# Merge and Bound: Direct Manipulations on Weights for Class Incremental Learning

**Taehoon Kim**[1]         **Donghwan Jang**[2]         **Bohyung Han**[3,4]

School of Informatics, University of Edinburgh[1]
Siebel School of Computing and Data Science, UIUC[2]
ECE[3] & IPAI[4], Seoul National University
{kthone, jh01120, bhhan}@snu.ac.kr

## Abstract

We present a novel training approach, named Merge-and-Bound (M&B) for Class Incremental Learning (CIL), which directly manipulates model weights in the parameter space for optimization. Our algorithm involves two types of weight merging: *inter-task weight merging* and *intra-task weight merging*. Inter-task weight merging unifies previous models by averaging the weights of models from all previous stages. On the other hand, intra-task weight merging facilitates the learning of current task by combining the model parameters within current stage. For reliable weight merging, we also propose a bounded update technique that aims to optimize the target model with minimal cumulative updates and preserve knowledge from previous tasks; this strategy reveals that it is possible to effectively obtain new models near old ones, reducing catastrophic forgetting. M&B is seamlessly integrated into existing CIL methods without modifying architecture components or revising learning objectives. We extensively evaluate our algorithm on standard CIL benchmarks and demonstrate superior performance compared to state-of-the-art methods.

## Introduction

Despite the remarkable achievements of recent deep neural networks (DNNs) (Radford et al. 2021; Ho, Jain, and Abbeel 2020; Brown et al. 2020), training under continually shifting data distributions encounters a significant challenge called *catastrophic forgetting*—significant performance degradation on previously learned data. Since the real-world environments, in which DNNs are deployed, dynamically change over time, addressing this issue becomes critical for enhancing the efficiency and applicability of DNNs.

Class Incremental Learning (CIL) (Douillard et al. 2020; Hou et al. 2019; Rebuffi et al. 2017; Simon, Koniusz, and Harandi 2021; Kang, Park, and Han 2022) is a kind of continual learning, which deals with a sequential influx of tasks, typically composed of disjoint class sets. Given the constraint of permitting to store only a few or no examples from previous tasks during training for new ones, catastrophic forgetting becomes a primary obstacle for CIL. Prior approaches have attempted to tackle this issue through methods including knowledge distillation (Hinton, Vinyals, and Dean 2015; Romero et al. 2015; Zagoruyko and Komodakis 2017; Kang, Park, and Han 2022), architecture expansion (Rusu et al. 2016; Yoon et al. 2018; Liu, Schiele, and

Sun 2021; Yan, Xie, and He 2021; Abati et al. 2020), or parameter regularization (Aljundi et al. 2018; Kirkpatrick et al. 2017; Zenke, Poole, and Ganguli 2017). However, these methods have inherent limitations, such as dependency on data from previous tasks for distilling knowledges of previous tasks, the need for additional network components, and poor performance, which hinder their wide-range applications.

We propose a novel training approach for CIL, referred to as Merge-and-Bound (M&B), which can be easily integrated into existing CIL methods without any modifications on architecture or loss function. Motivated by recent studies (Wortsman et al. 2022a; Rame et al. 2022; Wortsman et al. 2022b; Izmailov et al. 2018), demonstrating the benefits of weight averaging in aggregating the capabilities of multiple models, we introduce two types of weight merging techniques tailored for complex task dynamics of CIL— inter-task weight merging and intra-task weight merging— which respectively enhance the stability and plasticity of a CIL model. In order to preserve all the knowledge acquired up to the current task, *inter-task weight merging* averages the parameters of the models learned from all incremental stages in an online fashion to form a base model. The base model serves as an initialization point for a subsequent task. On the other hand, *intra-task weight merging* improves the model's adaptivity to new tasks. This technique improves the model's generalization ability for each new task by averaging multiple checkpoints along the training trajectory within the current task.

Furthermore, we incorporate *bounded model update* (Tian et al. 2023; Gouk, Hospedales, and Pontil 2020) for training models in each task, which constrains weight updates within the neighborhood of the initial model in each incremental stage. By preventing model parameters from deviating excessively from the base model, this strategy preserves knowledge gained from previous tasks. When this strategy is employed in conjunction with weight averaging, we anticipate a more reliable weight merging results by reducing the variation of merged models and leading them to stay within the same basin of the objective function with respect to the previous tasks. Throughout the integration of these techniques, we maintain a balance between stability, preserving knowledge from previous tasks, and plasticity, adapting to new tasks in CIL scenarios.

The contributions of this paper are summarized as follows:

- We propose two weight merging techniques specialized to CIL, inter- and intra-task weight merging. These techniques enhance the model's stability and plasticity by averaging model parameters across tasks (inter-task) and within each task (intra-task), respectively.

- We introduce a bounded model update strategy that constrains the total amount of model updates within each task. By enforcing the new models to remain close to the old ones, our approach alleviates the catastrophic forgetting of previously acquired knowledge and encourage reliable weight merging.

- Our algorithm, which can be conveniently integrated into existing CIL methods, substantially and consistently improves performance on multiple benchmarks with marginal extra computational complexity.

## Related Works

### Class Incremental Learning (CIL)

CIL is a challenging problem that aims to learn a model with the number of classes increasing stage-by-stage without forgetting the previously learned classes. We organize CIL methods into five groups based on their main strategies: parameter regularization, architecture expansion, bias correction, knowledge distillation, and rehearsal methods.

