

# Improving Plasticity in Online Continual Learning via Collaborative Learning

Anonymous CVPR submission

Paper ID 7545

## Abstract

001 *Online Continual Learning (CL) solves the problem of*  
002 *learning the ever-emerging new classification tasks from a*  
003 *continuous data stream. Unlike its offline counterpart, in*  
004 *online CL, the training data can only be seen once. Most*  
005 *existing online CL research regards catastrophic forgetting*  
006 *(i.e., model stability) as almost the only challenge. In this*  
007 *paper, we argue that the model’s capability to acquire new*  
008 *knowledge (i.e., model plasticity) is another challenge in*  
009 *online CL. While replay-based strategies have been shown*  
010 *to be effective in alleviating catastrophic forgetting, there is*  
011 *a notable gap in research attention toward improving model*  
012 *plasticity. To this end, we propose Collaborative Contin-*  
013 *ual Learning (CCL), a collaborative learning based strat-*  
014 *egy to improve the model’s capability in acquiring new con-*  
015 *cepts. Additionally, we introduce Distillation Chain (DC),*  
016 *a novel collaborative learning scheme to boost the training*  
017 *of the models. We adapted CCL-DC to existing represen-*  
018 *tative online CL works. Extensive experiments demonstrate*  
019 *that even if the learners are well-trained with state-of-the-*  
020 *art online CL methods, our strategy can still improve model*  
021 *plasticity dramatically, and thereby improve the overall per-*  
022 *formance by a large margin. The source code is included in*  
023 *the supplementary material and will be publicly available*  
024 *upon acceptance.*

## 025 1. Introduction

026 Continual Learning (CL) [11, 14, 35, 47] aims to learn a  
027 sequence of tasks incrementally and encourage the neural  
028 network to gain more performance on the tasks at hand,  
029 without forgetting heretofore learned knowledge. CL can  
030 be done in two different manners [4, 47]: *offline* continual  
031 learning and *online* continual learning. In offline CL, the  
032 learner can have infinite access to all the training data of  
033 the current task it trains on and may go through the data  
034 for any epoch. Contrary to offline CL, in online CL, the  
035 training data for each task also comes continually in a data  
036 stream, and the learner can only see the training data once.  
037 Apart from the learning manner, there are also three dif-

ferent CL scenarios [25, 32, 45]: Task-incremental Learn- 038  
ing (TIL), Domain-incremental Learning (DIL), and Class- 039  
incremental learning (CIL). In this paper, we focus on the 040  
CIL setting in online CL. 041

Various online CL methods [6, 7, 21, 22, 38, 48] have 042  
been proposed to help the models learn continually on one- 043  
epoch data stream, with alleviated forgetting. Among them, 044  
replay-based methods have shown remarkable success, and 045  
current state-of-the-art methods rely heavily on memory re- 046  
play to mitigate catastrophic forgetting [19, 33]. However, 047  
while most existing online CL research almost only focuses 048  
on improving model stability (i.e., alleviating catastrophic 049  
forgetting) in pursuit of better overall accuracy, the impor- 050  
tance of model plasticity (i.e., the capability to acquire new 051  
knowledge) is greatly overlooked. Contrary to offline CL, 052  
where it is possible to gain high plasticity by iterating sev- 053  
eral epochs on the current task before proceeding to the sub- 054  
sequent task, the plasticity in online CL is more important 055  
because the training data can only be seen once. As shown 056  
in Fig. 1, compared to learning without memory replay, the 057  
replay-based methods implicitly alleviate the low plasticity 058  
issue to some extent. Also, it is possible to improve the plas- 059  
ticity with multiple updates trick on incoming samples [3]. 060  
However, the combination of memory replay and multiple 061  
updates does not bridge the plasticity gap, and multiple up- 062  
dates trick will also lead to higher catastrophic forgetting. 063  
Overall, the plasticity gap hinders the performance of on- 064  
line CL methods. 065

In this paper, we claim that besides stability, the model’s 066  
ability to acquire new knowledge (i.e., model plasticity) 067  
is also vital in order to have a good overall accuracy. To 068  
shed light on how model plasticity and stability will impact 069  
the overall performance, we propose a quantitative link be- 070  
tween plasticity, stability, and final accuracy, showing that 071  
both plasticity and stability play crucial roles in the overall 072  
performance. 073

Guided by the quantitative relationship, we focus our- 074  
selves on the former-overlooked plasticity perspective. In- 075  
spired by the ability of collaborative learning to accelerate 076  
the convergence in non-continual scenarios [5], we incor- 077  
porated collaborative learning in online CL and observed 078

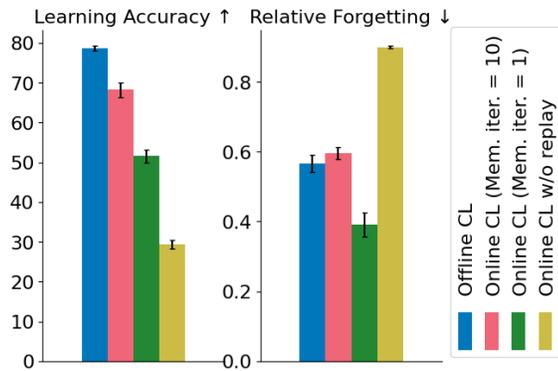


Figure 1. The plasticity (learning accuracy) and stability (relative forgetting, our metric proposed in Sec. 3) comparison of ER under different settings on CIFAR-100. For experiments with memory replay, the size of the memory buffer is set to 2,000. We can witness a plasticity gap between offline CL and online CL, even with memory replay and multiple update trick (memory iteration  $> 1$ ).

a similar phenomenon. To this end, we propose Collaborative Continual Learning with Distillation Chain (CCL-DC), a collaborative learning scheme that can be adapted to existing online CL methods. CCL-DC comprises two key components: Collaborative Continual Learning (CCL) and Distillation Chain (DC).

CCL involves two peer continual learners to learn from the data stream simultaneously in a peer teaching manner, and it enables us to have more parallelism in optimization and provides more maneuverability to the continual learners. To the best of our knowledge, CCL is the first to involve collaborative learning techniques in online CL research. Moreover, to fully exploit the potential of collaborative learning in online CL scenarios, we proposed DC, an entropy regularization based optimization strategy explicitly designed for online CL.

The main contribution of this paper can be summarized as follows.

1. We identify two important challenges in training online continual learners: plasticity and stability. Moreover, we propose a quantitative link between plasticity, stability, and final performance. Based on this, we find that plasticity is an important obstacle in online CL, which was greatly overlooked in the previous research;
2. To overcome the plasticity issue, we propose CCL-DC, a collaborative learning based strategy that can be seamlessly integrated into the existing methods and improve their performance by enhancing plasticity;
3. Extensive experiments show that CCL-DC can enhance the performance of existing methods by a large margin.

## 2. Related Work

**Continual Learning.** Continual Learning methods can be classified into three different categories: regularization-

based methods, parameter-isolation-based methods, and replay-based methods. Regularization-based methods [2, 9, 26, 29, 50] add extra regularization terms to balance the old and new tasks. Parameter-isolation-based methods [1, 17, 39–41] solve the problem explicitly by dynamically allocating task-specific parameters. Replay-based methods [6, 7, 10, 15, 21, 22, 34, 38, 48] maintain a small memory buffer that stores a few old training samples.

