# CENTROID DISTANCE DISTILLATION FOR EFFECTIVE REHEARSAL IN CONTINUAL LEARNING

Daofeng Liu<sup>1\*</sup>, Fan Lyu<sup>2\*</sup>, Linyan Li<sup>3</sup>, Zhenping Xia<sup>1</sup>, Fuyuan Hu<sup>1\*\*</sup>

<sup>1</sup>Suzhou University of Science and Technology, <sup>2</sup>Tianjin University, <sup>3</sup>Suzhou Institute of Trade & Commerce

## ABSTRACT

Rehearsal, retraining on a stored small data subset of old tasks, has been proven effective in solving catastrophic forgetting in continual learning. However, due to the sampled data may have a large bias towards the original dataset, retraining them is susceptible to driving continual domain drift of old tasks in feature space, resulting in forgetting. In this paper, we focus on tackling the continual domain drift problem with centroid distance distillation. First, we propose a centroid caching mechanism for sampling data points based on constructed centroids to reduce the sample bias in rehearsal. Then, we present a centroid distance distillation that only stores the centroid distance to reduce the continual domain drift. The experiments on four continual learning datasets show the superiority of the proposed method, and the continual domain drift can be reduced. Our code is available at https://github.com/Daofeng-liu/CDD-R.

Index Terms- continual learning, distillation, centroids

### **1. INTRODUCTION**

Continual Learning (CL) is used to enable a machine learning system to learn from a sequence of tasks like humans [1], which has been applied to many applications, such as the recommender systems [2], medical research [3, 4] and clinical management decisions [5]. However, CL suffers from a wellknown obstacle in neural networks called catastrophic forgetting [6], which is the inability to effectively retain old knowledge after learning a new task. The objective of CL is to improve the adaptative ability to new knowledge over time without forgetting past knowledge. To address catastrophic forgetting, the rehearsal method [7, 8, 9, 10], storing a small portion of data for replay, is proven simple and effective against other methods including regularization methods [11, 12, 13] and parameter-isolation methods [14, 15, 16].

To reduce forgetting, as shown in Fig. 1(a), the existing rehearsal may select biased samples for replaying and storing the corresponding features for distillation. However, the biased samples have a poor capacity to represent the original



**Fig. 1**. (a) Previous rehearsal may select biased samples, which results in continual domain drift. (b) Our rehearsal selects representative data based on centroids and suppresses continual domain drift through centroid distance distillation.

domain. Retraining on these samples leads to unpredictable drifts in feature space, namely, continual domain drifts (a.k.a. semantic drifts) [17]. Moreover, distillation on the biased samples is memory-costly when the number of tasks grows, and the biased samples mislead the relative relationship between tasks, which results in over-fitting and indistinguishability of old tasks and slow learning of new tasks.

In this paper, we propose a simple yet effective Centroid Distance Distillation (CDD) for Rehearsal to mitigate the forgetting raised by continual domain drift. As shown in Fig. 1(b), CDD contains two main steps. (1) Centroidbased sampling: we first select representative samples via a proposed centroid caching mechanism, which builds an auto-updated cache for each centroid for storing the most representative samples and the caches can also help the update of centroids. (2) Centroid distance distillation: in contrast to distilling storage-costly features, we propose to only store the relative relationship, *i.e.*, the pairwise centroid distance, which is distilled to guide the replay of old tasks for less forgetting with negligible storage. We demonstrate that our method can better alleviate the catastrophic forgetting raised by continual domain drift with only stored centroid distance through the experiments on four popular CL datasets.

#### 2. METHOD

## 2.1. Preliminary: Continual domain drift in CL

CL can be formulated as learning from a sequence of datasets  $\{\mathcal{D}_1, \cdots, \mathcal{D}_T\}$  in order, where  $\mathcal{D}_t = \{(x_i, y_i)\}_{i=1}^{N_t}$  is the

<sup>\*</sup> Co-first author.

<sup>\*\*</sup> Corresponding author: fuyuanhu@mail.usts.edu.cn.

Supported by the Natural Science Foundation of China (No. 61876121).



