# SAME ACCURACY, TWICE AS FAST: CONTINUAL LEARNING SURPASSES RETRAINING FROM SCRATCH

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027 028 029

030

038

Paper under double-blind review

#### Abstract

Continual learning aims to enable models to adapt to new datasets without losing performance on previously learned data, often assuming prior data is no longer available. However, in many practical scenarios, both old and new data are accessible. In such cases, good performance on both datasets is typically achieved by abandoning the model trained on the previous data and re-training a new model from scratch on both datasets. This training from scratch is computationally expensive. In contrast, methods that leverage the previously trained model are worthy of investigation as they could significantly reduce computational costs. Our evaluation framework quantifies the computational savings of such methods while maintaining or exceeding the performance of training from scratch. We identify key optimization aspects – initialization, regularization, data selection, and hyper-parameters – that can each contribute to reducing computational costs. For each aspect, we propose effective first-step methods that already yield substantial computational savings. By combining these methods, we achieve up to 2.7x reductions in computation time across various computer vision tasks, highlighting the potential for further advancements in this area.

#### 1 INTRODUCTION

In modern deep learning, the available data tends to expand and evolve over time, often resulting in a model already trained on a large dataset (*'old data'*) to be further adapted to perform well also on a new, smaller dataset (*'new data'*). This new data typically introduces slightly different distributions, such as additional classes, new domains, or corner cases. A common approach is to retrain the model from scratch using both the old and new datasets, but as model and dataset sizes increase, it becomes computationally expensive. This highlights the need for more efficient methods that can continuously train models without the cost of full retraining.

A large part of continual learning research has tried to approach this problem in a resourceconstrained setting, aiming to reach the highest possible performance on both old and new data 040 under some constraints (Verwimp et al., 2024). Such constraints can lead to suboptimal perfor-041 mance, in part due to catastrophic forgetting (French, 1999). In e.g. industry applications, it is often 042 undesirable to sacrifice performance to save computational costs, and retraining from a randomly 043 initialized model ('from scratch') is preferred over using older models (Huyen, 2022). This is a 044 wasteful approach; there is a model available that performs well on the old data, so why not use 045 it? In practice, it has been shown to be difficult to continue to train previously trained models, even without storage constraints on past examples (Ash & Adams, 2020). 046

In this paper, we focus on the cost of training a model on new data, when a model that performs well on the old data is available. The cost of training this old model is treated as a sunk cost (Kahneman & Tversky, 1972), its training happened in the past and the price for this training has already been paid. In contrast, future costs can be reduced or mitigated, and we show that using old models can be an effective way of doing so. Typically, in computational problems there is both a memory and computational aspect, but when models get large, the computational cost tends to outweigh memory costs. For example, the hard disks to store ImageNet21k are about 50 times cheaper than training a large vision transformer on the same data once (see Appendix A.1 for details of this estimate).

054

056

060

061

062

063

064



Figure 1: Test accuracy on CIFAR100 (70+30) with a model pre-trained on 70 classes. The 'scratch' method starts from random initialization, while the 'naive' approach uses the pre-trained model without modification. 'Ours' modifies the optimization process (see Section 3) and matches the scratch performance with  $2.7 \times$  lower computational cost.

While there are good reasons to work in memory-restricted settings (Verwimp et al., 2024), in this paper the focus is on reducing the computational cost when new data becomes available. Particularly, we investigate techniques that can make a difference when an old model is available, instead of general optimization techniques like e.g. mixed precision calculations (Micikevicius et al., 2018).

Instead of the wasteful from-scratch approach, we propose to focus on continuous solutions, which leverage the existing model when new data is available. The simplest of these solutions is to continue training the model on the full dataset, using the previous model's weights for initialization. While this solution uses pre-trained initialization, it converges at a similar speed and often underperforms compared to training from scratch, offering little computational gain. Figure 1 illustrates the limitations of this naive approach.

Here, we introduce an evaluation framework to measure the computational gain of continuous methods. In this framework, we assume access to a previously trained model on the old data, and both the entire new and the old datasets. The goal is to improve the efficiency and accuracy by leveraging this previously trained model, compared to training from scratch on the entire dataset, as well as naive continuous training. We measure the advantage of an approach by the reduction in computational cost required to reach the same accuracy, on the full dataset, as a model trained from scratch.

085 In Section 3, we start from the canonical SGD update rule and show that all of its aspects – the initialization, the objective function, the sampling process, and the hyper-parameters - can signif-087 icantly reduce computational costs. We outline and evaluate several strategies - inspired by the literature – that act on these aspects. Each of them presents a promising avenue for future research to accelerate the convergence of continuous models. In Section 4, we show that while these methods 090 individually reduce computational costs, their effects are to some extent complementary; combining 091 them yields even greater reductions. Finally, we further demonstrate that these methods improve 092 training efficiency across a variety of image classification datasets, multi-task settings, and domain incremental scenarios, highlighting the robustness and potential of our method to reduce the computational burden of re-training machine learning models. 094

096 1.1 OUR CONTRIBUTION

095

098

099

102

- We propose a novel continual learning evaluation framework, allowing models full access to previous data and measuring their computational cost rather than only their accuracy. This approach shifts the research focus towards more practical, real-world challenges.
- We discuss the different areas in the optimization process that can be explored to accelerate convergence. Each area is broad enough to be the focus of a different method.
- For each area, we introduce a first-step method based on the literature, which already significantly enhances the learning speed of continual models.
- We demonstrate that improvements in these areas are complementary, meaning advancements in one area do not preclude further gains in others.
- Through an empirical study, we show that our combined methods consistently accelerate convergence across multiple datasets and tasks.

