# HYPERPARAMETERS IN CONTINUAL LEARNING : A REALITY CHECK WITH CLASS-INCREMENTAL LEARNING

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Continual learning (CL) aims to train a model on a sequence of tasks (*i.e.*, a CL scenario) while balancing the trade-off between plasticity (learning new tasks) and stability (retaining prior knowledge). The dominantly adopted conventional evaluation protocol for CL algorithms selects the best hyperparameters within a given scenario and then evaluates the algorithms using these hyperparameters in the same scenario. However, this protocol has significant shortcomings: it overestimates the CL capacity of algorithms and relies on unrealistic hyperparameter tuning, which is not feasible for real-world applications. From the fundamental principles of evaluation in machine learning, we argue that the evaluation of CL algorithms should focus on assessing the generalizability of their CL capacity to unseen scenarios. Based on this, we propose the Generalizable Two-phase Evaluation Protocol consisting of a hyperparameter tuning phase and an evaluation phase. Both phases share the same scenario configuration (e.g., number of tasks) but are generated from different datasets. Hyperparameters of CL algorithms are tuned in the first phase and applied in the second phase to evaluate the algorithms. We apply this protocol to class-incremental learning, both with and without pretrained models. Across more than 8,000 experiments, our results show that most state-of-the-art algorithms fail to replicate their reported performance, highlighting that their CL capacity has been significantly overestimated in the conventional evaluation protocol.

033

000

001

002

004

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

#### 1 INTRODUCTION

In recent years, extensive research has been conducted on continual learning (CL) with the goal of effectively learning knowledge from a sequence of tasks (Wang et al., 2023). A neural network 035 model in such CL scenarios faces a crucial trade-off between learning new knowledge from novel tasks (plasticity) and maintaining knowledge on previous tasks (stability) (Mermillod et al., 2013). 037 To address this inherent trade-off, numerous algorithms have been proposed for successful CL in various domains (Wang et al., 2023). In these domains, many CL studies have focused on classification, primarily concentrating on class-incremental learning (class-IL) (Masana et al., 040 2020) without or with pretrained models (Zhou et al., 2024a). However, deploying CL algorithms 041 requires careful hyperparameter tuning. Figure 1 illustrates the conventional evaluation protocol 042 (including hyperparameter tuning) dominantly employed in both offline and online class-incremental 043 learning (Zhou et al., 2022; Boschini et al., 2022; Zhou et al., 2024b; Smith et al., 2023; Seo et al., 044 2024). Additionally, similar evaluation protocols are also widely applied across other CL domains for semantic segmentation (Cha et al., 2021b; Yuan & Zhao, 2024), test-time adaptation (Yoo et al., 2024; Lee et al., 2024), federated learning (Piao et al., 2024), self-supervised learning (Fini et al., 046 2022; Cha et al., 2024) and large language models (Ke et al., 2023; Wu et al., 2024). 047

Many algorithms have been considered state-of-the-art based on performance validated through the
 conventional evaluation protocol. However, this raises two issues: First, the hyperparameter tuning
 method used in this protocol is not applicable to real-world CL scenarios. Second, it results in
 evaluation overfitting to a given scenario and dataset, which in turn leads to an overestimation of their
 CL capacity. In other words, this protocol only assesses performance in a seen scenario but fails to
 evaluate generalizability to new, unseen ones—an essential aspect for real-world applications. While
 several alternative evaluation protocols and hyperparameter tuning methods have been proposed, they

056

060

061

062

096

098

099

102

103

104

**Conventional Evaluation Protocol for CL CL** Scenario Hyper Parameter Tuning Evaluation using  $\mathcal{H}^{2}$ Ϊ... Task 0 Task 1 Task T Train  $\theta_{1:T}$ Train  $\theta_{1:T}$  $(\boldsymbol{D}_{tr}^{1:T}, \boldsymbol{D}_{tr}^{1:T})$  $(\boldsymbol{D_{tr}^{1:T}}, \boldsymbol{D_{te}^{1:T}})$  $D_{t_1}^0$  $D^1$  $D_{ta}^T$ CL Alg.(H) Eval. Eval. CL Alg.( $\mathcal{H}^*$ ) Select best  $\mathcal{H}$ Reinit  $\boldsymbol{\theta}$ Generate Sample  $\mathcal{H}$ **Report Evaluation Result** н a CL scenario Dataset

Figure 1: This figure illustrates the conventional evaluation protocol. First, a CL scenario is constructed using a benchmark dataset, where each task has its own training, validation, and test sets. Second, to find the best hyperparameters  $\mathcal{H}^*$ , a model is sequentially trained up to the final task using the sampled hyperparameters. After training for each task t, the model  $\theta_t$  is evaluated using the validation dataset. This process is repeated for various hyperparameters, and the best hyperparameters are selected based on performance. Finally, the model is trained using the CL algorithm with the best hyperparameters  $\mathcal{H}^*$  in the same CL scenario, and report the evaluation result on the test dataset. Note that in many studies, the results are reported using only  $D_{val}$  without separating  $D_{te}$  (*i.e.*,  $D_{te} = D_{val}$ ).

also have limitations: 1) they require to tune additional hyperparameters for their methods (Delange et al., 2021; Liu et al., 2023), or 2) they are only applied to a few old algorithms, and have not gained widespread acceptance (Chaudhry et al., 2018b; Chen et al., 2023; Bornschein et al., 2023). As a result, the issues with the conventional evaluation protocol have been largely ignored, and it remains the dominant evaluation protocol for evaluating CL algorithms until now.

 In this paper, we aim to reveal the limitation of the conventional evaluation protocol by revisiting the fundamental principles of evaluation in machine learning. From this perspective, we argue that the evaluation of CL algorithms should prioritize assessing the generalizability of each algorithm's CL capacity across unseen scenarios.

To achieve this goal, we propose a revised evalua-081 tion protocol, the Generalizable Two-phase Evalua-082 tion Protocol (GTEP), which consists of two phases: the hyperparameter tuning phase and the evaluation phase. Both phases share the same CL scenario con-084 figuration (*e.g.*, the number of tasks and classes per 085 task) yet leverage distinct datasets. In the hyperparameter tuning phase, a model is trained using various 087 hyperparameters of an algorithm, and the best hyper-088 parameters are selected based on performance. These best hyperparameters are then applied directly to train 090 the model using the algorithm in the evaluation phase, 091 where the measured performance serves as a reliable 092 benchmark for the algorithm's CL capacity in unseen





scenarios. As an initial application of this protocol, we focus on the most actively researched domain
 of CL—class-incremental learning (class-IL)—both with and without pretrained models (Wang et al., 2023). From approximately 8,000 experiments, we derive the following key findings:

- First, as shown in Figure 2, most state-of-the-art class-IL algorithms achieve superior performance in the hyperparameter tuning (HT) phase, which is almost identical the conventional evaluation protocol. However, they reveal limited generalizability in their CL capacity in the evaluation (Eval) phase. This tendency is particularly pronounced in the recently proposed algorithms.
  - Second, further analysis shows that these algorithms are limited by long training times, a large number of required parameters, or significant performance variance, suggesting they are less efficient than expected.

Based on extensive experimental results with the proposed evaluation protocol, we highlight major shortcomings of the conventional approach, which consistently overestimates the CL capacity of algorithms. In conclusion, we advocate for a fundamental reconsideration of the evaluation protocol across all domains to drive meaningful progress in CL research.

# 108 2 RELATED WORK

109

110 **Continual learning** Continual learning (CL) research has been conducted in various domains (Wang 111 et al., 2023; Parisi et al., 2019; Delange et al., 2021; Masana et al., 2020). In the beginning, the 112 CL research focus on task-incremental learning (Parisi et al., 2019; Delange et al., 2021), exploring 113 diverse approaches (Li & Hoiem, 2017; Aljundi et al., 2018; Chaudhry et al., 2018a; Cha et al., 114 2021a; Yoon et al., 2017). As the field progressed, attention shifted to the more challenging scenario, class-incremental learning (class-IL) (Masana et al., 2020). This shift leads to the investigation of 115 116 exemplar-based methods, involving the effective utilization of exemplar memory storing a subset of the dataset from previous tasks (Rebuffi et al., 2017; Zhao et al., 2020; Cha et al., 2023a). Since 117 then, using the exemplar memory has become standard, with several methods building on this 118 foundation. Regularization-based methods, which overcome catastrophic forgetting by introducing a 119 novel regularization (Wu et al., 2019; Douillard et al., 2020), and model expansion-based methods, 120 which dynamically expand model capacity to balance the trade-off between stability and plasticity, 121 have become the most powerful approach, achieving state-of-the-art performance (Wang et al., 2022b; 122 Yan et al., 2021; Zhou et al., 2022; Wang et al., 2022a). 123

Class-IL using pretrained models has recently gained considerable attention for achieving strong 124 performance without relying on the exemplar memory (Zhou et al., 2024a). Prompt-based meth-125 ods enable class-IL through prompt learning while keeping the pretrained model frozen. These 126 approaches have evolved over time, incorporating techniques such as using prompt pool (Wang 127 et al., 2022d), prompt combination (Wang et al., 2022c), decomposed prompt (Smith et al., 2023), 128 and prompt generation (Jung et al., 2023). Additionally, representation-based methods derive class 129 prototypes from the pretrained model and use them for classification (Zhou et al., 2023b). To enhance 130 the separability of these prototypes, several recent methods have focused on reducing class-wise 131 correlation (McDonnell et al., 2024; Zhou et al., 2024b).

- 132 **Evaluation and hyperparameter tuning of CL** Several papers have proposed new evaluation 133 metrics and protocols for the proper assessment of CL algorithms in classification. Traditionally, 134 accuracy-based metrics (e.g., final and average accuracy) have been used as the primary metrics 135 of evaluating performance of CL algorithms (Parisi et al., 2019; Masana et al., 2020; Chaudhry 136 et al., 2018a). However, recent studies have highlighted limitations of these metrics, particularly 137 regarding computational costs (Prabhu et al., 2023) and learned representations (Cha et al., 2023b). 138 Delange et al. (2021) introduced a hyperparameter tuning method for task-incremental learning, which involves first conducting a maximum plasticity search and then selecting the best hyperparameters 139 using stability decay. Similarly, Liu et al. (2023) proposed a hyperparameter selection method for 140 class-IL based on a bandit algorithm. However, both approaches entail additional training costs and 141 the need to tune extra hyperparameters. Other studies have proposed evaluation protocols similar to 142 ours (Chaudhry et al., 2018b; Chen et al., 2023; Bornschein et al., 2023). However, these protocols 143 have only been applied to a limited number of older algorithms in specific domains, which fails to 144 fully uncover the limitations of the conventional evaluation protocol. In addition to these efforts, 145 despite discussions on proper CL evaluation (Mundt et al., 2022), the conventional evaluation protocol 146 has continued to dominate the assessment of state-of-the-art CL algorithms across various domains.
- 147 148

149 150

## 3 TOWARDS EVALUATING THE GENERALIZABILITY OF THE CL CAPACITY

### 151 3.1 MOTIVATION: IMPROPER HYPERPARAMETER TUNING

152 As shown in Figure 1, the primary flaw of the conventional evaluation protocol is that it optimizes an 153 algorithm's hyperparameters in a given CL scenario and then evaluates the algorithm using those 154 same hyperparameters. Surprisingly, many studies have reported their results by directly tuning 155 hyperparameters on test data without considering separate validation sets (*i.e.*, set  $D_{te}^{HT} = D_{val}^{HT}$ ), as 156 seen in studies such as Wu et al. (2019); Douillard et al. (2020); Zhao et al. (2020); Yan et al. (2021); 157 Wang et al. (2022b); Zhou et al. (2022); Wang et al. (2022a;d); Zhou et al. (2023b; 2024b), and 158 others. Note that this approach is only feasible in experimental scenarios where all task data is always 159 available. Consequently, this hyperparameter tuning method fails to capture the real challenges of CL and is not applicable to real-world situations. While many studies partially address this limitation 160 by reporting robust performance across various experiments with some fixed or minimally adjusted 161 hyperparameters (Wang et al., 2022a;d; Zhou et al., 2024b), these evaluations are still based on given

Task T

 $D_{tr}^T$ 

CL Scenario

Task 1

 $D_{tr}^1$ 

Task 0

 $D_{t}^{0}$ 

164 165

162

163



167

- 169 170
- 171

172 173

174 175

176



 $(\boldsymbol{D}_{tr}^{1:T}, \boldsymbol{D}_{val}^{1:T})$ 

Hyperparameter Tuning Phase

Train  $\theta_{1:T}$ 

Result

Р

64

 ${\mathcal H}$ 

 $\mathcal{H}_1$ 

Figure 3: Illustration of the proposed evaluation protocol. Both phases share the same CL scenario configuration (*e.g.*, the number of tasks and number of classes in each task) but they are generated from distinct datasets ( $D^{HT}$  and  $D^E$ ). Best hyperparameters are selected in the hyperparameter tuning phase. Then, the evaluation phase access a CL algorithm by training a model using them. Note that evaluating an algorithm solely based on the results from the hyperparameter tuning phase is identical to the conventional evaluation protocol without using  $D_{te}$ .

scenarios (*i.e.*, seen scenarios), making it challenging to assess whether the algorithms would perform
 equally well in unseen scenarios. Nevertheless, this conventional protocol remains the predominant
 evaluation protocol for assessing algorithms across most CL domains.

