TOWARDS REALISTIC HYPERPARAMETER OPTIMIZATION IN CONTINUAL LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

In continual learning (CL)—where a learner trains on a stream of data—standard hyperparameter optimisation (HPO) cannot be applied, as a learner does not have access to all of the data at the same time. This has prompted the development of CL-specific HPO frameworks. The most popular way to tune hyperparameters in CL is to repeatedly train over the whole data stream with different hyperparameter settings. However, this *end-of-training* HPO is unrealistic as in practice a learner can only see the stream once. Hence, there is an open question: *what HPO framework should a practitioner use for a CL problem in reality?* This paper answers this question by comparing several realistic HPO frameworks. We find that none of the HPO frameworks considered, including end-of-training HPO, perform consistently better than the rest on popular CL benchmarks. We therefore arrive at a twofold conclusion: a) on the popular CL benchmarks examined, a CL practitioner should select the HPO framework based on other factors, for example compute efficiency and b) to be able to discriminate between HPO frameworks there is a need to move beyond the current most commonly used CL benchmarks.

1 INTRODUCTION

027 028

025 026

000

001

002 003 004

010 011

012

013

014

015

016

017

018

019

021

029 Sequentially updating deep learning systems on a non-stationary data stream is a challenging problem which *continual learning* (CL) methods aim to address. The standard CL setup is when a learner 031 sees a sequence of tasks one-by-one and at the end of learning is evaluated on how well it performs 032 across all tasks. There have been many methods (Delange et al., 2021; Parisi et al., 2019; Wang 033 et al., 2023) designed for this problem and CL scenarios proposed (Hsu et al., 2018; Antoniou et al., 034 2020; van de Ven & Tolias, 2019). A key decision when using a CL method is selecting hyperparameter settings—learning rates, regularisation coefficients, etc. (Feurer & Hutter, 2019; Delange 035 et al., 2021; Wistuba et al., 2023). The most common way to fit hyperparameters for CL is end-of-036 training hyperparameter optimisation (HPO) (Delange et al., 2021; Buzzega et al., 2020)—shown in 037 Figure 1. This is when the hyperparameters are fit by training over the whole data stream with each hyperparameter configuration and then selecting the configuration that has the best end-of-training performance on a held-out validation set. However, end-of-training HPO is unrealistic as in the real 040 world a learner can only train over the data stream once and must select hyperparameters only using 041 the data it can currently access. Therefore, determining the best *realistic* way to perform HPO for 042 CL is currently an open problem. 043

In this work, we address the problem of deciding what realistic HPO framework to use in CL. To do 044 this, we benchmark a variety of approaches for performing HPO across different CL methodologies (ER (Chaudhry et al., 2020), ER-ACE (Caccia et al., 2021), iCaRL (Rebuffi et al., 2017), ESMER 046 (Sarfraz et al., 2023) and DER++ (Buzzega et al., 2020)). We investigate both fixed HPO frameworks 047 where the hyperparameters are kept constant throughout training and dynamic HPO frameworks 048 where hyperparameters are adapted throughout learning. For fixed HPO we examine (i) end-oftraining HPO as well as (ii) a first-task HPO framework where we fit the hyperparameters only using data from the first task (see Figure 1), a realistic and computationally efficient method. For 051 dynamic HPO, we consider (i) using data from the current task, (ii) using data stored in memory, and (iii) using validation sets from previous tasks to perform HPO for each new task. By comparing these 052 different HPO frameworks we shed light on what validation signal is sufficient to fit hyperparameters in CL and whether hyperparameters need to be adapted during training.



Our experiments highlight that to be able to better compare and develop CL HPO frame works there is a need to move beyond the current most popular CL benchmarks.

108 2 PRELIMINARIES AND RELATED WORK

110

111 CL is a large research area where many different settings have been looked at. In this work we look 112 at the most common CL setting which is known as standard CL, or sometimes offline CL (Prabhu 113 et al., 2020). In *Standard CL*, the learner sees a non-stationary sequence of data chunks called *tasks* one-by-one, such that it only has access to one chunk at a time and cannot access previously seen 114 or future chunks. Each task consists of examples which are data instance and label pairs (e.g. pairs 115 of images and their class) sampled from a subset of the classes. For example, the first task might be 116 examples of cows and sheep and the second task could be formed of examples of dogs and cats. The 117 goal of the learner is to classify new examples accurately after training on the whole data stream. 118 There are two common ways to evaluate a CL learner, task and class incremental learning. Task-119 *incremental* learning is when, at test-time, the learner knows which task a data instance comes from 120 and so only needs to distinguish between classes within that task. While, *class-incremental* learning 121 is when the learner is not given what task a data instance belongs to at test time and must distinguish 122 between all classes from all the tasks. An important part of the standard CL setting is the assumption 123 of memory constraints, which is why a learner cannot solve CL by storing previous data chunks in 124 memory. The memory constraints take the form of only allowing a learner to store a small amount 125 of previous data in memory and in constraining its use of memory for storing additional networks or parts of networks (Delange et al., 2021; Wang et al., 2023). 126

127 There have been many methods proposed for CL (Delange et al., 2021; Parisi et al., 2019; Wang 128 et al., 2023). One of the most popular and performant approaches to standard CL are replay meth-129 ods (Wang et al., 2023). This is especially true for class-incremental learning, where they are com-130 monly the best performing methods (van de Ven & Tolias, 2019; Wu et al., 2022; Mirzadeh et al., 131 2020; Lee & Storkey, 2024). Replay methods use a memory buffer to store a set of examples from previous tasks to regularise the updates on new tasks such that the learner does not forget previous 132 task knowledge. For example, the stereotypical replay method is *experience replay* (ER) (Chaudhry 133 et al., 2020; 2019b; Aljundi et al., 2019a) which for each learning step appends a sample of data from 134 the replay buffer to the batch of current task data to be trained on. More complex replay methods 135 often use a form of knowledge distillation on a sample of data from the replay buffer. For example, 136 DER++ (Buzzega et al., 2020), ESMER (Sarfraz et al., 2023) and iCaRL (Rebuffi et al., 2017) are 137 replay methods which use a method-specific knowledge distillation term. For each of these methods 138 the most common hyperparameters that are tuned are the learning rate and regularisation coeffi-139 cients, which need to be tuned to get good performance (see Appendix B). While other potential 140 hyperparameters are often not tuned in CL, e.g. momentum (Buzzega et al., 2020). 141

While the most common HPO framework used in standard CL is end-of-training HPO, there have 142 been several other HPO frameworks suggested (Kilickaya & Vanschoren, 2023; Parisi et al., 2019; 143 Cai et al., 2021). For example, Delange et al. (2021) propose a dynamic HPO framework. The 144 method adapts the hyperparameters for each task by first training with the hyperparameter config-145 uration which is assumed to have the least impact on previous task performance. Then the method 146 incrementally changes hyperparameter values to improve performance on the current task to a pre-147 specified value, while decreasing performance on previous tasks. However, this method assumes that the direction to change hyperparameters to increase performance on the current task is known 148 and that the interaction between different hyperparameters is understood. In this work we look at 149 a similar HPO framework, current-task HPO, which does not need the above assumptions. Also, 150 for the online CL scenario-which is different to standard CL-another HPO framework has been 151 proposed whereby end-of-training HPO is used on the first (or first k) tasks and then the hyper-152 parameters are fixed after that (Chaudhry et al., 2019a). To the best of our knowledge, this HPO 153 framework has been rarely used in standard CL up to this point. Here, we look at it in the form 154 of the first-task HPO framework and examine how it performs in the commonly used standard CL 155 setting. There has also been work on making dynamic HPO frameworks more efficient by sampling 156 fewer HPO configurations, for example using bandit methods (Liu et al., 2023) and analysis of vari-157 ance techniques (Semola et al., 2024). However, for simplicity, we only look at the more expensive 158 dynamic HPO frameworks which are an upper bound to the performance of these more efficient methods. While as shown above there has been work on HPO for CL, to the best of our knowledge 159 there has not been a comprehensive comparison between the main HPO frameworks proposed. This 160 is one of the key contributions of this work, shedding light on the relative performances of HPO 161 frameworks for CL.



