LUNCH: ADAPTIVE BALANCING OF CONTINUAL LEARNING VIA HYPERPARAMETER UNCERTAINTY

Anonymous authors

Paper under double-blind review

Abstract

Continual learning (CL) is characterized by learning sequentially arriving tasks and behaving as if they were observed simultaneously. In order to prevent catastrophic forgetting of old tasks when learning new tasks, representative CL methods usually employ additional loss terms to balance their contributions (e.g., regularization and replay), modulated by deterministic hyperparameters. However, this strategy struggles to accommodate real-time changes in data distributions and is also lack of robustness to subsequent unseen tasks, especially in online scenarios where CL is performed with a one-pass data stream. Inspired by adaptive weighting in multi-task learning, we propose an innovative approach named Learning UNCertain Hyperparameters (LUNCH) for adaptive balancing of task contributions in CL. Specifically, we formulate each CL-relevant hyperparameter as a function of optimizable uncertainty under homoscedastic assumption and ensure its training stability through the exponential moving average of network parameters. We further devise an evaluation protocol that moderately adjusts the hyperparameter values and reports their impact on performance, so as to analyze the sensitivity of these sub-optimal values in realistic applications. We perform extensive experiments to demonstrate the effectiveness and robustness of our approach, which significantly improves online CL in a plug-in manner (e.g., up to 11.26% and 5.64% on Split CIFAR-100 and Split Mini-ImageNet, respectively) as well as offline CL.¹

032

004

005

010 011

012

013

014

015

016

017

018

019

021

023

025

026

027

1 INTRODUCTION

033 The ability of continual learning (CL) is critical for artificial intelligence systems to accommodate 034 real-world changes, yet limited by catastrophic forgetting of old tasks when learning new tasks 035 (Wang et al., 2024a; McClelland et al., 1995). In order to strike an appropriate balance between task contributions within the same parameter space, representative CL methods often employ additional 037 loss terms to preserve previously learned knowledge, such as regularization of parameter changes (Kirkpatrick et al., 2017; Zenke et al., 2017) and replay of a few old training samples (Buzzega et al., 038 2020; Rebuffi et al., 2017). In general, the strength of these loss terms is regulated by deterministic hyperparameters obtained from a grid search (Chaudhry et al., 2018; Wang et al., 2023a). However, 040 this strategy is sub-optimal in performance as it struggles to adapt to real-time changes of data 041 distributions within the observed task sequence, and is also lacking robustness to subsequent unseen 042 tasks. These critical challenges tend to be more significant in online CL where each task is learned 043 from a one-pass data stream (Fini et al., 2020; Zhang et al., 2022). 044

In this regard, we analyze in depth the role of CL-relevant hyperparameters in balancing task con-045 tributions. We first formulate representative CL methods with a shared mathematical form of the 046 loss function. Besides a loss term for learning the current task, the loss function typically involves 047 additional loss terms that preserve previously learned knowledge in terms of the parameter space 048 and the output space with corresponding hyperparameters. These loss terms amount to approxi-049 mate multi-task learning (MTL) for all tasks ever seen, i.e., the upper bound of CL, while largely avoiding the use of old training samples. In MTL, adaptive weighting of task contributions in hy-051 perparameters has been shown to be an effective strategy compared to fixed weighting (i.e., using 052 deterministic hyperparameters) (Kendall et al., 2018; Liu et al., 2019; 2022; Lin et al., 2021), but

⁰⁵³

¹Our code is included in Supplementary Materials for examination and will be released upon acceptance.

remains under-explored and highly non-trivial for CL due to the dynamic and unpredictable nature of data distribution.

Based on the above analysis, we present Learnable UNCertain Hyperparameters (LUNCH), an inno-057 vative approach that enables adaptive balancing of task contributions in CL. Specifically, we formulate each CL-relevant hyperparameter as a function of optimizable uncertainty, which is initialized 059 high and then decreases during the learning of changes in data distributions. Under the homoscedas-060 tic uncertainty assumption, we derive probabilistic implementations for the loss terms of the param-061 eter space and output space, corresponding to regression and classification problems, respectively. 062 Whenever a new task is introduced, the uncertain hyperparameters need to be refreshed to re-balance 063 the contributions, resulting in a performance degradation known as the "stability gap" (De Lange 064 et al., 2022). In this regard, we perform exponential moving average of network parameters along the training trajectory, so as to stabilize training upon reinitialization. 065

066 We perform extensive experiments to evaluate our approach. Beyond the widely-used average accu-067 racy for overall performance, we consider two additional evaluation metrics including the average 068 anytime accuracy and the worst-case accuracy for real-time changes in data distributions. We further 069 evaluate the sensitivity of sub-optimal hyperparameter values through analyzing their impact under moderate adjustments. Our approach demonstrates outstanding performance with significant im-070 provements in effectiveness and robustness across various online CL benchmarks, benefiting recent 071 strong baselines in a plug-in manner (e.g., up to 11.26% and 5.64% on Split CIFAR-100 and Split 072 Mini-ImageNet, respectively) and also remarkably facilitate offline CL. 073

Our contributions can be summarized as follows: (1) We perform an in-depth analysis of CL-relevant hyperparameters under a unified framework of representative CL methods and task balancing strategies in MTL; (2) We propose an innovative approach that incorporates optimizable uncertainty into CL-relevant hyperparameters for adaptive balancing of task contributions, coupled with exponential moving average of network parameters to address the stability gap; and (3) Our approach significantly improves the effectiveness and robustness of CL, validated by extensive experiments.

080 081

082

2 RELATED WORK

083 **Continual Learning (CL)**, also known as incremental learning or lifelong learning, aims to over-084 come catastrophic forgetting of old tasks when learning new tasks (Wang et al., 2024a; McClelland 085 et al., 1995). Numerous efforts have been devoted into addressing this challenging issue. A majority of representative methods attempt to strike an appropriate balance between task contributions within 087 the same parameter space. For example, regularization-based methods employ explicit regulariza-088 tion terms to stabilize network parameters and simulate behaviors of the old model (Kirkpatrick et al., 2017; Buzzega et al., 2020; Li & Hoiem, 2017). Meanwhile, replay-based methods approx-089 imate and recover the old data distributions through preserving a small memory buffer or learning 090 a generative model (Buzzega et al., 2020; Shin et al., 2017; Aljundi et al., 2019). Other methods 091 that optimize network parameters in different parameter spaces are often collectively referred to as 092 architecture-based methods (Serra et al., 2018; Kang et al., 2022; Rusu et al., 2016), which explicitly 093 avoid the problem of balancing task contributions. However, this kind of method typically requires 094 the oracle of task identity at test time in order to select an appropriate parameter space, and is there-095 fore not prioritized in this work. Based on the availability of training samples, the widely-used CL 096 setups can be categorized into online CL and offline CL (detailed in Section 3.1), with the former 097 being considered realistic yet much more challenging.

098 Hyperparameter of CL. For representative CL methods, an appropriate management of hyperparameters (e.g., learning rates, regularization strengths, memory buffer sizes, etc.) is critical for 100 achieving outstanding performance across tasks. A primary consideration is task balancing, which 101 ensures that the model maintains performance on old tasks while learning new ones (Cha & Cho, 102 2024; Yildirim et al., 2023). As the upper bound of CL, multi-task learning (MTL) attempts to 103 address this problem through various weighting strategies: fixed weighting assigns constant impor-104 tance to each task, while adaptive weighting adjusts the importance based on task difficulty or model 105 performance (Kendall et al., 2018; Liu et al., 2019; 2022; Lin et al., 2021). Despite the effectiveness in MTL, adaptive weighting of CL-relevant hyperparameters is remarkably challenging due to the 106 dynamic and unpredictable nature of data distribution, and therefore remains largely under-explored 107 in literature.

