
One-Shot Domain Incremental Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Domain incremental learning (DIL) has been discussed in previous studies on deep
2 neural network models for classification. In practice, however, we may encounter
3 a situation where we need to perform DIL under the constraint that the samples
4 on the new domain are observed only infrequently. In this study, we consider the
5 extreme case where we have only one sample from the new domain, which we call
6 one-shot DIL (ODIL). In simulation experiments on ODIL, we observed that the
7 accuracy on both the new domain and the original domain deteriorated even on
8 applying existing DIL methods. We analyzed the reason for this problem through
9 various investigations and discovered that the cause would be the statistics of the
10 batch normalization layers. According to our analysis, we propose a new technique
11 regarding these statistics and demonstrate the effectiveness of the proposed method
12 in ODIL through experiments on open datasets.

13 1 Introduction

14 In recent years, deep learning has been widely used in image recognition, speech recognition, and
15 natural language processing [18]. There is a need to update a trained neural network model so that
16 the model can classify samples correctly on a new, untrained input distribution (domain) [6]. The
17 additional training for the new domain is called domain incremental learning (DIL) [32]. In DIL, we
18 assume that the input domains increase in each class, while the number of classes remains constant.
19 We need to correctly classify inputs on both the new and original domains.

20 In practice, although previous studies on DIL assume many samples from the new domain, we often
21 face situations in which we have few samples from the new domain.¹ Therefore, in this study, we
22 explore one-shot DIL (ODIL), which is supposed to have only one sample from the new domain.
23 Figure 1 shows an example of the new and original domains in ODIL, arranged using CIFAR10 [17].
24 In this example, no truck images are included in the original 9-class data, and such images are not
25 used in the first training. Suppose that one truck image is given and added to the automobile class
26 as a new domain after the first training. We assume that the model obtained by the first training
27 misclassifies the given truck image. Therefore, we attempt to update this model with the given truck
28 image as the second training, so that the model classifies truck images into the automobile class.

29 Previous studies proposed DIL methods with mechanisms to improve the accuracy on the new domain
30 while maintaining the accuracy on the original domain. Examples of such DIL methods are Elastic
31 Weight Consolidation (EWC) [16] and Gradient Episodic Memory (GEM) [21]. However, in our
32 experiments simulating ODIL, we could not maintain the accuracy on the original domain even after
33 applying EWC or GEM directly under the ODIL setting. Furthermore, though the goal of ODIL is to
34 improve the accuracy on the new domain, it worsened in ODIL with EWC or GEM. We carefully
35 investigated the reason for this deterioration and found that the statistics maintained in the batch
36 normalization layers [13] are more drastically affected under the ODIL setting compared to the

¹In Appendix B, we give an example of situations in which we have few samples from the new domain.

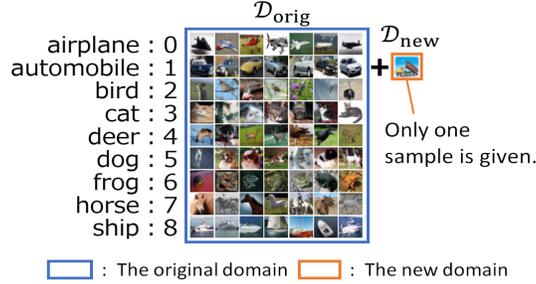


Figure 1: Example of the original domain and the new domain in one-shot domain incremental learning (ODIL) with CIFAR10 [17]. In this example, trucks are added to the “automobile” class as the new domain. However, only one sample is added.

37 general DIL setting. Therefore, in this paper, we propose a new technique regarding the statistics in
 38 the batch normalization layers to prevent the degradation of accuracy. We report the experimental
 39 results indicating that our technique is necessary even if we apply EWC or GEM under the ODIL
 40 setting. Note that the related work of our study is explained in Appendix A.

41 2 Problem Description

42 In this section, we define ODIL. Let $\mathcal{X}(\subseteq \mathbb{R}^d)$ be an input space. Let $\mathcal{Y} = \{1, \dots, K\}$ be the set of
 43 classes, where $K \in \mathbb{N}$ is the number of classes. Let $\mathcal{D}_{\text{orig}} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ ($\mathbf{x}_n \in \mathcal{X}, y_n \in \mathcal{Y}; n =$
 44 $1, \dots, N$) be a labeled dataset for classification, where $\mathbf{x}_n \sim p_{\text{orig}}(\mathbf{x}|y = k)$ ($k = 1, \dots, K$).
 45 In this work, $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Y}$ parametrized by $\theta \in \Theta(\subseteq \mathbb{R}^p)$ denotes a neural network model.
 46 We suppose that a neural network model $f_{\hat{\theta}_{\text{orig}}}$ ($\hat{\theta}_{\text{orig}} \in \Theta$) trained with $\mathcal{D}_{\text{orig}}$ is given. Under
 47 this premise, we suppose that a new labeled dataset \mathcal{D}_{new} is given. We assume that $|\mathcal{D}_{\text{new}}| = 1$
 48 and represent $\mathcal{D}_{\text{new}} = \{(\mathbf{x}_0, y_0)\}$ ($\mathbf{x}_0 \in \mathcal{X}, y_0 \in \mathcal{Y}$). The new sample (\mathbf{x}_0, y_0) is based on the
 49 following assumptions: 1) The trained model $f_{\hat{\theta}_{\text{orig}}}$ misclassifies the new sample (\mathbf{x}_0, y_0) . 2) The
 50 sample (\mathbf{x}_0, y_0) is not included in the original dataset $\mathcal{D}_{\text{orig}}$. 3) $\mathbf{x}_0 \sim p_{\text{new}}(\mathbf{x}|y = y_0)$, where
 51 $p_{\text{orig}}(\mathbf{x}|y = y_0)$ and $p_{\text{new}}(\mathbf{x}|y = y_0)$ are not equivalent.

52 We define ODIL as updating the trained model $f_{\hat{\theta}_{\text{orig}}}$ with \mathcal{D}_{new} to obtain a new model $f_{\hat{\theta}_{\text{new}}}$ ($\hat{\theta}_{\text{new}} \in$
 53 Θ) that can correctly classify the inputs on the new domain $p_{\text{new}}(\mathbf{x}|y = y_0)$. Note that the model
 54 is required to maintain the classification accuracy on the original domain $p_{\text{orig}}(\mathbf{x}|y = k)$ as much
 55 as possible. As in previous studies [21, 4], we allow a subset \mathcal{M} of the examples in the original
 56 data $\mathcal{D}_{\text{orig}}$ to be stored. We can use the subset \mathcal{M} in ODIL to prevent the deterioration of the accuracy
 57 on the original domain.