Parameter regularization methods (Aljundi et al. 2018; Kirkpatrick et al. 2017; Zenke, Poole, and Ganguli 2017) measure the importance of network parameters and adjust their flexibility to mitigate catastrophic forgetting. However, these methods suffer from unsatisfactory generalization performance in CIL scenarios (van de Ven and Tolias 2019; Hsu et al. 2018). Architecture expansion methods (Rusu et al. 2016; Yoon et al. 2018; Liu, Schiele, and Sun 2021; Yan, Xie, and He 2021; Abati et al. 2020) dynamically expand the network capacity to handle incoming tasks by adding new neurons or layers. However, they introduce computational burdens due to additional network components. Bias correction methods (Hou et al. 2019; Wu et al. 2019) address the bias towards new classes caused by the class imbalances in CIL by introducing scale and shift parameters or matching the scale of weight vectors. Knowledge distillation methods (Hinton, Vinyals, and Dean 2015; Romero et al. 2015; Zagoruyko and Komodakis 2017; Kang, Park, and Han 2022) encourage models to preserve previous task knowledge by mimicking the representations of old models. Several approaches match output distributions or attention maps to preserve important information. Rehearsal-based methods (Rebuffi et al. 2017; Ostapenko et al. 2019; Shin et al. 2017) store representative examples or employ generative models to mitigate forgetting. Examples include maintaining class centroids (Rebuffi et al. 2017) or using generative adversarial networks (Goodfellow et al. 2014; Liu et al. 2020a; Odena, Olah, and Shlens 2017) to generate synthetic examples.

In contrast, we propose a novel weight manipulation approach which is compatible and easy to integrate with the existing CIL methods, without requiring any changes to the network architectures or loss functions.

### Robust Transfer Learning

Robust transfer learning aims to adapt a pre-trained model to a new domain or task without losing its generalization ability. This is similar to continual learning, which seeks to prevent catastrophic forgetting of previous knowledge while learning new tasks sequentially.

A common approach to robust transfer learning is to freeze or constrain the weights of the pre-trained model (Kumar et al. 2022; Lee et al. 2023) while fine-tuning. This prevents feature distortion from adaptation and preserves the original knowledge. Another line of work explores the idea of weight averaging or interpolation to enhance the performance of fine-tuned models (Neyshabur, Sedghi, and Zhang 2020; Wortsman et al. 2022a; Ramé et al. 2022). However, these methods primarily concentrate on a single downstream task, neglecting the sequential domain shifts that are inherent in continual learning.

## Proposed Approaches

### Problem Formulation

CIL is a learning framework to handle a sequence of tasks, $\mathcal{T}_{1:K} = \{T_1, \cdots, T_K\}$. Each task $T_k$ consists of a labeled dataset $\mathcal{D}_k$ whose label set, $\mathcal{C}_k$, is disjoint to the ones defined in the past, i.e., $(\mathcal{C}_1 \cup \cdots \cup \mathcal{C}_{k-1}) \cap \mathcal{C}_k = \emptyset$. At the $k^{\text{th}}$ incremental stage, the current model $M_k(\cdot)$ is trained on integrated dataset $\mathcal{D}'_k = \mathcal{D}_k \cup \mathcal{B}_{k-1}$, where $\mathcal{B}_{k-1}$ is a set of representative exemplars that belong to all the previously learned classes. The performance of a CIL algorithm is evaluated using a test set comprising test data from all the stages.

### Motivation

Inspired by recent studies in transfer learning (Wortsman et al. 2022a; Rame et al. 2022; Wortsman et al. 2022b; Izmailov et al. 2018) that combine the weights of multiple models, we propose two unique model averaging strategies for CIL in the presence of complex dynamics in data distribution: inter-task weight merging and intra-task weight merging. Although both techniques concentrate on integrating the competency of multiple models, they have distinct objectives. Inter-task weight merging aims to consolidate knowledge from all the previously learned tasks and construct a comprehensive model by computing the moving average of the previous model parameters. Meanwhile, intra-task weight merging boosts the model's generalization capability on new task by merging the multiple models along the training trajectory of a current task.

To boost the effectiveness and reliability of both weight merging techniques, we introduce bounded model update scheme. Recent studies on model averaging (Neyshabur, Sedghi, and Zhang 2020; Wortsman et al. 2022a; Ramé et al. 2022) highlight that weight averaging techniques are effective for the models associated with the same loss basin, leading to identify a flat minimum within the basin. However, in the scenarios that data from previous tasks are scarce during incremental learning stages, training on the current task
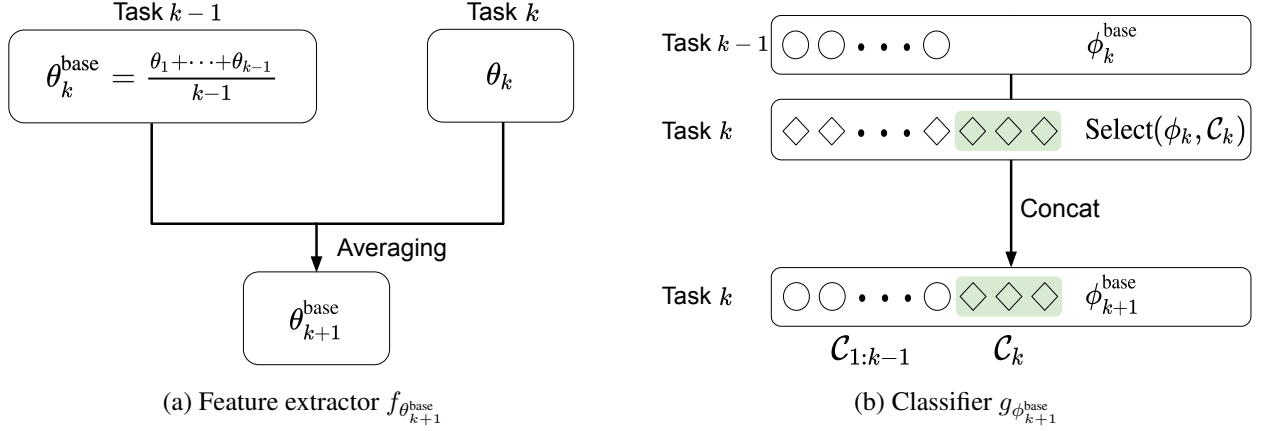
(a) Feature extractor $f_{\theta_{k+1}^{\text{base}}}$