Among these methods, replay-based strategies have gained huge success due to their impressive performance and simplicity. ER [38] is the fundamental replay-based method that leverages Cross-Entropy loss for classification and a random replay buffer. DER++ [6] stores the logits in the memory buffer and extends ER with the distillation of old stored logits. ER-ACE [7] extends ER with Asymmetric Cross-Entropy loss for classification to suppress the drift of old class representations. OCM [21] leverages a replay-based strategy by maximizing the mutual information between old and new class representations. GSA [22] solves cross-task class discrimination with replay-based strategy and Gradient Self Adaption. OnPro [48] uses online prototype learning to address shortcut learning and alleviate catastrophic forgetting.

These replay-based methods propose different strategies for alleviating catastrophic forgetting and improving the model stability. However, the importance of the model plasticity is greatly neglected in their research, despite their success in terms of final performance. In our work, these methods serve as the baselines and we adapted our strategy to these baselines to show the efficiency of our proposed approach.

**Collaborative Learning.** Collaborative learning [5, 20, 42, 51, 52] orients from online knowledge distillation (KD). Different from the conventional KD methods, online KD trains a cohort of deep networks from scratch in a peer-teaching manner. During the training process, the model imitates their peers and guides the training of other models simultaneously. DML [51] suggests peer student models learn from each other through the logit distillation between the probability distributions. Codistillation [5] is similar to DML and suggests the ensemble of peer networks can further improve the performance. More importantly, Codistillation shows that online KD can help the model converge faster on non-continual scenarios.

Despite the success of collaborative learning in non-continual scenarios, due to the lack of focus on plasticity, the research on collaborative learning in CL is still limited. To the best of our knowledge, there is no existing research using the collaborative learning technique to boost the training of online CL. Moreover, in our work, we propose DC, an entropy regularization based optimization strategy, which is designed to exploit the full potential of collaborative learning in online CL scenarios.

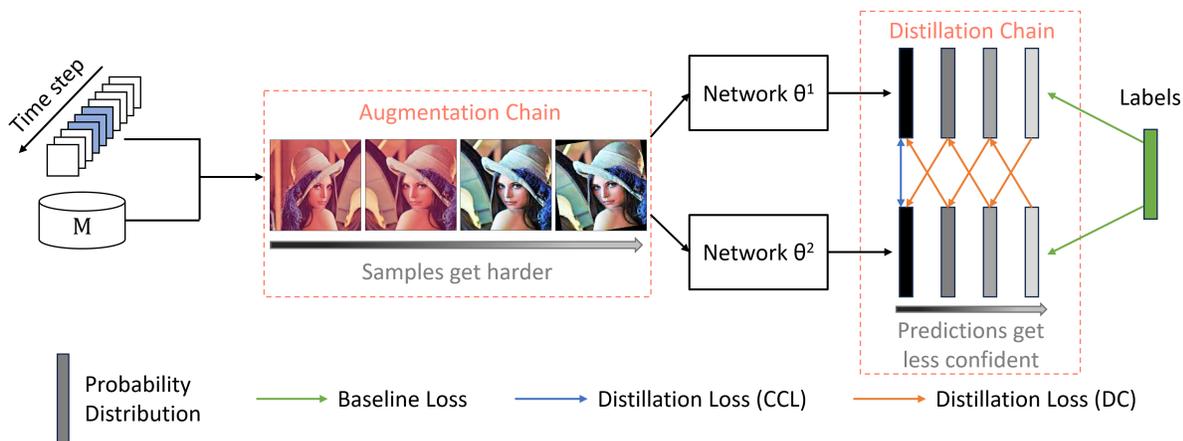


Figure 2. Overview of the proposed CCL-DC framework applied to a baseline online CL method. The proposed CCL-DC framework has two main components. The first one is CCL, which involves two peer continual learners that simultaneously learn from the data stream in a peer teaching manner. The second component is DC, which generates a chain of samples with varying levels of difficulty and feeds them to models to obtain a chain of logit distribution of different confidence levels. Then, in a collaborative learning approach, DC conducts distillation from *less confident* predictions to *more confident* predictions, to serve as a learned entropy regularization.

### 165 3. Plasticity and Stability in online CL

166 In this section, we revise the metric for model plasticity  
167 and propose a novel metric for model stability. In addition,  
168 we quantitatively derived the impact of model plasticity and  
169 stability on the final performance.

#### 170 3.1. Model Plasticity

171 The model plasticity measures the learner’s capability to  
172 learn new knowledge when a new task arrives. Several dif-  
173 ferent metrics have been proposed to measure the model  
174 plasticity [9, 30, 37, 47]. In our work, we evaluate the model  
175 plasticity with Learning Accuracy (LA) [37]. Formally, the  
176 Learning Accuracy for the  $j$ -th task is defined as:

$$177 \quad l_j = a_j^j, \quad (1)$$

178 where  $a_j^i$  is the accuracy evaluated on the test set of task  $j$   
179 after training the network from task 1 to task  $i$ . For an over-  
180 all metric normalized against all tasks, the averaged Learn-  
181 ing Accuracy is written as  $LA = \frac{1}{T} \sum_{j=1}^T l_j$ , and  $T$  is the  
182 number of tasks in total.

#### 183 3.2. Model Stability

184 The stability measures how much the model forgets given  
185 its current state. The most commonly used metric in previ-  
186 ous CL research is the Forgetting Measure (FM) [9]. Intu-  
187 itively, FM for the  $j$ -th task  $fm_j^k$  reveals how much per-  
188 formance the model loses on a given task  $j$ , after training on  
189 task  $k$ , compared with its maximum performance obtained  
190 in the past:

$$191 \quad fm_j^k = \max_{i \in \{1, \dots, k-1\}} (a_j^i - a_j^k), \forall j < k. \quad (2)$$

For the overall metric obtained across all tasks, FM can be  
expressed as:

$$FM = \frac{1}{T-1} \sum_{j=1}^{T-1} fm_j^T. \quad (3)$$

In our work, instead of using FM as the stability metric,  
we propose forgetting measure on a relative basis, which  
we call Relative Forgetting (RF). There are two reasons for  
shifting from absolute forgetting to relative forgetting:

1. RF is more fair for methods with higher plasticity. This  
is because methods with poor plasticity will never have  
a large FM, and FM is capped by the maximum perfor-  
mance obtained by the learner in the past. Moreover,  
even if two learners lose the same absolute performance,  
the more plastic learner can still be regarded as forget-  
ting less, because it has a higher peak performance and  
loses less proportion of its performance;
2. RF helps quantitatively derive the relationship between  
the model stability and final performance.

Intuitively, RF measures how much *proportion* of perfor-  
mance the model forgets. And RF for the  $j$ -th task after  
training on task  $k$ , can be defined as:

$$f_j^k = \max_{i \in \{1, \dots, k\}} \left( 1 - \frac{a_j^k}{a_j^i} \right), \forall j \leq k. \quad (4)$$

The overall metric averaged across all tasks can be written  
as:

$$RF = \frac{1}{T} \sum_{j=1}^T f_j^T. \quad (5)$$

### 216 3.3. Impact on the Overall Performance

217 For online CL, the model’s final average accuracy (AA) is  
218 the most vital metric. In this subsection, we try to show  
219 how the model plasticity and stability will impact the final  
220 performance quantitatively.

