
Defining Benchmarks for Continual Few-Shot Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

In recent years there has been substantial progress in few-shot learning, where a model is trained on a small labeled dataset related to a specific task, and in continual learning, where a model has to retain knowledge acquired on a sequence of datasets. However, the field has still to frame a suite of benchmarks for the hybrid setting combining these two paradigms, where a model is trained on several sequential few-shot tasks, and then tested on a validation set stemming from all those tasks. In this paper we propose such a setting, naming it *Continual Few-Shot Learning (CFSL)*. We first define a theoretical framework for CFSL, then we propose a range of flexible benchmarks to unify the evaluation criteria. As part of the benchmark, we introduce a compact variant of ImageNet, called *SlimageNet64*, which retains all original 1000 classes but only contains 200 instances of each one (a total of 200K data-points) downscaled to 64×64 pixels. We provide baselines for the proposed benchmarks using a number of popular few-shot and continual learning methods, exposing previously unknown strengths and weaknesses of those algorithms. The dataloader and dataset will be released with an open-source license.

1 Introduction

Two capabilities vital for an intelligent agent with finite memory are *few-shot learning*, the ability to learn from a handful of data-points, and *continual learning*, the ability to sequentially learn new tasks without forgetting previous ones. Taken individually these two areas have recently seen dramatic improvements mainly due to the introduction of proper benchmark tasks and datasets used to systematically compare different methods (Chen et al., 2019; Lesort et al., 2019a; Parisi et al., 2019). For the set-to-set few-shot setting (Vinyals et al., 2016) such benchmarks include Omniglot (Lake et al., 2015), CUB-200 (Welinder et al., 2010), Mini-ImageNet (Vinyals et al., 2016) and Tiered-ImageNet (Ren et al., 2018b). For the single-incremental-task continual setting (Maltoni and Lomonaco, 2019) and the multi-task continual setting (Zenke et al., 2017; Lopez-Paz and Ranzato, 2017) the benchmarks include permuted/rotated-MNIST (Zenke et al., 2017; Goodfellow et al., 2013), CIFAR10/100 (Krizhevsky et al., 2009), and CORE50 (Lomonaco and Maltoni, 2017). However, none of those benchmarks is particularly well suited for evaluating the hybrid setting of low-data sequential streams.

One of the main reasons behind the scarce consideration of the liaison between the two settings is that these problems have been often treated separately and handled by two distinct communities. Historically the research on continual learning has focused on the problem of avoiding the loss of previous knowledge when new tasks are presented to the learner, known as *catastrophic forgetting* (McCloskey and Cohen, 1989), without paying much attention to the low-data regime. On the other hand, the research on few-shot learning has mainly focused on achieving good generalization over

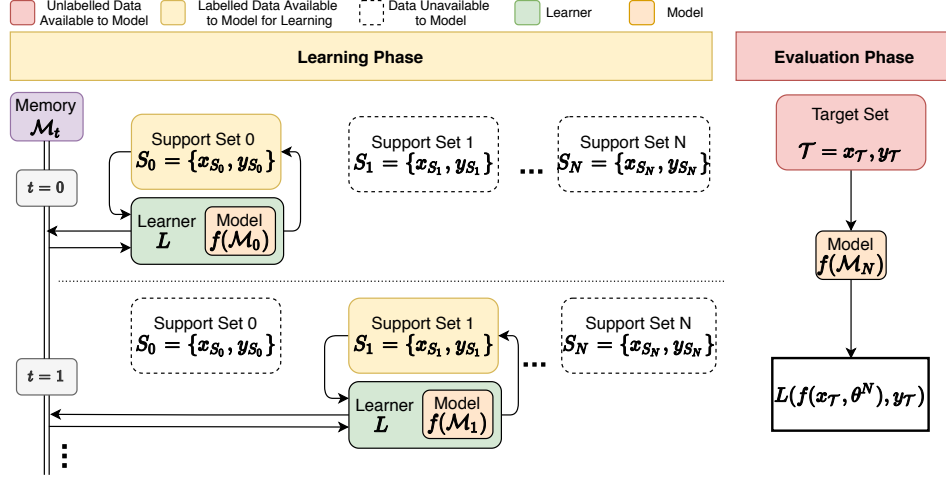


Figure 1: High level overview of the proposed benchmark. **Left block:** from the left, the learner acquires task-specific information from each set, one-by-one, without being allowed to view previous sets (memory constraint). The learner can store knowledge in a shared memory bank and use it in a classification model. On the rightmost side, future tasks are inaccessible to the learner. On the bottom, the same process is repeated on the second support set. Note that the first support set is now inaccessible. **Right block:** once the learner has viewed all support sets, it is given an evaluation set (target set) containing new examples of classes contained in the support-sets, and tasked with producing predictions for those samples. The evaluation procedure has access to the target set labels, and can establish a generalization measure for the model.

new tasks, without caring about possible future knowledge gain or loss. Scarce attention has been given to few-shot learning in the more practical continual learning scenario.

In this paper we propose to bridge the gap between the two settings by injecting the sequential component of continual learning into the framework of few-shot learning, calling this new paradigm *Continual Few-Shot Learning (CFSL)*. CFSL can be useful to the research community for a number of reasons: (i) as a framework for studying and improving the sample efficiency of mini-batch stochastic methods; (ii) as a minimal and efficient framework for studying and rectifying catastrophic forgetting providing insight into better learning dynamics, or by designing general methods to rectify the problem; (iii) as a framework for studying continual learning under memory constraints, and for testing meta-learning systems that are capable of continual learning (our settings fits on a single GPU with 11 GBs of memory). While we formally define the problem in Section 3, a high-level diagram is shown in Figure 1.