(b) Classifier $g_{\phi_{k+1}^{\text{base}}}$

Figure 1: Description of inter-task weight merging: Upon the completion of the $k^{\text{th}}$ incremental stage, we establish the base model $M_{k+1}^{\text{base}}(\cdot)$, which will serve as the *initialization point* for the $(k+1)^{\text{st}}$ stage. The model comprises a feature extractor $f_{\theta_{k+1}^{\text{base}}}(\cdot)$ and a classifier $g_{\phi_{k+1}^{\text{base}}}(\cdot)$, and they are constructed by the following procedures. (a) To construct the base feature extractor $f_{\theta_{k+1}^{\text{base}}}(\cdot)$, we set $\theta_{k+1}^{\text{base}}$ to the moving average of all the previous feature extractor weights, $\theta_1, \theta_2, \cdots, \theta_k$, which is easily computed with $\theta_k^{\text{base}}$ and $\theta_k$ in a recursive manner following Equation (1). (b) For learning the classifier $g_{\phi_{k+1}^{\text{base}}}(\cdot)$, we concatenate the weights of the current base classifier $\phi_k^{\text{base}}$ with the weights of the current classifier $\phi_k$ associated with the class set in the current task $\mathcal{C}_k$.

with a large number of training examples results in the convergence to different local minima from the one associated with earlier tasks. To tackle the issue, we impose a constraint on the magnitude of model updates and maintain the robustness of trained models to diverse tasks over many incremental stages. Through the bounded model update strategy for training on new tasks, trained models are less prone to stray from the loss basins of previous tasks. This approach stabilizes our weight merging techniques by reducing the variance of trained models across incremental stages and within each of an incremental task, facilitating the achievement of consistent performance.

In the rest of this section, we describe the details of each technical components, inter-task weight averaging, intra-task weight averaging, and bounded model update.

**Inter-task weight merging**

We merge the models from all the preceding tasks to the base model $M_k^{\text{base}}(\cdot)$, for consolidating the previously acquired knowledge. The base model is composed of a feature extractor $f_{\theta_k^{\text{base}}}(\cdot)$ and a classifier $g_{\phi_k^{\text{base}}}(\cdot)$, collectively parameterized by $\Theta_k^{\text{base}} = \{\theta_k^{\text{base}}, \phi_k^{\text{base}}\}$.

Upon the completion of the $k^{\text{th}}$ incremental stage, we combine the current feature extractor $f_{\theta_k}(\cdot)$ and classifier $g_{\phi_k}(\cdot)$ into the base model, $M_k^{\text{base}}(\cdot)$ at the corresponding stage, as illustrated in Figure 1. This process creates the new base model for the next stage, denoted by $M_{k+1}^{\text{base}}(\cdot)$. The new feature extractor is given by the moving average of all models as follows:

$$\theta_{k+1}^{\text{base}} = \frac{k-1}{k} \cdot \theta_k^{\text{base}} + \frac{1}{k} \cdot \theta_k. \qquad (1)$$

To obtain the classifier of the new base model at the $(k+1)^{\text{st}}$ stage, we concatenate the current base model classifier $g_{\phi_k^{\text{base}}}(\cdot)$ defined for the classes in $\mathcal{C}_{k-1}$ with the weights corresponding to the classes in $\mathcal{C}_k$ of the current classifier, $g_{\phi_k}(\cdot)$, which is expressed as

$$\phi_{k+1}^{\text{base}} = \text{Concat}(\phi_k^{\text{base}}, \text{Select}(\phi_k, \mathcal{C}_k)), \qquad (2)$$

where $\text{Select}(\phi_k, \mathcal{C}_k)$ extracts the weights corresponding to the classes in $\mathcal{C}_k$ of the current classifier $g_{\phi_k}(\cdot)$. At the $(k+1)^{\text{st}}$ incremental stage, we set the initial model as the base model, $M_{k+1}^{\text{base}}(\cdot)$, which encapsulates all the knowledge learned up to the current stage.

**Intra-task weight merging**

The primary goal of intra-task weight merging is to enhance generalization ability on the current task by averaging the multiple checkpoints along training trajectories. In the $k^{\text{th}}$ incremental stage, an intra-task merged model, $M_k^{\text{avg}}(\cdot)$, parameterized by $\Theta_k^{\text{avg}}$ is updated at every $e_a$ epochs as

$$\Theta_k^{\text{avg}} \leftarrow \frac{n \cdot \Theta_k^{\text{avg}} + \Theta_k}{n+1}, \qquad (3)$$

where $n$ denotes the number of models involved in the merging. Once the stage is completed, the final model, $M_k(\cdot)$, is replaced by the intra-task merged model, $M_k^{\text{avg}}(\cdot)$, *i.e.*, $\Theta_k \leftarrow \Theta_k^{\text{avg}}$, for inference and the computation of the base model, $M_{k+1}^{\text{base}}(\cdot)$.

For the models equipped with Batch Normalization (BN) (Ioffe and Szegedy 2015), an extra post-training data forwarding is required to estimate new mean and variance of the activations after model averaging (Garipov et al. 2018; Izmailov et al. 2018). To this end, instead of executing an additional forward pass after resetting the running statistics as in SWA (Izmailov et al. 2018), we conduct the forward pass on top of the running statistics of the current model