¹⁰⁸ 3 PRELIMINARIES

109 110 111

112

113

144

145

In this section, we first describe the problem formulation of CL with a unified framework of representative methods. We then analyze the selection of hyperparameter(s) in CL on the basis of the inherent connections between CL and MTL.

114 3.1 PROBLEM FORMULATION

116 Let us consider a sequence of tasks defined by a collection of respective training sets S =117 $\{\mathcal{D}_1, \cdots, \mathcal{D}_T\}$, where \mathcal{D}_t for each task t consists of several data-label pairs (x_t, y_t) with $x_t \in \mathcal{X}_t$ 118 and $y_t \in \mathcal{Y}_t$. The goal of CL is to learn a mapping $f_{\theta} : \bigcup_{t=1}^T \mathcal{X}_t \to \bigcup_{t=1}^T \mathcal{Y}_t$ parameterized by 119 trainable parameters θ with sequentially arriving \mathcal{D}_t , so as to achieve superior performance on all 120 observed tasks (Wang et al., 2024a). Since previous training samples are often unavailable at the 121 current training stage, it remains extremely challenging to strike an appropriate balance between old 122 and new tasks, resulting in catastrophic forgetting (i.e., f_{θ} abruptly and dramatically forget previ-123 ously learned knowledge upon new information). Regarding specific setups, training samples for each task can be reused for multiple epochs in offline CL, but arrive as a one-pass data stream in 124 online CL, which greatly adds to the challenge. 125

126 To alleviate catastrophic forgetting when optimizing θ within the *same* parameter space, many repre-127 sentative methods have been proposed for CL. These methods can be classified into regularization-128 based methods, which incorporate additional regularization term(s) to stabilize knowledge of the 129 parameter space and the output space (Kirkpatrick et al., 2017; Li & Hoiem, 2017), as well as replay-based methods (Aljundi et al., 2019; Zhang et al., 2022), which preserve some old train-130 ing samples with a small memory buffer. In particular, replay is often coupled with regularization 131 (Buzzega et al., 2020; Rebuffi et al., 2017) to encourage the current model f_{θ} to mimic the behaviors 132 of the old model f_{θ^*} with parameters θ^* when processing old training samples. 133

Following a recent work (Wang et al., 2024b), these two kinds of methods can be described as shared mathematical forms under a unified framework:

$$\mathcal{L}_{\rm CL} = \lambda_n \underbrace{\mathcal{L}_n(x, y)}_{\text{new task}} + \lambda_o \underbrace{\mathcal{L}_o(f_\theta(x), z)}_{\text{output space}} + \lambda_p \underbrace{\mathcal{L}_p(\theta, \theta^*)}_{\text{parameter space}},$$
(1)

where \mathcal{L}_n denotes the loss function for learning each new task. \mathcal{L}_o and \mathcal{L}_p restrict update rates in output space and parameter space, respectively. The example definitions of \mathcal{L}_o and \mathcal{L}_p will be described latter in Table 1. It can be seen that the contributions of new and old tasks are explicitly regulated by the hyperparameters $\{\lambda_n, \lambda_o, \lambda_p\}$.

3.2 HYPERPARAMETERS IN CONTINUAL LEARNING

146 With Eq. (1), we further analyze the selection of hyperparameters in representative CL methods. 147 Similar to regular machine learning methods, the optimal hyperparameter values for CL are usu-148 ally obtained by repeated iterations of the task sequence $S = \{\mathcal{D}_1, \cdots, \mathcal{D}_T\}$, which can be further 149 divided into two strategies. One is to run the first several tasks iteratively to determine the hyperpa-150 rameter values and use them to learn subsequent tasks (Chaudhry et al., 2018; Pham et al., 2021); the 151 other is to run the entire task sequence iteratively in different orders to determine the hyperparameter 152 values and provide a sensitivity analysis of them (Yildirim et al., 2023; Cha & Cho, 2024; Van de 153 Ven et al., 2022). Both strategies assume relatively stable changes in data distribution and employ deterministic hyperparameters for CL, making it difficult to adapt to real-world scenarios that are 154 highly dynamic and unpredictable (Semola et al., 2024; Cha & Cho, 2024). In addition, manually 155 adjusting these hyperparameters over time is both costly and impractical. 156

157 In fact, MTL is often considered to be the upper bound of CL. Both CL and MTL aim to achieve 158 the same objective, i.e., to find a solution θ that performs well for all observed tasks, with the main 159 difference being whether S is provided sequentially or simultaneously. Formally, the objective of 160 MTL can be defined as follows:

$$\mathcal{L}_{\mathrm{MTL}} = \sum_{i} w_i \mathcal{L}_i, \tag{2}$$

Table 1: Definition of representative CL methods that target the same parameter space. F is the 164 Fisher information matrix to approximate the importance of the network parameters. M denotes the memory buffer consisting of a few old training samples. x_{aug} is the augmentation of x. z is the 165 output logit of the old models in the training trajectory. 166

167	Method	Regularization Loss	Replay Loss
160	EWC (Kirkpatrick et al., 2017)	$(\theta - \theta^*)^\top F(\theta - \theta^*)$	-
170	ER (Shin et al., 2017)	-	$\mathbb{E}_{(x,y)\in M}(\mathcal{L}(x,y))$
171	MIR (Aljundi et al., 2019)	- (11.6.()) (12)	$\max \mathbb{E}_{(x,y)\in M}(\mathcal{L}(x,y))$
172	DER (Buzzega et al., 2020) DER $\pm \pm$ (Boschini et al. 2022)	$\mathbb{E}_{(x,y)\in M}(\ f_{\theta}(x) - z\ _{2}^{2})$ $\mathbb{E}_{(x,y)\in M}(\ f_{\theta}(x) - z\ _{2}^{2})$	$\mathbb{F}(x) = x(f(x,y))$
173	RAR (Zhang et al., 2022)	$ = (x,y) \in M (J\theta(x) - z _2) $	$\mathbb{E}_{(x,x_{\text{ang}},y)\in M}(\mathcal{L}(x,y))$ $\mathbb{E}_{(x,x_{\text{ang}},y)\in M}(\mathcal{L}(x_{\text{aug}},y))$
174		1	(,

162 163

> where \mathcal{L}_i corresponds to the loss function for learning each task, and w_i denotes the hyperparameter that regulates the task weight. The loss function of CL in Eq. (1) can be seen as an approximation of Eq. (2), with the use of old training samples largely avoided.

178 Although many MTL methods also employ deterministic hyperparameters selected from a grid 179 search, even simplified into an unweighted form $w_i = w_i$, adaptive balancing of task contributions has proven to be a superior strategy (Kendall et al., 2018; Liu et al., 2019; 2022). In particular, 181 the corresponding $\{w_i\}$ can be modeled as an optimizable function related to the relative confidence 182 between tasks (Liu et al., 2019; Kendall & Gal, 2017). However, these adaptive weighting strategies 183 remain under-explored and highly non-trivial for CL, due to the dynamic and unpredictable proper-184 ties of data distributions. To this end, we aim to provide an innovative approach to address the above 185 challenges, as detailed below.

