

# OrCo: Towards Better Generalization via Orthogonality and Contrast for Few-Shot Class-Incremental Learning

Noor Ahmed\*      Anna Kukleva\*      Bernt Schiele  
 {noahmed, akukleva, schiele}@mpi-inf.mpg.de  
 Max Planck Institute for Informatics, Saarland Informatics Campus

## Abstract

*Few-Shot Class-Incremental Learning (FSCIL) introduces a paradigm in which the problem space expands with limited data. FSCIL methods inherently face the challenge of catastrophic forgetting as data arrives incrementally, making models susceptible to overwriting previously acquired knowledge. Moreover, given the scarcity of labeled samples available at any given time, models may be prone to overfitting and find it challenging to strike a balance between extensive pretraining and the limited incremental data. To address these challenges, we propose the OrCo framework built on two core principles: features’ orthogonality in the representation space, and contrastive learning. In particular, we improve the generalization of the embedding space by employing a combination of supervised and self-supervised contrastive losses during the pretraining phase. Additionally, we introduce OrCo to address challenges arising from data limitations during incremental sessions. Through feature space perturbations and orthogonality between classes, the OrCo loss maximizes margins and reserves space for the following incremental data. This, in turn, ensures the accommodation of incoming classes in the feature space without compromising previously acquired knowledge. Our experimental results showcase state-of-the-art performance across three benchmark datasets, including mini-ImageNet, CIFAR100, and CUB datasets. Code is available at: <https://github.com/noorahmedds/OrCo>.*

## 1. Introduction

Real-world applications frequently encounter various challenges when acquiring data incrementally, with new information arriving in continuous portions. This scenario is commonly referred to as Class Incremental Learning (CIL) [2–4, 17, 18, 28, 29, 35, 37, 43, 47, 52]. Within

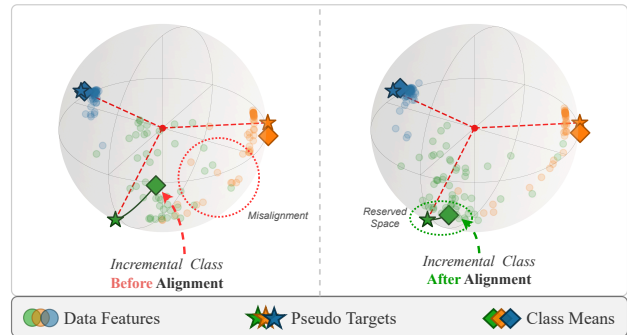


Figure 1. **PCA analysis on feature space before and after alignment.** Left: Before aligning incremental classes to orthogonal pseudo-targets. Right: After aligning incremental classes to assigned targets using **OrCo** loss. Our loss effectively reduces misalignment. Additionally, it enhances generalization for incoming classes by explicitly reserving space.

CIL, the foremost challenge lies in preventing catastrophic forgetting [13, 23, 31], where previously learned concepts are susceptible to being overwritten by the latest updates. However, in a Few-Shot Class-Incremental Learning (FSCIL) scenario [8, 16, 16, 22, 26, 30, 39, 44, 46, 48–50], characterized by the introduction of new information with only a few labeled samples, two additional challenges emerge: overfitting and intransigence [7, 38]. Overfitting arises as the model may memorize scarce input data and lose its generalization ability. On the other hand, intransigence involves maintaining a delicate balance, preserving knowledge from abundant existing classes while remaining adaptive enough to learn new tasks from a highly limited dataset. Advances in dealing with catastrophic forgetting, overfitting and intransigence are important steps toward improving the practical value of these methods.

Catastrophic forgetting is commonly tackled in CIL methods [18, 20, 35], which assume ample labeled training data. However, standard CIL methods struggle in scenarios with limited labeled data, such as FSCIL [39]. To

\* Equal Contribution

address the three challenges posed by FSCIL, recent approaches [46, 49] focus primarily on regularizing the feature space during incremental sessions, mitigating the risk of overfitting. These methods rely on a frozen backbone pretrained with standard cross-entropy on a substantial amount of data from the base session. However, we argue that achieving high performance on the pretraining dataset may not necessarily result in optimal generalization in subsequent incremental sessions with limited data. Therefore, in the first phase, we propose enhancing feature space generalization through contrastive learning, leveraging data from the base session.

In this work, we introduce the OrCo framework, a novel approach built on two fundamental pillars: features’ mutual orthogonality on the representation hypersphere and contrastive learning. During the first phase, we leverage supervised [21] and self-supervised contrastive learning [6, 14, 33] for pretraining the model. The interplay between these two learning paradigms enables the model to capture various types of semantic information that is particularly beneficial for the novel classes with limited data [1, 5, 19], implicitly addressing the *intransigence* issue. After the pretraining, we generate and fix mutually orthogonal random vectors, further referred to as pseudo-targets. In the second phase, we aim to align the fixed pretrained backbone to the pseudo-targets using abundant base data. The learning objective during this phase is our OrCo loss, which consists of three integral components: perturbed supervised contrastive loss (PSCL), loss term that enforces orthogonality of features in the embedding space, and standard cross-entropy loss. Notably, our PSCL leverages generated pseudo-targets to maximize margins between the classes and to preserve space for incremental data, enhancing orthogonality through contrastive learning (see figure 1). The third phase of our framework, which we apply in each subsequent incremental session, similarly aims to align the model with the pseudo-targets, but using only few-shot data from the incremental sessions. During the third phase, our PSCL addresses limited data challenges, mitigating the *overfitting* problem to the current incremental session and *catastrophic forgetting* of the previous sessions through margin maximization.

We summarize the contributions of this work as follows:

- We introduce the novel OrCo framework designed to tackle FSCIL, that is built on orthogonality and contrastive learning principles throughout both pretraining and incremental sessions.
- Our perturbed supervised contrastive loss introduces perturbations of orthogonal, data-independent vectors in the representation space. This approach induces increased margins between classes, enhancing generalization.
- We showcase robust performance on three datasets, outperforming previous state-of-the-art methods. Further-

more, we perform a thorough analysis to evaluate the importance of each component.

## 2. Related Work

**Few-Shot Learning (FSL).** In FSL, a model is trained on scarce data with just a few samples per class. Current literature can be divided into two predominant categories: optimisation-based [12, 32, 34] and metric-based methods [34, 38, 41]. Optimization-based approaches, such as MAML [12] and Reptile [32], find optimal parameters that can generalize quickly on other sets when subjected to fine-tuning. Conversely, metric-based methods utilize a pretrained model and compare support and query instances using similarity metrics. For example, Prototypical Networks [38] learns a metric space where the distance to class prototypes determine classification. And imprinting weights method [34] shows improved performance by using class means as strong initialisation for an evolving classifier and integrates principles from both categories.

**Class-Incremental Learning (CIL).** In the domain of CIL, a sequence of novel concepts must be learned without forgetting previously acquired knowledge. Recent works can be coarsely categorized in 3 groups. Foremost, there are knowledge distillation schemes [17, 28, 35, 43], which retain model behaviour across the adaptation process to avoid forgetting. Then, data-replay methods [2, 3, 18, 47, 52] show strong resistance to catastrophic forgetting by storing old class exemplars. Lastly, weight consolidation methods [4, 23, 29, 37] identify important weights and moderate training regimen.

**Few-Shot Class-Incremental Learning (FSCIL).** FSCIL paradigm demands rapid adaptation to novel classes with limited data. The methods can be divided into the following categories [40]: geometry preservation methods [39, 46], replay or distillation strategies [26, 49], metric learning methods [8, 22, 48, 50] and meta-learning [16, 30, 44], highlighting the breadth of methodologies in FSCIL. Parallel to our methodology, metric learning methods utilize tricks in the feature space, showcasing diverse approaches for accommodating incremental classes. FACT [50] creates virtual prototypes to reserve space and scale the model for incoming classes. NC-FSCIL [44], aligns class features with the classifier prototypes, which are formed as a simplex equiangular tight frame, using dot-regression loss. C-FSCIL [16] aligns class prototypes quasi-orthogonally to negate interference between classes. In contrast, our approach stands out for its use of contrastive learning with data agnostic pseudo-targets and margin maximization through perturbations in the embedding space improving generalization in incremental sessions.

