

# Learning Equi-angular Representations for Online Continual Learning

Minhyuk Seo<sup>1,†</sup> Hyunseo Koh<sup>1</sup> Wonje Jeung<sup>1</sup> Minjae Lee<sup>1</sup> San Kim<sup>1</sup> Hankook Lee<sup>2,3</sup>  
Sungjun Cho<sup>2</sup> Sungik Choi<sup>2</sup> Hyunwoo Kim<sup>4,\*</sup> Jonghyun Choi<sup>5,\*</sup>

<sup>1</sup>Yonsei Univ. <sup>2</sup>LG AI Research <sup>3</sup>Sungkyunkwan Univ. <sup>4</sup>Zhejiang Lab <sup>5</sup>Seoul National Univ.  
{dbd0508, specific0924, reccos1020, nasmik419}@yonsei.ac.kr  
khs8157@gmail.com, {sungik.choi, sungjun.cho}@lgresearch.ai  
hankook.lee@skku.edu, hwkim@zhejianglab.com, jonghyunchoi@snu.ac.kr

## Abstract

Online continual learning suffers from an underfitted solution due to insufficient training for prompt model update (e.g., single-epoch training). To address the challenge, we propose an efficient online continual learning method using the neural collapse phenomenon. In particular, we induce neural collapse to form a simplex equiangular tight frame (ETF) structure in the representation space so that the continuously learned model with a single epoch can better fit to the streamed data by proposing preparatory data training and residual correction in the representation space. With an extensive set of empirical validations using CIFAR-10/100, TinyImageNet, ImageNet-200, and ImageNet-1K, we show that our proposed method outperforms state-of-the-art methods by a noticeable margin in various online continual learning scenarios such as disjoint and Gaussian scheduled continuous (i.e., boundary-free) data setups. Code is available at <https://github.com/yonseivnl/earl>.

## 1. Introduction

A growing interest in continuous learning (CL) involves training the model using continuous data streams. Mostly, CL research has focused on the *offline* scenario that assumes the model can be trained in multiple epochs for the current task [9, 51, 64]. However, substantial storage and computational complexity are required to store all data to train a model for multiple epochs. Recently, there has been significant interest in *online* CL as a more realistic set-up with a lower computational overhead of allowing a single training pass through the data stream [1, 8, 34]. Learning a model in streamed data by a single training pass is challenging, since the temporal distribution at each intermediate

time point is likely imbalanced, even if the overall distribution of a dataset is balanced. Imbalanced data distributions would cause several problems, such as bias towards the major classes [28, 72] and the hindrance to generalization [63].

Recently, *minority collapse* [18], the phenomenon in which the angles between classifier vectors for minor classes (i.e., the classes that have a relatively small number of samples) become narrow, has been proposed as a fundamental issue in training with imbalanced data, making the classification of minor classes considerably more challenging. In contrast, for balanced data sets, it was proven that classifier vectors and the last layer activations for all classes converge into an optimal geometric structure, named the simplex *equiangular tight frame* (ETF) structure, where all pairwise angles between classes are equal and maximally widened when using cross entropy (CE) [27, 39, 62, 76] or mean squared error (MSE) [41, 50, 58, 75] loss. This phenomenon is called *neural collapse* [43].

Although neural collapse naturally occurs only in balanced training, several recent studies attempted to induce neural collapse in imbalanced training to address the minority collapse problem using a fixed ETF classifier [65, 73]. When employing a fixed ETF classifier, the output features are pulled toward the corresponding classifier vector during training, since the classifier is fixed. More recently, research has also been extended to induce neural collapse in offline CL using a fixed ETF classifier [66].

However, unlike offline CL, there are a number of challenges in inducing neural collapse in online CL, even with a fixed ETF classifier. The prerequisite for neural collapse is reaching the *terminal phase of training* (TPT) by sufficient training [43]. In offline CL, the model can reach the TPT phase for each task by multi-epoch training. In contrast, a single-pass training constraint often prevents the online CL from reaching TPT. As shown in Fig. 1, the offline CL quickly reaches TPT shortly after the arrival of novel task data, while online CL (vanilla ETF) does not.

<sup>†</sup>: Work done while interning at LG AI Research.

\*: Indicates corresponding authors.

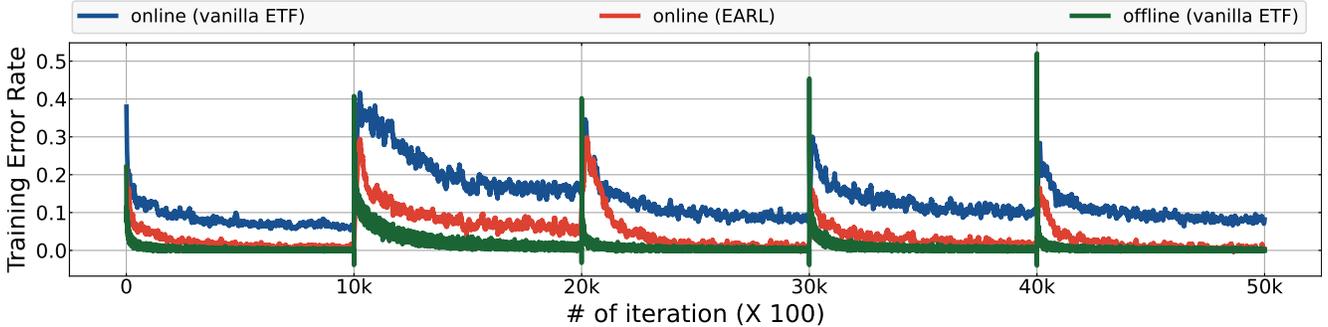


Figure 1. Comparison of training error rates between online CL and offline CL in the CIFAR-10 disjoint setup, where two novel classes are added every 10k samples. Vanilla ETF refers to a method where both preparatory data training and residual correction are removed from our proposed EARL.

Recently, the importance of anytime inference in online CL has been emphasized [7, 21, 34, 45], since a model should be available for inference not only at the end of a task, but also at any point during training to be practical for real-world applications. Hence, not only reaching TPT but also achieving faster convergence is necessary when using neural collapse in online CL.

However, we observe that the phenomenon in which the features of the new class become biased toward the features of the existing classes [6] hinders fast convergence of the last-layer features into the ETF structure, which we call the ‘bias problem.’ When features of old and new classes overlap (*i.e.*, biased) and are trained with the same objective, well-clustered features of old classes disperse (or perturb) [6], leading to destruction of the ETF structure formed by features of the old classes.

**Contributions.** To address the bias problem, we train the model to distinguish out-of-distribution (OOD) samples from existing classes’ samples. Specifically, we synthesize *preparatory data* that play the role of OOD data, by applying negative data augmentation [30, 57, 60] to samples from existing classes, refer to Sec. 4.2 for details. Since samples from new classes are OOD in the perspective of existing classes, training with preparatory data encourages the representation of new classes to be distinguished from existing classes in advance, thereby mitigating the bias problem. This promotes fast and stable convergence into the ETF structure.

Despite these efforts, however, the continuous stream of new samples cause ongoing data distribution shifts, which prevent the model from reaching the TPT and fully converging to the ETF structure. To address this, we propose to store the residuals between the target ETF classifier and the features during training. During inference, we correct the inference output using the stored residual to compensate for insufficient convergence in the training.

By accelerating convergence by preparatory data training and additional correction using residual, we noticeably

improve anytime inference performance than state of the arts. We name our method **Equi-Angular Representation Learning (EARL)**. We empirically demonstrate the effectiveness of our framework on CIFAR-10, CIFAR-100, Tiny-ImageNet, Imagenet-200, and ImageNet-1K. In particular, our framework outperforms various CL methods by a significant margin (+4.0% gain of  $A_{auc}$  in ImageNet-1K).

We summarize our contributions as follows:

- Proposing to induce neural collapse for online CL.
- Proposing ‘preparatory data training’ to address the ‘bias problem’ that the new classes are biased toward the existing classes, promoting faster induction of neural collapse.
- Proposing ‘residual correction’ scheme to compensate for not fully reaching neural collapse at inference to further improve anytime inference accuracy.

## 2. Related work

**Continual learning methods.** Various continual learning methods are being researched to prevent forgetting past tasks, broadly categorized into replay, parameter isolation, and regularization. Replay methods involve storing a small portion of data from previous tasks in episodic memory [1, 3, 24, 34, 70] or storing a generative model trained on data from previous tasks [48, 56]. By replaying the samples stored in episodic memory or generated from the stored generative model, the model prevents forgetting past tasks during subsequent learning of novel tasks. Furthermore, Boschini et al. [4], Buzzega et al. [5], Li and Hoiem [38], Wu et al. [64] used replay samples to distill information about past tasks.

Regularization methods [33, 37] apply a penalty to the change of important model parameters during the process of learning new tasks, allowing the model to retain information about previous tasks. Parameter isolation methods [11, 53, 74] expand the network by allocating specific layers for each task. This enables the network to store information about individual tasks and preserves their knowl-

edge without forgetting.

**Neural collapse.** Neural collapse is a phenomenon in which the activations of the last layer and the classifier vectors form a simplex *equiangular tight frame* (ETF) structure at the terminal phase of training (TPT) in a balanced dataset. [43]. Neural collapse has been demonstrated as the global optimum of balanced training using cross entropy (CE) loss [27, 39, 62, 76] and MSE [41, 50, 58, 75] loss functions, within a simplified model focused solely on last-layer optimization. Inducing neural collapse in imbalanced datasets poses challenges due to minority collapse [18] where minor classes are not well distinguished.

However, using a fixed ETF classifier, it is empirically and theoretically shown that neural collapse is induced even in imbalanced datasets [65]. Continual learning also needs to address imbalanced data since there is an imbalance in the data between novel classes and existing classes. Therefore, in offline CL, NC-FSCIL [47, 66] used a fixed ETF classifier to induce neural collapse. Meanwhile, online CL often fails to induce neural collapse compared to offline CL, *e.g.*, FSCIL, since it lacks sufficient training in multiple epochs, causing failure to reach TPT.

Please refer to the supplementary material for more comprehensive literature reviews including topics related to anytime inference.

### 3. Preliminaries

We here describe the problem statement for the online continual learning (Sec. 3.1), neural collapse (Sec. 3.2), and the equiangular tight frame classifier (Sec. 3.3).

#### 3.1. Problem Statement of Online CL

Continual learning (CL) aims to learn from a stream of data, rather than a fixed dataset as in standard learning. Formally, given a sequence of tasks  $\mathcal{T} = (T_1, T_2, \dots)$  where each task  $T_i$  is a training dataset  $D_i = \{(x_1^{(i)}, y_1^{(i)}), (x_2^{(i)}, y_2^{(i)}), \dots\}$ , an offline CL algorithm  $A_{\text{CL}}$  updates the model parameter  $\theta$  based on the current task  $D_k$ , *i.e.*,  $\theta_k = A_{\text{CL}}(\theta_{k-1}, D_k)$ , starting from the initial parameter  $\theta_0$ .

To avoid the forgetting issue in CL, many CL setups allow the use of episodic memory  $M_k$ , which is a limited-size subset of training data from previous tasks, *i.e.*,  $(\theta_k, M_k) = A_{\text{CL}}(\theta_{k-1}, M_{k-1}, D_k)$ . The objective is to minimize the error of  $\theta_k$  in all observed tasks  $\{T_i\}_{i=1}^k$ .

Unlike offline CL where all training samples in  $D_k$  are given as input, in online CL, the input is provided as a stream of samples  $(x_1^{(k)}, y_1^{(k)}), (x_2^{(k)}, y_2^{(k)}), \dots$ . Therefore, an online CL algorithm  $A_{\text{OCL}}$  is defined as:

$$(\theta_{k,t}, M_{k,t}) = A_{\text{OCL}}\left(\theta_{k,t-1}, M_{k,t-1}, (x_t^{(k)}, y_t^{(k)})\right), \quad (1)$$

with the same objective of minimizing the error of  $\theta_{k,t}$  on  $\{T_i\}_{i=1}^k$ . Since the algorithm does not have access to pre-

vious samples  $\{(x_s^{(k)}, y_s^{(k)}), s < t\}$  at time  $t$ , multi-epoch training with  $D_k$  is unavailable in online CL.