186

175

176

177

187 188

LEARNABLE UNCERTAIN HYPERPARAMETERS (LUNCH) 4

189 In this section, we design an innovative adaptive weighting strategy for CL that optimizes CL-190 relevant hyperparameters according to training progress. We incorporate optimizable uncertainty 191 into CL-relevant hyperparameters under the unified framework of representative CL methods, and 192 further rectify the "stability gap" introduced by uncertainty refresh.

193 We first define the specific forms of several representative CL methods with Eq. (1), as shown in 194 Table 1. These methods mainly focus on addressing the problem of catastrophic forgetting by lim-195 iting model updates in either parameter space or output space, so as to preserve previously learned 196 knowledge. For adaptive balancing of task contributions, we propose to incorporate optimizeable 197 parameters σ^2 (called "uncertainty") into the hyperparameter set $\{\lambda_n, \lambda_o, \lambda_p\}$ in Eq. (1). From a Bayesian perspective, such uncertainty can capture the model's confidence in different types of 199 tasks and accordingly adjust the CL process to balance new and old tasks (Guo et al., 2011). During the learning of each task in CL, the confidence in the predictive distribution $p(y|x,\theta)$ should grad-200 ually increase, and corresponding uncertainty σ^2 should gradually decrease from large to small and 201 eventually stabilize at a certain value (Kendall & Gal, 2017). 202

203 In particular, we observe that the loss terms \mathcal{L}_n , \mathcal{L}_o and \mathcal{L}_p correspond to addressing a regres-204 sion or classification problem (Bishop & Nasrabadi, 2006). Specifically, the classification problem 205 predicts discrete labels (e.g., DER employs a cross-entropy loss \mathcal{L}_{ρ} with a small memory buffer), 206 while the regression problem predicts continuous numerical value (e.g., EWC employs a weighted squared loss \mathcal{L}_p to stabilize parameter changes). Here we focus on homoscedastic aleatoric uncer-207 tainty (Kendall & Gal, 2017), whose corresponding hyperparameters $\{\lambda_n, \lambda_o, \lambda_p\}$ do not depend on 208 specific input data, rather stay constant for all inputs and vary only between different tasks. 209

210 For *regression* problems, the predictive distribution can be defined as a Gaussian distribution under 211 the Laplace approximation (Bishop & Nasrabadi, 2006), finding a Gaussian approximation to a continuous probability density: 212

$$p(y \mid f_{\theta}(x)) = \mathcal{N}\left(f_{\theta}(x), \sigma^2\right), \tag{3}$$

213 214

$$p(y \mid f_{\theta}(x)) = \mathcal{N}\left(f_{\theta}(x), \sigma^{2}\right), \qquad (3)$$

where σ^2 denotes the homoscedastic aleatoric uncertainty. Due to the space limit, the detailed 215 derivation can be found in Appendix A.1. When using the mean squared error as the loss function



Figure 1: Demonstration of our approach. Session 2 and Session 3 shows that the proposed LUNCH incorporates optimizable uncertainty into CL-relevant hyperparameters, enabling adaptive balancing of task contributions in both parameter space and output space.

 \mathcal{L} for regression problem, the corresponding log-likelihood becomes:

$$-\log p(y \mid f_{\theta}(x)) \propto \frac{1}{2\sigma^{2}} \|y - f_{\theta}(x)\|^{2} + \frac{1}{2}\log \sigma^{2}.$$
 (4)

For *classification* problems, we use a Gibbs distribution to capture the predictive distribution scaled by the learnable "temperature" σ^2 , which determines the flatness of discrete distribution (i.e., entropy) (Bishop & Nasrabadi, 2006):

$$p(y \mid f_{\theta}(x)) = \operatorname{Softmax}\left(\frac{1}{\sigma^2}f_{\theta}(x)\right),\tag{5}$$

where the σ^2 determines the degree of entropy divergence in the deterministic discrete distribution. Then, the corresponding log-likelihood becomes:

$$\log p\left(y=c \mid f_{\theta}(x)\right) = \frac{1}{\sigma^2} f_{\theta}(x) - \log\left[\sum_{c'=1} \exp\left(\frac{1}{\sigma^2} f_{\theta}^{c'}(x)\right)\right],\tag{6}$$

where $f_{\theta}^{c'}$ denotes the c'-th logit. To simplify the derivation in CL regarding its dynamic and unpredictable nature, the explicit trick $\frac{1}{\sigma^2} \sum_{c'=1}^{c} \exp\left(\frac{1}{\sigma^2} f_{\theta}^{c'}(x)\right) \approx \left(\sum_{c'=1}^{c} \exp\left(\frac{1}{\sigma^2} f_{\theta}^{c'}(x)\right)\right)^{\frac{1}{\sigma^2}}$ becomes an equality when $\sigma^2 \to 1$, which has been verified in empirically improving the performance. Due to the space limit, the derivation is detailed in Appendix A.2. Then, the overall log-likelihood of the predictive distribution is:

$$-\log p\left(y \mid f_{\theta}(x)\right) \approx \frac{1}{\sigma^2} \log \operatorname{Softmax}\left(y, f_{\theta}(x)\right) + \log \sigma^2.$$
(7)

Then, we can reformulate each CL-relevant hyperparameter in Eq. (1) as a function $\lambda(\sigma^2)$ of homoscedastic aleatoric uncertainty σ^2 based on the training progress. The overall objective function $\log p(y \mid x; \sigma_n^2, \sigma_o^2, \sigma_p^2)$ is defined as follows:

$$\log\left(p \mid x, y; \sigma_n^2, \sigma_o^2, \sigma_p^2\right) = \exp\left(-\log\sigma_n^2\right) \underbrace{\mathcal{L}_n(x, y)}_{\text{new task}} + \exp\left(-\log\sigma_o^2\right) \underbrace{\mathcal{L}_o\left(h_\theta(x), z\right)}_{\text{output space}} + \exp\left(-\log\sigma_p^2\right) \underbrace{\mathcal{L}_p\left(\theta, \theta^*\right)}_{\text{parameter space}} + \log\sigma_n^2 + \log\sigma_o^2 + \log\sigma_p^2, \tag{8}$$

where $\mathcal{L}_n(x, y)$, $\mathcal{L}_o(x, y)$ and $\mathcal{L}_p(x, y)$ represent the loss terms for CL. { $\lambda_n(\sigma_n^2) = \exp(-\log \sigma_n^2)$, $\lambda_o(\sigma_o^2) = \exp(-\log \sigma_o^2)$, $\lambda_p(\sigma_p^2) = \exp(-\log \sigma_p^2)$ } represent aforementioned uncertain hyperparameters. As uncertainty σ^2 increases, the relative weights of the loss function \mathcal{L} decrease, and vice versa. The additional term $\log \sigma^2$ discourages rapid changes in σ^2 , thus stabilizing the norm of relative weights. We also use the exponential mapping trick (Murphy, 2012) to



Figure 2: CL suffers substantial forgetting as each new task is introduced, followed by a phase of performance recovery. The experiment is performed with RAR on Split CIFAR-10.

ensure numerical stability, letting f_{θ} directly predict the log variance σ^2 to avoid division by zero. 285 As shown in Fig. 2 (left), when they gradually learn representative features in a new task, the corre-286 sponding uncertainty about the predictive distribution gradually decreases and ultimately stabilizes 287 at a fixed local optima (Kendall & Gal, 2017). 288