Our main contributions can be summarized as follows:

1. We formalize a highly general and flexible continual few-shot learning setting, taking into account recent considerations expressed in the literature.
2. We propose a novel benchmark and a compact dataset (SlimageNet64), releasing them under an open source license.
3. We compare recent state-of-the-art methods on our benchmark, showing how continual few-shot learning is effective in highlighting the strengths and weaknesses of those methods.

1.1 Motivation and applications

Here, we introduce the motivation behind CFSL with some practical examples. Consider typical user interfaces such as those used in online stores. The size of data points collected from each user is rather small (few-shot) and is generally stored in a sequential buffer or priority queue (continual). Suppose an underlying learning model has been deployed to enhance the user experience by suggesting new products that are likely to be of interest. This model should be able to rapidly adapt to each user (task) by accessing the sequential buffer while learning on the fly. There are multiple variants to take into account. For instance, if the user is unknown or previous data is not accessible (e.g. under privacy policies) the model has to rapidly infer the user preferences from a single task. On the other hand, if the user profile is known the model should retain knowledge about previous interactions without the

65 need of being retrained from scratch. Another example arise from human-robot interaction, where
 66 most of the applications require a robot to learn online by interacting with human teachers. For
 67 instance, in a manipulation task the human can provide a few trajectories representing the first task,
 68 then the second, the third, etc. The amount of trajectories for each task is usually rather limited and
 69 the tasks are learned sequentially. The robot should retain the knowledge of all tasks encountered so
 70 far, possibly by avoiding expensive training procedures that would overload the on-board hardware.

71 Note that neither the few-shot nor the continual setting are appropriate to deal with the aforementioned
 72 examples, since the former does not consider the sequential component and the latter does not account
 73 for the limited data size. Our CFSL formulation instead, can handle all these examples and other
 74 collateral variations, as discussed more thoroughly in Section 3.

75 2 Related Work

76 2.1 Few-Shot Learning

77 Progress in few-shot learning (FSL) was greatly accelerated after the introduction of the episodic
 78 few-shot training (Vinyals et al., 2016). This setting, for the first time, formalized few-shot learning
 79 as a well defined problem paving the way to the use of end-to-end differentiable algorithms that
 80 could be trained, tested, and compared. Among the first algorithms to be proposed there were meta-
 81 learned solutions, which here we group into three categories: metric-learning, optimization-based,
 82 hallucination (Chen et al., 2019). *Metric-learning* techniques are based on the idea of parameterizing
 83 embeddings via neural networks and then use distance metrics to match target points to support points
 84 in latent space (Vinyals et al., 2016; Edwards and Storkey, 2017; Snell et al., 2017). *Optimization-*
 85 *based* or *gradient-based* techniques are trained to perform a controlled optimization or parameter
 86 initialization to learn efficiently from a support set and generalize to a target set (Ravi and Larochelle,
 87 2016; Li et al., 2017; Finn et al., 2017; Antoniou et al., 2019; Rusu et al., 2019; Antoniou and Storkey,
 88 2019). *Hallucination* techniques utilize one or both the aforementioned methods in combination with
 89 a generative process to produce additional samples as a complement to the support set (Antoniou et al.,
 90 2017). There have been a number of methods that do not clearly fall in one of the previous categories
 91 (Santoro et al., 2017; Santurkar et al., 2018; Chen et al., 2019), including Bayesian approaches (Grant
 92 et al., 2018; Gordon et al., 2019; Patacchiola et al., 2019). For more detail, we refer the reader to the
 93 original work as well as a survey on few-shot learning (Chen et al., 2019).

94 2.2 Continual Learning

95 The problem of continual learning (CL), also called life-long learning, has been considered since
 96 the beginnings of artificial intelligence and it remains an open challenge in machine learning (Parisi
 97 et al., 2019). In standard offline supervised learning, algorithms can usually access any data point
 98 as many times as necessary during the training phase. In contrast, in CL data arrives sequentially
 99 and might only be ever seen once during the training process. Following the taxonomy of (Maltoni
 100 and Lomonaco, 2019), we group the continual learning methods into three classes: architectural,
 101 rehearsal, and regularization methods. Each category brings with it a different set of advantages and
 102 disadvantages under various resource constraints. *Architectural* approaches can be constrained on
 103 the amount of available RAM (Rusu et al., 2016; Mallya et al., 2018; Mallya and Lazebnik, 2018;
 104 Lesort et al., 2019a). Whereas, *rehearsal* strategies can become quickly bounded by the amount of
 105 available storage (Rebuffi et al., 2017; Lesort et al., 2018, 2019b). *Regularization* approaches can
 106 be free from resource constraints but incur in severe issues in the way they adapt model parameters
 107 (Kirkpatrick et al., 2017; Zenke et al., 2017; Lee, 2017; He and Jaeger, 2018; Mitchell et al., 2018).
 108 The mentioned strategies can often be intersected and combined to form even more powerful models
 109 (Rebuffi et al., 2017; Kemker et al., 2018; Maltoni and Lomonaco, 2019). Due to space constraints,
 110 we refer the reader to recent surveys on continual learning (Lesort et al., 2019c; Parisi et al., 2019).

