

Memory Efficient Data-Free Distillation for Continual Learning (Appendix)

Xiaorong Li^{a,1}, Shipeng Wang^{a,1}, Jian Sun^{a,*}, Zongben Xu^a

^a*School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an, Shaanxi, China*

1 **A. Related Work**

2 The existing continual learning methods can be roughly divided into five
3 categories: rehearsal-based, architecture-based, algorithm-based, regularization-
4 based and distillation-based approaches.

5 **Rehearsal-based approach** commonly relies on storing data of previous
6 tasks in a memory buffer, which is then replayed when learning a new task to
7 maintain network performance on previous tasks [1, 2]. Some works utilize
8 data from previous tasks to formulate a constrained optimization problem
9 to avoid the increases of training losses on previous tasks [3, 4, 5]. Gu et
10 al. [6] focus on the setting where the training data is accessed only once
11 during the training phase. They propose a sample selection strategy that
12 selects stored samples whose network parameter gradients are most inter-
13 ferred by new incoming samples. However, most of these works suffer from
14 class imbalance [7] between classes from previous and new tasks. Addition-

*Corresponding author.

Email addresses: lixiaorong@stu.xjtu.edu.cn (Xiaorong Li),
wangshipeng8128@stu.xjtu.edu.cn (Shipeng Wang), jiansun@xjtu.edu.cn
(Jian Sun), zbxu@xjtu.edu.cn (Zongben Xu)

¹Equal contribution.

15 ally, storing data of previous tasks introduces privacy and security concerns.
16 An alternative way is generating synthetic data of previous tasks based on
17 generative models [8, 9], whose performance depends on the quality of syn-
18 thetic data. Wei et al. [10] further apply knowledge distillation to mitigate
19 catastrophic forgetting of the generator for tackling a new problem named
20 incremental zero-shot learning, where the network incrementally learns the
21 knowledge of new classes and will be tested on all previously learned classes
22 and unseen classes. Nonetheless, generating high-quality synthetic data re-
23 mains a challenge.

24 **Architecture-based approach** concentrates on dynamically customiz-
25 ing the network’s structure for each task, such as expanding [11], pruning [12]
26 or masking [13] the neural connections, to alleviate catastrophic forgetting.
27 For example, Packnet [12] prunes the network and dedicates a specific sub-
28 part of the network for each task. Piggyback [13] learns masks of network
29 parameters to identify the task-specific part of the network. Learn-to-Grow
30 [11] employs architecture search to find the optimal structure for each task.
31 DRT [14] disentangles the latent features into class-disentangled and task-
32 disentangled features by two branches of networks. In contrast to them, our
33 method avoids modifying network architecture for each task, but designs a
34 novel regularizer in network training loss, which is easy to implement.

35 **Algorithm-based approach** concentrates on designing a network pa-
36 rameter update rule on the new task which alleviates the performance dete-
37 rioration on previous tasks. For example, OWM [15] projects the gradient
38 obtained at each training step into an orthogonal space of the space spanned
39 by input features of all network layers, where the projection matrix is up-

40 dated by recursive least squares. Adam-NSCL [16] considers the network
41 parameter update as the projection of gradient into the approximate null
42 space of uncentered feature covariance based on theoretical analysis. GPM
43 [17] forces the network parameter update to lie in the orthogonal space of
44 the space spanned by input features, which is dependent on storing the ba-
45 sis of the complement to orthogonal space obtained by employing singular
46 value decomposition on partial input features of previous tasks. However, to
47 obtain the network parameter update in each step, these works [15, 16, 17]
48 need to store the projection matrix or basis of subspace, where the memory
49 usages of them are larger than ours.

50 **Regularization-based approach** penalizes the variations of network
51 parameters to preserve the performance on previous tasks, where each net-
52 work parameter is associated with an importance weight. Previous works
53 focus on designing the importance weight in different ways [18, 19, 20, 21].
54 For example, MAS [19] aims to consider the sensitivity of the output function,
55 and finally implements this by estimating the sensitivity of the norm of pre-
56 dicted output w.r.t. the parameters to measure the importance weight. In-
57 spired by [19], MUC-MAS [21] proposes to integrate an ensemble of auxiliary
58 classifiers to estimate the importance weight, where the auxiliary classifiers
59 are trained on out-of-distribution data irrelevant to the current task. How-
60 ever, these works implicitly assume that the network parameters are inde-
61 pendent without considering the impact of the correlation among parameters
62 on the network performance, while our method takes this into consideration.