175 Figure 2: Depiction of current-task, seen-tasks (Mem) and seen-tasks (Val) HPO frameworks, which 176 dynamically adapt hyperparameters (HPs) for each task. Each methods splits the data of the current 177 task into train and validation sets. Then, current-task HPO uses this validation set to fit the HPs for 178 the current task. While, seen-tasks (Mem) and seen-tasks (Val) use a combination of this validation 179 set and either a sample of data from previous tasks stored in memory or validation sets of previous tasks, respectively. Then current-task and seen-tasks (Mem) HPO retrain on the combined validation and train sets to complete the learning process on that task. Seen-tasks (Val) does not retrain, instead 181 it takes the model fitted using the best found hyperparameters as the final model for the current 182 task. This is to ensure that the current task's validation set has not been trained on when fitting 183 hyperparameters for future tasks. 184

3 STANDARD CL

187 188

While the setting we look at, standard CL, is mentioned above, we describe it more formally here. 189 In standard CL a learner sees a sequence of tasks, D_1, \ldots, D_T , where each task consists of a chunk 190 of data. The chunks of data consist of a set of examples, where an example is a pair formed of 191 a data instance $\mathbf{x} \in X$ and label $y \in C$. Each task only contains examples from a given subset 192 of the classes, in other words for all $(\mathbf{x}, y) \in D_i$ we have that $y \in C_i$ and $C_i \subseteq C$ is the subset 193 of classes the examples of that task can belong to. In this work we look at the most common 194 setting, where no two tasks have examples from the same class. This means that for any two task 195 i and j we have that $C_i \cap C_j = \emptyset$. Additionally, learners can have a memory buffer of previous 196 examples which consists at task i of the set M_i . Training consists of the learner sequentially seeing 197 each task in order and it cannot access the data from previous or future tasks. For each task, its data chunk is split into training and validations sets, $Train_i \subseteq D_i$ and $Val_i \subseteq D_i$, to enable the use of HPO frameworks. Then after fitting the hyperparameters the learner usually retrains on the 199 combination of the training and validation sets, $D_i = \text{Train}_i \cup \text{Val}_i$. After training the learner is 200 tested by evaluating its performance on a held-out set of data which consists of an equal number 201 of examples from all the classes. We look at two evaluation scenarios, task-incremental learning 202 and class-incremental learning. Task-incremental learning is where the learner receives with each 203 test data instance the task it belongs to and therefore the subset of classes that the data instance 204 can belong to. While for *class-incremental* learning, no indication is given of what task a test data 205 instance belongs to. 206

207 208

209

4 HPO FRAMEWORKS FOR CL

In this work, we examine several HPO frameworks for CL to see which should be the preferred choice to use in CL. We look both at fixed HPO frameworks which keep the values of hyperparameters constant throughout training and dynamic HPO which adapts the hyperparameters per task.
The fixed HPO frameworks we look at are end-of-training HPO and first-task HPO and the dynamic HPO frameworks we look at are current-task HPO, seen-tasks HPO (Mem) and seen-tasks HPO (Val). Each of these frameworks are described in turn below and we present an overview of their advantages and disadvantages in Table 1.

Table 1: Advantages and disadvantages of different HPO frameworks. Where, for time complexity, *K* refers to the number of hyperparameter configurations looked at and *T* is the number of tasks in the data stream. The asterisk (*) for seen-tasks HPO (Val) denotes that, while it does not require knowledge of future tasks like end-of-training HPO, it does require additional storage compared to other methods. The additional memory is needed to store the validation sets of previous tasks.

HPO Framework	Realistic?	Efficient? (Time Complexity)
End-of-training HPO	×	$\checkmark (\mathcal{O}(T \times K))$
First-task HPO	1	$\checkmark (\mathcal{O}(T+K))$
Current-task HPO	1	$\checkmark (\mathcal{O}(T \times K))$
Seen-tasks HPO (Val)	√ *	$\checkmark (\mathcal{O}(T \times K))$
Seen-tasks HPO (Mem)	✓	$\checkmark (\mathcal{O}(T \times K))$

228 229

222

230 End-of-training HPO is the most common HPO framework for CL (shown in Figure 1). It selects 231 hyperparameters by first training each hyperparameter configuration on the whole data stream. Sec-232 ond, it evaluates the final model fitted using each hyperparameter configuration on a validation set 233 formed of each task's held-out validation set, and selects the configuration with the highest valida-234 tion performance. Last, it retrains using the selected configuration on the whole data stream where 235 the validation data for each task is added to the training data. The model fitted at the end of this 236 training run is the final model to be evaluated. This HPO framework is expensive as it needs to 237 perform a training run over all the data stream for each hyperparameter configuration looked at. Additionally, it is unrealistic as it requires running through the data stream multiple times, which is not 238 possible in many real-world settings. It might be thought that to make end-of-training more realistic 239 the learner could store a network for each hyperparameter configuration: updating each network on 240 every task and performing selection at the end of training. This idea would remove the requirement 241 of running through the data stream multiple times. However, it would also require a large amount of 242 extra memory. Additionally, the learner would have to store and not train on the validation data for 243 each previous task. Therefore, because of underlying constraints on memory usage in standard CL, 244 it is not possible to use such an idea. 245

First-task HPO is a fixed HPO framework which is illustrated in Figure 1. It selects hyperpa-246 rameters by training each hyperparameter configuration on the *first task*. Next, it measures the 247 performance of each configuration on the held-out validation set of the first task. The configuration 248 with the highest validation accuracy is then used to retrain on the first task using both the training 249 and validation data and thereafter for all of the future tasks. First-task HPO is computationally effi-250 cient as it trains using each hyperparameter configuration solely on the first task and then only trains 251 using one configuration for the rest of the tasks. This is much less costly than end-of-training HPO, 252 which for all tasks must train using each hyperparameter configuration. Additionally, first-task HPO 253 can be used in real-world settings as it only assumes access to data available at the start of training, 254 the first task, and not future tasks like end-of-training HPO.

Current-task HPO is a dynamic HPO framework which selects hyperparameters for each task using the validation set of the *current task* (shown in Figure 2). This is a greedy strategy, selecting the hyperparameters that maximise the validation performance of the current task. It is roughly as computationally expensive as end-of-training HPO, as it has to validate each hyperparameter configuration for each task. However, it is more realistic than end-of-training HPO as it only needs access to the current task's data.