111 2.3 Joining Meta-Learning and Continual Learning

112 Attempts to combine continual-learning with meta-learning produced a set of new research areas
 113 (Caccia et al., 2020). Continual Few-Shot Learning falls into *meta continual-learning* which can
 114 also be thought of as ‘learning to continually learn’ (Finn et al., 2019; He et al., 2019; Harrison
 115 et al., 2019). In contrast, *continual meta-learning* refers to ‘continually learning to learn’ which

attempts to make the process of meta-learning continuous as opposed to the standard meta-learning which is typically performed offline. Very recently Caccia et al. (2020) proposed a hybrid task, called Online faSt Adaptation and Knowledge Accumulation (OSAKA), linking continual-meta learning and meta continual-learning to study continual learning in the context of non-stationarity on shifting task distributions and unknown identities. Related to continual few-shot learning is the field of *incremental few-shot learning* (IFSL, (Gidaris and Komodakis, 2018; Ren et al., 2018a)). In contrast to standard few-shot learning and our work, in IFSL target set is composed of ‘novel’ classes (drawn from a never-seen-before dataset) as well as classes seen during the meta-learning phase (called ‘base classes’), and it does not consider continual updates to the novel classes. These methods are significantly different in terms of training and testing procedures. For this reason, we will not analyze this line of research any further.

Inconsistencies in the evaluation protocol In the literature, there are no established benchmarks that integrates few-shot and continual learning. Related tasks were introduced to prove the efficacy of a given system, making such tasks very restricted in terms of what methods they are applicable on, and how many aspects they could investigate. We found that tasks and datasets vary from paper to paper, making it challenging to know the actual performance of a given algorithm in comparison to others. For instance, the method proposed by Vuorio et al. (2018) has been tested exclusively on variants of MNIST. The method of Javed and White (2019) has been tested on Omniglot and incremental sine-waves. Spigler (2019) evaluated on MNIST, Le et al. (2019) on CIFAR100 and permuted MNIST, and Beaulieu et al. (2020) on Omniglot. It is evident how the problem of continual few-shot learning is not well defined, making it challenging to benchmark and compare the performance of algorithms.

3 Continual Few-Shot Learning

3.1 Definition of the problem

A continual few-shot learning (CFSL) task is composed of a sequence \mathcal{G} of small training sets (support sets) $\mathcal{G} = \{\mathcal{S}_n\}_{n=1}^{N_G}$, and a small evaluation set (target set) \mathcal{T} . Each support set $\mathcal{S} = \{(\mathbf{x}_n, y_n)\}_{n=1}^{N_S}$ is a set of input-label pairs just like in the standard few-shot learning setup (Chen et al., 2019). A target set $\mathcal{T} = \{(\mathbf{x}_n, y_n)\}_{n=1}^{N_T}$ is a set of input-label pairs containing previously unseen instances of classes stemming from \mathcal{G} . The objective of the learner is to perform well on the validation set \mathcal{T} having only temporal access to the labeled data contained in the support \mathcal{S} .

The size of the support set N_S is defined by the number of classes N_C (way) and by the *number of samples per class* K (shot), such that if we have a 5-way/1-shot setup we end up with $N_S = N_C \times K = 5 \times 1 = 5$. The *Number of Support Sets Per Task* (NSS) parameter determines N_G , the cardinality of \mathcal{G} . A *Class-Change Interval* (CCI) parameter dictates how often the classes within the sequence \mathcal{G} should change, expressed in numbers of support sets. This corresponds to assigning the elements in the support sets to a series of disjoint class sets $\bigcap_{i=1}^I \mathcal{C}_i = \{\emptyset\}$, where $I = \lceil NSS/CCI \rceil$ and \mathcal{C}_i is a set of unique classes of size N_C . For example, if $CCI=2$ then we will draw support sets whose classes change every 2 samples. As a result, support sets \mathcal{S}_1 and \mathcal{S}_2 will contain different instances of the same class set \mathcal{C}_1 , whereas \mathcal{S}_3 and \mathcal{S}_4 will contain different instances from the class set \mathcal{C}_2 . The process of generating CFSL tasks is also described in Algorithm 1 (Appendix E, supplementary material).

A *learner* is a process which extracts task-specific information and distills it into a classification model. The model can be generically defined as a function $f(\mathbf{x}, \theta)$ parameterized by a vector of weights θ . At evaluation time the learner is tested through a loss function

$$\mathcal{L} = \left(f(\mathbf{x}_{\mathcal{T}}, \theta), y_{\mathcal{T}} \right), \quad (1)$$

where $\mathbf{x}_{\mathcal{T}}$ and $y_{\mathcal{T}}$ are the input-output pairs belonging to the target set. Note that we intentionally provided a definition that is generic enough to fit into different methodologies and not restricted to the use of neural networks.

To remove the possibility of converting a continual learning task to a non-continual one, we introduce a restriction, which dictates that a support set \mathcal{S} is sampled from \mathcal{G} without replacement, and deleted once it has been used by the learner. The learner should never have access to more than one support set at a time, and should not be able to review a support set once it has moved to the next one. This restriction induces a strict sequentiality in the access of \mathcal{G} .

167 The setup we have described so far is very flexible, and it allows us to define a variety of different
 168 tasks and therefore to target different problems. In the following section we provide a description of
 169 those tasks and show that they are consistent with the continual learning literature.