221 The model’s final average accuracy can be calculated by:

$$222 \quad AA = \frac{1}{T} \sum_{j=1}^T a_j^T. \quad (6)$$

223 With the definition of our plasticity metric (LA) and stabil-  
224 ity metric (RF), we can easily find the relationship between  
225 learning accuracy, relative forgetting, and accuracy:

$$226 \quad a_j^i \geq l_j \times (1 - f_j^i), \quad (7)$$

227 where we take the equal sign when  $a_j^j = \max_{i \in \{1, \dots, i\}} a_j^i$ .  
228 When generalizing the class-wise final accuracy  $a_j^T$  to the  
229 average accuracy (AA), we need to take the dot product of  
230 the LA vector  $[l_1, \dots, l_T]$  with RF vector  $[f_1^T, \dots, f_T^T]$  which  
231 is trivial. More intuitively, in practice, we can make the  
232 approximation with:

$$233 \quad AA \gtrsim LA \times (1 - RF). \quad (8)$$

234 As indicated by Eq. 8, the lower bound of the final per-  
235 formance is proportional to  $LA$  and  $1 - RF$ , which suggests  
236 that both plasticity (LA) and stability (RF) play a crucial  
237 role in the final accuracy. Our findings reveal the impor-  
238 tance of the model plasticity which was neglected in the  
239 past. And it can serve as a good guide for future online CL  
240 research.

## 241 4. Proposed Method

242 In this section, we first justify our motivation with the find-  
243 ings in Sec. 3. Then, we introduce our proposed strategy:  
244 Collaborative Continual Learning and Distillation Chain.  
245 Finally, we show how to adapt our proposed strategy to the  
246 existing online CL methods and boost their plasticity.

### 247 4.1. Motivation Justification

248 Online continual learners aim to continuously adapt to non-  
249 stationary data streams, efficiently acquiring new knowl-  
250 edge while retaining previously learned information. In  
251 current online CL research, almost all of the efforts focus  
252 on alleviating catastrophic forgetting, and the importance  
253 of “learning capability” on new knowledge is greatly over-  
254 looked. However, our finding in Sec. 3 shows both plasticity  
255 and stability play important roles in achieving decent final  
256 performance.

257 While replay-based methods were originally designed  
258 to tackle forgetting issues, Fig. 1 demonstrates that they  
259 can implicitly mitigate the plasticity gap. Nonetheless, as

shown in Fig. 1, the limited plasticity is still a significant  
barrier to the performance, even with replay. To this end,  
we explicitly focus on the plasticity perspective.

The potential of collaborative learning to improve con-  
vergence in non-continual scenarios [5] positions it as a  
promising candidate for enhancing plasticity. With the ap-  
parent lack of focus on plasticity, collaborative learning has  
yet to be leveraged to boost convergence of online continual  
learners. In our research, we propose to exploit collabora-  
tive learning convergence properties for improving plastic-  
ity. We find that similar to non-continual scenarios, collab-  
orative learning strategy can boost convergence by allowing  
more parallelism in the training and more maneuverability  
of the continual learners. Moreover, to fully take advan-  
tage of collaborative learning, we also propose Distillation  
Chain (DC), an entropy regularization based optimization  
strategy in collaborative learning specifically designed for  
online CL.

### 278 4.2. Collaborative Continual Learning

The introduced Collaborative Continual Learning (CCL)  
enables more parallelism and flexibility in training online  
continual learners, and it is the key to improving the model  
plasticity and the final performance. As shown in Fig. 2,  
CCL involves two peer continual learners of the same ar-  
chitecture and optimizer setting training in a peer-teaching  
manner. In the training phase, networks are supervised with  
both the ground truth label and the predictions of their peers.  
In the inference phase, models can either make predictions  
collaboratively with ensemble methods [5] to get a better  
performance or predict independently for the sake of com-  
putation efficiency. If we denote two networks in CCL as  
 $\theta^1$  and  $\theta^2$ , we formulate our loss to network  $\theta^1$  as:

$$292 \quad \mathcal{L}_{CCL}^1 = \lambda_1 \cdot \mathcal{L}_{cls}(\theta^1(X), y) + \lambda_2 \cdot D_{KL}(\theta^1(X)/\tau, \theta^2(X)/\tau), \quad (9)$$

where  $(X, y)$  is the data-label pair,  $\mathcal{L}_{cls}(\cdot)$  is the classifica-  
tion loss in the baseline method CCL adapts to,  $D_{KL}(\cdot)$  is  
the Kullback-Leibler divergence,  $\lambda_1$  and  $\lambda_2$  are balancing  
hyperparameters and  $\tau$  is the temperature hyperparameter.  
Note that the network  $\theta^2$  should be trained with  $\mathcal{L}_{CCL}^2$ , re-  
spectively.

### 299 4.3. Distillation Chain

To fully take advantage of CCL, we propose Distillation  
Chain (DC), an entropy regularization based strategy explic-  
itly designed for online CL. As illustrated in Fig. 2, DC  
comprises two steps: (1) generating a chain of samples with  
different levels of difficulty [43] using data augmentation,  
and (2) distillation of logit distribution from *harder* samples  
to *easier* samples in a collaborative learning way.

The main motivation of DC originates from the idea  
of entropy regularization-based optimization strategies, like

309 label smoothing [44], knowledge distillation [49], and con-  
 310 fidence penalty [36], where we find that overconfidence will  
 311 hurt performance in non-continual scenarios. As shown  
 312 in the supplementary material, we observed a similar phe-  
 313 nomenon in online CL. To tackle the problem, DC uses data  
 314 augmentation strategies to generate samples with different  
 315 levels of difficulty and produces logit distribution with dif-  
 316 ferent confidence. The distillation from *less confident* pre-  
 317 dictions to *more confident* predictions weakens the overall  
 318 confidence of the network and benefits the performance by  
 319 improving the generalization capability.

320 In our work, we use a geometric distortion comprised of  
 321 RandomCrop and RandomHorizontalFlip as the first step of  
 322 DC augmentation. After that, we use RandAugment [13]  
 323 for the subsequent augmentations and we involve two hy-  
 324 perparameters  $N$  and  $M$  for RandAugment. We take three  
 325 augmentation steps and distill the logit distribution from the  
 326 teacher with *harder* samples to the student with *easier* sam-  
 327 ples. We formulate our loss with DC to network  $\theta^1$  as:

$$328 \mathcal{L}_{DC}^1 = \lambda_1 \sum_{i=1}^3 \mathcal{L}_{cls}(\theta^1(X_i), y) \quad (10)$$

$$+ \lambda_2 \sum_{i=1}^3 D_{KL}(\theta^1(X_{i-1})/\tau, \theta^2(X_i)/\tau),$$

329 where  $X_i$  is the augmentation of input sample  $X$  after  $i$   
 330 augmentation steps. More discussion about why and how  
 331 DC works can be found in the supplementary material.