**Fig. 2.** The framework of Centroid-Based Rehearsal. ① The model computes a set of centroids and candidate samples at the cache layer. ② After task training, memory is selected from the cache by centroid-based rehearsal sampling. ③ We compute the centroid relationship of memory data and distill the stored relationship to reduce continual domain drift.

dataset for the *t*-th task. We denote the model parameters as two main parts, the shared feature extraction layers  $h: x \to f$ and the task-specific linear classification layer  $g_t: f \to p$  for task *t*, where *f* and *p* are the deep feature and the predicted probability of the input image respectively. In CL, the model suffers from catastrophic forgetting due to the unavailability of data from past tasks. Rehearsal [7, 18], as an effective method to reduce forgetting, stores a small number of samples to memory  $\mathcal{M}$ , and the stored data will be retrained together with the current training to maintain the past task knowledge.

#### 2.2. Centroid-based Sampling for Rehearsal

Although the continual learning method based on rehearsal improves the ability to remember past knowledge, the memory size is far small compared to the original data set D, and the continual domain drift phenomenon will inevitably occur. Moreover, the existing may sample data with large biases such as outliers, which even worse the continual domain drift.

To sample representative samples, following Agg-Var [19, 20], we propose a centroid-based sampling method for rehearsal. A centroid c is the cluster center of the training domain [21, 22, 20]. Given a training data point (x, y), if the closest centroid c has a distance less than a given threshold  $\epsilon$ , this centroid is updated via:

$$\mathbf{c} = \frac{n \times \mathbf{c} + h(x)}{n+1},\tag{1}$$

where n is the number of the data points already represented by the centroid. Elsewise, a new centroid will be constructed directly using h(x). We consider that the data pass only once and then discard in a stream following [9, 18]. It is difficult to determine which sample should be selected if they have been discarded before the optimal centroid.

Thus, together with the centroid updating, we build centroid-aware caches as the medium to select the closest sample candidates for each centroid, and the cache will be discarded after the task end. For a centroid  $\mathbf{c}$ , its cache is

Algorithm 1 Update centroid and cac	che.
<b>Input:</b> Centroids $\{\mathbf{c}_1, \cdots, \mathbf{c}_N\}$ , Ce	entroid update numbers
$\{K_1, \cdots, K_N\}$ , model h, Cache	e $\mathcal{A}$ , Distance threshold
$\epsilon$ , Cache size $\gamma$ , data $(x, y)$	
<b>Output:</b> Updated $\{\mathbf{c}_1, \cdots, \mathbf{c}_{N'}\}$ and	d ${\mathcal A}$
1: $i^* = \arg\min_{i \in [1,N]} \ h(x) - \mathbf{c}_i\ $	<i>⊲Nearest centroid</i>
2: $d = \ h(x) - \mathbf{c}_{i^*}\ $	
3: if $d > \epsilon$ then	
4: $N' = N + 1, \mathbf{c}_{N'} = h(x)$	
5: <b>else</b>	
$6:  N' = N_{K}$	
7: $\mathbf{c}_{i^*} = \frac{\sum_{j=1}^{K} h(x_j \in \mathcal{A}_{\mathbf{c}_{i^*}}) + h(x)}{K + 1}$	<i>⊲ Update centroid</i>
8: end if	
9: if $K_{i^*} < \gamma$ then	
10: $\mathcal{A}_{\mathbf{c}_{i^*}} \leftarrow \mathcal{A}_{\mathbf{c}_{i^*}} \cup (x, y)$	$\triangleleft$ Update cache
11: <b>else</b>	
12: $j = \operatorname{argmax}_{i} \ \mathcal{A}_{\mathbf{c}_{i^{*}}} - \mathbf{c}_{i^{*}}\ $	
13: Remove $\mathcal{A}_{\mathbf{c},j}$ from $\mathcal{A}_{\mathbf{c}}$	
14: end if	

represented as  $\mathcal{A}_{\mathbf{c}} = \{(x_j, y_j)\}_{j=1}^{K}$ . The candidate data are the *K* nearest data in the data stream to each centroid and are replaced with nearer samples in the current batch. If the number of candidates for a centroid is less than a pre-defined  $\gamma$ , (x, y) is added to the cache directly. If the candidate number is equal to  $\gamma$ , the updated centroid compares the distance with  $\gamma + 1$  candidates and removes the farthest one.