58 For example, in Fig. 1, the original domain contains classes other than trucks. The new domain
 59 contains only trucks. The given truck image corresponds to \mathbf{x}_0 , and the label y_0 represents the
 60 automobile class. We expect that the new model $f_{\hat{\theta}_{\text{new}}}$ will correctly classify the trucks as well as the
 61 other data. Most previous studies about DIL assumed that the domains expand in all classes. Unlike
 62 them, we assume that the domain increases in only one class because the new dataset contains only
 63 one sample.

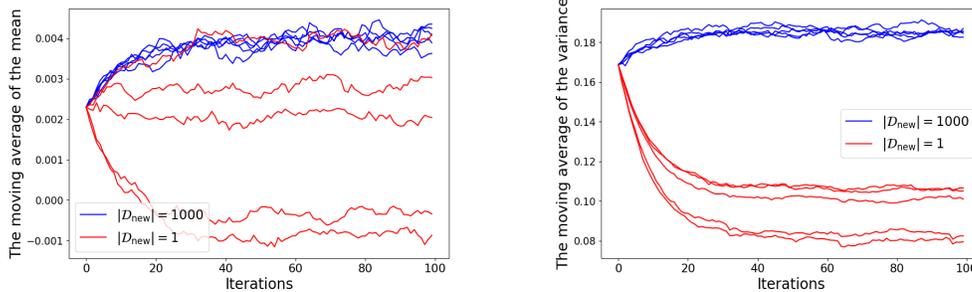
64 3 Method for ODIL

65 3.1 Problems with Batch Normalization in ODIL

66 EWC and GEM have mechanisms to improve accuracy on the new domain while maintaining
 67 accuracy on the original domain [16, 21]. Therefore, it is natural to expect that the same results can
 68 be achieved in ODIL as well. However, in experiments simulating ODIL, we could not maintain
 69 the accuracy on the original domain even when using EWC or GEM. Further, although ODIL was
 70 intended to improve the accuracy on the new domain, EWC and GEM actually worsened the accuracy
 71 of the new domain in ODIL, as shown in Table 1. Table 1 summarizes the test accuracy on the new

Table 1: Test accuracy on the original and new domains when $|\mathcal{D}_{\text{new}}| = 1000$ or $|\mathcal{D}_{\text{new}}| = 1$.

	$f_{\hat{\theta}_{\text{orig}}}$	$f_{\hat{\theta}_{\text{new}}}$					
		CE		CE+EWC		CE+GEM	
		1000	1	1000	1	1000	1
p_{new}	0.6940	0.9410	0.5765	0.9390	0.5955	0.9960	0.7315
p_{orig}	0.9550	0.9420	0.4936	0.9434	0.4647	0.9306	0.9071



(a) Transition of the moving averages of the mean. (b) Transition of the moving averages of the variance.

Figure 2: Transition of the moving averages of the statistics at the batch normalization layer closest to the input layer in ResNet18 [9]. We accumulated the moving averages of the statistics at every forward propagation and plotted their sequences. Since the input is normalized in parallel for each channel in the batch normalization layer, we computed the average for the channels.

72 domain $p_{\text{new}}(\mathbf{x}|y = y_0)$ and the original domain $p_{\text{orig}}(\mathbf{x}|y = k)(k = 1, \dots, K)$ in experiments
 73 simulating ODIL as shown in Fig. 1. We tried two cases, $|\mathcal{D}_{\text{new}}| = 1$ and $|\mathcal{D}_{\text{new}}| = 1000$. In the
 74 $|\mathcal{D}_{\text{new}}| = 1$ setting, we resampled data from \mathcal{D}_{new} with data augmentation [28] at each step of the
 75 optimization to increase the number of learnable samples for the new domain. The details of this data
 76 augmentation are explained in Appendix C. We used ResNet18 [9] as an image classifier. CE denotes
 77 the case where we performed optimization with simple cross-entropy loss. CE+EWC and CE+GEM
 78 are the same as CE except that EWC and GEM are applied, respectively. The setup of the experiments
 79 is explained in Appendix D. When $|\mathcal{D}_{\text{new}}| = 1$, the accuracy on both domains decreases after training
 80 with \mathcal{D}_{new} , in the case of CE and CE+EWC. This deterioration on both domains is a problem specific
 81 to ODIL as the deterioration on the new domain does not occur when $|\mathcal{D}_{\text{new}}| = 1000$.

82 We carefully analyzed the cause of this problem and discovered that the cause lies in the statistics in
 83 the batch normalization layers. Figure 2 shows the transitions of the moving averages of the mean
 84 and variance of a batch normalization layer as $f_{\hat{\theta}_{\text{orig}}}$ is updated with \mathcal{D}_{new} . We chose an example
 85 in CE. We ran five trials with varying \mathcal{D}_{new} for each of the $|\mathcal{D}_{\text{new}}| = 1000$ and $|\mathcal{D}_{\text{new}}| = 1$ settings.
 86 We plotted all results of the five trials as the five lines in the figures. In Fig. 2a, the moving average
 87 of the mean shifts in more widely different directions over the five trials in $|\mathcal{D}_{\text{new}}| = 1$ than in
 88 $|\mathcal{D}_{\text{new}}| = 1000$. When $|\mathcal{D}_{\text{new}}| = 1$, for some trials, a positive moving average is obtained, while for
 89 others, a negative moving average is obtained. This variation is caused by the fact that the mean value
 90 is calculated under a strong influence of the single input \mathbf{x}_0 . The mean value heavily influenced by
 91 the single input \mathbf{x}_0 causes the trainable weights and biases in the batch normalization layers to be
 92 updated in unexpected directions during training. Furthermore, in inference, the moving average
 93 that is heavily influenced by the single input \mathbf{x}_0 degrades the generalization ability. In Fig. 2b,
 94 the moving average of the variance shifts in the opposite direction between $|\mathcal{D}_{\text{new}}| = 1000$ and
 95 $|\mathcal{D}_{\text{new}}| = 1$. The moving average of the variance increases when $|\mathcal{D}_{\text{new}}| = 1000$, which is natural
 96 from the viewpoint that the number of domains increases. However, when $|\mathcal{D}_{\text{new}}| = 1$, the moving
 97 average of the variance decreases. This discrepancy is due to a lack of diversity in the new data in
 98 ODIL. This results in a reduced variance in the mini-batches, in turn causing the trainable weights
 99 and biases in the batch normalization layers to be updated in unexpected directions during training.
 100 Furthermore, in inference, the much smaller moving average of variance causes misclassification
 101 because the outputs of the batch normalization layers become much larger than the inputs. The lower
 102 accuracies in $|\mathcal{D}_{\text{new}}| = 1$ are due to the above problems. Although we have tried EWC and GEM,
 103 existing DIL methods other than EWC and GEM are also likely to show the same trend as long as