170 3.2 Task Types

171 In this section we define an empirical procedure under the form of specific task types. Our CFSL
 172 tasks fully cover the standard single-incremental task scenario (Lomonaco and Maltoni, 2017) while
 173 introducing an additional, super-class NI setting consistent with Domain-Incremental learning (van de
 174 Ven and Tolias, 2018). Specifically, task A, B, and D are equivalent to *New Instances (NI)*, *New*
 175 *Classes (NC)*, and *New Instances and Classes (NIC)* (Maltoni and Lomonaco, 2019), respectively.
 176 Task C captures the super-set NI setting where instances are sampled across super-classes, instead
 177 of being sampled from previously defined class categories. This task is most similar to Domain-
 178 Incremental learning (van de Ven and Tolias, 2018). A detailed comparison between our tasks and
 179 existing work is reported in Appendix D (supplementary material). Figure 2 showcases a high-level
 180 visual representation of the proposed tasks.

181 **A New Samples:** In this task type, each support set within a given task are sampled from the
 182 same set of preselected classes. As a result, each support set will share the same classes
 183 and labels but will contain previously unseen samples of those classes. To achieve this, we
 184 can set CCI to be equal to (or higher than) the number of support sets in a given task ($CCI \geq NSS$). For every support set we sample new instances from the same classes as seen in
 185 previous support sets of the same task. **Motivation:** The standard supervised regime can be
 186 considered as a particular case of this task, which can be used for studying mini-batch
 187 stochastic optimization models as well as meta-learning efficient algorithms for doing so.

189 **B New Classes:** In this task type, each support set has different set of classes from the other
 190 support sets within a given task ($CCI = 1$). Each class within each support set has a
 191 corresponding unique output unit in the model. **Motivation:** Class-incremental learning
 192 is a particular case of this task. Since this setting allows expanding the number of class
 193 descriptors, the agent is not forced to explicitly rewrite previous knowledge at the class-level,
 194 but it may have to rewrite representations at lower-levels.

195 **C New Classes with Overwrite:** This task is similar to task B in that each support set has
 196 different set of unique classes ($CCI = 1$). The difference is that each class in the overall task
 197 is grouped and assigned a new label by *overwriting* the true label. As a result the number
 198 of output units in the model is equivalent to the number of unique classes within a single
 199 support set $\tilde{\mathcal{C}}$. Intuitively, $\tilde{\mathcal{C}}$ could be the hierarchical categories of classes in $\mathcal{G}^y = \bigcup_{n=1}^{N_G} \mathcal{S}_n^y$,
 200 however, in our experiments we assign the hierarchical categories arbitrarily. **Motivation:**
 201 This setting emulates situations where an agent is tasked with learning data-streams while
 202 being limited in storing knowledge in a preset number of output labels. As a result the agent
 203 has to learn super classes, rewriting previous knowledge at the class-level.

204 **D New Classes with New Samples:** Combines tasks A and B, where each support set contains
 205 different instances of the same set of classes for some predefined class change interval
 206 ($1 < CCI < NSS$). The CCI defines when the classes (and labels) change to a different and
 207 disjoint set. **Motivation:** This setting emulates situations where an agent is tasked with
 208 both learning new class descriptors and updating such descriptors by observing new class
 209 instances. It sheds light on how agents can perform on a setting that mixes all previous
 210 settings into one.

211 3.3 Metrics

212 **Test Generalization Performance.** A proposed model should be evaluated on at least the test sets
 213 of Omniglot and SlimageNet, on all the tasks of interest. This is done by presenting the model with a
 214 number of previously unseen continual tasks sampled from these test sets, and then using the target
 215 set metrics as the task-level generalization metrics. To obtain a measure of generalization across the
 216 whole test set the model should be evaluated on a number of previously unseen and unique tasks.
 217 The mean and standard deviation of both accuracy and performance should be used as generalization
 218 measures to compare models.



Figure 2: Visual representation of the four continual few-shot task types. Each row corresponds to a task with Number of Support Sets, NSS=4, and a defined Class-Change Interval (CCI). Given a sequence of support sets, S_n , the aim is to correctly classify samples in the target set, T . Colored frames correspond to the associated support set labels.

219 **Multiply-Addition operations (MACs)** Measures the computational expense of both the learner
 220 and the model operations during learning and inference time. In other words, it measures the memory
 221 footprint that the model itself needs to execute during a cycle of inference, or a meta-learning cycle.

222 **Across Task Memory (ATM).** Even though we have imposed a restriction on the access to \mathcal{G} , the
 223 learner is still authorized to store in a local *memory bank* \mathcal{M} some representations of the inputs and/or
 224 output vectors (often implemented as embedding vectors or inner loop parameters)

$$\mathcal{M} = \{(\hat{\mathbf{x}}, \hat{y})_{S_1}, \dots, (\hat{\mathbf{x}}, \hat{y})_{S_{N_G}}\}, \quad (2)$$

225 where $\hat{\mathbf{x}}$ and \hat{y} are representations of \mathbf{x} and y obtained after a given learner has processed \mathbf{x} and y
 226 and stored some of their useful components. Most learners will be compressing a given support set,
 227 but this is not strictly the case.