289 With significant changes in data distribution over the course of training (e.g., switching from \mathcal{D}_{t-1} 290 to \mathcal{D}_t), we need to refresh the uncertain hyperparameters for adaptation (i.e., the uncertainty is reraised, and then gradually decreased): 291

$$\lambda_n\left(\sigma_n^2\right) = \lambda_n\left(\sigma_{n,\text{init}}^2\right); \lambda_o\left(\sigma_o^2\right) = \lambda_o\left(\sigma_{o,\text{init}}^2\right); \lambda_p\left(\sigma_p^2\right) = \lambda_p\left(\sigma_{p,\text{init}}^2\right).$$
(9)

This is accompanied by a critical challenge called the "stability gap" (De Lange et al., 2022): f_{θ} suffers substantial forgetting when starting to learn new tasks, followed by a phase of performance 295 recovery (see Fig. 2, right). 296

297 To alleviate the mismatch between the refreshed hyperparameters and the current training progress, 298 we employ the temporal ensemble (Laine & Aila, 2016) strategy to offset the bias. Specifically, we 299 collect models along the training trajectory with exponential moving average (EMA), denoted as 300 $f_{\theta^{\text{ema.}}}$. For the current model f_{θ} parameterized by trainable θ , the EMA at step t is defined as:

$$\theta_t^{\text{ema}} = \frac{\beta \theta_{t-1}^{\text{ema}} + (1-\beta)\theta_t}{1-\beta^t} = \beta^t \theta_0^{\text{ema}} + \sum_{i=1}^t (1-\beta)\beta^{t-i}\theta_i, \tag{10}$$

304 where β controls the strength of EMA, β^t denotes β raised to the power of t, and the parameter θ_i 305 is at the *i*-th step in the training trajectory before t step. The temporal ensemble of different models 306 can enhance the stability of online CL, which has also been observed in a recent study (Soutif-307 Cormerais et al., 2023). To stabilize training when refreshing hyperparameters, the EMA strategy 308 implicitly integrates hyperparameters from different training stages of the old tasks, providing an 309 adaptive benefit for implementing the optimizable uncertainty.

Finally, we describe the LUNCH training procedure with a pseudo-code in Algorithm 1. The gray 311 area highlights the key procedure. As can be seen, our approach is easy to implement and compatible 312 with many representative CL methods, such as ER (Rolnick et al., 2019), DER (Buzzega et al., 313 2020), RAR (Zhang et al., 2022), etc. 314

315 316

317 318

319 320

310

272

273

275

281

282

283 284

292 293

301 302

303

EXPERIMENTS 5

In this section, we briefly describe the experimental setups and then analyze the experimental results.

5.1 EXPERIMENTAL SETUPS

321 Benchmark. Here we consider several benchmarks that are commonly used in the literature of 322 online CL (Wang et al., 2023b). Specifically, CIFAR-10 dataset includes 10 classes of images sized 323 32×32 , randomly divided into 5 disjoint tasks with 2 classes each. CIFAR-100 dataset includes

ithm 1 Learnable UNCertain Hyperparameters (LUNCH) for CL
aput : Deep model f_{θ} with parameters θ ; uncertainty $\{\sigma_n^2, \sigma_o^2, \sigma_p^2\}$; datasets $S =$
$\mathcal{D}_1,\cdots,\mathcal{D}_T\}.$
nitialization : $\theta_0^{\text{ema}} = \theta_{\text{init}}$; $E = 1$ for online CL; $E > 1$ for offline CL.
or task $t = 1, \cdots, T$ do
if $t = 1$ then
for epoch $e = 1, \cdots, E$ do
Update the parameters θ of f_{θ} only with loss \mathcal{L}_n
end for
Save the parameters θ of f_{θ} as θ_1^{ema}
else
for epoch $e = 1, \cdots, E$ do
Calculate the unified CL loss by Eq. (8).
Update the parameters θ of f_{θ} and the uncertainty $\{\sigma_n^2, \sigma_o^2, \sigma_p^2\}$
Update θ_t^{ema} by Eq. (10)
end for
Refresh the uncertainty $\{\sigma_n^2, \sigma_o^2, \sigma_p^2\}$ in Eq. (9)
end if
ad for

347100 classes of images sized 32×32 , randomly divided into 20 disjoint tasks with 5 classes each.348Mini-ImageNet dataset includes 100 classes of images sized 84×84 , randomly split into 20 disjoint
tasks with 5 classes each. The common image resolution for the ImageNet-R dataset is 224×224 ,
randomly divided into 20 disjoint tasks with 10 classes each (Hendrycks et al., 2021).

Implementation. We use Empirical Replay (ER) (Rolnick et al., 2019), Dark Experience Replay (DER) (Buzzega et al., 2020) and Repeated Augmented Rehearsal (RAR) (Zhang et al., 2022) as the main baselines for CL. Following the implementation of RAR (Zhang et al., 2022), we train a ResNet-32 backbone with an SGD optimizer of learning rate 0.1 in all experiments. The batch size is set to 32 for Split CIFAR-10/100, and 64 for both Split Mini-ImageNet and Split ImageNet-R. The memory buffer is unified to maintain 2000 training samples in total. For offline CL, the number of epochs is set to 100 that is sufficient for convergence on each task.

Evaluation Metric. We consider multiple evaluation metrics to provide a comprehensive analysis of CL. We first define the Final Average Accuracy (FAA) to effectively evaluate the overall performance after learning the last task \mathcal{D}_T . Formally, \mathcal{D}_i represents the *i*-th task, while f_{θ_i} denotes the model parameterized by trainable θ_i at the *i*-th task:

$$FAA = \frac{1}{T} \sum_{i=1}^{T} A(\mathcal{D}_i, f_{\theta_T}).$$
(11)

However, FAA only provides a snapshot of the state after learning all tasks, ignoring performance during task transitions. To provide a more comprehensive assessment, we therefore consider two other metrics. One is the Average Anytime Accuracy (AAA), which measures the average performance on all observed tasks (De Lange et al., 2022).

$$AAA = \frac{1}{T} \sum_{j=1}^{T} \frac{1}{j} \sum_{i=1}^{j} A\left(\mathcal{D}_{i}, f_{\theta_{j}}\right).$$
(12)

Another is the Worst-Case Accuracy (WCA), which evaluates the performance of CL methods in
 the worst-case scenarios (De Lange et al., 2022). WCA first obtains the average minimum accuracy
 called minAcc in previous tasks and then combines it with the performance of current task (i.e., the