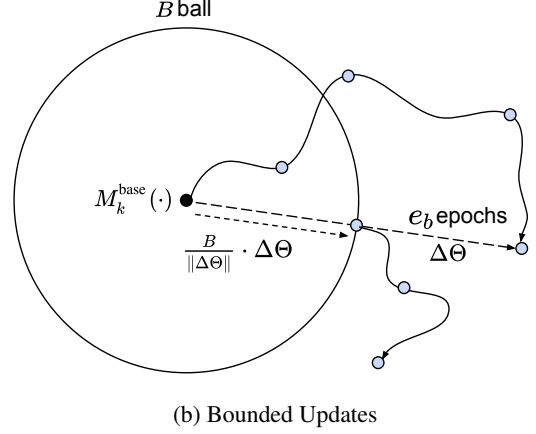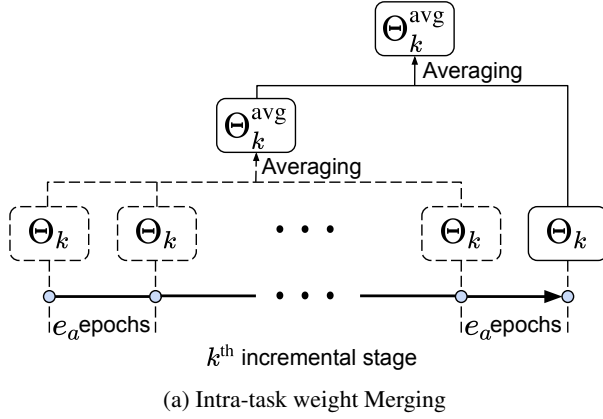| (a) Intra-task weight Merging | (b) Bounded Updates |

Figure 2: (a) Illustration of the intra-task weight merging: We introduce an intra-task weight merging by the moving average of weights in multiple models along the training trajectories, as described in Equation (3). Intra-task weight merged model is utilized for inference and for computing the next stage base model $M_k^{\text{base}}(\cdot)$. (b) Illustration of the bounded update technique: We constrain the weight updates around the base model denoted by $M_k^{\text{base}}(\cdot)$. This strategy is designed to preserve the knowledge in the base model but search for the unexplored space during optimization.

$M_k(\cdot)$ estimated before model merging, which incurs slight updates of the BN statistics. This strategy alleviates the bias towards the current task caused by sample deficiency of the classes introduced in the previous tasks.

## Bounded Model Update

We enforces a constraint for model update in CIL, which bounds the magnitude of weight updates from the base model at every $e_b$ epochs as

$$\Delta\Theta \leftarrow \begin{cases} B \cdot \frac{\Delta\Theta}{\|\Delta\Theta\|}, & \text{if } \|\Delta\Theta\| > B \\ \Delta\Theta, & \text{otherwise} \end{cases}, \qquad (4)$$

where $\Delta\Theta$ denotes the displacement of the current model from the base model, $M_k^{\text{base}}(\cdot)$, and $B$ indicates the limit of the total model update magnitude. Note that the proposed bounded model update is performed within an incremental stage jointly with intra-task weight merging. Given that the base model is presumed to hold the knowledge of previous tasks, this strategy aims to directly prevent model updates in the current task from diverging from the base model in the weight space, thereby preventing the loss of previously learned information. By combining bounded updates with inter-task and intra-task weight merging, this approach is expected to foster an reliable ensemble effect.

## Experiments

### Datasets and Evaluation Protocol

We conduct experiments on CIFAR-100 (Krizhevsky, Nair, and Hinton 2009) and ImageNet-100/1000 (Russakovsky et al. 2015). CIFAR-100 contains 50,000 and 10,000 training and validation images, respectively, in 100 classes. Following the standard protocols, we evaluate all the compared algorithms using three different class orders on CIFAR-100 and the class order specified in (Douillard et al. 2020) on ImageNet-100 and ImageNet-1000.

We incorporate Merge-and-Bound (M&B) into various state-of-the-art CIL algorithms based on knowledge distillation (PODNet (Douillard et al. 2020), AFC (Kang, Park, and Han 2022)), architecture expansion (FOSTER (Wang et al. 2022)), and virtual class augmentation (IL2A (Zhu et al. 2021)). Note that IL2A is employed to compare with non-exemplar-based methods.

As in the previous works (Douillard et al. 2020; Hou et al. 2019; Liu, Schiele, and Sun 2021), we first train the model using a half of the entire classes at the initial stage and split the remaining classes into 5/10/25/50 stages on CIFAR-100, 5/10/25 stages on ImageNet-100, and 5/10 stages on ImageNet-1000 to simulate CIL scenarios. We test models on all the seen classes at each incremental stage, and report the *average incremental accuracy* (Rebuffi et al. 2017; Douillard et al. 2020; Hou et al. 2019)—average accuracy over all incremental stages.

### Implementation Details

As our approach is a plug-in method, we follow the implementation settings of the existing methods (Douillard et al. 2020; Kang, Park, and Han 2022; Wang et al. 2022) in principle. The memory budget size is set to 20 per class unless specified otherwise. Detailed description about implementation details and hyperparameters are discussed in supplementary document.

### Results on CIFAR-100 and ImageNet-100/1000

**CIFAR-100**  Table 1 illustrates that the proposed algorithm, denoted by M&B, enhances the performance of the baseline models in all CIL scenarios, with notable margins. The performance gains of M&B appear modest in FOSTER (Wang et al. 2022). This is because of the design choice of FOSTER (Wang et al. 2022), which exploits enlarged model capacity and ensemble effects from multiple models. Such strategies tend to reduce the unique benefits of the proposed approach, especially when the number of incremental

Table 1: CIL performance (%) on CIFAR-100. The proposed training technique (M&B) consistently improves performance when plugged into the existing methods. Note that we run 3 experiments with 3 different orders and report the average result. Models with ∗ denotes the report of reproduced results. The bold-faced numbers indicate the best performance.