Table 2: Test accuracy on the new and original domains in the baseline or proposed (CIFAR10)

	$f_{\hat{\theta}_{\text{orig}}}$	$f_{\hat{\theta}_{\text{new}}}$					
		CE		CE+EWC		CE+GEM	
		baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.6940	0.5765 ± 0.2297	0.9595 ± 0.0430	0.5955 ± 0.2614	0.9420 ± 0.0586	0.7315 ± 0.2287	0.9485 ± 0.0417
p_{orig}	0.9550	0.4936 ± 0.1847	0.9335 ± 0.0440	0.4647 ± 0.2202	0.9332 ± 0.0192	0.9071 ± 0.2512	0.9404 ± 0.0250

Table 3: Test accuracy on the new and original domains in the baseline or proposed (RESISC45)

	$f_{\hat{\theta}_{\text{orig}}}$	$f_{\hat{\theta}_{\text{new}}}$					
		CE		CE+EWC		CE+GEM	
		baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.1071	0.4250 ± 0.1217	0.8250 ± 0.0862	0.3893 ± 0.1416	0.8214 ± 0.0868	0.3643 ± 0.1081	0.7786 ± 0.1020
p_{orig}	0.9580	0.9482 ± 0.0096	0.8912 ± 0.0712	0.9502 ± 0.0166	0.8879 ± 0.0728	0.9478 ± 0.0083	0.9294 ± 0.0549

104 batch normalization layers are present in the model. Therefore, we propose to use fixed statistics in
 105 the batch normalization layers in ODIL, as detailed in the following section.

106 3.2 Modifying Batch Normalization Statistics

107 As explained in Section 3.1, the shifted statistics in the batch normalization layers degrade the
 108 accuracy on the new and original domains under the ODIL setting. Therefore, we modify the statistics
 109 as follows: 1) We do not update the moving averages of the statistics and we fix them to the values
 110 they had before training the new domain. 2) We use these constant moving averages for batch
 111 normalization when calculating forward and backward propagation during training. 3) We also use
 112 the constant moving averages in inference. These improvements avoid the unexpected shift of the
 113 statistics shown in Section 3.1.

114 4 Experiments

115 To investigate the effectiveness of the proposed method, we performed experiments on the image
 116 datasets CIFAR10 [17] and RESISC45 [5]. As experiments using CIFAR10, the example in Fig. 1 was
 117 adopted. In addition, as experiments using RESISC45, we adopted the example given by replacing
 118 “truck” with “airplane” and “automobile” with “airport” in Fig. 1. We used ResNet18 [9] as an image
 119 classifier. We trained the models in ODIL with the following two settings and calculated the test
 120 accuracy on both the new domain $p_{\text{new}}(\mathbf{x}|y = y_0)$ and the original domain $p_{\text{orig}}(\mathbf{x}|y = k)(k =$
 121 $1, \dots, K)$. **baseline:** In the batch normalization layers, the statistics are calculated from the inputs
 122 as usual. Simultaneously, the moving averages of the statistics are updated. **proposed:** As explained
 123 in Section 3.2, the batch normalization layers constantly use the moving averages obtained before
 124 training the new domain. These moving averages are not updated. Note that there is no difference
 125 between “baseline” and “proposed” other than the batch normalization layers. The setup of the
 126 experiments is explained in Appendix D. Additional experiments are shown in Appendix E.

127 Tables 2 and 3 list the results of CIFAR10 and RESISC45, respectively. We calculated the median and
 128 standard deviation of the accuracies over 10 trials with varying \mathcal{D}_{new} (median \pm standard deviation).
 129 In the case of CIFAR10, some settings of the baseline degraded the accuracy on both the new and
 130 original domains through ODIL. However, the proposed method avoided the deterioration of accuracy
 131 and improved the accuracy on the new domain. The proposed method improved the accuracy not only
 132 on the new domain but also on the original domain in the experiments with CIFAR10. There is no
 133 example where the proposed method is worse than the baseline in Table 2. In the case of RESISC45,
 134 the baseline improved the accuracy on the new domain through ODIL; however, compared to the
 135 baseline, the proposed method achieved higher accuracy on the new domain. Although there is a
 136 trade-off between the accuracy on the new domain and that on the original domain, the decrease in
 137 the accuracy on the original domain is much smaller than the increase in the accuracy on the new
 138 domain.

References

- 139
- 140 [1] Antreas Antoniou, Massimiliano Patacchiola, Mateusz Ochal, and Amos Storkey. Defining benchmarks
141 for continual few-shot learning. *Workshop on Meta-Learning, 34th Conference on Neural Information
142 Processing Systems (NeurIPS 2020)*, 2020.
- 143 [2] Antreas Antoniou and Amos J Storkey. Learning to learn by self-critique. In *Advances in Neural
144 Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 145 [3] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMP-
146 STAT'2010*, pages 177–186, Heidelberg, 2010. Physica-Verlag HD.
- 147 [4] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and SIMONE CALDERARA. Dark
148 experience for general continual learning: a strong, simple baseline. In *Advances in Neural Information
149 Processing Systems*, volume 33, pages 15920–15930. Curran Associates, Inc., 2020.
- 150 [5] Gong Cheng, Junwei Han, and Xiaoqiang Lu. Remote sensing image scene classification: Benchmark and
151 state of the art. *Proceedings of the IEEE*, 105(10):1865–1883, 2017.
- 152 [6] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory
153 Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks.
154 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3366–3385, 2022.
- 155 [7] Xiang Gao, Ripon K. Saha, Mukul R. Prasad, and Abhik Roychoudhury. Fuzz testing based data
156 augmentation to improve robustness of deep neural networks. In *2020 IEEE/ACM 42nd International
157 Conference on Software Engineering (ICSE)*, pages 1147–1158, 2020.
- 158 [8] Vibhor Gupta, Jyoti Narwariya, Pankaj Malhotra, Lovekesh Vig, and Gautam Shroff. Continual learning
159 for multivariate time series tasks with variable input dimensions. In *2021 IEEE International Conference
160 on Data Mining (ICDM)*, pages 161–170, 2021.
- 161 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
162 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- 163 [10] Michael Hersche, Geethan Karunaratne, Giovanni Cherubini, Luca Benini, Abu Sebastian, and Abbas
164 Rahimi. Constrained few-shot class-incremental learning. In *Proceedings of the IEEE/CVF Conference on
165 Computer Vision and Pattern Recognition (CVPR)*, pages 9057–9067, June 2022.
- 166 [11] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier
167 incrementally via rebalancing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
168 Pattern Recognition (CVPR)*, June 2019.
- 169 [12] YenChang Hsu, YenCheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual learning
170 scenarios: A categorization and case for strong baselines. *Continual Learning Workshop, 32nd Conference
171 on Neural Information Processing Systems (NIPS 2018)*, 2018.
- 172 [13] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing
173 internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning*, pages
174 448–456, 2015.
- 175 [14] Fahdi Kanavati and Masayuki Tsuneki. Partial transfusion: on the expressive influence of trainable batch
176 norm parameters for transfer learning. In *Proceedings of the Fourth Conference on Medical Imaging with
177 Deep Learning*, pages 338–353, 2021.
- 178 [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:
179 1412.6980*, 2017.
- 180 [16] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu,
181 Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia
182 Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks.
183 *Proceedings of the National Academy of Sciences of the United States of America*, 114(13):3521–3526,
184 2017.
- 185 [17] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical
186 Report 0, University of Toronto, Toronto, Ontario, 2009.
- 187 [18] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–444, 2015.
- 188 [19] Yann LeCun, Corinna Cortes, and CJ Burges. MNIST handwritten digit database. *ATT Labs [Online]*.
189 Available: <http://yann.lecun.com/exdb/mnist>, Accessed on July 21, 2023, 2, 2010.