228 The potential compression rate is not directly correlated to the complexity of the model (e.g. number
 229 of parameters, FLOPs, etc). For instance, compression can be achieved by removing some of the
 230 dimensions of the input, or by using a lossless data compression algorithm, which may not require
 231 additional parameters or may have minimal impact on the execution time. In this regard, the concept
 232 of memory bank \mathcal{M} helps to disambiguate model complexity from any additional memory allocated
 233 for compressed representations of inputs. We can use the cardinality of \mathcal{M} , indicated as $|\mathcal{M}|$, to
 234 quantify the learner efficiency. Given two learners with their corresponding models $f(\mathbf{x}, \theta_1)$ and
 235 $f(\mathbf{x}, \theta_2)$, and assuming that the size of θ_1 is equal to the size of θ_2 with $\mathcal{L}_1 = \mathcal{L}_2$, then the learner
 236 with smaller cardinality $|\mathcal{M}|$ must be preferred.

237 In order to compare performances across different tasks and datasets, we relate the size of the stored
 238 task-specific representations (in bytes) $\mathcal{M}^{\hat{\mathbf{x}}}$ (e.g. embedding vectors in ProtoNets, and inner loop
 239 parameters for MAML) during task-specific information extraction to the size of vectors (in bytes) \mathbf{x}
 240 contained in the episode $\mathcal{G}^x = \cup_{n=1}^{N_G} S_n^x$. Recall that $\hat{\mathbf{x}}$ is a compressed version of \mathbf{x} and therefore
 241 $F < H$, where H and F are the vectors of dimensionality of \mathbf{x} and $\hat{\mathbf{x}}$, respectively. To reduce the
 242 notation burden we have only considered the inputs \mathbf{x} and not the targets y , since \mathbf{x} is significantly
 243 larger than y . Based on these considerations we define Across-Task Memory (ATM)

$$\text{ATM} = \frac{|\mathcal{M}^{\hat{\mathbf{x}}}|}{|\mathcal{G}^x|}, \quad (3)$$

244 where $\mathcal{M}^{\hat{\mathbf{x}}}$ is the stored representation of a series of support sets and \mathcal{G}^x is the size of the support
 245 sets. For each utilized floating point arithmetic unit we include a computation that takes into account
 246 the floating point precision level. For example, if both $\mathcal{M}^{\hat{\mathbf{x}}}$ and \mathcal{G}^x use the same floating point
 247 standard then it is divided out, but if the representational form uses a lower precision than the actual
 248 data-points then it becomes compressive. From a practical standpoint (image classification), the ATM
 249 can be estimated relating the total number of bytes stored in the memory bank (ATM numerator)
 250 with the total number of bytes over all the images in the episode (ATM denominator). Given the
 251 definition above: $\text{ATM} < 1$ for learners with efficient memory, $\text{ATM} = 0$ for learners with no
 252 memory, and $\text{ATM} > 1$ for learners with inefficient memory. Note that the ATM is undefined for
 253 empty episodes $\mathcal{G} = \{\emptyset\}$. ATM is task/dataset agnostic and can be used to compare various models
 254 (or the same model) across different settings.

3.4 Choice of datasets

The proposed CFSL benchmark has to be supplemented with appropriate datasets. The choice of these datasets should be considered decoupled from the benchmark itself. However, building some of the tasks requires datasets that meet specific desiderata: very high degree of diversity in terms of classes, high number of categories, fair (but not overabundant) number of samples per class, contained in size and memory footprint. The popular Omniglot (Lake et al., 2015) is a good first choice for a lower-difficulty dataset. Finding a higher complexity dataset (e.g. with RGB images) is arduous since existing datasets do not satisfy all the desiderata, being insufficient for robust and extensive benchmarking of CFSL algorithms (see detailed comparison in Appendix B).

SlimageNet64. To solve this issue, we propose a new variant of ImageNet64 \times 64 (Chrabaszcz et al., 2017), named *SlimageNet64* (derived from Slim and ImageNet) consisting of 200 instances from each of the 1000 object categories of the ILSVRC-2012 dataset (Krizhevsky et al., 2012; Russakovsky et al., 2015), for a total of 200K RGB images with a resolution of 64×64 pixels. In Appendix B (Table 2) we report a detailed comparison of all the datasets available, showing how SlimageNet64 is an optimal choice over multiple criteria.

SlimageNet64 vs Tiered-ImageNet. The closest alternative to SlimageNet64 is Tiered-ImageNet (Ren et al., 2018b), a subset of ILSVRC-12 with a total of 608 classes. We have found three main issues with Tiered-ImageNet: (i) it is not provided as a stand-alone archive but aggregated via a generating script (different scripts have been released leading to an incompatibility between versions); (ii) it has just 608 classes, meaning a reduced diversity in the distribution space; (iii) it has a large memory footprint (29 GB size) that makes it impractical for training a large number of models in a reasonable time. On the other hand, SlimageNet64 contains more classes (1000) and it has a lower memory footprint (9 GB size) due to the smaller resolution of the images. This makes SlimageNet64 more compact since both dataset and model can be easily loaded on a single GPU.

4 Experiments

For the purposes of establishing baselines in the CFSL tasks outlined in this paper we chose to use eight existing methods: (1) randomly initializing a convolutional neural network, and fine tuning on incoming tasks; (2) pretraining a convolutional neural network on all training set classes, then fine-tune on sequential tasks (Chen et al., 2019); (3) Elastic Weight Consolidation (EWC, Kirkpatrick et al. 2017) on a network which has been initialized and then finetuned (EWC-Init) as a baseline for CL methods; (4) EWC on a network which has been pretrained and finetuned (EWC-Pretrain); (5) Prototypical Networks (Snell et al., 2017) as baseline for metric-based FSL methods; (6) the Improved Model Agnostic Meta-Learning or MAML++ L (Low-End) model (Antoniou et al., 2019) as baseline for optimization based FSL methods; (7) MAML++ H (High-End) model (Antoniou and Storkey, 2019) (dense-net backbone, squeeze excite attention, mid-tier baseline); and (8) the Self-Critique and Adapt model (SCA) (Antoniou and Storkey, 2019), as a state-of-the-art algorithm for FSL (high-tier baseline). Additional details are reported in Appendix A (supplementary material). Table 1 and Figure 3 show the results obtained in the experimental section.