261 Seen-tasks HPO (Mem) and seen-tasks HPO (Val) are dynamic HPO frameworks (shown in 262 Figure 2). They select hyperparameters for each task by a validation set formed of current task 263 validation data along with some historic data from the stream. We consider two ways to integrate 264 historic task data. Seen-tasks HPO (Mem) uses a sample of data from the current memory buffer. 265 Seen-tasks HPO (Val) uses the validation sets of previous tasks. So, unlike current-task HPO, the 266 hyperparameters are fit using both current and previous task data. This should aid the HPO procedure in selecting hyperparameters that ensure previous tasks are not forgotten. Like current task 267 HPO, both seen-tasks HPO (Mem) and seen-tasks HPO (Val) are as computationally expensive as 268 end-of-training HPO. Seen-tasks HPO (Val) assumes it is possible to access the validation sets of 269 previous tasks which makes it less realistic than current or first task HPO. This is unlike seen-tasks

HPO (Mem) which does not assume this as it uses data stored in the memory buffer to measure performance on the previous tasks. But, this comes at the cost of biasing its validation performance as the data in the memory buffer has been trained on in previous tasks.

For seen-tasks HPO (Mem), three additional details are important to mention. First, to ensure we are not training on validation data, the sample from memory used in the validation set is not trained on for the current task. Second, as the memory buffer contains different amounts of data for each task, we sample the same proportion of examples from each task to add to the validation set. Last, unlike for the other HPO frameworks, the validation set combined with the sample from memory might be class imbalanced. Therefore, unlike other methods which use validation accuracy as the performance metric, for seen-tasks HPO (Mem) we use the median of per class accuracies to reduce the impact of class imbalance.

281 282

283 284

5 EXPERIMENTS

285 **Benchmarks** In our experiments we look at two settings, the commonly used split-task setting 286 (Buzzega et al., 2020; Delange et al., 2021) and the heterogeneous task setting. We look at these 287 settings using the datasets CIFAR-10, CIFAR-100, CORe50 and Tiny ImageNet (Krizhevsky, 2009; Lomonaco & Maltoni, 2017; Wu et al., 2015). We chose to use these datasets and the split-task 288 setting due to their commonplace use in the CL literature (Wang et al., 2023) and hence to maximise 289 the insights our results can have on current practice. In the split-task setting, each task has the same 290 number of classes associated with it and no two tasks share a class. For CIFAR-10 and CORe50, the 291 dataset is split into five tasks, each containing the data from two or ten of the classes, respectively. 292 For CIFAR-100 and Tiny ImageNet, the datasets are split into ten tasks, where each task contains 293 the data of 10 or 20 classes, respectively. In the heterogeneous task setting, instead of each task 294 having the same number of classes associated with it they have a varying amount, from two to ten, 295 but still no two tasks share a class (see Appendix A for more details). This is to make the tasks 296 have differing amounts of data and difficulty. We only look at CIFAR-100 and Tiny ImageNet for 297 the heterogeneous task setting due to computational cost. Additionally, for the heterogeneous task 298 setting we divide the datasets into twenty tasks to test how HPO frameworks perform on longer task 299 sequences. For both settings, if required by the HPO framework, we split the data of the task into train and validations sets, where the validation set contains 10% of the task's data evenly sampled 300 from each class associated with the task. 301

We evaluate the methods at the end of training using a standard performance metric for CL, average accuracy (Chaudhry et al., 2019a). The average accuracy of a method is the mean accuracy over each task on a held-out test set which contains an equal amount of data from each task. For classincremental learning, the learner must classify between all classes at test time as it is not told what task a test data instance comes from. For task-incremental learning, the learner knows what task each test data instance comes from, meaning only classes from that task will be predicted.

308 **CL methods** To evaluate how well each HPO framework performs we look at applying them to 309 fit the hyperparameters of several common and well performing CL methods. More specifically, we utilise the CL methods: ER (Chaudhry et al., 2020), ER-ACE (Caccia et al., 2021), iCaRL (Rebuffi 310 et al., 2017), ESMER (Sarfraz et al., 2023) and DER++ (Buzzega et al., 2020). For these methods we 311 fit the learning rate and any regularisation coefficients they have using each HPO framework. While 312 all HPO frameworks looked at can be used with any underlying sampler/selector of hyperparameter 313 configurations, for simplicity and to be consistent with common practice in CL (Buzzega et al., 314 2020; Boschini et al., 2022; Sarfraz et al., 2023) we use grid search. We look at the combination 315 of ten different learning rate values and for each regularisation coefficient three different values. 316 This means for DER++ we search across 90 different hyperparameter configurations (learning rate 317 and two regularisation coefficients) and for ESMER we search across 30 different configurations 318 (learning rate and the loss margin coefficient). While, for ER, iCaRL and ER-ACE we look at 10 319 different configurations as they have no regularisation coefficients to fit. The hyperparameter grid 320 used is very similar to the ones looked at in several popular works on CL (Buzzega et al., 2020; 321 Boschini et al., 2022) and is given in full in Appendix A. Moreover, for each method we use: a ResNet18 (He et al., 2016) as the underlying backbone network; random crop and horizontal flip 322 data augmentations when training; and a memory buffer of size 5120, in common with previous 323 work (Buzzega et al., 2020).

324 Table 2: Results of using different HPO frameworks for ER, iCaRL, ER-ACE, ESMER and DER++ 325 on the standard split-task CIFAR-10 and CIFAR-100 benchmarks. We report mean average accuracy 326 over three runs with their standard errors and, to highlight effect size, bold results which are greater by +0.5% average accuracy than any other for that CL method. The table shows that all HPO 327 frameworks perform similarly; none perform consistently better than the rest. 328