# 108 1.2 RELATED WORK

110 **Continual learning.** Continual learning concerns cases where data is not available all at once (see 111 e.g. (De Lange et al., 2021; Wang et al., 2024) for surveys). Most studies have focused on settings 112 with strong constraints on the amount of old data that can be stored (Verwimp et al., 2024). Replay methods (Chaudhry et al., 2019; Buzzega et al., 2020) aim to use this stored data as effectively as 113 possible. However, even with replay mechanisms, learning new data often leads to (catastrophic) 114 forgetting of previously learned examples. Other approaches modify the model architectures (e.g. 115 Yan et al., 2021) and regularization losses (e.g. Li & Hoiem, 2017) to reduce forgetting. In contrast, 116 some works have looked at settings where computational cost is restricted rather than memory costs. 117 In these cases, standard replay outperforms other methods (Prabhu et al., 2023). Later works (Smith 118 et al., 2024; Harun et al., 2023) improved replay techniques in these settings. Our work, however, 119 imposes no memory restrictions and instead focuses on accelerating the learning compared to models 120 trained from scratch, a challenge that is in some cases harder than outperforming standard replay on 121 a continuous model, see Section 3.1.

122

123 **Faster optimization.** At its core, machine learning involves solving a challenging and compu-124 tationally expensive optimization problem, and many approaches have been proposed to stream-125 line this process (Sun et al., 2019). First-order methods, such as stochastic gradient descent 126 (SGD)(Robbins & Monro, 1951), are the most widely used, where the full gradient is typically 127 approximated using small batches. However, SGD can suffer from slow convergence, especially in 128 high-variance settings(Johnson & Zhang, 2013), which can be mitigated by techniques like Nesterov 129 momentum (Sutskever et al., 2013). Adaptive approaches such as AdaGrad (Duchi et al., 2011) and 130 Adam (Kingma, 2014) help by adjusting the learning rate dynamically, though explicit learning rate 131 scheduling often enhances their performance (Loshchilov & Hutter, 2022; Smith & Topin, 2019). Second-order methods, despite their promise, are hampered by the computational expense of esti-132 mating the Hessian (Martens, 2016). Goyal et al. (2017) also highlights the intricate relationship 133 between batch size and learning rate in neural network optimization. Additionally, regularization 134 techniques like batch normalization (Ioffe & Szegedy, 2015) and weight normalization (Salimans & 135 Kingma, 2016) can further improve convergence. These considerations become especially impor-136 tant when optimizing from a pre-trained model, as is the case in this work, rather than from random 137 initialization (Narkhede et al., 2022).

138 139

**Warm start.** Starting from non-random initialization is most commonly used in transfer learning, 140 where a pre-trained model (typically on ImageNet) serves as a starting point to kick-start training 141 downstream tasks (Zhuang et al., 2020). Contrary to our work, these works are often not concerned 142 with performance on the pre-training (i.e. old) data. While beneficial, pre-training may hurt down-143 stream performance (Zoph et al., 2020). This may be explained by a loss of plasticity in trained 144 networks (Dohare et al., 2024; Abbas et al., 2023). Ash & Adams (2020) showed that when con-145 tinuously training on the same data source, lower performance is reached than when starting from 146 scratch. Gupta et al. (2023); Parmar et al. (2024) study how to continually pre-train NLP models 147 which is strongly related to our setting. They examine learning rate scheduling but do not consider loss of plasticity, regularization, and data selection as done in this paper. 148

149 150

## **2 PROBLEM DESCRIPTION**

151 152

We consider the following setup: an existing dataset, denoted as  $\mathcal{D}_{old} = (\mathcal{X}_{old}, \mathcal{Y}_{old})$ , where  $\mathcal{X}_{old}$ represents the input examples and  $\mathcal{Y}_{old}$  the corresponding labels. A model  $f_{old} : \mathcal{X}_{old} \to \mathcal{Y}_{old}$ , has already been trained on this dataset. We are then provided with a new data,  $\mathcal{D}_{new} = (\mathcal{X}_{new}, \mathcal{Y}_{new})$ , which may include new classes. The objective is to get a model  $f : (\mathcal{X}_{old} \cup \mathcal{X}_{new}) \to (\mathcal{Y}_{old} \cup \mathcal{Y}_{new})$  that performs well on both the new and the old datasets,  $\mathcal{D} = \mathcal{D}_{old} \cup \mathcal{D}_{new}$ , with minimal computational cost.

To ensure a fair comparison between models, we use the same architecture when comparing the computation costs of different training methods. This, and keeping a fixed batch size, allows us to quantify computational cost through the number of training iterations. Let  $f^i$  represent the model fafter *i* training iterations. We define the speed *s* of a model *f* in achieving a target accuracy, *a*, as 162 the first iteration in which f reached or surpassed that accuracy. Formally: 163

$$s(f,a) = \min_{i} \left\{ i \in \mathbb{N} : \frac{1}{|\mathcal{D}|} \sum_{j}^{|\mathcal{D}|} \mathbb{1}[f^{i}(x_{j}) = y_{j}] \ge a \right\}$$
(1)

(2)

We use the relative speed-up  $L_r$  when comparing the performance of different models against a baseline model. The baseline model  $f_{scratch}$  is trained from scratch on the combined dataset  $\mathcal{D}$  and achieves an accuracy of  $a_{scratch}$ . Then we can define  $L_r$  as:

 $L_r(f, f_{scratch}) = \frac{s(f_{scratch}, a_{scratch})}{s(f, r/100 \cdot a_{scratch})}$ 

168

174 175

176

177

178

179

181 182

183

or the relative number of iterations that a model f requires to reach the same  $(L_{100})$  or a fraction of (e.g.  $L_{99}$ ) the final accuracy of a model trained from scratch. For instance,  $L_{99} = 2$  would indicate that model f attains an accuracy of  $0.99 \cdot a_{scratch}$  with a 2 times lower computational cost than was required to train the full model from scratch to attain ascratch, or equivalently, half the number of iterations. To reduce notation complexity, we will simply use  $L_r$  in the remainder of the paper when the models can be inferred from context. Note that this measure only works when each iteration has the same computational cost, but the idea can easily be extended to when this is not the case.