#### 332 4.4. Apply CCL-DC to online CL methods

333 The overall loss to network  $\theta^1$  when adapting CCL-DC can  
 334 be written as:

$$335 \mathcal{L}^1 = \mathcal{L}_{Baseline} + \mathcal{L}_{CCL}^1 + \mathcal{L}_{DC}^1, \quad (11)$$

336 where  $\mathcal{L}_{Baseline}$  is the loss function of the baseline model  
 337 CCL-DC adapts to. Note that the model  $\theta^2$  should be  
 338 trained similarly. In Algorithm 1, we provide a Pytorch-  
 339 like pseudo-code demonstrating how to incorporate CCL-  
 340 DC into a given baseline. For simplicity, we only show the  
 341 loss function for model  $\theta^1$ . Also, we omitted the memory  
 342 buffer in the pseudo-code. However, the training should  
 343 be consistent with the baseline, using both streaming and  
 344 memory data.

## 345 5. Experiments

### 346 5.1. Experimental Setup

347 **Datasets.** We use four image classification benchmark  
 348 datasets to evaluate the effectiveness of our method, includ-  
 349 ing CIFAR-10 [27], CIFAR-100 [27], TinyImageNet [28],  
 350 and ImageNet-100 [24]. More detailed information about

**Algorithm 1** PyTorch-like pseudo-code of CCL-DC to in-  
 351 tegrate to other baselines.

```

352 # model1: student model
353 # model2: teacher model
354 # optim1: optimizer for student model
355 # cls: classification loss in baseline
356 for x, y in dataloader:
357     # Baseline loss
358     loss_baseline = criterion_baseline(model1, x, y)
359
360     # DC Augmentation
361     x1 = geometric_distortion(x)
362     x2 = RandAugment(x1, N, M)
363     x3 = RandAugment(x2, N, M)
364
365     # CCL-DC loss
366     ls, ls1, ls2, ls3 = model1(x, x1, x2, x3)
367     lt, lt1, lt2, lt3 = model2(x, x1, x2, x3) # no grad
368
369     loss_cls = cls(ls, y) + cls(ls1, y) + cls(ls2, y) +
370     ↪ cls(ls3, y)
371     loss_ccl = kl_div(ls/t, lt/t) # temperature t
372     loss_dc = kl_div(ls/t, lt1/t) + kl_div(ls1/t, lt2/t) +
373     ↪ kl_div(ls2/t, lt3/t)
374
375     loss_ours = lam1*loss_cls + lam2*(loss_ccl + loss_dc)
376     loss = loss_baseline + loss_ours
377
378     optim1.zero_grad()
379     loss.backward()
380     optim1.step()
  
```

the dataset split and task allocation is given in the supple-  
 351 mentary material. 352

**Baselines.** To show the effectiveness of our strategy, we  
 353 applied CCL-DC to six typical and state-of-the-art online  
 354 CL methods, including ER [38], DER++ [6], ER-ACE [7],  
 355 OCM [21], GSA [22], and OnPro [48]. 356

**Implementation details.** We use full ResNet-18 [23] (not  
 357 pre-trained) as the backbone for every method. For each  
 358 baseline method, we perform a hyperparameter search on  
 359 CIFAR-100,  $M=2k$ , and apply the hyperparameter to all of  
 360 the settings. For fair comparison, we use the same opti-  
 361 mizer and hyperparameter setting when adapting CCL-DC  
 362 to the baselines. For hyperparameters unique to CCL-DC,  
 363 we conduct another hyperparameter search as stated in the  
 364 supplementary material. We set the streaming batch size to  
 365 10 and the memory batch size to 64. We do not use the  
 366 multiple update trick as described in [3]. More detailed in-  
 367 formation about data augmentation, hyperparameter search,  
 368 and hardware environments is given in the supplementary  
 369 material. 370

### 371 5.2. Results and Analysis

**Final average accuracy.** Table 1 presents the results of  
 372 average accuracy (AA) at the end of the training on four  
 373 datasets. As indicated in Sec. 4, to fully take advantage  
 374 of collaborative learning, we show the results with the en-  
 375 semble of two models, with the independent model perfor-  
 376

Dataset	CIFAR10		CIFAR100			Tiny-ImageNet			ImageNet-100
	500	1000	1000	2000	5000	2000	5000	10000	5000
ER [38]	56.68±1.89	62.32±4.13	24.47±0.72	31.89±1.45	39.41±1.81	10.82±0.79	19.16±1.42	24.71±2.52	33.30±1.74
ER + Ours	<b>66.43±2.48</b>	<b>74.10±1.71</b>	<b>33.43±1.06</b>	<b>44.45±1.04</b>	<b>53.81±1.16</b>	<b>16.56±1.63</b>	<b>29.39±1.23</b>	<b>37.73±0.85</b>	<b>43.11±1.49</b>
DER++ [6]	58.04±2.30	64.02±1.92	25.09±1.41	32.33±2.66	38.31±2.28	8.73±1.58	17.95±2.49	19.40±3.71	34.75±2.23
DER++ + Ours	<b>68.79±1.42</b>	<b>74.25±1.10</b>	<b>34.36±0.89</b>	<b>43.52±1.35</b>	<b>52.95±0.86</b>	<b>10.99±1.39</b>	<b>21.68±1.94</b>	<b>28.01±2.46</b>	<b>45.70±1.32</b>
ER-ACE [7]	53.26±3.04	59.94±2.40	28.36±1.99	34.21±1.53	39.39±1.31	13.56±1.00	20.84±0.43	25.92±1.07	38.37±1.20
ER-ACE + Ours	<b>70.08±1.38</b>	<b>75.56±1.14</b>	<b>37.20±1.15</b>	<b>45.14±1.00</b>	<b>53.92±0.48</b>	<b>18.32±1.49</b>	<b>26.22±2.01</b>	<b>32.23±1.70</b>	<b>45.15±1.94</b>
OCM [21]	68.19±1.75	73.15±1.05	28.02±0.74	35.69±1.36	42.22±1.06	18.36±0.95	26.74±1.02	31.94±1.19	23.67±2.36
OCM + Ours	<b>74.14±0.85</b>	<b>77.66±1.46</b>	<b>35.00±1.15</b>	<b>43.34±1.51</b>	<b>51.43±1.37</b>	<b>23.36±1.18</b>	<b>33.17±0.97</b>	<b>39.25±0.88</b>	<b>43.19±0.98</b>
GSA [22]	60.34±1.97	66.54±2.28	27.72±1.57	35.08±1.37	41.41±1.65	12.44±1.17	19.59±1.30	25.34±1.43	41.03±0.99
GSA + Ours	<b>68.91±1.68</b>	<b>75.78±1.16</b>	<b>35.56±1.39</b>	<b>44.74±1.32</b>	<b>55.39±1.09</b>	<b>16.70±1.66</b>	<b>28.11±1.70</b>	<b>37.13±1.75</b>	<b>44.28±1.16</b>
OnPro [48]	70.47±2.12	74.70±1.51	27.22±0.77	33.33±0.93	41.59±1.38	14.32±1.40	21.13±2.12	26.38±2.18	38.75±1.03
OnPro + Ours	<b>74.49±2.14</b>	<b>78.64±1.42</b>	<b>34.76±1.12</b>	<b>41.89±0.82</b>	<b>50.01±0.85</b>	<b>21.81±1.02</b>	<b>32.00±0.72</b>	<b>38.18±1.02</b>	<b>47.93±1.26</b>

Table 1. Average Accuracy (% , higher is better) on four benchmark datasets with difference memory buffer size  $M$ , with and without our proposed CCL scheme. The result of our method is given by the ensemble of two peer models. All values are averages of 10 runs.