| | CIFAR-100 | | | |
| Number of tasks | 5 | 10 | 25 | 50 |
|---|---|---|---|---|
| iCaRL (Rebuffi et al. 2017) | 58.08 | 53.78 | 50.60 | 44.20 |
| BiC (Wu et al. 2019) | 56.86 | 53.21 | 48.96 | 47.09 |
| Mnemonics (Liu et al. 2020b) | 63.34 | 62.28 | 60.96 | - |
| GeoDL* (Simon, Koniusz, and Harandi 2021) | 65.34 | 63.61 | 60.21 | 52.28 |
| UCIR (Hou et al. 2019) | 64.01 | 61.22 | 57.57 | 49.30 |
| PODNet* (Douillard et al. 2020) | 64.83±0.62 | 62.75±0.74 | 60.73±0.62 | 58.37±0.83 |
| PODNet (Douillard et al. 2020) + M&B | **67.69**±0.50 | **66.49**±0.50 | **64.93**±0.42 | **63.29**±0.37 |
| AFC* (Kang, Park, and Han 2022) | 66.11±0.60 | 64.77±0.74 | 63.68±0.74 | 61.94±0.60 |
| AFC (Kang, Park, and Han 2022) + M&B | **67.96**±0.70 | **67.10**±0.68 | **66.12**±0.67 | **65.38**±0.41 |
| FOSTER* (Wang et al. 2022) | 72.23±0.51 | 69.12±0.67 | 65.45±0.90 | 59.60±0.85 |
| FOSTER (Wang et al. 2022) + M&B | **73.03**±0.61 | **70.58**±0.64 | **66.79**±0.84 | **62.47**±0.97 |

Table 2: CIL performance (%) on ImageNet-100/1000. M&B demonstrates significant performance gains when integrated into existing methods, even in the large-scale benchmarks for CIL.

| | ImageNet-100 | | | ImageNet-1000 | |
| Number of tasks | 5 | 10 | 25 | 5 | 10 |
|---|---|---|---|---|---|
| iCaRL (Rebuffi et al. 2017) | 65.56 | 60.90 | 54.56 | 51.36 | 46.72 |
| BiC (Wu et al. 2019) | 68.97 | 65.14 | 59.65 | 45.72 | 44.31 |
| Mnemonics (Liu et al. 2020b) | 72.58 | 71.37 | 69.74 | 64.54 | 63.01 |
| GeoDL* (Simon, Koniusz, and Harandi 2021) | 73.87 | 73.55 | 71.72 | 65.23 | 64.46 |
| UCIR (Hou et al. 2019) | 71.04 | 70.71 | 62.94 | 64.34 | 61.18 |
| PODNet* (Douillard et al. 2020) | 74.06 | 71.51 | 67.31 | 68.18 | 65.58 |
| PODNet (Douillard et al. 2020) + M&B | **75.98** | **74.08** | **70.70** | **69.53** | **67.76** |
| AFC* (Kang, Park, and Han 2022) | 76.91 | 75.26 | 73.65 | 68.06 | 66.39 |
| AFC (Kang, Park, and Han 2022) + M&B | **77.05** | **76.35** | **74.35** | **70.28** | **69.51** |
| FOSTER* (Wang et al. 2022) | 80.22 | 78.15 | 71.74 | – | – |
| FOSTER (Wang et al. 2022) + M&B | **80.59** | **79.20** | **73.24** | – | – |

stages is small. However, the performance of the proposed method stands out again when the number of incremental stages becomes 50, where the benefits of FOSTER (Wang et al. 2022) diminish abruptly due to lack of model capacity.

**ImageNet-100/1000** Table 2 shows the results on large-scale benchmarks, ImageNet-100/1000. While the proposed method consistently boosts the baseline algorithms, Table 2 clearly illustrates that the performance gains are particularly remarkable for ImageNet-1000, which is the most challenging benchmark for CIL. Additionally, across all datasets, the proposed algorithm demonstrates enhanced performance gains with an increasing number of tasks. This characteristic is highly desirable for CIL in practical scenarios, which have no restrictions or information on the number of incoming tasks. Note that we were unable to reproduce the ImageNet-1000 experiments for FOSTER (Wang et al. 2022) since the training configuration was not released publicly.

## Ablation Studies

We perform various ablation studies to validate the effectiveness of the proposed training technique. All the experiments are performed on CIFAR-100 with 50 incremental stages unless specified otherwise. To measure the stability and plasticity, we compute the average new accuracy, which is the av-

Table 3: Component analysis of our algorithm for inter-task weight merging, intra-task weight merging, and bounded model update.

| | Forgetting ↓ | Avg. new acc. ↑ | Overall acc.↑ |
|---|---|---|---|
| (a) M&B | 15.38 | 59.35 | **65.38** |
| (b) w/o inter-task | 21.77 | 62.74 | 61.81 |
| (c) w/o intra-task | **13.10** | 51.60 | 65.08 |
| (d) w/o bounded | 18.72 | **64.14** | 64.21 |

erage accuracy of new classes over incremental stages, and the forgetting metric (Lee et al. 2019), which is the average of the performance degradation for each class.

**Variations of the proposed method** Table 3 presents the results obtained from different combinations of components in M&B. Each component contributes to the overall improvements, albeit in different ways. Comparing M&B to w/o inter-task and w/o bounded, we observe a decrease in forgetting, indicating that inter-task weight merging and bounded updates effectively help retain previous knowledge. However, the low average new accuracy in w/o intra-task suggests that integrating bounded updates solely with inter-task weight merging leads to significant degradation in adaptation. On the other hand, the increase in average new ac-

Table 4: CIL performance (%) with a limited memory budget on CIFAR-100: 1 exemplar memory per class for POD-Net (Douillard et al. 2020) and AFC (Kang, Park, and Han 2022), and no memory for IL2A (Zhu et al. 2021).

| | CIFAR-100 | | | |
|---|---|---|---|---|
| Number of tasks | 5 | 10 | 25 | 50 |
| PODNet* (Douillard et al. 2020) | 43.70 | 34.19 | 26.58 | 14.78 |
| PODNet (Douillard et al. 2020) + M&B | **52.27** | **45.63** | **36.81** | **20.84** |
| AFC* (Kang, Park, and Han 2022) | 49.82 | 42.78 | 35.51 | 23.59 |
| AFC (Kang, Park, and Han 2022) + M&B | **53.41** | **48.92** | **46.08** | **35.25** |
| IL2A* (Zhu et al. 2021) | 65.70 | 58.14 | 54.40 | 20.42 |
| IL2A (Zhu et al. 2021) + M&B | **67.46** | **61.80** | **58.25** | **43.54** |