#### 2.1IMPLEMENTATION DETAILS

Datasets. We conducted experiments on a variety of image classification datasets, including CIFAR-10 (Krizhevsky et al., 2009), CIFAR-100, subsets of ImageNet (ImageNet-100 and 185 ImageNet-200) (Deng et al., 2009), and Adaptiope (Ringwald & Stiefelhagen, 2021). For continuous training, each dataset was divided into disjoint subsets, and training proceeded in a cumulative 187 manner. For example, in CIFAR-100 (80+10+10), classes are split into three groups: 80 classes, 188 followed by 10, and 10. The model was trained sequentially, first on the 80 classes, then on the 189 combined 80 + 10, and finally on all 100 classes. Class splits were randomized, where the specific 190 seeds are available in the attached code.

191

199

201

203 204

205

192 **Training and baselines.** Unless otherwise specified, we used ResNet-18 with a cosine annealing 193 scheduler (Loshchilov & Hutter, 2016), the Adam optimizer, and a learning rate of 0.001. All models 194 are trained with standard cropping and horizontal flipping augmentations. Additional model and 195 hyperparameter details can be found in the attached code and Appendix A.2. The *scratch* baseline 196 indicates a model that is trained from a random initialization on both old and new data together. The *naive* baseline represents a continuous model that simply continues training from the old model, 197 without any modification.

**Experimental details.** Every experiment shown in this paper is repeated five times. The results 200 shown are the averages of these experiments, accompanied by their standard error in plots. The results presented are always obtained by using the combined test sets of the old and the new datasets. 202

#### 3 METHOD

206 In deep learning, optimization is typically performed using variants of stochastic gradient descent, 207 where model parameters  $\theta_i$  are updated by subtracting an estimate of the gradient of the objective function, based on a small batch of samples. The standard minibatch SGD update rule is: 208

$$\theta_{i+1} = \theta_i - \frac{\eta}{N} \sum_{j \in B_i} \nabla L(\theta_i, (x_j, y_j))$$
(3)

210 211 212

209

213 where  $\eta$  is the learning rate, L is the objective function, and  $B_i$  a batch of N examples sampled from the full dataset  $\mathcal{D}$ . All components of this update – model initialization ( $\theta_0$ ), batch composi-214 tion, learning rate adjustments, and modifications to the objective function – play a critical role in 215 the optimization trajectory and convergence speed. For other optimizers like e.g. Adam (Kingma,



Figure 2: **Initialization**. Naive continuous training is slower and less accurate than retraining from scratch. Re-introducing plasticity with shrink-and-perturb improves both speed and accuracy, surpassing scratch training.

Figure 3: **Objective function**. Regularizing the objective function with L2-losses is beneficial in both from scratch and continuous learning, yet the latter outperforms the former when using L2-init regularization.

2014), the simple average here would be replaced by a momentum-based average, but the idea remains the same. In the following sections, we explore how each of these elements can significantly accelerate continuous model training. CIFAR-100 (70+30) serves as a case study for comparing various strategies.

238 3.1 INITIALIZATION

226

227

228

229

230 231 232

233

234

235

236 237

251

252 253

260

261

When training models from scratch, careful random initialization offers several advantageous properties (Narkhede et al., 2022). However, during continuous training, many of these advantages are lost as training does not start from a random set of weights. Several works have shown that this potentially leads to reduced plasticity, i.e. not being able to learn new information as fast and accurately as a model trained from scratch (Ash & Adams, 2020; Dohare et al., 2023). This issue is similarly observed in continuous training, as shown in Figure 2. The naive benchmark is both worse and converges slower than retraining from scratch.

To restore plasticity during warm-start training without distribution shifts, Ash & Adams (2020) introduced the 'shrink and perturb' method. Rather than using the previously trained model's weights directly, the weights are shrunk by a factor  $\alpha$  and combined with a small portion of randomly initialized parameters  $\theta_{random}$ . The resulting initialization  $\theta_{init}$  is computed as:

$$\theta_{init} = \alpha \theta_{old} + \beta \theta_{random} \tag{4}$$

In our experiments, we used  $\alpha = 0.4$  and  $\beta = 0.001$  without tuning, as proposed by the original authors. However, a broad range of  $\alpha$  and  $\beta$  values yielded qualitatively similar results (see Appendix A.3). In Figure 2, we compare ResNet-18 models trained continuously on CIFAR-100 (70+30) with and without the shrink-and-perturb method. The results show that re-introducing plasticity can not only accelerate convergence but also potentially improve final accuracy compared to both scratch training and continuous training without it.

#### 3.2 OBJECTIVE FUNCTION AND REGULARIZATION

An alternative approach to tackle the reduced plasticity problem is to prevent it from arising in the first place by changing the objective function L. Dohare et al. (2023) proposed maintaining plasticity through continual backpropagation, while Kumar et al. (2023) introduced a simplified version that uses an L2 regularizer towards the initial random weights  $\theta_0$ , rather than towards the origin as is typically done in L2 regularization. In our setting, this involves modifying the objective function L to include this L2-init regularizer:

$$L_{reg}(\theta_i, (x_j, y_j)) = L(\theta_i, (x_j, y_j)) + \lambda \|\theta_i - \theta_0\|^2$$
(5)



Figure 4: Batch composition. 'Old / new' sam- Figure 5: Hyperparameters. Shortening the pling balances old and new examples in each learning rate scheduler allows for faster conver-281 batch, unlike the naive baseline, which uses pro- gence but at the cost of lower final accuracy. On 282 portional sampling. 'Easy / hard' sampling re- its own, changing the scheduler does not reach 283 duces the inclusion of the easiest and hardest old the required accuracy, yet when combined with 284 examples, significantly improving performance. the other aspects, it becomes important (see Sec-285 (Naive and 'old/new' results nearly overlap.)

tion 4.1)

In our experiments, we used  $\lambda = 0.01$  without tuning, as proposed by the original authors. However, 288 a broad range of  $\lambda$  values yielded qualitatively similar results (see Appendix A.4). In Figure 3, we 289 train ResNet-18 models on CIFAR-100 (70+30) and compared the results of models trained from 290 scratch with those trained continuously, both with the L2-init regularizer. The results indicate that 291 adding this regularization term accelerates the convergence of the continuously trained models, more 292 so than it improves models trained from scratch. Since regularization can also benefit models trained 293 from scratch, we use this baseline here, as well as the more standard L2 regularization (i.e.  $\theta_0 = 0$ ) which is not as effective as L2-init regularization.