Results: accuracy. (i) *Omniglot*. The results on Omniglot in the New Classes without Overwrite Setting (B) MAML++ Low-End is inferior to ProtoNets, whilst in the New Classes with Overwrite Settings (C) this result is reversed. We infer that embedding-based methods are better at retaining information from previously seen classes, assuming that each new class remains distinct. However, when overwriting is enabled this trend is overturned because ProtoNet prototypes are shared by a number of super-classes containing categories that are harder to semantically disentangle. Gradient based methods such as MAML++ dominate in this case, since they can update their weights towards new tasks, and therefore achieve a better disentanglement of super-classes. SCA and High-End MAML++ (which utilize both embeddings and gradient-based optimization) produce the best performance across all settings. In the New Samples Setting (A), gradient based methods tend to outperform embedding-based methods while hybrid methods produce the best results. Furthermore, in the New Classes and Samples Setting (D), embedding-based methods outperform gradient-based methods, whilst hybrid methods continue to produce the best performing models. Regarding CL methods, we notice that the performance of EWC is close to the baseline conditions (Init & Tune, Pretrain & Tune), showing that EWC is not particularly effective in the CFSL setting. (ii) *SlimageNet*. Results on SlimageNet show that ProtoNets consistently outperform Low-End MAML++, even in New Classes with Overwrite

Table 1: Accuracy and standard deviation (percentage) on the test set for the proposed benchmarks. Best in bold.

Task Type	B	C	A	D	B	C	A	D	B	C	A	
NSS	3	3	3	4	5	5	5	8	10	10	10	
CCI	1	1	3	2	1	1	1	2	1	1	10	
Overwrite	False	True	True	False	False	True	True	False	False	True	True	
Omniglot	Init & Tune	10.87±0.01	27.51±0.01	44.76±0.01	8.74±0.01	6.15±0.01	24.52±0.01	45.30±0.01	3.93±0.01	3.12±0.01	22.16±0.01	45.64±0.01
	Pre. & Tune	9.97±0.14	26.75±0.27	32.44±0.29	7.91±0.15	6.02±0.02	24.51±0.06	31.89±1.10	3.86±0.06	3.13±0.03	22.30±0.06	33.17±0.39
	EWC-Init	10.75±0.00	27.76±0.00	42.55±0.00	8.72±0.00	6.23±0.00	24.49±0.00	42.18±0.00	3.91±0.00	3.02±0.00	22.22±0.01	41.82±0.00
	EWC-Pre.	9.55±0.00	26.23±0.00	30.27±0.00	7.73±0.00	5.93±0.00	24.20±0.00	30.15±0.00	3.71±0.00	3.23±0.02	22.10±0.00	30.14±0.00
	ProtoNets	95.30±0.12	45.44±0.19	98.73±0.02	48.98±0.03	91.52±0.20	35.10±0.09	98.73±0.12	48.44±0.03	83.72±0.19	27.39±0.17	98.65±0.14
	MAML++L	38.18±0.14	46.12±0.15	99.38±0.07	28.87±0.07	22.69±0.07	35.76±0.14	99.41±0.04	14.29±0.05	11.30±0.02	27.82±0.03	99.44±0.01
	MAML++H	96.14±0.02	96.77±0.08	99.73±0.04	49.44±0.02	92.70±0.03	93.47±0.05	99.80±0.01	49.00±0.04	85.56±0.10	86.38±0.14	99.86±0.01
	SCA	96.84±0.04	97.38±0.02	99.82±0.01	49.71±0.01	93.81±0.02	94.08±0.45	99.88±0.03	49.51±0.01	86.07±0.03	87.29±0.19	99.88±0.01
SlimageNet64	Init & Tune	8.4±0.01	21.3±0.01	24.4±0.01	6.1±0.01	4.5±0.01	20.8±0.01	24.7±0.01	3.0±0.01	2.4±0.01	20.5±0.01	24.9±0.01
	Pre. & Tune	8.7±0.03	21.9±0.11	24.2±0.17	6.4±0.01	4.9±0.02	21.2±0.05	24.5±0.23	3.3±0.03	2.7±0.03	20.7±0.10	24.4±0.20
	EWC-Init	8.51±0.00	21.66±0.00	24.81±0.00	6.11±0.00	4.60±0.00	20.98±0.00	24.87±0.00	2.93±0.00	2.38±0.00	20.39±0.00	24.67±0.00
	EWC-Pre.	8.83±0.09	22.21±0.04	23.47±0.09	6.42±0.08	5.18±0.01	21.17±0.06	23.44±0.35	3.26±0.02	2.66±0.01	20.74±0.04	24.40±0.42
	ProtoNets	24.1±0.05	25.9±0.23	43.1±0.24	15.1±0.03	18.2±0.14	22.7±0.09	43.3±0.03	10.4±0.12	12.3±0.09	21.0±0.06	43.7±0.15
	MAML++L	13.6±0.04	25.5±0.23	42.7±0.10	10.2±0.11	7.9±0.13	22.6±0.03	43.0±0.12	5.0±0.08	3.6±0.14	20.8±0.09	43.0±0.42
	MAML++H	27.2±0.25	33.8±0.16	61.2±0.36	16.8±0.18	21.0±0.21	30.4±0.51	68.6±0.47	12.3±0.11	14.4±0.12	25.7±0.10	75.6±0.10
	SCA	27.9±0.16	34.0±0.23	65.3±0.15	17.3±0.07	22.0±0.18	30.1±0.36	72.0±0.36	12.7±0.08	14.6±0.07	26.3±0.13	77.4±0.06