Algorithm 1: The Generalizable Two-phase Evaluation Protocol

**Input** : A CL algorithm  $\mathcal{A}$ , a model  $\theta$ , the dataset for the hyperparameter tuning phase  $D^{HT}$ , the dataset for the evaluation phase  $D^E$ , the number of random samplings R, the number of trials S, and the number of hyperparameters K. **Output** : Final evaluation result ( $P^E$ ) for a CL algorithm  $\mathcal{A}$  in the evaluation phase

1.  $\{(\mathcal{H}_i, P_i^{HT})\}_{i=1}^R \leftarrow \texttt{HyperparameterTuning}(\theta, \mathcal{A}, D^{HT}, R, S, K)$ 

191 192

186

187 188

189

190

193

194

196

197

2.  $\mathcal{H}^* \leftarrow \text{SelectBestHyperparameter}(\{(\mathcal{H}_i, P_i^{HT})\}_{i=1}^R)$ 

3.  $P^E \leftarrow \text{Evaluation}(\theta, \mathcal{A}, D^E, \mathcal{H}^*, S)$ 

3.2 GENERALIZABLE TWO-PHASE EVALUATION PROTOCOL (GTEP) FOR CL

Given the previously discussed issues with the conventional evaluation protocol, the key question 199 becomes: What hyperparameter tuning and evaluation protocol should be used to properly assess CL 200 algorithms? Note that effective evaluation in machine learning should prioritize realistic methods 201 tailored to each learning scenario, rather than rigidly adhering to assumptions (e.g., i.i.d.) for 202 theoretical convenience. In this regards, we argue that evaluating the generalizability of each 203 algorithm's CL capacity is essential. For example, consider a real-world CL scenario where an 204 algorithm is applied to a CL scenario consisting of a sequence of tasks. Since the entire task data 205 would not be fully accessible at once, the conventional hyperparameter tuning method cannot be 206 applied. In such cases, a reasonable approach is to construct a simulated CL scenario, reflecting the 207 expected actual CL scenario, using a benchmark or available dataset. This involves identifying the best hyperparameters in the simulated scenario and then applying them to the actual CL scenario. 208 In other words, one of the basic evaluation protocols-consistent with the fundamental principles 209 of evaluation in machine learning—is to tune hyperparameters in seen scenarios (e.g., simulated 210 scenarios) and test them in unseen scenarios (e.g., actual scenarios). 211

Building on the above concept, we propose a revised evaluation protocol consisting of two phases, the Generalizable Two-phase Evaluation Protocol (GTEP): hyperparameter tuning and evaluation. Figure 3 and Algorithm 1 outlines the overall process. The key idea is that CL scenarios for the hyperparameter tuning and evaluation phases are generated from different datasets (*i.e.*,  $D^{HT} \neq D^E$ ) but share the same scenario configuration (*e.g.*, the number of tasks and classes per task), based 216 on expectations on the actual scenario. In the hyperparameter tuning phase, the goal is to identify 217 the best hyperparameters for the CL algorithm. In the evaluation phase, these hyperparameters are 218 applied to assess the algorithm's CL capacity in unseen scenarios, providing a more realistic measure 219 of its generalizability.

220 The pseudo algorithm of the hyperparameter tuning phase is outlined in Algorithm 2 of the Appendix. 221 First, we randomly sample hyperparameters  $h_k$  from a predefined set  $h_k^{Set}$  and build a list of selected 222 hyperparameters  $\mathcal{H}_r$ . Next, we generate a predefined CL scenario using the function  $\mathcal{F}$  with shuffled 223 class orderings. Afterward, the model  $\theta$  is trained using the selected hyperparameters  $\mathcal{H}_r$ , the CL algorithm  $\mathcal{A}$ , and the training dataset  $D_{tr}^{HT}$ . Performance  $(P^{HT})$  is then measured on the validation dataset  $D_{val}^{HT}$ . This phase returns multiple sets of hyperparameters and their corresponding 224 225 226 performance. Next, using the SelectBestHyperparameter function in Algorithm 1, we select the best hyperparameters, denoted as  $\mathcal{H}^*$ . Note that the hyperparameter tuning phase is identical to 227 the conventional evaluation protocol. However, we only use the results from this phase to select the 228 best hyperparameters. 229

230 In the evaluation phase (shown in Algorithm 3 of the Appendix), we train a model  $\theta$  using the CL 231 algorithm with the best hyperparameters  $\mathcal{H}^*$ . The trained model is then tested on the validation 232 dataset  $D_{val}^{E}$ . The final performance metric is the averaged performance  $(P^{E})$  of the trained model across multiple class orderings, which serves as the evaluation criterion for the CL algorithm. 233

234 To find the best hyperparameters for each algorithm, we optimize both algorithm-specific hyperpa-235 rameters (e.g., regularization strength) and general hyperparameters (e.g., learning rate and batch 236 size). During the hyperparameter tuning phase, we train the model with R sets of randomly selected 237 hyperparameters and account for S task orderings per set. In the evaluation phase, we assess the 238 performance across S task orderings as well. In this paper, we set R = 30 and S = 5 for all experi-239 ments. We also take into account various similarity scenarios between the hyperparameter tuning dataset  $(D^{HT})$  and the evaluation dataset  $(D^E)$ . High similarity indicates that the characteristics of 240 the dataset used in the actual scenario are somewhat predictable, allowing us to generate a scenario 241 for the hyperparameter tuning phase using a similar dataset. Conversely, low similarity suggests 242 unpredictability, indicating that the datasets used to generate scenarios in both phases differ signifi-243 cantly. Evaluating each algorithm under both similarity cases offers a comprehensive understanding 244 of the generalizability of its CL capacity. Furthermore, these efforts towards accurate evaluation 245 highlight the methodological differences from previously proposed evaluation protocols (Chaudhry 246 et al., 2018b; Chen et al., 2023; Bornschein et al., 2023; Mundt et al., 2022), as the revised evaluation 247 protocol. Additionally, note that the high-level concept of the proposed protocol can be applied to 248 various CL domains by considering the specific characteristics of these domains (e.g., imbalanced 249 classes per task, blurred task boundaries, or entirely different domains such as semantic segmentation) 250 in the CL scenario generation process (denoted as  $\mathcal{F}$  in Algorithms 2 and 3) for both phases.

251

253 254

255

256

257 258

259 260

#### 4 **EXPERIMENTAL RESULTS**

In this section, we present extensive experimental results using our proposed protocol in the most actively researched domain of continual learning (CL) (Wang et al., 2023), class-incremental learning (class-IL) both without and with pretrained models (Masana et al., 2020; Zhou et al., 2024a; 2023a).

### 4.1 CLASS-INCREMENTAL LEARNING WITHOUT PRETRAINED MODELS

**Experimental** settings We conduct the hyperparameter tuning and evaluation phases 261 using benchmark datasets, as shown in Table 1. From 262 ImageNet-1k (Deng et al., 2009), we derive two 263 subsets, ImageNet-100-1 and ImageNet-100-2, each 264 containing 100 randomly selected non-overlapping 265 classes. To account for varying dataset similari-266 ties, we further divide CIFAR-100 (Krizhevsky et al., 2009) and ImageNet-100-1 into disjoint classes, gen-267 erating CIFAR-50-1, CIFAR-50-2, ImageNet-50-1, 268

Table 1: Scenarios and datasets.

Scenario	$D^{HT}$	$D^E$
10 Tasks (C10×T10)	ImageNet-100-1	ImageNet-100-2
6 Tasks (C50×T1 + C10×T5)		
10 Tasks (C5×T10)	ImageNet-50-1, CIFAR-50-1	ImageNet-50-2, CIFAR-50-2
$\begin{array}{c} 6 \text{ Tasks} \\ (\text{C25}{\times}\text{T1} + \text{C5}{\times}\text{T5}) \end{array}$		CH111C 50 2

and ImageNet-50-2. We focus on two primary class-incremental learning (class-IL) scenarios (Masana 269 et al., 2020): 10 Tasks, where the model learns an equal number of classes from each task, and 6

293

295

296

297

298

299

300

301

270 Tasks, where the model learns half of the total classes in the first task then evenly distributes the 271 remaining classes evenly across subsequent tasks. Note that evaluating using both scenarios has 272 been widely considered the proper assessment of each algorithm (Masana et al., 2020; Zhou et al., 273 2023a) The table presents the configuration of the number classes (C) and tasks (T) for each scenario. 274 We conduct experiments using ResNet (He et al., 2016). We employ two key performance metrics commonly used for evaluating class-IL algorithms (Masana et al., 2020): Acc is final classification 275 accuracy for the entire validation dataset after training the final task, and AvgAcc =  $\frac{1}{T} \sum_{t=1}^{T} Acc_t$ , 276 where  $Acc_t$  denotes accuracy on the validation data up to task t. The hyperparameters that yield the 277 highest **harmonic mean** of Acc and AvgAcc are selected during the hyperparameter tuning phase. 278

279 **Baselines** We evaluate nine major class-IL algorithms, including replay-based methods (Replay, 280 iCaRL (Rebuffi et al., 2017), and WA (Zhao et al., 2020)) and regularization-based methods (BiC (Wu 281 et al., 2019) and PODNet (Douillard et al., 2020)) and expansion-based methods (DER (Yan et al., 282 2021), FOSTER (Wang et al., 2022b), and BEEF (Wang et al., 2022a)). Note that we use the partially implemented DER, as neither PyCIL nor the official DER code includes the implementation details 283 for masking and pruning. Replay serves as a naive baseline, where a model is fine-tuned using both 284 the exemplar memory and the current task's dataset. Note that these algorithms have demonstrated 285 progressively improved performance in the order of their publication. Among them, FOSTER, 286 BEEF, and MEMO are recognized as the current *state-of-the-art*, reporting superior performance that 287 surpasses DER by a small margin. We conduct experiments using the implementation code proposed 288 in PyCIL (Zhou et al., 2023a). The size of the exemplar memory is set to 2000 for ImageNet-100, and 289 1000 for ImageNet-50 and CIFAR-50 variants. More details on settings, predefined hyperparameter 290 sets and selected hyperparameters are presented in Section B.1 of the Appendix. 291



Figure 4: Experimental results (AvgAcc) on the 10 Tasks scenario using ImageNet-100-1 for  $D^{HT}$ and ImageNet-100-2 for  $D^E$  (high similarity). The term 'Original' and  $\mathcal{H}^*$  refer to the use of reported hyperparameters and hyperparameters selected from our protocol, respectively. BEEF constantly returns NaN in training loss at specific seeds so we could not report its performance.

Experiments using original and selected hyperparameters To demonstrate whether the hyperpa-307 rameters identified during the hyperparameter tuning phase achieve better performance than those 308 previously reported, we conduct experiments with both sets of hyperparameters. Figure 4(a) presents 309 results on  $D^{HT}$  = ImageNet-100-1, showing that using the best hyperparameters ( $\mathcal{H}^*$ ) generally out-310 performs the original ones across all algorithms except BEEF. Note that the performance differences 311 among DER, FOSTER, and MEMO are within their respective standard deviations. Considering the 312 hyperparameter tuning phase aligns with the conventional evaluation protocol, this graph indicates 313 that each algorithm reflects the performance trends observed in their respective papers, gradually 314 improving over time in accordance with the order of publication. On the other hand, we confirm that 315 BEEF is significantly sensitive to hyperparameters, as it occasionally results in NaN (Not a Number) in training loss for specific seeds, even when using the original hyperparameters. 316

In the evaluation phase, we apply the best hyperparameters to assess the CL capacity in unseen scenarios generated by  $D^E$ . Note that, due to differences in the datasets between these phases, the final performance may vary across phases, even when using identical hyperparameters for each algorithm. Figure 4(b) presents experimental results. The graph shows that the CL capacity of the state-of-the-art algorithms (*i.e.*, FOSTER, BEEF, and MEMO) is significantly inferior to that of older algorithms, such as WA, BiC and PODNet. Additionally, BEEF again produces NaN values for certain seeds. In contrast, DER demonstrates superior generalizability of its CL capacity, consistently maintaining strong performance in both phases.

357

358

359

360

361

362

364

366

367

368



Figure 5: Bar graphs depict the experimental results from the evaluation phase. The Y-axis represents final classification accuracy (Acc). The parentheses next to each algorithm indicate the publication year. The bar graphs in the first row show the experimental results using the best hyerparameters selected in the hyperparameter tuning phase with  $D^{HT} = \text{CIFAR-50-1}$ , while the graphs in the second row show the results using  $D^{HT} = \text{ImageNet-50-1}$ . In cases of using ImageNet-50-1 or ImageNet-50-2, we encountered challenges in obtaining results for BEEF due to NaN issues.

347 Experiments across diverse similarity cases To broadly assess the generalizability of each algorithm's CL capacity, we conduct experiments across various similarity cases. The bar graphs 348 in the first row of Figure 5 display results for both high and low similarity cases, using the best 349 hyperparameters selected during the hyperparameter tuning phase using  $D^{HT} = CIFAR-50-1$ . In 350 most cases, iCaRL performs worse than Replay, and BiC also struggles in some cases (e.g., 6 tasks 351 in both high and low similarity settings). Additionally, both WA and PODNet outperform other 352 regularization-based methods, with PODNet particularly excelling in the 6 Tasks. Lastly, the current 353 state-of-the-art methods—FOSTER, BEEF, and MEMO—exhibit lower performance compared to 354 DER, with BEEF again showing significant sensitivity, especially on ImageNet-50-2. 355

The second row of Figure 5 presents results using the best hyperparameters selected based on  $D^{HT}$  = ImageNet-50-1. The trends are consistent with previous experiments: DER maintains superior performance in most cases, although FOSTER surpasses DER in the low similarity case (6 tasks). Additionally, BEEF suffers from NaN issues in training loss for certain seeds.





Additional analysis Figure 6(a) shows the evaluation results for each task t in the evaluation phase,
 with shaded areas representing the standard deviation across 5 trials. From these graphs, it is evident
 that DER consistently outperforms current state-of-the-art algorithms (*i.e.*, FOSTER, BEEF and
 MEMO). Considering the standard deviation, the performances of FOSTER and MEMO are nearly
 indistinguishable. Among the remaining algorithms, WA demonstrates relatively better performance
 while BEEF performs similarly to the order algorithms.

378 Recent studies have increasingly focused on evaluating CL algorithms based on their training costs, 379 particularly in terms of GPU usage and energy consumption (Prabhu et al., 2023; Chavan et al., 380 2023). However, these evaluations were often conducted by either limiting the number of training 381 iterations or comparing costs under a fixed number of training epochs. Building on this, we extend 382 the analysis by examining the final model size and total training time for each algorithm, using their best hyperparameters to ensure a fair and comprehensive comparison of efficiency. Figures 6(b) and 6(c) present scatter plots showing achieved accuracy, total parameter counts, and training times. DER 384 performs the best and requires relatively less training time. Nevertheless, it exhibits considerable 385 inefficiency in the total number of parameters, which increases linearly with the number of tasks, 386 raising concerns about its actual cost-efficiency as a CL algorithm. On the other hand, BiC, BEEF, 387 and MEMO fail to demonstrate superior performance while requiring similar or longer training times 388 compared to DER, highlighting their serious inefficiency.