#### 3.2. Neural Collapse

Neural collapse [43] is a phenomenon of the penultimate features after the convergence of training on a balanced dataset. When neural collapse (NC) occurs, the collection of  $K$  classifier vectors  $\mathbf{W}_{\text{ETF}} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K] \in \mathbb{R}^{d \times K}$  forms a simplex equiangular tight frame (ETF), which satisfies:

$$\mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1, & i = j \\ -\frac{1}{K-1}, & i \neq j \end{cases}, \quad \forall i, j \in [1, \dots, K], \quad (2)$$

where  $K - 1 \leq d$ , and the penultimate feature of a training sample collapses into an ETF vector  $\mathbf{w}_i$ .

#### 3.3. Equiangular Tight Frame (ETF) Classifier

Inspired by neural collapse as described in Sec. 3.2, a fixed ETF classifier has been utilized for inducing neural collapse in imbalanced datasets [65, 66, 76]. Here, the classifier is initialized by the ETF structure  $\mathbf{W}_{\text{ETF}}$  at the beginning of training and fixed during training to induce the penultimate feature  $f(x)$  to converge to the ideal balanced scenario. For training  $f(x)$ , it is only required to attract  $f(x)$  to the corresponding classifier vector for convergence, since the classifier is fixed during training. Therefore, following Yang et al. [65], we use the *dot regression* (DR) loss as a training objective, as it shows to outperform cross entropy (CE) loss when using a fixed ETF classifier in imbalanced datasets [65].<sup>1</sup> The DR loss can be written as follows:

$$\mathcal{L}_{\text{DR}}(\hat{f}(x), y; \mathbf{W}_{\text{ETF}}) = \frac{1}{2} \left( \mathbf{w}_y \hat{f}(x) - 1 \right)^2, \quad (3)$$

where  $\hat{f}(x) = f(x) / \|f(x)\|_2$  is the  $L_2$  normalized feature of the model  $f$ ,  $y$  is the label of the input  $x$ , and  $\mathbf{w}_y$  is a classifier vector in  $\mathbf{W}_{\text{ETF}}$  for the label  $y$ .

### 4. Approach

Despite the success of the fixed ETF classifier in both imbalanced training [65, 73] and offline CL [66], the ETF classifier has not yet been explored for online CL due to the necessity of sufficient training for neural collapse. To be specific, streamed data are trained only once in online CL, which makes it harder to induce neural collapse than in offline CL which supports multi-epoch training.

To learn a better converged model without multi-epoch training for online CL, we propose two novel methods, each for the training phase and the inference phase, respectively. In the training phase, we accelerate convergence by proposing *preparatory data training* (Sec. 4.2). In the inference

<sup>1</sup>If the training objective includes a contrastive term between different classes like cross entropy, it could cause incorrect gradients [65].

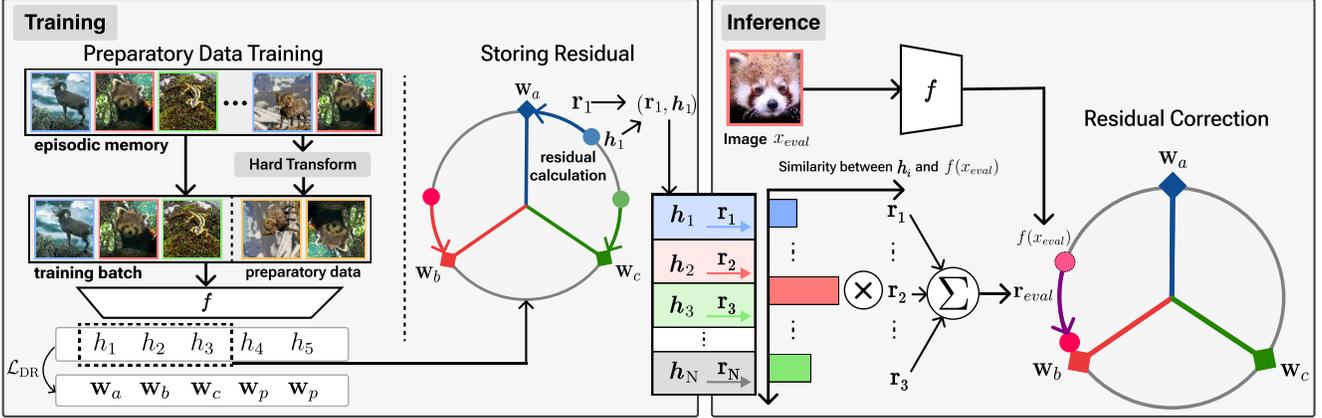


Figure 2. Overview of EARL.  $w_i$  denotes the ETF classifier vector for class  $i$ .  $h_a$  denotes the output of the model. The colors of the data denote the class to which the data belong. The arrow  $r_i$  denotes the residual between the last layer activation  $h_i$  and the classifier vector  $w_i$  for class  $i$ . During training, both memory and preparatory data are used for replaying, and the residuals between  $h_i$  and  $w_i$  are stored in feature-residual memory. During inference, using the similarity between  $f(x_{eval})$  and  $h_i$  in feature-residual memory,  $r_{eval}$  is obtained by a weighted sum of  $r_i$ 's. Finally, by adding  $r_{eval}$ ,  $f(x_{eval})$  is corrected. The purple arrow indicates ‘residual correction’ (Sec. 4.3).

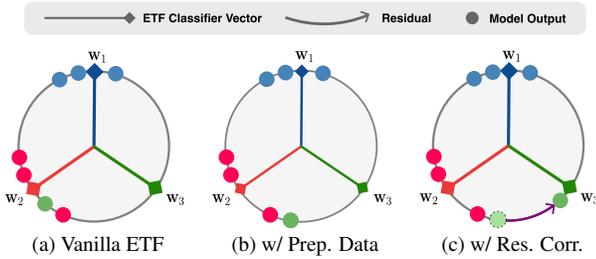


Figure 3. Illustrative effects for each component of EARL. (a) In online CL, features of novel classes are biased towards the features of the previous class. (b) By training with preparatory data (**w/ Prep. Data**, Sec. 4.2), we address the bias problem. (c) In inference, for features that do not fully converge to an ETF classifier, we add residuals (**w/ Res. Corr.**, Sec. 4.3) to features that have not yet reached the corresponding classifier vectors, making features aligned with them. Purple arrow: the ‘residual correction’, Colors: classes.

phase, we propose to correct the remaining discrepancy between the classifiers and the features using *residual correction* (Sec. 4.3). We illustrate an overview of EARL in Fig. 2 and the effect of EARL in Fig. 3.

#### 4.1. Inducing Neural Collapse for Online CL

For online CL, we first try to induce neural collapse with a fixed ETF classifier:  $\mathbf{W}_{ETF} \in R^{d \times K}$  where  $d$  is the dimension of the embedding space and  $K$  is the number of ETF vectors. While prior works [65, 66] use prior knowledge of the total number of classes that will be encountered for setting  $K$ , it is unrealistic to know the knowledge (*i.e.*, the exact number of classes) as it evolves continuously over time under realistic CL scenarios.

To address the challenge, we propose using the maximum possible number of classifier vectors for  $K$ . Based

on the simplex ETF property that the maximum number of ETF vectors is  $d + 1$  [20, 65], we define our ETF classifier  $\mathbf{W}_{ETF} \in R^{d \times (d+1)}$  by letting  $K := d + 1$ .

#### 4.2. Preparatory data at training

Since novel classes continuously arrive in CL, both the data from previous tasks ( $x_{old}$ ) and the current task ( $x_{new}$ ) coexist. While  $\hat{f}(x_{old})$  are placed closer to their corresponding ETF classifier vectors  $w_{new}$  by training,  $\hat{f}(x_{new})$  are biased toward the cluster of  $\hat{f}(x_{old})$ , as we can see in Fig. 4-(a), where  $\hat{f}(x_{new})$  and  $\hat{f}(x_{old})$  are the outputs of the model for inputs  $x_{new}$  and  $x_{old}$ , respectively. We call this ‘*bias problem*.’ When  $\hat{f}(x_{new})$  and  $\hat{f}(x_{old})$  overlap due to bias and are optimized with the same objective function, training on the new class interferes with the representations of the old classes, dispersing the well-clustered  $\hat{f}(x_{old})$  [6]. It destroys the ETF structure formed by  $\hat{f}(x_{old})$ , which hinders convergence towards neural collapse. Although the bias problem exists even without a fixed ETF classifier [6, 13, 64], it poses a greater problem when using a fixed ETF classifier, since the dispersed feature of old classes has to be restored to the corresponding fixed classifier vector. In contrast, when using a learnable classifier, the dispersed representation of the old classes does not have to be restored to the original position, since the learnable classifier vector will be re-optimized from the dispersed feature. However, the angles between the classifier vector may become narrow, *i.e.*, minority collapse occurs.

To accelerate convergence in the ETF structure, we propose to prevent novel classes from being biased towards existing classes when introduced, mitigating the bias problem. Specifically, we train the model to avoid making predictions in favor of the existing classes for images that do not be-

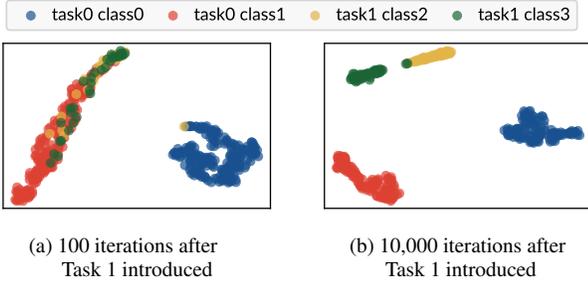


Figure 4. t-SNE visualization of ‘bias-problem’ in data distributions (class 0 to 3). (a) Only after 100 iterations of training after task 1 appears, learning is likely insufficient, and we can see that the features of new classes (class 2, 3) are biased towards the feature cluster of the existing class (*i.e.*, class 1). (b) With more training iterations (10,000 iter), the features are well clustered by class.

long to them. To this end, we propose to use preparatory data  $x_p$  that is different from existing classes, obtained by transforming the samples in the episodic memory. By training with the preparatory data so that their representations differ from existing classes, we prevent biased predictions towards existing classes when a new class arrives.

For preparatory data to have a different representation from existing classes when trained, their image semantics should be distinguishable from existing classes, provided that they contain enough semantic information to be considered as an image (*i.e.*, not noise). There are various approaches to synthesize images with modified semantics, such as using generative models [12, 23, 68] or negative transformations [57, 60]. We use negative transformations, as generative models are expensive in computation and storage [15], which is undesirable for CL scenarios with limited memory and computations. On the contrary, negative transformations are relatively efficient in computation and storage [42] and is reported to create images with semantics different from the original image [60]. We empirically observed that data generated through negative transformations have a distinct semantic class from the original class while preserving its semantic information as an image (*i.e.*, not as random noise), and we summarize the results in the supplementary material.

Formally, for a set of existing classes  $Y$  and the set of possible transformations  $G$ , we randomly select  $y \in Y$  and  $g \in G$ , and randomly retrieve a sample of class  $y$  from memory and apply the transformation  $g$  to obtain *preparatory data*  $x_p$ . We assign labels of the unseen classes to the preparatory data by a mapping function  $m : Y \times G \rightarrow Y'$  where  $Y'$  denotes a set of unseen classes, *i.e.*,  $Y' = \{y | y \notin Y, 1 \leq y \leq K\}$  and  $K$  is the total number of classifier vectors in  $\mathbf{W}_{\text{ETF}} \in \mathbb{R}^{d \times K}$ . Thus, preparatory data  $x_p$  from class  $y$  and transformation  $g$  is pulled towards the classifier vector  $\mathbf{w}_p$ , where  $p = m(y, g)$ . When a new class  $y_{\text{new}}$

is added to  $Y$ , we update  $m$  by randomly assigning a new mapping  $m(y_{\text{new}}, g)$  for  $y_{\text{new}}$  and  $g \in G$ .