Table 5: Analysis of inter-task weight merging. EMA denotes exponential moving average, where the numbers in parentheses indicate the smoothing factor; a higher smoothing factor assigns more weights to the most recent tasks.

| | Averaging factor | | | |
|---|---|---|---|---|
| | Avg (Ours) | EMA (0.9) | EMA (0.5) | EMA (0.1) |
| Forgetting ↓ | **15.38** | 19.12 | 18.28 | 17.41 |
| Avg. new acc. ↑ | 59.35 | 61.34 | **61.74** | 59.10 |
| Overall acc. ↑ | **65.38** | 64.09 | 64.10 | 64.01 |

Table 6: Results by varying the weight averaging periods and the BN statistics computation methods in intra-class weight merging. R means resetting the running statistics and NC indicates no change of the BN statistics after the weight merging. † denotes our choice for reporting the main results.

| | Weight averaging period | | | | BatchNorm | | |
|---|---|---|---|---|---|---|---|
| | 1† | 5 | 10 | 15 | Ours | R | NC |
| Forgetting ↓ | **15.38** | 16.12 | 15.59 | 15.58 | **15.38** | 36.67 | 25.60 |
| Avg. new acc. ↑ | 59.35 | **59.41** | 58.50 | 58.95 | 59.35 | 39.75 | **60.47** |
| Overall acc. ↑ | 65.38 | 64.85 | 65.19 | **65.51** | **65.38** | 25.60 | 64.57 |

Table 7: Results by varying the bounding period and the size of the bound.

| | Bounding period | | | Bounding threshold | | |
|---|---|---|---|---|---|---|
| | 5 | 10 | 15† | 5 | 10† | 15 |
| Forgetting ↓ | 15.89 | 16.45 | **15.38** | 16.53 | **15.38** | 16.33 |
| Avg. new acc. ↑ | 57.92 | 58.12 | **59.35** | 58.66 | 59.35 | **59.66** |
| Overall acc. ↑ | 64.81 | 64.73 | **65.38** | 64.50 | **65.38** | 64.99 |

curacy in ours compared to w/o intra-task indicates that the inclusion of intra-task weight merging, in combination with bounded updates, helps mitigate this problem. Thus, the combination of these components enhances the performance of our method in retaining previous knowledge while facilitating adaptation to new tasks.

**Results on limited memory budgets** We evaluate the performance of M&B with exemplar-based methods such as PODNet (Douillard et al. 2020) and AFC (Kang, Park, and Han 2022) when only one exemplar is available. We also verify the effectiveness of M&B when it is combined with non-exemplar-based method, IL2A (Zhu et al. 2021). As shown in Table 4, our algorithm significantly outperforms existing methods when operating under limited memory budgets by exploiting model merging techniques and constraining the amount of model updates. This property is desirable conceptually because CIL may have a large number of stages and even a small number of exemplars may be difficult to hold in practice.

**Variations in inter-task weight merging** Table 5 shows the results from different strategies of inter-task weight merging. Our moving average technique provides the same solution with the offline averaging, which outperforms Exponential Moving Averaging (EMA) with various smoothing factors. The EMA methods favor recent models and lead to forgetting the previous knowledge.

**Variations in intra-task weight merging** We explore the characteristics of intra-task weight merging by varying the weight averaging periods and the BN statistics computation strategies. According to Table 6, the proposed intra-task weight merging technique is robust to the changes of averaging period but the results are affected by the meth-

ods to compute the BN statistics. Due to the distinct characteristics of CIL that it should perform well on the previous classes, resetting the running statistics (R) gives severe performance degradation since the computed running statistics will be highly biased towards current tasks in computation procedure after the reset. Also, not forwarding the additional data path (NC) shows the degraded performance since there is discrepancy between the running statistics and intra-task merged model since the statistics for the merged model are not computed.
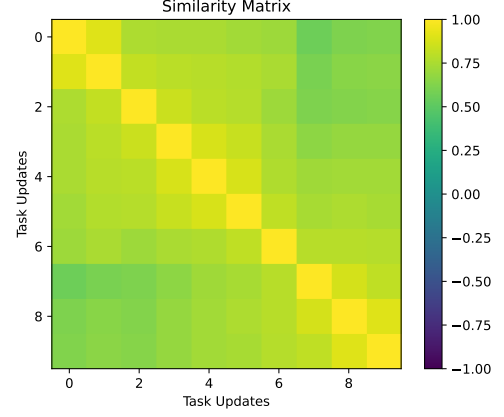
**Variations in bounded model updates** We conducted experiments by varying the frequency and the allowed size of the bounded model updates. As illustrated in Table 7, the overall accuracy is robust to the changes in the bounding frequency while applying the bounded model update results in clear advantage. For the bounding threshold, we observe that larger threshold results in larger adaptivity since it allows model to change in large magnitude for the new task.

**Effects of M&B on weights updates** To better comprehend the impact of our method, M&B, within the weight space, we conducted an analysis of the cosine similarity among task updates. Task updates are defined by the weight changes from the beginning to the end of each task. As depicted in Figure 3, the absence of M&B results in task updates that are largely independent, and in some cases, exhibit a negative correlation. Conversely, the integration of M&B leads to a significant positive correlation among task updates. These aligned task updates enable stable integration of different tasks and reduce the catastrophic forgetting by preventing the learning of new tasks from disrupting previous task updates.