Dataset	CIFAR10		CIFAR100			Tiny-ImageNet			ImageNet-100
	500	1000	1000	2000	5000	2000	5000	10000	5000
ER	83.13±1.60	78.15±3.60	53.77±1.51	51.53±1.66	50.79±0.71	68.15±1.47	64.99±1.22	64.44±1.45	53.95±1.51
ER + Ours	<b>90.60±1.50</b>	<b>89.99±1.50</b>	<b>72.38±0.66</b>	<b>70.86±0.72</b>	<b>68.84±1.05</b>	<b>85.24±0.53</b>	<b>81.75±0.83</b>	<b>79.54±0.74</b>	<b>68.73±1.21</b>
DER++	77.14±2.96	78.00±2.16	56.13±3.75	55.33±3.26	56.32±3.44	70.01±1.83	66.87±1.30	70.28±2.42	60.65±2.97
DER++ + Ours	<b>88.85±1.88</b>	<b>89.00±1.67</b>	<b>72.85±1.37</b>	<b>71.54±1.99</b>	<b>69.52±2.37</b>	<b>82.83±1.27</b>	<b>78.80±1.62</b>	<b>77.79±0.86</b>	<b>70.16±1.03</b>
ER-ACE	57.66±4.16	61.59±3.35	38.53±1.61	39.95±2.00	41.56±1.44	5.60±1.45	4.83±0.78	4.92±0.95	49.82±1.05
ER-ACE + Ours	<b>88.37±1.39</b>	<b>88.40±1.15</b>	<b>69.47±0.88</b>	<b>68.39±1.32</b>	<b>66.63±0.90</b>	<b>21.91±5.16</b>	<b>21.88±4.39</b>	<b>18.88±3.12</b>	<b>68.52±0.82</b>
OCM	78.71±3.66	81.33±2.06	40.87±1.60	42.00±1.48	42.43±1.80	18.56±2.87	15.86±2.01	15.03±2.02	20.77±1.88
OCM + Ours	<b>82.39±2.23</b>	<b>84.53±1.63</b>	<b>48.89±2.04</b>	<b>49.83±2.01</b>	<b>49.94±2.16</b>	<b>31.69±1.81</b>	<b>29.54±2.35</b>	<b>28.10±2.28</b>	<b>48.20±1.38</b>
GSA	79.87±3.26	77.09±4.55	58.16±1.58	55.13±1.81	50.34±1.73	20.46±1.59	15.86±1.26	14.50±0.63	62.59±1.17
GSA + Ours	<b>91.69±1.11</b>	<b>90.98±1.33</b>	<b>73.73±1.03</b>	<b>72.68±0.98</b>	<b>70.36±1.07</b>	<b>80.36±1.22</b>	<b>74.77±1.66</b>	<b>70.71±1.19</b>	<b>73.71±1.12</b>
OnPro	84.23±2.00	85.60±1.56	41.34±1.63	42.59±1.65	42.92±1.00	20.84±1.47	16.73±1.27	15.82±1.04	39.60±0.86
OnPro + Ours	<b>90.39±1.59</b>	<b>90.18±1.58</b>	<b>46.30±1.10</b>	<b>47.13±1.01</b>	<b>47.27±1.81</b>	<b>25.87±1.91</b>	<b>21.40±1.52</b>	<b>19.75±1.22</b>	<b>52.55±2.18</b>

Table 2. Learning Accuracy (% , higher is better) on four benchmark datasets with difference memory buffer size  $M$ , with and without our proposed CCL scheme. The result of our method is given by the ensemble of two peer models. All values are averages of 10 runs.

377 mance available in Sec. 6. Generally, the ensemble method  
378 provides about 1% additional accuracy compared to inde-  
379 pendent inference. For all datasets, memory size  $M$ , and  
380 baseline methods, applying CCL-DC constantly improves  
381 the performance by a large margin. Notably, even for state-  
382 of-the-art methods like GSA and OnPro, we can still gain  
383 significant performance when incorporating CCL-DC.

384 More interestingly, for almost all settings with different  
385 memory buffer sizes  $M$ , the performance gain tends to be  
386 a constant on a relative basis. For example, CCL-DC can  
387 boost the performance of ER on Tiny-ImageNet from 10.82  
388 to 16.56 when  $M=2k$ , which is a 53.0% performance gain  
389 on a relative basis. The performance gain is 53.4% and  
390 52.7% when  $M=5k$  and  $M=10k$  respectively. This indi-  
391 cates that we can achieve a decent performance gain regard-  
392 less of the memory buffer size, and it shows the scalability  
393 of our method to different resource conditions.

**Plasticity and stability metric.** As mentioned in Sec. 3,  
we evaluate the plasticity and stability of different continual  
learners with Learning Accuracy and Relative Forgetting,  
respectively. Table 2 shows the plasticity metric on four  
datasets. For all settings, CCL-DC constantly improves the  
model plasticity by a large margin. For model stability, as  
indicated by RF in Table 3, models trained with CCL-DC  
are comparable with the baselines under most cases. ER-  
ACE is an exception as its plasticity is unexpectedly low,  
especially on TinyImagenet. Also, the stability of ER-ACE  
is compromised when incorporating CCL-DC. We will ex-  
plain the reason for this unexpected phenomenon in the sup-  
plementary material.