381				$\frac{1}{4}\lambda$			λ			4λ	
000	Benchmark	Method	$FAA(\uparrow)$	$WCA(\uparrow)$	$AAA(\uparrow)$	$FAA(\uparrow)$	$WCA(\uparrow)$	$AAA(\uparrow)$	$FAA(\uparrow)$	WCA (\uparrow)	AAA (\uparrow)
382		ER	$43.51{\scriptstyle \pm 0.40}$	$38.58{\scriptstyle\pm 5.01}$	57.12 ± 4.06	34.77±8.01	$32.97{\scriptstyle\pm 6.47}$	$54.98{\scriptstyle\pm2.67}$	34.78±3.22	$32.97{\scriptstyle\pm2.23}$	$43.26{\scriptstyle\pm3.24}$
383		w/ Ours	47.50±0.58	40.59±4.96	59.73±3.38	45.64±1.14	39.24±3.04	58.22±3.68	47.64±3.99	39.24±2.87	53.07±3.65
20/		$\Delta(\uparrow)$	3.99	2.01	2.61	10.87	6.27	3.24	12.86	6.27	9.81
304	G I' CIEAD 10	DER	30.35±0.22	36.05±0.22	35.69±0.23	4/.03±4.81	41.46±8.15	50.5/±0.84	15.92±0.84	15.91±0.84	30.92±0.54
385	Split CIFAR-10	w/ Ours	39.02±1.77	38.9/±1.83	42.70±3.36	34.98±3.55	45./8±6.25	38.40±0.79	19.52±1.42	18.00±1.39	34.39±2.07
386		$\Delta(1)$	2.47	2.92	/.01	21.07	4.52	26.67	3.00	2.75	3.07
000		KAK	11 72 + 2.81	27.21±4.36	40.13±2.49	31.97±5.01	24.34±3.73	42 71 Laur	33.33±3.01	31.09±6.47	41.40±2.67
387		$\Delta(\uparrow)$	6.13	5 1 5	5 96	2 05	4 55	42.71±2.15	10.69	7 55	12 00
388		ER	9 99 _{+2 38}	5.22+1.32	12.91+251	9.66+1.89	5 27+1 32	12.84+1.74	9.66+1.89	5 21+1 32	12.84+1.74
000		w/ Ours	15 33+563	$779_{\pm 3.69}$	16 68+573	20.92 + 3.29	10.95+2.10	20.66+2.18	20.90 ± 3.29	10.95+210	20.67+2.18
389		$\Delta(\uparrow)$	5.34	2.57	3.77	11.26	5.68	7.82	11.24	5.74	7.83
390		DER	4.83±0.13	5.83±0.13	7.24±0.18	7.83±0.12	7.13±0.13	$7.07_{\pm 0.18}$	7.83±0.12	7.33 ± 0.12	9.79 ± 0.18
201	Split CIFAR-100	w/ Ours	$9.59_{\pm 0.43}$	$7.59_{\pm 0.43}$	$10.34_{\pm 0.48}$	12.59±0.43	$9.59_{\pm 0.43}$	$14.44_{\pm 0.47}$	13.59±0.33	$12.59_{\pm 0.44}$	$14.34{\scriptstyle\pm0.48}$
221	•	$\Delta(\uparrow)$	4.76	1.76	3.10	4.76	2.46	7.37	5.76	5.26	4.55
392		RAR	4.46 ± 0.42	3.34 ± 0.35	$9.49_{\pm 1.72}$	8.42±2.38	4.82 ± 1.15	$10.49{\scriptstyle\pm1.48}$	7.69±0.16	7.97 ± 0.11	7.51 ± 1.05
303		w/ Ours	6.01±0.19	$5.83{\scriptstyle \pm 0.17}$	$11.29{\scriptstyle \pm 0.67}$	15.54±4.10	7.82 ± 1.78	15.71 ± 3.13	14.27±0.11	$11.32{\scriptstyle \pm 0.12}$	13.51 ± 1.04
000		$\Delta(\uparrow)$	1.55	2.49	1.80	7.12	3.00	5.22	6.58	3.35	6.00
394		ER	2.13 ± 0.94	3.11 ± 0.93	8.71±1.11	2.75±0.39	$2.63{\scriptstyle \pm 0.23}$	$8.63{\scriptstyle \pm 0.66}$	2.91±0.26	2.76 ± 0.26	6.78±1.56
395		w/ Ours	5.83±0.13	5.78±0.12	10.61±1.19	7.01±0.16	6.18±0.16	11.38±0.33	7.21±0.92	5.97±0.86	11.37 ± 1.86
000		$\Delta(\uparrow)$	3.70	2.67	1.90	4.26	3.55	2.75	4.30	3.21	4.59
396	0 P. M. 17	DER	3.42±1.31	2.53±0.11	4.21±0.69	2.02 ± 0.48	3.13±0.46	4.84±0.29	3.77±0.39	3.26±0.37	4.28±0.73
397	Split Mini-ImageNet	w/ Ours	0.58±0.71	5.55±0.73	9.73±0.02	7.00±0.77	6./4±0.77	8.19±0.09	0.42±0.18	6.12±0.31	7.05±0.37
200		$\Delta(\uparrow)$	3.10	3.02	3.52	5.04	3.01	5.55	2.05	2.80	2.78
220		KAK	5.08±0.22	5.48±0.19 7.31⊥0.12	0.29±1.02	2.23±0.46	2.99±0.45	0.30±0.86	5.29±0.97	5.15±0.89	10 53 Loco
399		$\Lambda(\uparrow)$	3.14	3.83	3 39	5.42	3.48	3.21	3.68	3 33	2.89

Table 2: Overall performance of *online* CL. The results of all methods are averaged over five runswith different random seeds and task orders.

minimum accuracy retained after learning the current task):

r

$$\operatorname{ninAcc}_{t} = \frac{1}{t-1} \sum_{i=1}^{t-1} \min_{i < j \le t} \operatorname{Acc}\left(\mathcal{D}_{i}, f_{\theta_{j}}\right),$$

$$\operatorname{WCA} = \frac{1}{T} \operatorname{Acc}\left(\mathcal{D}_{T}, f_{\theta_{T}}\right) + \left(1 - \frac{1}{T}\right) \operatorname{minAcc}_{T}.$$
(13)

Evaluation Protocol. To evaluate the practical performance of each method, we devise a novel evaluation protocol that reflects the impact of *sub-optimal* hyperparameter values. Specifically, we evaluate CL models across different hyperparameter scales, i.e., $\{\frac{1}{k}\lambda, \lambda, k\lambda\}, k > 1$, where λ represents the optimal hyperparameter values obtained from a grid search, and k is a scale factor. By varying the scale of λ , we are able to examine whether the CL model generalizes well beyond the desired conditions of hyperparameters. In practice, we set k = 4 for $\{\frac{1}{k}\lambda, \lambda, k\lambda\}$, which usually provides a sufficient and meaningful range for examination.

416 417

400 401

402 403 404

378

5.2 EXPERIMENTAL RESULTS

418 **Overall Performance**. Table 2 shows the overall performance of *online* CL, where λ represents 419 the optimal value of each hyperparameter obtained from a search on the grid. It is evident that the 420 improvement in FAA is consistently significant (e.g., ER is improved by $\{5.34\%, 11.26\%, 11.24\%\}$ 421 on Split CIFAR-100), indicating the remarkable benefit of LUNCH plugin through adaptive task 422 weighting. It also shows significant improvements in AAA and WCA (e.g., RAR is improved by {5.52%, 3.35%, 2.78%} in AAA on Split Mini-ImageNet and by {5.15%, 4.55%, 7.55%} in WCA 423 on Split CIFAR-10), suggesting that LUNCH improves the performance of deep models across the 424 training trajectory and in worst-case scenarios. Additionally, each CL method shows considerable 425 variation with different hyperparameters initializations, thus the results from our protocol can be 426 considered as the expected performance of each CL method when applied in real-world scenarios to 427 comprehensively evaluate the generalizability of CL methods. 428

Additionally, LUNCH significantly improves the performance of the baseline CL methods on different scale benchmarks. Fig. 3 shows the average results and standard deviations of five runs on Split CIFAR-100 and Split ImageNet-R. *RAR w/ Ours* consistently outperforms the original *RAR* in terms of AAA in a series of continual learning tasks. Due to the page limit, additional empirical results in

433

434

435

436

437 438

439

440

441

442 443

444 445

446

452

453

476

477

0.50 RAR RAR 0.6 RAR w/ Ours RAR w/ Ours 0.45 0.40 0.5 €^{0.35} (%) 0.4 WW _{0.3} ¥ 0.30 0.20 0.2 0.15 0.1 0.10 15.0 17.5 20.0 10.0 15.0 17.5 20.0 0.0 2.5 5.0 10.0 12.5 25 12.5 7.5 0.0 5.0 7.5 Tasks Tasks (a) AAA on Split CIFAR-100 (b) AAA on Split ImageNet-R

Figure 3: Performance curves of AAA on benchmark of different scales.