- 190 [20] Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting batch normalization for
191 practical domain adaptation. *Workshop in International Conference on Learning Representations*, 2017.
- 192 [21] David Lopez-Paz and Marc' Aurelio Ranzato. Gradient episodic memory for continual learning. In
193 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 194 [22] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International
195 Conference on Learning Representations*, 2017.
- 196 [23] Sébastien Marcel and Yann Rodriguez. Torchvision the machine-vision package of torch. In *Proceedings
197 of the 18th ACM international conference on Multimedia*, pages 1485–1488, 2010.
- 198 [24] SylvestreAlvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. iCaRL: Incre-
199 mental classifier and representation learning. In *Proceedings of the IEEE Conference on Computer Vision
200 and Pattern Recognition (CVPR)*, July 2017.
- 201 [25] Xuhong Ren, Bing Yu, Hua Qi, Felix Juefei-Xu, Zhuo Li, Wanli Xue, Lei Ma, and Jianjun Zhao. Few-shot
202 guided mix for dnn repairing. In *2020 IEEE International Conference on Software Maintenance and
203 Evolution (ICSME)*, pages 717–721, 2020.
- 204 [26] Pierre Sermanet and Yann LeCun. Traffic sign recognition with multi-scale convolutional networks. In
205 *The 2011 International Joint Conference on Neural Networks*, pages 2809–2813, 2011.
- 206 [27] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting
207 with hard attention to the task. In *Proceedings of the 35th International Conference on Machine Learning*,
208 pages 4548–4557, 2018.
- 209 [28] Connor Shorten and Taghi M. Khoshgoftaar. A survey on image data augmentation for deep learning.
210 *Journal of Big Data*, 6(60), 2019.
- 211 [29] Jeongju Sohn, Sungmin Kang, and Shin Yoo. Arachne: Search-based repair of deep neural networks. *ACM
212 Trans. Softw. Eng. Methodol.*, 32(4), May 2023.
- 213 [30] Domen Tabernik and Danijel Skočaj. Deep learning for large-scale traffic-sign detection and recognition.
214 *IEEE Transactions on Intelligent Transportation Systems*, 21(4):1427–1440, 2020.
- 215 [31] Xiaoyu Tao, Xiaopeng Hong, Xinyuan Chang, Songlin Dong, Xing Wei, and Yihong Gong. Few-shot
216 class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
217 Recognition (CVPR)*, June 2020.
- 218 [32] Gido M. van de Ven and Andreas S. Tolias. Generative replay with feedback connections as a general
219 strategy for continual learning. arXiv preprint arXiv: 1809.10635, 2018.
- 220 [33] Moslem Yazdanpanah, Aamer Abdul Rahman, Muawiz Chaudhary, Christian Desrosiers, Mohammad
221 Havaei, Eugene Belilovsky, and Samira Ebrahimi Kahou. Revisiting learnable affines for batch norm in
222 few-shot transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
223 Recognition (CVPR)*, pages 9109–9118, June 2022.
- 224 [34] Bing Yu, Hua Qi, Qing Guo, Felix Juefei-Xu, Xiaofei Xie, Lei Ma, and Jianjun Zhao. DeepRepair: Style-
225 guided repairing for deep neural networks in the real-world operational environment. *IEEE Transactions
226 on Reliability*, pages 1–16, 2021.
- 227 [35] Chi Zhang, Nan Song, Guosheng Lin, Yun Zheng, Pan Pan, and Yinghui Xu. Few-shot incremental learning
228 with continually evolved classifiers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
229 Pattern Recognition (CVPR)*, pages 12455–12464, June 2021.
- 230 [36] Hao Zhang and W.K. Chan. Apricot: A weight-adaptation approach to fixing deep learning models. In *2019
231 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 376–387,
232 2019.
- 233 [37] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing
234 He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2021.
- 235 [38] Stop sign. Available Online: https://en.wikipedia.org/wiki/Stop_sign, Accessed on January
236 16, 2023.

237 **A Related Work**

238 **A.1 Continual Learning**

239 DIL is a part of continual learning [6]. There are three types of scenarios in continual learning: DIL,
240 class incremental learning (CIL), and task incremental learning (TIL) [32]. In DIL, several datasets
241 with different domains are observed over time. All the datasets have the same number of classes, and
242 a new domain is added to each class as a new dataset is observed. Both the original and new domains
243 must be classifiable by the training model.

244 In CIL, each dataset that is observed over time contains an exclusive subset of the classes from a
245 whole set of data [24, 11]. Therefore, the observation of a new dataset means the addition of new
246 classes. The training model must classify all classes. Each time a new dataset is observed, the number
247 of units in the output layer of neural network models is expanded for accommodating new classes.
248 Note that this scenario assumes that different datasets have different domains.

249 In TIL, the output spaces of all the datasets observed over time are disjoint [27, 8]. For example,
250 consider the case where the original dataset deals with a 10-class classification problem, while the
251 new dataset deals with a single output regression task [12]. Neural network models often have a
252 multi-head output layer in TIL, wherein each head is responsible for a specific task. Note that in TIL,
253 we assume that different datasets have different domains and different output distributions.

254 **A.2 Few-Shot Continual Learning**

255 Recent studies on computer vision discussed few-shot CIL (FCIL) [35, 10]. FCIL is similar to CIL
256 except that it involves a limited number of training samples. In many cases, we assume that the first
257 training dataset has a large sample size, while the subsequent datasets have smaller sizes [31]. The
258 new dataset has new classes that are not included in the previous datasets. A training model must
259 discriminate between the new and original classes, even though the new classes have few samples.