#### 296 3.3 BATCH COMPOSITION 297

Estimating the full gradient of a dataset is expensive in the deep learning case, as many iterations 298 are required to reach a good solution. To overcome this, minibatches containing a small part of the 299 entire dataset are used to estimate the gradient. Typically, examples are sampled from the full dataset 300 with uniform probability. In continual learning, it is a common practice to balance examples from 301 old and new data in a batch (Rolnick et al., 2019), which either increases or decreases the sampling 302 probability of old data depending on whether there is either more or less old than new data. 303

304 In Figure 4, we show that having an equal number of new and old examples ('old / new sampling') does not improve the results compared to the naive approach, which balances old and new examples 305 in a batch according to their ratio in the full dataset (i.e. 70% old examples and 30% new ones 306 in this example). Katharopoulos & Fleuret (2018) show that the optimal sampling distribution is 307 proportional to the gradient norms of the individual examples. In Appendix A.7, we show that at the 308 very start of training the gradient norms of the old examples are indeed smaller, but after about 100 309 iterations, there is no noticeable difference on the class level, which explains the results. 310

While the ratio of old and new examples does not have an immediate effect, not all data in the replay 311 memory is equally useful. Often in continual learning, access to the old dataset  $\mathcal{D}_{old}$  is limited 312 to a small fraction. In our case, the replay memory has infinite capacity. We build on the work 313 of Hacohen & Tuytelaars (2024), who propose a sampling strategy that reduces the importance of 314 very easy and very difficult examples in the memory. More formally, they define learning speed ls 315 of a sample  $x_i$  as the epoch e in which sample is classified correctly: 316

317 318

287

295

$$ls(x_j, y_j) = \frac{1}{E} \sum_{i=1}^{E} \mathbb{1}[f^i(x) = y]$$
(6)

319 320

321 with E the total number of epochs and  $f^i$  the model at epoch i. 322

The learning speed is used to order the old examples from easy to hard, given that the necessary 323 information is recorded during training of the old model (For more details, see Appendix A.5). Using this order, we reduce the sampling probability of the 10% easiest and highest examples to one-tenth of the other examples. In the easy examples, there is no information left, and the hardest ones are never learned, and thus as useful. The results in Figure 4 show that this approach ('easy / hard sampling') is helpful and speeds up training in the continuous case. In Appendix A.5 we show that this is robust to the exact hyperparameters and that, while sacrificing some convergence speed, the easiest and hardest examples can be removed entirely.

330

# 331 3.4 LEARNING RATE SCHEDULING332

The learning rate controls the step size in stochastic gradient descent, directly influencing convergence speed. In deep learning, the learning rate typically starts high to allow faster progress toward optimal solutions and decreases gradually to avoid overshooting (Zeiler, 2012). Common scheduling strategies include 'multistep' scheduling (Zagoruyko & Komodakis, 2016), which reduces the learning rate at set intervals, and cosine annealing (Loshchilov & Hutter, 2017), which decays it more smoothly over time.

When training a model continuously, since the model is already trained on parts of the data, a more aggressive learning rate scheduling can accelerate convergence. While the learning rate still needs to decrease over time, this can happen more rapidly over fewer iterations, as the model has already learned key information. As shown in Figure 15 in the Appendix, the loss curve for the continuous model reaches a plateau much faster compared to scratch training, indicating that the model approaches a local optimum more quickly. This supports the need for a faster reduction in the learning rate (Zeiler, 2012).

346 We train ResNet-18 models on CIFAR-100 (70+30 split), using different lengths of cosine learning 347 rate schedulers. The most aggressive scheduler used only 25% of the total iterations, while others used 50% and 75%. The learning curves, depicting the mean test errors from these experiments, 348 are shown in Figure 5. The results indicate that more aggressive schedulers lead to faster conver-349 gence of the continuous model to a performance level comparable to the scratch solution. However, 350 overly aggressive scheduling hurts the final performance, as the networks fail to achieve 100% of 351 the scratch model's final accuracy. On its own, reducing the learning length is not helpful, yet when 352 combined with the other aspects, it becomes important (See Section 4.1). Similar qualitative re-353 sults were observed with multistep scheduling, see Appendix A.6. For the remainder of this paper, 354 any adjustments to the learning rate scheduler for specific methods are explicitly mentioned, and a 355 comparison to the unmodified variant is provided. 356

## 4 Results

In the previous section, we explored several key aspects of SGD optimization, and each can serve as
 the foundation for methods that reduce the computational cost of continuous training. By leveraging
 insights from various works in the literature, we proposed a method for each aspect, demonstrating
 the effectiveness of each method individually.

In this section, we begin by showing that these methods are complementary, with their combination further accelerating convergence beyond what is achieved by applying them separately. We show that although the same techniques can lead to better convergence when training from scratch, the benefit is larger in the continuous setting. We then extend this combined approach to various datasets, scenarios, and data splits, demonstrating its applicability across image classification tasks.