Table 3: Overall performance of *offline* CL with FAA as the evaluation metric. The results of all methods are averaged over five runs with different random seeds and task orders.

Dataset	ER	ER w/ Ours	$\Delta(\uparrow)$	DER	DER w/ Ours	$\Delta(\uparrow)$	RAR	RAR w/ Ours	$\Delta \uparrow$
Split CIFAR-10	65.26±0.48	69.49±0.99	4.23	71.53±0.49	74.15±0.38	2.62	73.52 ± 0.65	77.12±0.59	3.60
Split CIFAR-100	45.14 ± 0.65	$48.21 \pm \textbf{0.73}$	3.07	57.85±0.57	$61.07 {\scriptstyle \pm 0.43}$	3.22	60.64 ± 1.01	$63.61{\scriptstyle \pm 0.95}$	2.97

WCA and FAA are provided in Appendix B.1, and these empirical results are consistent with those in AAA.

On the other hand, Table 3 shows the mean and standard deviation of five runs after applying our approach to all baselines (e.g., up to $\sim 4\%$ in FAA) in the *offline* scenarios. From Table 2 and Table 3, we observe that *RAR w/ Ours* often achieves superior performance gains compared to other baselines, mainly because it effectively utilizes advanced data augmentation to capture useful visual information (Zhang et al., 2022). Therefore, we will use RAR as the primary baseline for further analysis in the sequel.

460 Ablation Study We conduct an ab-

lation study to demonstrate the ef-462 fectiveness of each component, in-463 cluding the uncertain hyperparam-464 eters and the temporal ensembling 465 strategy EMA. We compared several 466 variants to analyze their individual contributions. As shown in Table 4, 467 RAR w/ Unc introduces optimizable 468 uncertainty into CL-relevant hyper-469 parameters without temporal ensem-470 bling, RAR w/ EMA only incorporates 471

Table 4: Ablation study of LUNCH with RAR as baseline.

Dataset	Metric	RAR	w/ EMA	w/ Unc	w/ Ours
	FAA	42.07	46.79	49.25	55.28
Split CIFAR-10	WCA	37.41	37.81	45.81	49.97
-	AAA	58.23	60.05	59.23	61.06
	FAA	9.19	9.43	9.33	11.77
Split CIFAR-100	WCA	5.21	5.93	5.55	7.39
-	AAA	11.78	12.37	12.17	14.78
	FAA	3.21	4.74	4.43	7.42
Split Mini-ImageNet	WCA	3.26	3.32	4.69	7.12
	AAA	10.62	11.62	12.99	16.12

the temporal ensembling strategy with deterministic hyperparameters, and *RAR w/ Ours* represents the full version of LUNCH that leverages both components. It is obvious that both the learnable uncertain hyperparameters and the temporal ensembling strategy contribute to CL, and removing either results in performance degradation.

Table 5: Impact of sub-optimal hyperparameter values in terms of batch size, learning rate decay and weight decay. We report FAA averaged over five runs with different random seeds and task orders.

Dataset	ER	ER w/ Ours	$\Delta(\uparrow)$	DER	DER w/ Ours	$\Delta(\uparrow)$	RAR	RAR w/ Ours	$\Delta(\uparrow)$
Split CIFAR-10	23.11±1.21	27.38 ± 1.13	4.27	24.11±1.37	29.84 ± 0.84	5.73	33.11±1.96	38.49 ± 2.04	5.38
Split CIFAR-100	$8.74_{\pm 1.09}$	$14.99{\scriptstyle\pm1.41}$	3.25	10.31 ± 1.11	13.42 ± 0.49	3.11	11.47 ± 0.73	15.21 ± 1.12	3.74

Additionally, we evaluate the computational costs of different methods integrating LUNCH (i.e., *w/ Ours*) or not (i.e., *w/o Ours*). As shown in Fig. 4, it is evident that there is minimal difference in time between the *w/o Ours* and *w/ Ours* cases. Our method enhances performance without significantly increasing the extra computational costs, thus making it practical for CL applications with limited computational resources.

486 Sensitivity Analysis. We further extend our evalua-487 tion protocol to verify the effectiveness of LUNCH 488 in improving the robustness of other hyperparame-489 ters, such as the learning rate decay, batch size, and 490 weight decay. These hyperparameters are also relevant to CL yet in an implicit manner (Mirzadeh 491 et al., 2020; Cha & Cho, 2024). For ease of imple-492 mentation, we randomly sample a set of sub-optimal 493 hyperparameters from predefined values (i.e., mod-494 erately rescaling optimal hyperparameters in Ap-495 pendix B.2), and average multiple runs to evalu-496 ate the practical performance of each method, thus 497 avoiding overestimation of CL capabilities in prac-498 tical use. As shown in Table 5, we observe that our 499 approach achieves better performance compared to



Figure 4: Comparison of extra computational time with and without LUNCH (Ours).

the corresponding baselines, and in particular the improvement of RAR is more significant. Aligning
with the results of previous analysis, LUNCH can significantly alleviate hyperparameter sensitivity,
especially with more advanced methods.

Of note, the above results pose the particular challenge of hyperparameter sensitivity in CL. A majority of CL methods are developed with the "optimal" hyperparameter values obtained from a grid search, and thus may fail to produce desirable results based on such hyperparameter values when there exist large changes in data distribution. In contrast, our protocol of evaluating hyperparameter sensitivity can faithfully reflect this issue, and the proposed LUNCH can effectively address it as a plug-in manner.

509 510

6 DISCUSSION AND CONCLUSION

511 512

In this work, we propose an innovative approach for adaptive weighting of task contributions in CL, 513 which optimizes CL-relevant hyperparameters according to training progress and thus alleviates 514 catastrophic forgetting. In particular, we formulate each CL-relevant hyperparameter as a function 515 of learnable uncertainty under homoscedastic assumption, and ensure training stability via the ex-516 ponential moving average of network parameters along the training trajectory. Extensive empirical 517 results demonstrate the benefits of our approach in improving the effectiveness and robustness of the corresponding baselines in a variety of CL scenarios. We hope that this work will inspire more 518 explorations of adaptive weighting in CL scenarios and shed light on building general, practical, and 519 effective CL methods in this field. 520

521 This work also has some potential limitations. First, we follow the implementations of represen-522 tative online CL methods, which mainly employ ResNet-based backbone. We leave examining the 523 effectiveness of uncertain hyperparameters with ViT-based backbone as a potential future work. Sec-524 ond, we focus on CL methods within the same parameter space. Therefore, the proposed uncertain 525 hyperparameters are not applicable to many architecture-based CL methods that focus on multiple 526 parameter spaces. Since this work is essentially a fundamental research in machine learning, the 527 potential negative impact is not obvious at the current stage.