260 Another study explored few-shot DIL [1]. The few-shot DIL called New Classes with Overwrite
261 Settings in [1] is similar to ours. The difference lies in the number of classes in the new domain.
262 In the earlier work [1], the number of classes in the new domain is the same as that in the original
263 domain. In contrast, in the present work, the new domain contains only one class regardless of
264 the number of classes in the original domain, and only one sample belongs to that one class. The
265 meta-learning methods used in [1] assume that there are multiple classes and multiple samples in
266 each domain. The Self-Critique and Adapt model (SCA) [2], which can achieve the best accuracy
267 in [1], needs two datasets—a support set and a target set—for each domain. The two datasets must be
268 different because the target set is used for validation. Hence, it is difficult to apply this method to our
269 ODIL.

270 **A.3 Neural Network Software Repair**

271 In our settings, a new sample from the new domain is assumed to be misclassified before training the
272 new domain. Therefore, we consider ODIL as the process of repairing or debugging neural network
273 models. Some studies have discussed neural network software repair in software engineering [36, 29].
274 These studies aimed to reduce the number of misclassified samples by updating the weight parameters
275 of a trained model. Note that the predictions of correctly classified samples cannot be changed when
276 the weight parameters are updated. Some studies assumed 1000 samples as the number of samples
277 to be repaired [34, 25]. Other studies set the number of samples to be repaired depending on the
278 accuracy of a pre-trained model [36, 7]. Few studies discussed neural network software repair with
279 limited training samples. In contrast, the present study considers the situation of only one additional
280 class with only one sample. Therefore, our study is valuable not only for continual learning but also
281 for neural network software repair.

282 **A.4 Transfer Learning**

283 Transfer learning [37] is similar to continual learning and neural network software repair. However,
284 transfer learning differs from these techniques in that it allows us to forget the original data pattern.
285 Some studies on transfer learning discussed batch normalization. Kanavati and Tsuneki [14] claimed
286 that fine-tuning only the trainable weights and biases of the batch normalization layers yields

287 performance similar to that achieved by fine-tuning all the weights of the neural network models.
 288 Meanwhile, Yazdanpanah et al. [33] showed that shifting and scaling normalized features by the
 289 weights and biases of the batch normalization layers is detrimental in few-shot transfer learning. They
 290 demonstrated that the removal of batch normalization weights and biases can have a positive impact
 291 on performance.

292 In addition, a study explored not only trainable weights and biases but also the statistics that are
 293 computed by the inputs in the batch normalization layers. Li et al. [20] achieved a deep adaptation
 294 effect to the new domain by modulating the statistics from the source domain to the target domain in
 295 the batch normalization layers. In the present work, we proposed a technique different from that used
 296 in the earlier study [20] for the statistics in the batch normalization layers. We changed the setting of
 297 statistics in the batch normalization layers depending on the number of samples on the new domain.

298 B Applications

299 Our study contributes to many cases where samples from new domains are rarely observed over time,
 300 but the model must adapt to the new domains. For example, consider the traffic sign recognition
 301 requirements of self-driving cars. For classifying traffic signs into categories such as stop, one-way,
 302 and speed limit, based on images obtained from cameras, automatic classification is made possible by
 303 using neural network models trained on a dataset containing images of traffic signs [26, 30]. However,
 304 even if we train a large image dataset of traffic signs for this classification, we may encounter a new,
 305 unprecedented type of traffic sign that will be misclassified by the trained model in the test drives. In
 306 the category of stop signs, many different types of stop signs are in use [38]. Each time we find a
 307 new type of stop sign, we must update the trained model. In particular, even if there is only one stop
 308 sign in the world that looks like the one we found, we still need to update the model so long as the
 309 self-driving cars are using it. In addition, the updated model is expected to successfully recognize the
 310 same stop sign in different situations.

311 C Algorithm for ODIL

Algorithm 1: Our ODIL algorithm with improved batch normalization.

Input: \mathcal{M} , $\mathcal{D}_{\text{new}} = \{(\mathbf{x}_0, y_0)\}$, $f_{\hat{\theta}_{\text{orig}}}$, \mathcal{L} : loss function, ν : learning rate,
 transform: the random operator of augmentation,
 judge : $\Theta \rightarrow \{\text{True}, \text{False}\}$: the judgment for termination of iteration

Output: $f_{\hat{\theta}_{\text{new}}}$

```

1  $\theta_0 \leftarrow \hat{\theta}_{\text{orig}}$ 
2  $t \leftarrow 0$ 
3 while judge( $\theta_t$ ) = False do
4   Sample a mini-batch  $\mathcal{B}$  from  $\mathcal{M}$ .
5   for  $m = 1, 2, \dots, M$  do
312 6   |  $\mathbf{z}_m \leftarrow \text{transform}(\mathbf{x}_0)$ 
7   end
8    $\mathcal{C} \leftarrow \{(\mathbf{z}_1, y_0), (\mathbf{z}_2, y_0), \dots, (\mathbf{z}_M, y_0)\}$ 
9    $\mathcal{B}' \leftarrow \mathcal{B} \cup \mathcal{C}$ 
10  Compute the forward and backward propagation of the loss  $\mathcal{L}(\theta_t; \mathcal{B}')$  with the proposed method in
    Section 3.2.
11   $\theta_{t+1} \leftarrow \theta_t - \nu \nabla_{\theta_t} \mathcal{L}(\theta_t; \mathcal{B}')$ 
12   $t \leftarrow t + 1$ 
13 end
14  $\hat{\theta}_{\text{new}} \leftarrow \theta_t$ 
15 return  $f_{\hat{\theta}_{\text{new}}}$ 

```

313 Algorithm 1 shows the details of our algorithm for ODIL. We assume that the neural network
 314 model $f_{\hat{\theta}_{\text{orig}}}$ trained on the original data $\mathcal{D}_{\text{orig}}$ is given, as in Section 2. In ODIL, the new domain
 315 $p_{\text{new}}(\mathbf{x}|y = y_0)$ becomes the training target in addition to the original domain $p_{\text{orig}}(\mathbf{x}|y = k)(k =$
 316 $1, \dots, K)$. However, we have only one sample on the new domain, denoted as (\mathbf{x}_0, y_0) . This makes
 317 training on the new domain difficult. Therefore, we apply data augmentation [28] to \mathbf{x}_0 , as described
 318 in Lines 5-8 in Alg. 1. We replicate the new data \mathbf{x}_0 into M samples ($M \in \mathbb{N}$) and transform them

319 randomly and differently in each iteration.² The vectors $z_1, z_2, \dots, z_M \in \mathcal{X}$ denote the new samples
 320 obtained by transforming x_0 . We assume that the transformed samples z_1, z_2, \dots, z_M belong to the
 321 same class y_0 as x_0 . However, if all the samples in the mini-batches belong to the class y_0 , updating θ
 322 with these mini-batches may result in a model that classifies all the inputs into y_0 . To avoid this
 323 problem, we concatenate the mini-batch \mathcal{B} , which is the mini-batch of \mathcal{M} , with $\mathcal{C} = \{(z_m, y_0)\}_{m=1}^M$
 324 (Line 9).