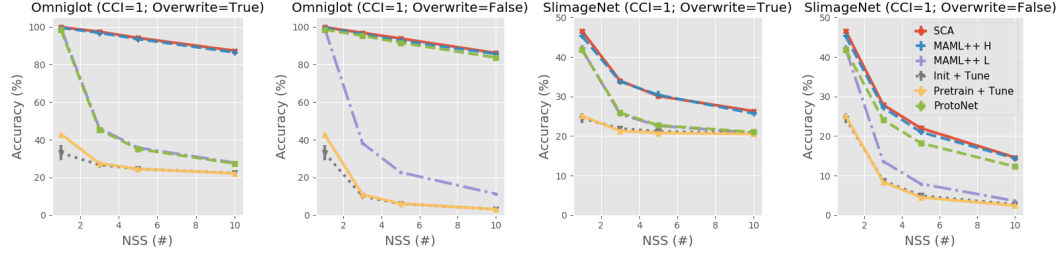


Figure 3: Accuracy (percentage) on the Omniglot and SlimageNet test set for different values of Number of Support Sets Per Task (NSS), Class Change Interval (CCI) equal 1, and with/without overwrite.

(C) where it was previously inferior. This might indicate that in SlimageNet retaining information about previously seen tasks is more important than disentangling complicated super-classes. Overall models that use both embedding-based and gradient-based methods, seem to outperform methods that do just one of the two. In New Classes and Samples (D), embedding-based methods outperform gradient-based ones by a significant margin, while hybrid approaches consistently generate the best performing models. Interestingly, in New Samples (A) the embedding-based and gradient-based methods produce very similar results, whereas in Omniglot gradient-based methods dominated.

Results: memory. Figure 4 in Appendix C (supp. material) shows the ATM and MAC costs for $NSS \in [1, 640]$. ProtoNets are the most efficient by two orders magnitude. Other methods (e.g. Low-End MAML++) starts off cheaper but increase with time, due to accumulation of features for each new learned class. In terms of ATM it is worth noting that methods such as MAML++ H and SCA tend to become incrementally cheaper than MAML++ L as the number of support sets increases. Whereas in terms of MACs MAML++ H and SCA are the most expensive by an order of magnitude or more compared to MAML++ L and ProtoNets.

5 Conclusion

In this paper, we have introduced a novel benchmark for Continual Few-shot Learning and a new variant of ImageNet, called SlimageNet64, that contains all of 1000 ImageNet classes, but only 200 samples from each class, downscaled to 64×64 . Furthermore, we have run experiments on the proposed benchmark, utilizing a number of popular few-shot learning models and baselines. We have found that embedding-based models tend to perform better when incoming tasks contain different classes from one another, whereas gradient-based methods tend to perform better when the task-classes form super-classes of randomly combined categories. Methods utilizing both embedding-based and gradient-based methods (i.e. High-End MAML++ and SCA) outperform methods that use either of the two. In conclusion, we are confident that the proposed benchmark and dataset, will help increasing the rate of progress and the understanding of the behavior of systems trained in a continual and data-limited setting.

Broader Impact

The main motivation of this work has been to propose a benchmark for the continual few-shot learning setting. The ability to learn online from a reduced amount of data is crucial if we want to have systems that are able to deal with concrete real-world problems. Applications include (but are not limited to): online classification under constrained computational resources, learning in robotic systems, facial identification (with obstructions and continual features change), learning under privacy policy restrictions (with data removed or encrypted after a short period).

These examples of applications present evident challenges which require to be handled carefully. While we have provided the framework and clear metrics, it is also important for researchers to train their methods on domain-specific data targeted on the final application. In particular, it is necessary to consider if the data and the system are biased towards some categories. In our proposed dataset (SlimageNet64) there may be implicit biases that have not been investigated enough. We think that our benchmark can be of benefit to the research community in general, but we invite the individual researchers to consider the aforementioned problems before the final deployment.