389 390 391

417

421

422

423

424

425

#### 4.2 CLASS-INCREMENTAL LEARNING WITH PRETRAINED MODELS

392 **Experimental details** We conduct both the hyperparameter tuning and evaluation phases 393 using widely used datasets in class-incremental 394 learning (class-IL) with pretrained models, in-395 cluding CUB-200 (Wah et al., 2011), ImageNet-396 R (Hendrycks et al., 2021a), and ImageNet-397 A (Hendrycks et al., 2021b), all of which contain 200 classes. To explore diverse similarity cases, we 398 divide these datasets into two disjoint subsets, as out-399 lined in Table 2. Following Sun et al. (2023), we 400 consider two major class-IL scenarios: 20 Tasks and 401 10 Tasks, where the model learns an equal number of 402

Table 2	: Scenarios an	d datasets.
onario	$D^{HT}$	$D^E$

Scenario	$D^{n_1}$	$D^{L}$
$\begin{array}{r} \hline 20 \text{ Tasks} \\ (C10 \times T20) \\ \hline 10 \text{ Tasks} \\ (C20 \times T10) \\ \hline \end{array}$	CUB-200, ImageNet-R	ImageNet-R, ImageNet-A
$\frac{(C20\times110)}{20 \text{ Tasks}}$ $(C5\timesT20)$	CUB-100-1, ImageNat P. 1	CUB-100-2, ImagaNat P 2
10 Tasks (C10×T10)	inageivet-K-1	ImageNet-A-2, ImageNet-A-2

classes in each task. Note that the 20 Tasks scenario has been commonly regarded as the standard for 403 better evaluating algorithm performance due to the need to handle more tasks. For all experiments, 404 we employ the ViT B16 model, which is pretrained on ImageNet (Dosovitskiy et al., 2020). The best 405 hyperparameters are selected based on the same metrics: the **harmonic mean** of **Acc** and **AvgAcc**. 406

**Baselines** We select six major algorithms: prompt-based methods (L2P (Wang et al., 2022d), 407 DualPrompt (Wang et al., 2022c) and CODA-Prompt (Smith et al., 2023)) and representation-based 408 methods (Adam-Adapter (Zhou et al., 2023b), Ranpac (McDonnell et al., 2024) and EASE (Zhou 409 et al., 2024b)). Within each category, CODA-Prompt and EASE represent current state-of-the-art 410 algorithms. Although DAP (Jung et al., 2023) reports better performance within the prompt-based 411 method category, we exclude it due to fairness issues in comparison, as mentioned in Zhou et al. 412 (2024a). All experiments are conducted using code implemented in PILOT (Sun et al., 2023). Details 413 on experimental settings, predefined hyperparameter sets the best hyperparameters are proposed in 414 Section B.2 of the Appendix.





**Experiments using original and selected hyperparameters** To verify best hyperparameters 429 selected in the hyperparameter tuning phase, we conduct experiments on  $D^{HT} = CUB-200$  using 430 both the reported and selected hyperparameters of each algorithm. Figure 7(a) demonstrates that 431 using the selected hyperparameters leads to better performance across all algorithms. Additionally,



448  $D^{E} = \text{ImageNet-R-2}$   $D^{E} = \text{ImageNet-R-2}$   $D^{E} = \text{CUB100-2}$   $D^{E} = \text{ImageNet-A-2}$   $D^{E} = \text{ImageNet-A-2}$ 449 Figure 8: Bar graphs depict the experimental results from the evaluation phase. The Y-axis represents 450 final accuracy (Acc). In the legend, the parentheses next to each algorithm indicate the publication 451 year. The bar graphs in the first row show the experimental results using the best hyerparameters 452 selected in the hyperparameter tuning phase with  $D^{HT} = \text{CUB100-1}$ , while the graphs in the second 453 row display the results using  $D^{HT} = \text{ImageNet-R-1}$ .

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

we observe that the performance of each algorithm gradually improves in accordance with their
 publication order, as reported in the respective papers. Note that Ranpac and EASE demonstrate
 similar performance, with differences falling within their standard deviations.

458 Following our evaluation protocol, we apply the best hyperparameters for each algorithm in the 459 evaluation phase. We conduct experiments for two evaluation phases using ImageNet-R and ImageNet-A as  $D^E$  and Figure 7(b) shows experimental results. From these results, we can confirm the following 460 461 observations: First, among the prompt-based algorithms (solid lines), DualPrompt exhibits degraded performance compared to L2P in both evaluation phases. Additionally, CODA-Prompt demonstrates 462 superior performance in all cases, although it shows nearly identical performance to L2P in the 463 ImageNet-R. In the case of the representation-based algorithms (dashed lines), Ranpac consistently 464 maintains its superiority across all datasets. However, EASE, recognized as the current state-of-the-art, 465 shows significantly poorer performance in both evaluation phases. 466

**Experiments across diverse similarity cases** Figure 8 presents the experimental results evaluated 467 in the evaluation phase. Similar to trends reported in Zhou et al. (2024a), representation-based 468 methods generally outperform prompt-based methods. However, significant differences are observed 469 under the proposed evaluation protocol: First, the prompt-based methods have reported substantial 470 performance improvements over previous algorithms (e.g., 7-10% increases on the CUB200 dataset 471 for each algorithm (Zhou et al., 2024a)). However, the proposed evaluation protocol reveals either no 472 significant performance difference between them (e.g., low similarity (20 tasks) using ImageNet-R-2 473 in the first row of the graph) or cases where an order algorithm outperforms a newer one (e.g., high 474 similarity (20 tasks) using CUB100-2 in the first row of the graph). Second, the current state-of-the-art 475 representation-based method, EASE, often underperforms compared to Ranpac, especially in low 476 similarity cases (e.g., low similarity (10 tasks) using ImageNet-R-2 in the first row of the graph). 477 Lastly, while Ranpac achieves the best performance in most cases, it exhibits significantly degraded performance in several low similarity cases (e.g., low similarity (20 tasks) using ImageNet-A-2 in 478 the first row of the graph). This degradation is attributed to considerable performance instability in 479 certain tasks. 480

Additional analysis As we already confirmed in the previous experiments, Figure 9(a) illustrates that Ranpac suffers from significant instability in certain tasks, resulting in a substantial increase in standard deviation (shaded area). Furthermore, we observe that the state-of-the-art algorithms, EASE and CODA-Prompt in their respective categories, do not consistently outperform baseline algorithms like ADAM and DualPrompt in many cases, highlighting a lack of generalizability in their CL capacity.



Figure 9: Experimental analyses in the evaluation phase. All experimental results are obtained by first identifying the best hyperparameters using ImageNet-R-1 (20 Tasks) in the hyperparameter tuning phase, then evaluating each algorithm using ImageNet-A-2 (20 Tasks) in the evaluation phase. (b) and (c) show the results after training up to the final task.

Figures 9(b) and 9(c) display the number of trainable parameters and training times with the best hyperparameters. For prompt-based algorithms, training times are comparable; however, CODA-Prompt requires more parameters while delivering lower performance compared to DualPrompt. Among representation-based methods, the oldest algorithm (*i.e.*, ADAM) achieves the best performance with minimal costs in terms of trainable parameters and training time.

504 In Section C of the Appendix, we present additional experimental results, including training graphs 505 and numerical data related to the results discussed in the manuscript.

# 506 5 CONCLUDING REMARKS

507 **Problems with the conventional evaluation protocol** The conventional evaluation protocol, which 508 is predominantly used to assess continual learning (CL) algorithms, has significant flaws because it 509 fails to consider real-world situations. In particular, the hyperparameter tuning method, which relies 510 on repeated training in a given scenario is not only inapplicable to real-world CL scenarios but also 511 tends to overestimate the CL capacity of each algorithm. According to the fundamental evaluation 512 principles of machine learning, the evaluation of CL algorithms should prioritize assessing the generalizability of their CL capacity to unseen scenarios. In this regard, we propose the Generalizable 513 Two-phase Evaluation Protocol (GTEP), which involves tuning hyperparameters in seen scenarios 514 and then applying them to unseen scenarios, considering various similarity cases. 515

516 Summary of experimental findings Our experiments across various similarity cases provide several 517 key insights: First, the CL capacity of many algorithms, especially recent ones, has been significantly 518 overestimated. Although most state-of-the-art algorithms perform well in seen scenarios, their CL capacity to unseen scenarios is often lacking. Second, we found that some of these algorithms are 519 highly sensitive to hyperparameters, resulting in instances where they fail to learn specific task orders 520 or exhibit significant performance variance on certain tasks. These two findings indicate that they 521 have reported overfitted results to the seen scenarios under the conventional evaluation protocol, 522 raising serious questions about their real-world applicability. Finally, even algorithms that perform 523 relatively well in the proposed protocol often incur excessive costs (e.g., training time and number of 524 parameters), undermining one of the key objectives of continual learning: cost-efficiency. Although 525 we reported the experimental results in class-incremental learning, we argue that these issues can 526 naturally be inferred to occur in other CL domains that use the same conventional evaluation protocol. 527

How should we evaluate going forward? – Key Takeaways We believe that the proposed evaluation protocol provides a fundamental approach to assess the generalizability of each algorithm's CL capacity, taking into account both the fundamental evaluation principles in machine learning and its real-world applications. Therefore, to make meaningful progress in CL research, we suggest that future evaluations across all CL domains should at least check the following:

- Does the proposed algorithm outperform baseline algorithms when the best hyperparameters selected from the hyperparameter tuning phase are applied to the evaluation phase?
- In the evaluation phase, is the proposed algorithm more efficient in terms of training costs (*e.g.*, total parameters and training time) compared to baseline algorithms? Additionally, does it avoid significant instability?
- 538

533

534

536

There are several limitations to our work and we discuss both limitations and future work in Section D of the Appendix.

#### 540 REFERENCES 541

560

561

562

563 564

565

566

567

568

569

571

576

583

590

542 Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In Proceedings of the European Conference 543 on Computer Vision (ECCV), pp. 139–154, 2018. 544

- Jorg Bornschein, Alexandre Galashov, Ross Hemsley, Amal Rannen-Triki, Yutian Chen, Arslan 546 Chaudhry, Xu Owen He, Arthur Douillard, Massimo Caccia, Qixuan Feng, et al. Nevis' 22: A 547 stream of 100 tasks sampled from 30 years of computer vision research. Journal of Machine 548 Learning Research, 24(308):1-77, 2023. 549
- Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class-550 incremental continual learning into the extended der-verse. IEEE transactions on pattern analysis 551 and machine intelligence, 45(5):5497-5512, 2022. 552
- 553 Sungmin Cha, Hsiang Hsu, Taebaek Hwang, Flavio Calmon, and Taesup Moon. {CPR}: Classifier-554 projection regularization for continual learning. In International Conference on Learning Repre-555 sentations, 2021a. URL https://openreview.net/forum?id=F2v4aqEL6ze. 556
- Sungmin Cha, YoungJoon Yoo, Taesup Moon, et al. Ssul: Semantic segmentation with unknown 558 label for exemplar-based class-incremental learning. Advances in neural information processing systems, 34:10919-10930, 2021b. 559
  - Sungmin Cha, Sungjun Cho, Dasol Hwang, Sunwon Hong, Moontae Lee, and Taesup Moon. Rebalancing batch normalization for exemplar-based class-incremental learning. In Proceedings of the *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20127–20136, 2023a.
  - Sungmin Cha, Jihwan Kwak, Dongsub Shim, Hyunwoo Kim, Moontae Lee, Honglak Lee, and Taesup Moon. Towards more objective evaluation of class incremental learning: Representation learning perspective, 2023b.
  - Sungmin Cha, Kyunghyun Cho, and Taesup Moon. Regularizing with pseudo-negatives for continual self-supervised learning. In Forty-first International Conference on Machine Learning, 2024.
- 570 Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In Proceedings of the 572 European Conference on Computer Vision (ECCV), pp. 532–547, 2018a. 573
- 574 Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. arXiv preprint arXiv:1812.00420, 2018b. 575
- Vivek Chavan, Paul Koch, Marian Schlüter, and Clemens Briese. Towards realistic evaluation of 577 industrial continual learning scenarios with an emphasis on energy consumption and computational 578 footprint. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pp. 579 11506-11518, 2023. 580
- 581 Jiefeng Chen, Timothy Nguyen, Dilan Gorur, and Arslan Chaudhry. Is forgetting less a good inductive 582 bias for forward transfer? arXiv preprint arXiv:2303.08207, 2023.
- Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg 584 Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification 585 tasks. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2021. 586
- 587 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale 588 hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, 589 pp. 248–255. Ieee, 2009.
- 591 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An 592 image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.

594 595	Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled
595	outputs distillation for small-tasks incremental learning. In Computer Vision–ECCV 2020: 16th
090 507	European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16, pp. 86–102.
597	Springer, 2020.
590	Enrico Fini Victor G Turrisi Da Costa Xavier Alameda-Pineda Elisa Ricci Karteek Alahari and
599	Julien Mairal. Self-supervised models are continual learners. In <i>Proceedings of the IEEE/CVF</i>
601	Conference on Computer Vision and Pattern Recognition, pp. 9621–9630, 2022.
602	
603	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
604	recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition,
605	pp. 770–778, 2016.
606	Dan Hendrycks, Steven Basart, Norman Mu, Sauray Kadayath, Frank Wang, Evan Dorundo, Rahul
607	Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical
608	analysis of out-of-distribution generalization. In Proceedings of the IEEE/CVF international
609	conference on computer vision, pp. 8340–8349, 2021a.
610	
611	Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial
612	nn 15262–15271 2021b
613	pp. 15262–15271, 26216.
614	Dahuin Jung, Dongyoon Han, Jihwan Bang, and Hwanjun Song. Generating instance-level prompts
615	for rehearsal-free continual learning. In Proceedings of the IEEE/CVF International Conference
616	on Computer Vision, pp. 11847–11857, 2023.
617	Zimmen Ka Ville Chan Harmei Lin Tetenne Kanishi Cambah Kim and Ding Lin. Cantinual
618	Dixuali Ke, Tijia Shao, Haowel Lill, Tatsuya Kollishi, Gyunak Killi, and Ding Liu. Continual pre-training of language models. In <i>The Flowerth International Conference on Learning Represen</i>
619	tations 2023 URL https://openreview.net/forum?id=m_GDIItaI30
620	
621	Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
622	Daeun Lee, Jaehong Yoon, and Sung Ju Hwang, Becotta: Input-dependent online blending of experts
023	for continual test-time adaptation. In <i>Forty-first International Conference on Machine Learning</i> .
625	2024.
626	
627	Zhizhong Li and Derek Hoiem. Learning without forgetting. <i>IEEE transactions on pattern analysis</i>
628	and machine intelligence, 40(12):2935–2947, 2017.
629	Yaoyao Liu, Yingying Li, Bernt Schiele, and Ojanru Sun. Online hyperparameter optimization for
630	class-incremental learning. arXiv preprint arXiv:2301.05032, 2023.
631	
632	Marc Masana, Xialei Liu, Bartlomiej Twardowski, Mikel Menta, Andrew D Bagdanov, and Joost
633	van de weijer. Class-incremental learning: survey and performance evaluation on image classifica-
634	uon. <i>arxiv preprint arxiv.2010.13277</i> , 2020.
635	Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasneiad, and Anton van den Hengel.
636	Ranpac: Random projections and pre-trained models for continual learning. Advances in Neural
637	Information Processing Systems, 36, 2024.
638	Mandal Manualla di Annella Decadale de l'Decidi Decidi Marcia (1996) e 1997 de la 1997 de 1997 de 1997
639	Martial Mermillod, Aurelia Bugaiska, and Patrick Bonin. The stability-plasticity dilemma: Inves-
640	nsychology 4.504 2013
641	рыусногоду, т.50т, 2015.
642	Martin Mundt, Steven Lang, Quentin Delfosse, and Kristian Kersting. CLEVA-compass: A continual
643	learning evaluation assessment compass to promote research transparency and comparability. In
644	International Conference on Learning Representations, 2022. URL https://openreview.
040	net/iorum?id=rHMaBYbkkRJ.
040	

German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual 647 lifelong learning with neural networks: A review. Neural Networks, 113:54-71, 2019.