We compose the transformation set  $G$  as a negative transformation that can modify semantic information (*i.e.*, change the label), used in the self-supervised literature [19, 22] and out-of-distribution (OOD) detection [30, 57, 60]. Specifically, we use rotation by 90, 180, and 270 degrees [19, 22] as our negative transformation. Since the vertical information of the image is more important than the horizontal information [26, 59], *e.g.*, images of a car facing left or right are common, but a car flipped upside down is very rare compared to a car standing right, rotation by a large angle causes loss of semantic information of images [60]. Thus, we can use rotated data as semantically different images from images of existing classes, *i.e.*, suitable for preparatory data. A more detailed analysis of various negative transformations is provided in the Supple.

We jointly train the model using real data (*i.e.*, samples of existing classes) and preparatory data. Formally, we write the objective for the model parameter  $\theta$  as:

$$\arg \min_{\theta} \left[ \mathbb{E}_{(x,y) \sim \mathcal{M}} \mathcal{L}_{\text{DR}} \left( \hat{f}_{\theta}(x), y; \mathbf{W}_{\text{ETF}} \right) + \lambda \mathbb{E}_{(x',y') \sim \mathcal{M}, g \sim G} \mathcal{L}_{\text{DR}} \left( \hat{f}_{\theta}(g(x')), m(y', g); \mathbf{W}_{\text{ETF}} \right) \right], \quad (4)$$

where  $\mathcal{M}$  denotes episodic memory,  $(x, y)$  and  $(x', y')$  denote the samples in  $\mathcal{M}$ ,  $G$  denotes the set of possible transformations,  $g$  denotes a transformation in  $G$ ,  $m$  denotes the mapping function, and  $\lambda$  is a hyperparameter for balancing real data and preparatory data.

### 4.3. Residual correction at inference

Despite preparatory data training that accelerates convergence towards ETF, in online CL, new samples in a data stream hinder models from reaching the Terminal Phase of Training (TPT) and fully converging to the ETF during single-epoch training. When the output of the model  $f(x)$  does not fully converge to the ETF classifier, the model would not perform well at all times, leading to poor anytime inference accuracy ( $A_{\text{AUC}}$ ).

To address this issue, we want to correct the residual between  $f(x)$  and the corresponding classifier vector  $\mathbf{w}_y$  during inference, where  $y$  is label of input  $x$ .

However, the discrepancy between the prediction on test data  $x_{\text{eval}}$  and the ground-truth ETF classifier vector is not available in the inference phase. In a similar situation, ResMem [67] stores the residual obtained from training data and uses them during inference, but they assume two-stage learning algorithms, *i.e.*, standard training and inference rather than incremental learning and anytime inference.

Inspired by ResMem [67], which uses residuals from the training in the inference stage, we propose to use the residual obtained from the training process to compensate for the

Method	CIFAR-10				CIFAR-100			
	Disjoint		Gaussian-Scheduled		Disjoint		Gaussian-Scheduled	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$						
EWC (Kirkpatrick <i>et al.</i> , 2017)	75.25±0.78	60.80±2.20	59.62±0.31	64.24±1.97	52.08±0.83	41.55±0.85	38.22±0.50	42.52±0.58
ER (Rolnick <i>et al.</i> , 2019)	75.94±0.86	63.56±1.32	60.13±0.56	64.81±2.70	52.95±1.25	42.82±0.05	41.12±0.56	42.74±1.09
ER-MIR (Aljundi <i>et al.</i> , 2019)	75.89±1.02	61.93±0.93	60.39±0.48	61.64±3.86	52.93±1.44	42.47±0.13	41.19±0.63	42.93±1.18
REMIN (Hayes <i>et al.</i> , 2020)	69.55±0.91	53.34±1.01	58.01±0.72	59.27±1.86	40.87±0.76	36.17±1.83	23.40±2.25	28.78±1.71
DER++ (Buzzega <i>et al.</i> , 2020)	74.78±0.72	59.20±0.95	59.44±0.28	66.11±1.80	38.16±1.57	38.55±2.11	29.38±2.58	38.20±3.13
SCR (Mai <i>et al.</i> , 2021)	75.61±0.93	56.52±0.52	60.62±0.43	58.41±2.39	41.84±0.74	36.00±0.83	31.33±0.41	32.11±0.39
ODDL (Ye <i>et al.</i> , 2022)	75.03±1.00	61.61±3.55	65.46±0.46	66.19±2.08	40.26±0.50	41.88±4.52	38.82±0.49	41.35±1.08
MEMO (Zhou <i>et al.</i> , 2023)	73.21±0.49	62.47±3.38	59.26±0.90	62.01±1.17	40.60±1.11	39.87±0.46	23.41±1.63	32.74±2.11
X-DER (Boschini <i>et al.</i> , 2023)	77.59±0.62	65.40±4.79	58.28±1.17	59.76±4.28	52.80±1.61	43.73±0.86	41.94±0.57	46.00±1.04
EARL (Ours)	<b>78.31±0.72</b>	<b>65.71±2.26</b>	<b>69.52±0.19</b>	<b>70.41±1.97</b>	<b>57.12±1.22</b>	<b>44.40±0.68</b>	<b>47.89±0.61</b>	<b>46.09±0.26</b>

Method	TinyImageNet				ImageNet-200			
	Disjoint		Gaussian-Scheduled		Disjoint		Gaussian-Scheduled	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$						
EWC (Kirkpatrick <i>et al.</i> , 2017)	37.95±0.93	27.50±0.80	25.29±0.81	26.06±0.52	41.84±0.64	31.57±0.80	30.71±0.27	33.33±0.98
ER (Rolnick <i>et al.</i> , 2019)	37.43±1.05	27.47±0.63	26.37±0.89	25.79±0.44	41.51±0.76	30.87±0.72	32.39±0.36	33.09±0.37
ER-MIR (Aljundi <i>et al.</i> , 2019)	37.81±1.06	26.72±0.86	26.22±0.69	25.11±1.04	38.28±0.38	33.12±0.73	32.17±0.44	33.85±0.93
REMIN (Hayes <i>et al.</i> , 2020)	28.37±0.13	27.68±0.45	10.19±0.60	14.90±1.49	39.25±0.93	31.98±0.84	30.23±0.62	33.98±0.09
DER++ (Buzzega <i>et al.</i> , 2020)	39.38±0.60	29.36±0.81	27.23±1.94	31.53±0.80	43.50±0.31	34.56±0.50	35.22±0.26	38.38±0.97
SCR (Mai <i>et al.</i> , 2021)	34.65±1.08	22.18±0.32	25.86±0.94	22.54±0.59	41.90±0.40	28.92±0.40	33.24±0.32	30.98±0.28
MEMO (Zhou <i>et al.</i> , 2023)	27.36±0.61	27.57±0.52	10.82±1.23	18.03±1.36	41.55±0.23	34.19±1.47	32.54±0.39	36.11±1.06
X-DER (Boschini <i>et al.</i> , 2023)	35.15±2.12	26.67±0.52	29.71±0.86	28.10±0.50	43.41±0.47	34.14±0.98	36.31±0.17	38.61±0.55
EARL (Ours)	<b>41.77±1.26</b>	<b>29.65±0.20</b>	<b>35.08±0.70</b>	<b>32.49±1.21</b>	<b>44.88±0.29</b>	<b>34.79±0.55</b>	<b>39.14±0.47</b>	<b>38.83±0.35</b>

Table 1. Comparison of online CL methods on Disjoint and Gaussian Scheduled Setup for CIFAR-10, CIFAR-100, TinyImageNet and ImageNet-200.

remaining discrepancy between the classifiers and the features at inference.

To select which of the stored residuals to use during inference, we not only store the residual, but also  $f(x)$  to choose the stored residual closest to  $f(x_{infer})$ . Therefore, we retain ‘feature-residual’ pairs in a ‘feature-residual memory’  $M = \{(\hat{h}_i, \mathbf{r}_i)\}_{i=1}^N$ , where  $\hat{h}_i = \hat{f}(x_i)$ ,  $\mathbf{r}_i = \mathbf{w}_{y_i} - \hat{f}(x_i)$ , where  $N$  is the size of feature-residual memory, and  $\mathbf{w}_{y_i}$  is the classifier vector for class  $y_i$ .

During inference, we select the  $k$  nearest neighbor  $\hat{h}$ ’s, *i.e.*,  $\{\hat{h}_{n_1}, \hat{h}_{n_2}, \dots, \hat{h}_{n_k}\}$  from  $\{\hat{h}_i\}_{i=1}^N$ , since using only the nearest residual for correction may lead to incorrect inference predictions if a wrong residual is selected from a different class. Finally, following ResMem [67], we calculate the residual correction term  $\mathbf{r}_{eval}$  by a weighted average of the corresponding residuals  $\{\mathbf{r}_{n_1}, \mathbf{r}_{n_2}, \dots, \mathbf{r}_{n_k}\}$ , with weights  $\{s_1, s_2, \dots, s_k\}$  that are inversely proportional to the distance from  $\hat{f}(x_{eval})$ , as:

$$\mathbf{r}_{eval} = \sum_{i=1}^k s_i \mathbf{r}_i, \quad s_i = \frac{e^{-(\hat{f}(x_{eval}) - \hat{h}_i)/\tau}}{\sum_{j=1}^k e^{-(\hat{f}(x_{eval}) - \hat{h}_j)/\tau}}, \quad (5)$$

where  $\tau$  is a temperature hyperparameter. We add the residual-correcting term  $\mathbf{r}_{eval}$  to the model output  $\hat{f}(x_{eval})$

to obtain the corrected output  $\hat{f}(x_{eval})_{corrected}$  as:

$$\hat{f}(x_{eval})_{corrected} = \hat{f}(x_{eval}) + \mathbf{r}_{eval}. \quad (6)$$

## 5. Experiments

### 5.1. Experimental setup

We perform experiments on four datasets: CIFAR-10, CIFAR-100, TinyImageNet, ImageNet-200, and ImageNet-1K. For all datasets, our experiments are conducted on both a disjoint setup [44] and a Gaussian scheduled setup [55, 61]. We report the average and standard deviation results in three different seeds, except ImageNet-1k due to computational resource constraints [3, 34].

To evaluate anytime inference performance, we use the area under the curve accuracy ( $A_{auc}$ ) [7, 34], which measures the area under the accuracy curve. We also use last accuracy ( $A_{last}$ ) which measures accuracy at the end of training. For detailed information about the experimental setup, refer to the Supplementary.

**Baselines.** We compare EARL with various CL methods in different categories: regularization (EWC [33]), network expansion (MEMO [74]), replay (ER [52], ER-MIR [1], REMIND [24], SCR [40], ODDL [4]), and distillation

Method	ImageNet-1K			
	Disjoint		Gaussian-Scheduled	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$
EWC (Kirkpatrick <i>et al.</i> , 2017)	30.17	18.78	17.09	20.15
ER (Rolnick <i>et al.</i> , 2019)	30.18	18.96	17.17	18.39
ER-MIR (Aljundi <i>et al.</i> , 2019)	31.68	19.87	19.37	16.30
REMIND (Hayes <i>et al.</i> , 2020)	32.47	<b>24.59</b>	17.42	19.79
DER++ (Buzzega <i>et al.</i> , 2020)	31.08	18.98	23.73	23.14
SCR (Mai <i>et al.</i> , 2021)	26.72	13.60	19.82	15.84
MEMO (Zhou <i>et al.</i> , 2023)	30.45	19.08	14.06	20.12
X-DER (Boschini <i>et al.</i> , 2023)	31.67	18.18	25.18	12.51
EARL (Ours)	<b>34.33</b>	23.19	<b>30.53</b>	<b>29.48</b>

Table 2. Comparison of online CL methods on Disjoint and Gaussian Scheduled Setup for ImageNet-1K.