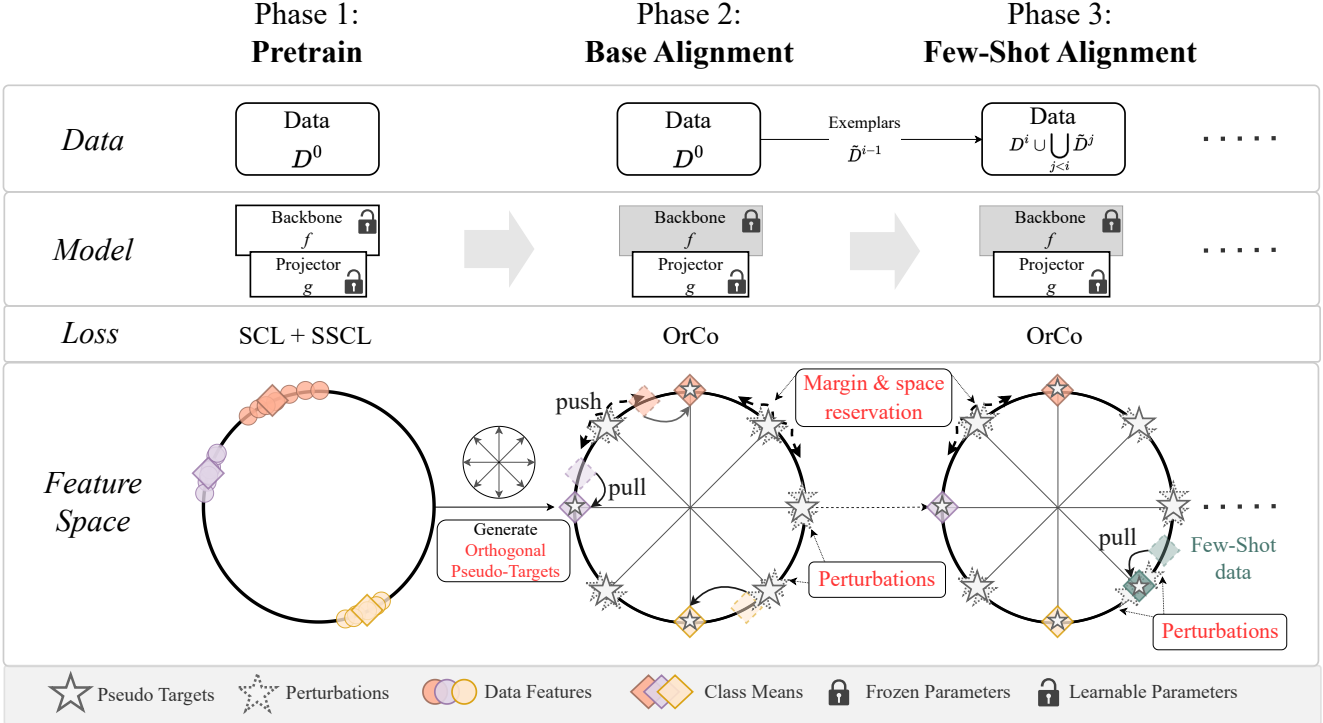


Figure 2. **Overview of OrCo framework.** Our OrCo framework is a three-phase approach for FSCIL. Phase 1 (Pretrain): We pretrain both backbone and projection head with SCL and SSCL on base dataset  $D^0$ . Before the next phase, we generate mutually orthogonal pseudo-targets. Phase 2 (Base Alignment): We aim to align the base dataset  $D^0$  to the pseudo-targets through our OrCo loss. This involves pulling class features towards the nearest pseudo-targets and pushing forces based on perturbations around unassigned pseudo-targets (grey stars without assigned colored class means) to increase the margin and preserve space for incoming classes. Phase 3 (Few-Shot Alignment): Phase 3, employed in each subsequent incremental session, is similar to Phase 2 and assigns pseudo-targets to incremental class means with further alignment using our OrCo loss.

### 3. OrCo Framework

We begin with necessary preliminaries in section 3.1, followed by the description of our OrCo framework in section 3.2 and OrCo loss in section 3.3.

#### 3.1. Preliminaries

**FSCIL Setting.** FSCIL consists of multiple incremental sessions. An initial 0-th session is often reserved to learn a generalisable representation on an abundant base dataset. This is followed by multiple few-shot incremental sessions with limited data. To formalise, an M-Session N-Way and K-Shot FSCIL task consists of  $D_{seq} = \{D^0, D^1, \dots, D^M\}$ . These are all the datasets written in sequence where  $D^i = \{(x_i, y_i)\}_{i=1}^{|D^i|}$  is the dataset for the  $i$ -th session. The 0-th session dataset  $D^0$ , also referred to as base dataset, consists of  $C^0$  classes, each with a large number of samples. The training set for each following few-shot incremental session ( $i > 0$ ) has  $N$  classes. Each of these classes has  $K$  samples, typically ranging from 1 to 5 samples per class. Taking the  $i$ -th session as an example, the model’s performance is as-

essed on validation sets from the current ( $i$ -th) and all previously encountered datasets ( $< i$ ). The entire FSCIL task comprises a total of  $C$  classes. In our OrCo framework, we use base dataset  $D^0$  for pretraining the model during phase 1 and for base alignment during phase 2.

**Target Generation.** We employ a Target Generation loss, similarly as in [27], to generate mutually orthogonal vectors across the representation hypersphere with a dimensionality of  $d$ . First, we define a set of random vectors  $T = \{t_i\}$  where  $\{t_i\} \in \mathbb{R}^d$ . The optimization of the following loss with respect to these random vectors maximizes the angle between any pair of vectors  $t_i, t_j \in T$ , thereby ensuring their mutual orthogonality:

$$\mathcal{L}_{TG}(T) = \frac{1}{|T|} \sum_{i=1}^{|T|} \log \sum_{j=1}^{|T|} e^{t_i \cdot t_j / \tau_o} \quad (1)$$

where  $\tau_o$  is the temperature parameter. These optimized vectors, which we further refer to as pseudo-targets, remain fixed throughout our training process.

**Contrastive Loss.** The objective of contrastive representation learning is to create an embedding space where similar sample pairs are in close proximity, while dissimilar pairs are distant. In this work, we adopt the InfoNCE loss [21, 33] as our contrastive objective. With positive set  $P_i$  and negative set  $N_i$  defined for each data sample  $z_i$ , called anchor, this loss aims to bring any  $z_j \in P_i$  closer to its anchor  $z_i$  and push any  $z_k \in N_i$  further away from the anchor  $z_i$ :

$$\mathcal{L}_{CL}(i; \theta) = \frac{-1}{|P_i|} \sum_{z_j \in P_i} \log \frac{\exp(z_i \cdot z_j / \tau)}{\sum_{z_k \in N_i} \exp(z_i \cdot z_k / \tau)}, \quad (2)$$

where  $\tau$  is the temperature parameter. In the classical self-supervised contrastive learning (SSCL) scenario, where labels for individual instances are unavailable, the positive set comprises augmentations of the anchor, and all other instances are treated as part of the negative set [33]. In contrast, for the supervised contrastive loss (SCL), the positive set includes all instances from the same class as the anchor, while the negative set encompasses instances from all other classes [21]. To enhance clarity, we denote the supervised contrastive loss as  $\mathcal{L}_{SCL}$  and self-supervised contrastive loss as  $\mathcal{L}_{SSCL}$ . We consider SCL and SSCL as the cornerstone guiding our work due to their discriminative nature, robustness, and extendability. We leave formal definition of the losses for the supplement.

### 3.2. OrCo Framework

**Overview.** Our OrCo framework (see figure 2) for FSCIL begins with a pretraining of the model in the first phase, focusing on learning representations, which are transferable to the new tasks. To achieve this, we leverage both supervised and self-supervised contrastive losses [5, 19]. Before the second phase, we generate a set of mutually orthogonal vectors, which we term as pseudo-targets. In the subsequent second phase, referred to as base alignment in figure 2, we allocate pseudo-targets to class means and ensure alignment through our OrCo loss, using abundant base data  $D^0$ . The third phase, implemented in each subsequent incremental session, similarly focuses on assigning incoming but few-shot data to unassigned pseudo-targets, followed by alignment through our OrCo loss. Our OrCo loss comprises three key components: cross-entropy, orthogonality loss, and our novel perturbed supervised contrastive loss (PSCL). The cross-entropy loss aligns incremental data with assigned fixed orthogonal pseudo-targets, the orthogonality loss enforces a geometric constraint on the entire feature space to mimic the pseudo-targets distribution, and our PSCL enhances crucial robustness for FSCIL tasks through margin maximization and space reservation, leveraging mutual orthogonality of pseudo-targets.

**Phase 1: Pretrain.** In the first pretraining phase, we

learn an encoder that accumulates knowledge and generates distinctive features. Using a combination of supervised contrastive loss (SCL) and self-supervised contrastive loss (SSCL), we enhance feature separation within classes, improving model transferability to incremental sessions [5, 19]. To this end, we train the model encoder  $f$  and MLP projection head  $g$  using base data  $D^0$ , mapping input images to  $\mathcal{R}^d$  feature space. The pretraining loss is then defined as:

$$\mathcal{L}_{pretrain}(D^0; f, g) = (1 - \alpha) * \mathcal{L}_{SCL} + \alpha * \mathcal{L}_{SSCL}, \quad (3)$$

where  $\alpha$  controls the contribution of each contrastive loss.

**Pseudo-targets.** During the first phase, we do not employ any explicit class vectors that can be used for linear classification. Therefore, we generate data-independent mutually orthogonal pseudo-targets  $T = \{t_j\} \in \mathcal{R}^d$  on the hypersphere by optimizing loss shown in equation 1, where  $|T| \geq C$ . Further, these pseudo-targets are fixed and assigned to classes, which, in turn, maximize margins between the classes and improve generalization.

**Phase 2: Base Alignment.** In addition to the pretraining phase, we introduce the second phase based on the base dataset  $D^0$ . This phase initiates alignment between the projection head  $g$  and the set of generated pseudo-targets  $T$ . More specifically, we create class means by averaging features with the same labels. Then, we employ a one-to-one matching approach, utilizing the Hungarian algorithm [25], to assign class means with the most fitting set of pseudo-targets  $T^0$ , where  $|T^0| = |C^0|$ . Note that each class  $y_j \in C^0$  is then associated with the respective pseudo-target  $t_j^0 \in T^0$  and we denote the remaining unassigned pseudo-targets as  $T_u^0 = T \setminus T^0$ . Further, we use pseudo-targets  $T^0$  as base class representations for classification. Despite the optimal initial assignment, we enhance the alignment of the projection head  $g$  and the respective pseudo-targets  $T^0$  through the optimization of our OrCo loss. Further insights into the specifics and motivation behind our OrCo loss, we elaborate on in section 3.3.