After the training of the current task, we select the sample data from the cache into memory and delete the whole cache. To prevent bias in the centroids and the updated model, in Eq. (1), we also consider using cache to update centroids:

$$\mathbf{c} = \frac{\sum_{x' \in \mathcal{A}_{\mathbf{c}}} h(x') + h(x)}{|\mathcal{A}_{\mathbf{c}}| + 1},$$
(2)

where the centroid c is updated by the average of the features of the candidate data corresponding to this centroid (closest and distance less than  $\epsilon$ ). Easy to know, in Eq. (2), important centroids are updated more frequently than some outliers.

Based on the cache, the memory buffer is obtained by

$$\mathcal{M}^{t} = \bigcup_{i}^{|\mathcal{M}^{t}|} \left\{ (x_{i}, y_{i}) \sim P^{t}((x, y)) \right\}, \qquad (3)$$

where  $P_i^t = \frac{m_i}{\sum_{k=1}^n m_k} \in [0, 1]$  is the sampling probability of the *i*-th centroid and  $m_i$  is the total update frequency of centroid *i*. It is worth noting why we do not only sample from caches with larger confidence because we need to also keep the diversity to some extent. Otherwise, the stored samples will have larger biases than random selection

## 2.3. Distillation on centroid distance

Only selecting samples with diversity and representativeness for rehearsal based on the centroids is not enough to effectively solve the continual domain drift problem. Because in the process of continual learning, domain drift will blur the

Methods		Permuted MNIST			Split CIFAR	
Wethous	$A_T(\%)$	$F_T$	LTR	$A_T(\%)$	$F_T$	LTR
A-GEM [11]	$89.32 \pm 0.46$	$0.07\pm0.004$	$0.367 \pm 0.013$	$61.28 \pm 1.88$	$0.09 \pm 0.018$	$0.643 \pm 0.124$
ER [10]	$90.47 \pm 0.14$	$0.03 \pm 0.001$	$0.367 \pm 0.013$	$63.97 \pm 1.30$	$0.06 \pm 0.006$	$0.451 \pm 0.333$
MEGA [23]	$91.21 \pm 0.10$	$0.05\pm0.001$	$0.524 \pm 0.017$	$66.12 \pm 1.94$	$0.06\pm0.015$	$0.356 \pm 0.114$
DER [24]	$92.03 \pm 0.19$	$0.04\pm0.001$	$0.402 \pm 0.012$	$68.49 \pm 1.45$	$0.06\pm0.009$	$0.371 \pm 0.087$
ASER [9]	-	-	-	$65.53 \pm 1.89$	$0.07\pm0.007$	$0.544 \pm 0.133$
SCR [25]	$91.74 \pm 0.63$	$0.05\pm0.004$	$0.492 \pm 0.041$	$67.99 \pm 1.89$	$0.05\pm0.004$	$0.258 \pm 0.024$
MDMTR [17]	$91.97 \pm 0.23$	$0.05\pm0.002$	$0.521 \pm 0.022$	$66.38 \pm 1.63$	$0.05\pm0.006$	$0.377 \pm 0.076$
MDMTR+FD [17]	$93.97 \pm 0.15$	$0.03\pm0.002$	$0.283 \pm 0.019$	$69.20 \pm 1.60$	$0.04\pm0.010$	$0.283 \pm 0.099$
Ours	$92.22 \pm 0.22$	$0.04\pm0.002$	$0.484 \pm 0.019$	$69.65 \pm 1.55$	$0.04 \pm 0.013$	$0.192 \pm 0.094$
Ours+FD	$94.12 \pm 0.11$	$0.01 \pm 0.007$	$0.013 \pm 0.009$	$70.69 \pm 2.33$	$0.03 \pm 0.012$	$0.119 \pm 0.065$
Methods		Split CUB			Split AWA	
Wethods	$A_T(\%)$	$F_T$	LTR	$A_T(\%)$	$F_T$	LTR
A-GEM [11]	$61.82 \pm 3.72$	$0.08\pm0.021$	$0.456 \pm 0.174$	$44.95 \pm 2.97$	$0.05\pm0.014$	$0.178 \pm 0.082$
ER [10]	$73.63 \pm 0.52$	$0.01\pm0.005$	$0.001 \pm 0.001$	$53.27 \pm 4.05$	$0.02\pm0.030$	$0.014 \pm 0.015$
MEGA [23]	$80.58 \pm 1.94$	$0.01\pm0.017$	$0.002\pm0.002$	$54.28 \pm 4.84$	$0.05\pm0.040$	$0.070 \pm 0.114$
DER [24]	$76.56 \pm 2.48$	$0.01\pm0.015$	$0.025 \pm 0.018$	$50.70 \pm 4.91$	$0.04\pm0.040$	$0.063 \pm 0.094$
ASER [9]	$75.58 \pm 3.72$	$0.02\pm0.010$	$0.037 \pm 0.029$	$46.72 \pm 3.20$	$0.05\pm0.006$	$0.171 \pm 0.021$
SCR [25]	$81.43 \pm 1.97$	$0.01\pm0.007$	$0.007\pm0.009$	$54.35 \pm 2.68$	$0.02\pm0.012$	$0.022\pm0.010$
MDMT [17]	$83.06 \pm 4.39$	$0.20\pm0.028$	$0.015\pm0.023$	$58.20 \pm 2.51$	$0.02\pm0.011$	$0.035\pm0.025$
MDMT+FD [17]	$83.98 \pm 2.35$	$0.01\pm0.015$	$0.021 \pm 0.018$	$61.26 \pm 3.36$	$0.02\pm0.027$	$0.002 \pm 0.002$
Ours	$84.85 \pm 2.46$	$0.01\pm0.012$	$0.008 \pm 0.010$	$59.26 \pm 4.72$	$0.03 \pm 0.037$	$0.034 \pm 0.065$
Ours+FD	$85.75 \pm 1.99$	$0.01 \pm 0.004$	$0.004 \pm 0.006$	$61.92 \pm 2.94$	$0.02 \pm 0.027$	$0.013 \pm 0.002$