		CIFA	R-10	CIFA	R-100
CL Method	HPO Framework	Class-IL.	Task-IL.	Class-IL.	Task-IL.
ER	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 83.55_{\pm 0.44} \\ \textbf{84.38}_{\pm 0.45} \\ 82.10_{\pm 2.21} \\ 83.67_{\pm 0.73} \\ 79.49_{\pm 0.63} \end{array}$	$\begin{array}{c} 97.18_{\pm 0.14} \\ 96.82_{\pm 0.17} \\ 96.39_{\pm 0.50} \\ 96.84_{\pm 0.21} \\ 95.93_{\pm 0.09} \end{array}$	$\begin{array}{c} 51.03 {\scriptstyle \pm 0.43} \\ 49.61 {\scriptstyle \pm 0.34} \\ 50.64 {\scriptstyle \pm 0.40} \\ 51.46 {\scriptstyle \pm 0.36} \\ 47.39 {\scriptstyle \pm 0.24} \end{array}$	$\begin{array}{c} 85.68_{\pm 0.29} \\ 84.97_{\pm 0.19} \\ 85.47_{\pm 0.18} \\ 85.65_{\pm 0.06} \\ 84.83_{\pm 0.22} \end{array}$
iCaRL	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 77.79_{\pm 0.23} \\ 77.83_{\pm 0.22} \\ 76.15_{\pm 0.75} \\ 77.58_{\pm 0.49} \\ 76.67_{\pm 0.44} \end{array}$	$\begin{array}{c} \textbf{98.52}_{\pm 0.03} \\ 95.31_{\pm 0.12} \\ 93.29_{\pm 0.61} \\ 94.32_{\pm 1.01} \\ 95.41_{\pm 0.28} \end{array}$	$\begin{array}{c} 54.30_{\pm 0.36}\\ 52.56_{\pm 0.10}\\ 54.26_{\pm 0.02}\\ 51.89_{\pm 0.39}\\ 49.16_{\pm 0.23}\end{array}$	$\begin{array}{c} 85.74_{\pm 0.45}\\ 84.60_{\pm 0.09}\\ 85.74_{\pm 0.06}\\ 84.02_{\pm 0.68}\\ 82.43_{\pm 0.23}\end{array}$
ER-ACE	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 82.34_{\pm 0.30}\\ 83.20_{\pm 0.79}\\ \textbf{83.99}_{\pm 0.22}\\ 81.94_{\pm 1.55}\\ 81.61_{\pm 0.15}\end{array}$	$\begin{array}{c} 96.74_{\pm 0.01} \\ 96.67_{\pm 0.18} \\ 96.58_{\pm 0.15} \\ 95.90_{\pm 0.51} \\ 96.40_{\pm 0.13} \end{array}$	$\begin{array}{c} 55.58 {\pm} 0.39 \\ 56.36 {\pm} 0.29 \\ 56.46 {\pm} 0.36 \\ 54.37 {\pm} 0.25 \\ 53.76 {\pm} 0.21 \end{array}$	$\begin{array}{c} 85.73 {\scriptstyle \pm 0.09} \\ 86.11 {\scriptstyle \pm 0.154} \\ 86.35 {\scriptstyle \pm 0.02} \\ 85.02 {\scriptstyle \pm 0.14} \\ 84.56 {\scriptstyle \pm 0.31} \end{array}$
ESMER	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 80.73 {\scriptstyle \pm 0.15} \\ 77.89 {\scriptstyle \pm 0.46} \\ 81.69 {\scriptstyle \pm 0.25} \\ 81.29 {\scriptstyle \pm 0.03} \\ 70.95 {\scriptstyle \pm 0.94} \end{array}$	$\begin{array}{c} 96.50_{\pm 0.01} \\ 96.15_{\pm 0.12} \\ 96.03_{\pm 0.05} \\ 96.46_{\pm 0.06} \\ 95.79_{\pm 0.14} \end{array}$	$\begin{array}{c} 56.16_{\pm 0.54} \\ 56.61_{\pm 0.20} \\ 55.11_{\pm 0.13} \\ 53.81_{\pm 0.44} \\ \textbf{57.50}_{\pm 0.14} \end{array}$	$\begin{array}{c} 88.69_{\pm 0.35} \\ 89.05_{\pm 0.10} \\ 88.96_{\pm 0.08} \\ 87.26_{\pm 0.13} \\ 89.27_{\pm 0.16} \end{array}$
DER++	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 84.40_{\pm 0.94}\\ 85.22_{\pm 0.08}\\ 84.90_{\pm 0.11}\\ 85.44_{\pm 0.38}\\ 82.18_{\pm 0.26}\end{array}$	$\begin{array}{c} 95.75_{\pm 0.33} \\ 96.14_{\pm 0.10} \\ 95.92_{\pm 0.11} \\ 96.22_{\pm 0.15} \\ 94.75_{\pm 0.28} \end{array}$	$\begin{array}{c} 56.04_{\pm 3.67} \\ 55.20_{\pm 0.78} \\ 55.00_{\pm 1.21} \\ 56.59_{\pm 0.64} \\ 56.94_{\pm 0.66} \end{array}$	$\begin{array}{c} 83.13_{\pm 2.69} \\ 81.68_{\pm 0.66} \\ 83.14_{\pm 0.76} \\ 83.61_{\pm 0.42} \\ 83.08_{\pm 0.21} \end{array}$

357 358

359 360

5.1 RESULTS

361 362 For the split-task setting, the results of our experiments show that none of the HPO frameworks looked at perform much better than the rest. The results are presented in Table 2 and 3 and we have 364 bolded the results which are better by +0.5% than any of the other HPO frameworks results for a given CL method. The reason we chose to bold results in this way is to be able to draw attention 366 to and reference observed effect sizes. We want to do this because if the observed effect sizes are 367 small it suggests that no method performs much better than any other and hence that other factors 368 become more important when selecting a HPO framework, e.g. compute cost. In Table 2 there are 369 few bolded results and for those that exist, the HPO framework which achieves it varies. This shows 370 that, for the datasets shown in Table 2, there is only a small difference in performance between HPO 371 frameworks. While in Table 3 there are more bolded results indicating a slightly greater variance in 372 the performance of HPO frameworks-perhaps due to the greater complexity of the datasets looked 373 at. However, as in Table 2, in Table 3 the HPO framework that performs the best differs across 374 datasets and CL methods. These results show that no HPO framework performs consistently better 375 than the rest. For instance, on CIFAR-100, no HPO framework improves accuracy over the other methods by more than +0.5% for all CL methods but ESMER in class-incremental learning. This 376 suggest that for the split-task setting there is no general advantage in using one HPO framework over 377

another in terms of predictive performance.

Table 3: Results of using different HPO frameworks for ER, iCaRL, ER-ACE, ESMER and DER++ on the standard split-task CORe50 and Tiny ImageNet benchmarks. We report mean average accu-racy over three runs with their standard errors and, to highlight effect size, bold results which are greater by +0.5% average accuracy than any other for that CL method. The table shows that all HPO frameworks perform similarly; none perform consistently better than the rest.