(DER++ [5], X-DER [4]). For more details on the implementation of these methods and a comparison with NC-FSCIL [66], which attempts to induce neural collapse in the context of few-shot class incremental learning, please refer to the supplementary material.

**Implementation details.** We use three components to architect our model: a backbone network  $g(\cdot)$ , a projection MLP  $p(\cdot)$ , and a fixed ETF classifier  $\mathbf{W}_{ETF}$  (*i.e.*, our model  $f$  can be defined as  $f(x) = p \circ g(x)$ ). For the projection layer, we attach an MLP projection layer  $p_{\theta_p}$  to the output of the backbone network  $g_{\theta_g}$ , where  $\theta_g$  and  $\theta_p$  denote the parameters of the backbone network and the projection layer, respectively (*i.e.*, the model  $f$  can be defined as  $f(x) = p \circ g(x)$ ), following [10, 46, 66]. For all methods, we use ResNet-18 [25] as the backbone network.

Following [34, 69], we employ memory-only training, where a random batch is selected from memory at each iteration. Furthermore, for episodic memory sampling, EARL uses the Greedy Balanced Sampling strategy [49]. We describe the details about hyperparameters and the pseudocode of EARL in Supplementary for the sake of space.

## 5.2. Results

We first compare the accuracy of online continual learning methods, including EARL, and summarize the results in Table 1. As shown in the table, EARL outperforms other baselines on all benchmarks, both in disjoint and Gaussian-scheduled setups, except  $A_{last}$  in ImageNet-1K disjoint setup. In particular, high  $A_{AUC}$  suggests that EARL outperforms other methods for all the time that the data stream is provided to the model, which implies that it can be used for inference at anytime.

Furthermore, EARL, SCR, ER, MIR, DER, and ODDL do not use task boundary information during training, *i.e.*, *task-free*, in contrast to EWC, X-DER, MEMO, and REMIND which use task-boundary information, *i.e.*, *task-aware*. Nevertheless, EARL outperforms these methods even in the disjoint setup where utilizing task boundary information is advantageous due to abrupt distribution shifts

at the task boundaries, except MEMO in ImageNet-1K disjoint. Since MEMO not only uses the task boundary information but also expands the network per task, this extra advantage of using a larger model is emphasized in large-scale datasets such as ImageNet which require a large learning capacity of the model. Please refer to the supplementary material for details of the computational budget comparison among the baseline methods.

**Memory budget.** We consider not only the size of episodic memory, but also the size of the model, logits, and other components, following MEMO [74]. In the case of the memory budget for CIFAR-10 and the architecture of ResNet-18 (size of 42.6MB), EARL incurs an additional memory cost of 0.2MB to store feature-residual pairs. EWC requires an additional cost of 85.2MB due to the need to store Fisher Information per parameter and the previous model. MEMO, which expands the network per task, incurs an additional cost of 32MB per task. DER also requires an additional cost of 0.1MB to store logits during distillation.

## 5.3. Ablation study

We conduct an ablation study on the two components of EARL, preparatory data training and residual correction, and summarize the results in Table 3. Although both components contribute to performance improvement, preparatory data training shows a larger gain in performance. Training with preparatory data improves the baseline by 2.7%  $\sim$  3.3% in  $A_{auc}$  and 2.1%  $\sim$  4.4% in  $A_{last}$ . When combined with residual correction, we observe further gains across all metrics in both disjoint and Gaussian-scheduled setups.

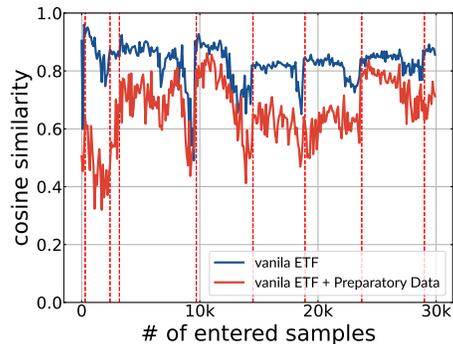


Figure 5. Average similarity between features of the most recently added class’s samples and the closest classifier vectors of the old classes (CIFAR-10, Gaussian Scheduled). Baseline is a vanilla ETF model trained only using episodic memory.

For further analysis, Fig. 6 visualizes the cosine similarity between the output features of 50 randomly selected test set samples of novel class ‘4’ and the classifier vectors  $\mathbf{w}_i$ . In baseline (a), the features of the class 4 are strongly biased towards the classifier vectors of the old classes (*i.e.*,  $\mathbf{w}_0, \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3$ ), rather than the correct classifier,  $\mathbf{w}_4$ . (c)

METHOD	CIFAR-10				CIFAR-100			
	Disjoint		Gaussian-Scheduled		Disjoint		Gaussian-Scheduled	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$						
EARL (Ours)	<b>78.61±0.72</b>	<b>66.01±2.26</b>	<b>69.62±0.19</b>	<b>70.91±1.97</b>	<b>57.42±1.24</b>	<b>44.60±0.65</b>	<b>48.19±0.61</b>	<b>46.10±0.26</b>
(-) RC	77.78±0.52	65.52±1.63	68.37±0.10	70.17±1.08	56.28±1.40	44.29±1.04	47.05±0.71	46.07±0.52
(-) RC & PDT	74.77±0.77	61.10±3.12	65.60±0.25	67.39±0.78	52.91±1.05	41.08±0.60	44.34±0.55	44.45±0.48

Table 3. Ablation Study. RC and PDT refer to preparatory data training (Sec. 4.2) and the residual correction (Sec. 4.3), respectively.

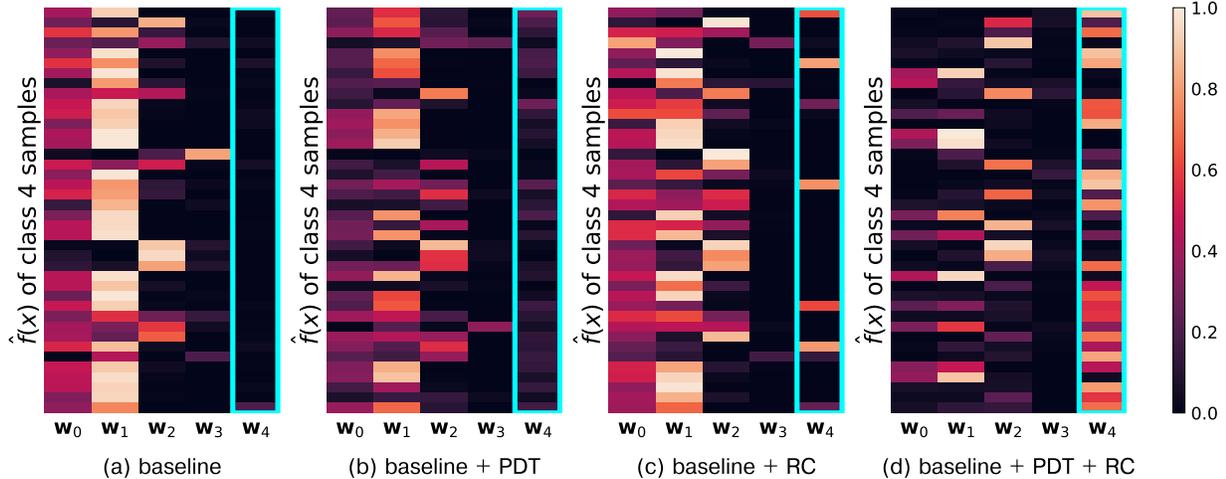


Figure 6. Cosine similarity between features  $\hat{f}(x)$  for class 4 and the ETF classifier vectors  $w_i$  at the 50<sup>th</sup> iteration after the introduction of class 4 in the Gaussian Scheduled CIFAR-10 setup. As we can see in the cyan highlighted box, EARL promotes the convergence of  $\hat{f}(x)$  for class 4 toward the ground truth classifier vector  $w_4$ .

When residual correction is used, some samples show high similarity with  $w_4$  compared to the baseline.

However, since incorrect residuals can be added due to the bias problem, more samples have high similarity with  $w_2$  than the baseline (*i.e.*, wrong residuals are added). (b) When using preparatory data training, the bias toward  $w_0$  and  $w_1$  classes significantly decreases compared to the baseline. Fig. 5 also shows the effect of preparatory data, which reduces the similarity between the novel class features and existing classes. (d) Using both residual correction and preparatory data training shows a remarkable alignment with the ground truth classifier  $w_4$ . A more detailed analysis of the ablation results is provided in the Supple.

## 6. Conclusion

To better learn online data in a continuous data stream without multiple epoch training, we propose to induce neural collapse, which aligns last layer activations to the corresponding classifier vectors in the representation space. Unlike in offline CL, it is challenging to induce neural collapse in online CL due to insufficient training epochs and continuously streamed new data. We first observe that the bias of the new class towards existing classes slows the convergence of features toward neural collapse.

To mitigate the bias, we propose synthesizing preparatory data for unseen classes by transforming the samples

of existing classes. Using the preparatory data for training, we accelerate neural collapse in an online CL scenario. Additionally, we propose residual correction to resolve the remaining discrepancy toward neural collapse at inference, which arises due to the continuous stream of new data. In our empirical evaluations, the proposed methods outperform state-of-the-art online CL methods in various datasets and setups, especially with high performance on anytime inference.

**Limitations and Future Work.** Since our work uses ETF structure, it has an inherent limitation that the number of possible classifier vectors in the ETF classifier is limited by the dimension of the embedding space. Considering life-long learning, where the number of new classes goes to infinity, it is interesting to explore the idea of dynamically expanding the ETF structure so that the model can continually learn the ever-increasing number of concepts in the real world.

**Acknowledgement.** This work was partly supported by the NRF grant (No.2022R1A2C4002300, 20%) and IITP grants (No.2022-0-00077 (10%), No.2022-0-00113 (25%), No.2022-0-00959 (10%), No.2022-0-00871 (10%), No.2020-0-01361 (5%, Yonsei AI), No.2021-0-01343 (5%, SNU AI), No.2021-0-02068 (5%, AI Innov. Hub), No.2022-0-00951 (5%)) funded by the Korea government (MSIT) and CARAI grant funded by DAPA and ADD (UD230017TD) 5%. This work was done while J. Choi was with Yonsei University and partially while H. Lee, M. Seo and H. Kim was with LG AI Research.