Table 1. Comparisons on four datasets, where the mean and std are over 5 seeds.

decision boundary of the old task and cause forgetting of the old knowledge. A naive way is to anchor the memory feature as mentioned via distillation [17, 24]. However, constraining the memory feature move may undermine the new task learning and the stored features are a large storage burden.

In this paper, we propose to only store the relative relationships among centroids, i.e., the Cosine similarity. In specific, we first calculate the similarity between any two centroids of the current task (intra- and inter-class), and the pairwise distances can be represented by a matrix  $\mathbf{W}^t$ , where

$$\mathbf{W}_{ij}^{t} = \frac{\mathbf{c}_{i} \cdot \mathbf{c}_{j}}{||\mathbf{c}_{i}|| \, ||\mathbf{c}_{j}||}.$$
(4)

The storage cost of this matrix is much smaller than the raw samples and their features.

In the learning of new tasks, we recalculate the centroid distance matrix  $\mathbf{W}^{t'}$  based on only the memory and distill it with the stored  $\mathbf{W}^{t}$  with a Centroid Distillation (CD) loss

$$\mathcal{L}_{\rm CD} = \sum_{t=1}^{k-1} \left\| \mathbf{W}^t - \mathbf{W}^{t'} \right\|^2.$$
 (5)

The centroid distance matrix is distilled in the following tasks to keep the decision boundary from blurring and drifting. The CD loss takes advantage of the information between the centroids to further suppress the continual domain drift, where diversity and representativeness will be kept. The relative relationships among centroids also have the classes of old tasks discriminative from each other in the current tasks.