325 Then, in forward and backward propagation (Line 10), we perform batch normalization using the
 326 technique proposed in Section 3.2. The moving averages of the statistics are not updated. Note that
 327 we did not use the proposed technique in the experiments explained in Section 3.1, and used it as
 328 “proposed” in the experiments explained in Section 4. We continue the training iteration as long as
 329 a function “judge” returns “False” and stop it when “judge” returns “True” (Line 3). We define
 330 “judge” as the function that returns “True” if the element corresponding to y_0 in the softmax output
 331 of $f_\theta(x_0)$ is greater than a threshold $\delta \in [0, 1]$. “judge” returns “False” if it is less than δ .

332 When using EWC in general DIL, the loss for the new domain is computed only with new data [16].
 333 In ODIL, however, we compute the loss for the new domain with \mathcal{B}' in Alg. 1 to prevent the trained
 334 model from classifying all the inputs into y_0 . The Fisher information needed for EWC is computed
 335 with \mathcal{M} . For GEM, we compute the loss for the new domain with \mathcal{B}' , similar to EWC. The loss for
 336 the original domain is computed with a mini-batch \mathcal{B} randomly sampled from \mathcal{M} .

337 D Details of Experiments

338 D.1 Datasets

339 Although we only show the results when we used CIFAR10 and RESISC45 in the main paper, we
 340 also show the results when we used MNIST [19] in Appendix E. CIFAR10 and MNIST classify
 341 images into 10 classes, while RESISC45 classifies images into 45 classes. Each image is assigned a
 342 label representing the class to which it belongs. To extract the datasets corresponding to the subset
 343 $\mathcal{M}(\subset \mathcal{D}_{\text{orig}})$ and the new data $\mathcal{D}_{\text{new}} = \{(x_0, y_0)\}$ from CIFAR10, MNIST, and RESISC45, we
 344 perform the following operations:

345 **Step 1.** Randomly split a dataset into training, validation, and test datasets, which are denoted by
 346 $\mathcal{D}^{(\text{train})}$, $\mathcal{D}^{(\text{val})}$, and $\mathcal{D}^{(\text{test})}$, respectively.

347 **Step 2.** Exclude one common class from each of the three datasets. The excluded datasets are
 348 denoted by $\mathcal{D}_{\text{new}}^{(\text{train})}$, $\mathcal{D}_{\text{new}}^{(\text{val})}$, and $\mathcal{D}_{\text{new}}^{(\text{test})}$. The remaining datasets, which have 9 (or 44)
 349 classes, are denoted by $\mathcal{D}_{\text{orig}}^{(\text{train})}$, $\mathcal{D}_{\text{orig}}^{(\text{val})}$, and $\mathcal{D}_{\text{orig}}^{(\text{test})}$. Let y'_0 be the class label of samples
 350 included in the datasets $\mathcal{D}_{\text{new}}^{(\text{train})}$, $\mathcal{D}_{\text{new}}^{(\text{val})}$, and $\mathcal{D}_{\text{new}}^{(\text{test})}$.

351 **Step 3.** Change y'_0 to one of the remaining 9 (or 44) classes. Let y_0 be the class label after the
 352 change.

353 **Step 4.** Suppose that $\mathcal{D}_{\text{orig}}^{(\text{train})}$ represent the original data $\mathcal{D}_{\text{orig}}$ on $p_{\text{orig}}(x|y = k)$. These data are
 354 used for training to obtain the pre-trained model $f_{\hat{\theta}_{\text{orig}}}$. Note that $f_{\hat{\theta}_{\text{orig}}}$ is a 9-class (or
 355 44-class) classifier.

356 **Step 5.** Classify the samples in $\mathcal{D}_{\text{new}}^{(\text{train})}$ with $f_{\hat{\theta}_{\text{orig}}}$, and extract the misclassified samples. The
 357 dataset containing the extracted samples is denoted by $\mathcal{D}'_{\text{new}}$.

358 **Step 6.** Randomly select one sample from $\mathcal{D}'_{\text{new}}$. We assume that the selected sample is the given
 359 new data $\mathcal{D}_{\text{new}} = \{(x_0, y_0)\}$ on $p_{\text{new}}(x|y = y_0)$.

360 **Step 7.** Randomly select 1000 samples from the dataset $\mathcal{D}_{\text{orig}}^{(\text{train})}$. We assume that the dataset
 361 including the selected samples is the subset \mathcal{M} .

362 **Step 8.** Run Alg. 1 with the given \mathcal{M} , $f_{\hat{\theta}_{\text{orig}}}$, and \mathcal{D}_{new} .

363 While CIFAR10 and MNIST are originally 10-class classification problems, the above process reduces
 364 the number of classes by one to a 9-class classification problem. In the case of RESISC45, the number

²The transformations methods are described in Appendix D

Table 4: Settings of the original and new domains in y_0 .

dataset	Set 1		Set 2	
	y'_0	y_0	y'_0	y_0
CIFAR10	“truck”	“automobile”	“dog”	“cat”
MNIST	“9”	“8”	“1”	“0”
RESISC45	“airplane”	“airport”	“overpass”	“intersection”

of classes is reduced to 44. We set the new domain $p_{\text{new}}(\mathbf{x}|y = y_0)$ by changing y'_0 , which is not included in these 9 (or 44) classes, to one of these classes. Table 4 lists the settings of y'_0 and y_0 in our experiments with CIFAR10, MNIST, and RESISC45. We adopted two settings per dataset, which are named Set 1 and Set 2. The results of Set 1 of CIFAR10 and RESISC45 are shown in Section 4. Other results are presented in Appendix E. In the experiments, we randomly selected 10 samples from $\mathcal{D}'_{\text{new}}$ in Step 6 and performed ODIL 10 times, assuming that each sample was the given new sample (\mathbf{x}_0, y_0) . Over the 10 trials, we calculated the median and standard deviation of the accuracy.

Besides the new data \mathcal{D}_{new} and the original data $\mathcal{D}_{\text{orig}}$, the following datasets were used in the experiments.