		CORe50		Tiny ImageNet	
CL Method	HPO Framework	Class-IL.	Task-IL.	Class-IL.	Task-IL.
ER	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 37.37_{\pm 1.03} \\ 38.37_{\pm 0.38} \\ 35.97_{\pm 0.24} \\ \textbf{39.12}_{\pm 0.64} \\ 36.10_{\pm 1.15} \end{array}$	$\begin{array}{c} 55.51_{\pm 0.41} \\ 56.95_{\pm 0.62} \\ 53.40_{\pm 1.01} \\ 57.32_{\pm 0.63} \\ 54.28_{\pm 0.77} \end{array}$	$\begin{array}{c} 28.01_{\pm 0.09} \\ 28.51_{\pm 0.18} \\ 25.79_{\pm 0.21} \\ 28.45_{\pm 0.28} \\ \textbf{29.58}_{\pm 0.25} \end{array}$	$\begin{array}{c} 68.17_{\pm 0.06} \\ \textbf{68.72}_{\pm 0.13} \\ 66.96_{\pm 0.15} \\ 68.16_{\pm 0.26} \\ 68.02_{\pm 0.14} \end{array}$
iCaRL	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 54.30_{\pm 0.36}\\ 52.56_{\pm 0.10}\\ 54.26_{\pm 0.02}\\ 51.89_{\pm 0.39}\\ 49.16_{\pm 0.23}\end{array}$	$\begin{array}{c} 85.74_{\pm 0.45} \\ 84.60_{\pm 0.09} \\ 85.74_{\pm 0.06} \\ 84.02_{\pm 0.68} \\ 82.43_{\pm 0.23} \end{array}$	$\begin{array}{c} 37.09_{\pm 0.27} \\ 36.42_{\pm 0.22} \\ 37.17_{\pm 0.28} \\ 34.81_{\pm 0.42} \\ 36.79_{\pm 0.13} \end{array}$	$\begin{array}{c} 70.37_{\pm 0.36} \\ 70.11_{\pm 0.13} \\ 70.67_{\pm 0.03} \\ 68.42_{\pm 0.41} \\ 70.46_{\pm 0.08} \end{array}$
ER-ACE	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 39.33 {\scriptstyle \pm 0.79} \\ 37.81 {\scriptstyle \pm 0.71} \\ 43.59 {\scriptstyle \pm 0.09} \\ 44.32 {\scriptstyle \pm 0.69} \\ 37.60 {\scriptstyle \pm 0.69} \end{array}$	$\begin{array}{c} 58.14_{\pm 1.29} \\ 56.02_{\pm 0.60} \\ 61.33_{\pm 0.33} \\ \textbf{62.28}_{\pm 0.51} \\ 56.01_{\pm 1.17} \end{array}$	$\begin{array}{c} \textbf{38.94}_{\pm 0.47} \\ 36.94_{\pm 0.67} \\ 37.63_{\pm 0.38} \\ 36.06_{\pm 0.37} \\ 32.37_{\pm 0.34} \end{array}$	$\begin{array}{c} \textbf{70.18}_{\pm 0.23} \\ 68.16_{\pm 0.30} \\ 68.25_{\pm 0.41} \\ 67.69_{\pm 0.26} \\ 64.37_{\pm 0.47} \end{array}$
ESMER	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 45.08_{\pm1.06} \\ \textbf{47.07}_{\pm1.18} \\ 46.01_{\pm0.90} \\ 43.29_{\pm1.11} \\ 42.15_{\pm1.24} \end{array}$	$\begin{array}{c} 62.05_{\pm 0.45} \\ 63.69_{\pm 0.95} \\ 63.32_{\pm 0.59} \\ 60.77_{\pm 0.80} \\ 58.78_{\pm 1.10} \end{array}$	$\begin{array}{c} \textbf{47.33}_{\pm 0.30} \\ 46.69_{\pm 0.56} \\ 45.20_{\pm 0.53} \\ 44.82_{\pm 0.16} \\ 44.26_{\pm 0.20} \end{array}$	$\begin{array}{c} 76.18_{\pm 0.22} \\ 75.72_{\pm 0.24} \\ 74.93_{\pm 0.29} \\ 74.27_{\pm 0.11} \\ 74.54_{\pm 0.31} \end{array}$
DER++	End-of-training HPO First-task HPO Current-task HPO Seen-tasks HPO (Val) Seen-tasks HPO (Mem)	$\begin{array}{c} 51.87_{\pm 0.44} \\ 46.07_{\pm 1.58} \\ 51.58_{\pm 0.77} \\ 49.19_{\pm 0.37} \\ 41.08_{\pm 1.91} \end{array}$	$\begin{array}{c} 63.48_{\pm 0.61} \\ 58.07_{\pm 1.18} \\ \textbf{64.19}_{\pm 046} \\ 62.10_{\pm 0.65} \\ 54.73_{\pm 2.16} \end{array}$	$\begin{array}{c} \textbf{39.89}_{\pm 0.27} \\ 35.98_{\pm 0.63} \\ 36.64_{\pm 0.33} \\ 31.88_{\pm 5.36} \\ 33.54_{\pm 0.13} \end{array}$	$\begin{array}{c} \textbf{70.41}_{\pm 0.17} \\ 65.86_{\pm 0.37} \\ 66.43_{\pm 0.49} \\ 64.20_{\pm 3.00} \\ 63.68_{\pm 0.17} \end{array}$

In the heterogeneous task setting we also see that none of the HPO frameworks perform consistently better than the rest. The results for this setting are presented in Table 4 and we have again bolded the results which are better by +0.5% than any of the other HPO frameworks for a given CL method. Like the results for the split-task setting, there are many columns for each CL method which have no bolded result and for the three which do the HPO framework which achieves it is different. Therefore, we conclude that in the heterogeneous task setting it is also the case that there is no one best HPO framework. The reason we look at the heterogeneous task setting is because we expected a greater benefit from adapting hyperparameters per task, given that unlike the split-task setting each task is quite different. However, our results show that this is not the case and that it is possible to use the same hyperparameters across all the tasks and still perform well.

Performance of first-task HPO Our results show that all of the HPO frameworks tested perform similarly. Therefore, we conclude that other factors should be used when deciding what realistic HPO framework to use on these common CL benchmarks. For example, taking computational cost into account would mean that first-task HPO would be a good method to use as it is the most com-putationally efficient. Given this, we describe here in more detail its relative performance compared to the other HPO frameworks tested. In the split-task setting, we see from Table 2 and 3, that for ER some of its results are bolded. Thus, first-task HPO sometimes achieves the best performance. Ad-ditionally, for the spilt task setting, there is an average performance difference from end-of-training HPO to first-task HPO of -0.62% in class-incremental learning and -0.91% in task-incremental learning. While, for the heterogeneous tasks setting there is an average performance difference from Table 4: Results of using different HPO frameworks for ER, iCaRL, ER-ACE, ESMER and DER++ on heterogeneous task benchmarks. We report mean average accuracy over three runs with their standard errors and, to highlight effect size, bold the results which are greater by +0.5% accuracy than any other for that CL method. The table shows that no HPO framework is consistently better than the rest.

		Hetero-CIFAR-100	Hetero-TinyImg
CL Method	HPO Framework	Class-IL.	Class-IL.
	End-of-training HPO First-task HPO	$50.41_{\pm 0.21}$ 50.33 + 0.50	$39.41_{\pm 0.57}$ 40.77 + 0.24
FR	Current-task HPO	$49.77_{\pm 0.91}$	40.65 ± 0.07
Lit	Seen-tasks HPO (Val)	$51 70 \pm 0.22$	40.55 ± 0.97
	Seen-tasks HPO (Mem)	$45.52_{\pm 0.41}$	$44.62_{\pm 0.18}$
	End-of-training HPO	$51.54_{\pm 0.38}$	$37.17_{\pm 0.48}$
	First-task HPO	$49.81_{\pm 0.10}$	$37.47_{\pm 0.26}$
iCaRL	Current-task HPO	$51.34_{\pm 0.32}$	$37.07_{\pm 0.07}$
	Seen-tasks HPO (Val)	$48.15_{\pm 0.09}$	$35.70_{\pm 0.23}$
	Seen-tasks HPO (Mem)	$47.87_{\pm 0.15}$	$35.27_{\pm 1.12}$
	End-of-training HPO	51.96 ± 0.60	$45.47_{\pm 0.42}$
	First-task HPO	$51.37_{\pm 0.16}$	$43.62_{\pm 1.09}$
ER-ACE	Current-task HPO	$51.78_{\pm 0.30}$	$43.87_{\pm 0.20}$
	Seen-tasks HPO (Val)	$51.94_{\pm 0.12}$	$43.15_{\pm 0.63}$
	Seen-tasks HPO (Mem)	$48.15_{\pm 0.28}$	$42.19_{\pm 0.84}$
	End-of-training HPO	$50.54_{\pm 0.16}$	$44.87_{\pm 0.26}$
	First-task HPO	$50.43_{\pm 0.34}$	$45.84_{\pm 0.50}$
ESMER	Current-task HPO	$50.68_{\pm 0.31}$	$44.50_{\pm 0.31}$
	Seen-tasks HPO (Val)	$47.96_{\pm 0.61}$	$42.18_{\pm 0.22}$
	Seen-tasks HPO (Mem)	$50.56_{\pm 0.40}$	$46.00_{\pm 0.43}$
	End-of-training HPO	$54.12_{\pm 0.70}$	$46.41_{\pm 0.77}$
	First-task HPO	$54.87_{\pm 0.39}$	$43.45_{\pm 3.55}$
DER++	Current-task HPO	$55.10_{\pm 0.52}$	$45.95_{\pm 0.93}$
	Seen-tasks HPO (Val)	$54.67_{\pm 0.57}$	$46.51_{\pm 0.49}$
	Seen-tasks HPO (Mem)	$49.06_{\pm 3.90}$	$25.78_{\pm 7.40}$