**Phase 3: Few-Shot Alignment.** In the third phase of our framework, applied in each subsequent incremental session, our goal remains to align incoming data with the pseudo-targets. The following sessions introduce few shot incremental data  $D^i$  for the  $i$ -th session. Building on previous methods [3, 8, 26, 43], we maintain some exemplars from previously seen classes, constituting a joint set  $D_{joint}^i = \{D^i \cup \{\cup_{j=0}^{i-1} \tilde{D}^j\}\}$ , where  $\tilde{D}^j$  denotes saved exemplars from earlier sessions. Keeping random exemplars from previous sessions serves to mitigate both overfitting and catastrophic forgetting issues. Similarly to the second phase, we assign pseudo-targets to incremental class means. To be specific, we determine the optimal assignment between  $T_u^{i-1}$  and the current class means, resulting in the optimal

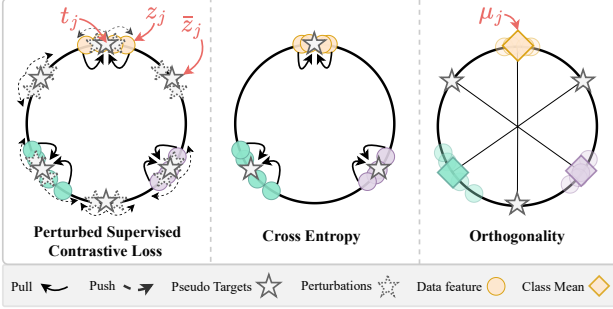


Figure 3. **OrCo loss** consists of three components: our proposed perturbed supervised contrastive loss (PSCL), cross-entropy loss (CE), and orthogonality loss (ORTH).  $z_j$  denotes the real data anchor point for a contrastive loss,  $\bar{z}_j$  denotes the unassigned pseudo-target anchor, and  $t_j$  denotes an additional positive sample for yellow class in the form of an assigned pseudo-target.  $\mu_j$  represents the within-batch mean features.

assignment set of pseudo-targets  $T^i$ . Respectively, the unassigned set of pseudo-targets becomes  $T_u^i = T_u^{i-1} \setminus T^i$ . During incremental session  $i$ , we optimize our OrCo loss given data  $D_{joint}^i$ , pseudo-targets  $T = \{T^i\}_0^i \cup T_u^i$  and the assignment between  $T = \{T^i\}_0^i$  and the respective classes.

### 3.3. OrCo Loss

During the second and the third phases, we optimize the parameters of the projection head  $g$  with our OrCo loss. This loss comprises three integral components: our novel perturbed supervised contrastive loss (PSCL), cross-entropy loss (CE) and orthogonality loss (ORTH), see figure 3. The aim of optimizing the OrCo loss is to align classes with their assigned pseudo-targets, simultaneously maximizing the margins between classes. This, in turn, enhances overall generalization performance.

To maximize the margins in the representation space, we introduce uniform perturbations of the pseudo-targets  $T$ , resulting in perturbed pseudo-targets  $\tilde{T} = \{\tilde{t}_j\}_0^{|T|}$  defined as

$$\tilde{t}_j = t_j + \mathcal{U}(-\lambda, \lambda), \quad (4)$$

where  $\mathcal{U}$  stands for uniform distribution and  $\lambda$  defines sampling boundaries. To utilize the introduced perturbations, we redefine positive  $P_j$  and negative  $N_j$  sets for the contrastive loss for the anchor  $z_j$  in equation 3. Note that during incremental session  $i$  the positive set  $P_j^i$  in the standard SCL contains all  $z_k \in D_{joint}^i$  such that the label  $y_k$  is equal to anchor label  $y_j$ , e.g. in figure 3, all yellow circles belong to the positive set for yellow anchor  $z_j$ . And the negative set consists of remaining samples  $N_j^i = D_{joint}^i \setminus P_j^i$ , in figure 3, the negative set is composed of all other colors.

To adapt standard SCL to PSCL, we expand the definition of the positive set. The anchor  $z_j$ , with its assigned

pseudo-target  $t_j \in T$ , becomes an additional positive pair, see figure 3. Furthermore, considering the previously defined pseudo-target perturbations, we incorporate them into the positive set, resulting in  $\tilde{P}_j^i = P_j^i \cup t_j \cup \tilde{t}_j$ . In figure 3, the positive set for yellow anchor  $z_j$  contains all yellow circles and additionally the pseudo-target  $t_j$  with its perturbations. This extension of the positive set introduces additional pushing forces for incremental classes and, therefore, enables the maximization of margins between classes. We show that this approach proves to be especially advantageous in scenarios with limited samples, as it mimics augmentations in the feature space.

On the other hand, we expand the anchor definition. In standard SCL, each anchor  $z_j$  belongs to the set of real training data  $D_{joint}^i$ , e.g. in figure 3, anchors for standard SCL are only circles. However, we propose to use anchors from both real data and unassigned pseudo-targets (circles and unassigned stars in figure 3), specifically  $z_j \in D_{joint}^i$  and  $\bar{z}_j \in T_u^i$ . The positive set for the anchor  $\bar{z}_j \in T_u^i$  (unassigned pseudo-target) contains only corresponding perturbed pseudo-targets  $\tilde{P}_j^i = \{\tilde{t}_j\}$  (dashed stars around  $\bar{z}_j$  in figure 3), while the negative set  $\tilde{N}_j^i = \{D_{joint}^i \cup \tilde{T}_u^i\} \setminus \{\tilde{t}_j\}$  includes all real data and other perturbed pseudo-targets. This approach ensures that each unassigned pseudo-target pushes all other classes away, thereby promoting space preservation for the following incremental sessions.

To complement PSCL, we employ cross-entropy loss for sample  $z_j$  that pulls class features to their assigned targets during the few-shot incremental session  $i$ :

$$\mathcal{L}_{CE}(z_j) = - \sum_{c \in \{C^I\}_1^i} y_c \log \frac{\exp(z_j t_c^T)}{\sum_{k \in \{C^I\}_1^i} \exp(z_k t_c^T)}. \quad (5)$$

We further employ the orthogonality loss ( $\mathcal{L}_{ORTH}$ ) defined similarly as in equation 1. It differs, however, in that it uses the mean class features  $\mu_j$  (see figure 3) as input and enforces an intrinsic geometric constraint on the feature landscape. See section 10 in supplementary for details.

Finally, our OrCo loss is a combination of the three losses introduced above:

$$\mathcal{L}_{OrCo} = \mathcal{L}_{PSCL} + \mathcal{L}_{CE} + \mathcal{L}_{ORTH}. \quad (6)$$

During testing, a sample is assigned a label based on the nearest assigned pseudo target.

## 4. Experimental Results

In section 4.1, dataset and evaluation protocol are introduced. Then we compare with state-of-the-art methods on 3 popular benchmarks in section 4.2. Finally, we validate the effectiveness of each of the components in section 4.3.

Method	Base Acc	Session-wise Harmonic Mean (%) $\uparrow$								aHM	$\Delta$ aHM
		1	2	3	4	5	6	7	8		
IW [34]	83.10	49.49	45.09	45.98	46.30	44.67	42.48	43.26	45.65	45.36	<b>+12.76</b>
FACT [50]	75.78	27.20	27.84	27.94	25.17	22.46	20.54	20.88	21.25	24.16	<b>+33.96</b>
CEC [46]	72.17	31.91	31.84	30.98	30.74	28.14	26.78	26.96	27.42	29.35	<b>+28.78</b>
C-FSCIL [16]	76.60	9.74	20.53	28.68	31.91	34.85	35.05	37.72	37.92	29.55	<b>+28.57</b>
LIMIT [51]	73.27	40.34	33.58	31.81	31.74	29.32	29.11	29.57	30.28	31.97	<b>+26.15</b>
LCwoF [26]	64.45	41.24	38.96	39.08	38.67	36.75	35.47	34.71	35.02	37.49	<b>+20.63</b>
BiDist [49]	74.67	42.42	43.86	43.87	40.34	38.97	38.01	36.85	38.47	40.35	<b>+17.77</b>
NC-FSCIL [44]	<b>84.37</b>	62.34	61.04	55.93	53.13	49.68	47.08	46.22	45.57	52.62	<b>+5.50</b>
OrCo	83.30	<b>68.71</b>	<b>63.87</b>	<b>60.94</b>	<b>57.98</b>	<b>55.27</b>	<b>52.41</b>	<b>52.68</b>	<b>53.12</b>	<b>58.12</b>	

Table 1. **Sota comparison on mini-ImageNet.** aHM denotes the average of the harmonic mean across all sessions. IW [34] is evaluated based on the model learning in our pretrain phase. Detailed results of the individual sessions are in the supplement.