63 **Distillation-based approach** is inspired by knowledge distillation [22].
64 To preserve the performance on previous tasks in continual learning, a distil-

65 lation term is utilized to penalize the variations of outputs between teacher
66 and student networks, where the teacher and student networks are respec-
67 tively set as the network learned from the previous tasks and the network
68 being trained on the current task. Ideally, the data associated with the dis-
69 tillation term should be the original data of previous tasks. However, the
70 full datasets of previous tasks are inaccessible in continual learning, thus
71 existing works mainly focus on how to substitute the datasets of previous
72 tasks [23, 24, 25]. For example, methods in [24, 26] employ the training data
73 of the current task for knowledge distillation. iCaRL [23] stores a few sam-
74 ples of previous tasks as a coreset for knowledge distillation. Besides using
75 coreset, the method in [25] additionally leverages a large stream of unlabeled
76 data in the wild which is assumed to be available at any time for knowledge
77 distillation. Commonly, the dataset they adopted for knowledge distillation
78 may fail to reflect the knowledge of the full original datasets of previous
79 tasks, as a result of which, these methods may not well distill full knowledge
80 of the previous tasks. Moreover, the teacher networks need to be saved in
81 their methods. Different from them, we make use of the whole knowledge
82 of datasets of previous tasks encoded by gradients for knowledge distillation,
83 and we do not need to store teacher networks.

84 **B. Proof of Theorem 1**

85 In this section, we begin by introducing the notation used in the proof,
86 followed by presenting Lemma 1 which forms the basis of proof for Theorem 1.
87 Then we recall Theorem 1 and provide its proof.

88 **Notation.** To mitigate catastrophic forgetting, we employ a distillation

89 loss penalizing the output variations between teacher and student networks
90 on previous task data X_i of task \mathcal{T}_i , i.e., $\min_w \|f(X_i, w) - f(X_i, w_i^*)\|_2^2$, where
91 $f(\cdot, w_i^*)$ is the teacher network and $f(\cdot, w)$ is the student network. Since the
92 previous task data X_i are not available when learning task \mathcal{T}_t , we approximate
93 $f(X_i, w)$ with its first-order Taylor expansion at $w = w_i^*$, i.e., $f(X_i, w) \approx$
94 $f(X_i, w_i^*) + G_i^\top (w - w_i^*)$, where $G_i = \frac{\partial f(X_i, w_i^*)}{\partial w_i^*} \in \mathbb{R}^{|w| \times C}$ with $|w|$ and C as the
95 dimension of w and the number of classes respectively. The gradients G_i are
96 further compressed and recovered for memory efficiency, and the recovered
97 gradients are denoted as \tilde{G}_i .

98 For clarity in notation, we omit the subscript i denoting the task index of
99 gradients G_i . The gradient matrix $G \in \mathbb{R}^{|w| \times C}$ on task \mathcal{T}_i is the concatenation
100 of layer-wise gradient matrix $G^l \in \mathbb{R}^{|w^l| \times C}$ such that $|w| = \sum_{l=1}^L |w^l|$ with
101 L denoting the number of layers of the network. We additionally introduce
102 a notation $\bar{G} \in \mathbb{R}^{|w| \times C}$ as the concatenation of matrix $EXP(\bar{G}^l) \in \mathbb{R}^{|w^l| \times C}$
103 ($l = 1, \dots, L$) for proving Lemma 1, where $EXP(\cdot)$ denotes the operation
104 that copies each element of the input object to the elements corresponding to
105 the convolutional kernel. For the l -th convolutional layer of the network with
106 n_{out}^l 3D convolutional kernels of size $n_{in}^l \times k^l \times k^l$, the gradients of the l -th
107 layer $G^l \in \mathbb{R}^{|w^l| \times C}$ are gradients of C elements of network output w.r.t. the
108 j -th convolutional kernel, where $G^l = [\mathbf{g}_1^\top, \dots, \mathbf{g}_{n_{out}^l}^\top]^\top$, $\mathbf{g}_j \in \mathbb{R}^{(n_{in}^l \times k^l \times k^l) \times C}$
109 ($j = 1, \dots, n_{out}^l$) and $|w^l| = n_{out}^l \times n_{in}^l \times k^l \times k^l$.