**Similarity between representations** To assess the impact of the proposed method on feature representations, we measure the similarity of representations. To this end, we extract the feature representations of the test examples from the final layer of the models trained on different tasks, *e.g.*, baseline
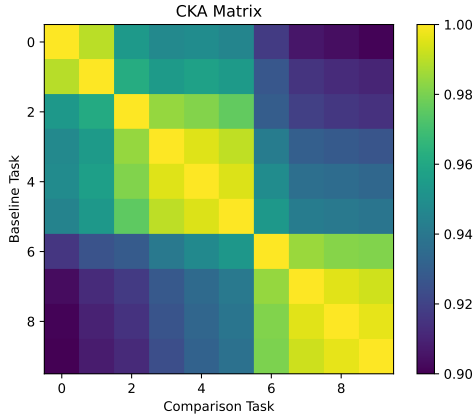
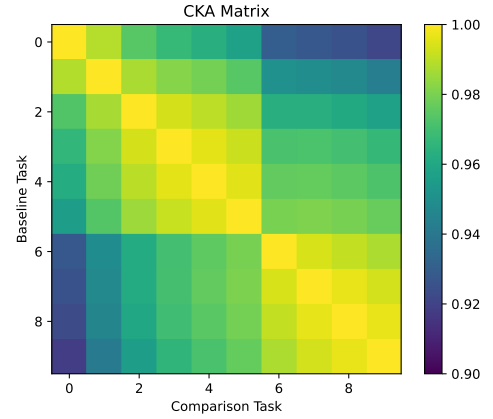(a) Task similarity of AFC (Kang, Park, and Han 2022)



(b) Task similarity of AFC (Kang, Park, and Han 2022) + M&B

Figure 3: We measure the cosine similarities between all pairs of the model update vectors occurred in each stage. The model update vectors become positively correlated when M&B is incorporated.



(a) CKA of AFC (Kang, Park, and Han 2022)



(b) CKA of AFC (Kang, Park, and Han 2022) + M&B

Figure 4: CKA between models after training individual incremental stages. We visualize the similarity between pairs of models obtained from two different tasks—baseline task and comparison task—by measuring CKA of the representations of test examples of all classes learned up to baseline task extracted from the two models.

task and comparison task using Centered Kernel Alignment (CKA) (Cortes, Mohri, and Rostamizadeh 2012; Kornblith et al. 2019) using test data of all classes learned up to the baseline task. Figure 4 clearly illustrates that M&B enhances feature similarities across the models trained in different incremental stages, which implies that M&B alleviates catastrophic forgetting at representation level.

**Computational cost** We provide precise training time overhead induced by our method. Based on a single NVIDIA RTX-8000 GPU with a ResNet-32 backbone, we observe that the inter- and intra-task weight merging only take 0.003 seconds each, whereas the bounded model update operation requires 0.011 seconds. Given that these operations only occur intermittently, *e.g.*, per task or every several epochs, the additional training time is marginal. Note that one epoch of additional forwarding is required for computing BN statistics per task. Because M&B is a training technique, it incurs no extra overhead for inference. These results demonstrates

that the proposed method only requires negligible computational cost and shows the great performance gains.

## Conclusion

We introduce an unique Class Incremental Learning (CIL) approach that incorporates weight ensemble techniques to counteract catastrophic forgetting. By redefining CIL as a sequential transfer learning process, innovative inter-task and intra-task weight averaging, as well as bounded update techniques were introduced. The proposed strategy is seamlessly integrated into existing CIL strategies without necessitating changes to the network architecture or loss function. Additionally, it exceeded the performance of contemporary methods on numerous benchmarks, specifically in low memory budget scenarios. This paper contributes significantly to CIL research by elucidating how model weights can be effectively employed to maintain and transfer knowledge across tasks.

# References

Abati, D.; Tomczak, J.; Blankevoort, T.; Calderara, S.; Cucchiara, R.; and Bejnordi, B. E. 2020. Conditional Channel Gated Networks for Task-Aware Continual Learning. In *CVPR*.

Aljundi, R.; Babiloni, F.; Elhoseiny, M.; Rohrbach, M.; and Tuytelaars, T. 2018. Memory Aware Synapses: Learning what (not) to forget. In *ECCV*.

Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. In *NeurIPS*.

Cortes, C.; Mohri, M.; and Rostamizadeh, A. 2012. Algorithms for learning kernels based on centered alignment. In *JMLR*.

Douillard, A.; Cord, M.; Ollion, C.; and Robert, T. 2020. PODNet: Pooled Outputs Distillation for Small-Tasks Incremental Learning. In *ECCV*.

Garipov, T.; Izmailov, P.; Podoprikhin, D.; Vetrov, D. P.; and Wilson, A. G. 2018. Loss surfaces, mode connectivity, and fast ensembling of dnns. In *NeurIPS*.

Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative Adversarial Nets. In *NIPS*.

Gouk, H.; Hospedales, T. M.; and Pontil, M. 2020. Distance-based regularisation of deep networks for fine-tuning. In *arXiv preprint arXiv:2002.08253*.

Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the Knowledge in a Neural Network. *arXiv preprint arXiv:1503.02531*.

Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. In *NeurIPS*.

Hou, S.; Pan, X.; Loy, C. C.; Wang, Z.; and Lin, D. 2019. Learning a Unified Classifier Incrementally via Rebalancing. In *CVPR*.

Hsu, Y.-C.; Liu, Y.-C.; Ramasamy, A.; and Kira, Z. 2018. Re-evaluating Continual Learning Scenarios: A Categorization and Case for Strong Baselines. *arXiv preprint arXiv:1810.12488*.

Ioffe, S.; and Szegedy, C. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*.

Izmailov, P.; Podoprikhin, D.; Garipov, T.; Vetrov, D.; and Wilson, A. G. 2018. Averaging weights leads to wider optima and better generalization. In *UAI*.

Kang, M.; Park, J.; and Han, B. 2022. Class-incremental learning by knowledge distillation with adaptive feature consolidation. In *CVPR*.

Kirkpatrick, J.; Pascanu, R.; Rabinowitz, N.; Veness, J.; Desjardins, G.; Rusu, A. A.; Milan, K.; Quan, J.; Ramalho, T.; Grabska-Barwinska, A.; et al. 2017. Overcoming Catastrophic Forgetting in Neural Networks. *Proceedings of the national academy of sciences*.