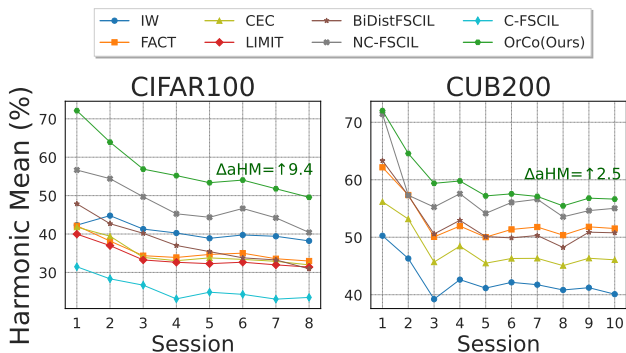


Figure 4. **Sota comparisons on CIFAR100 and CUB200 datasets.** Performance curves, that measure harmonic mean, of our method comparing to recent sota methods. Left: CIFAR100. Right: CUB200.  $\Delta$ aHM denotes the average harmonic mean improvement over the runner-up method.

#### 4.1. Datasets and Evaluation

We conduct evaluation of our OrCo framework on three FSCIL benchmark datasets: mini-ImageNet [36], CIFAR100 [24] and CUB200 [42]. In the setting formalised by [39] mini-ImageNet and CIFAR100 are organized into 60 base classes and 40 incremental classes structured in a 5-way, 5-shot FSCIL scenario for a total of 8 sessions. CUB200, a dataset for fine-grained bird species classification, contains equal number of base and incremental classes for a total of 200 classes. The dataset is organised as a 10-way, 5-shot FSCIL task and presents a rigorous challenge.

**Performance measure** Commonly used FSCIL datasets all have a quantity bias towards the base classes. mini-ImageNet and CIFAR100 have both 60% of the data in the base classes and CUB200 with 50%. Consequently, standard accuracy measures like Top-1 accuracy will be skewed

in favour of the base-classes. For instance, a method which has a base accuracy  $A_{base} = 100\%$  on CUB200 and performs weakly on the first incremental session  $A_{inc}^1 = 10\%$  would produce a Top-1 average accuracy  $A_{cls}^1 = 91.82\%$ . At first glance, this accuracy may not entirely represent inherent biases in a method, though such measures are commonly used to benchmark performance. To tackle this, harmonic mean has risen as a robust evaluation measure in FSCIL [26, 49, 50]. In the given scenario, the harmonic mean would penalise the method aggressively resulting in a metric score of  $A_{hm}^1 = 18.18\%$ , accurately indicating bias. More concretely, we compute harmonic mean by combining base class accuracy and incremental session accuracy:  $A_{hm}^j = (2 \times A_{base} \times A_{inc}^j) / (A_{base} + A_{inc}^j)$ . In addition to this, we propose average harmonic mean (aHM) which is simply averaging the harmonic mean scores from all sessions for a consolidated view.

**Implementation Details** Our model is optimised using LARS [45] for the pretraining phase and SGD with momentum for phase 2 and 3. For CUB200 dataset, we skip the pretraining following [39, 44, 46] and initialize the model with ImageNet pretrained weights. For the second and third phase, we finetune only the projection head. For the PSCL loss we choose a perturbation magnitude  $\lambda = 1e-2$ . We train the projection head for 10 epochs during the second phase and 100 epochs for the third phase. Cosine scheduling is employed with a maximum learning rate set to 0.1. Augmentations include, random crop, random horizontal flip, random grayscale and a random application of color jitter. Details can be found in the supplement section 11.

#### 4.2. Comparison to state-of-the-art

In this section, we conduct a comparative analysis of our proposed OrCo with recent state-of-the-art approaches. Ta-

PSCL	CE	ORTH	Session-wise Harmonic Mean (%) $\uparrow$								aHM
			1	2	3	4	5	6	7	8	
	✓		65.46	56.29	44.12	36.96	26.90	21.11	18.90	16.19	35.74
	✓	✓	65.30	56.21	43.96	37.30	28.31	22.01	19.66	16.64	36.17
✓			50.70	45.42	42.68	39.84	38.71	37.94	36.26	35.87	40.93
✓		✓	52.34	47.24	43.79	41.62	41.15	39.68	38.69	37.34	42.73
✓	✓		68.04	<b>63.94</b>	60.22	<b>58.00</b>	<b>55.44</b>	51.51	51.88	52.74	57.72
✓	✓	✓	<b>68.71</b>	63.87	<b>60.94</b>	57.98	55.27	<b>52.41</b>	<b>52.68</b>	<b>53.12</b>	<b>58.12</b>

Table 2. **Influence of OrCo loss components.** PSCL denotes perturbed supervised contrastive loss, CE denotes cross-entropy, ORTH denotes orthogonality loss. See figure 3 for visualization of each component. Ablation study on mini-ImageNet.

ble 1 presents the results obtained on the mini-ImageNet dataset, while figure 4 illustrates the evaluation results on the CUB200 and CIFAR100 datasets. Our method demonstrates superior performance across all three datasets, surpassing previous state-of-the-art methods by a significant margin, particularly achieving improvements of 9.4% and 5.5% on CIFAR100 and mini-ImageNet, respectively. Notably, the effectiveness of OrCo is consistently evident across all incremental sessions.

In addition to reporting results for the standard FSCIL methods, we also present the performance of the Imprinted Weights method (IW) [34] based on our model pretrained during Phase 1. The robust performance of this method indicates the efficacy of our pretraining strategy in facilitating effective transferability to downstream tasks, such as incremental few-shot learning sessions, thereby addressing the intransigence problem. We present a detailed breakdown of each session in section 7 of the supplement.

### 4.3. Analysis

To validate the effectiveness of each component of our framework, in this section, we show an analysis based on the mini-ImageNet dataset.

**OrCo loss.** We assess the efficacy of the components comprising our OrCo loss in table 2. The OrCo loss consists of three integral components illustrated in figure 3: the cross-entropy loss (CE), the orthogonality loss (ORTH), and the perturbed supervised contrastive loss (PSCL). We observe that CE struggles to generalize on underrepresented incremental classes. On the contrary, PSCL enhances the robust SCL approach with pseudo target perturbations and provides better class separation. PSCL, on its own, shows steady generalization, with only a 14.83% drop in harmonic mean. CE, however, while starting strong, ultimately becomes biased towards base classes, leading to a significant 49.26% drop in harmonic mean. By integrating the dynamic yet fundamentally discerning features of CE with the stability offered by PSCL, a significant enhancement in har-

Sampling	aHM
Rand	57.23
Orth	<b>58.12</b>

Table 3. **Importance of explicit orthogonality loss for pseudo-target generation.** Rand denotes random sampling from normal distribution.

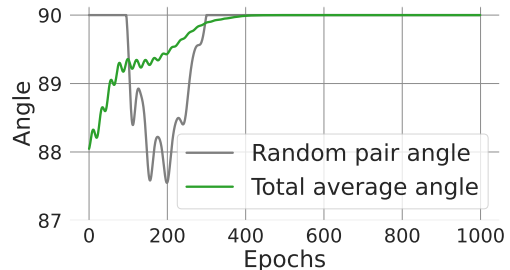


Figure 5. **Measurement of angle during orthogonality optimization.** The green curve corresponds to the evolution of the average angle between all pairs during the optimization. The gray curve shows measurements of random pairs at each epoch.

monic mean is achieved. Moreover, the orthogonality loss (ORTH) consistently improves performance (+0.4%). Its integration into our loss formulation further underscores the significance of orthogonality.

**Influence of mutually orthogonal pseudo-targets.** We note that independent randomly sampled vectors from a Gaussian distribution  $\mathcal{N}(0, 1)$  are theoretically orthogonal on the surface of a unit sphere (refer to the supplement for a discussion). However, in practice, we observe only a near-orthogonal behavior, as illustrated in our training curve for orthogonal pseudo-target generation in figure 5. To examine the impact of a perfectly orthogonal target space, we assess our method using both randomly sampled targets from a Gaussian distribution and our generated orthogonal targets, as presented in table 3. We observe a 0.89% improvement in performance when explicit orthogonality constraints are applied. This finding suggests that an aligned and orthogonal feature space is more effective in addressing data imbalances between base and incremental sessions. Consequently, we incorporate orthogonality as a fundamental principle in our framework, recognizing its significant role in enhancing the overall effectiveness of our model.

**Pseudo-targets perturbations.** OrCo relies on perturba-

Perturbation	$FP_{inc} \downarrow$	$Sim_{cls} \downarrow$	$Sim_{cls \rightarrow target} \downarrow$	$HM_8 \uparrow$
w/o	66.5	0.105	0.013	20.14
$\mathcal{N}$	54.6	<b>0.002</b>	<b>0.006</b>	50.23
$\mathcal{U}$	<b>52.5</b>	0.011	<b>0.006</b>	<b>53.12</b>

Table 4. **Influence of perturbations in PSCL.** Comparison of our PSCL loss with or without perturbations of pseudo-targets.  $\mathcal{N}$ ,  $\mathcal{U}$  denotes Gaussian and Uniform distributions, respectively, from which  $\lambda$ , as in equation 4, sampled during training.  $FP_{inc}$  refers to the False Positive rate among all incremental classes.  $Sim_{cls}$  computes the average pairwise cosine similarity between all class pairs.  $Sim_{cls \rightarrow target}$  indicates the pairwise cosine similarity between classes and unassigned target pairs over all sessions.  $HM_8$  refers to the 8-th and final session harmonic mean.

tions of fixed pseudo-targets to introduce a margin between previously encountered and incoming classes. We compare OrCo against a variant where the contrastive loss does not receive any pseudo-targets’ perturbations (w/o). In contrast to this, our perturbation schemes with sampling  $\lambda$ , as in equation 4, from Gaussian ( $\mathcal{N}$ ) and uniform ( $\mathcal{U}$ ) distributions consistently enhance the final session harmonic mean ( $HM_8$ ) by over 30%, as shown in table 4.