110 Recall the proposed compression approach, we first average the gradients
111 of each convolutional kernel with size $n_{in}^l \times k^l \times k^l$, which results in a compact
112 matrix $\bar{G}^l \in \mathbb{R}^{n_{out}^l \times C}$. Specifically, we obtain \bar{G}^l by $\bar{G}^l = [\bar{\mathbf{g}}_1^\top, \dots, \bar{\mathbf{g}}_{n_{out}^l}^\top]^\top$,
113 where $\bar{\mathbf{g}}_j = \frac{1}{n_{in}^l \times k^l \times k^l} (\mathbf{1}^\top \mathbf{g}_j)^\top \in \mathbb{R}^C$ and $\mathbf{1} \in \mathbb{R}^{n_{in}^l \times k^l \times k^l}$ is a column vector

114 with its all elements as 1. We then apply SVD to the compact matrix $\bar{G}^l \in$
 115 $\mathbb{R}^{n_{out}^l \times C}$ and obtain its SVD decompositions, i.e., $U^l \in \mathbb{R}^{n_{out}^l \times r}$, $V^l \in \mathbb{R}^{C \times r}$
 116 and $\Lambda^l \in \mathbb{R}^{r \times r}$ ($r < C$), which are used to approximate the compact matrix
 117 \bar{G}^l . The resulting approximate compact matrix is denoted as $\tilde{G}^l \in \mathbb{R}^{n_{out}^l \times C}$,
 118 which is obtained by $\tilde{G}^l = U^l \Lambda^l V^{l\top}$.

119 We now introduce Lemma 1 which analyzes the approximation error be-
 120 tween the gradients G_i and the recovered gradients \tilde{G}_i . Lemma 1 forms the
 121 basis of the proof of Theorem 1.

122 **Lemma 1.** *Given the gradients G of the network output $f(X_i, w_i^*)$ w.r.t.*
 123 *the learned network parameters w_i^* , i.e., $G = \frac{\partial f(X_i, w_i^*)}{\partial w_i^*}$, and the recovered*
 124 *gradients \tilde{G} , then we have the following bound on the error between G and*
 125 *\tilde{G} :*

$$\|G - \tilde{G}\|_2^2 \leq \sum_{l=1}^L \left[\left(\sum_{j=1}^{n_{out}^l} \sum_{c=1}^C (n_{in}^l \times k^l \times k^l) \times \sigma_{j,c}^2 \right) + (n_{in}^l \times k^l \times k^l) \left(\sum_{j=r+1}^{C-r} s_j^2 \right) \right], \quad (1)$$

126 where $\sigma_{j,c}^2$ is the variance of the gradients of the c -th element of the net-
 127 work output $f(X_i, w_i^*)$ w.r.t. the j -th convolutional kernel, r is the number
 128 of selected singular values, s_j represents the j -th singular value that is not
 129 selected.

Proof.

$$\begin{aligned} \|G - \tilde{G}\|_2^2 &= \|G - \bar{G} + \bar{G} - \tilde{G}\|_2^2 \\ &\leq \|G - \bar{G}\|_2^2 + \|\bar{G} - \tilde{G}\|_2^2 \end{aligned}$$

$$\begin{aligned}
&= \sum_{l=1}^L \left[\|G^l - EXP(\bar{G}^l)\|_2^2 + \|EXP(\bar{G}^l) - EXP(\tilde{G}^l)\|_2^2 \right] \\
&= \sum_{l=1}^L \left[\left(\sum_{j=1}^{n_{out}^l} \|\mathbf{g}_j - EXP(\bar{\mathbf{g}}_j)\|_2^2 \right) + (n_{in}^l \times k^l \times k^l) \|\bar{G}^l - \tilde{G}^l\|_2^2 \right] \\
&= \sum_{l=1}^L \left[\left(\sum_{j=1}^{n_{out}^l} \sum_{c=1}^C (n_{in}^l \times k^l \times k^l) \times \sigma_{j,c}^2 \right) + (n_{in}^l \times k^l \times k^l) \left(\sum_{j=r+1}^{C-r} s_j^2 \right) \right]. \quad (2)
\end{aligned}$$