Kornblith, S.; Norouzi, M.; Lee, H.; and Hinton, G. 2019. Similarity of neural network representations revisited. In *ICML*.

Krizhevsky, A.; Nair, V.; and Hinton, G. 2009. Learning Multiple Layers of Features from Tiny Images. Technical report.

Kumar, A.; Raghunathan, A.; Jones, R.; Ma, T.; and Liang, P. 2022. Fine-Tuning Can Distort Pretrained Features and Underperform Out-of-Distribution. In *ICLR*.

Lee, K.; Lee, K.; Shin, J.; and Lee, H. 2019. Overcoming catastrophic forgetting with unlabeled data in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 312–321.

Lee, Y.; Chen, A. S.; Tajwar, F.; Kumar, A.; Yao, H.; Liang, P.; and Finn, C. 2023. Surgical Fine-Tuning Improves Adaptation to Distribution Shifts. In *ICLR*.

Liu, X.; Wu, C.; Menta, M.; Herranz, L.; Raducanu, B.; Bagdanov, A. D.; Jui, S.; and de Weijer, J. v. 2020a. Generative Feature Replay for Class-Incremental Learning. In *CVPR Workshops*.

Liu, Y.; Schiele, B.; and Sun, Q. 2021. Adaptive Aggregation Networks for Class-Incremental Learning. In *CVPR*.

Liu, Y.; Su, Y.; Liu, A.-A.; Schiele, B.; and Sun, Q. 2020b. Mnemonics Training: Multi-Class Incremental Learning without Forgetting. In *CVPR*.

Neyshabur, B.; Sedghi, H.; and Zhang, C. 2020. What Is Being Transferred in Transfer Learning? In *NeurIPS*.

Odena, A.; Olah, C.; and Shlens, J. 2017. Conditional Image Synthesis with Auxiliary Classifier Gans. In *ICML*.

Ostapenko, O.; Puscas, M.; Klein, T.; Jahnichen, P.; and Nabi, M. 2019. Learning to Remember: A Synaptic Plasticity Driven Framework for Continual Learning. In *CVPR*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *ICML*.

Ramé, A.; Ahuja, K.; Zhang, J.; Cord, M.; Bottou, L.; and Lopez-Paz, D. 2022. Model Ratatouille: Recycling Diverse Models for Out-of-Distribution Generalization. *arXiv preprint arXiv:2212.10445*.

Rame, A.; Kirchmeyer, M.; Rahier, T.; Rakotomamonjy, A.; Cord, M.; et al. 2022. Diverse Weight Averaging for Out-of-Distribution Generalization. In *NeurIPS*.

Rebuffi, S.-A.; Kolesnikov, A.; Sperl, G.; and Lampert, C. H. 2017. iCaRL: Incremental Classifier and Representation Learning. In *CVPR*.

Romero, A.; Ballas, N.; Kahou, S. E.; Chassang, A.; Gatta, C.; and Bengio, Y. 2015. FitNets: Hints For Thin Deep Nets. In *ICLR*.

Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *IJCV*.

Rusu, A. A.; Rabinowitz, N. C.; Desjardins, G.; Soyer, H.; Kirkpatrick, J.; Kavukcuoglu, K.; Pascanu, R.; and Hadsell, R. 2016. Progressive Neural Networks. *arXiv preprint arXiv:1606.04671*.

Shin, H.; Lee, J. K.; Kim, J.; and Kim, J. 2017. Continual Learning with Deep Generative Replay. In *NIPS*.

Simon, C.; Koniusz, P.; and Harandi, M. 2021. On Learning the Geodesic Path for Incremental Learning. In *CVPR*.

Tian, J.; Dai, X.; Ma, C.-Y.; He, Z.; Liu, Y.-C.; and Kira, Z. 2023. Trainable Projected Gradient Method for Robust Fine-tuning. In *CVPR*.

van de Ven, G. M.; and Tolias, A. S. 2019. Three Scenarios for Continual Learning. *arXiv preprint arXiv:1904.07734*.

Wang, F.-Y.; Zhou, D.-W.; Ye, H.-J.; and Zhan, D.-C. 2022. Foster: Feature boosting and compression for class-incremental learning. In *ECCV*.

Wortsman, M.; Ilharco, G.; Gadre, S. Y.; Roelofs, R.; Gontijo-Lopes, R.; Morcos, A. S.; Namkoong, H.; Farhadi, A.; Carmon, Y.; Kornblith, S.; et al. 2022a. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *ICML*.

Wortsman, M.; Ilharco, G.; Kim, J. W.; Li, M.; Kornblith, S.; Roelofs, R.; Lopes, R. G.; Hajishirzi, H.; Farhadi, A.; Namkoong, H.; et al. 2022b. Robust fine-tuning of zero-shot models. In *CVPR*.

Wu, Y.; Chen, Y.; Wang, L.; Ye, Y.; Liu, Z.; Guo, Y.; and Fu, Y. 2019. Large Scale Incremental Learning. In *CVPR*.

Yan, S.; Xie, J.; and He, X. 2021. DER: Dynamically Expandable Representation for Class Incremental Learning. In *CVPR*.

Yoon, J.; Yang, E.; Lee, J.; and Hwang, S. J. 2018. Lifelong Learning with Dynamically Expandable Networks. In *ICLR*.

Zagoruyko, S.; and Komodakis, N. 2017. Paying More Attention to Attention: Improving the Performance of Convolutional Neural Networks via Attention Transfer. In *ICLR*.

Zenke, F.; Poole, B.; and Ganguli, S. 2017. Continual Learning Through Synaptic Intelligence. In *ICML*.

Zhu, F.; Cheng, Z.; Zhang, X.-y.; and Liu, C.-l. 2021. Class-incremental learning via dual augmentation. In *NeurIPS*.