- $\mathcal{D}_{\text{orig}}^{(\text{val})}$ is used for tuning the learning rate.
- $\mathcal{D}_{\text{new}}^{(\text{test})}$ is used to calculate the test accuracy on the new domain $p_{\text{new}}(\mathbf{x}|y = y_0)$.
- $\mathcal{D}_{\text{orig}}^{(\text{test})}$ is used to calculate the test accuracy on the original domain $p_{\text{orig}}(\mathbf{x}|y = k)$.

$\mathcal{D}_{\text{new}}^{(\text{val})}$ was not used to tune the learning rate because of the restriction that only one new data must be available for ODIL. Using the aforementioned datasets, we calculated the accuracy on the 9 (or 44) classes, excluding y'_0 .

D.2 Setup for Model Training

First Training. We followed Step 1 to Step 3 in Appendix D.1 and performed the first training with only the original data $\mathcal{D}_{\text{orig}}$ in Step 4. In this training, we used SGD [3] as the optimizer. The learning rate was scheduled by CosineAnnealing [22] from 0.1 to 0.0, and there were 200 epochs. The momentum was set to 0.9 and the weight decay was set to 0.0005. We used the same setup regardless of the new data size $|\mathcal{D}_{\text{new}}|$, DIL methods (CE, CE+EWC, or CE+GEM), and batch normalization (“baseline” or “proposed”). We used the same pre-trained model $f_{\hat{\theta}_{\text{orig}}}$ for all settings.

Second Training in the $|\mathcal{D}_{\text{new}}| = 1000$ setting. For the $|\mathcal{D}_{\text{new}}| = 1000$ setting, we skipped Steps 5 and 6 in Appendix D.1 and randomly selected 1000 samples from $\mathcal{D}_{\text{new}}^{(\text{train})}$. We assumed that the dataset including the selected 1000 samples is the new data \mathcal{D}_{new} . We extracted the subset \mathcal{M} as in Step 7, and $|\mathcal{M}| = 1000$. In Step 8, we updated $f_{\hat{\theta}_{\text{orig}}}$ by the Adam optimizer [15] with the fixed learning rate 10^{-5} and the number of iterations 100. The weight decay was set to 0.0001. In this setting, we did not adopt Alg. 1. We concatenated \mathcal{D}_{new} and \mathcal{M} prior to training iterations and performed training on the concatenated dataset using the standard approach.

Second Training in the $|\mathcal{D}_{\text{new}}| = 1$ setting. We followed Steps 5 and 6. We extracted the subset \mathcal{M} as in Step 7, and $|\mathcal{M}| = 1000$. In Step 8, we updated $f_{\hat{\theta}_{\text{orig}}}$ with the Adam optimizer. We calculated the accuracy of the validation set of the original domain for the fixed learning rate settings of 10^{-8} , 10^{-7} , 10^{-6} , 10^{-5} , 10^{-4} , 10^{-3} , 10^{-2} , and 10^{-1} . Then, we selected the settings that satisfied the constraint that the training iterations were terminated (“judge” returns “True”) within 100 iterations. Moreover, we adopted the learning rate with the highest accuracy among the selected settings. The weight decay was set to 0.0001. For the transformations from \mathbf{x}_0 to $\mathbf{z}_1, \dots, \mathbf{z}_M$ in Alg. 1, we used RandomRotation, RandomResizedCrop, RandomAffine, and RandomPerspective, which were implemented in Torchvision [23]. For “judge” in Alg. 1, we set $\delta = 0.99$. In the experiments whose results are shown in the main paper, we set $|\mathcal{B}| = |\mathcal{C}| = 32$ when running Alg. 1. In the following section, we show additional results for a different setting.

Note that we set the regularization term coefficient of EWC to 0.5 in the above two settings.

406 **E Additional Experiments**

407 We presented the results of Set 1 of CIFAR10 and RESISC45 in the main paper. In this section, we
408 also show the results of Set 2 of CIFAR10 and RESISC45 and those of Sets 1 and 2 of MNIST [19].
409 In addition, although in the main paper, we only showed the cases in the $|\mathcal{B}| = |\mathcal{C}| = 32$ setting, in
410 this section, we also show the cases in the $|\mathcal{B}| = 63, |\mathcal{C}| = 1$ setting.

411 Tables 5 and 6 show the results when we used CIFAR10. Some settings of the baseline degraded
412 the accuracy on both the new and original domains through ODIL. However, the proposed method
413 avoided the deterioration of accuracy and improved the accuracy on the new domain. A comparison
414 of the proposed method with the baseline shows that the accuracy on the new domain increased by
415 9%-38%, and the accuracy on the original domain increased by 0%-46%. The proposed method
416 improved the accuracy not only on the new domain but also on the original domain in the experiments
417 with CIFAR10. There is no example where the proposed method is worse than the baseline in Table 3.

418 Tables 7 and 8 show the results when we used MNIST. The baseline improved the accuracy on the
419 new domain through ODIL; however, compared to the baseline, the proposed method achieved higher
420 accuracy on the new domain. A comparison of the proposed method with the baseline shows that the
421 accuracy on the new domain increased by 0.8%-15% while maintaining the reduction in accuracy
422 on the original domain below 0.6%. Although there is a trade-off between the accuracy on the new
423 domain and that on the original domain, the decrease in the accuracy on the original domain is much
424 smaller than the increase in the accuracy on the new domain.

425 RESISC45 results show similar trends to MNIST. From Tables 9 and 10, a comparison of the proposed
426 method with the baseline shows that the accuracy on the new domain increased by 37%-57% while
427 maintaining the reduction in accuracy on the original domain below 7%. Similar to MNIST, the
428 decrease in the accuracy on the original domain is much smaller than the increase in the accuracy on
429 the new domain.

430 Note that the difference between the baseline and the proposed method lies only in the statistics of the
431 batch normalization layers. These results show that good accuracy can be achieved just by modifying
432 the batch normalization statistics. As explained in Section 3.1, the statistics shift unexpectedly in the
433 baseline, and avoiding this unexpected shift has a strong effect on the accuracy in ODIL. Moreover,
434 from the results of CE, we can predict that other DIL methods with cross-entropy loss show the same
435 trend.