130

□

131 We now recall Theorem 1 and provide its proof. The proof of Theorem 1
132 is based on the conclusion of Lemma 1.

133 **Theorem 1** (Bound on Relative Approximation Error). *Given the network*
134 *output $f(X_i, w)$ on previous task data X_i , we approximate it by $f(X_i, w) \approx$*
135 *$f(X_i, w_i^*) + \tilde{G}_i^\top (w - w_i^*)$, where \tilde{G}_i is an approximation of the gradients G_i*
136 *of $f(X_i, w_i^*)$ w.r.t. w_i^* , then we have the following bound on relative approx-*
137 *imation error between $f(X_i, w)$ and $f(X_i, w_i^*) + \tilde{G}_i^\top (w - w_i^*)$:*

$$\frac{\|f(X_i, w) - f(X_i, w_i^*) - \tilde{G}_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2} \leq e_1 + e_2 \quad (3)$$

with

$$e_1 = \frac{\|f(X_i, w) - f(X_i, w_i^*) - G_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2}, \quad (4)$$

$$\begin{aligned}
e_2 &= \frac{\|(w - w_i^*)\|_2^2 \alpha^2 \sum_{l=1}^L \sum_{j=1}^{n_{out}^l} \sum_{c=1}^C (n_{in}^l \times k^l \times k^l) \times \sigma_{j,c}^2}{\|f(X_i, w_i^*)\|_2^2} \\
&+ \frac{\|(w - w_i^*)\|_2^2 \alpha^2 \sum_{l=1}^L (n_{in}^l \times k^l \times k^l) (\sum_{j=r+1}^{C-r} s_j^2)}{\|f(X_i, w_i^*)\|_2^2}, \quad (5)
\end{aligned}$$

138 where $\alpha = \max(\{\cos\theta_1, \dots, \cos\theta_C\})$ is the maximum cosine value with θ_c
139 ($c = 1, \dots, C$) as the angle between the c -th column of $G_i - \tilde{G}_i$ and $w - w_i^*$,
140 $\sigma_{j,c}^2$ is the variance of the gradients of the c -th element of the network output
141 $f(X_i, w_i^*)$ w.r.t. the j -th convolutional kernel, r is the number of selected
142 singular values, s_j represents the j -th singular value that is not selected.

143 According to Theorem 1, the bound on relative approximation error is
144 determined by e_1 and e_2 . e_1 measures the relative truncation error of the
145 Taylor expansion, while e_2 quantifies the approximation error of compressing
146 the gradients. Specifically, the first term and second term of e_2 are influenced
147 by our average operation and SVD respectively.

148 *Proof.* Considering the relative approximation error $\frac{\|f(X_i, w) - f(X_i, w_i^*) - \tilde{G}_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2}$,
149 we have

$$\begin{aligned}
& \frac{\|f(X_i, w) - f(X_i, w_i^*) - \tilde{G}_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2} \\
&= \frac{\|f(X_i, w) - f(X_i, w_i^*) - G_i^\top (w - w_i^*) + G_i^\top (w - w_i^*) - \tilde{G}_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2} \\
&\leq \frac{\|f(X_i, w) - f(X_i, w_i^*) - G_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2} + \frac{\|(G_i - \tilde{G}_i)^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2} \\
&= e_1 + \frac{\sum_{c=1}^C [G_i - \tilde{G}_i]_c^2 (w - w_i^*)^2 \cos^2 \theta_c}{\|f(X_i, w_i^*)\|_2^2} \\
&\leq e_1 + \frac{\sum_{c=1}^C [G_i - \tilde{G}_i]_c^2 (w - w_i^*)^2 \alpha^2}{\|f(X_i, w_i^*)\|_2^2}
\end{aligned}$$

$$= e_1 + \frac{\|(G_i - \tilde{G}_i)\|_2^2 \|(w - w_i^*)\|_2^2 \alpha^2}{\|f(X_i, w_i^*)\|_2^2}, \quad (6)$$

150 where $[\cdot]_c$ denotes the c -th column of the matrix, $\alpha = \max(\{\cos\theta_1, \dots, \cos\theta_C\})$
 151 is the maximum cosine value and θ_c ($c = 1, \dots, C$) is the angle between the
 152 c -th column of $G_i - \tilde{G}_i$ and $w - w_i^*$.