648	Hongming Piao, Yichen Wu, Dapeng Wu, and Ying Wei. Federated continual learning via prompt-
649	based dual knowledge transfer. In Forty-first International Conference on Machine Learning,
650	2024.
651	

- Ameya Prabhu, Hasan Abed Al Kader Hammoud, Puneet K Dokania, Philip HS Torr, Ser-Nam Lim, Bernard Ghanem, and Adel Bibi. Computationally budgeted continual learning: What does matter? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3698–3707, 2023.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:
   Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- Minhyuk Seo, Hyunseo Koh, Wonje Jeung, Minjae Lee, San Kim, Hankook Lee, Sungjun Cho, Sungik Choi, Hyunwoo Kim, and Jonghyun Choi. Learning equi-angular representations for online continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23933–23942, 2024.
- James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf
   Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual decomposed
   attention-based prompting for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11909–11919, 2023.
- Hai-Long Sun, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Pilot: A pre-trained model-based
   continual learning toolbox. *arXiv preprint arXiv:2309.07117*, 2023.
- Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd
   birds-200-2011 dataset. 2011.
- Fu-Yun Wang, Da-Wei Zhou, Liu Liu, Han-Jia Ye, Yatao Bian, De-Chuan Zhan, and Peilin Zhao.
   Beef: Bi-compatible class-incremental learning via energy-based expansion and fusion. In *The Eleventh International Conference on Learning Representations*, 2022a.
- Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *European conference on computer vision*, pp. 398–414. Springer, 2022b.
- Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning:
   Theory, method and application. *arXiv preprint arXiv:2302.00487*, 2023.
- Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *European Conference on Computer Vision*, pp. 631–648. Springer, 2022c.
- Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent
   Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 139–149, 2022d.
- Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari.
   Continual learning for large language models: A survey. *arXiv preprint arXiv:2402.01364*, 2024.
- Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large
   scale incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 374–382, 2019.
- Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3014–3023, 2021.
- Jayeon Yoo, Dongkwan Lee, Inseop Chung, Donghyun Kim, and Nojun Kwak. What how and when should object detectors update in continually changing test domains? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23354–23363, 2024.

702 703	Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. <i>arXiv preprint arXiv:1708.01547</i> , 2017.
704	r · · · · · · · · · · · · · · · · · · ·
705 706	Bo Yuan and Danpei Zhao. A survey on continual semantic segmentation: Theory, challenge, method and application. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2024.
707	Denne Zhao Vi Vice Cracing Can Die Zhang and Shu Teo Vie Meinteining discrimination and
708	bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. Maintaining discrimination and foirness in class incremental learning. In <i>Proceedings of the IEEE/CVE conference on computer</i>
709	vision and pattern recognition, pp. 13208–13217, 2020.
710	Da-Wei Zhou, Oi-Wei Wang, Han-Iia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards
711	memory-efficient class-incremental learning. arXiv preprint arXiv:2205.13218, 2022.
713	Da-Wei Zhou, Fu-Yun Wang, Han-Jia Ye, and De-Chuan Zhan. Pycil: a python toolbox for class-
714 715	incremental learning. SCIENCE CHINA Information Sciences, 66(9):197101-, 2023a. doi: https://doi.org/10.1007/s11432-022-3600-y.
716	
717	Da-Wei Zhou, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Revisiting class-incremental learn-
718	ing with pre-trained models: Generalizability and adaptivity are all you need. arXiv preprint
719	arXiv:2303.0/338, 2023b.
720	Da-Wei Zhou, Hai-Long Sun, Jingyi Ning, Han-Jia Ye, and De-Chuan Zhan. Continual learning with
721	pre-trained models: A survey. <i>arXiv preprint arXiv:2401.16386</i> , 2024a.
700	Da-Wei Zhou, Hai-Long Sun, Han-Jia Ye, and De-Chuan Zhan. Expandable subspace ensemble for
723	pre-trained model-based class-incremental learning. arXiv preprint arXiv:2403.12030, 2024b.
724	
725	
726	
/2/	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

# 756 A ALGORITHM TABLES

The function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : { $\{(\mathcal{H}_i, P_i^{HT})\}_{i=1}^R$ 1. Result $\leftarrow$ {} 2. for $r \leftarrow 1$ to $R$ do 3. for $k \leftarrow 1$ to $K$ do 4. $h_k \leftarrow \text{RandomSample}(h_k^{\text{Set}})$ 5. $\mathcal{H}_r \leftarrow (h_0, h_1, \cdots, h_K)$ 6. for $s \leftarrow 1$ to $S$ do 7. Initialize $\theta$ 8. $D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^S P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(\mathcal{P}^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^S P_s^E$	number	of random samplings $R$ , the number of trials $S$ , the number of hyperparameters $K$ , and
1. Result $\leftarrow$ {} 2. for $r \leftarrow 1$ to $R$ do 3. for $k \leftarrow 1$ to $K$ do 4. $h_k \leftarrow \text{RandomSample}(h_k^{\text{Set}})$ 5. $\mathcal{H}_r \leftarrow (h_0, h_1, \dots, h_K)$ 6. for $s \leftarrow 1$ to $S$ do 7. Initialize $\theta$ 8. $D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	Output	$\{(\mathcal{H}_i, P_i^{HT})\}_{i=1}^{R}$
2. for $r \leftarrow 1$ to $R$ do 3. for $k \leftarrow 1$ to $K$ do 4. $h_k \leftarrow \text{RandomSample}(h_k^{\text{Set}})$ 5. $\mathcal{H}_r \leftarrow (h_0, h_1, \cdots, h_K)$ 6. for $s \leftarrow 1$ to $S$ do 7. Initialize $\theta$ 8. $D_{tT}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result	1.	$\operatorname{Result} \leftarrow \{\}$
3. for $k \leftarrow 1$ to $K$ do 4. $h_k \leftarrow \text{RandomSample}(h_k^{\text{Set}})$ 5. $\mathcal{H}_r \leftarrow (h_0, h_1, \cdots, h_K)$ 6. for $s \leftarrow 1$ to $S$ do 7. Initialize $\theta$ 8. $D_{tT}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tT}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	2.	for $r \leftarrow 1$ to $R$ do
4. $h_k \leftarrow \text{RandomSample}(h_k^{\text{Set}})$ 5. $\mathcal{H}_r \leftarrow (h_0, h_1, \cdots, h_K)$ 6. <b>for</b> $s \leftarrow 1$ to $S$ <b>do</b> 7. Initialize $\theta$ 8. $D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. <b>for</b> $s \leftarrow 1$ to $S$ <b>do</b> 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	3.	for $k \leftarrow 1$ to K do
5. $\mathcal{H}_r \leftarrow (h_0, h_1, \dots, h_K)$ 6. <b>for</b> $s \leftarrow 1$ to $S$ <b>do</b> 7. Initialize $\theta$ 8. $D_{tT}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. <b>for</b> $s \leftarrow 1$ to $S$ <b>do</b> 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	4.	$h_k \leftarrow RandomSample(h_k^{Set})$
6. for $s \leftarrow 1$ to $S$ do 7. Initialize $\theta$ 8. $D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	5.	$\mathcal{H}_r \leftarrow (h_0, h_1, \cdots, h_K)$
7. Initialize $\theta$ 8. $D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT})))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	6.	for $s \leftarrow 1$ to $S$ do
8. $D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$ 9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result <b>Algorithm 3:</b> Pseudo algorithm of the evaluation phase <b>Input</b> : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . <b>Output</b> : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	7.	Initialize $\theta$
9. $P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$ 10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result Algorithm 3: Pseudo algorithm of the evaluation phase Input : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . Output : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	8.	$D_{tr}^{HT}, D_{val}^{HT} \leftarrow \mathcal{F}(\text{Shuffle}(D^{HT}))$
10. $P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$ 11. Add $(\mathcal{H}_r, P_r^{HT})$ to Result Algorithm 3: Pseudo algorithm of the evaluation phase Input : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . Output : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	9.	$P_s^{HT} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{HT}, D_{val}^{HT}, \theta, \mathcal{H}_r)$
<ol> <li>Add (<i>H<sub>r</sub></i>, <i>P<sub>r</sub><sup>HT</sup></i>) to Result</li> <li>Algorithm 3: Pseudo algorithm of the evaluation phase</li> <li>Input : A CL Algorithm <i>A</i>, a model <i>θ</i>, the dataset for the Eval phase <i>D<sup>E</sup></i>, the best hyperparameter value <i>H</i>*, the number of trials <i>S</i>, the number of hyperparameters <i>K</i>, and the function that generates a CL scenario <i>F</i>.</li> <li>Output : Final evaluation result (<i>P<sup>E</sup></i>) for <i>A</i> <ol> <li>for <i>s</i> ← 1 <i>to S</i> do</li> <li>Initialize <i>θ</i></li> <li><i>D<sup>E</sup><sub>tr</sub></i>, <i>D<sup>E</sup><sub>val</sub></i> ← <i>F</i>(Shuffle(<i>D<sup>E</sup></i>))</li> <li><i>P<sup>E</sup><sub>s</sub></i> ← TrainCL(<i>A</i>, <i>D<sup>E</sup><sub>tr</sub></i>, <i>D<sup>E</sup><sub>val</sub></i>, <i>θ</i>, <i>H</i>*)</li> <li><i>P<sup>E</sup><sub>s</sub></i> ← <sup>1</sup>/<sub>S</sub> ∑<sup>S</sup><sub>s=1</sub> <i>P<sup>E</sup><sub>s</sub></i></li> </ol> </li> </ol>	10.	$P_r^{HT} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^{HT}$
Algorithm 3: Pseudo algorithm of the evaluation phase         Input : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best         hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ .         Output : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do         2. Initialize $\theta$ 3. $D^E_{tr}, D^E_{val} \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P^E_s \leftarrow \text{TrainCL}(\mathcal{A}, D^E_{tr}, D^E_{val}, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P^E_s$	11.	Add $(\mathcal{H}_r, P_r^{HT})$ to Result
Algorithm 3: Pseudo algorithm of the evaluation phase Input : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . Output : Final evaluation result ( $P^E$ ) for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D^E_{tr}, D^E_{val} \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P^E_s \leftarrow \text{TrainCL}(\mathcal{A}, D^E_{tr}, D^E_{val}, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P^E_s$		
Input : A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best hyperparameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the function that generates a CL scenario $\mathcal{F}$ . Output : Final evaluation result $(P^E)$ for $\mathcal{A}$ 1. for $s \leftarrow 1$ to $S$ do 2. Initialize $\theta$ 3. $D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(\text{Shuffle}(D^E))$ 4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	Algorith	<b>m 3:</b> Pseudo algorithm of the evaluation phase
2. Initialize $\theta$ 3. $D_{tr}^{E}, D_{val}^{E} \leftarrow \mathcal{F}(\text{Shuffle}(D^{E}))$ 4. $P_{s}^{E} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{E}, D_{val}^{E}, \theta, \mathcal{H}^{*})$ 5. $P^{E} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_{s}^{E}$	Input : A hyperpa function Output :	A CL Algorithm $\mathcal{A}$ , a model $\theta$ , the dataset for the Eval phase $D^E$ , the best rameter value $\mathcal{H}^*$ , the number of trials $S$ , the number of hyperparameters $K$ , and the n that generates a CL scenario $\mathcal{F}$ . Final evaluation result $(P^E)$ for $\mathcal{A}$ for $s \leftarrow 1$ to $S$ do
3. $D_{tr}^{E}, D_{val}^{E} \leftarrow \mathcal{F}(\text{Shuffle}(D^{E}))$ 4. $P_{s}^{E} \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^{E}, D_{val}^{E}, \theta, \mathcal{H}^{*})$ 5. $P^{E} \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_{s}^{E}$	1.	Initialize $\theta$
4. $P_s^E \leftarrow \text{TrainCL}(\mathcal{A}, D_{tr}^E, D_{val}^E, \theta, \mathcal{H}^*)$ 5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P_s^E$	1. 2.	
5. $P^E \leftarrow \frac{1}{S} \sum_{s=1}^{S} P^E_s$	1. 2. 3.	$D_{tr}^E, D_{val}^E \leftarrow \mathcal{F}(Shuffle(D^E))$
	1. 2. 3. 4.	$D_{tr}^{E}, D_{val}^{E} \leftarrow \mathcal{F}(Shuffle(D^{E}))$ $P_{s}^{E} \leftarrow TrainCL(\mathcal{A}, D_{tr}^{E}, D_{val}^{E}, \theta, \mathcal{H}^{*})$

827 828

835

836

837 838

839

840

841 842

843

844

# B Additional Details on Experimental Settings

#### 812 B.1 CLASS-INCREMENTAL LEARNING WITHOUT A PRETRAINED MODEL 813

**Experimental details** We conduct all experiments using PyCIL (Zhou et al., 2023a) in the following environment: Python 3.8, PyTorch 1.13.1, and CUDA 11.7. We use ResNet-18 and ResNet-32 architectures for our experiments. For class-incremental learning without a pretrained model, we employ the SGD optimizer with a momentum of 0.9 across all methods, consistent with their respective implementations. Other hyperparameters, however, are sampled during the hyperparameter tuning phase.

Table 3: Hyperparameters for training the first task.

Hyperparameters	Values
Init epochs	200
Init learning rate	0.1
Init milestones	[60, 120, 170] (Only applied when 'StepLR' is selected)
Init learning rate decay	0.1
Init weight decay	0.0005

Predefiend hyperparameters Recent studies have demonstrated that newer algorithms perform
better when trained for more epochs on the first task and fewer epochs on subsequent tasks (Masana
et al., 2020). Additionally, it is known that performance on the first task significantly impacts overall
performance (Cha et al., 2023b). To apply this approach consistently across all algorithms, we train a
model on the first task using the hyperparameters listed in Table 3. Subsequently, we train that model
with randomly sampled hyperparameters starting from the second task.



Figure 10 shows the number of hyperparameters for each algorithm. We consider both algorithm-845 specific and general hyperparameters in the hyperparameter tuning phase. Table 4 presents the sets of 846 predefined hyperparameters considered for each algorithm. Note that 'Epoch', 'Num milestones', 'LR 847 decay', 'Batch size', 'Weight decay', and 'LR scheduler' are commonly considered hyperparameters 848 for all algorithms. Additionally, both 'Num milestones' and 'Lr decay' are applicable only when 849 'StepLR' is selected as a scheduler. The others are specific hyperparameters of each algorithm. We 850 consider all the hyperparameters necessary for implementing each algorithm. For instance, even if a 851 specific algorithm uses the same value for a particular hyperparameter across all experiments (e.g., 852 fixing the strength of an additional regularization to 1), we aimed to find the best hyperparameter for 853 it (e.g., setting the strength as  $\alpha$  and finding the best value of it in the hyperparameter tuning phase). 854 We determine the range of values for the predefined hyperparameters based on the following criteria. 855 First, for general hyperparameters, we establish the range to include all optimal values reported by each algorithm. For specific hyperparameters related to each algorithm, we not only include the 856 optimal values report in the papers but also considered the full range of values that were explored 857 during their hyperparameter searches. 858

When the LR scheduler is set to StepLR, the milestones must be determined. To achieve this, we generalize the process of random sampling based on the milestones used in existing algorithms. First, we randomly sample num\_milestones. Based on this sampling, the milestones for the StepLR are set according to the following rule: For example, if Num\_milestones is set to 2, the milestones are defined as [epoch\*(2/5), epoch\*(4/5)]. If set to 3, the milestones become [epoch\*(2/7), epoch\*(4/7), epoch\*(6/7)]. Similarly, for 4 milestones, the values are [epoch\*(2/9), epoch\*(4/9), epoch\*(6/9), epoch\*(8/9)]. However, note that the num\_milestones is ignored when another LR schduler is selected. 