end-of-training HPO to first-task HPO of -0.39%. These results indicate, compared to standard
practice, that by using first-task framework it is possible to perform realistic HPO for much less
computation with only a small expected cost to performance. However, it is important to point out
that first-task HPO has a failure case of when the first task is not informative for the hyperparameter
choices of subsequent tasks. This failure case does not happen on the standard CL benchmarks used
in this work nor in the heterogeneous task setting where the tasks are designed to be more different.
Therefore, it is an open question whether such a failure case will arise if the standard CL benchmarks
used by the community change to be different, hopefully more realistic, data streams.

One of the potential reasons that the performance is similar between HPO frameworks is that there is little variation between the performance of different hyperparameter configurations. To see whether this is the case, we have plotted in Figure 3 histograms of the performance of using different fixed HPO configurations for DER++. The histograms show that hyperparameter configurations achieve a wide range of average accuracies. Therefore, the performance of different HPO configurations is not the reason why the HPO frameworks have similar results. Additionally, in Appendix B, we examine whether using default hyperparameters performs as well as selecting hyperparameters using HPO. We found that using default hyperparameters in most cases performed worse than using a HPO framework. Hence, our results suggest that HPO is necessary but that out of the HPO frameworks tested there is no one best performing method.



Figure 3: Histograms of the validation accuracy at the end of training for each hyperparameter setting searched over for DER++. We look at standard CL benchmarks and heterogeneous task benchmarks, which are identified by having a 'Hetero' in their name. The histograms show that different hyperparameter settings give a varying range of performances and only a few achieve near to the top performance.

505

506

507

508

6 CONCLUSIONS

512 513

521

522 523

525

In this paper we have benchmarked several hyperparameter optimisation (HPO) frameworks for CL 514 which are more realistic than the currently commonly used end-of-training HPO framework. We 515 benchmarked both fixed HPO frameworks, which fix the hyperparameters throughout training, and 516 dynamic HPO frameworks that continually adapt the hyperparameters. Our results show for com-517 monly used CL benchmarks that all the HPO frameworks achieve similar performances and none 518 consistently outperforms the others. Because of this, we recommend that practitioners using these 519 benchmarks should select a realistic HPO framework using other factors-for example compute 520 cost, for which *first-task* HPO is a good choice. Our results also suggest that future work on HPO for CL should move towards the use of new benchmarks where a difference in performances across HPO frameworks could arise.

524 **REPRODUCIBILITY STATEMENT**

To make our experiments reproducible, we provide in Section 5 a description of the setup used in 526 this work, which is in common with many other works in CL (Buzzega et al., 2020), and provide 527 more specific experimental details in Appendix A. Additionally, we provide the code used in the 528 supplementary material. 529

530 531

532

534

535

REFERENCES

- Rahaf Aljundi, Lucas Caccia, Eugene Belilovsky, Massimo Caccia, Min Lin, Laurent Charlin, and Tinne Tuytelaars. Online Continual Learning with Maximal Interfered Retrieval. In Proceedings of the 33rd Conference on the Advances in Neural Information Processing Systems, pp. 11849– 11860, 2019a.
- 536
- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient Based Sample Selection 538 for Online Continual Learning. In Proceedings of the 33rd Conference on the Advances in Neural Information Processing Systems, pp. 11816–11825, 2019b.

560

561

562

563

567

568

569

570

- Antreas Antoniou, Massimiliano Patacchiola, Mateusz Ochal, and Amos Storkey. Defining Benchmarks for Continual Few-shot Learning. *arXiv preprint arXiv:2004.11967*, 2020.
- James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for Hyper-Parameter
 Optimization. In *Proceeding of the 25th Conference on the Advances in Neural Information Processing Systems*, 2011.
- Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class Incremental Continual Learning into the Extended DER-verse. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5497–5512, 2022.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark Experience for General Continual Learning: a Strong, Simple Baseline. In *Proceedings of the 33rd Conference on the Advances in Neural Information Processing Systems*, pp. 15920–15930, 2020.
- Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky.
 New Insights on Reducing Abrupt Representation Change in Online Continual Learning. In *Proceedings of the 10th International Conference on Learning Representations*, 2021.
- Zhipeng Cai, Ozan Sener, and Vladlen Koltun. Online continual learning with natural distribu tion shifts: An empirical study with visual data. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8281–8290, 2021.
 - Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient Lifelong Learning with A-GEM. In *Proceedings of the 7th International Conference on Learning Representations*, 2019a.
- Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K
 Dokania, Philip HS Torr, and Marc'Aurelio Ranzato. On Tiny Episodic Memories in Contin ual Learning. *arXiv preprint arXiv:1902.10486*, 2019b.
 - Arslan Chaudhry, Naeemullah Khan, Puneet Dokania, and Philip Torr. Continual Learning in Lowrank Orthogonal Subspaces. In *Proceeding of the 34th Conference on the Advances in Neural Information Processing Systems*, pp. 9900–9911, 2020.
- Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg
 Slabaugh, and Tinne Tuytelaars. A Continual Learning Survey: Defying Forgetting in Classification Tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3366–3385,
 2021.
- Matthias Feurer and Frank Hutter. Hyperparameter Optimization. Automated machine learning: Methods, systems, challenges, pp. 3–33, 2019.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Sigrid Passano Hellan, Huibin Shen, François-Xavier Aubet, David Salinas, and Aaron Klein.
 Obeying the Order: Introducing Ordered Transfer Hyperparameter Optimisation. *arXiv preprint arXiv:2306.16916*, 2023.
- Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating Continual Learning Scenarios: A Categorization and Case for Strong Baselines. In *Proceedings of the 3rd Continual Learning Workshop, at the 32nd Conference on the Advances in Neural Information Processing Systems*, 2018.
- 589 Mert Kilickaya and Joaquin Vanschoren. What can AutoML do for Continual Learning? *arXiv* 590 *preprint arXiv:2311.11963*, 2023.
- Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. *Preprint*, 2009.
- ⁵⁹³ Thomas L Lee and Amos Storkey. Chunking: Forgetting Matters in Continual Learning even without Changing Tasks. *arXiv preprint arXiv:2310.02206*, 2023.