369 370

371

357 358

359

4.1 Ablations

In the previous section, we introduced several methods aimed at speeding up the convergence of continuous models, each targeting a different aspect of the optimization process. We now examine whether these improvements are complementary or if the benefits of one method negate the effectiveness of others. Table 1 presents results from training five ResNet-18 models and compares all combinations. While not fully cumulative, each combination offers benefits, either in convergence speed or in final accuracy. Additionally, we applied the same techniques to a model trained from scratch, to rule out the possibility that the techniques are *only* improving overall learning capabilities. 378 Table 1: Ablation results of the four different aspects studied in Section 3 under 'continuous' on 379 CIFAR100 (70+30). All speed-ups are relative to the model trained from scratch without any mod-380 ification, '/' indicates the required accuracy was not reached. Each of the aspects individually establishes a speed-up, and combining them improves the results further, albeit not fully cumulative. 381 On the right hand side of the table, 'scratch' shows the result of a model trained from scratch using 382 the same techniques. Some of them improve the training speed, but to a lesser extent than in the 383 continuous case. 384

				Co	ontinuous		Scratch		
Initialization	Regularization	Data	Scheduler	Max Acc	$L_{99}$	$L_{100}$	Max Acc	$L_{99}$	$L_{100}$
				65.26	/	/	65.92	$\times 1.33$	$\times 1.00$
~				66.43	$\times 1.49$	$\times 1.20$	65.06	/	/
	$\checkmark$			68.60	$\times 1.94$	$\times 1.66$	67.90	$\times 1.69$	$\times 1.51$
		$\checkmark$		66.65	$\times 1.60$	$\times 1.31$	66.61	$\times 1.54$	$\times 1.31$
			$\times 0.25$	64.15	/	/	63.72	/	/
$\checkmark$	$\checkmark$			68.91	$\times 2.19$	$\times 1.82$	68.01	$\times 1.75$	$\times 1.54$
$\checkmark$		$\checkmark$		66.76	$\times 2.92$	$\times 2.53$	66.61	$\times 1.54$	$\times 1.31$
	$\checkmark$	$\checkmark$		68.86	$\times 1.96$	$\times 1.69$	68.30	$\times 1.67$	$\times 1.49$
$\checkmark$	$\checkmark$	$\checkmark$		69.47	$\times 2.19$	$\times 1.84$	68.31	$\times 1.69$	$\times 1.52$
$\checkmark$	$\checkmark$	$\checkmark$	$\times 0.25$	68.26	imes 5.73	imes 5.32	63.74	/	/

The results show that while some optimization aspects have an effect on the scratch models as well 398 (most notably regularization), their influence is always stronger in the continuous case. This is 399 especially clear when shortening the learning rate scheduler, which allows the continuous model to 400 learn faster, yet greatly reduces the final accuracy of the from-scratch model. These results indicate 401 that the proposed techniques are effective in leveraging the benefits of having a model trained on 402 the old data available. Future research targeting different aspects of the continuous learning process 403 could yield further gains in training speed, as multiple optimization strategies can be effectively 404 integrated to speed up the learning. 405

#### 406 4.2 MULTIPLE TASKS 407

408 We used CIFAR100 (70+30) as a case study in Section 3, adding 30 new classes. In many continual 409 learning scenarios, new data is not only added once, but repeatedly. In Figure 6 we show that 410 applying the same aspects on consecutive tasks continues to work as in the one-task case. For each 411 task, we also train a model from scratch on all classes available up to this point. In Figure 6, the 412 yellow dots indicate when the continuous model reaches the same performance  $(L_{100})$  as the fromscratch model for that task. For every task  $L_{100} > 1$  and gets progressively larger, indicating that 413 the continuous model benefits more than it has trained for longer in the past. 414

415 These results highlight an important insight: in total, the continuous models have trained for more 416 iterations than the models that are trained from scratch, which also explains why their final accuracy 417 may surpass that of training from scratch. However, when accumulating past costs, the total cost of 418 the continuous model is significantly lower than the sum of costs for models trained from scratch for every task. 419

420 421

422

388

391 392

396 397

## 4.3 DOMAIN ADAPTATION

All previous experiments had new classes in the 'new' data, often referred to as class-incremental 423 learning in continual learning (De Lange et al., 2021). Here, we use the Adaptiope dataset (Ringwald 424 & Stiefelhagen, 2021) to show that our method also works when new data contains no new classes, 425 but the same classes from a different domain. In particular, the 'old' data consists of 123 categories 426 of product images from a shopping website and the 'new' data are real-life images of the same 427 products captured by users. The objective is to achieve strong performance across both domains, 428 ideally with faster convergence than re-training the model from scratch on both datasets. 429

Figure 7 shows how our method outperforms both training from scratch and the naive continuous 430 approach when including a new domain. The relative speed-up compared to the naive baselines is 431 smaller than in the class-incremental case, but the final accuracy is considerably improved.



442 Figure 6: Test accuracy on the full CIFAR-100 when training a continuous model with our 443 method on CIFAR100  $(50 + \sum_{i=1}^{10} 5)$ . Yellow 444 dots mark the  $L_{100}$  iteration, where the continu-445 ous model outperforms the from-scratch model. 446 The shaded area indicates maximum possible 447 performance based on data availability (e.g., 55%448 during the first task 50 + 5). 449



Figure 7: Domain adaptation results using the Adaptiope dataset. The 'old' data contains product images, the 'new' data has real life images. The test accuracy is on both domains. Although naively continuing to train is nearly as fast as our method in this setting, our final test accuracy is considerably higher.

Table 2: Results of our method across various datasets and splits. We report final test accuracy and relative speedup to reach similar performance compared to the scratch solution. The multiplier after the algorithm name indicates learning rate scheduling aggressiveness (e.g., Ours ( $\times 0.5$ ) means the minimum learning rate is achieved at half the iterations of Ours ( $\times$ 1)). '/' indicates that the accuracy level was not reached. Table (a) presents results on different datasets, while Table (b) provides an in-depth analysis of CIFAR-100 with various class splits.

(a) More image classification benchmarks

(b) Different ratios of old and new classes.