## References

- [1] Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning. In *CVPR*, 2019. 1, 2, 6
- [2] Soumya Banerjee, Vinay K Verma, Avideep Mukherjee, Deepak Gupta, Vinay P Nambodiri, and Piyush Rai. Verse: Virtual-gradient aware streaming lifelong learning with anytime inference. *arXiv preprint arXiv:2309.08227*, 2023. 12
- [3] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, 2021. 2, 6
- [4] Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class-incremental continual learning into the extended der-verse. *IEEE TPAMI*, 45(5): 5497–5512, 2022. 2, 6, 7
- [5] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *NeurIPS*, 2020. 2, 7
- [6] Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky. New insights on reducing abrupt representation change in online continual learning. In *ICLR*, 2021. 2, 4
- [7] Lucas Caccia, Jing Xu, Myle Ott, Marcaurelio Ranzato, and Ludovic Denoyer. On anytime learning at macroscale. In *CoLLAs*. PMLR, 2022. 2, 6
- [8] Zhipeng Cai, Ozan Sener, and Vladlen Koltun. Online continual learning with natural distribution shifts: An empirical study with visual data. In *ICCV*, 2021. 1
- [9] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *ECCV*, 2018. 1
- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*. PMLR, 2020. 7
- [11] Brian Cheung, Alexander Terekhov, Yubei Chen, Pulkit Agrawal, and Bruno Olshausen. Superposition of many models into one. In *NeurIPS*, 2019. 2
- [12] Florent Chiaroni, Ghazaleh Khodabandelou, Mohamed-Cherif Rahal, Nicolas Hueber, and Frederic Dufaux. Counter-examples generation from a positive unlabeled image dataset. In *CVPR*. Elsevier, 2020. 5
- [13] Aristotelis Chrysakis and Marie-Francine Moens. Online bias correction for task-free continual learning. In *ICLR*, 2023. 4
- [14] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 702–703, 2020. 14
- [15] Mehmet Dedeoglu, Sen Lin, Zhaofeng Zhang, and Junshan Zhang. Continual learning of generative models with limited data: From wasserstein-1 barycenter to adaptive coalescence. *IEEE TNNLS*, 2023. 5
- [16] Thang Doan, Seyed Iman Mirzadeh, and Mehrdad Farajtabar. Continual learning beyond a single model. *arXiv preprint arXiv:2202.09826*, 2022. 12
- [17] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE international conference on computer vision*, pages 1422–1430, 2015. 12
- [18] Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training. *PNAS*, 118(43): e2103091118, 2021. 1, 3
- [19] Zeyu Feng, Chang Xu, and Dacheng Tao. Self-supervised representation learning by rotation feature decoupling. In *CVPR*, 2019. 5
- [20] Matthew Fickus, John Jasper, Emily J King, and Dustin G Mixon. Equiangular tight frames that contain regular simplices. *Linear Algebra and its applications*, 555:98–138, 2018. 4
- [21] Yasir Ghunaim, Adel Bibi, Kumail Alhamoud, Motasem Alfarrar, Hasan Abed Al Kader Hammoud, Ameya Prabhu, Philip HS Torr, and Bernard Ghanem. Real-time evaluation in online continual learning: A new hope. In *CVPR*, 2023. 2
- [22] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018. 5
- [23] Gabriele Graffieti, Davide Maltoni, Lorenzo Pellegrini, and Vincenzo Lomonaco. Generative negative replay for continual learning. *Neural Networks*, 162:369–383, 2023. 5
- [24] Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. Remind your neural network to prevent catastrophic forgetting. In *ECCV*. Springer, 2020. 2, 6
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 7
- [26] Dong-Hwan Jang, Sanghyeok Chu, Joonhyuk Kim, and Bohyung Han. Pooling revisited: Your receptive field is suboptimal. In *CVPR*, 2022. 5
- [27] Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J Su. An unconstrained layer-peeled perspective on neural collapse. In *ICLR*, 2021. 1, 3
- [28] Haeyong Kang, Thang Vu, and Chang D Yoo. Learning imbalanced datasets with maximum margin loss. In *IEEE ICIP*. IEEE, 2021. 1
- [29] Chris Dongjoo Kim, Jinseo Jeong, Sangwoo Moon, and Gunhee Kim. Continual learning on noisy data streams via self-purified replay. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 537–547, 2021. 12
- [30] Gyuhak Kim, Sepideh Esmailpour, Changnan Xiao, and Bing Liu. Continual learning based on ood detection and task masking. In *CVPR*, 2022. 2, 5
- [31] Jun Kimata, Tomoya Nitta, and Toru Tamaki. Objectmix: Data augmentation by copy-pasting objects in videos for action recognition. In *Proceedings of the 4th ACM International Conference on Multimedia in Asia*, pages 1–7, 2022. 12
- [32] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, San Diego, CA, USA, 2015. 14

- [33] 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. *PNAS*, 114(13):3521–3526, 2017. 2, 6
- [34] Hyunseo Koh, Dahyun Kim, Jung-Woo Ha, and Jonghyun Choi. Online continual learning on class incremental blurry task configuration with anytime inference. In *ICLR*, 2022. 1, 2, 6, 7
- [35] Hyunseo Koh, Minhyuk Seo, Jihwan Bang, Hwanjun Song, Deokki Hong, Seulki Park, Jung-Woo Ha, and Jonghyun Choi. Online boundary-free continual learning by scheduled data prior. In *The Eleventh International Conference on Learning Representations*, 2023. 14
- [36] Sanghyeok Lee, Minkyu Jeon, Injae Kim, Yongyang Xiong, and Hyunwoo J Kim. Sagemix: Saliency-guided mixup for point clouds. *Advances in Neural Information Processing Systems*, 35:23580–23592, 2022. 12
- [37] Timothée Lesort, Andrei Stoian, and David Filliat. Regularization shortcomings for continual learning. *arXiv preprint arXiv:1912.03049*, 2019. 2
- [38] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE TPAMI*, 40(12):2935–2947, 2017. 2
- [39] Jianfeng Lu and Stefan Steinerberger. Neural collapse with cross-entropy loss. *arXiv preprint arXiv:2012.08465*, 2020. 1, 3
- [40] 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 *CVPRW*, 2021. 6
- [41] Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *arXiv preprint arXiv:2011.11619*, 2020. 1, 3
- [42] Sergey Pankov. Learning image transformations without training examples. In *ISVC*. Springer, 2011. 5
- [43] Vardan Pappayan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *PNAS*, 117(40):24652–24663, 2020. 1, 3
- [44] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural networks*, 113:54–71, 2019. 6
- [45] Lorenzo Pellegrini, Gabriele Graffieti, Vincenzo Lomonaco, and Davide Maltoni. Latent replay for real-time continual learning. In *IROS*. IEEE, 2020. 2
- [46] Can Peng, Kun Zhao, Tianren Wang, Meng Li, and Brian C Lovell. Few-shot class-incremental learning from an open-set perspective. In *ECCV*, pages 382–397. Springer, 2022. 7
- [47] Federico Pernici, Matteo Bruni, Claudio Bacchi, Francesco Turchini, and Alberto Del Bimbo. Class-incremental learning with pre-allocated fixed classifiers. In *ICPR*. IEEE, 2021. 3
- [48] Jary Pomponi, Simone Scardapane, and Aurelio Uncini. Continual learning with invertible generative models. *Neural Networks*, 164:606–616, 2023. 2
- [49] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*. Springer, 2020. 7
- [50] Akshay Rangamani and Andrzej Banburski-Fahey. Neural collapse in deep homogeneous classifiers and the role of weight decay. In *ICASSP*. IEEE, 2022. 1, 3
- [51] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017. 1
- [52] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In *NeurIPS*, 2019. 6
- [53] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016. 2
- [54] Minhyuk Seo, Hyunseo Koh, and Jonghyun Choi. Budgeted online continual learning by adaptive layer freezing and frequency-based sampling. 2023. 12
- [55] Murray Shanahan, Christos Kaplanis, and Jovana Mitrović. Encoders and ensembles for task-free continual learning. *arXiv preprint arXiv:2105.13327*, 2021. 6
- [56] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *NeurIPS*, 2017. 2
- [57] Abhishek Sinha, Kumar Ayush, Jiaming Song, Burak Uzkent, Hongxia Jin, and Stefano Ermon. Negative data augmentation. *arXiv preprint arXiv:2102.05113*, 2021. 2, 5
- [58] Tom Tirer and Joan Bruna. Extended unconstrained features model for exploring deep neural collapse. In *ICML*. PMLR, 2022. 1, 3
- [59] Alice Barbora Tumpach and Peter Kán. Temporal alignment of human motion data: A geometric point of view. In *GSI*. Springer, 2023. 5
- [60] Chenyang Wang, Junjun Jiang, Xiong Zhou, and Xianming Liu. Resmooth: Detecting and utilizing ood samples when training with data augmentation. *IEEE TNNLS*, 2022. 2, 5
- [61] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaohu Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *CVPR*, 2022. 6
- [62] Stephan Wojtowysch. On the emergence of simplex symmetry in the final and penultimate layers of neural network classifiers. In *MSML*, pages 1–21. PMLR, 2021. 1, 3
- [63] Ou Wu. Rethinking class imbalance in machine learning. *arXiv preprint arXiv:2305.03900*, 2023. 1
- [64] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *CVPR*, pages 374–382, 2019. 1, 2, 4
- [65] Yibo Yang, Shixiang Chen, Xiangtai Li, Liang Xie, Zhouchen Lin, and Dacheng Tao. Inducing neural collapse in imbalanced learning: Do we really need a learnable classifier at the end of deep neural network? In *NeurIPS*, 2022. 1, 3, 4
- [66] Yibo Yang, Haobo Yuan, Xiangtai Li, Zhouchen Lin, Philip Torr, and Dacheng Tao. Neural collapse inspired feature-

- classifier alignment for few-shot class incremental learning. In *ICLR*, 2023. 1, 3, 4, 7
- [67] Zitong Yang, Michal Lukasik, Vaishnavh Nagarajan, Zonglin Li, Ankit Singh Rawat, Manzil Zaheer, Aditya Krishna Menon, and Sanjiv Kumar. Resmem: Learn what you can and memorize the rest. *arXiv preprint arXiv:2302.01576*, 2023. 5, 6
- [68] Zhi Yang, Jiwei Qin, Chuan Lin, Yanping Chen, Ruizhang Huang, and Yongbin Qin. Ganrec: A negative sampling model with generative adversarial network for recommendation. *Expert Systems with Applications*, 214:119155, 2023. 5
- [69] Fei Ye and Adrian G Bors. Task-free continual learning via online discrepancy distance learning. In *NeurIPS*, 2022. 7
- [70] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju Hwang. Online coreset selection for rehearsal-based continual learning. In *ICLR*, 2021. 2
- [71] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 12
- [72] Bowen Zhao, Chen Chen, Qi Ju, and ShuTao Xia. Energy aligning for biased models. *arXiv preprint arXiv:2106.03343*, 2021. 1
- [73] Zhisheng Zhong, Jiequan Cui, Yibo Yang, Xiaoyang Wu, Xiaojuan Qi, Xiangyu Zhang, and Jiaya Jia. Understanding imbalanced semantic segmentation through neural collapse. In *CVPR*, 2023. 1, 3
- [74] Da-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards memory-efficient class-incremental learning. In *ICLR*, 2022. 2, 6, 7
- [75] Jinxin Zhou, Xiao Li, Tianyu Ding, Chong You, Qing Qu, and Zhihui Zhu. On the optimization landscape of neural collapse under mse loss: Global optimality with unconstrained features. In *ICML*. PMLR, 2022. 1, 3
- [76] Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. In *NeurIPS*, 2021. 1, 3

Note: Blue characters denote the reference of the main paper.

## 7. Anytime Inference in Online CL (L164)

In online CL, new data continuously arrive in a stream rather than in a large chunk (e.g., task unit). Several previous works [2][29] train the model only after a large chunk of new data accumulates, which leads to poor inference performance on new data during data accumulation, since the model is not trained during the accumulation period [6, 27].

However, inference queries can occur at any time including the data accumulation, i.e., *anytime inference*, whose importance has been emphasized by recent research [2, 16] [17, 27, 36]. Consequently, we focus not only on  $A_{last}$ , which is measured after the learning has finished for all data, but also on  $A_{AUC}$ , which measures the average accuracy during training, i.e., ‘anytime inference’ performance [27].

## 8. Analysis about Negative Transformation (L332)

Negative Transformation	Gaussian-Scheduled		Disjoint	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$
Patch Permutation	66.37±0.36	69.04±1.39	75.34±0.88	63.73±2.05
Negative Cutmix	65.26±0.22	64.48±1.07	72.86±0.13	56.56±3.25
Gaussian Noise	64.83±0.23	66.97±1.40	74.18±0.22	63.32±1.52
Negative Rotation	<b>69.52±0.13</b>	<b>70.28±1.77</b>	<b>77.86±0.71</b>	<b>69.50±1.94</b>

Table 4. Comparison of various negative transformations to obtain preparatory data in CIFAR-10

Negative transformation is a transformation that modifies semantic information of original images and is widely used in self-supervised learning [4, 5] and out-of-distribution (OOD) detection [25, 47, 50]. We consider various negative transformations such as negative rotation (i.e., rotation by 90, 180, and 270 degrees) [15, 18], patch permutation [17], negative CutMix, and Gaussian noise. Negative Cutmix is a kind of Cutmix [71] that finely divides an image into small patches and mixes it with another finely divided image. Gaussian noise refers to filling all the image pixels with a Gaussian random noise. Even though Gaussian noise is not a kind of transformation, it is a straightforward approach to generate synthetic input that can be distinguished from the existing images; thus, we also consider this.