For a detailed evaluation of our method, we employ false positive and cosine similarity analyses. By measuring the false positive rate within only incremental classes ( $FP_{inc}$ ), we observe improved separation between the few-shot classes with the perturbed objective.

Subsequently, we calculate the average inter-class cosine similarity ( $Sim_{cls}$ ) for all base and few-shot incremental classes, providing an indication of the spread of each class on the unit sphere. A lower value suggests more compact representations. Notably, we observe values at least 10 times lower for the training with perturbations. Lastly, we assess the availability of space around unassigned pseudo-targets ( $Sim_{cls \rightarrow target}$ ) by computing the average similarity of all data features with respect to all unassigned targets. A higher average similarity corresponds to smaller margins between the features and the unassigned pseudo-targets. Table 4 illustrates that perturbations indeed increase the margin around unassigned pseudo-targets. Further discussions can be found in the supplement.

**Influence of pretraining.** To evaluate our pretraining strategy, we compare it against cross entropy (CE) and standard supervised contrastive loss (SCL) [21]. As shown in table 5, the addition of self-supervised contrastive loss (SCL+SSCL) to the pretraining session significantly enhances generalization on unseen data, showcasing improved transfer capabilities, which aligns with previous findings [5, 19]. Additionally, we present the accuracy on the validation set for the base classes  $D^0$  immediately after the pretrain phase for each strategy.

Pretrain Strategy	Phase 1:Accuracy	aHM
CE	85.70	55.20
SCL	85.18	57.38
SCL + SSCL (Ours)	<b>85.95</b>	<b>58.12</b>

Table 5. **Influence of pretraining.** aHM denotes average harmonic mean. CE is cross-entropy, SCL is supervised contrastive loss, and SSCL is self-supervised contrastive loss.

Fine-tuned params	Performance Decay $\downarrow$	Base Decay $\downarrow$
$f, g$	28.94	20.52
$g$	<b>26.99</b>	<b>17.83</b>

Table 6. **Influence of frozen parameters.** Analysing of catastrophic forgetting when 1) fine-tuning the entire model ( $f, g$ ) and 2) fine-tuning only the projection head ( $g$ ).

While all strategies exhibit close to 85% accuracy on the base validation set, our approach yields a 0.74% higher average harmonic mean compared to *SCL* and a notable 2.92% improvement over *CE*. The significance of this lies in the fact that our frozen backbone network, maintained during incremental sessions, is capable of producing strong and unique features even for unseen classes.

**Frozen parameters.** Table 6 illustrates how OrCo effectively addresses catastrophic forgetting by adopting a strategy of freezing the backbone and training only the projection head. The observed overall performance decay, along with a 2.7% greater loss of base accuracy across all sessions, demonstrates favorable outcomes for decoupling the learning process after Phase 1.

## 5. Conclusion

This paper introduced OrCo method to boost the performance of FSCIL by addressing its inherent challenges: catastrophic forgetting, overfitting, and intransigence. The OrCo framework is a novel approach that tackles these issues by leveraging features’ mutual orthogonality on the representation hypersphere and contrastive learning. By combining supervised and self-supervised contrastive learning during pretraining, the model captures diverse semantic information crucial for novel classes with limited data, implicitly addressing the intransigence challenge. Employing the proposed OrCo loss during subsequent incremental sessions ensures alignment with the generated fixed pseudo-targets, maximizing margins between classes and preserving space for incremental data. This comprehensive approach not only enhances feature space generalization but also mitigates overfitting and catastrophic forgetting, marking steps toward improving the practical value of incremental learning methods in real-world applications.



## References

- [1] Touqeer Ahmad, Akshay Raj Dhamija, Steve Cruz, Ryan Rabinowitz, Chunchun Li, Mohsen Jafarzadeh, and Terrance E. Boult. Few-shot class incremental learning leveraging self-supervised features. In *CVPR Workshops*, 2022. 2
- [2] Eden Belouadah and Adrian Popescu. Il2m: Class incremental learning with dual memory. In *ICCV*, 2019. 1, 2
- [3] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *ECCV*, 2018. 2, 4
- [4] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *ECCV*, 2018. 1, 2
- [5] Mayee Chen, Daniel Y Fu, Avanika Narayan, Michael Zhang, Zhao Song, Kayvon Fatahalian, and Christopher Ré. Perfectly balanced: Improving transfer and robustness of supervised contrastive learning. In *ICML*. PMLR, 2022. 2, 4, 8
- [6] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*. PMLR, 2020. 2
- [7] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019. 1
- [8] Ali Cheraghian, Shafin Rahman, Pengfei Fang, Soumava Kumar Roy, Lars Petersson, and Mehrtaash Harandi. Semantic-aware knowledge distillation for few-shot class-incremental learning. In *CVPR*, 2021. 1, 2, 4
- [9] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation policies from data. *arXiv preprint arXiv:1805.09501*, 2018. 14
- [10] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPR workshops*, 2020. 14
- [11] Matthias De Lange and Tinne Tuytelaars. Continual prototype evolution: Learning online from non-stationary data streams. In *ICCV*, 2021. 12, 13
- [12] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*. PMLR, 2017. 2
- [13] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013. 1
- [14] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 2010. 2
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 13
- [16] Michael Hersche, Geethan Karunaratne, Giovanni Cherubini, Luca Benini, Abu Sebastian, and Abbas Rahimi. Constrained few-shot class-incremental learning. In *CVPR*, 2022. 1, 2, 6, 14, 15, 16
- [17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 1, 2
- [18] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, 2019. 1, 2
- [19] Ashraf Islam, Chun-Fu Richard Chen, Rameswar Panda, Leonid Karlinsky, Richard Radke, and Rogerio Feris. A broad study on the transferability of visual representations with contrastive learning. In *ICCV*, 2021. 2, 4, 8
- [20] Ronald Kemker and Christopher Kanan. Farnet: Brain-inspired model for incremental learning. In *ICLR*, 2018. 1
- [21] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *NIPS*, 2020. 2, 4, 8
- [22] Do-Yeon Kim, Dong-Jun Han, Jun Seo, and Jaekyun Moon. Warping the space: Weight space rotation for class-incremental few-shot learning. In *ICLR*, 2022. 1, 2
- [23] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 2017. 1, 2
- [24] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 6
- [25] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 1955. 4
- [26] A Kukleva, H Kuehne, and B Schiele. Generalized and incremental few-shot learning by explicit learning and calibration without forgetting. in 2021 iee. In *ICCV*, 2021. 1, 2, 4, 6, 14, 15
- [27] Tianhong Li, Peng Cao, Yuan Yuan, Lijie Fan, Yuzhe Yang, Rogerio S Feris, Piotr Indyk, and Dina Katabi. Targeted supervised contrastive learning for long-tailed recognition. In *CVPR*, pages 6918–6928, 2022. 3
- [28] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE TPAMI*, 2017. 1, 2
- [29] Xialei Liu, Marc Masana, Luis Herranz, Joost Van de Weijer, Antonio M Lopez, and Andrew D Bagdanov. Rotate your networks: Better weight consolidation and less catastrophic forgetting. In *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018. 1, 2
- [30] Pratik Mazumder, Pravendra Singh, and Piyush Rai. Few-shot lifelong learning. In *AAAI*, 2021. 1, 2
- [31] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*. Elsevier, 1989. 1
- [32] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms. *arXiv preprint arXiv:1803.02999*, 2018. 2
- [33] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 2, 4

- [34] Hang Qi, Matthew Brown, and David G Lowe. Low-shot learning with imprinted weights. In *CVPR*, 2018. 2, 6, 7, 12, 14, 15, 16
- [35] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017. 1, 2
- [36] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115, 2015. 6
- [37] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *ICML*. PMLR, 2018. 1, 2
- [38] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *NIPS*, 30, 2017. 1, 2
- [39] Xiaoyu Tao, Xiaopeng Hong, Xinyuan Chang, Songlin Dong, Xing Wei, and Yihong Gong. Few-shot class-incremental learning. In *CVPR*, 2020. 1, 2, 6
- [40] Songsong Tian, Lusi Li, Weijun Li, Hang Ran, Xin Ning, and Prayag Tiwari. A survey on few-shot class-incremental learning. *Neural Networks*, 2024. 2
- [41] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *NIPS*, 29, 2016. 2
- [42] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 6
- [43] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *CVPR*, 2019. 1, 2, 4
- [44] 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, 2, 6, 12, 13, 14, 15, 16
- [45] Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017. 6
- [46] Chi Zhang, Nan Song, Guosheng Lin, Yun Zheng, Pan Pan, and Yinghui Xu. Few-shot incremental learning with continually evolved classifiers. In *CVPR*, 2021. 1, 2, 6, 14, 15, 16
- [47] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. Maintaining discrimination and fairness in class incremental learning. In *CVPR*, 2020. 1, 2
- [48] Hanbin Zhao, Yongjian Fu, Mintong Kang, Qi Tian, Fei Wu, and Xi Li. Mgsvf: Multi-grained slow vs. fast framework for few-shot class-incremental learning. *IEEE TPAMI*, 2021. 1, 2
- [49] Linglan Zhao, Jing Lu, Yunlu Xu, Zhanzhan Cheng, Dashan Guo, Yi Niu, and Xiangzhong Fang. Few-shot class-incremental learning via class-aware bilateral distillation. In *CVPR*, 2023. 2, 6, 12, 13, 14, 15, 16
- [50] Da-Wei Zhou, Fu-Yun Wang, Han-Jia Ye, Liang Ma, Shiliang Pu, and De-Chuan Zhan. Forward compatible few-shot class-incremental learning. In *CVPR*, 2022. 1, 2, 6, 13, 14, 15, 16
- [51] Da-Wei Zhou, Han-Jia Ye, Liang Ma, Di Xie, Shiliang Pu, and De-Chuan Zhan. Few-shot class-incremental learning by sampling multi-phase tasks. *IEEE TPAMI*, 2022. 6, 12, 14, 15, 16
- [52] Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu. Prototype augmentation and self-supervision for incremental learning. In *CVPR*, 2021. 1, 2

# OrCo: Towards Better Generalization via Orthogonality and Contrast for Few-Shot Class-Incremental Learning

## Supplementary Material

Within the supplement, we provide additional ablation studies in section 6, detailed breakdown tables and confusion matrices in section 7, an extended discussion on theory of orthogonality in section 8, formulation of contrastive losses in section 9 and additional implementation details in section 11.