153 By substituting inequality (2) in the second term of r.h.s. of inequality
 154 (6), we have

$$\frac{\|(G_i - \tilde{G}_i)\|_2^2 \|(w - w_i^*)\|_2^2 \alpha^2}{\|f(X_i, w_i^*)\|_2^2} \leq e_2. \quad (7)$$

Combining inequalities (6) and (7), we can conclude that

$$\frac{\|f(X_i, w) - f(X_i, w_i^*) - \tilde{G}_i^\top (w - w_i^*)\|_2^2}{\|f(X_i, w_i^*)\|_2^2} \leq e_1 + e_2.$$

155

□

156 C. Details on Compressing Gradients of BN layer

157 For the batch normalization (BN) layer with n feature channels, each
 158 channel is equipped with an affine transformation with two parameters, i.e.,
 159 scaling weight and bias. Therefore, the gradients of scaling weight and bias
 160 at BN layer are both with the size of $n \times C$, where C is the number of classes.
 161 We respectively compress the gradients of scaling weight and bias by SVD
 162 without using average operation.

163 D. Algorithm

Algorithm 1 Data-free distillation for continual learning.

Inputs: datasets $\{X_t, Y_t\}$ for task $\mathcal{T}_t \in \{\mathcal{T}_1, \mathcal{T}_2, \dots\}$, network $f(\cdot, w)$ of depth L , hyperparameter λ .

Output: learned network $f(\cdot, w^*)$.

```

1: Memory buffer  $\mathcal{M} = \emptyset$ .
2: for task  $\mathcal{T}_t \in \{\mathcal{T}_1, \mathcal{T}_2, \dots\}$  do
3:   if  $t = 1$  then
4:     Initialize  $w$  randomly.
5:      $w_t^* \leftarrow \arg \min_w \mathcal{L}(\hat{Y}_t(w), Y_t)$ .
6:   else
7:     # Reconstruction stage as illustrated in Sec.3.2
8:     Take  $\{U_i^l, \Lambda_i^l, V_i^l\}_{l=1}^L$  from  $\mathcal{M}$  for  $i = 1, \dots, t-1$ .
9:      $G_i = \text{Reconstruct}(U_i^1, \Lambda_i^1, V_i^1, \dots, U_i^L, \Lambda_i^L, V_i^L)$  for  $i = 1, \dots, t-1$ .
10:    Initialize  $w$  with  $w_{t-1}^*$ .
11:     $w_t^* \leftarrow \arg \min_w \mathcal{L}(\hat{Y}_t(w), Y_t) + \lambda \sum_{i=1}^{t-1} (w^\top G_i G_i^\top w - 2w^\top G_i G_i^\top w_i^*)$ .
12:   end if
13:   Obtain the gradient  $G_t = \frac{\partial f(X_t, w_t^*)}{\partial w_t^*}$ .
14:   # Compression stage as illustrated in Sec. 3.2
15:    $U_t^l, \Lambda_t^l, V_t^l = \text{Compress}(G_t^l)$ .
16:    $\mathcal{M} = \mathcal{M} \cup \{U_t^l, \Lambda_t^l, V_t^l\}_{l=1}^L$ .
17: end for

```

164 In this Appendix, we summarize the algorithm for the proposed DFD in
165 Alg. 1. Given a task sequence $\{\mathcal{T}_1, \mathcal{T}_2, \dots\}$, The basic pipeline of DFD is it-
166 eratively optimizing problem (5) of the manuscript to acquire new knowledge
167 while reserving the performance on previous tasks.