594 595 596	Thomas L Lee and Amos Storkey. Approximate Bayesian Class-Conditional Models under Con- tinuous Representation Shift. In <i>Proceedings of the 27th International Conference on Artificial</i> <i>Intelligence and Statistics</i> , 2024.
597 598 599 600	Yaoyao Liu, Yingying Li, Bernt Schiele, and Qianru Sun. Online hyperparameter optimization for class-incremental learning. In <i>Proceedings of the 37th AAAI Conference on Artificial Intelligence</i> , pp. 8906–8913, 2023.
601 602	Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. In <i>Conference on robot learning</i> , pp. 17–26. PMLR, 2017.
603 604 605 606	Seyed Iman Mirzadeh, Mehrdad Farajtabar, Razvan Pascanu, and Hassan Ghasemzadeh. Under- standing the Role of Training Regimes in Continual Learning. In <i>Proceedings of the 33rd confer-</i> <i>ence on the Advances in Neural Information Processing Systems</i> , pp. 7308–7320, 2020.
607 608	German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual Lifelong Learning with Neural Networks: A review. <i>Neural Networks</i> , 113:54 – 71, 2019.
609 610 611 612	Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. GDumb: A Simple Approach that Questions our Progress in Continual Learning. In <i>Proceeding of the 16th European Conference on Computer Vision</i> , pp. 524–540, 2020.
613 614 615	Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. ICARL: Incremental Classifier and Representation Learning. In <i>Proceedings of the IEEE conference on</i> <i>Computer Vision and Pattern Recognition</i> , pp. 2001–2010, 2017.
616 617 618	Fahad Sarfraz, Elahe Arani, and Bahram Zonooz. Error Sensitivity Modulation based Experience Replay: Mitigating Abrupt Representation Drift in Continual Learning. In <i>Proceedings of the</i> <i>Eleventh International Conference on Learning Representations</i> , 2023.
620 621	Rudy Semola, Julio Hurtado, Vincenzo Lomonaco, and Davide Bacciu. Adaptive Hyperparameter Optimization for Continual Learning Scenarios. <i>arXiv preprint arXiv:2403.07015</i> , 2024.
622 623 624	Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical Bayesian Optimization of Machine Learning Algorithms. In <i>Proceeding of the 26th Conference on the Advances in Neural Information Processing Systems</i> , 2012.
625 626 627	Gido M van de Ven and Andreas S Tolias. Three Scenarios for Continual Learning. <i>arXiv preprint arXiv:1904.07734</i> , 2019.
628 629	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A Comprehensive Survey of Continual Learning: Theory, Method and Application. <i>arXiv preprint arXiv:2302.00487</i> , 2023.
630 631 632	Martin Wistuba, Martin Ferianc, Lukas Balles, Cédric Archambeau, and Giovanni Zappella. Renate: A Library for Real-World Continual Learning. <i>arXiv preprint arXiv:2304.12067</i> , 2023.
633 634	Jiayu Wu, Qixiang Zhang, and Guoxi Xu. Tiny Imagenet Challenge (cs231n), http://tiny- imagenet.herokuapp.com/. Technical report, Stanford, 2015.
636 637 638 639	Tongtong Wu, Massimo Caccia, Zhuang Li, Yuan-Fang Li, Guilin Qi, and Gholamreza Haffari. Pretrained Language Models in Continual Learning: A Comparative Study. In <i>Proceedings of the</i> 10th International Conference on Learning Representations, 2022.
640 641 642	
644 645	
645 647	

648 A ADDITIONAL EXPERIMENTAL DETAILS

650 While we have aimed to include all the main experimental details in the main paper there are a few 651 others to mention here. First, we mostly follow the experimental setup of Buzzega et al. (2020) and 652 Boschini et al. (2022) and use the Mammoth library produced by those works as the base of our 653 code. Second, we use as our optimiser SGD with no momentum or weigh decay, as is done in other 654 works (Aljundi et al., 2019b; Buzzega et al., 2020; Chaudhry et al., 2019a; Lee & Storkey, 2023). Third, in the heterogeneous tasks setting we look at tasks sequences where each task in order has the 655 following number of classes associated with it [9, 2, 7, 3, 4, 9, 8, 3, 3, 7, 4, 4, 5, 9, 4, 5, 2, 8, 2, 2] and 656 all the data of a class is contained in the task associated with it. For Tiny ImageNet we only use the 657 first 100 classes in the heterogeneous tasks setting to reduce runtime and to make it more comparable 658 to CIFAR-100 in that setting. In the heterogeneous tasks setting each task has a variable amount of 659 data. For example, using CIFAR-100, the first task contains nine classes and so it will contain in 660 total 4500 examples (500 examples per task) while the second task contains two classes so will only 661 contain 1000 examples. Also, as in each task the learner needs to discriminate between a varying 662 number of classes the difficultly should vary between tasks. Additionally, in the heterogeneous 663 tasks setting we only look at class-incremental learning. Finally, the CORe50 dataset consists of 664 data drawn from multiple different background and lighting environments called sessions and the 665 test data consists of data from different sessions than the training data. Therefore, to insure that we more accurately model the covariate shift from the training to test data in our validation signal, 666 we construct the validation sets for CORe50 differently from the other datasets where it is sampled 667 randomly. Specifically, we use the data of session 2 contained in the task as the validation data for 668 that task. 669

670 We record here the hyperparameter grid that we sample over when performing HPO. We look at 671 learning rates in the set {0.2, 0.15, 0.1, 0.075, 0.05, 0.03, 0.01, 0.0075, 0.005, 0.0025}. For DER++, we perform HPO over both regularisation coefficients where we sample α in the set {0.2, 0.5, 1.0} 672 and β in the set {0.2, 0.5, 1.0}. For ESMER, we perform HPO over the loss margin coefficient 673 where we sample over the set $\{1.5, 1.2, 1.0\}$. We sample all possible combinations of learning rates 674 and regularisation coefficients in each of our HPO frameworks. This grid contains the ones used 675 in the popular works Buzzega et al. (2020), Boschini et al. (2022) and Sarfraz et al. (2023), where 676 we add additional learning rate settings and, for some datasets, regularisation coefficients settings. 677 We note here that while we use grid search in this paper to align with common practice in CL 678 (Buzzega et al., 2020; Delange et al., 2021), any hyperparameter sampling/selecting method can be 679 used with each of the HPO frameworks looked at. For example, tree-structured Parzen estimators 680 are a common Bayesian HPO method to sample hyperparameter configurations for neural networks 681 (Bergstra et al., 2011). Additionally, Gaussian process based HPO methods are also commonly used (Snoek et al., 2012) and have been looked at in settings related to online learning (Hellan et al., 682 2023). 683

684 685

686

B EXPERIMENTS USING DEFAULT HYPERPARAMETER VALUES

To test whether HPO is needed in CL and if instead using default hyperparameters is sufficient, we perform experiments using default hyperparameters. The experimental setup is the same as the main paper and we use for the default learning rate the default given by PyTorch, 0.001, and use 1.0 as the default for regularisation coefficients. The results are presented in Tables 5 and 7. The tables show that using default hyperparameters leads to worse performance than using HPO. Additionally, for some dataset and CL method combinations the default hyperparameters perform very badly showing the need to adapt hyperparameters to the dataset and CL method used.

- 694
- 695
- 696 697
- 609
- 699
- 700
- 701

Table 5: Comparison of using default hyperparameters versus using a HPO framework on split-task CIFAR-10 and CIFAR-100, where we only present the most common HPO framework (End-of-training HPO) and the most efficient (First-task HPO) for readability. We report mean average accuracies over three runs with their standard errors. The table shows that using default HPs leads to worse performance than using HPO for standard CL benchmarks.