Among them, we choose to use rotations 90, 180, and 270 degrees as the negative transform set  $T$ , because the rotation transformation is simple to implement, widely applicable from low- to high-resolution images, and also outperforms other transformations in our empirical evaluations, as we can see in Tab. 4. Unlike rotation transformation, which preserves image continuity while modifying semantic information, patch permutation or Negative Cutmix could cause discontinuities at the boundary of patches [31, 36], leading to loss of image features.

Furthermore, to compare the effect of various negative transformations, we compare the cosine similarities between the features and the corresponding ground truth classifier as shown in Fig. 7. As we can see in the cyan highlighting box, the rotation transformation promotes the convergence of  $\hat{f}(x)$  for class 0 towards the ground truth classifier vector  $w_0$ .

## 9. Effect of Negative Rotation Transformation (L304)

We use negative rotation, which rotates images by 90, 180, and 270 degrees [15, 18]. Since the vertical information of the image is more important than the horizontal information [22, 49] as mentioned in Sec.4.2, rotation by a large angle changes the semantic information [50]. To empirically verify that the negative rotation transformation alters the semantic information of the image, we generate four datasets by rotating the original data by 0, 90,

180, and 270 degrees and train the model using each dataset. Subsequently, we performed inference on all four datasets to check whether they have different semantics.

Specifically, we used CIFAR-10-R90, CIFAR-10-R180, and CIFAR-10-R270, created by applying rotation transformations of 90, 180, and 270 degrees to CIFAR-10 dataset, respectively. We summarize the results in Tab. 5.

Train Dataset	Evaluation Dataset			
	CIFAR-10	CIFAR-10-R90	CIFAR-10-R180	CIFAR-10-R270
CIFAR-10	<b>92.9</b>	31.7	34.7	31.8
CIFAR-10-R90	32.1	<b>92.6</b>	31.9	35.3
CIFAR-10-R180	33.7	30.9	<b>92.5</b>	32.3
CIFAR-10-R270	32.4	34.3	31.4	<b>92.4</b>

Table 5. Comparison of the inference accuracy of the rotated datasets after training with each rotated dataset.

## 10. Prevent Forgetting with ETF Classifier

We compare the performance and forgetting between a fixed ETF classifier and a learnable classifier. We summarize the results in Tab. 6. The fixed ETF classifier not only effectively prevents forgetting, but also has high accuracy.

Classifier	CIFAR-10		CIFAR-100	
	$A_{AUC} \uparrow$	Forgetting $\downarrow$	$A_{AUC} \uparrow$	Forgetting $\downarrow$
Learnable	66.02±0.18	6.54±1.02	44.37±0.84	9.04±1.45
Fixed ETF	<b>69.61±0.35</b>	<b>3.75±2.21</b>	<b>47.78±0.69</b>	<b>7.41±1.02</b>

Table 6. Comparison of Forgetting and  $A_{AUC}$  between a learnable classifier and a classifier with fixed ETF structure on CIFAR-10/100.

## 11. Comparison of Baselines with Computational Constraint

Following [54][17], we measure the number of FLOPs required for each method to complete one iteration. We then calculate the relative FLOPs with respect to the ER [43], which is the simplest method, and adjust the total training iteration to align the total FLOPs used for the entire training process. With this computational constraint, we compare online CL methods on disjoint and Gaussian scheduled setup for CIFAR-10, CIFAR-100, TinyImageNet and ImageNet-200. Furthermore, we measure the Average Online Accuracy (AOA)[7, 17], which uses the newly encountered streaming data for evaluation before incorporating it into the training process. We plot the online accuracy on TinyImageNet in Fig. 8.

## 12. Comparison of Baselines on CLEAR-10/100

To evaluate the setups in which the domain of classes changes over time, we compare CL methods with the CLEAR benchmark and summarize the result in Tab. 8. In our experiments, we address a more realistic scenario, called the domain-class-IL setup, where both novel classes and domain shifts take place. This contrasts with domain-IL setups, where all classes are initially provided and only the data distribution changes over time.

EARL significantly outperforms the baselines in CLEAR-10/100 benchmarks and achieves high performance. We perform experiments under the computational constraints mentioned in Sec. 11.

## 13. Comparison of Computational Budget (L)

EARL requires an additional computation for residual correction, which calculates  $k$  nearest features, and FLOPs of additional cost can be formu-

†: Work done while interning at LG AI Research.

\*: Indicates corresponding authors.

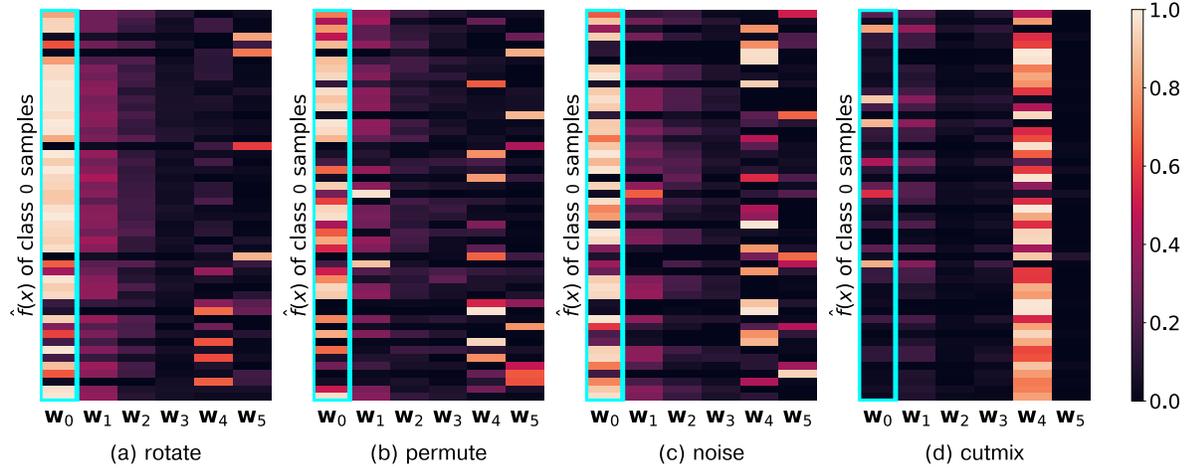


Figure 7. Cosine similarity between the features  $\hat{f}(x)$  for class 0 and the ETF classifier vectors  $w_i$  at the  $10000^{th}$  iteration after the introduction of class 0 in the Gaussian Scheduled CIFAR-10 setup.

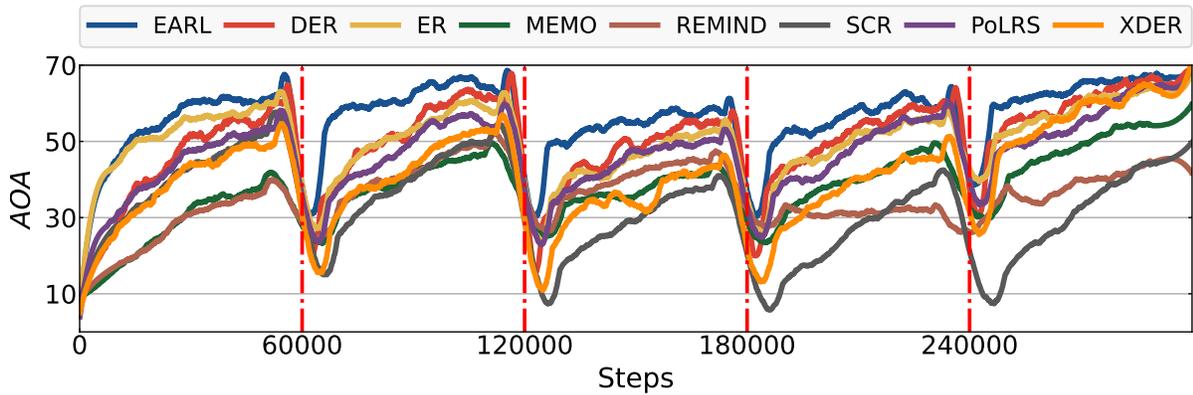


Figure 8. Online accuracy in TinyImageNet Disjoint. Red lines denote task boundary

Method	CIFAR-10						CIFAR-100					
	Disjoint			Gaussian-Scheduled			Disjoint			Gaussian-Scheduled		
	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$
EWC (Kirkpatrick et al., 2017)	77.30±0.91	62.78±1.25	88.70±0.89	60.91±0.30	63.56±3.04	83.22±0.16	52.94±1.00	43.27±0.41	62.04±0.52	40.89±0.40	42.80±0.89	55.29±0.18
ER (Rolnick et al., 2019)	76.89±0.96	63.92±2.36	88.47±0.78	60.71±0.10	66.71±2.49	82.47±0.17	53.13±1.28	41.97±0.21	62.06±0.73	41.30±0.22	44.10±0.14	55.28±0.33
ER-MIR (Aljundi et al., 2019)	75.08±0.10	63.13±3.32	87.27±1.52	57.25±0.71	59.34±2.12	79.29±0.69	50.17±0.91	42.17±0.50	57.74±0.31	35.52±0.57	41.29±0.71	46.44±0.55
REMIND (Hayes et al., 2020)	67.73±0.42	49.33±1.30	81.30±1.17	55.94±1.06	51.73±4.11	77.41±1.27	40.76±0.46	37.64±1.13	48.18±0.98	22.55±1.58	30.27±1.94	33.02±1.19
DER++ (Buzzegga et al., 2020)	75.71±1.46	58.49±3.16	<b>88.93±1.00</b>	59.14±1.15	66.16±3.83	82.14±0.76	41.39±0.80	42.03±0.81	57.00±0.30	28.63±1.93	37.10±2.23	47.39±1.38
SCR (Mai et al., 2021)	75.89±0.47	57.90±2.45	87.91±1.82	60.38±0.25	63.02±3.71	78.73±0.48	36.74±1.72	30.54±0.99	34.56±1.26	26.53±1.00	27.59±0.96	29.88±0.93
ODDL (Ye et al., 2022)	75.91±0.87	61.89±4.47	88.61±1.50	61.02±0.60	65.25±1.01	83.58±0.69	53.16±1.18	42.89±0.64	61.19±0.78	42.46±0.30	44.26±1.26	56.33±0.36
MEMO (Zhou et al., 2023)	72.59±0.18	63.29±5.07	86.33±1.81	57.88±1.30	61.67±2.32	79.04±0.82	39.48±0.73	37.53±0.64	50.06±0.13	22.46±1.40	33.27±2.39	36.02±1.21
X-DER (Boschini et al., 2023)	75.99±1.22	65.69±4.05	83.28±2.05	57.73±0.76	66.59±2.45	76.63±0.32	47.22±1.21	43.72±0.56	50.73±0.78	36.81±0.56	46.20±0.53	47.59±0.33
EARL (Ours)	<b>78.54±0.48</b>	<b>66.16±0.84</b>	88.10±0.58	<b>69.77±0.05</b>	<b>71.46±1.84</b>	<b>86.00±0.17</b>	<b>57.65±1.30</b>	<b>45.15±0.05</b>	<b>67.09±0.69</b>	<b>48.05±0.67</b>	<b>47.27±0.69</b>	<b>63.69±0.76</b>