528 529

530

7 REPRODUCIBILITY

The source code of our experiments is included in the supplymentary materials. The theoretical results, including key assumptions and proofs, are provided in Appendix A. The specific hyperparameters used in our experiments can be found in Section 5.1 and the associated yml configuration files are provided in the source code. These resources allow for a comprehensive replication of our results.

536

537 REFERENCES

539 Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. In H. Wal-

540 541 542	lach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, pp. 11849–11860. 2019.
543 544	Christopher M Bishop and Nasser M Nasrabadi. <i>Pattern recognition and machine learning</i> , volume 4. Springer, 2006.
545 546 547 548	Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class- incremental continual learning into the extended der-verse. <i>IEEE transactions on pattern analysis</i> <i>and machine intelligence</i> , 45(5):5497–5512, 2022.
549 550 551	Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. <i>Advances in Neural Information Processing Systems</i> , 33:15920–15930, 2020.
552 553 554	Sungmin Cha and Kyunghyun Cho. Hyperparameters in continual learning: a reality check. <i>arXiv</i> preprint arXiv:2403.09066, 2024.
555 556	Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. In <i>International Conference on Learning Representations</i> , 2018.
557 558 559	Matthias De Lange, Gido van de Ven, and Tinne Tuytelaars. Continual evaluation for lifelong learning: Identifying the stability gap. <i>arXiv preprint arXiv:2205.13452</i> , 2022.
560 561 562	Enrico Fini, Stéphane Lathuiliere, Enver Sangineto, Moin Nabi, and Elisa Ricci. Online continual learning under extreme memory constraints. In <i>European Conference on Computer Vision</i> , pp. 720–735. Springer, 2020.
563 564 565	Shengbo Guo, Onno Zoeter, and Cedric Archambeau. Sparse bayesian multi-task learning. Advances in Neural Information Processing Systems, 24, 2011.
566 567 568	Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. The many faces of robustness: A critical analysis of out-of-distribution generalization. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 8340–8349, 2021.
569 570 571 572 573	Haeyong Kang, Rusty John Lloyd Mina, Sultan Rizky Hikmawan Madjid, Jaehong Yoon, Mark Hasegawa-Johnson, Sung Ju Hwang, and Chang D Yoo. Forget-free continual learning with winning subnetworks. In <i>International Conference on Machine Learning</i> , pp. 10734–10750. PMLR, 2022.
574 575	Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? <i>Advances in neural information processing systems</i> , 30, 2017.
576 577 578 579	Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 7482–7491, 2018.
580 581 582 583	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. <i>Proceedings of the National Academy of Sciences</i> , 114(13):3521–3526, 2017.
584 585 586	Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. <i>arXiv preprint arXiv:1610.02242</i> , 2016.
587 588	Zhizhong Li and Derek Hoiem. Learning without forgetting. <i>IEEE Transactions on Pattern Analysis</i> and Machine Intelligence, 40(12):2935–2947, 2017.
589 590 591	Baijiong Lin, Feiyang Ye, Yu Zhang, and Ivor W Tsang. Reasonable effectiveness of random weight- ing: A litmus test for multi-task learning. <i>arXiv preprint arXiv:2111.10603</i> , 2021.
592 593	Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 1871–1880, 2019.

614

631

635

- 594 Shikun Liu, Stephen James, Andrew J Davison, and Edward Johns. Auto-lambda: Disentangling dynamic task relationships. arXiv preprint arXiv:2202.03091, 2022. 596
- James L McClelland, Bruce L McNaughton, and Randall C O'Reilly. Why there are complementary 597 learning systems in the hippocampus and neocortex: insights from the successes and failures of 598 connectionist models of learning and memory. Psychological Review, 102(3):419, 1995.
- 600 Seyed Iman Mirzadeh, Mehrdad Farajtabar, Razvan Pascanu, and Hassan Ghasemzadeh. Under-601 standing the role of training regimes in continual learning. Advances in Neural Information Pro-602 cessing Systems, 33:7308–7320, 2020.
- Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012. 604
- 605 Quang Pham, Chenghao Liu, and HOI Steven. Continual normalization: Rethinking batch normal-606 ization for online continual learning. In International Conference on Learning Representations, 607 2021. 608
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: 609 Incremental classifier and representation learning. In Proceedings of the IEEE Conference on 610 Computer Vision and Pattern Recognition, pp. 2001–2010, 2017. 611
- 612 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience 613 replay for continual learning. Advances in neural information processing systems, 32, 2019.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray 615 Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. arXiv preprint 616 arXiv:1606.04671, 2016. 617
- 618 Rudy Semola, Julio Hurtado, Vincenzo Lomonaco, and Davide Bacciu. Adaptive hyperparameter 619 optimization for continual learning scenarios. arXiv preprint arXiv:2403.07015, 2024.
- 620 Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic 621 forgetting with hard attention to the task. In International Conference on Machine Learning, pp. 622 4548-4557. PMLR, 2018. 623
- 624 Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. Advances in Neural Information Processing Systems, 30, 2017. 625
- 626 Albin Soutif-Cormerais, Antonio Carta, and Joost Van de Weijer. Improving online continual learn-627 ing performance and stability with temporal ensembles. In Conference on Lifelong Learning 628 Agents, pp. 828-845. PMLR, 2023. 629
- Gido M Van de Ven, Tinne Tuytelaars, and Andreas S Tolias. Three types of incremental learning. 630 *Nature Machine Intelligence*, 4(12):1185–1197, 2022.
- 632 Liyuan Wang, Xingxing Zhang, Qian Li, Mingtian Zhang, Hang Su, Jun Zhu, and Yi Zhong. Incor-633 porating neuro-inspired adaptability for continual learning in artificial intelligence. arXiv preprint 634 arXiv:2308.14991, 2023a.
- Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual 636 learning: Theory, method and application. arXiv preprint arXiv:2302.00487, 2023b. 637
- 638 Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual 639 learning: Theory, method and application. IEEE Transactions on Pattern Analysis and Machine 640 Intelligence, 2024a.
- Zhenyi Wang, Yan Li, Li Shen, and Heng Huang. A unified and general framework for continual 642 learning. arXiv preprint arXiv:2403.13249, 2024b. 643
- 644 Elif Ceren Gok Yildirim, Murat Onur Yildirim, Mert Kilickaya, and Joaquin Vanschoren. Adaptive 645 regularization for class-incremental learning. arXiv preprint arXiv:2303.13113, 2023. 646
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. 647 In International Conference on Machine Learning, pp. 3987–3995. PMLR, 2017.