168 Concretely, to train the network acquiring knowledge of the first task \mathcal{T}_1 ,
169 we minimize the cross-entropy loss to obtain the learned network parameter
170 w_1^* as illustrated in line 9. When training the network on task \mathcal{T}_t ($t > 1$), we
171 first recover the gradients as described in lines 13 to 14. Then we optimize
172 problem (5) of the manuscript as described in line 18. After obtaining the
173 learned network parameters for task \mathcal{T}_t , we update the memory buffer as

174 described in line 26 following the compression stage in line 25.

175 E. Details on the Hyperparameters

176 The setting of hyperparameters of our method is listed in Table 1, where
177 “lr”, “wd” and “bs” denote the initial learning rate, weight decay and batch
178 size respectively. The network is trained for 80 epochs in total for all experi-
179 ments. The initial learning rate decays at epochs 30 and 60 with a multiplier
180 of 0.5 for all experiments.

Table 1: The settings of hyperparameters of our method on all experiments.

Experiments	lr	wd	bs	λ
10-split CIFAR-100	1×10^{-3}	5×10^{-5}	32	20
20-split CIFAR-100	1×10^{-4}	5×10^{-4}	32	10
25-split TinyImageNet	5×10^{-5}	1×10^{-5}	16	50
20-split miniImageNet	5×10^{-5}	5×10^{-5}	32	20
10-split SubImageNet	1×10^{-4}	5×10^{-5}	16	200

181 F. Effect of r on the Bound on Relative Approximation Error

182 We now study the impact of the number of stored top singular values,
183 denoted by r , on the bound of the relative approximation error. According
184 to Theorem 1 of the paper, r is irrelevant to e_1 but affects the value of e_2
185 which captures the approximation error of compressing the gradients. To
186 investigate the effect of r on the bound of error, we report the values of the
187 bound for different values of r on 10-split CIFAR-100 in Table 2, where we
188 set $w = w_{10}^*$ and $w_i^* = w_1^*$. As shown in Table 2, we observe that the value of
189 e_1 is dominant and the values of e_2 remain stable across different choices of r .

Table 2: The effect of r on the bound of relative approximation error in 10-split CIFAR-100.

The number of stored top singular values	e_1	e_2	α in e_2	$\ (w_{10}^* - w_1^*)\ _2^2 \alpha^2$ in e_2
$r = 1$	0.32	8.9×10^{-5}	4.9×10^{-3}	1.3×10^{-3}
$r = 2$	0.32	3.7×10^{-5}	6.2×10^{-3}	7.1×10^{-4}
$r = 3$	0.32	1.2×10^{-5}	4.6×10^{-3}	2.8×10^{-4}
$r = 4$	0.32	1.8×10^{-5}	6.7×10^{-3}	5.3×10^{-4}
$r = 5$	0.32	1.4×10^{-5}	8.0×10^{-3}	5.2×10^{-4}
$r = 6$	0.32	8.0×10^{-6}	6.8×10^{-3}	3.5×10^{-4}
$r = 7$	0.32	8.3×10^{-6}	8.1×10^{-3}	4.5×10^{-4}
$r = 8$	0.32	9.1×10^{-6}	9.5×10^{-3}	6.0×10^{-4}

190 This finding is consistent with our observation that the network performance
 191 is robust to the changes of r , as demonstrated in Figure 4 of the paper.

192 References

- 193 [1] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A
 194 simple approach that questions our progress in continual learning. In
 195 *Eur. Conf. Comput. Vis.*, pages 524–540, 2020.
- 196 [2] Chen Zhuang, Shaoli Huang, Gong Cheng, and Jifeng Ning. Multi-
 197 criteria selection of rehearsal samples for continual learning. *Pattern*
 198 *Recognition*, 132:108907, 2022.
- 199 [3] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory
 200 for continual learning. In *Adv. Neural Inform. Process. Syst.*, pages
 201 6467–6476, 2017.