### 6. More Ablations

**What to pull & what to perturb in OrCo loss.** Our OrCo loss comprises both pull and push components, influencing the distribution over the hypersphere. The pull effect is driven by cross-entropy loss (CE), where data features align with their assigned pseudo-targets. As illustrated in table 7, we show the advantage of aligning data features and pseudo-targets specifically from incremental sessions during the third phase. Introducing pseudo-targets assigned to the base classes to CE loss results in a performance degradation of approximately 1% in  $HM_8$ , due to an increased bias towards base classes. Next, we study the impact of perturbations, which create additional pushing forces, on different subsets of pseudo-targets. Our findings indicate that perturbing both incremental- and base-assigned pseudo-targets consistently hampers performance compared to perturbing only those assigned to incremental classes, resulting in about 9% improvement in  $HM_8$ . Higher base accuracy indicates that perturbations of both base- and incremental-assigned pseudo-targets provide more room for prevalent base classes, hindering the learning of novel classes and favouring base-class bias.

**Pseudo-targets assignment strategy.** In table 8, we highlight the crucial role of optimal initial alignment between pseudo-targets and class means. We compare a random assignment strategy to a Hungarian matching algorithm. Hungarian matching allows to find an optimal assignment based on distances between class means and pseudo-targets. We identify two optimal assignment strategies within hungarian matching 1) Reassignment and 2) Greedy Assignment. For the former, class means are reassigned to closest pseudo-targets at the beginning of each session whereas the later, carries forward the assignment from previous sessions.

We find that the random assignment strategy leads to a notable degradation in accuracy, particularly evident after the second phase for the base classes, amounting to approximately 8%. Greedy assignment performed better than reassignment. Despite reassignment being theoretically optimal, in practice we observe a performance drop likely due

CE	Perturbed Pseudo-Targets	Base Acc	Inc Acc	$HM_8$	aHM	aACC
<b>Inc</b>	<b>Inc</b>	67.60	<b>43.80</b>	<b>53.12</b>	<b>58.12</b>	67.14
	Base+Inc	78.13	30.86	44.26	50.44	69.17
Base+Inc	Inc	69.65	41.85	52.30	57.76	67.90
	Base+Inc	<b>78.90</b>	29.53	43.00	48.77	<b>69.25</b>

Table 7. **What to pull & what to perturb.** CE denotes cross-entropy that pulls data features to the pseudo-targets; Inc denotes that only assigned to incremental sessions pseudo-targets participate in the CE loss, Base+Inc denotes both base- and incremental-assigned pseudo-targets. The choice of perturbed pseudo-targets can include incremental assigned pseudo-targets with unassigned pseudo-targets (Inc), or all assigned pseudo-targets with unassigned pseudo-targets (Base+Inc). Base/Inc Acc denotes accuracy from the last 8<sup>th</sup> session. aACC denotes average accuracy over all sessions. Results on mini-ImageNet.

Assignment	Base Acc ↑	Base Decay ↓	aHM ↑	aACC ↑
Random	75.75	20.40	54.40	59.42
Reassignment	<b>83.30</b>	29.65	55.49	62.74
Greedy	<b>83.30</b>	<b>15.72</b>	<b>58.12</b>	<b>67.14</b>

Table 8. **Pseudo-targets assignment strategy.** Comparing our optimal assignment strategy against random assignment of pseudo-targets.

to noisy few-shot classes appearing geometrically closer to previously assigned pseudo-targets hence causing a shift of previously seen assigned classes and causing misalignment. This can be clearly seen in the loss of generalisation given a base decay of 29.65% vs 15.62% for best case.

Overall accuracy is substantially improved, demonstrating the critical contribution of the optimal assignment approach in addressing forgetting and achieving better alignment.

**Number of exemplars.** Due to the memory constraints inherent in FSCIL, it is common to utilize a constrained number of exemplars from the previous task. To investigate this, we conducted tests with 0, 1, and 5 exemplars, and the results are presented in table 9. We note that even with just 1 exemplar, our model achieves a performance improvement of 2.84% compared to our strong baseline, the IW method.

#	Session-wise Harmonic Mean (%) $\uparrow$								aHM	aACC
	1	2	3	4	5	6	7	8		
0	69.3	47.9	42.3	34.9	31.2	28.0	24.8	24.4	37.8	54.6
1	64.4	53.4	48.7	48.9	45.4	40.5	40.7	43.5	48.2	64.9
5	<b>68.7</b>	<b>63.9</b>	<b>60.9</b>	<b>58.0</b>	<b>55.3</b>	<b>52.4</b>	<b>52.7</b>	<b>53.1</b>	<b>58.1</b>	<b>67.1</b>

Table 9. Number of saved exemplars (#) for incremental sessions.

## 7. More Results

**Base and incremental accuracy breakdown.** We show our SOTA results with a base and incremental session accuracy breakdown for all sessions in tables 12 to 14. Note that, for Imprinted Weights (IW) [34] we use the implementation of a decoupled learning strategy from [11] and we initialise the method with our model pretrained during Phase 1. We note that LIMIT [51] has been unintentionally left out from figure 4 (main) for CUB200. Furthermore, we report the average accuracy metric to provide a comprehensive overview of our results. Note that the considerable data imbalance between base and incremental classes has an impact on this metric. The improved accuracy especially in the base classes, as illustrated in the breakdown tables, contributes to the overall enhancement of this measure.

**Confusion matrices.** In figure 7, we compare session-wise confusion matrices for a) OrCo, b) NC-FSCIL [44], and c) BiDist [49] during the final session of mini-ImageNet. Our benchmark involves assessing OrCo against its two closest competitors. OrCo plays a crucial role in finding a delicate balance between preserving knowledge of base classes and efficiently learning new ones, showcasing significantly enhanced learning capabilities in incremental classes. Notably, other methods exhibit a strong bias towards the base classes due to low transferability, while our pretraining session establishes a robust backbone. Moreover, our space reservation scheme, along with strong separation using perturbed targets and robust contrastive learning, enables us to learn a highly performant learner.

## 8. Theory of Orthogonality

This section covers the mathematical theory of orthogonality of independent vectors in high dimensional space. Let  $X_1, X_2, \dots, X_n$  be independent and identically distributed (i.i.d.) random vectors sampled from a normal distribution with mean 0 and variance 1. These vectors are in  $\mathbb{R}^n$ . The claim is that these vectors are mutually orthogonal on the unit sphere. To prove this, let’s first establish that the vectors are normalized to have a length of 1.

Given a vector  $X_i = (X_{i1}, X_{i2}, \dots, X_{in})$ , its length is given by:

$$\|X_i\| = \sqrt{X_{i1}^2 + X_{i2}^2 + \dots + X_{in}^2}$$

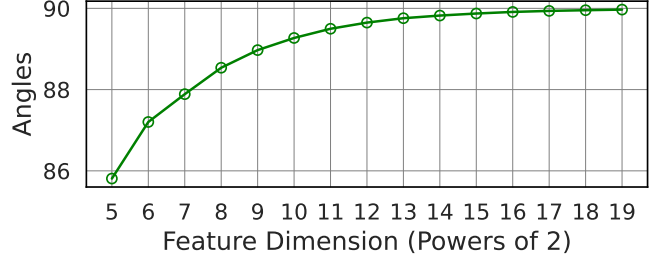


Figure 6. Practical effect of dimensions on average pair-wise angle of a 100 independent random vectors.