### 5.3. Ablation Studies

**Effect of multiview learning.** As mentioned in Sec. 4,  
CCL-DC benefits from multiview learning with data aug-

Dataset	CIFAR10		CIFAR100			Tiny-ImageNet			ImageNet-100
	500	1000	1000	2000	5000	2000	5000	10000	5000
Memory Size $M$									
ER	31.63 $\pm$ 3.81	20.63 $\pm$ 8.32	55.71 $\pm$ 2.24	39.11 $\pm$ 3.87	23.05 $\pm$ 3.69	85.00 $\pm$ 1.30	71.62 $\pm$ 2.18	62.43 $\pm$ 3.83	39.26 $\pm$ 3.21
ER + Ours	<b>26.74<math>\pm</math>3.99</b>	<b>17.58<math>\pm</math>2.71</b>	<b>54.34<math>\pm</math>2.22</b>	<b>37.67<math>\pm</math>2.16</b>	<b>21.98<math>\pm</math>2.59</b>	<b>81.13<math>\pm</math>1.93</b>	<b>64.79<math>\pm</math>1.32</b>	<b>53.18<math>\pm</math>0.99</b>	<b>37.78<math>\pm</math>2.18</b>
DER++	23.60 $\pm$ 3.64	17.71 $\pm$ 2.18	55.65 $\pm$ 4.36	41.27 $\pm$ 4.93	31.72 $\pm$ 3.95	87.79 $\pm$ 2.35	73.28 $\pm$ 3.88	72.51 $\pm$ 5.53	42.97 $\pm$ 5.89
DER++ + Ours	<b>22.62<math>\pm</math>3.03</b>	<b>16.43<math>\pm</math>3.36</b>	<b>53.45<math>\pm</math>1.40</b>	<b>39.39<math>\pm</math>2.71</b>	<b>23.71<math>\pm</math>3.39</b>	<b>87.16<math>\pm</math>1.60</b>	<b>73.15<math>\pm</math>2.15</b>	<b>64.48<math>\pm</math>3.08</b>	<b>35.32<math>\pm</math>2.80</b>
ER-ACE	<b>12.25<math>\pm</math>3.84</b>	<b>9.92<math>\pm</math>2.83</b>	<b>25.88<math>\pm</math>4.10</b>	<b>17.68<math>\pm</math>1.90</b>	<b>10.62<math>\pm</math>2.08</b>	57.41 $\pm$ 2.38	44.48 $\pm$ 1.96	37.83 $\pm$ 3.12	<b>23.92<math>\pm</math>2.05</b>
ER-ACE + Ours	20.62 $\pm$ 2.26	14.32 $\pm$ 2.58	46.78 $\pm$ 1.91	34.19 $\pm$ 2.40	19.01 $\pm$ 0.94	<b>56.56<math>\pm</math>4.16</b>	<b>42.20<math>\pm</math>3.94</b>	<b>31.13<math>\pm</math>3.52</b>	34.43 $\pm$ 3.60
OCM	13.05 $\pm$ 4.37	11.00 $\pm$ 3.11	31.16 $\pm$ 2.69	17.90 $\pm$ 3.73	6.85 $\pm$ 2.25	56.66 $\pm$ 2.53	40.59 $\pm$ 1.55	30.80 $\pm$ 2.29	<b>4.55<math>\pm</math>1.60</b>
OCM + Ours	<b>10.75<math>\pm</math>2.52</b>	<b>8.45<math>\pm</math>2.63</b>	<b>29.65<math>\pm</math>4.00</b>	<b>17.02<math>\pm</math>3.01</b>	<b>6.16<math>\pm</math>1.35</b>	<b>51.58<math>\pm</math>2.81</b>	<b>35.58<math>\pm</math>2.54</b>	<b>27.24<math>\pm</math>1.60</b>	15.33 $\pm$ 2.28
GSA	25.02 $\pm$ 2.83	<b>16.56<math>\pm</math>4.02</b>	53.42 $\pm$ 3.12	<b>37.29<math>\pm</math>2.60</b>	<b>20.50<math>\pm</math>4.33</b>	<b>66.87<math>\pm</math>3.31</b>	<b>53.42<math>\pm</math>3.84</b>	<b>43.44<math>\pm</math>3.81</b>	<b>35.44<math>\pm</math>2.42</b>
GSA + Ours	<b>24.96<math>\pm</math>3.27</b>	16.59 $\pm$ 2.09	<b>52.29<math>\pm</math>2.06</b>	38.76 $\pm$ 2.41	21.36 $\pm$ 2.36	80.08 $\pm$ 1.97	63.85 $\pm$ 1.78	49.73 $\pm$ 2.10	40.46 $\pm$ 2.54
OnPro	<b>16.47<math>\pm</math>4.23</b>	12.93 $\pm$ 3.02	35.03 $\pm$ 4.45	24.26 $\pm$ 2.31	12.04 $\pm$ 2.11	64.69 $\pm$ 3.36	50.47 $\pm$ 4.20	42.81 $\pm$ 4.63	<b>14.44<math>\pm</math>2.08</b>
OnPro + Ours	17.54 $\pm$ 4.15	<b>12.90<math>\pm</math>2.77</b>	<b>27.64<math>\pm</math>3.29</b>	<b>17.78<math>\pm</math>1.39</b>	<b>8.41<math>\pm</math>2.62</b>	<b>56.03<math>\pm</math>2.96</b>	<b>38.70<math>\pm</math>1.88</b>	<b>29.24<math>\pm</math>1.33</b>	15.72 $\pm$ 3.29

Table 3. Relative Forgetting (% , lower is better) on four benchmark datasets with difference memory buffer size  $M$ , with and without our proposed CCL scheme. The result of our method is given by the ensemble of two peer models. All values are averages of 10 runs.

Method	Acc. $\uparrow$	LA $\uparrow$
ER	31.89 $\pm$ 1.45	51.53 $\pm$ 1.66
ER + Multiview	38.18 $\pm$ 1.46	64.02 $\pm$ 1.12
ER + Ours (CCL only)	41.05 $\pm$ 1.21	68.76 $\pm$ 0.79
ER + Ours	44.45 $\pm$ 1.04	70.86 $\pm$ 0.72
ER-ACE	34.21 $\pm$ 1.53	39.95 $\pm$ 2.00
ER-ACE + Multiview	38.61 $\pm$ 1.48	47.45 $\pm$ 1.88
ER-ACE + Ours (CCL only)	40.90 $\pm$ 1.08	50.91 $\pm$ 1.63
ER-ACE + Ours	45.14 $\pm$ 1.00	68.39 $\pm$ 1.32

Table 4. Ablation studies on CIFAR-100 ( $M=2k$ ). We report the ensemble performance for methods incorporating CCL.

mentation in DC. For fair comparison, we explore how multiview learning will impact the performance of the baselines. We apply the classification loss part of CCL-DC to the baselines. Table 4 demonstrates that multiview learning can improve both AA and LA of baselines. However, those performance gains are still inferior to CCL-DC.

**Effect of CCL.** We evaluate how CCL alone can improve the baselines. In the experiments, we remove multiview learning and DC, and we train the continual learner pair with the loss illustrated in Eq. 9. Table 4 shows the performance gain for ER and ER-ACE. We can see that CCL alone can provide significant gains in both final accuracy and plasticity. Also, when combining CCL with DC, the performance can be further improved.

**Distillation scheme of DC.** We also evaluate the effectiveness of DC’s strategy of distilling from harder samples to easier samples in collaborative learning manner. As shown in Table 5, we compared it with other distillation strategies. The result shows that the distillation scheme of DC constantly outperforms other schemes. Extra experiments explaining the working mechanism of DC can be found in the supplementary material.

Method	Distillation scheme	Acc. $\uparrow$	LA $\uparrow$
ER	Easy to hard	40.95 $\pm$ 0.97	60.03 $\pm$ 0.98
ER	Same difficulty	43.64 $\pm$ 1.09	69.49 $\pm$ 0.78
ER	Hard to easy (Ours)	44.45 $\pm$ 1.04	70.86 $\pm$ 0.72
ER-ACE	Easy to hard	38.46 $\pm$ 1.51	39.00 $\pm$ 1.03
ER-ACE	Same difficulty	43.81 $\pm$ 1.28	55.37 $\pm$ 1.54
ER-ACE	Hard to easy (Ours)	45.14 $\pm$ 1.00	68.39 $\pm$ 1.32

Table 5. Comparison of different distillation schemes in DC on CIFAR-100 ( $M=2k$ ).