#### **3. EXPERIMENTS**

#### 3.1. Dataset and experimental details

**Dataset.** *Permuted MNIST* is a variant of MNIST [26] dataset where the input pixels for each task have different random permutation as different tasks. *Split CIFAR* is a split version of the original CIFAR100 dataset [27], which splits 100 classes into 20 disjoint tasks, each containing 5 classes. *Split CUB* is a random splitting of the 200 classes of the CUB

Table 2. Ablation study on Split CIFAR.					
CD	FD	$A_{\rm T}(\%)$	$F_{\mathrm{T}}$	LTR	
-	-	$66.38 \pm 1.63$	$0.052 \pm 0.006$	$0.377 \pm 0.076$	
-	$\checkmark$	$68.97 \pm 2.21$	$0.040\pm0.009$	$0.254 \pm 0.031$	
$\checkmark$	-	$69.65 \pm 1.55$	$0.035\pm0.013$	$0.192 \pm 0.094$	
$\checkmark$	<b>√</b>	$70.69 \pm 2.33$	$0.027 \pm 0.012$	$0.119 \pm 0.065$	

dataset [28] into 20 disjoint tasks, each containing 10 classes. *Split AWA* is an incremental version of the AWA dataset [29] that splits 50 animal categories into 20 tasks, each with 5 different categories, but the categories are repeatable across tasks. We follow previous works [11, 10, 17] to conduct experiments on the above four datasets.

**Evaluation metric.** Average Accuracy  $(A_T)$  is the average of the accuracy of all tasks after the model is trained on the last task. Forgetting Measure  $(F_T)$  represents the the accuracy drop of old tasks after the model is trained on all tasks. Long-Term Remembering (LTR) [23] computes the accuracy drop of each task relative to when the task was first trained.

**Implementation detail.** Following [11, 17], we implement our methods with different backbones. For Permuted MNIST, we use a fc network with two hidden layers of 256 ReLU units. For Split CIFAR, we use a reduced resnet18 [30]. For Split CUB and Split AWA, we use a standard resnet18. The model is optimized using stochastic gradient descent with a mini-batch size of 10. For Permuted MNIST, Split CIFAR, Split CUB, Split AWA,  $\epsilon$  and  $\gamma$  are respectively set to (7,35), (8,20), (11,10), (7,35) via grid searching.

## 3.2. Experimental Results

**Main comparisons.** As shown in Table 1, our method compares with other SOTAs[31, 11, 10, 23, 24, 17, 9, 25] in three metrics. For  $A_T$ , our method shows a clear improvement compared to other methods on all four datasets. This indicates that our method can better suppress forgetting and reduce domain drift on all seen tasks. For  $F_T$  and LTR, our method also has an advantageous position, which shows that



Fig. 3. T-SNE of task 1 on Permuted MNIST after the learning of task 1 and 6. Trangles are the centroids of each class.

ĺ	<u> Table 3.</u>	The effect of	<u>the</u>	<u>centroid</u>	distance	$\epsilon$

$\epsilon$	$A_{\rm T}(\%)$	$F_{\mathrm{T}}$	-	LTR		
6 6	$9.76 \pm 1.59$	$0.034 \pm 0.034$	007 0.175	$2 \pm 0.044$		
7 7	$0.16 \pm 1.96$	$0.029 \pm 0.029$	012 0.148	$8 \pm 0.072$		
8 7	$0.69 \pm 2.33$	$0.027 \pm 0$	.012 <b>0.11</b>	$9 \pm 0.065$		
9 6	$9.81 \pm 1.76$	$0.030 \pm 0.030$	008 0.143	$3 \pm 0.048$		
Table 4. Sampling comparisons on Split CIFAR.						
Methods (w	ithout distillation)	$  A_T(\%)$	$F_{\rm T}$	LTR		
Ri	ng buffer	$66.38 \pm 1.63$	$0.052\pm0.006$	$0.377\pm0.076$		
MoF		$66.58 \pm 1.75$	$0.053 \pm 0.010$	$0.359 \pm 0.106$		
	GSS	$62.06 \pm 3.58$	$0.115\pm0.021$	$0.912 \pm 0.183$		
	Ours	$69.18 \pm 0.74$	$0.039 \pm 0.008$	$0.236 \pm 0.061$		

our method can gain long-term memory by reducing the domain drift. Our method achieves the best results compared to SOTAs without FD (Feature Distillation) loss [17]. With FD loss, our method can be further improved.