Algorithm	Hyperparameter Name	$h^{Set}$
	Epoch	[30, 70, 120, 160, 200]
	ĹR	[0.05, 0.1, 0.15, 0.2, 0.3]
	Num	[2, 2, 4]
All algorithms	milestones	[2, 5, 4]
	LR	[0 1 0 3 0 5]
	decay	[0.1, 0.5, 0.5]
	Batch	[32 64 128 256 512]
	size	[52, 04, 120, 250, 512]
	Weigh	[0.0001_0.0005_0.001_0.005]
	decay	[0.0001; 0.0005; 0.001; 0.005]
	LR	['StepLR' 'Cosine']
	Scheduler	[ StepErk , Cosine ]
iCaRL BiC WA and EOSTER	Т	[0 5 1 1 5 2 2 5]
	(KD)	[0.0, 1, 1.0, 2, 2.0]
BiC. WA and FOSTER	$\lambda$	[05, 1, 15, 2, 3]
	(KD)	[010, 1, 110, 2, 0]
BiC	Split	[0.05, 0.1, 0.15, 0.2, 0.3]
	ratio	[0.00, 0.11, 0.12, 0.2, 0.0]
iCaRL PODNet DER and MEMO	$\lambda$	[05, 1, 15, 2, 3]
	(Aux)	[0.0, 2, 2.0, 2, 2]
FOSTER	$\lambda$	[0.5, 1, 1.5, 2, 3]
	(FE)	
FOSTER	$\beta_1$	[0.93, 0.95, 0.97, 0.99]
FOSTER	$\beta_2$	[0.93, 0.95, 0.97, 0.99]
PODNet	Num	[10 20 30 50 100]
102100	proxy	
PODNet_FOSTER and BEEF	Post FT	[5, 10, 20, 30, 50]
T ODTAN, T OD TERA AND DEER	epochs	7 [30, 70, 120, 160, 200] (FOSTER and BEEF)
PODNet	Post FT	[0.001.0.003.0.005.0.007.0.01]
Toblet	LR	[0.001, 0.003, 0.003, 0.007, 0.01]
PODNet	Adaptive factor	[True, False]
BEFE	Energy	[0.001_0.005_0.01_0.02_0.05]
	weight	[0.001, 0.005, 0.01, 0.02, 0.05]
BEEF	Logit	[1114172023]
	alignment	[1.1, 1.1, 1.7, 2.0, 2.0]
MEMO	Exemplar	[16, 32, 64, 128, 256]
in Enito	batch size	[10, 52, 01, 120, 250]

Table 4: The predefined set of hyperparametes for class-IL without a pretrained model.

Original hyperparameters The following shows the original hyperparameters of each algorithm reported in PyCIL.

921	• Replay: ep_70_lr_0.1_lr_decay_0.1_batch_128_w_decay_0.0002_scheduler_steplr
922	• BiC: en 170 lr 0.1 lr decay 0.1 batch 128 w decay 0.0002 scheduler stenlr
923	T 2 lambda kd 0 split ratio 0.1
924	$POPN_{t} = 160 \text{ miletare } 2 \text{ is } 0.1 \text{ is decay } 0.1 \text{ betch } 120 \text{ miletare } 0.0005  selected a second$
925	• PODNet: ep_160_milestone_2_ir_0.1_ir_decay_0.1_batch_128_w_decay_0.0005_scheduler_cosine lambda c 5 lambda f 1.0 nb proxy 10 ft epochs 20 ft lrate 0.005 adaptive factor True
926	• EOSTER: en 170 lr 0.1 lr decay 0.1 batch 128 w decay 0.0005 scheduler cosine
927	T 2 lambda kd 1 fe 1 beta 0.96 0.97 comp en 130
928	$MEMO_{1} = 170 \text{ milestone} 2 \text{ is } 0.1 \text{ is decay} 0.1 \text{ hetch } 120 \text{ we decay} 0.0002 \text{ scheduler steple}$
929 930	lambda_aux_1_examplar_bs_64
931	• iCaRL: ep 170 lr 0.1 lr decay 0.1 batch 128 w decay 0.0002 scheduler steplr
932	T_2_lambda_aux_1
933	• WA: ep_170_milestone_3_lr_0.1_lr_decay_0.1_batch_128_w_decay_0.0002_scheduler_steplr
934	T_2.0_lambda_kd_0
935 936	• DER: ep_170_lr_0.1_lr_decay_0.1_batch_128_w_decay_0.0002_scheduler_steplr
937	Tambua_aux_1
938	• BEEF: ep_170_milestone_4_lr_0.1_lr_decay_0.1_batch_128_w_decay_0.0005_scheduler_cosine
939	fusion_ep_60_energy_w_0.01_logits_align_1.7
940	Note that setting 'lambda $kd = 0$ ' for both BiC and WA indicates the use of their adaptive rule
941	The that setting fullout_ke = 0 for both bie and with indicates the use of their adaptive fullo.
942	
943	
944	
945	
946	
947	
948	
949	
950	
951	
952	
953	
954	
955	
956	
957	
958	
959	
960	
961	
962	
963	
904	
302	
900	
907	
969	
970	
971	
511	

Best hyperparameters (ImageNet-100, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-100 (10 Tasks).

975	
976	• Replay: ep_70_milestone_3_lr_0.2_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_steplr
977	• BiC: ep 120 milestone 3 lr 0.1 lr decay 0.1 batch 32 w decay 0.0001 scheduler steplr
978	$T_1_ambda_kd_3.0_split_ratio_0.1$
979	• PODNet en 30 milestone 4 lr 0.05 lr decay 0.1 batch 64 w decay 0.0001 scheduler steplr
980	lambda c 3 lambda f 1.5 nb proxy 20 ft epochs 5 ft lrate 0.005 adaptive factor False
981	• FOSTER: en 30 milestone 4 lr 0.05 lr decay 0.1 hatch 64 w decay 0.0001 scheduler stenlr
982	T_1_lambda_kd_1.5_fe_1_beta_0.93_0.97_comp_ep_160
984 985	• MEMO: ep_120_milestone_3_lr_0.15_lr_decay_0.1_batch_512_w_decay_0.001_scheduler_steplr ambda_aux_0.5_examplar_bs_32
986 987	<ul> <li>iCaRL: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine</li> <li>T 2.5 lambda aux 2</li> </ul>
988 989	<ul> <li>WA: ep_120_milestone_3_lr_0.1_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_steplr</li> <li>T_1_lambda_kd_3.0_split_ratio_0.1</li> </ul>
990 991	<ul> <li>DER: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine lambda_aux_3</li> </ul>
992	$\mathbf{P}_{\mathbf{r}} = \mathbf{P}_{\mathbf{r}} + $
993	• BEEF: ep_120_ninestone_2_ir_0.2_ir_decay_0.5_batch_128_w_decay_0.0001_scheduler_steph fusion_ep_30_energy_w_0.02_logits_align_2.3
994	$\mathbf{D}_{\mathbf{r}}$
995 996	of each algorithm selected in the hyperparameter tuning phase using ImageNet-100 (6 Tasks).
997	• Replay: ep. 70 milestone 3 lr 0.2 lr decay 0.1 batch 64 w decay 0.0001 scheduler steplr
998	DiCt or 20 milestone 4 la 0.05 la decay 0.1 betch (4 m decay 0.0001 scheduler steph
999 1000	T_1_lambda_kd_1.5_split_ratio_0.1
1001 1002	• PODNet: ep_30_milestone_2_lr_0.15_lr_decay_0.1_batch_128_w_decay_0.001_scheduler_steplr lambda_c_9_lambda_f_0.5_nb_proxy_100_ft_epochs_10_ft_lrate_0.007_adaptive_factor_False
1003 1004	• FOSTER: ep_70_milestone_3_lr_0.05_lr_decay_0.1_batch_512_w_decay_0.0001_scheduler_cosine T_2.5_lambda_kd_0.5_fe_3_beta_0.95_0.93_comp_ep_30
1005 1006	• MEMO: ep_120_milestone_3_lr_0.15_lr_decay_0.1_batch_512_w_decay_0.001_scheduler_steplr lambda_aux_0.5_examplar_bs_32
1007 1008	• iCaRL: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine T 2.5 lambda aux 2
1009 1010	<ul> <li>WA: ep_170_lr_0.1_lr_decay_0.1_batch_128_w_decay_0.0002_scheduler_steplr</li> <li>T 2 lambda kd 1</li> </ul>
1011	<ul> <li>DER: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine lambda_aux_3</li> </ul>
1013 1014 1015	• BEEF: ep_30_milestone_4_lr_0.05_lr_decay_0.1_batch_128_w_decay_0.0001_scheduler_steplr fusion ep 70 energy w 0.01 logits align 1.4
1016 1017 1018	<b>Best hyperparameters (CIFAR-50, 10 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using CIFAR-50 (10 Tasks).
1019	• Replay: ep_160_milestone_3_lr_0.15_lr_decay_0.3_batch_32_w_decay_0.0001_scheduler_cosine
1020 1021	• BiC: ep_200_milestone_2_lr_0.1_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_cosine T 0.5 lambda kd 0.5 split ratio 0.2
1022 1023	<ul> <li>PODNet: ep_70_milestone_2_lr_0.1_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_steplr lambda_c_1_lambda_f_1_nb_proxy_10_ft_epochs_30_ft_lrate_0.007_adaptive_factor_False</li> </ul>
1024	• FOSTER: ep_120_milestone_3_lr_0.1_lr_decay_0.5_batch_32_w_decay_0.0001_scheduler_steplr T_2_lambda_kd_1.5_fe_0.5_beta_0.97_0.93_comp_ep_160

1026 1027	• MEMO: ep_120_milestone_4_lr_0.05_lr_decay_0.3_batch_32_w_decay_0.0005_scheduler_steplr lambda_aux_0.5_examplar_bs_16
1028 1029 1030	• iCaRL: ep_70_milestone_3_lr_0.05_lr_decay_0.3_batch_32_w_decay_0.001_scheduler_cosine T_2.5_lambda_aux_1
1031 1032	• WA: ep_160_milestone_4_lr_0.05_lr_decay_0.1_batch_64_w_decay_0.001_scheduler_cosine T_2_lambda_kd_3
1033 1034	• DER: ep_200_milestone_3_lr_0.2_lr_decay_0.1_batch_256_w_decay_0.001_scheduler_cosine lambda_aux_2
1035 1036	• BEEF: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_128_w_decay_0.0001_scheduler_cosine fusion_ep_200_energy_w_0.02_logits_align_2.3
1037 1038 1039	<b>Best hyperparameters (CIFAR-50, 6 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using CIFAR-50 (6 Tasks).
1040	• Replay: ep_70_milestone_2_lr_0.05_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_cosine
1042 1043	• BiC: ep_120_milestone_2_lr_0.05_lr_decay_0.3_batch_32_w_decay_0.0001_scheduler_cosine T_2.5_lambda_kd_1.5_split_ratio_0.3
1044 1045	• PODNet: ep_30_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.0005_scheduler_cosine lambda_c_1_lambda_f_3_nb_proxy_30_ft_epochs_50_ft_lrate_0.003_adaptive_factor_False
1046 1047	• FOSTER: ep_70_milestone_2_lr_0.05_lr_decay_0.1_batch_64_w_decay_0.0005_scheduler_steplr T_1.5_lambda_kd_1_fe_3_beta_0.97_0.93_comp_ep_200
1048 1049 1050	• MEMO: ep_160_milestone_4_lr_0.05_lr_decay_0.1_batch_32_w_decay_0.001_scheduler_cosine lambda_aux_0.5_examplar_bs_256
1051 1052	• iCaRL: ep_120_milestone_2_lr_0.05_lr_decay_0.1_batch_32_w_decay_0.0005_scheduler_steplr T_1_lambda_aux_1
1053 1054	• WA: ep_160_milestone_3_lr_0.05_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine T_2_lambda_kd_1.5
1055 1056	• DER: ep_120_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001_scheduler_cosine lambda_aux_1.5
1057 1058 1059	• BEEF: ep_30_milestone_4_lr_0.05_lr_decay_0.1_batch_128_w_decay_0.0001_scheduler_steplr fusion_ep_70_energy_w_0.01_logits_align_1.4
1060 1061 1062	<b>Best hyperparameters (ImageNet-50, 10 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-50 (10 Tasks).
1063	• Replay: ep_70_milestone_3_lr_0.2_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_steplr
1064 1065	• BiC: ep_120_milestone_3_lr_0.1_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_steplr T_1_lambda_kd_3_split_ratio_0.1
1066 1067	• PODNet: ep_30_milestone_4_lr_0.2_lr_decay_0.3_batch_64_w_decay_0.0005_scheduler_cosine lambda_c_7_lambda_f_1_nb_proxy_50_ft_epochs_20_ft_lrate_0.007_adaptive_factor_True
1069 1070	• FOSTER: ep_120_milestone_3_lr_0.1_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_steplr T_1_lambda_kd_3_fe_1_beta_0.99_0.93_comp_ep_160
1071 1072	• MEMO: ep_120_milestone_3_lr_0.1_lr_decay_0.1_batch_32_w_decay_0.0001_scheduler_steplr lambda_aux_1_examplar_bs_256
1073 1074	• iCaRL: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine T_2.5_lambda_aux_2
1075 1076	• WA: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine T_2.5_lambda_kd_2
1077 1078	• DER: ep_200_milestone_3_lr_0.15_lr_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine lambda_aux_3
1079	- DEEE NAN