	U				. ,				
Dataset	Algorithm	Max Acc	$L_{99}$	$L_{100}$	Dataset	Algorithm	Max Acc	$L_{99}$	L <sub>100</sub>
CIFAR10 (8+2)	Scratch Continuous	83.32 81.37	×1.55 /	×1 /	CIFAR100	Scratch + L2-init Scratch ( $\times 0.5$ )	$\begin{array}{c} 67.90 \\ 67.14 \end{array}$	×1.29 /	×1.00
	Ours ( $\times$ 1) Ours ( $\times$ 0.5)	83.95 <b>84.33</b>	$\substack{\times 1.83 \\ \times 3.61}$	$\begin{array}{c} \times 1.29 \\ \times 2.92 \end{array}$	90+10	Ours $(\times 1.0)$ Ours $(\times 0.5)$	$69.42 \\ 69.07$	$\times 1.56$ $\times 2.95$	$\times 1.41 \times 2.65$
Adaptiope-PI (100+23)	Scratch	74.59	×1.14	×1	20110	Ours $(\times 0.25)$	68.42	$\times 5.05$	×4.47
	Continuous Ours ( $\times$ 1) Ours ( $\times$ 0.5)	74.93 <b>77.51</b> 76.75	$\times 1.17 \times 1.55 \times 2.36$	$1.02 \\ \times 1.42 \\ \times 2.24$	70+30	Ours ( $\times 1.0$ ) Ours ( $\times 0.5$ ) Ours ( $\times 0.25$ )	69.47 68.77 68.26	$\times 1.61 \\ \times 2.84 \\ \times 4.74$	$\times 1.45 \\ \times 2.56 \\ \times 4.34$
ImageNet-100 (80+20)	Scratch Continuous Ours (×1) Ours (×0.5)	62.34 62.21 <b>67.41</b> 66.26	$\times 1.18 \\ \times 1.23 \\ \times 2.16 \\ \times 3.39$	$\times 1 / \\ \times 1.97 \\ \times 3.09$	50+50	Ours (×1.0) Ours (×0.5) Ours (×0.25)	69.09 68.42 67.57	$\times 1.50 \\ \times 2.70 \\ \times 4.34$	×1.35 ×2.37 /
ImageNet-200 (180+20)	Scratch Continuous Ours (×1)	55.16 53.5 <b>59.62</b>	×1.18 / ×1.69	×1 / ×1.64	30+70	Ours (×1.0) Ours (×0.5) Ours (×0.25)	$68.81 \\ 68.10 \\ 66.93$	$\times 1.49 \\ \times 2.56 \\ /$	×1.31 ×2.30 /
. /	Ours $(\times 0.5)$	58.54	$\times 2.84$	$\times 2.75$					

472 473 474

475

450 451

452

453

454

455

456 457

#### 4.4 OTHER DATASETS AND SCENARIOS

476 In Section 3, we demonstrated the effectiveness of each method using the CIFAR-100 (70+30)477 dataset for consistency across experiments. Table 2a extends these results to a range of datasets, including CIFAR-10, ImageNet-100, ImageNet-200, and the product images of Adaptiope. For each 478 dataset, we compare the performance of models trained from scratch, a naive continuous model, 479 and a continuous model using our full method. We report the relative speed-up to reach 99% and 480 100% of the scratch-trained model's accuracy, along with the final accuracy. Across all the different 481 datasets we tested, our method consistently enhances the speed of convergence both for the  $L_{99}$  and 482  $L_{100}$  cases, often also surpassing the final accuracy of the scratch model. 483

Table 2b presents results for different class split ratios, including CIFAR-100 splits of (90+10), 484 (70+30), (50+50), and (30+70). Our methods consistently accelerate training, with speed-up being 485 more significant when the second data split is smaller. This aligns with intuition: when the initial 'old' data is larger, the old model already has more knowledge, requiring fewer updates when exposed to new data, compared to training from scratch, which treats both splits equally.

## 5 CONCLUSION

489

490 491

500 501

502

519

526

527

528

529 530

531

532

Throughout this paper, we have shown that the computational cost that is paid to train models on old 492 data should not be treated as something that is lost, but rather as a starting point for further training 493 whenever new data becomes available. We studied the four main components of the traditional SGD 494 update rule – initialization, objective function, data selection and hyperparameters – and showed that 495 each of them individually can contribute to faster convergence on old and new data combined when 496 starting from an old model. These aspects are thoroughly tested to be robust and effective across a 497 wide range of scenarios. Our proposals should be seen as starting points to work further towards 498 solutions that are even more effective in re-using old models, and saving computational costs across 499 the board in machine learning development and applications.

#### References

- Zaheer Abbas, Rosie Zhao, Joseph Modayil, Adam White, and Marlos C Machado. Loss of plasticity
   in continual deep reinforcement learning. In *Conference on Lifelong Learning Agents*, pp. 620–
   636. PMLR, 2023.
- Jordan Ash and Ryan P Adams. On warm-starting neural network training. Advances in neural information processing systems, 33:3884–3894, 2020.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.
- Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, P Dokania, P Torr, and M Ranzato. Continual learning with tiny episodic memories. In *Workshop on Multi-Task and Lifelong Reinforcement Learning*, 2019.
- Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory
  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.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009.
- Shibhansh Dohare, J Fernando Hernandez-Garcia, Parash Rahman, A Rupam Mahmood, and
  Richard S Sutton. Maintaining plasticity in deep continual learning. *arXiv preprint arXiv:2306.13812*, 2023.
  - Shibhansh Dohare, J Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A Rupam Mahmood, and Richard S Sutton. Loss of plasticity in deep continual learning. *Nature*, 632(8026): 768–774, 2024.
  - John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4):128–135, 1999.
- Google. Pricing Cloud Storage Google Cloud cloud.google.com. https://cloud.
   google.com/storage/pricing, 2024a. [Accessed 01-10-2024].
- Google. GPU pricing Compute Engine: Virtual Machines (VMs) —no Google Cloud cloud.google.com. https://cloud.google.com/compute/gpus-pricing, 2024b. [Accessed 01-10-2024].