648 649	Yaqian Zhang, Bernhard Pfahringer, Eibe Frank, Albert Bifet, Nick Jin Sean Lim, and Yunzhe Jia. A simple but strong baseline for online continual learning: Repeated augmented rehearsal. <i>Advances</i>
650	in Neural Information Processing Systems, 35:14771–14783, 2022.
651	
652	
653	
654	
655	
656	
657	
658	
659	
660	
661	
662	
663	
664	
665	
666	
667	
668	
669	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
604	
605	
606	
607	
608	
699	
700	
701	

A UNCERTAINTY DERIVATION

A.1 REGRESSION TASKS

The derivation of uncertainty in regression is shown in followings.

$$p\left(y \mid f_{\theta}(x), \sigma^{2}\right) = \operatorname{Lap}\left(f_{\theta}(x), \sigma^{2}\right),$$

$$= \frac{1}{2\sigma^{2}} \exp\left(-\frac{\|y - f_{\theta}(x)\|^{2}}{\sigma^{2}}\right)$$

$$\log p\left(y \mid f_{\theta}(x), \sigma^{2}\right) = -\frac{1}{\sigma^{2}} |y - f_{\theta}(x)| - \log 2\sigma^{2},$$

$$-\log p\left(y \mid f_{\theta}(x), \sigma^{2}\right) = \frac{1}{\sigma^{2}} |y - f_{\theta}(x)| + \log 2\sigma^{2},$$

$$\propto \frac{1}{\sigma^{2}} |y - f_{\theta}(x)| + \log \sigma^{2}.$$
(14)

Thus, we can derive the joint distribution of multiple regression tasks:

$$p(y_1, y_2, \dots \mid f_{\theta}(x), \sigma_1^2, \sigma_2^2, \dots) = p(y_1 \mid f_{\theta}(x), \sigma_1^2) \cdot p(y_2 \mid f_{\theta}(x), \sigma_2^2) \cdot \dots$$

= Lap(y_1; f_{\theta}(x), \sigma_1^2) \cdot Lap(y_2; f_{\theta}(x), \sigma_2^2) \cdot \dots (15)

$$-\log p(y_1, y_2, \dots | f_{\theta}(x), \sigma_1^2, \sigma_2^2 \dots) \propto \frac{1}{\sigma_1^2} ||y_1 - f_{\theta}^1(x)||^2 + \frac{1}{\sigma_2} ||y_2 - f_{\theta}^2(x)||^2 + \log \sigma_1^2 + \log \sigma_2^2 + \dots$$
(16)

A.2 CLASSIFICATION TASKS

We use the maximum likelihood estimation as our objective function:

$$p\left(y=c \mid f_{\theta}(x), \sigma^{2}\right) = \frac{\exp\left(f_{\theta}^{c}(x)/\sigma^{2}\right)}{\sum_{c'} \exp\left(f_{\theta}^{c'}(x)/\sigma^{2}\right)},$$

$$\log p\left(y=c \mid f_{\theta}(x), \sigma^{2}\right) = \frac{1}{\sigma^{2}}f_{\theta}^{c}(x) - \log \sum_{c'} \exp\left(\frac{1}{\sigma^{2}}f_{\theta}^{c'}(x)\right).$$
(17)

For classification problems, we use the cross entropy loss function:

$$p(y = c \mid x, \theta) = -\log\left(\operatorname{Softmax}\left(y = c, f_{\theta}(x)\right)\right),$$

$$\log p(y = c \mid x\theta) = \log\left(\sum_{c'} \exp(f_{\theta}^{c'}(x))\right) - f_{\theta}^{c}(x).$$
 (18)

We can rewrite and derive the corresponding forms:

$$\log p\left(y=c \mid f_{\theta}(x), \sigma^{2}\right) = \frac{1}{\sigma^{2}} f_{\theta}^{c}(x) - \log \sum_{c'} \exp\left(\frac{1}{\sigma^{2}} f_{\theta}^{c'}(x)\right) + \frac{1}{\sigma^{2}} \log \sum_{c'} \exp\left(f_{\theta}(x)\right) - \frac{1}{\sigma^{2}} \log \sum_{c'} \exp\left(f_{\theta}(x)\right) = \frac{1}{\sigma^{2}} f^{c_{\theta}}(x) - \frac{1}{\sigma^{2}} \log \sum_{c'} \exp\left(f_{\theta}(x)\right) + \log\left(\sum_{c'} \exp\left(f_{\theta}(x)\right)\right)^{\frac{1}{\sigma^{2}}} - \log \sum_{c'} \exp\left(\frac{1}{\sigma^{2}} f_{\theta}^{c'}(x)\right) = -\frac{1}{\sigma^{2}} f^{c_{\theta}}(x) - \log\left(\sum_{c'} \exp\left(f_{\theta}(x)\right)\right)^{\frac{1}{\sigma^{2}}} - \log\left(\sum_{c'} \exp\left(\frac{1}{\sigma^{2}} f_{\theta}^{c'}(x)\right)\right)$$
(19)

$$= -\frac{1}{\sigma^2} \mathcal{L}(y = c, \theta) - \log \frac{\sum_{c'} \exp\left(\frac{1}{\sigma^2} f_{\theta}(x)\right)}{\left(\sum_{c'} \exp\left(f'_{\theta}(x)\right)\right)^{\frac{1}{\sigma^2}}}.$$

Under the assumption $\left(\sum_{c'} \exp\left(f_{\theta}^{c'}(x)\right)\right)^{\frac{1}{\sigma^2}} \approx \frac{1}{\sigma^2} \sum_{c'} \exp\left(\frac{1}{\sigma^2} f_{\theta}^{c'}(x)\right)$ that allows us to simplify the objective function:

$$-\log p\left(y=c \mid f_{\theta}(x), \sigma^{2}\right) \approx \frac{1}{\sigma^{2}} p(y=c \mid x, \theta) + \log \sigma^{2}.$$
(20)

Thus, we can derive the joint distribution of multiple classification tasks (take two tasks as examples):

$$p\left(\theta, \sigma_1^2, \sigma_2^2\right) \propto \frac{1}{\sigma_1^2} \mathcal{L}_1(\theta) + \frac{1}{\sigma_2^2} \mathcal{L}_2(\theta) + \log \sigma_1^2 + \log \sigma_2^2 + \dots,$$
(21)

$$p\left(y_1, y_2 \mid f_{\theta}(x), \sigma_1^2, \sigma_2^2\right) = \operatorname{Softmax}\left(y_1; \frac{1}{\sigma_1^2} f_{\theta}(x)\right) \cdot \operatorname{Softmax}\left(y_2; \frac{1}{\sigma_2^2} f_{\theta}(x)\right), \quad (22)$$

$$-\log p\left(y_1, y_2 \mid f_{\theta}(x), \sigma_1^2, \sigma_2^2\right) = \log \operatorname{Softmax}\left(y_1; \frac{1}{\sigma_1^2} f_{\theta}(x)\right) + \log \operatorname{Softmax}\left(y_2; \frac{1}{\sigma_2^2} f_{\theta}(x)\right)$$
(23)

$$-\log p\left(y_1, y_2 \dots \mid f_{\theta}(x), \sigma_1^2, \sigma_2^2 \dots\right) = \frac{1}{\sigma_1^2} \mathcal{L}_1(\theta) + \frac{1}{\sigma_2^2} \mathcal{L}_2(\theta) + \frac{1}{2}\log \sigma_1^2 + \frac{1}{2}\log \sigma_2^2 + \dots$$
(24)

В MORE EXPERIMENTAL DETAILS

B.1 MORE EXPERIMENTAL RESULTS ABOUT PERFORMANCE CURVES

Fig. 5 shows the average results and standard deviations of five runs in WCA and AAA on Split CIFAR-100 and Split ImageNet-R. RAR w/ Ours achieves superior performance in both WCA and FAA compared to RAR, aligning with the main results of AAA in Fig. 3.



Figure 5: The performance curves of WCA and FAA for RAR with LUNCH plugin.

810 B.2 HYPERPARAMETER SETS

All hyperparameters for the sensitivity analysis in our proposed protocol were sampled from the hyperparameter set in Table 6.

Table 6: The pre-defined set of hyperparameter values.

Hyperparameters	Value sets
Learning rate decay	[0.025, 0.05, 0.1, 0.2, 0.4]
Batch size	[32, 64, 128, 256, 512]
Weight decay	[0.0005, 0.001, 0.002, 0.004]

817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858