Jin Young Choi. 2019. Backbone cannot be trained at once: Rolling back to pre-trained network for person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 8859–8867.
- [30] Andrew M Saxe, James L McClelland, and Surya Ganguli. 2013. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120* (2013).
- [31] Jake Snell, Kevin Swersky, and Richard S Zemel. 2017. Prototypical networks for few-shot learning. *arXiv preprint arXiv:1703.05175* (2017).
- [32] Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales. 2018. Learning to compare: Relation network for few-shot learning. In *CVPR*. 1199–1208.
- [33] Alex Tamkin, Trisha Singh, Davide Giovanardi, and Noah Goodman. 2020. Investigating transferability in pretrained language models. *arXiv preprint arXiv:2004.14975* (2020).
- [34] Yonglong Tian, Yue Wang, Dilip Krishnan, Joshua B Tenenbaum, and Phillip Isola. 2020. Rethinking few-shot image classification: a good embedding is all you need?. In *ECCV*. 266–282.
- [35] Hung-Yu Tseng, Hsin-Ying Lee, Jia-Bin Huang, and Ming-Hsuan Yang. 2020. Cross-Domain Few-Shot Classification via Learned Feature-Wise Transformation. In *ICLR*.
- [36] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. 2016. Matching networks for one shot learning. *NeurIPS* 29 (2016), 3630–3638.
- [37] Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadhi Bagheri, and Ronald M Summers. 2017. ChestX-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *CVPR*. 2097–2106.
- [38] Yaqing Wang, Quanming Yao, James T Kwok, and Lionel M Ni. 2020. Generalizing from a few examples: A survey on few-shot learning. *ACM Computing Surveys (CSUR)* 53, 3 (2020), 1–34.
- [39] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. 2014. How transferable are features in deep neural networks? *arXiv preprint arXiv:1411.1792* (2014).
- [40] Chiyuan Zhang, Samy Bengio, and Yoram Singer. 2019. Are all layers created equal? *arXiv preprint arXiv:1902.01996* (2019).
- [41] Tianyi Zhang, Felix Wu, Arzo Katiyar, Kilian Q Weinberger, and Yoav Artzi. 2020. Revisiting few-sample BERT fine-tuning. *arXiv preprint arXiv:2006.05987* (2020).
- [42] Luisa Zintgraf, Kyriacos Shiarli, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. 2019. Fast context adaptation via meta-learning. In *ICML*. 7693–7702.