548

549

550

553

554

555

563

564

565

569

574

580

540	Priva Goval Piotr Dollár Ross Girshick Pieter Noordhuis Lukasz Wesolowski Aano Kurola An-
	Triya Obyai, Tibu Donai, Ross Olisinek, Teter Robrandis, Eukasz Wesolowski, Aapo Rytola, Al-
541	drew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet
542	in 1 hour. arXiv preprint arXiv:1706.02677, 2017.
543	

- Kshitij Gupta, Benjamin Thérien, Adam Ibrahim, Mats Leon Richter, Quentin Gregory Anthony, 544 Eugene Belilovsky, Irina Rish, and Timothée Lesort. Continual pre-training of large language 545 models: How to re-warm your model? In Workshop on Efficient Systems for Foundation Models 546 @ ICML2023, 2023. URL https://openreview.net/forum?id=pg7PUJe0Tl. 547
  - Guy Hacohen and Tinne Tuytelaars. Forgetting order of continual learning: Examples that are learned first are forgotten last. arXiv preprint arXiv:2406.09935, 2024.
- Md Yousuf Harun, Jhair Gallardo, Junyu Chen, and Christopher Kanan. Grasp: a rehearsal policy 551 for efficient online continual learning. arXiv preprint arXiv:2308.13646, 2023. 552
  - Chip Huyen. Real-time machine challenges learning: and solutions, Jan 2022. URL https://huyenchip.com/2022/01/02/ real-time-machine-learning-challenges-and-solutions.html# towards-continual-learning. Online; accessed 29-August-2024.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by 558 reducing internal covariate shift. CoRR, abs/1502.03167, 2015. URL http://arxiv.org/ 559 abs/1502.03167.
- 561 Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance 562 reduction. Advances in neural information processing systems, 26, 2013.
  - Daniel Kahneman and Amos Tversky. Subjective probability: A judgment of representativeness. Cognitive psychology, 3(3):430-454, 1972.
- 566 Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning with 567 importance sampling. In International conference on machine learning, pp. 2525–2534. PMLR, 568 2018.
- Diederik P Kingma. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 570 2014. 571
- 572 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 573 Technical Report, 2009.
- Saurabh Kumar, Henrik Marklund, and Benjamin Van Roy. Maintaining plasticity in continual learn-575 ing via regenerative regularization. In Proceedings of The 3th Conference on Lifelong Learning 576 Agents, 2023. 577
- 578 Zhizhong Li and Derek Hoiem. Learning without forgetting. IEEE transactions on pattern analysis 579 and machine intelligence, 40(12):2935-2947, 2017.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. arXiv 581 preprint arXiv:1608.03983, 2016. 582
- 583 Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In In-584 ternational Conference on Learning Representations, 2017. URL https://openreview. 585 net/forum?id=Skq89Scxx.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In International Conference on Learning Representations, 2022. 588
- 589 James Martens. Second-order optimization for neural networks. University of Toronto (Canada), 2016.
- С McCallum. John Price and performance changes of computer 592 techhttps://ourworldindata.org/grapher/ nology with time. historical-cost-of-computer-memory-and-storage, 2023.

594 595 596 597	Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed precision training. In <i>International Conference on Learning Representations</i> , 2018. URL https://openreview.net/forum?id=r1gs9JgRZ.
598 599 600	Meenal V Narkhede, Prashant P Bartakke, and Mukul S Sutaone. A review on weight initialization strategies for neural networks. <i>Artificial intelligence review</i> , 55(1):291–322, 2022.
601 602 603	Nvidia. NVIDIA A100 GPUs Power the Modern Data Center — nvidia.com. https://www.nvidia.com/en-us/data-center/a100/, 2024. [Accessed 01-10-2024].
604 605 606	Jupinder Parmar, Sanjev Satheesh, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Reuse, don't retrain: A recipe for continued pretraining of language models. <i>arXiv preprint</i> <i>arXiv:2407.07263</i> , 2024.
607 608 609 610 611	Ameya Prabhu, Hasan Abed Al Kader Hammoud, Puneet K Dokania, Philip HS Torr, Ser-Nam Lim, Bernard Ghanem, and Adel Bibi. Computationally budgeted continual learning: What does mat- ter? In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 3698–3707, 2023.
612 613 614	Tobias Ringwald and Rainer Stiefelhagen. Adaptiope: A modern benchmark for unsupervised do- main adaptation. In <i>Proceedings of the IEEE/CVF winter conference on applications of computer</i> <i>vision</i> , pp. 101–110, 2021.
615 616 617	Maxime Rio. Tech Insights: A behind-the-scenes look at rolling out new GPU resources for NZ researchers — nesi.org.nz. https://www.nesi.org.nz/case-studies/
618 610	2023. [Accessed 01-10-2024].
620 621	Herbert Robbins and Sutton Monro. A stochastic approximation method. <i>The annals of mathematical statistics</i> , pp. 400–407, 1951.
622 623 624	David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. <i>Advances in neural information processing systems</i> , 32, 2019.
625 626	Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. <i>Advances in neural information processing systems</i> , 29, 2016.
627 628 629 630	James Seale Smith, Lazar Valkov, Shaunak Halbe, Vyshnavi Gutta, Rogerio Feris, Zsolt Kira, and Leonid Karlinsky. Adaptive memory replay for continual learning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 3605–3615, 2024.
631 632 633	Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In <i>Artificial intelligence and machine learning for multi-domain operations applications</i> , volume 11006, pp. 369–386. SPIE, 2019.
634 635 636	Shiliang Sun, Zehui Cao, Han Zhu, and Jing Zhao. A survey of optimization methods from a machine learning perspective. <i>IEEE transactions on cybernetics</i> , 50(8):3668–3681, 2019.
637 638 639	Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initial- ization and momentum in deep learning. In <i>International conference on machine learning</i> , pp. 1139–1147. PMLR, 2013.
640 641	A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
642 643 644 645	Eli Verwimp, Rahaf Aljundi, Shai Ben-David, Matthias Bethge, Andrea Cossu, Alexander Gepperth, Tyler L Hayes, Eyke Hüllermeier, Christopher Kanan, Dhireesha Kudithipudi, et al. Continual learning: Applications and the road forward. <i>Transactions on Machine Learning Research</i> , 2024.
646 647	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: theory, method and application. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 2024.

Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3014–3023, 2021.

- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Matthew D Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1): 43–76, 2020.
- Barret Zoph, Golnaz Ghiasi, Tsung-Yi Lin, Yin Cui, Hanxiao Liu, Ekin Dogus Cubuk, and Quoc Le. Rethinking pre-training and self-training. *Advances in neural information processing systems*, 33: 3833–3845, 2020.