Since each component  $X_{ij}$  is independently sampled from a normal distribution with mean 0 and variance 1, the expected value of  $X_{ij}^2$  is 1. Therefore, the expected value of the length squared is:

$$\mathbb{E}[\|X_i\|^2] = \mathbb{E}[X_{i1}^2 + X_{i2}^2 + \dots + X_{in}^2] = n$$

This means that  $\frac{1}{\sqrt{n}}X_i$  has a length of 1 in expectation. Now, let’s consider the inner product of two different vectors  $X_i$  and  $X_j$  (where  $i \neq j$ ):

$$\mathbb{E}[X_i \cdot X_j] = \sum_{k=1}^n \sum_{k'=1}^n \mathbb{E}[X_{ik} \cdot X_{jk'}]$$

Since  $X_{ik}$  and  $X_{jk}$  are independent for  $i \neq j$ , the cross-terms in the summation will have an expected value of 0, and the only non-zero terms will be the ones where  $k = k'$ , resulting in:

$$\mathbb{E}[X_i \cdot X_j] = \sum_{k=1}^n \mathbb{E}[X_{ik} \cdot X_{jk}] = \sum_{k=1}^n \mathbb{E}[X_{ik}^2] \delta_{ij} = \delta_{ij} \cdot n$$

where  $\delta_{ij}$  is the Kronecker delta. Therefore, the expected value of the inner product is  $n$  if  $i = j$  and 0 otherwise. This means that the vectors are orthogonal in expectation. It is important to note that this property specifically holds for vectors drawn from a normal distribution with mean 0 and variance 1.

The figure 6 shows the practical effects of the above theory which yields near orthogonality but not perfect orthogonality. Only near feature dimension =  $2^{15}$  do we generate nearly orthogonal vectors. Which would lead to the projection head having  $\sim 68$  million parameter making our framework incomparable to other methods. Explicit orthogonality as we have shown previously in table 3 (main) yields better results which is why we explore this practical constraint.

## 9. Contrastive losses

**Supervised Contrastive Loss.** Given a set of sample-label pairs  $(x_i, y_i) \in Z^{SCL}$ , we define the positive set  $P_i^{SCL}$  for

$x_i$  as the collection of pairs  $(x_j, y_j)$  where  $j$  varies over all instances such that  $y_j = y_i$ . Correspondingly, the negative set  $N_i^{SCL}$  is defined as  $Z^{SCL} \setminus P_i^{SCL}$ . Then, supervised contrastive loss (SCL) is defined as:

$$\mathcal{L}_{SCL}(i; \theta) = \frac{-1}{|P_i^{SCL}|} \sum_{x_j \in P_i^{SCL}} \log \frac{\exp(x_i \cdot x_j / \tau)}{\sum_{x_k \in N_i^{SCL}} \exp(x_i \cdot x_k / \tau)}.$$

**Self-Supervised Contrastive Loss.** Given a set of samples  $x_i \in Z^{SSCL}$ , we define a positive for  $x_i$  as  $A(x_i)$  where  $A(\cdot)$  is a random transformation. Then, self-supervised contrastive loss (SSCL) loss is defined as:

$$\mathcal{L}_{SSCL}(i; \theta) = -1 \cdot \log \frac{\exp(x_i \cdot A(x_i) / \tau)}{\sum_{x_k \in Z^{SSCL}, i \neq k} \exp(x_i \cdot x_k / \tau)}.$$

## 10. Orthogonality Loss

In this section, we further expand on the orthogonality loss in our framework. We employ the orthogonality loss as an implicit geometric constraint on the set  $O$  predicated on the current batch.  $O$  contains the following: mean features for all classes within batch  $\mu_j$ , assigned targets not represented within batch and all unassigned targets  $T_u^i$ .

Formally, let us assume the session  $i$  with data  $D^i$  and classes  $C^i$ . In order to define a set  $O$  we compute some preliminaries. Firstly, for every training batch  $B$  we compute the within-batch mean for all data features. This is computed as:

$$\mu_j = \frac{1}{|C^j|} \sum_{k=0}^{|C^j|} z_k, \forall j \in C_B^i \quad (7)$$

where  $C_B^i \in C^i$  refers to all classes appearing in this particular batch. The combined set of all means can then be termed  $M_B$ . For the classes that did not appear in this batch we define as  $\neg C_B^i = C^i \setminus C_B^i$ . Subsequently we define a mapping function from seen class labels to the assigned pseudo target.

$$h : C^i \rightarrow T^i \quad (8)$$

We incorporate the remaining real data by adding the following set of assigned pseudo targets as  $\neg T_B^i = h(\neg C_B^i)$ . For completeness, we combine the above with the unassigned targets  $T_u^i$  leading to the following definition of  $O$ :

$$O = \{M_B \cup \neg T_B^i \cup T_u^i \mid B\} \quad (9)$$

Finally, the orthogonality loss takes the form:

$$\mathcal{L}_{ORTH}(O) = \frac{1}{|O|} \sum_{i=1}^{|O|} \log \sum_{j=1}^{|O|} e^{o_i \cdot o_j / \tau_o}, o_i, o_j \in O \quad (10)$$

In essence, the orthogonality loss introduces a subtle geometric constraint between real class features and the pseudo-targets. Additionally the batch-wise construction helps regularise the loss function.

Although the improvement from the Orthogonality loss are not prominent like in PSCL, it remains measurable and consistent across settings (e.g. in table 2 line 1 vs line 2, line 3 vs line 4), and thus contributes to our results. Additionally, we would like to highlight that each loss term in equation 6 incorporates orthogonality constraints either implicitly or explicitly. E.g. the orthogonality enforced on the pseudo-targets implies an implicit orthogonality among the features as they are pushed to these targets.

## 11. Further Implementation Details

**Model parameters.** Our representation learning framework is composed of:

- Encoder/Backbone Network ( $f$ ) We use the ResNet18 and 12 [15] variant for our experiments. Depending on the variant of the encoder, the representation vector has output dimensions  $D_E = 512$  or  $640$ , respectively. In table 10 we compare our choice of architecture against other FSCIL studies.
- Projection Network (projector)( $g$ ). Following the encoder, a projection MLP maps the representation vector from  $f$  to the contrastive subspace. We use a two layer projection head with a hidden dimension of size 2048. By convention output projections are normalised to the hypersphere. For simplicity, the dimension of the projected hyper sphere is initialised as  $d = 2 \lceil \log(C) \rceil$  where  $C$  is the total number of classes for our FSCIL task. Following convention, we assume normalised feature vectors.

**Training details.** For CUB200 dataset, we skip the pre-training phase given that it is common in literature to use pretrained ImageNet weights for the backbone [11, 44, 49, 50]. For any incremental sessions ( $> 0$ ) we finetune only the projection head from phase 2. For the PSCL loss we choose a perturbation magnitude  $\lambda_{pert} = 1e-2$  for our experiments and perturb only the incremental targets and unassigned targets. We train the projection head for 10 epochs in the 0-th incremental session and 100 epochs for all following sessions. Cosine scheduling is employed with warmup for a few epochs for all our phases with a maximum learning rate set to 0.4 in phase 1, 0.25 in phase 2 and 0.1 in phase 3. Given the equation 4 (main), we double our batch size by over sampling target perturbations such that the number of perturbed targets is always the same as the original training batch size. Orthogonality loss is applied batch wise. More concretely, the orthogonality loss takes as input, the within class average features inside a batch along with any other pseudo targets not in the batch. Cross entropy is applied exclusively to incremental class data as

Method	Model		
	mini-ImageNet	CIFAR100	CUB200
IW [34]	ResNet18	ResNet12	ResNet18
FACT [50]	ResNet18	ResNet20	ResNet18
CEC [46]	ResNet18	ResNet20	ResNet18
C-FSCIL [16]	ResNet12	ResNet12	-
LIMIT [51]	ResNet18	ResNet20	ResNet18
LCwoF [26]	ResNet18	-	-
BiDist [49]	ResNet18	ResNet18	ResNet18
NC-FSCIL [44]	ResNet12	ResNet12	ResNet18
OrCo	ResNet18	ResNet12	ResNet18

Table 10. ResNet architectures used in FSCIL literature

base classes are already aligned.

For pseudo target generation we employ an SGD optimiser, with 1e-2 learning rate for 2000 epochs to minimize the loss. For CUB200 the dimension of the projection head is higher which constitutes longer training cycle to fully orthogonalize the pseudo targets. CUB200 was most susceptible to forgetting for which reason we ensured that base classes were incorporated inside the CE loss component of the loss function. With the lack of a pretrain step in CUB200, we must also finetune the backbone during phase 2 to align base classes while also capturing specific representation which is important for learning effectively on a fine-grained dataset. For pretraining, we use RandAug [10] for mini-Imagenet and AutoAugment policy [9] for CIFAR100. Additionally, following the implementation of [50] we apply auto augment policy for CIFAR100 during the incremental sessions as well.

Architecture	Parameter Count (million)	aHM	Base Acc
ResNet-18	<b>12.49</b>	58.12	83.30
ResNet-12	15.06	<b>59.30</b>	<b>83.65</b>

Table 11. Comparing ResNet-12 to ResNet-18

## 12. ResNet 12 with Mini-ImageNet

In this section we measure the efficacy of our method on a different backbone. More specifically we train a ResNet-12 backbone used by [44] with our method. In table 11 we show the results. Our reported results with ResNet-18 are 58.12, while the results with ResNet-12 are 59.30 indicating an improvement with the wider ResNet-12 architecture. ResNet-18 remains as our elected architecture due to

its prominence in prior works, low parameter count, while still maintaining state-of-the art performance.