## 6. Discussions

In this section, we analyze some properties of CCL-DC.

**Improving plasticity.** One of the important advantages of CCL-DC is that it can improve the plasticity of continual learners. This can be evident by plasticity metrics like LA. Moreover, we have observed that the plasticity of CCL-DC facilitates the model to converge faster and descend to a deeper loss. Figure 3 illustrates the classification loss (cross-entropy) curve of the model. To obtain the loss curve, we take a snapshot of the model every 10 iterations and compute the cross-entropy over all the training samples on the *current* task. We plot the curve on the logarithm scale so that it is easy to observe that CCL-DC helps the model descend deeper at the end of each task.

**Improving feature discrimination.** Another advantage of CCL-DC is its ability to enhance the feature discrimination of continual learners. Fig. 4 illustrates the t-SNE visualization [46] of the memory data’s embedding space at the end of the training. We can see that the feature representation of the method with CCL-DC is more discriminative compared with the baseline.

Moreover, we can evaluate the feature discrimination using the clustering methods. Following [31], we remove the

Method	NCM Acc. $\uparrow$	Logit Acc. $\uparrow$
ER	36.56 $\pm$ 0.60	31.89 $\pm$ 1.45
ER + Ours	44.76 $\pm$ 0.55	44.45 $\pm$ 1.04
ER-ACE	34.91 $\pm$ 1.02	34.21 $\pm$ 1.53
ER-ACE + Ours	45.62 $\pm$ 1.04	45.14 $\pm$ 1.00
OnPro	34.32 $\pm$ 0.95	33.33 $\pm$ 0.93
OnPro + Ours	42.82 $\pm$ 0.67	41.89 $\pm$ 0.82

Table 6. Final average accuracy on CIFAR-100 (M=2k), with and without NCM classifier.

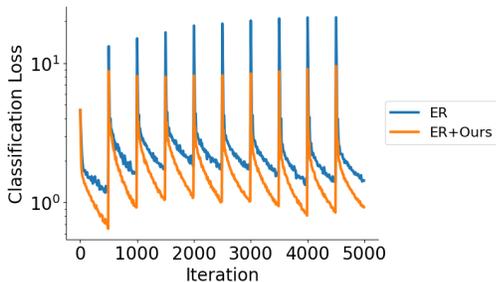


Figure 3. Classification loss curve of ER on CIFAR-100 (M=2k). The curve is calculated on all training samples of the *current* task. Since there are 10 tasks in total, the curve has 10 peaks.

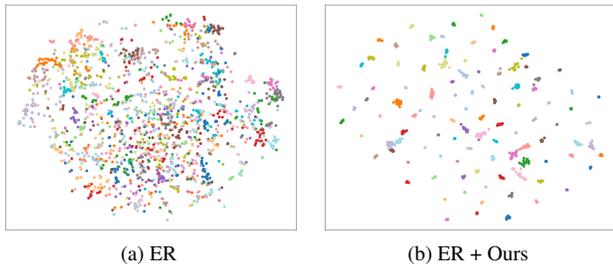


Figure 4. T-SNE visualization of memory data at the end of training on CIFAR-100 (M=2k).

455 final FC classifier and use Nearest-Class-Mean (NCM) [31]  
 456 classifier with intermediate representations. Table 6 demon-  
 457 strates that CCL-DC can greatly enhance the NCM accu-  
 458 racy, which evidences the capability of CCL-DC in improv-  
 459 ing feature discrimination.

460 **Alleviating shortcut learning.** Shortcut learning [18] is  
 461 another commonly observed issue that hinders the general-  
 462 ization capability of continual learners [48]. In Fig. 5, we  
 463 use GradCAM++ [8] on the training set of ImageNet-100  
 464 (M=5k) at the end of the training of ER and GSA. Although  
 465 both ER and GSA make correct predictions, we observed  
 466 that they focus on irrelevant objects, which indicates a ten-  
 467 dency toward shortcut learning. Also, we can see that by  
 468 integrating CCL-DC, the shortcut learning can be greatly  
 469 alleviated.

470 **Independent network performance.** Although the en-  
 471 semble method gives extra performance at inference time,  
 472 by averaging the logit output of two networks in CCL-DC,

Method	Ind. Acc. $\uparrow$	Ens. Acc. $\uparrow$
ER + Ours	43.58 $\pm$ 1.05	44.45 $\pm$ 1.04
DER++ + Ours	42.79 $\pm$ 1.38	43.52 $\pm$ 1.35
ER-ACE + Ours	44.15 $\pm$ 1.05	45.14 $\pm$ 1.00
OCM + Ours	42.39 $\pm$ 1.36	43.34 $\pm$ 1.51
GSA + Ours	43.84 $\pm$ 1.34	44.74 $\pm$ 1.32
OnPro + Ours	41.18 $\pm$ 0.83	41.89 $\pm$ 0.82

Table 7. Comparison of the final average accuracy achieved through independent inference and the use of the ensemble method on CIFAR-100 (M=2k). Independent accuracy (Ind. Acc.) is calculated by averaging the accuracy of two networks in CCL-DC. All values are averaged over 10 runs.

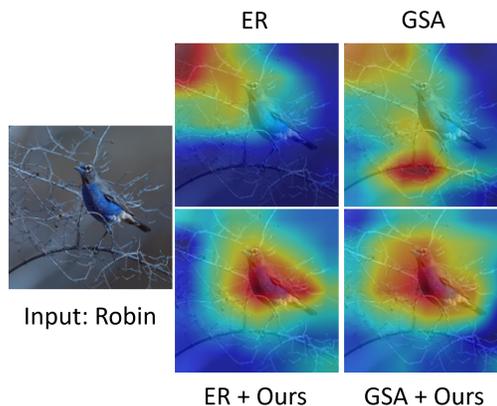


Figure 5. GradCAM++ visualization on the training set of ImageNet-100 (M=5k). Shortcut learning exists in the baseline methods despite making correct predictions.

it also doubles the computation. In some cases, compu- 473  
 tational efficiency becomes more crucial during inference. 474  
 Continual learners trained with CCL-DC are also able to do 475  
 inference independently, albeit with a slight performance 476  
 drop compared with ensemble inference. Table 7 illustrates 477  
 the accuracy achieved through independent inference. It is 478  
 evident that the performance loss in independent inference, 479  
 when compared to ensemble inference, is minimal (approx- 480  
 imately 1%). 481

## 7. Conclusion 482

In this paper, we highlight the significance of plasticity 483  
 in online CL, which has been neglected in prior re- 484  
 search. We also establish the quantitative link between 485  
 plasticity, stability, and final accuracy. The quanti- 486  
 tative relationship sheds light on the future direction 487  
 of online CL research. Based on this, we introduce 488  
 collaborative learning into online CL and propose CCL- 489  
 DC, a strategy that can be seamlessly integrated into 490  
 existing online CL methods. Extensive experiments 491  
 show the effectiveness of CCL-DC in boosting plas- 492  
 ticity and subsequently improving the final performance. 493  
 494