## A APPENDIX

A.1 COST ESTIMATION

To train a machine learning model there is both a memory (storage) and compute cost. Here we will
work out an example for training a Vision Transformer (ViT) (Vaswani, 2017) on ImageNet 21k.
This dataset is about 1.13 terabytes. Cloud storage for common providers averages around 0.023\$
per GB each month, which would be around 30.13\$ per month (Google, 2024a). However, local
storage is much cheaper than cloud storage. Today hard disk storage is about 11\$ per TB (McCallum, 2023), which can last for 10 years or more.

ViT models were trained on 8 Nividia P100 gpus for 3.5 days (Vaswani, 2017), which cost \$1.46
per hour on Google Cloud, which would be \$981 in total for the entire training (Google, 2024b).
Buying the GPUs would be much more expensive, with prices for a single A100 GPU above \$10.000 (Nvidia, 2024). Modern A100 GPUs are about 3.5 times faster on real-life work loads (Rio, 2023), but more expensive to rent. At 4\$4.05 per hour, the total price would still be \$777.5.

685 686

687

651 652

653

654 655 656

657

658 659

660

661

662 663

664

665

670 671

672 673

## A.2 HYPERPARAMETERS

Unless specified otherwise, all experiments are trained with the Adam optimizer (Kingma, 2014), with default settings of  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and no weight decay. The starting learning rate is equal to 0.001 and the batch size is consistently 128 in all experiments. Unless specified differently, a cosine scheduler is used where the minimal learning rate 1e - 6 is reached after 39.100 iterations, which is equal to 50 epochs using the complete CIFAR100 dataset. All experiments use cropping and random horizontal flip augmentations.

694

A.3 INITIALIZATION

696

Figure 8 shows different values of  $\alpha$  in the shrink and perturb update rule. In genreal, shrinking more gives better results, although it slows down learning results at the start of learning. Shrinking enough is necessary and not doing so may lead to suboptimal accuracies because of loss of plasticity. In Figure 9 different values of  $\beta$  are tested, which add various levels of noise to the initialization. There is no visible effect of the size of this parameter, which is roughly in line with (Ash & Adams, 2020).



Figure 8: Grid test of the  $\alpha$  (shrink) hyperparam-Figure 9: Grid test of the  $\beta$  (perturb) hyperpaeter in the shrink and perturb update rule.



rameter in the shrink and perturb update rule

#### A 4 REGULARIZATION

Figure 10 shows the influence of the  $\lambda$  parameter in the L2-init regularization, which controls the strength of the regularization. A too low value will not have any effect, while setting this value too high, will lead to slower than necessary learning speeds.



Figure 10: Ablation of the  $\lambda$  parameter in the L2-init regularization.

#### A.5 BATCH COMPOSITION

Figure 11 shows how fast examples are learned during the first task of CIFAR100 (70+30). The x-axis shows the epochs, while each row indicates a single sample. The examples on the bottom are green almost from the start of training, these are the ones that are very easy: the model almost has no difficulty learning them, which also means they do not carry a lot of information. The ones on the top are almost completely red: they are never learned. These examples are equally not very useful: they are so hard the model can not learn them anyway (which may be a result of bad labeling, bad images etc.). 

In Figure 12, we ablate two of the parameters of the 'easy / hard' sampling process. c indicates what proportion of the easy and hard examples is influenced. e.g. c = 0.2 means that 20% procent of the examples is affected, or the 10% easiest and the 10% hardest. r indicates the probability that the easy and hard examples are sampled in a batch, relative to the examples in the middle. e.g. with r = 0.1, an easy or hard example is 10 times less likely to be in a minibatch. The result with r = 0 shows that we can even get rid of these examples, reducing the memory requirements for this method. When r = 0, this becomes the method proposed by (Hacohen & Tuytelaars, 2024), on which this idea is based.



Figure 11: How fast examples are learned during the initial learning phase of training on the
old data of CIFAR100 70+30 (i.e. the first 70
classes). Some examples are almost learned instantly (bottom), others never (top). The cutoffs
indicate both the 10% easiest and hardest examples.

Figure 12: Test accuracy of CIFAR100 on the 70 + 30 benchmark. Removing the easiest and hardest examples (r = 0) barely has an influence. Sampling them with a low probability (r = 0.1) allows to learn a little faster on the continuous model, without such an effect on the from scratch model.

### A.6 SCHEDULERS

Figure 13 shows an experiment on CIFAR100 (70+30), with a multistep learning rate scheduler (which reduces the learning rate by a fixed factor at fixed steps) rather than a cosine learning rate scheduler. This scheduler reaches the same conclusion, although in this case the iterations where the scheduler kicks in has a much larger influence than in the cosine scheduler case. In the continuous case, the learning rate is kept too high for too long, preventing the model to learn the last details, while the from scratch model is still learning.



Figure 13: When using multi-step schedulers the qualitative result stays the same, although cosine schedulers are better suited to this task of continuous learning.

#### A.7 GRADIENT NORMS

Figure 14 shows the average gradient norm of examples of old and new data in the CIFAR100 (70+30) baseline during training of the naive baseline. While at the very start the gradient of new data is considerably higher, this difference has completely disappeared after more than 100 iterations, which is nearly immediately considering 40.000 iterations in total. This explains mostly why it is not useful to oversample new data, which would indicate that new data is more important to learn than remembering the old data, which is not the goal.



Figure 14: Gradient norms of 'old' and 'new' data during training of the naive baseline on CI-FAR100 (70+30).

A.8 LOSS PLATEAU

Figure 15 shows how the loss of the naive solution plateaus earlier than when training from scratch, thus allowing to schedule learning rates earlier. However, it is not because such a loss plateau is reached that the final result will be better, it merely indicates that the results won't get any better when training longer than from the point the plateau is reached. In fact, as can be seen in Figure 15, the loss starts to increase again as overfitting takes place. In this sense, it might even be necessary to schedule early enough to obtain the best possible results.



Figure 15: The continuous naive baseline loss plateaus earlier than that of the from scratch baseline.