		CIFAR-10		CIFAR-100	
CL Method	HPO Framework	Class-IL.	Task-IL.	Class-IL.	Task-IL.
ER	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 83.55_{\pm 0.44} \\ 84.38_{\pm 0.45} \\ 74.60_{\pm 0.79} \end{array}$	$\begin{array}{c} 97.18_{\pm 0.14} \\ 96.82_{\pm 0.17} \\ 94.53_{\pm 0.13} \end{array}$	$\begin{array}{c} 51.03_{\pm 0.43} \\ 49.61_{\pm 0.34} \\ 35.39_{\pm 0.36} \end{array}$	$\begin{array}{c} 85.68_{\pm 0.29} \\ 84.97_{\pm 0.19} \\ 72.83_{\pm 0.24} \end{array}$
iCaRL	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 77.79_{\pm 0.23} \\ 77.83_{\pm 0.22} \\ 68.34_{\pm 0.49} \end{array}$	$\begin{array}{c} 98.52_{\pm 0.03} \\ 95.31_{\pm 0.12} \\ 92.98_{\pm 0.21} \end{array}$	$\begin{array}{c} 54.30_{\pm 0.36} \\ 52.56_{\pm 0.10} \\ 11.54_{\pm 0.25} \end{array}$	$\begin{array}{c} 85.74_{\pm 0.45} \\ 84.60_{\pm 0.09} \\ 41.66_{\pm 0.54} \end{array}$
ER-ACE	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 82.34_{\pm 0.30} \\ 83.20_{\pm 0.79} \\ 75.46_{\pm 0.21} \end{array}$	$\begin{array}{c} 96.74_{\pm 0.01} \\ 96.67_{\pm 0.18} \\ 94.71_{\pm 0.06} \end{array}$	$\begin{array}{c} 55.58_{\pm 0.39} \\ 56.36_{\pm 0.29} \\ 42.65_{\pm 0.57} \end{array}$	$\begin{array}{c} 85.73_{\pm 0.09} \\ 86.11_{\pm 0.154} \\ 76.28_{\pm 0.19} \end{array}$
ESMER	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 80.73_{\pm 0.15} \\ 77.89_{\pm 0.46} \\ 68.86_{\pm 1.06} \end{array}$	$\begin{array}{c} 96.50_{\pm 0.01} \\ 96.15_{\pm 0.12} \\ 93.54_{\pm 0.20} \end{array}$	$\begin{array}{c} 56.16_{\pm 0.54} \\ 56.61_{\pm 0.20} \\ 42.94_{\pm 0.61} \end{array}$	$\begin{array}{c} 88.69_{\pm 0.35} \\ 89.05_{\pm 0.10} \\ 79.64_{\pm 0.36} \end{array}$
DER++	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 84.40_{\pm 0.94} \\ 85.22_{\pm 0.08} \\ 77.59_{\pm 0.45} \end{array}$	$\begin{array}{c} 95.75_{\pm 0.33} \\ 96.14_{\pm 0.10} \\ 93.83_{\pm 0.40} \end{array}$	$\begin{array}{c} 56.04_{\pm 3.67} \\ 55.20_{\pm 0.78} \\ 46.11_{\pm 1.16} \end{array}$	$\begin{array}{c} 83.13_{\pm 2.69} \\ 81.68_{\pm 0.66} \\ 78.14_{\pm 1.28} \end{array}$

Table 6: Comparison of using default hyperparameters versus using a HPO framework on split-task
Tiny ImageNet, where we only present the most common HPO framework (End-of-training HPO)
and the most efficient (First-task HPO) for readability. We report mean average accuracies over three
runs with their standard errors. The table shows that using default HPs leads to worse performance
than using HPO for standard CL benchmarks.

		TinyImageNet	
CL Method	HPO Framework	Class-IL.	Task-IL.
ER	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 28.01_{\pm 0.09} \\ 28.51_{\pm 0.18} \\ 16.27_{\pm 0.20} \end{array}$	$\begin{array}{c} 68.17_{\pm 0.0} \\ 68.72_{\pm 0.1} \\ 50.99_{\pm 0.4} \end{array}$
iCaRL	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 37.09_{\pm 0.27} \\ 36.42_{\pm 0.22} \\ 5.30_{\pm 0.03} \end{array}$	$\begin{array}{c} 70.37_{\pm 0.3} \\ 70.11_{\pm 0.1} \\ 23.97_{\pm 0.1} \end{array}$
ER-ACE	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 38.94_{\pm 0.47} \\ 36.94_{\pm 0.67} \\ 25.84_{\pm 0.26} \end{array}$	$\begin{array}{c} 70.18_{\pm 0.2} \\ 68.16_{\pm 0.3} \\ 56.25_{\pm 0.1} \end{array}$
ESMER	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 47.33_{\pm 0.30} \\ 46.69_{\pm 0.56} \\ 33.11_{\pm 0.39} \end{array}$	$\begin{array}{c} 76.18_{\pm 0.2} \\ 75.72_{\pm 0.2} \\ 63.15_{\pm 0.1} \end{array}$
DER++	End-of-training HPO First-task HPO Default HPs	$39.89_{\pm 0.27} \\ 35.98_{\pm 0.63} \\ 25.66_{\pm 0.16}$	$70.41_{\pm 0.1}$ $65.86_{\pm 0.3}$ $59.14_{\pm 0.5}$

Table 7: Comparison of using default hyperparameters versus using a HPO framework on heterogeneous task benchmarks, where we only present the most common HPO framework (End-of-training
HPO) and the most efficient (First-task HPO) for readability. We report mean average accuracies
over three runs with their standard errors. The table shows that using default HPs leads to worse
performance than using HPO for heterogeneous task benchmarks.

			Hetero-CIFAR-100	Hetero-TinyImg
CL M	lethod	HPO Framework	Class-IL.	Class-IL.
ER		End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 50.41_{\pm 0.21} \\ 50.33_{\pm 0.50} \\ 33.76_{\pm 0.78} \end{array}$	$\begin{array}{c} 39.41_{\pm 0.57} \\ 40.77_{\pm 0.34} \\ 26.88_{\pm 0.45} \end{array}$
iCaRl		End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 51.54_{\pm 0.38} \\ 49.81_{\pm 0.10} \\ 12.23_{\pm 0.19} \end{array}$	$\begin{array}{c} 37.17_{\pm 0.48} \\ 37.47_{\pm 0.26} \\ 10.6_{\pm 0.26} \end{array}$
ER-A	CE	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 51.96_{\pm 0.60} \\ 51.37_{\pm 0.16} \\ 38.11_{\pm 0.80} \end{array}$	$\begin{array}{c} 45.47_{\pm 0.42} \\ 43.62_{\pm 1.09} \\ 32.37_{\pm 0.53} \end{array}$
ESMI	ER	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 50.54_{\pm 0.16} \\ 50.43_{\pm 0.34} \\ 37.92_{\pm 0.30} \end{array}$	$\begin{array}{c} 44.87_{\pm 0.26} \\ 45.84_{\pm 0.50} \\ 34.22_{\pm 0.41} \end{array}$
DER-	++	End-of-training HPO First-task HPO Default HPs	$\begin{array}{c} 54.12_{\pm 0.70} \\ 54.87_{\pm 0.39} \\ 44.43_{\pm 0.51} \end{array}$	$\begin{array}{c} 46.41_{\pm 0.77} \\ 43.45_{\pm 3.55} \\ 30.21_{\pm 1.53} \end{array}$