Ablation study. We then explore the impact of each configuration. We show the ablation experiments on the Split CIFAR dataset in Table 2. The first row of the table shows the performance of our experimental baseline, without adding any configuration. When we add CD loss or FD loss, the experimental results can be clearly improved. The experimental performance is further improved when CD loss and FD loss are used together. In Table 3, We also conduct an experimental study on the choice of the hyper-parameter centroid distance  $\epsilon$ . The larger the centroid distance, the less the class trains out of the centroid, and the smaller the centroid distance, the more the centroid is obtained. Easy to observe, when a suitable centroid distance is chosen, the experimental effect can reach optimal performance. In Fig. 4, we set three memory sizes on the Split CIFAR dataset to explore the effect of changing the cache size  $\gamma$  on the experimental results. We can see that for different memory sizes, the optimal cache sizes are different. Moreover, a larger memory size always means better performance with an appropriate cache size.

**Sampling strategy comparisons.** In Table 4, we compare with other sampling methods. Ring buffer is classical online random sampling methods. Mean-of-Feature (MoF) [7] samples the data closest to the mean by calculating the mean of the features for all data in each class, which means larger stor-



Fig. 4. The effect of the cache size  $\gamma$  with different memory size settings. Memory size is the number of samples per class.

age is needed. Gradient-Based Sample Selection (GSS) [32] diversifies the gradients of the samples in the memory buffer. In contrast, our centroid-based sampling achieves the best results without much storage.

**Continual domain drift observation** We explore the continual domain drift phenomenon in continual learning by the t-SNE visualization. In Fig. 3, we show the features distribution of Task 1 on Permuted MNIST at the end of Task 1 phase and Task 6 phase. The DER and SCR methods do not suppress the occurrence of continual domain drift and thus lead to the forgetting of past knowledge. Our method can effectively mitigate the continual domain drift by centroid-based sampling and centroid distance distillation, and the feature distribution remains relatively stable after the end of Task 6. When our method is combined with FD loss, the continual domain drift phenomenon is further reduced.

## 4. CONCLUSION

In this paper, we tackled catastrophic forgetting from the perspective of continual domain drift. We first constructed centroids for each class in an online fashion. Then, with guidance from centroids, we stored representative data points to reduce the dataset bias. We also stored the relative centroid distance, which is used to distill for long-term remembering. The experimental results on four continual learning datasets show the superiority of our method.

#### 5. REFERENCES

- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter, "Continual lifelong learning with neural networks: A review," *Neural Networks*, vol. 113, pp. 54–71, 2019.
- [2] Ivens Portugal, Paulo Alencar, and Donald Cowan, "The use of machine learning algorithms in recommender systems: A systematic review," *Expert Systems with Applications*, vol. 97, pp. 205–227, 2018.
- [3] Jeffrey De Fauw, Joseph R Ledsam, Bernardino Romera-Paredes, Stanislav Nikolov, Nenad Tomasev, Sam Blackwell, Harry Askham, Xavier Glorot, Brendan O'Donoghue, Daniel Visentin, et al., "Clinically applicable deep learning for diagnosis and referral in retinal disease," *Nature medicine*, vol. 24, no. 9, pp. 1342–1350, 2018.
- [4] Michael D Abràmoff, Philip T Lavin, Michele Birch, Nilay Shah, and James C Folk, "Pivotal trial of an autonomous aibased diagnostic system for detection of diabetic retinopathy in primary care offices," *NPJ digital medicine*, vol. 1, no. 1, pp. 1–8, 2018.
- [5] Mohammad M Ghassemi, Tuka Alhanai, M Brandon Westover, Roger G Mark, and Shamim Nemati, "Personalized medication dosing using volatile data streams," in *AAAIW*, 2018.
- [6] Michael McCloskey and Neal J Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in *Psychology of learning and motivation*, vol. 24, pp. 109–165. Elsevier, 1989.
- [7] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert, "icarl: Incremental classifier and representation learning," in *CVPR*, 2017, pp. 2001–2010.
- [8] Qing Sun, Fan Lyu, Fanhua Shang, Wei Feng, and Liang Wan, "Exploring example influence in continual learning," *NIPS*, 2022.
- [9] Dongsub Shim, Zheda Mai, Jihwan Jeong, Scott Sanner, Hyunwoo Kim, and Jongseong Jang, "Online class-incremental continual learning with adversarial shapley value," vol. 35, no. 11, pp. 9630–9638, 2021.
- [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.
- [11] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny, "Efficient lifelong learning with agem," arXiv preprint arXiv:1812.00420, 2018.
- [12] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al., "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, vol. 2, no. 7, 2015.
- [13] Kaile Du, Fan Lyu, Fuyuan Hu, Linyan Li, Wei Feng, Fenglei Xu, and Qiming Fu, "Agcn: augmented graph convolutional network for lifelong multi-label image recognition," in *ICME*. IEEE, 2022, pp. 01–06.
- [14] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang, "Lifelong learning with dynamically expandable networks," arXiv preprint arXiv:1708.01547, 2017.