## References

- 495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550
- [1] Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert gate: Lifelong learning with a network of experts. In *CVPR*, 2017. 2
- [2] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *ECCV*, 2018. 2
- [3] Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. In *NeurIPS*. 2019. 1, 5
- [4] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *NeurIPS*, 2019. 1
- [5] Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert Ormandi, George E Dahl, and Geoffrey E Hinton. Large scale distributed neural network training through online distillation. *arXiv preprint arXiv:1804.03235*, 2018. 1, 2, 4
- [6] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *NeurIPS*, 2020. 1, 2, 5, 6
- [7] Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky. New insights on reducing abrupt representation change in online continual learning. *arXiv preprint arXiv:2104.05025*, 2021. 1, 2, 5, 6
- [8] Aditya Chattopadhyay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-cam++: Generalized gradient-based visual explanations for deep convolutional networks. In *WACV*, 2018. 8
- [9] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *ECCV*, 2018. 2, 3
- [10] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019. 2
- [11] Zhiyuan Chen and Bing Liu. Lifelong machine learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 12(3):1–207, 2018. 1
- [12] Yubei Chen Chun-Hsiao Yeh. IN100pytorch: Pytorch implementation: Training resnets on imagenet-100. <https://github.com/danielchye/Imagenet-100-Pytorch>, 2022.
- [13] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPR-W*, 2020. 5
- [14] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3366–3385, 2021. 1
- [15] Cyprien de Masson D’Autume, Sebastian Ruder, Lingpeng Kong, and Dani Yogatama. Episodic memory in lifelong language learning. In *NeurIPS*, 2019. 2
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 555
- [17] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. *arXiv preprint arXiv:1701.08734*, 2017. 2
- [18] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. 8
- [19] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013. 1
- [20] Qiushan Guo, Xinjiang Wang, Yichao Wu, Zhipeng Yu, Ding Liang, Xiaolin Hu, and Ping Luo. Online knowledge distillation via collaborative learning. In *CVPR*, 2020. 2
- [21] Yiduo Guo, Bing Liu, and Dongyan Zhao. Online continual learning through mutual information maximization. In *ICML*, 2022. 1, 2, 5, 6
- [22] Yiduo Guo, Bing Liu, and Dongyan Zhao. Dealing with cross-task class discrimination in online continual learning. In *CVPR*, 2023. 1, 2, 5, 6
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [24] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, 2019. 5
- [25] Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. *arXiv preprint arXiv:1810.12488*, 2018. 1
- [26] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017. 2
- [27] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [28] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015. 5
- [29] Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang. Overcoming catastrophic forgetting by incremental moment matching. In *NeurIPS*, 2017. 2
- [30] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *NeurIPS*, 2017. 3
- [31] Zheda Mai, Ruiwen Li, Hyunwoo Kim, and Scott Sanner. Supervised contrastive replay: Revisiting the nearest class mean classifier in online class-incremental continual learning. In *CVPR*, 2021. 7, 8
- 551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608

- 609 [32] Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel  
610 Menta, Andrew D Bagdanov, and Joost Van De Weijer.  
611 Class-incremental learning: survey and performance evalu-  
612 ation on image classification. *IEEE Transactions on Pattern  
613 Analysis and Machine Intelligence*, 45(5):5513–5533, 2022.  
614 1
- 615 [33] Michael McCloskey and Neal J Cohen. Catastrophic inter-  
616 ference in connectionist networks: The sequential learning  
617 problem. In *Psychology of Learning and Motivation*, pages  
618 109–165. 1989. 1
- 619 [34] Nicolas Michel, Giovanni Chierchia, Romain Negrel, and  
620 Jean-François Bercher. Learning representations on the unit  
621 sphere: Application to online continual learning. *arXiv  
622 preprint arXiv:2306.03364*, 2023. 2
- 623 [35] German I Parisi, Ronald Kemker, Jose L Part, Christopher  
624 Kanan, and Stefan Wermter. Continual lifelong learning with  
625 neural networks: A review. *Neural Networks*, 113:54–71,  
626 2019. 1
- 627 [36] Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz  
628 Kaiser, and Geoffrey Hinton. Regularizing neural networks  
629 by penalizing confident output distributions. *arXiv preprint  
630 arXiv:1701.06548*, 2017. 5
- 631 [37] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu,  
632 Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn  
633 without forgetting by maximizing transfer and minimizing  
634 interference. *arXiv preprint arXiv:1810.11910*, 2018. 3
- 635 [38] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lil-  
636 licrap, and Gregory Wayne. Experience replay for continual  
637 learning. In *NeurIPS*, 2019. 1, 2, 5, 6
- 638 [39] Amir Rosenfeld and John K Tsotsos. Incremental learn-  
639 ing through deep adaptation. *IEEE Transactions on Pattern  
640 Analysis and Machine Intelligence*, 42(3):651–663, 2018. 2
- 641 [40] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins,  
642 Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Raz-  
643 van Pascanu, and Raia Hadsell. Progressive neural networks.  
644 *arXiv preprint arXiv:1606.04671*, 2016.
- 645 [41] Joan Serra, Didac Suris, Marius Miron, and Alexandros  
646 Karatzoglou. Overcoming catastrophic forgetting with hard  
647 attention to the task. In *ICML*, 2018. 2
- 648 [42] Guocong Song and Wei Chai. Collaborative learning for  
649 deep neural networks. In *NeurIPS*, 2018. 2
- 650 [43] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu  
651 Sebe. Curriculum learning: A survey. *International Journal  
652 of Computer Vision*, 130(6):1526–1565, 2022. 4
- 653 [44] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon  
654 Shlens, and Zbigniew Wojna. Rethinking the inception ar-  
655 chitecture for computer vision. In *CVPR*, 2016. 5
- 656 [45] Gido M Van de Ven and Andreas S Tolias. Three scenar-  
657 ios for continual learning. *arXiv preprint arXiv:1904.07734*,  
658 2019. 1
- 659 [46] Laurens Van der Maaten and Geoffrey Hinton. Visualizing  
660 data using t-sne. *Journal of Machine Learning Research*, 9  
661 (11), 2008. 7
- 662 [47] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A  
663 comprehensive survey of continual learning: Theory, method  
664 and application. *arXiv preprint arXiv:2302.00487*, 2023. 1,  
665 3
- [48] Yujie Wei, Jiaxin Ye, Zhizhong Huang, Junping Zhang, and  
Hongming Shan. Online prototype learning for online con-  
tinual learning. In *ICCV*, 2023. 1, 2, 5, 6, 8
- [49] Li Yuan, Francis EH Tay, Guilin Li, Tao Wang, and Jiashi  
Feng. Revisiting knowledge distillation via label smoothing  
regularization. In *CVPR*, 2020. 5
- [50] Friedemann Zenke, Ben Poole, and Surya Ganguli. Contin-  
ual learning through synaptic intelligence. In *ICML*, 2017.  
2
- [51] Ying Zhang, Tao Xiang, Timothy M Hospedales, and  
Huchuan Lu. Deep mutual learning. In *CVPR*, 2018. 2
- [52] Xiatian Zhu, Shaogang Gong, et al. Knowledge distillation  
by on-the-fly native ensemble. In *NeurIPS*, 2018. 2