• BEEF NaN

1080 1081	<b>Best hyperparameters (ImageNet-50, 6 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-50 (6 Tasks).
1082	• Replay: ep 200 milestone 2 lr 0.2 lr decay 0.3 batch 32 w decay 0.0001 scheduler steplr
1084	• BiC on 120 milestone 3 lr 0.1 lr decay 0.1 batch 32 w decay 0.0001 scheduler steplr
1085	T_1_lambda_kd_3_split_ratio_0.1
1086	• PODNet: ep_30_milestone_4_lr_0.2_lr_decay_0.3_batch_64_w_decay_0.0005_scheduler_cosine
1088	lambda_c_7_lambda_f_1_nb_proxy_50_ft_epochs_20_ft_lrate_0.007_adaptive_factor_True
1089	• FOSTER: ep_30_milestone_4_lr_0.2_lr_decay_0.3_batch_64_w_decay_0.0005_scheduler_cosine T_2_lambda_kd_1_fe_2_beta_0.97_0.99_comp_ep_120
1090	• MEMO: ep_120_milestone_3_lr_0.15_lr_decay_0.1_batch_512_w_decay_0.001_scheduler_steplr
1092	$\frac{1}{1000} = \frac{1}{1000} = \frac{1}{1000} = \frac{1}{1000} = \frac{1}{10000} = \frac{1}{10000000000000000000000000000000000$
1093 1094	• iCaRL: ep_200_milestone_3_ir_0.15_ir_decay_0.1_batch_64_w_decay_0.0001_scheduler_cosine T_2.5_lambda_aux_2
1095 1096	• WA: ep_120_milestone_4_lr_0.1_lr_decay_0.5_batch_64_w_decay_0.0005_scheduler_steplr T_1_lambda_kd_1
1097	$DED = 120 \text{ milestars} (1 + 0.2) \ln decen (0.5) heath (120 + 0.0005) established as in$
1098	• DER: ep_120_Innesione_4_ir_0.5_ir_decay_0.5_batch_128_w_decay_0.0005_scheduler_cosine lambda_aux_0.5
1099	• DEEE: NoN
1100	• BEEF. Indin
1101	
1102	
1103	
1104	
1105	
1106	
1107	
1108	
1109	
1110	
1110	
1112	
1114	
1115	
1116	
1117	
1118	
1119	
1120	
1121	
1122	
1123	
1124	
1125	
1126	
1127	
1128	
1129	
1130	
1132	
1133	

B.2 EXPERIMENTAL SETTINGS FOR CLASS-INCREMENTAL LEARNING WITH A PRETRAINED MODEL 

Experimental details For experiments using the proposed evaluation protocol on class-incremental learning algorithms with a pretrained model, we employ the PILOT (Sun et al., 2023) code for each algorithm. The experimental setup closely followed PILOT's environment, using Python 3.8, PyTorch 2.0.1, and CUDA 11.7. 



**Pretrained hyperparameters** The process of selecting hyperparameters for algorithms using a pretrained model is similar to the previous experiments. We comprehensively consider both general hyperparameters and algorithm-specific ones, finding the best hyperparameters during the tuning phase. Figure 11 shows the number of hyperparameters for each algorithm. The predefined hyperparameters used for this process are listed in Table 5. Using the selected hyperparameters, we train each algorithm across the entire CL scenario. The range of each hyperparameter is set based on values reported in previous work for each type of algorithm. Unlike the algorithms without pretrained models, which use the same optimizer (*i.e.*, SGD), different optimizers have been used across algorithms in this case, so we also perform sampling for the optimizer. For hyperparameters of the optimizer that were not sampled, we use the default values provided in PyTorch.

Algorithm	**	
	Hyperparameter Name	$h^{Set}$
(		[3, 5, 10, 15, 20, 25]
	Epoch	(for L2P, DualPrompt. CODA-Pormpt)
	Epoon	/ [5, 10, 15, 20, 25, 30]
		(for Adam-Adapter, Ranpac, EASE)
All algorithms		[0.000875, 0.001575, 0.001875, 0.002575, 0.0025]
An argonanis	LR	(101 L21, Dual rompt. CODA-rompt)
		(for Adam-Adapter, Ranpac, EASE)
	Num	
	milestones	[2, 3, 4]
	LR	[0 1 0 3 0 5]
	decay	
	size	[6, 10, 24, 46, 04, 126] (for L 2P DualPrompt CODA-Prompt Adam-Adapter
	5120	[0, 0.0001, 0.0005]
	Weigh	(for L2P, DualPrompt, CODA-Prompt)
	decay	/ [0.0001, 0.0005, 0.001, 0.005]
		(for Adam-Adapter, Ranpac, EASE)
	LR	['steplr' 'cosine' 'constant']
	Scheduler	
LODD	Optimizer	['sgd', 'adam', 'adamw']
L2P, DualPrompt	M Size	
L2P L2P	$Length (L_p)$	
L2F L 2P DualPrompt		
DualPrompt	Prompt length of $g(L_{r})$	[5, 10, 15, 20, 30]
DualPrompt	$\frac{11000 \text{ Jeres for g } (L_g)}{\text{Length } (L_e)}$	[5, 10, 15, 20, 30]
CODA-Prompt	Pool size	[30, 50, 100, 200, 300]
CODA-Prompt	Prompt length	[4, 8, 16, 24, 32]
CODA-Prompt	Orthogonality Mu	[0.2, 0.1, 0.01, 0.001, 0]
Adam-Adapter, Ranpac,	EASE FFN num	[4,8,16,32,64]
Ranpac	М	[5000, 10000, 15000, 20000]
Ranpac	Prompt token num	[3, 5, 10, 20, 30, 50]
<b>Driginal hyperparan</b> eported in PILOT.	neters The following show	rs the original hyperparameters of each algorithm
Driginal hyperparan eported in PILOT. • L2P_ep_10_1 scheduler_co	neters The following show nilestone_3_lr_0.001875_h nstant_optimizer_adam_siz	r_decay_0_batch_32_w_decay_0 e_10_length_5_top_k_5_lamb_0.1
<ul> <li>Driginal hyperparan eported in PILOT.</li> <li>L2P_ep_10_1 scheduler_co</li> <li>DualPrompt_ scheduler_co</li> </ul>	neters The following show nilestone_3_lr_0.001875_li nstant_optimizer_adam_siz ep_10_milestone_4_lr_0.00 nstant_optimizer_adam_siz	rs the original hyperparameters of each algorithm decay_0_batch_32_w_decay_0 e_10_length_5_top_k_5_lamb_0.1 01_lr_decay_0.0_batch_24_w_decay_0.0 e_10_L_e_5_L_g_5_top_k_1_lamb_0.1
<ul> <li>Driginal hyperparan eported in PILOT.</li> <li>L2P_ep_10_1 scheduler_co</li> <li>DualPrompt_ scheduler_co</li> <li>CODA-Prom scheduler_co</li> </ul>	neters The following show milestone_3_lr_0.001875_ln nstant_optimizer_adam_siz ep_10_milestone_4_lr_0.00 nstant_optimizer_adam_siz pt_ep_50_milestone_2_lr_0 sine_optimizer_adam_e_po	s the original hyperparameters of each algorithm decay_0_batch_32_w_decay_0 e_10_length_5_top_k_5_lamb_0.1 01_lr_decay_0.0_batch_24_w_decay_0.0 e_10_L_e_5_L_g_5_top_k_1_lamb_0.1 0.001_lr_decay_0.0_batch_128_w_decay_0.0 ol_size_100_e_p_length_8_ortho_mu_0.0
<ul> <li>Driginal hyperparan eported in PILOT.</li> <li>L2P_ep_10_1 scheduler_co</li> <li>DualPrompt_ scheduler_co</li> <li>CODA-Prom scheduler_co</li> <li>Adam_ep_10 scheduler_co</li> </ul>	neters The following show nilestone_3_lr_0.001875_lr nstant_optimizer_adam_siz ep_10_milestone_4_lr_0.00 nstant_optimizer_adam_siz pt_ep_50_milestone_2_lr_0 sine_optimizer_adam_e_po _milestone_3_lr_0.05_lr_do nstant_optimizer_sgd_ffn_r	r_decay_0_batch_32_w_decay_0 e_10_length_5_top_k_5_lamb_0.1 01_lr_decay_0.0_batch_24_w_decay_0.0 e_10_L_e_5_L_g_5_top_k_1_lamb_0.1 0.001_lr_decay_0.0_batch_128_w_decay_0.0 ol_size_100_e_p_length_8_ortho_mu_0.0 ecay_0.0_batch_16_w_decay_0.005 num_100
<ul> <li>Driginal hyperparam eported in PILOT.</li> <li>L2P_ep_10_1 scheduler_co</li> <li>DualPrompt_ scheduler_co</li> <li>CODA-Prom scheduler_co</li> <li>Adam_ep_10 scheduler_co</li> <li>Ranpac_ep_1 scheduler_co</li> </ul>	neters The following show milestone_3_lr_0.001875_lr nstant_optimizer_adam_siz ep_10_milestone_4_lr_0.00 nstant_optimizer_adam_siz pt_ep_50_milestone_2_lr_0 sine_optimizer_adam_e_po _milestone_3_lr_0.05_lr_donstant_optimizer_sgd_ffn_r 0_milestone_2_lr_0.05_lr_ nstant_optimizer_sgd_ffn_r	The original hyperparameters of each algorithm $c_decay_0_batch_32_w_decay_0$ $e_10_length_5_top_k_5_lamb_0.1$ $01_lr_decay_0.0_batch_24_w_decay_0.0$ $e_10_L_e_5_L_g_5_top_k_1_lamb_0.1$ $0.001_lr_decay_0.0_batch_128_w_decay_0.0$ $ol_size_100_e_p_length_8_ortho_mu_0.0$ $ecay_0.0_batch_16_w_decay_0.005$ $num_100$ $decay_0.0_batch_16_w_decay_0.005$ $num_64_M_10000_pt_num_30$
<ul> <li>Driginal hyperparam eported in PILOT.</li> <li>L2P_ep_10_1 scheduler_co</li> <li>DualPrompt_ scheduler_co</li> <li>CODA-Prom scheduler_co</li> <li>Adam_ep_10 scheduler_co</li> <li>Ranpac_ep_1 scheduler_co</li> <li>EASE_ep_20 scheduler_co</li> </ul>	neters The following show milestone_3_lr_0.001875_ln nstant_optimizer_adam_siz ep_10_milestone_4_lr_0.00 nstant_optimizer_adam_siz pt_ep_50_milestone_2_lr_0 sine_optimizer_adam_e_po _milestone_3_lr_0.05_lr_d nstant_optimizer_sgd_ffn_r 0_milestone_2_lr_0.05_lr_ nstant_optimizer_sgd_ffn_rd i_milestone_4_lr_0.05_lr_d sine_optimizer_sgd_ffn_nu	The provided provide
<ul> <li>Driginal hyperparam eported in PILOT.</li> <li>L2P_ep_10_1 scheduler_co</li> <li>DualPrompt_ scheduler_co</li> <li>CODA-Prom scheduler_co</li> <li>Adam_ep_10 scheduler_co</li> <li>Ranpac_ep_1 scheduler_co</li> <li>EASE_ep_20 scheduler_co</li> <li>Best hyperparameter ach algorithm selecte</li> </ul>	neters The following show milestone_3_lr_0.001875_li nstant_optimizer_adam_siz ep_10_milestone_4_lr_0.00 nstant_optimizer_adam_siz pt_ep_50_milestone_2_lr_( sine_optimizer_adam_e_po _milestone_3_lr_0.05_lr_dd nstant_optimizer_sgd_ffn_r 0_milestone_2_lr_0.05_lr_ nstant_optimizer_sgd_ffn_r _milestone_4_lr_0.05_lr_d sine_optimizer_sgd_ffn_nu s (CUB-200, 20 Tasks) Th d in the hyperparameter tun	rs the original hyperparameters of each algorithm -decay_0_batch_32_w_decay_0 e_10_length_5_top_k_5_lamb_0.1 01_lr_decay_0.0_batch_24_w_decay_0.0 e_10_L_e_5_L_g_5_top_k_1_lamb_0.1 0.001_lr_decay_0.0_batch_128_w_decay_0.0 ol_size_100_e_p_length_8_ortho_mu_0.0 ecay_0.0_batch_16_w_decay_0.005 hum_100 decay_0.0_batch_16_w_decay_0.005 hum_64_M_10000_pt_num_30 ecay_0.0_batch_16_w_decay_0.005 m_64_alpha_0.1 e following represents the best hyperparameters of ing phase using CUB-200 (20 Tasks).

### Table 5: The predefined set of hyperparametes for class-IL with a pretrained model.

• DualPrompt: ep\_25\_milestone\_3\_lr\_0.000875\_lr\_decay\_0.1\_batch\_48\_w\_decay\_0.0005 scheduler\_constant\_optimizer\_adamw\_size\_20\_L\_e\_5\_L\_g\_30\_top\_k\_1\_lamb\_0.3