Method	TinyImageNet						ImageNet-200					
	Disjoint			Gaussian-Scheduled			Disjoint			Gaussian-Scheduled		
	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$
EWC (Kirkpatrick et al., 2017)	37.57±0.83	27.43±0.89	49.58±0.42	26.35±0.84	25.61±0.23	37.77±0.25	41.85±0.64	31.57±0.76	64.10±0.26	32.22±0.27	32.71±0.94	54.68±0.31
ER (Rolnick et al., 2019)	37.29±0.81	27.04±0.40	49.66±0.39	26.37±0.90	25.92±0.38	37.66±0.50	41.65±0.62	32.18±0.22	64.22±0.23	32.43±0.51	32.85±0.27	54.72±0.50
ER-MIR (Aljundi et al., 2019)	37.73±0.84	27.40±0.22	48.05±0.15	24.00±0.73	25.10±0.41	32.34±0.46	39.88±0.44	33.27±0.66	59.87±0.14	28.65±0.26	33.30±1.84	48.36±0.44
REMIND (Hayes et al., 2020)	29.05±0.83	27.30±0.59	36.43±0.83	10.22±0.93	16.59±0.70	16.85±1.17	38.32±0.73	31.60±0.46	39.80±1.85	30.46±0.37	33.36±0.82	38.13±1.17
DER++ (Buzzegga et al., 2020)	39.24±0.80	29.77±1.16	49.30±0.81	27.27±2.13	31.40±0.89	40.19±1.47	44.67±0.17	32.65±0.54	65.74±0.21	37.11±0.28	37.85±0.74	59.81±0.33
SCR (Mai et al., 2021)	34.07±1.16	24.24±0.43	33.28±0.93	25.62±0.90	24.28±0.24	29.36±0.41	41.94±0.31	28.53±0.37	61.97±0.40	33.27±0.46	31.41±0.31	54.12±0.48
MEMO (Zhou et al., 2023)	27.84±0.38	27.52±0.52	37.35±0.26	9.45±0.60	17.30±0.65	17.50±0.49	38.68±0.39	31.70±0.53	58.92±0.53	32.09±0.26	35.95±1.40	50.88±0.10
X-DER (Boschini et al., 2023)	35.68±0.38	27.15±1.66	41.36±1.41	23.51±2.43	24.80±4.06	32.17±3.71	44.03±0.37	33.58±1.17	56.40±0.08	34.92±0.24	38.29±1.12	53.75±0.23
EARL (Ours)	<b>42.19±1.18</b>	<b>29.50±0.78</b>	<b>56.41±0.08</b>	<b>34.79±0.67</b>	<b>31.68±0.21</b>	<b>49.77±0.13</b>	<b>45.01±0.36</b>	<b>34.25±0.69</b>	<b>67.08±0.34</b>	<b>39.04±0.22</b>	<b>38.95±0.50</b>	<b>60.80±0.13</b>

Table 7. Comparison of online CL methods on Disjoint and Gaussian Scheduled Setup for CIFAR-10, CIFAR-100, TinyImageNet and ImageNet-200 with computational constraint.

Method	CLEAR10			CLEAR100		
	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$AOA \uparrow$
EWC	70.88±1.15	69.46±2.40	76.96±1.16	45.74±0.40	47.61±0.54	46.15±0.28
ER	70.70±1.22	68.86±3.00	76.65±1.12	45.59±0.91	47.89±1.11	46.05±0.18
ER-MIR	68.21±0.94	65.58±2.33	74.53±1.10	43.21±1.03	46.60±1.21	42.24±0.47
REMIND	66.48±1.93	66.91±1.19	72.34±1.78	36.67±0.76	47.90±0.58	35.49±0.09
DER++	71.93±0.90	70.41±2.67	77.51±1.47	47.34±0.63	49.63±0.75	47.02±0.33
SCR++	73.32±0.85	70.81±1.69	77.90±1.11	44.67±0.77	44.70±0.72	46.15±0.28
MEMO	65.04±1.72	62.64±2.67	71.94±1.80	44.35±0.54	46.47±1.58	41.24±0.12
PoLRS	65.65±1.88	61.06±6.17	71.67±1.53	41.17±1.83	41.62±2.56	40.15±1.50
X-DER	69.77±0.85	68.67±3.04	74.76±1.23	44.76±1.21	49.01±1.31	43.60±0.65
EARL	<b>77.85±0.96</b>	<b>76.51±1.97</b>	<b>81.51±1.17</b>	<b>56.97±0.25</b>	<b>59.03±1.09</b>	<b>55.35±0.06</b>

Table 8. Comparison of online CL methods on CLEAR-10/100 with computational constraint.

lated as  $d \times N \times 3 + k \times N$ , where  $d$  is the dimension of stored features,  $N$  is the total number of feature-residual pairs, and  $k$  denotes  $\text{top}_k$  in  $k$ -NN. The first term is for calculating the distances between the stored features and the feature of the inference image and is multiplied by 3 since it involves subtraction, squaring, and addition operations. The second term is the cost of selecting  $\text{top}_k$  features among  $N$  features. Compared to the cost incurred by the naive inference process, which requires forward flops of the model, it involves a very minimal amount of additional cost. Taking the example of ImageNet-200 with the ResNet-18 architecture that has  $\frac{1}{3}$  GFLOPs in the model forward, only 0.5% of additional cost is consumed, since we use  $N = 2,000 (= 10 \times 200)$  and  $k = 15$ .

## 14. Details About Experiment Setup (L404)

This paper focuses on online class-incremental learning in two types of setups: disjoint [35] and Gaussian scheduled [35] [45, 51]. In the disjoint setup, each class is assigned to a specific task, *i.e.*, tasks do not share any classes. On the other hand, in the Gaussian scheduled setup, the arrival time of each class follows a Gaussian distribution  $\mathcal{N}(\mu_i, \sigma)$ . Since the class distribution is shifted every time step, the Gaussian scheduled setup is a boundary-free setup. We set  $\sigma$  to 0.1 and  $\mu_i$ , mean of the class  $i$ , to  $\frac{i}{N}$ , where  $N$  is the number of classes.

## 15. Effect of $k$

$k$ , which is a hyperparameter that determines the number of the nearest features from the inference image used to weight-sum to calculate the residual, has a negligible impact, except for  $k = 1$  as incorrect residual largely affects predictions. We illustrate the consistency of accuracy at various values of  $k$  in Fig. 9.

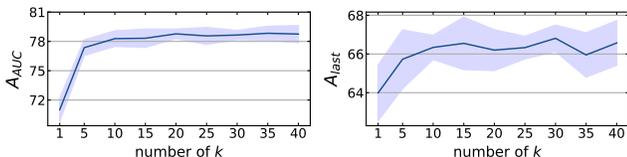


Figure 9. Performance variations of  $A_{AUC}$  and  $A_{last}$  with respect to the change in  $k$  in the CIFAR-10 disjoint setup.

## 16. Hyperparameters (L428)

For all methods, we use Adam optimizer [32] with a constant learning rate (LR) of 0.0003. For data augmentation, we use RandAugment [14]. For hyperparameters such as iteration, memory size, and number of tasks for each dataset, we follow prior works [2, 27, 40]. Specifically, we set the number of iterations as 1, 3, 3, 0.25, and 0.25, and memory sizes as 500,

2,000, 4,000, 4,000, and 10000 for the CIFAR-10, CIFAR-100, TinyImageNet, ImageNet-200, and ImageNet-1k datasets, respectively. To ensure a fair comparison among methods, we use the same batch size in all methods. Specifically, we use 16, 32, 32, 256, and 256 for CIFAR-10, CIFAR-100, TinyImageNet, ImageNet-200, and ImageNet-1k datasets, respectively. Note that the number of preparatory data is included in the batch size for EARL, *i.e.*, batch size = the number of retrieved samples from memory + the number of preparatory data.

To ensure that EARL does not depend on a specific dataset, we use the following hyperparameters for all datasets:  $\tau = 0.9$ ,  $k = 15$ ,  $d = 4096$  and  $\lambda = 1$ , where  $\tau$  refers to the temperature parameter during residual correction,  $k$  refers to the number of samples for  $\text{top}_k$  in  $k$ -NN during residual correction,  $d$  refers to the output dimension of the projection layer, and  $\lambda$  refers to the loss balancing parameter between real data training and preparatory data training. In addition, for  $N$ , the total number of stored feature-residual pairs in EARL, we used  $N = 10 * |C|$  for all dataset, where  $|C|$  is the number of classes seen so far.

## 17. Implementation of Baselines (L587)

Some of the baselines assumed multi-epoch training, *i.e.*, offline CL, so we modified them to be used in online CL for comparisons with our proposed method and other baselines.

**REMIND.** Remind is a feature-replay method and freezes front layers after base initialization, offline training phase using a subset of the data during the initial stage of training. We modified the base initialization of REMIND [20] to be suitable for use in online CL. In offline CL setup, REMIND freezes 7 layers after base initialization. However, as shown in Table 9, in online CL, it is not suitable to freeze many layers compared to offline CL, as the model is not sufficiently trained. We studied various numbers of frozen layers and revealed that freezing 6 layers (*i.e.*, 3 blocks) is the optimal number for freezing in ResNet-18 in CIFAR-10 and CIFAR-100. Compared to freezing 7 layers, which is the original proposed freezing criterion, REMIND performed comparatively well when freezing 6 layers in the online CL, despite the decrease in the number of stored features due to the larger sizes of stored features. Regarding the hyperparameters used in REMIND, such as the size of the codebook and the number of codebooks for product quantization, we performed additional hyperparameter search experiments with CIFAR-100, as shown in Table 10, and used the same hyperparameter in all remaining datasets.

For CIFAR-10, CIFAR-100, and TinyImageNet, which have a relatively small number of images, we increased the proportion of the initialization sample base from 10% to 60% of the total samples, since 10% images are not enough for the lower level layers to represent highly transferable features, leading to low performance. In ImageNet-200 and ImageNet-1k, we followed the original setting (*i.e.*, 10% base initialization samples and freezing 7 layers after base initialization), following [20].

Frozen Layers	CIFAR-10				CIFAR-100			
	Disjoint		Gaussian-Scheduled		Disjoint		Gaussian-Scheduled	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$						
4 Layers	66.63±1.24	46.24±0.12	54.55±0.71	45.56±1.54	38.73±0.27	30.19±0.66	23.76±0.86	26.94±0.37
5 Layers	<b>70.11±0.66</b>	51.01±0.79	56.47±0.96	52.00±1.94	41.87±0.05	<b>38.81±0.41</b>	24.79±1.73	32.12±2.81
6 Layers	69.55±0.91	47.28±3.92	<b>57.15±0.71</b>	<b>53.40±0.70</b>	<b>41.92±0.06</b>	37.64±1.09	<b>25.56±1.10</b>	<b>34.08±1.02</b>
7 Layers	68.59±0.10	<b>53.71±2.27</b>	55.37±0.60	52.53±1.98	40.52±0.17	37.42±0.81	23.94±1.30	33.28±1.43

Table 9. REMIND performance as a function of the number of frozen layers in ResNet-18 with CIFAR-10 and CIFAR-100. Rather than freezing all the layers except the last layer (*i.e.*, freezing 7 Layers), continuously updating the last two layers and fixing the rest (*i.e.*, freezing 6 Layers) shows the best performance due to the limitation of online CL.

# of Codebooks	CIFAR-100				Codebook Size	CIFAR-100			
	Disjoint		Gaussian-Scheduled			Disjoint		Gaussian-Scheduled	
	$A_{AUC} \uparrow$	$A_{last} \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$		$A_{AUC} \uparrow$	$A_{last} \uparrow$	$A_{AUC} \uparrow$	$A_{last} \uparrow$
8	41.07±0.27	36.97±0.27	25.17±1.02	<b>34.14±1.18</b>	256	<b>41.92±0.06</b>	<b>37.64±1.09</b>	25.56±1.10	<b>34.08±1.02</b>
16	41.84±0.29	36.48±0.34	25.37±0.99	33.30±0.26	512	41.48±0.19	36.98±0.51	<b>25.85±1.39</b>	33.01±1.99
32	<b>41.92±0.06</b>	<b>37.64±1.09</b>	<b>25.56±1.10</b>	34.08±1.02	1024	41.28±0.18	36.48±0.34	25.13±0.90	31.83±0.19
64	40.13±0.15	33.50±0.46	24.39±0.85	29.48±0.04	2048	41.11±0.21	36.10±0.33	25.05±0.90	31.33±0.11

Table 10. REMIND performance as a function of different codebook sizes and number of codebooks with CIFAR-100. Original hyperparameters (codebook size: 256, number of codebooks: 32) consistently show the best performance in the online CL setting. The same hyperparameters were used uniformly for all datasets.