Table 5: Test accuracy on the new and original domains in the baseline or proposed method (Set 1 of CIFAR10)

domain	$f_{\theta_{\text{orig}}}$	batch size	$f_{\theta_{\text{new}}}$					
			CE		CE+EWC		CE+GEM	
			baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.6940	$ \mathcal{B} = \mathcal{C} = 32$	0.5765 ± 0.2297	0.9595 ± 0.0430	0.5955 ± 0.2614	0.9420 ± 0.0586	0.7315 ± 0.2287	0.9485 ± 0.0417
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.6895 ± 0.1466	0.9260 ± 0.0425	0.7815 ± 0.1008	0.9375 ± 0.0456	0.8015 ± 0.0901	0.9220 ± 0.0452
p_{orig}	0.9550	$ \mathcal{B} = \mathcal{C} = 32$	0.4936 ± 0.1847	0.9335 ± 0.0440	0.4647 ± 0.2202	0.9332 ± 0.0192	0.9071 ± 0.2512	0.9404 ± 0.0250
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9413 ± 0.0054	0.9424 ± 0.0091	0.9394 ± 0.0070	0.9414 ± 0.0164	0.9399 ± 0.0062	0.9431 ± 0.0137

Table 6: Test accuracy on the new and original domains in the baseline or proposed method (Set 2 of CIFAR10)

domain	$f_{\theta_{\text{orig}}}$	batch size	$f_{\theta_{\text{new}}}$					
			CE		CE+EWC		CE+GEM	
			baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.7800	$ \mathcal{B} = \mathcal{C} = 32$	0.7595 ± 0.3077	0.9335 ± 0.0257	0.7665 ± 0.2721	0.9330 ± 0.0241	0.8081 ± 0.1010	0.9240 ± 0.0255
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.7740 ± 0.0771	0.9050 ± 0.0110	0.7305 ± 0.0526	0.9035 ± 0.0260	0.8075 ± 0.0663	0.8985 ± 0.0223
p_{orig}	0.9672	$ \mathcal{B} = \mathcal{C} = 32$	0.9363 ± 0.2308	0.9528 ± 0.0267	0.9360 ± 0.2210	0.9522 ± 0.0259	0.9556 ± 0.0313	0.9556 ± 0.0201
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9544 ± 0.0051	0.9580 ± 0.0023	0.9580 ± 0.0057	0.9587 ± 0.0204	0.9532 ± 0.0051	0.9584 ± 0.0041

Table 7: Test accuracy on the new and original domains in the baseline or proposed method (Set 1 of MNIST)

domain	$f_{\theta_{\text{orig}}}$	batch size	$f_{\theta_{\text{new}}}$					
			CE		CE+EWC		CE+GEM	
			baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.1021	$ \mathcal{B} = \mathcal{C} = 32$	0.7671 ± 0.1880	0.8826 ± 0.0715	0.8414 ± 0.1454	0.8930 ± 0.0611	0.7800 ± 0.1108	0.8677 ± 0.0688
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.7522 ± 0.2107	0.8251 ± 0.0832	0.7190 ± 0.1416	0.8746 ± 0.0778	0.7507 ± 0.1017	0.8593 ± 0.0863
p_{orig}	0.9962	$ \mathcal{B} = \mathcal{C} = 32$	0.9940 ± 0.0203	0.9904 ± 0.0046	0.9953 ± 0.0068	0.9903 ± 0.0050	0.9958 ± 0.0043	0.9909 ± 0.0035
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9959 ± 0.2806	0.9921 ± 0.0026	0.9959 ± 0.0013	0.9917 ± 0.0022	0.9957 ± 0.0019	0.9928 ± 0.0047

Table 8: Test accuracy on the new and original domains in the baseline or proposed method (Set 2 of MNIST)

domain	$f_{\theta_{\text{orig}}}$	batch size	$f_{\theta_{\text{new}}}$					
			CE		CE+EWC		CE+GEM	
			baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.0555	$ \mathcal{B} = \mathcal{C} = 32$	0.8639 ± 0.0763	0.9313 ± 0.1016	0.8656 ± 0.1085	0.9256 ± 0.0972	0.8819 ± 0.1008	0.9198 ± 0.0968
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9110 ± 0.0692	0.9194 ± 0.1089	0.8863 ± 0.0847	0.9167 ± 0.0783	0.8903 ± 0.0751	0.9022 ± 0.0676
p_{orig}	0.9951	$ \mathcal{B} = \mathcal{C} = 32$	0.9933 ± 0.0049	0.9931 ± 0.0056	0.9931 ± 0.0377	0.9932 ± 0.0055	0.9936 ± 0.0015	0.9937 ± 0.0026
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9946 ± 0.0005	0.9936 ± 0.0010	0.9941 ± 0.0004	0.9937 ± 0.0026	0.9941 ± 0.0011	0.9941 ± 0.0008

Table 9: Test accuracy on the new and original domains in the baseline or proposed method (Set 1 of RESISC45)

domain	$f_{\hat{\theta}_{\text{orig}}}$	batch size	$f_{\hat{\theta}_{\text{new}}}$					
			CE		CE+EWC		CE+GEM	
			baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.1071	$ \mathcal{B} = \mathcal{C} = 32$	0.4250 ± 0.1217	0.8250 ± 0.0862	0.3893 ± 0.1416	0.8214 ± 0.0868	0.3643 ± 0.1081	0.7786 ± 0.1020
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.2750 ± 0.1067	0.7179 ± 0.1433	0.2750 ± 0.0637	0.7179 ± 0.1433	0.3321 ± 0.1187	0.7393 ± 0.1559
p_{orig}	0.9580	$ \mathcal{B} = \mathcal{C} = 32$	0.9482 ± 0.0096	0.8912 ± 0.0712	0.9502 ± 0.0166	0.8879 ± 0.0728	0.9478 ± 0.0083	0.9294 ± 0.0549
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9515 ± 0.0117	0.9403 ± 0.0117	0.9539 ± 0.0023	0.9402 ± 0.0117	0.9546 ± 0.0231	0.9410 ± 0.0124

Table 10: Test accuracy on the new and original domains in the baseline or proposed method (Set 2 of RESISC45)

domain	$f_{\hat{\theta}_{\text{orig}}}$	batch size	$f_{\hat{\theta}_{\text{new}}}$					
			CE		CE+EWC		CE+GEM	
			baseline	proposed	baseline	proposed	baseline	proposed
p_{new}	0.0429	$ \mathcal{B} = \mathcal{C} = 32$	0.1500 ± 0.0584	0.7179 ± 0.1591	0.1393 ± 0.0776	0.7179 ± 0.1591	0.1429 ± 0.0616	0.5786 ± 0.1512
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.0429 ± 0.0450	0.5250 ± 0.1475	0.0500 ± 0.0421	0.5250 ± 0.1475	0.0929 ± 0.0463	0.4714 ± 0.1187
p_{orig}	0.9584	$ \mathcal{B} = \mathcal{C} = 32$	0.9505 ± 0.0105	0.9265 ± 0.0455	0.9481 ± 0.0118	0.9265 ± 0.0456	0.9499 ± 0.0108	0.9356 ± 0.0217
		$ \mathcal{B} = 63, \mathcal{C} = 1$	0.9535 ± 0.0025	0.9388 ± 0.0172	0.9543 ± 0.0022	0.9389 ± 0.0171	0.9563 ± 0.0006	0.9425 ± 0.0113