1242	• CODA-Prompt: ep_25_milestone_2_lr_0.000875_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_steplr_optimizer_sgd_e_pool_size_30_e_p_length_4_ortho_mu_0.01
1244 1245	<ul> <li>Adam: ep_15_milestone_4_lr_0.05_lr_decay_0.5_batch_48_w_decay_0.0005</li> <li>acheduler_conjug_activities_cod_ffrunn_8</li> </ul>
1246	scheduler_cosine_optimizer_sgu_nn_num_8
1247 1248	• Ranpac: ep_30_milestone_4_lr_0.01_lr_decay_0.1_batch_8_w_decay_0.0005 scheduler_cosine_optimizer_sgd_ffn_num_32_M_20000_pt_num_5
1249	• EASE: ep_15_milestone_4_lr_0.02_lr_decay_0.5_batch_128_w_decay_0.001 scheduler_cosine_optimizer_sgd_ffn_num_8_alpha_0.01
1251 1252	<b>Best hyperparameters (CUB-200, 10 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using CUB-200 (10 Tasks).
1253	
1254 1255	• L2P: ep_25_milestone_2_lr_0.001875_lr_decay_0.3_batch_64_w_decay_0.0005 scheduler_cosine_optimizer_adamw_size_10_length_6_top_k_6_lamb_0.5
1256 1257	• DualPrompt: ep_25_milestone_3_lr_0.0025_lr_decay_0.5_batch_128_w_decay_0.0005 scheduler_steplr_optimizer_sgd_size_20_L_e_10_L_g_10_top_k_1_lamb_0.5
1258 1259	• CODA-Prompt: ep_25_milestone_3_lr_0.0025_lr_decay_0.3_batch_64_w_decay_0 scheduler_cosine_optimizer_adamw_e_pool_size_100_e_p_length_8_ortho_mu_0
1260 1261	• Adam: ep_20_milestone_3_lr_0.04_lr_decay_0.3_batch_8_w_decay_0.0005 scheduler steplr optimizer sgd ffn num 32
1262	• Ranpac: ep 30 milestone 4 lr 0.02 lr decay 0.3 batch 16 w decay 0.0001
1263	scheduler_steplr_optimizer_sgd_ffn_num_64_M_10000_pt_num_3
1265	• EASE: ep_15_milestone_4_lr_0.01_lr_decay_0.3_batch_64_w_decay_0.0005
1266	scheduler_steplr_optimizer_sgd_ffn_num_8_alpha_0.05
1267 1268 1269	<b>Best hyperparameters (ImageNet-R, 20 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (20 Tasks).
1270 1271	<ul> <li>L2P_ep_25_milestone_3_lr_0.000875_lr_decay_0.5_batch_64_w_decay_0 scheduler_steplr_optimizer_adam_size_10_length_10_top_k_4_lamb_0.5</li> </ul>
1272 1273	<ul> <li>DualPrompt: ep_15_milestone_4_lr_0.001875_lr_decay_0.5_batch_128_w_decay_0 scheduler_steplr_optimizer_adam_size_20_L_e_30_L_g_5_top_k_1_lamb_0.5</li> </ul>
1274 1275	CODA Presents on 15 milestone 2 in 0.002275 in descer 0.1 betch 48 m descer 0.0001
1070	• CODA-Prompt: ep_15_milestone_2_ir_0.002375_ir_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001
1276 1277	<ul> <li>CODA-Prompt: ep_15_milestone_2_ir_0.002375_ir_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> </ul>
1276 1277 1278 1279 1280	<ul> <li>CODA-Prompt: ep_15_milestone_2_lr_0.002375_lr_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> </ul>
1276 1277 1278 1279 1280 1281 1282	<ul> <li>CODA-Prompt: ep_15_milestone_2_ir_0.002375_ir_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> </ul>
1276 1277 1278 1279 1280 1281 1282 1283 1284 1285	<ul> <li>CODA-Prompt: ep_15_milestone_2_lr_0.002375_lr_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> <li>Best hyperparameters (ImageNet-R, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (10 Tasks).</li> </ul>
1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287	<ul> <li>CODA-Prompt: ep_15_milestone_2_lr_0.002375_lr_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> <li>Best hyperparameters (ImageNet-R, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (10 Tasks).</li> <li>L2P: ep_25_milestone_3_lr_0.001375_lr_decay_0.5_batch_128_w_decay_0 exheduler_constant_optimizer_sd_ffn_decay_0.5_batch_128_w_decay_0 exheduler_constant_optimizer_sd_ffn_decay_0.5_batch_128_w_decay_0</li> </ul>
1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288	<ul> <li>CODA-Prompt: ep_15_milestone_2_ir_0.002375_ir_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> <li>Best hyperparameters (ImageNet-R, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (10 Tasks).</li> <li>L2P: ep_25_milestone_3_lr_0.001375_lr_decay_0.5_batch_128_w_decay_0 scheduler_constant_optimizer_adamw_size_20_length_6_top_k_10_lamb_0.3</li> </ul>
1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290	<ul> <li>CODA-Prompt: ep_13_milestone_2_ir_0.002375_ir_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> <li>Best hyperparameters (ImageNet-R, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (10 Tasks).</li> <li>L2P: ep_25_milestone_3_lr_0.001375_lr_decay_0.5_batch_128_w_decay_0 scheduler_constant_optimizer_adamw_size_20_length_6_top_k_10_lamb_0.3</li> <li>DualPrompt: ep_25_milestone_2_lr_0.001375_lr_decay_0.3_batch_128_w_decay_0.0005</li> </ul>
1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290 1291 1292	<ul> <li>CODA-Prompt: ep_15_milestone_2_lr_0.002375_lr_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> <li>Best hyperparameters (ImageNet-R, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (10 Tasks).</li> <li>L2P: ep_25_milestone_3_lr_0.001375_lr_decay_0.5_batch_128_w_decay_0 scheduler_constant_optimizer_adamw_size_20_length_6_top_k_10_lamb_0.3</li> <li>DualPrompt: ep_25_milestone_2_lr_0.001375_lr_decay_0.3_batch_128_w_decay_0.0005 scheduler_constant_optimizer_adamw_size_30_L_e_20_L_g_20_top_k_1_lamb_0.3</li> <li>CODA-Prompt: ep_20_milestone_2_lr_0.001375_lr_decay_0.1_batch_48_w_decay_0 scheduler_steplr_optimizer_adam_size_300_e_p_length_8_ortho_mu_0</li> </ul>
1276 1277 1278 1279 1280 1281 1282 1283 1284 1285 1286 1287 1288 1289 1290 1291 1292 1293 1294	<ul> <li>CODA-Prompt: ep_15_milestone_2_lr_0.002375_lr_decay_0.1_batch_48_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_300_e_p_length_32_ortho_mu_0.001</li> <li>Adam: ep_25_milestone_3_lr_0.05_lr_decay_0.5_batch_64_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_64</li> <li>Ranpac: ep_20_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_15000_pt_num_20</li> <li>EASE: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_24_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_16_alpha_0.15</li> <li>Best hyperparameters (ImageNet-R, 10 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R (10 Tasks).</li> <li>L2P: ep_25_milestone_3_lr_0.001375_lr_decay_0.5_batch_128_w_decay_0 scheduler_constant_optimizer_adamw_size_20_length_6_top_k_10_lamb_0.3</li> <li>DualPrompt: ep_25_milestone_2_lr_0.001375_lr_decay_0.3_batch_128_w_decay_0.0005 scheduler_constant_optimizer_adamw_size_30_L_e_20_L_g_20_top_k_1_lamb_0.3</li> <li>CODA-Prompt: ep_20_milestone_2_lr_0.001375_lr_decay_0.1_batch_48_w_decay_0 scheduler_steplr_optimizer_adam_e_pool_size_300_e_p_length_8_ortho_mu_0</li> <li>Adam: ep_30_milestone_2_lr_0.05_lr_decay_0.1_batch_64_w_decay_0.001</li> </ul>

• EASE: ep_30_milestone_4_lr_0.05_lr_decay_0.3_batch_128_w_decay_0.001 scheduler_cosine_optimizer_adam_ffn_num_16	
<ul> <li>Best hyperparameters (CUB-100-1, 20 Tasks) The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using CUB-100-1 (20 Tasks)</li> </ul>	5
• L2P: ep_20_milestone_3_lr_0.002375_lr_decay_0.3_batch_128_w_decay_0.0005	
<ul> <li>scheduler_constant_optimizer_adamw_size_20_length_8_top_k_4_lamb_0.5</li> <li>DualPrompt: ep_25_milestone_4_lr_0.001375_lr_decay_0.1_batch_128_w_decay_0</li> </ul>	
<ul> <li>scheduler_constant_optimizer_adam_size_15_L_e_15_L_g_20_top_k_1_lamb_0.5</li> <li>CODA-Prompt: ep_10_milestone_4_lr_0.0025_lr_decay_0.3_batch_64_w_decay_0</li> </ul>	
<ul> <li>scheduler_constant_optimizer_adam_e_pool_size_200_e_p_length_4_ortho_mu_0.001</li> <li>Adam: ep_5_milestone_2_lr_0.05_lr_decay_0.5_batch_16_w_decay_0.0005</li> </ul>	
<ul> <li>Ranpac: ep_15_milestone_4_lr_0.04_lr_decay_0.1_batch_128_w_decay_0.0001</li> </ul>	
<ul> <li>EASE: ep_10_milestone_4_lr_0.01_lr_decay_0.1_batch_16_w_decay_0.001</li> <li>scheduler_constant_optimizer_sgd_ffn_num_32_alpha_0.05</li> </ul>	
Best hyperparameters (CUB-100-1, 10 Tasks) The following represents the best hyperparameters	5
of each algorithm selected in the hyperparameter tuning phase using CUB-100-1 (10 Tasks).	
<ul> <li>b21: cp_25_initestone_2_in_0.0025_in_decay_0.5_batch_128_w_decay_0</li> <li>scheduler_cosine_optimizer_adam_size_20_length_4_top_k_6_lamb_0.5</li> <li>DuelPrempt: on 25_milestone_4_in_0.002275_in_decay_0.2_batch_128_w_decay_0.0005</li> </ul>	
<ul> <li>Dual Flohipt: ep_25_intestone_4_in_0.002375_in_decay_0.5_bach_128_w_decay_0.0003 scheduler_cosine_optimizer_adamw_size_15_L_e_30_L_g_15_top_k_1_lamb_0.5</li> <li>CODA Brownth on 25_milestone_2 in 0.001275_in_decay_0.1 hoteh 64_w_decay_0.0001</li> </ul>	
• CODA-Prompt: ep_25_initestone_5_ir_0.001575_ir_decay_0.1_batch_04_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_50_e_p_length_4_ortho_mu_0	-
• Adam: ep_25_milestone_2_lr_0.05_lr_decay_0.3_batch_24_w_decay_0.0001 scheduler_steplr_optimizer_sgd_ffn_num_32	
• Ranpac: ep_25_milestone_3_lr_0.02_lr_decay_0.3_batch_16_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_16_M_10000_pt_num_20	
• EASE: ep_15_milestone_3_lr_0.03_lr_decay_0.5_batch_128_w_decay_0.0001 scheduler_constant_optimizer_sgd_ffn_num_64_alpha_0.05	
<b>Best hyperparameters (ImageNet-R-1, 20 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R-1 (20 Tasks).	-
<ul> <li>L2P: ep_15_milestone_2_lr_0.002375_lr_decay_0.5_batch_48_w_decay_0 scheduler_cosine_optimizer_adamw_size_25_length_6_top_k_10_lamb_0.3</li> </ul>	
• DualPrompt: ep_20_milestone_3_lr_0.001875_lr_decay_0.1_batch_128_w_decay_0 scheduler_steplr_optimizer_adam_size_10_L_e_30_L_g_30_top_k_1_lamb_0.1	
<ul> <li>CODA-Prompt: ep_15_milestone_2_lr_0.002375_lr_decay_0.5_batch_64_w_decay_0.0001 scheduler_cosine_optimizer_adamw_e_pool_size_100_e_p_length_4_ortho_mu_0.01</li> </ul>	L
• Adam: ep_25_milestone_3_lr_0.04_lr_decay_0.3_batch_24_w_decay_0.001 scheduler_constant_optimizer_sgd_ffn_num_16	
• Ranpac: ep_10_milestone_2_lr_0.02_lr_decay_0.3_batch_8_w_decay_0.0001 scheduler_cosine_optimizer_sgd_ffn_num_8_M_20000_pt_num_10	
• EASE: ep_15_milestone_4_lr_0.03_lr_decay_0.5_batch_16_w_decay_0.0005 scheduler_constant_optimizer_sgd_ffn_num_64_alpha_0.05	
<b>Best hyperparameters (ImageNet-R-1, 10 Tasks)</b> The following represents the best hyperparameters of each algorithm selected in the hyperparameter tuning phase using ImageNet-R-1 (10 Tasks).	- )

1350	L 2D
1351	scheduler_cosine_optimizer_adamy_size_20_length_8_top_k_10_lamb_0.5
1352	scheduler_cosine_optimizer_adamw_size_zo_iengui_6_top_k_10_tamo_0.5
1353	DualPrompt: ep_25_milestone_4_lr_0.000875_lr_decay_0.3_batch_64_w_decay_0
1354	scheduler_cosine_optimizer_adam_size_15_L_e_30_L_g_20_top_k_1_lamb_0.1
1355	CODA-Prompt: ep_15_milestone_3_lr_0.001375_lr_decay_0.5_batch_64_w_decay_0.0005
1356	scheduler_constant_optimizer_adamw_e_pool_size_300_e_p_length_4_ortho_mu_0.01
1357	Adam: ep 25 milestone 3 lr 0.04 lr decay 0.3 batch 24 w decay 0.001
1358	scheduler_constant_optimizer_sgd_ffn_num_16
1359	Rannac: en 25 milestone 4 lr 0.05 lr decay 0.1 batch 24 w decay 0.0001
1360	scheduler constant optimizer sgd ffn num 64 M 20000 pt num 10
1361	$\Sigma A \Sigma \Sigma a = 10$ milastana $A = 0.04$ in decay 0.1 betab 24 w decay 0.001
1362	scheduler_constant_ontimizer_sgd_ffn_num_64_alpha_0.2
1363	seneduler_constant_optimizer_sgd_nn_nunn_04_arpna_0.2
1364	
1365	
1366	
1367	
1368	
1360	
1370	
1370	
1372	
1372	
1373	
1375	
1375	
1377	
1378	
1370	
1375	
1381	
1382	
1382	
1384	
1385	
1386	
1387	
1388	
1389	
1300	
1391	
1392	
1393	
1394	
1395	
1396	
1397	
1398	
1399	
1400	
1401	
1402	
1403	

# 1404 C ADDITIONAL EXPERIMENTAL RESULTS ON THE EVALUATION PHASE

# 1406 C.1 RESULT TABLES

### 1408 Class-IL without a pretrained model ( $D^{HT}$ = ImageNet-100-1)

Table 6: The experimental results of class-IL without a pretrained model (using original hyperparameters) The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^{HT}$ = ImageNet-100
Replay	41.21(1.06) / 59.82(1.48)
iCaRL	40.50(1.19) / 60.12(1.41)
BiC	39.61(2.39) / 64.27(1.59)
WA	53.34(1.39) / 68.92(1.54)
PODNet	46.66(1.11) / 64.13(1.20)
DER	61.96(1.04) / 72.10(1.41)
FOSTER	60.68(0.71) / 69.97(1.70)
BEEF	NaN
MEMO	59.59(1.29) / 70.04(1.62)

Table 7: The experimental results of class-IL without a pretrained model (using  $D^{HT}$  = ImageNet-100-1) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

	10 Tasks (Acc / AvgAcc)	$D^{HT}$ = ImageNet-100-1	$D^E = $ ImageNet-100-2
	Replay	44.78(1.19) / 59.85(0.95)	44.27(1.05) / 61.49(0.87)
	iCaRL	42.58(1.06) / 61.27(1.26)	42.44(1.50) / 63.39(1.18)
ĺ	BiC	54.22(1.27) / 67.31(0.74)	58.77(0.96) / 71.81(1.42)
	WA	54.67(0.60) / 69.54(1.41)	59.89(1.18) / 72.93(1.94)
	PODNet	55.35(0.93) / 68.74(1.52)	57.48(0.94) / 71.76(1.62)
	DER	63.31(0.42) / 72.93(0.87)	70.23(0.46) / 77.12(1.20)
ĺ	FOSTER	58.36(0.85) / 71.99(0.98)	61.46(0.98) / 68.41(1.23)
	BEEF	NaN	NaN
ĺ	MEMO	57.91(0.54) / 71.25(1.41)	61.94(0.78) / 71.35(2.17)

Table 8: The experimental results of class-IL without a pretrained model (using  $D^{HT}$  = ImageNet-100-1) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

6 Tasks (Acc / AvgAcc)	$D^{HT}$ = ImageNet-100	$D^E = $ ImageNet-200
Replay	42.93(2.41) / 53.81(1.72)	43.26(1.38) / 49.28(0.53)
iCaRL	46.62(1.54) / 57.27(0.73)	45.64(1.49) / 59.18(0.54)
BiC	37.14(1.62) / 36.42(1.89)	38.43(2.53) / 40.89(3.07)
WA	58.72(1.02) / 65.58(1.55)	60.58(1.35) / 69.47(1.71)
PODNet	67.22(0.67) / 75.05(1.16)	65.51(1.83) / 75.82(1.03)
DER	72.20(0.51) / 77.68(1.08)	75.83(0.64) / 81.19(0.70)
FOSTER	69.48(0.50) / 74.59(1.18)	71.62(1.08) / 78.29(1.14)
BEEF	74.67(0.14) / 78.92(0.54)	75.09(0.29) / 81.31(0.50)
MEMO	59.91(0.87) / 67.22(1.63)	62.80(3.16) / 68.77(6.26)

#### 1458 Class-IL without a pretrained model ( $D^{HT}$ = CIFAR-50-1)

Table 9: The experimental results of class-IL without a pretrained model (using  $D^{HT}$  = CIFAR-50-1) in the hyperparameter tuning phase.) The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^E = CIFAR-50-2$	$D^E$ = ImageNet-50-2
Replay	45.42(2.19) / 65.88(1.97)	42.51(0.47) / 60.72(1.58)
iCaRL	47.12(2.80) / 66.71(2.07)	42.44(1.00) / 61.55(1.64)
BiC	52.83(2.83) / 69.16(2.30)	49.52(1.16) / 67.09(1.74)
WA	54.89(2.13) / 69.85(2.32)	53.64(1.47) / 67.75(1.90)
PODNet	51.20(1.76) / 69.47(0.13)	51.70(1.19) / 67.86(1.67)
DER	63.51(1.98) / 75.04(1.24)	63.40(1.02) / 72.67(1.62)
FOSTER	60.00(2.72) / 72.29(2.09)	62.09(1.83) / 70.24(1.50)
BEEF	57.24(1.48) / 72.26(2.05)	NaN
MEMO	60.72(2.41) / 73.78(1.99)	54.91(1.59) / 68.06(2.10)

1475<br/>1476<br/>1477Table 10: The experimental results of class-IL without a pretrained model (using  $D^{HT}$  = CIFAR-50-1)<br/>in the hyperparameter tuning phase.) The values in parentheses represent the standard deviation.