- 202 [4] Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mo-
203 hamed Elhoseiny. Efficient lifelong learning with a-GEM. In *Int. Conf.*
204 *Learn. Represent.*, 2019.
- 205 [5] Yunhui Guo, Mingrui Liu, Tianbao Yang, and Tajana Rosing. Improved
206 schemes for episodic memory-based lifelong learning. In *Adv. Neural*
207 *Inform. Process. Syst.*, pages 1023–1035, 2020.
- 208 [6] Yanan Gu, Xu Yang, Kun Wei, and Cheng Deng. Not just selection,
209 but exploration: Online class-incremental continual learning via dual
210 view consistency. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages
211 7442–7451, 2022.
- 212 [7] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. Main-
213 taining discrimination and fairness in class incremental learning. In
214 *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 13208–13217, 2020.
- 215 [8] Huaiyu Li, Weiming Dong, and Bao-Gang Hu. Incremental concept
216 learning via online generative memory recall. *IEEE Trans. Neural Net-*
217 *works Learn. Syst.*, 32(7):3206–3216, 2021.
- 218 [9] Qicheng Lao, Mehrzad Mortazavi, Marzieh Tahaei, Francis Dutil,
219 Thomas Fevens, and Mohammad Havaei. Focl: Feature-oriented con-
220 tinual learning for generative models. *Pattern Recognition*, 120:108127,
221 2021.
- 222 [10] Kun Wei, Cheng Deng, Xu Yang, and Dacheng Tao. Incremental zero-
223 shot learning. *IEEE Trans. Cyber.*, 52(12):13788–13799, 2022.

- 224 [11] Xilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong.
225 Learn to grow: A continual structure learning framework for overcoming
226 catastrophic forgetting. In *Int. Conf. Mach. Learn.*, pages 3925–3934,
227 2019.
- 228 [12] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks
229 to a single network by iterative pruning. In *IEEE Conf. Comput. Vis.*
230 *Pattern Recog.*, pages 7765–7773, 2018.
- 231 [13] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapt-
232 ing a single network to multiple tasks by learning to mask weights. In
233 *Eur. Conf. Comput. Vis.*, pages 67–82, 2018.
- 234 [14] Kun Wei, Da Chen, Yuhong Li, Xu Yang, Cheng Deng, and Dacheng
235 Tao. Incremental embedding learning with disentangled representation
236 translation. *IEEE Trans. Neural Networks Learn. Syst.*, pages 1–13,
237 2022.
- 238 [15] Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. Continual learning
239 of context-dependent processing in neural networks. *Nat. Mach. Intell.*,
240 1(8):364–372, 2019.
- 241 [16] Shipeng Wang, Xiaorong Li, Jian Sun, and Zongben Xu. Training net-
242 works in null space of feature covariance for continual learning. In *IEEE*
243 *Conf. Comput. Vis. Pattern Recog.*, pages 184–193, 2021.
- 244 [17] Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection mem-
245 ory for continual learning. In *Int. Conf. Learn. Represent.*, 2021.

- 246 [18] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guil-
247 laume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ram-
248 malho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic
249 forgetting in neural networks. *Proc. Natl. Acad. Sci. USA*, 114(13):3521–
250 3526, 2017.
- 251 [19] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus
252 Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning
253 what (not) to forget. In *Eur. Conf. Comput. Vis.*, pages 139–154, 2018.
- 254 [20] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning
255 through synaptic intelligence. In *Int. Conf. Mach. Learn.*, pages 3987–
256 3995, 2017.
- 257 [21] Yu Liu, Sarah Parisot, Gregory Slabaugh, Xu Jia, Ales Leonardis, and
258 Tinne Tuytelaars. More classifiers, less forgetting: A generic multi-
259 classifier paradigm for incremental learning. In *Eur. Conf. Comput.*
260 *Vis.*, pages 699–716, 2020.
- 261 [22] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge
262 in a neural network. In *Adv. Neural Inform. Process. Syst. Worksh.*,
263 2015.
- 264 [23] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and
265 Christoph H Lampert. icarl: Incremental classifier and representation
266 learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2001–2010,
267 2017.

- 268 [24] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE*
269 *Trans. Pattern Anal. Mach. Intell.*, 40(12):2935–2947, 2017.
- 270 [25] Kibok Lee, Kimin Lee, Jinwoo Shin, and Honglak Lee. Overcoming
271 catastrophic forgetting with unlabeled data in the wild. In *Int. Conf.*
272 *Comput. Vis.*, pages 312–321, 2019.
- 273 [26] Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu.
274 Prototype augmentation and self-supervision for incremental learning.
275 In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5871–5880, 2021.