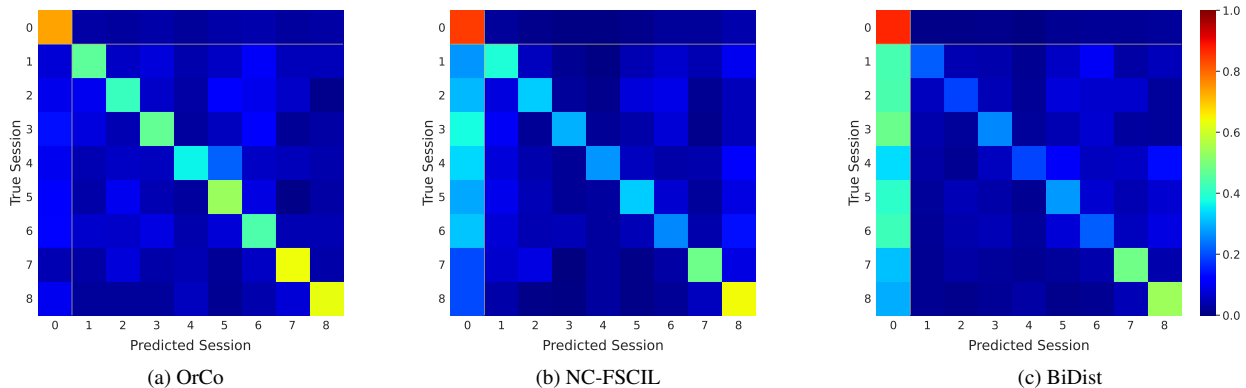


Figure 7. Visualising the session-wise confusion matrix for mini-ImageNet using a) OrCo, b) NC-FSCIL [44], and c) BiDist [49]. Each matrix demonstrates the predictive accuracy for base and incremental sessions, separated by yellow lines. High values on the diagonal (indicative of correct session predictions) are contrasted with low off-diagonal values (representing misclassifications). The first column in each matrix highlights potential prediction bias towards base classes. Our method’s performance, as illustrated, demonstrates both high diagonal accuracy and a balanced approach in reducing base class bias, as compared to the results of the competing methods.

Method	Class Group	Session-wise Accuracy (%)										Means	aACC
		0	1	2	3	4	5	6	7	8			
IW [34]	Base	<b>83.10</b>	<b>81.17</b>	<b>80.58</b>	<b>79.93</b>	<b>79.55</b>	<b>78.88</b>	<b>78.38</b>	<b>78.08</b>	<b>77.25</b>	<b>79.66</b>	68.77	
	Incremental	-	35.60	31.30	32.27	32.65	31.16	29.13	29.91	32.40	31.80		
FACT [50]	Base	75.78	75.22	74.83	74.47	74.30	74.05	73.82	73.47	73.37	74.37	60.86	
	Incremental	-	16.60	17.10	17.20	15.15	13.24	11.93	12.17	12.43	14.48		
CEC [46]	Base	72.17	70.77	70.05	69.53	69.27	68.95	68.65	68.25	67.95	69.51	59.57	
	Incremental	-	20.60	20.60	19.93	19.75	17.68	16.63	16.80	17.18	18.65		
C-FSCIL [16]	Base	76.60	76.15	74.70	73.82	73.18	71.10	70.08	68.27	67.58	72.39	61.21	
	Incremental	-	5.20	11.90	17.80	20.40	23.08	23.37	26.06	26.35	19.27		
LIMIT [51]	Base	73.27	70.43	69.47	68.68	68.18	67.73	67.30	67.07	66.68	68.76	58.04	
	Incremental	-	29.80	24.00	22.73	21.85	19.72	19.40	19.71	20.18	22.17		
LCwoF [26]	Base	64.45	57.33	53.31	52.87	51.38	48.25	47.60	47.51	47.73	52.27	46.80	
	Incremental	-	32.20	30.70	31.00	31.12	29.68	28.27	27.34	27.65	29.75		
BiDist [49]	Base	74.67	73.63	72.50	71.03	70.63	70.37	68.70	67.98	69.25	70.97	61.25	
	Incremental	-	32.60	31.40	30.33	30.30	25.48	25.23	27.09	25.10	28.44		
NC-FSCIL [44]	Base	84.37	78.25	76.00	75.73	74.80	75.42	75.52	75.13	74.77	76.67	67.82	
	Incremental	-	51.80	51.00	44.33	41.20	37.04	34.20	33.37	32.77	40.71		
OrCo	Base	83.30	76.40	74.10	72.00	71.20	70.50	69.20	68.10	67.60	72.49	67.14	
	Incremental	-	<b>62.40</b>	<b>56.10</b>	<b>52.80</b>	<b>48.90</b>	<b>45.40</b>	<b>42.20</b>	<b>42.90</b>	<b>43.80</b>	<b>49.31</b>		

Table 12. Base and Incremental accuracy shown per session for mini-ImageNet.

Method	Class Group	Session-wise Accuracy (%)										Means	aACC
		0	1	2	3	4	5	6	7	8			
IW [34]	Base	78.58	75.45	75.15	74.65	74.38	74.07	73.80	73.55	73.03	74.74	64.05	
	Incremental	-	29.40	31.90	28.53	27.60	26.36	27.20	26.91	25.88	27.97		
C-FSCIL [16]	Base	77.35	76.70	76.17	75.52	75.35	74.22	73.92	73.63	72.87	75.08	61.42	
	Incremental	-	19.80	17.40	16.20	13.65	14.92	14.53	13.66	14.00	15.52		
LIMIT [51]	Base	79.63	75.40	74.47	73.70	73.22	72.52	72.22	72.02	71.32	73.83	61.66	
	Incremental	-	27.20	24.60	21.47	21.00	20.76	21.10	20.54	20.13	22.10		
CEC [46]	Base	72.93	72.13	71.42	70.72	70.12	69.20	68.67	68.43	67.75	70.15	59.53	
	Incremental	-	29.60	27.00	22.60	21.80	22.40	22.13	21.66	21.08	23.53		
FACT [50]	Base	78.72	76.23	75.30	74.63	73.90	73.07	72.58	72.28	71.73	74.27	62.55	
	Incremental	-	29.80	25.60	21.20	20.70	20.24	22.33	21.69	21.95	22.94		
BiDist [49]	Base	69.68	68.45	67.55	66.47	65.80	64.87	64.77	64.27	64.50	66.26	56.91	
	Incremental	-	36.80	31.20	28.80	25.75	24.36	22.87	22.43	20.35	26.57		
NC-FSCIL [44]	Base	<b>82.52</b>	<b>79.55</b>	<b>78.63</b>	<b>77.98</b>	<b>77.60</b>	<b>75.98</b>	<b>74.45</b>	<b>75.18</b>	<b>73.98</b>	<b>77.32</b>	67.50	
	Incremental	-	44.00	41.60	36.47	31.95	31.32	33.97	31.31	29.30	34.99		
OrCo	Base	80.08	67.37	68.12	63.30	63.40	63.93	61.45	61.08	58.22	65.22	62.11	
	Incremental	-	<b>77.60</b>	<b>60.20</b>	<b>51.67</b>	<b>48.90</b>	<b>45.80</b>	<b>48.23</b>	<b>44.94</b>	<b>43.15</b>	<b>52.56</b>		

Table 13. Base and Incremental accuracy shown per session for CIFAR100.

Method	Class Group	Session-wise Accuracy (%)											Means	aACC
		0	1	2	3	4	5	6	7	8	9	10		
IW [34]	Base	67.53	67.07	66.83	66.55	66.48	66.31	66.13	65.99	65.85	65.85	65.75	66.39	59.72
	Incremental	-	29.03	27.74	25.00	25.95	26.61	26.51	25.45	24.34	26.08	26.93	26.36	
CEC [46]	Base	75.64	74.27	73.88	73.64	72.66	72.31	71.75	71.09	70.98	70.64	70.46	72.48	61.33
	Incremental	-	45.16	41.52	33.1	36.34	33.1	34.17	34.34	32.96	34.41	34.16	35.93	
BiDist [49]	Base	75.98	74.23	73.71	73.85	73.08	72.35	71.68	71.65	71.68	71.12	70.36	72.70	62.91
	Incremental	-	55.20	46.11	38.19	41.67	37.62	38.11	39.28	36.81	39.73	39.83	41.26	
FACT [50]	Base	77.23	75.04	74.83	74.79	74.41	74.20	73.60	73.46	73.22	72.87	72.84	74.22	64.42
	Incremental	-	53.05	47.17	38.08	40.21	38.37	39.66	40.25	38.30	40.11	39.93	41.51	
LIMIT [51]	Base	79.63	<b>79.02</b>	<b>78.81</b>	<b>78.77</b>	<b>78.39</b>	<b>78.04</b>	<b>77.72</b>	<b>77.51</b>	<b>77.24</b>	<b>76.68</b>	73.46	<b>77.75</b>	65.49
	Incremental	-	49.82	44.88	37.38	39.35	37.35	38.40	40.01	38.47	39.92	42.15	40.77	
NC-FSCIL [44]	Base	<b>80.45</b>	76.89	77.62	78.63	77.23	77.13	76.85	76.29	76.61	75.94	<b>76.19</b>	77.26	67.29
	Incremental	-	66.67	45.41	42.59	45.88	41.72	44.11	44.99	41.16	42.66	43.07	45.83	
OrCo	Base	75.59	65.54	63.79	66.76	64.91	65.05	65.50	65.01	66.24	66.27	66.62	66.48	62.36
	Incremental	-	<b>79.93</b>	<b>65.37</b>	<b>53.47</b>	<b>55.41</b>	<b>51.03</b>	<b>51.31</b>	<b>50.90</b>	<b>47.69</b>	<b>49.70</b>	<b>49.25</b>	<b>55.41</b>	

Table 14. Base and Incremental accuracy shown per session for CUB200.