- [15] Arun Mallya, Dillon Davis, and Svetlana Lazebnik, "Piggyback: Adapting a single network to multiple tasks by learning to mask weights," in *ECCV*, 2018, pp. 67–82.
- [16] Ju Xu and Zhanxing Zhu, "Reinforced continual learning," *NIPS*, vol. 31, 2018.
- [17] Fan Lyu, Shuai Wang, Wei Feng, Zihan Ye, Fuyuan Hu, and Song Wang, "Multi-domain multi-task rehearsal for lifelong learning," in AAAI, 2021, vol. 35, pp. 8819–8827.
- [18] Rishabh Tiwari, Krishnateja Killamsetty, Rishabh Iyer, and Pradeep Shenoy, "Gcr: Gradient coreset based replay buffer selection for continual learning," in CVPR, 2022, pp. 99–108.
- [19] Ali Ayub and Alan R Wagner, "Centroid based concept learning for rgb-d indoor scene classification," *arXiv preprint arXiv:1911.00155*, 2019.
- [20] Ali Ayub and Alan R Wagner, "Cognitively-inspired model for incremental learning using a few examples," in *CVPRW*, 2020, pp. 222–223.
- [21] Xiaoyu Tao, Xiaopeng Hong, Xinyuan Chang, Songlin Dong, Xing Wei, and Yihong Gong, "Few-shot class-incremental learning," in *CVPR*, 2020, pp. 12183–12192.
- [22] Jian Liang, Ran He, Zhenan Sun, and Tieniu Tan, "Distant supervised centroid shift: A simple and efficient approach to visual domain adaptation," in CVPR, 2019, pp. 2975–2984.
- [23] Yunhui Guo, Mingrui Liu, Tianbao Yang, and Tajana Rosing, "Learning with long-term remembering: Following the lead of mixed stochastic gradient," *ICLR*, 2020.
- [24] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara, "Dark experience for general continual learning: a strong, simple baseline," *NIPS*, vol. 33, pp. 15920–15930, 2020.
- [25] Zheda Mai, Ruiwen Li, Hyunwoo Kim, and Scott Sanner, "Supervised contrastive replay: Revisiting the nearest class mean classifier in online class-incremental continual learning," pp. 3589–3599, 2021.
- [26] Yann LeCun, "The mnist database of handwritten digits," http://yann. lecun. com/exdb/mnist/, 1998.
- [27] Alex Krizhevsky, Geoffrey Hinton, et al., "Learning multiple layers of features from tiny images," 2009.
- [28] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie, "The caltech-ucsd birds-200-2011 dataset," 2011.
- [29] Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling, "Learning to detect unseen object classes by between-class attribute transfer," in *CVPR*. IEEE, 2009, pp. 951–958.
- [30] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [31] David Lopez-Paz and Marc'Aurelio Ranzato, "Gradient episodic memory for continual learning," *NIPS*, vol. 30, 2017.
- [32] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio, "Gradient based sample selection for online continual learning," *NIPS*, vol. 32, 2019.