6 Tasks (Acc / AvgAcc)	$D^E = CIFAR-50-2$	$D^E = $ ImageNet-50-2
Replay	48.00(1.98) / 59.86(1.03)	46.30(1.31) / 55.67(0.64)
iCaRL	46.09(1.51) / 59.14(1.39)	46.21(1.72) / 57.79(1.06)
BiC	58.22(1.20) / 68.16(1.96)	46.26(3.26) / 59.07(3.87)
WA	61.37(1.02) / 70.56(0.51)	61.47(0.72) / 69.67(0.63)
PODNet	62.62(0.39) / 72.62(0.75)	64.30(0.78) / 73.56(1.01)
DER	67.98(1.34) / 75.88(0.78)	70.68(0.75) / 76.56(0.95)
FOSTER	66.45(0.55) / 73.93(0.77)	69.86(0.45) / 75.27(0.83)
BEEF	65.51(1.29) / 72.98(0.50)	NaN
MEMO	64.64(1.54) / 73.50(0.83)	51.40(3.39) / 62.11(3.33)

1491 Class-IL without a pretrained model ( $D^{HT}$  = ImageNet-50-1)

Table 11: The experimental results of class-IL without a pretrained model (using  $D^{HT}$  = ImageNet-50-1) in the hyperparameter tuning phase.) The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^E$ = ImageNet-50-2	$D^E = CIFAR-50-2$
Replay	43.71(0.81) / 58.75(1.60)	44.19(2.17) / 63.57(1.50)
iCaRL	39.41(1.46) / 59.51(1.70)	41.59(3.10) / 62.42(2.85)
BiC	51.26(1.39) / 65.33(2.48)	51.22(3.67) / 66.41(2.92)
WA	51.85(0.79) / 67.23(1.79)	57.72(1.92) / 71.39(2.00)
PODNet	51.31(1.24) / 67.28(1.53)	48.19(1.17) / 65.77(1.29)
DER	64.89(1.16) / 74.15(1.56)	63.64(1.32) / 75.32(1.21)
FOSTER	61.57(0.70) / 72.38(1.20)	58.64(2.15) / 72.89(1.81)
BEEF	NaN	NaN
MEMO	57.56(1.24) / 68.36(2.27)	58.99(1.01) / 72.43(1.81)

1512	Table 12: The experimental results of class-IL without a pretrained model (using $D^{HT}$ =
1513	ImageNet-50-1) in the hyperparameter tuning phase.) The values in parentheses represent the
1514	standard deviation.

1516	6 Tasks	DE = Imaga Nat 50.2	$D^E - CIEAR 50.2$
1517	(Acc / AvgAcc)	$D^{-} = \text{ImageIvel}{-30-2}$	$D^{-} = CIFAR-30-2$
1518	Replay	42.82(1.43) / 53.50(1.54)	42.28(0.71) / 52.18(1.31)
1519	iCaRL	42.47(1.73) / 54.65(1.85)	40.24(2.64) / 52.89(2.14)
1520	BiC	44.68(2.81) / 54.19(2.93)	39.65(1.32) / 49.49(1.46)
1521	WA	55.68(0.07) / 64.69(0.72)	56.14(1.99) / 64.08(1.60)
1522	PODNet	64.10(0.80) / 72.50(0.81)	61.33(0.54) / 71.27(1.07)
1523	DER	70.28(0.98) / 76.14(1.00)	64.76(1.06) / 72.89(1.28)
1524	FOSTER	68.40(1.08) / 75.02(0.94)	65.31(0.26) / 73.80(0.68)
1525	BEEF	NaN	NaN
1526	MEMO	50.92(1.25) / 60.93(1.67)	50.58(2.62) / 60.66(2.65)
1527		•	

## 1528 Class-IL with a pretrained model ( $D^{HT}$ = CUB-200)

Table 13: The experimental results of class-IL with a pretrained model (using original hyperparameters) The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^{HT} = \text{CUB-200}$
L2P	72.32(0.62) / 76.82(0.30)
DualPrompt	68.74(0.54) / 74.39(0.68)
CODA-Prompt	75.19(0.33) / 80.27(0.93)
Adam	71.21(1.06) / 77.52(1.24)
Ranpac	78.27(0.57) / 83.24(0.44)
EASE	77.07(0.19) / 82.65(0.68)

Table 14: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = CUB-200) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^E = \text{ImageNet-R}$	$D^E = $ ImageNet-A
L2P	71.86(0.66) / 77.42(0.92)	45.13(1.25) / 53.57(0.92)
DualPrompt	66.33(0.42) / 73.03(0.60)	39.97(2.32) / 52.58(0.70)
CODA-Prompt	72.86(0.44) / 78.49(0.99)	51.63(0.50) / 61.00(0.47)
Adam	72.68(0.77) / 79.09(0.89)	57.03(0.47) / 66.50(1.22)
Ranpac	79.59(0.29) / 84.46(0.41)	66.14(0.40) / 73.63(1.05)
EASE	61.96(0.06) / 67.74(0.67)	49.32(0.48) / 58.30(0.86)

Table 15: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = CUB-200) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

## 

#### 1578 Class-IL with a pretrained model ( $D^{HT}$ = ImageNet-R)

Table 16: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = ImageNet-R) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

20 Tasks (Acc / AvgAcc)	$D^E = \text{CUB-200}$	$D^E = $ ImageNet-A
L2P	63.76(1.81) / 76.59(1.48)	36.97(1.31) / 46.78(0.71)
DualPrompt	68.78(0.78) / 79.67(1.04)	47.54(0.79) / 55.91(0.84)
CODA-Prompt	67.92(2.11) / 79.65(1.93)	50.07(0.29) / 59.76(0.58)
Adam	85.38(0.19) / 90.87(0.90)	53.86(1.44) / 63.99(2.61)
Ranpac	89.86(0.22) / 93.44(0.78)	38.53(31.11) / 67.65(3.37)
EASE	79.89(1.22) / 87.58(1.19)	53.99(1.05) / 64.11(0.78)

Table 17: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = ImageNet-R) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^E = \text{CUB-200}$	$D^E$ = ImageNet-A
L2P	69.75(1.79) / 79.92(1.24)	43.50(0.99) / 50.06(1.18)
DualPrompt	71.74(1.01) / 82.22(1.10)	39.47(0.79) / 50.63(0.94)
CODA-Prompt	72.30(1.11) / 83.00(1.35)	52.39(0.38) / 61.87(1.01)
Adam	85.90(0.17) / 90.93(0.89)	56.63(0.78) / 65.94(1.45)
Ranpac	89.99(0.29) / 93.36(0.83)	63.78(1.52) / 71.70(1.88)
EASE	74.00(0.78) / 83.69(0.74)	54.76(1.36) / 66.14(1.65)

#### 1620 Class-IL with a pretrained model ( $D^{HT}$ = CUB-100-1)

Table 18: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = CUB-100-1) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

20 Tasks (Acc / AvgAcc)	$D^E = \text{CUB-100-2}$	$D^E$ = ImageNet-R-2	ImageNet-A-2
L2P	54.12(3.59) / 68.33(3.73)	66.01(0.74) / 72.17(1.04)	28.08(2.38) / 39.18(2.75)
DualPrompt	59.83(1.63) / 73.54(2.68)	65.51(0.32) / 71.58(0.68)	33.90(2.26) / 44.84(2.25)
CODA-Prompt	58.16(1.88) / 71.05(2.68)	66.73(0.61) / 73.06(0.46)	30.62(0.82) / 41.70(1.70)
Adam	85.95(0.08) / 90.56(0.24)	67.77(0.84) / 74.53(1.74)	43.93(0.09) / 55.63(2.69)
Ranpac	89.52(0.35) / 90.52(2.96)	74.53(0.28) / 79.80(0.81)	30.30(22.41) / 45.87(4.57)
EASE	85.19(0.49) / 89.91(0.74)	67.17(0.29) / 73.61(0.75)	44.11(0.29) / 55.42(2.83)

Table 19: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = CUB-100-1) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^E = \text{CUB-100-2}$	$D^E = $ ImageNet-R-2	ImageNet-A-2
L2P	66.15(1.41) / 76.68(1.49)	70.11(0.53) / 75.61(0.87)	34.96(0.92) / 44.98(2.26)
DualPrompt	67.20(2.59) / 78.28(1.68)	68.29(0.49) / 74.32(0.89)	38.43(1.52) / 49.15(2.43)
CODA-Prompt	68.37(2.71) / 78.93(2.57)	70.35(0.81) / 75.59(0.90)	37.23(1.87) / 47.48(1.85)
Adam	86.76(0.21) / 90.75(0.46)	72.73(0.27) / 79.42(0.59)	44.81(0.85) / 55.08(2.22)
Ranpac	90.60(0.36) / 93.08(0.65)	80.40(0.3) / 85.00(0.47)	49.56(2.52) / 57.60(1.96)
EASE	85.86(0.10) / 90.11(0.26)	63.36(0.03) / 69.36(0.95)	43.88(0.15) / 54.49(2.64)

## **Class-IL** with a pretrained model $(D^{HT} = \text{ImageNet-R-1})$

Table 20: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = ImageNet-R-1) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

20 Tasks (Acc / AvgAcc)	$D^E = \text{ImageNet-R-2}$	$D^E = \text{CUB-100-2}$	ImageNet-A-2
L2P	66.15(0.85) / 71.93(1.13)	51.04(1.45) / 66.04(1.71)	25.13(2.27) / 34.21(2.51)
DualPrompt	65.77(0.78) / 71.83(1.17)	57.13(3.40) / 71.15(2.25)	31.96(2.49) / 41.71(1.76)
CODA-Prompt	66.44(0.66) / 72.62(0.36)	57.24(1.90) / 71.27(1.95)	30.48(1.62)/41.30(2.56)
Adam	70.69(0.73) / 77.86(0.51)	86.35(0.14) / 90.83(0.56)	44.25(0.86) / 55.84(2.75)
Ranpac	76.15(0.93) / 81.68(0.94)	73.73(31.52) / 89.58(2.03)	35.06(15.86) / 47.04(6.07)
EASE	75.16(0.68) / 81.68(0.71)	76.36(2.61) / 84.35(2.55)	42.49(1.76) / 54.40(3.21)

Table 21: The experimental results of class-IL with a pretrained model (using  $D^{HT}$  = ImageNet-R-1) in the hyperparameter tuning phase. The values in parentheses represent the standard deviation.

10 Tasks (Acc / AvgAcc)	$D^E = \text{ImageNet-R-2}$	$D^E = \text{CUB-100-2}$	ImageNet-A-2
L2P	70.35(0.64) / 75.66(0.30)	63.71(2.33) / 74.62(1.61)	29.10(1.24) / 38.80(1.44)
DualPrompt	69.97(0.25) / 75.93(0.62)	66.66(1.12) / 78.11(1.43)	32.42(0.68) / 42.31(2.02)
CODA-Prompt	72.17(0.46) / 77.80(0.50)	66.98(1.3) / 78.70(0.98)	37.04(1.49) / 46.47(2.45)
Adam	72.84(0.67) / 79.69(0.86)	85.26(0.41) / 89.77(0.45)	37.36(2.72) / 48.62(4.07)
Ranpac	80.70(0.50) / 85.28(0.46)	91.09(0.51)/91.63(3.51)	41.98(19.61) / 58.79(4.70)
EASE	78.33(0.41) / 83.82(0.71)	79.70(1.47) / 86.23(1.59)	42.49(0.69) / 53.69(2.61)

## 1674 C.2 TRAINING GRAPHS







<sup>1782</sup> Class-IL with a pretrained model ( $D^{HT}$  = ImageNet-R-1,  $D^E$  = ImageNet-R-2)

#### D LIMITATIONS AND FUTURE WORK

Our study has limitations. First, evaluating each algorithm using the proposed evaluation protocol requires a substantial number of training trials. Although we believe that our protocol serves as a basic method for more accurately assessing CL algorithms, it is not a perfect evaluation protocol. Consequently, developing more efficient protocols for accurately evaluating CL algorithms remains a

significant and interesting research direction.

Second, we did not account for unpredictable CL scenarios, such as varying task numbers or class distributions. Our study assumes CL scenarios are predictable, but each phase's dataset differs. This is because, in real-world situations, some level of predictability is possible, and evaluating algorithms in completely unpredictable scenarios would be too harsh. Nevertheless, we believe that it is essential to explore evaluation methods for unpredictable scenarios in broader CL research, potentially through adaptive algorithms that can adjust hyperparameters for each task.

Finally, our evaluation focused solely on offline class-incremental learning algorithms. We think that similar challenges associated with the conventional evaluation protocol also exist in other CL domains, such as online class-incremental learning, class-incremental semantic segmentation, continual self-supervised learning, and continual reinforcement learning. As part of our future work, we intend to first apply the proposed protocol to online class-incremental learning algorithms, followed by its implementation in other domains.