**MEMO.** MEMO [63] retrieves samples that are relatively close to the class mean feature at every task boundary and stores them in episodic memory for replay. However, in online CL, the model can not access all data for the current task, *i.e.*, it continuously updates the memory using stream data from the current task. Therefore, the sampling strategy of MEMO is replaced by class-balanced random sampling, which is used to replace the sampling strategy of RM [2] for online CL in [27].

## 18. Comparison between NC-FSCIL and Vanilla ETF

NC-FSCIL [56] is a recently proposed offline CL method that attempts to induce neural collapse in few-shot class incremental learning (FSCIL), a CL setup with few training samples per class. NC-FSCIL uses the fixed ETF classifier, a backbone network  $f$ , and a projection layer. It freezes the backbone network after training the base task and further fine-tunes the projection layer in the following incremental tasks by using the replay memory  $\mathcal{M}^{(t)}$ , which stores the class-wise mean features  $\mathbf{h}_c$  retrieved from the backbone network for each old class  $c$  as follows:

$$\mathcal{M}^{(t)} = \left\{ \mathbf{h}_c | c \in \bigcup_{j=0}^{t-1} C^{(j)} \right\}, \quad 1 \leq t \leq T, \quad (7)$$

where  $\mathbf{h}_c = \text{Avg}_i \{ f(x_i, \theta) | y_i = c \}$ ,  $T$  refers to the total number of tasks, and  $C^{(j)}$  refers to the set of classes for task  $j$ .

Considering that our setup is neither an offline CL nor a few-shot class incremental learning setup, for a fair comparison with online CL methods, we modify the freezing strategy of NC-FSCIL. Instead of freezing the whole backbone, we freeze only 6 layers, which achieved the best performance in online CL as shown in Table 9.

Despite the high performance of NC-FSCIL in offline CL, it has a lower performance than that of vanilla ETF (*i.e.*, baseline of EARL without preparatory data training and residual correction) in online CL, as shown in Tab. 11.

## 19. Pseudocode for the Our Method (L428)

Algorithm 1 and Algorithm 2 provide detailed pseudocode for EARL.

## 20. Detailed analysis of ablation results (L509)

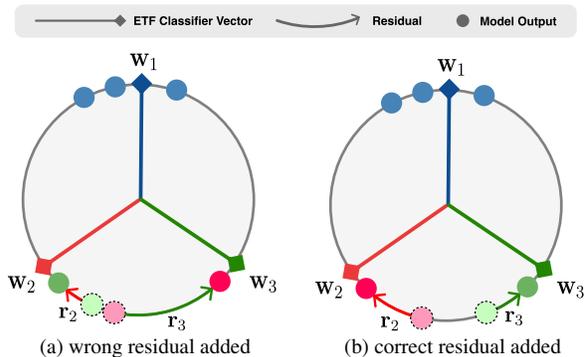


Figure 10. Without preparatory data training and relying solely on residual correction, incorrect residuals can be added due to bias, which can potentially lead to decreased performance. On the other hand, the combination of preparatory data training and joint training leads to the addition of correct residuals by addressing the bias problem.

Note that exclusively relying on residual correction without the training of preparatory data may result in the addition of incorrect residuals due to bias problem, as illustrated in Fig. 10-(a). The bias in CL causes the features of novel classes and old classes to overlap. When the features of multiple classes are clustered together, as in Fig. 4-(a), residuals of one class can be added to the residual-correcting term of other classes in the cluster, since we select the residuals using  $k$  nearest-neighbors of corresponding features. Thus, the residuals of old classes are often added to novel class samples and vice versa, hurting the accuracy of both old and novel classes.

In this context, the use of preparatory data not only accelerates the convergence of ETF during training, but also promotes accurate residual

Methods	CIFAR-10				CIFAR-100			
	Disjoint		Gaussian-Scheduled		Disjoint		Gaussian-Scheduled	
	$A_{\text{AUC}} \uparrow$	$A_{\text{last}} \uparrow$						
NC-FSCIL	68.09±0.69	48.05±1.69	53.12±0.77	46.11±0.36	39.54±0.79	34.83±1.27	30.91±1.46	35.73±1.13
Vanilla ETF	<b>75.27±0.77</b>	<b>62.10±4.12</b>	<b>65.80±0.25</b>	<b>67.79±0.78</b>	<b>52.91±1.05</b>	<b>41.08±0.60</b>	<b>33.34±0.55</b>	<b>44.45±0.48</b>

Table 11. Comparison between NC-FSCIL and Vanilla ETF (EARL w/o Preparatory data training and Residual correction) on CIFAR-10 and CIFAR-100.

---

### Algorithm 1 Training Phase

---

```

1: Input model  $f_\theta$ , Memory  $\mathcal{M}$ , Residual Memory  $\mathcal{M}_{\text{RES}}$ , Training data stream  $\mathcal{D}$ , ETF classifier  $\mathbf{W}$ , Negative transformation  $\mathcal{R}_r$ , Learning rate  $\mu$ 
2: for  $(x, y) \in \mathcal{D}$  do ▷ Sample arrives from training data stream  $\mathcal{D}$ 
3:   Update  $\mathcal{M} \leftarrow \text{ClassBalancedSampler}(\mathcal{M}, (x, y))$  ▷ Update memory
4:   Sample  $(X, Y) \leftarrow \text{RandomRetrieval}(\mathcal{M})$  ▷ Get batch  $(X, Y)$  from memory
5:   Sample  $(X', Y') \leftarrow \text{RandomRetrieval}(\mathcal{M})$  ▷ Get batch  $(X', Y')$  to make preparatory data
6:    $(X_p, Y_p) \leftarrow \mathcal{R}_r(X', Y')$  ▷ Negative transformation for preparatory data training
7:    $\hat{f}_\theta(X) = \frac{f_\theta(X)}{|f_\theta(X)|}, \hat{f}_\theta(X_p) = \frac{f_\theta(X_p)}{|f_\theta(X_p)|}$  ▷ Normalize model output
8:    $\mathbf{r} = \mathbf{W}_Y - \hat{f}_\theta(X)$  ▷ Calculate Residuals
9:   Update  $\mathcal{M}_{\text{RES}} \leftarrow (\hat{f}_\theta(X), \mathbf{r})$  ▷ Update feature-residual memory
10:   $\mathcal{L}(X, Y, X_p, Y_p; \theta, \mathbf{W}) = L_{DR}(\hat{f}_\theta(X), \mathbf{W}_Y) + L_{DR}(\hat{f}_\theta(X_p), \mathbf{W}_{Y_p})$  ▷ Calculate dot-regression loss
11:  Update  $\theta \leftarrow \theta - \mu \cdot \nabla_\theta \mathcal{L}(X, Y, X_p, Y_p; \theta, \mathbf{W})$  ▷ Update model
12: end for
13: Output  $f_\theta$ 

```

---

### Algorithm 2 Inference Phase

---

```

1: Input model  $f_\theta$ , inference input  $x_{\text{eval}}$ , Residual Memory  $\mathcal{M}_{\text{RES}}$ , ETF classifier  $\mathbf{W}$ , number of nearest neighbors  $k$ , softmax temperature  $\tau$ 
2:  $\{(\hat{h}_i, \mathbf{r}_i)\}_{i=1}^{|\mathcal{M}_{\text{RES}}|} \leftarrow \mathcal{M}_{\text{RES}}$  ▷ Get residual and features from residual memory
3:  $\hat{f}_\theta(x_{\text{eval}}) = \frac{f_\theta(x_{\text{eval}})}{|f_\theta(x_{\text{eval}})|}$  ▷ Normalize model output
4:  $\{n_1, n_2, \dots, n_k\} \leftarrow k\text{-arg min}_i \left( \left| \hat{f}_\theta(x_{\text{eval}}) - \hat{h}_i \right| \right)$  ▷ Calculate  $k$  nearest neighbor features
5:  $s_{n_i} = \frac{e^{-(\hat{f}_\theta(x_{\text{eval}}) - \hat{h}_{n_i})/\tau}}{\sum_{j=1}^k e^{-(\hat{f}_\theta(x_{\text{eval}}) - \hat{h}_{n_j})/\tau}}$  ▷ calculate residual weights
6:  $\mathbf{r} = \sum_{i=1}^k s_{n_i} \mathbf{r}_{n_i}$  ▷ Calculate residual-correcting term
7:  $\hat{f}_\theta(x_{\text{eval}})_{\text{corrected}} \leftarrow \hat{f}_\theta(x_{\text{eval}}) + \mathbf{r}$  ▷ Add residual on features
8:  $y_{\text{pred}} = \arg \max_y (\text{CosineSimilarity}(\mathbf{W}_y, \hat{f}_\theta(x_{\text{eval}})_{\text{corrected}}))$  ▷ Predict class
9: Output  $y_{\text{pred}}$ 

```

---

addition during inference. In conclusion, the combination of residual correction and preparatory data training effectively aligns the model output with the corresponding ETF classifier, as demonstrated in Fig. 10-(b).

## 21. Properties of Neural Collapse

**(NC1) Collapse of Variability:** The last layer feature output of each data point collapses toward the class mean feature of its respective class. In other words,  $h_{k,i}$ , last layer feature of sample  $i$  in class  $k$ , collapse to  $\mu_k = \sum_{i=1}^{n_k} h_{k,i}$  for  $\forall k \in [1, K]$  where  $n_k$  is the number of samples for class

$k$ . By considering within-class covariance and between-class covariance

$$\begin{aligned} \Sigma_W &= \frac{1}{K} \sum_{k=1}^K \frac{1}{n_k} \left( \sum_{i=1}^{n_k} (h_{k,i} - \mu_k) \right), \\ \Sigma_B &= \frac{1}{K} \sum_{k=1}^K (\mu_k - \mu_G), \end{aligned} \quad (8)$$

where  $\mu_G = \sum_{k=1}^K \mu_k$ , empirical variability can be measured as

$$NC1 := \frac{1}{K} \text{trace}(\Sigma_W \Sigma_B^\dagger). \quad (9)$$

**(NC2) Convergence to simplex equiangular tight frame (ETF):** Class means  $\mu_k (k \in [1, K])$  centered by the global mean  $\mu_G$  converge to

vertices of a simplex ETF structure, i.e., matrix  $M = [m_1 \ m_2 \ \cdots \ m_K]$  where  $m_k = \frac{\mu_k - \mu_G}{\|\mu_k - \mu_G\|^2}$  satisfies the following equation:

$$\mathbf{M}\mathbf{M}^T = \frac{1}{K-1}(KI_K - \mathbf{1}_K\mathbf{1}_K^T). \quad (10)$$

The degree of convergence can be measured using:

$$\frac{\mathbf{M}\mathbf{M}^T}{\|\mathbf{M}\mathbf{M}^T\|_F} - \frac{1}{\sqrt{K-1}}(I_K - \frac{1}{K}\mathbf{1}_K\mathbf{1}_K^T). \quad (11)$$

**(NC3) Convergence to self-duality:** Classifier  $\mathbf{W}$  converges to the simplex ETF  $\mathbf{M}$  formed by recentered feature mean, and during this convergence, the classifier vector  $w_k$  aligns with their corresponding feature mean  $m_k$  where  $w_k$  means classifier weight for class  $k$ ,  $k \in [1, K]$ , i.e.,

$$\frac{\mathbf{M}}{\|\mathbf{M}\|_F} = \frac{\mathbf{W}}{\|\mathbf{W}\|_F} \quad (12)$$

Duality can be measured by measuring:

$$\frac{\mathbf{W}\mathbf{M}^T}{\|\mathbf{W}\mathbf{M}^T\|_F} - \frac{1}{\sqrt{K-1}}(I_K - \frac{1}{K}\mathbf{1}_K\mathbf{1}_K^T). \quad (13)$$