References

- Antoniou, A., Edwards, H., and Storkey, A. (2019). How to train your MAML. *In International Conference on Learning Representations*.
- Antoniou, A. and Storkey, A. (2019). Learning to Learn by Self-Critique. *Neural Information Processing Systems, NeurIPS*.
- Antoniou, A., Storkey, A., and Edwards, H. (2017). Data Augmentation Generative Adversarial Networks. *arXiv preprint arXiv:1711.04340*.
- Beaulieu, S., Frati, L., Miconi, T., Lehman, J., Stanley, K. O., Clune, J., and Cheney, N. (2020). Learning to Continually Learn. *arXiv preprint arXiv:2002.09571*.
- Caccia, M., Rodriguez, P., Ostapenko, O., Normandin, F., Lin, M., Caccia, L., Laradji, I., Rish, I., Lacoste, A., Vazquez, D., and Charlin, L. (2020). Online fast adaptation and knowledge accumulation: a new approach to continual learning. *arXiv preprint arXiv:2003.05856*.
- Chen, W.-Y., Liu, Y.-C., Kira, Z., Wang, Y.-C., and Huang, J.-B. (2019). A Closer Look at Few-Shot Classification. *In International Conference on Learning Representations*.
- Chrabaszcz, P., Loshchilov, I., and Hutter, F. (2017). A Downsampled Variant of Imagenet as an Alternative to the CIFAR datasets. *Computing Research Repository*.
- Edwards, H. and Storkey, A. (2017). Towards a neural statistician. *In International Conference on Learning Representations*.
- Finn, C., Abbeel, P., and Levine, S. (2017). Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks. *arXiv preprint arXiv:1703.03400*.
- Finn, C., Rajeswaran, A., Kakade, S., and Levine, S. (2019). Online Meta-Learning. *arXiv preprint arXiv:1902.08438*.
- Gidaris, S. and Komodakis, N. (2018). Dynamic Few-Shot Visual Learning without Forgetting. *In Computer Vision and Pattern Recognition*.
- Goodfellow, I. J., Mirza, M., Xiao, D., Courville, A., and Bengio, Y. (2013). An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*.
- Gordon, J., Bronskill, J., Bauer, M., Nowozin, S., and Turner, R. (2019). Meta-learning probabilistic inference for prediction. *In International Conference on Learning Representations*.
- Grant, E., Finn, C., Levine, S., Darrell, T., and Griffiths, T. (2018). Recasting gradient-based meta-learning as hierarchical bayes. *arXiv preprint arXiv:1801.08930*.
- Harrison, J., Sharma, A., Finn, C., and Pavone, M. (2019). Continuous Meta-Learning without Tasks. *arXiv preprint arXiv:1912.08866*, pages 1–21.

381 He, X. and Jaeger, H. (2018). Overcoming Catastrophic Interference using Conceptor-Aided Back-
382 propagation. *International Conference on Learning Representations*.

383 He, X., Sygnowski, J., Galashov, A., Rusu, A. A., Teh, Y. W., and Pascanu, R. (2019). Task Agnostic
384 Continual Learning via Meta Learning. *arXiv preprint arXiv:1906.05201*.

385 Hsu, Y.-C., Liu, Y.-C., Ramasamy, A., and Kira, Z. (2018). Re-evaluating Continual Learning
386 Scenarios: A Categorization and Case for Strong Baselines. *Continual Learning Workshop, 32nd
387 Conference on Neural Information Processing Systems (NIPS 2018)*.

388 Javed, K. and White, M. (2019). Meta-learning representations for continual learning. *arXiv preprint
389 arXiv:1905.12588*.

390 Kemker, R., McClure, M., Abitino, A., Hayes, T. L., and Kanan, C. (2018). Measuring catastrophic
391 forgetting in neural networks. In *Association for the Advancement of Artificial Intelligence*.

392 Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan,
393 J., Ramalho, T., Grabska-Barwinska, A., et al. (2017). Overcoming catastrophic forgetting in
394 neural networks. *Proceedings of the national academy of sciences*.

395 Krizhevsky, A., Hinton, G., et al. (2009). Learning multiple layers of features from tiny images.
396 Technical report, Citeseer.

397 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolu-
398 tional neural networks. In *Neural Information Processing Systems*.

399 Lake, B. M., Salakhutdinov, R., and Tenenbaum, J. B. (2015). Human-level concept learning through
400 probabilistic program induction. *Science*.

401 Le, C., Wei, X., Wang, B., Zhang, L., and Chen, Z. (2019). Learning Continually from Low-shot
402 Data Stream. *arXiv preprint arXiv:1908.10223*.

403 LeCun, Y. (1998). The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>.

404 Lee, S. (2017). Toward continual learning for conversational agents. *arXiv preprint arXiv:1712.09943*.

405 Lesort, T., Caselles-Dupré, H., Garcia-Ortiz, M., Stoian, A., and Filliat, D. (2019a). Generative
406 models from the perspective of continual learning. In *International Joint Conference on Neural
407 Networks*.

408 Lesort, T., Díaz-Rodríguez, N., Goudou, J.-F., and Filliat, D. (2018). State representation learning for
409 control: An overview. *Neural Networks*.

410 Lesort, T., Gepperth, A., Stoian, A., and Filliat, D. (2019b). Marginal replay vs conditional replay for
411 continual learning. In *International Conference on Artificial Neural Networks*. Springer.

412 Lesort, T., Lomonaco, V., Stoian, A., Maltoni, D., Filliat, D., and Díaz-Rodríguez, N. (2019c).
413 Continual learning for robotics: Definition, framework, learning strategies, opportunities and
414 challenges. *arXiv preprint arXiv:1907.00182*.

415 Li, Z., Zhou, F., Chen, F., and Li, H. (2017). Meta-SGD: Learning to Learn Quickly for Few Shot
416 Learning. *arXiv preprint arXiv:1707.09835*.

417 Lomonaco, V. and Maltoni, D. (2017). Core50: A New Dataset and Benchmark for Continuous
418 Object Recognition. *arXiv preprint arXiv:1705.03550*.

419 Lopez-Paz, D. and Ranzato, M. (2017). Gradient Episodic Memory for Continual Learning. In
420 *Advances in Neural Information Processing Systems*.

421 Mallya, A., Davis, D., and Lazebnik, S. (2018). Piggyback: Adapting a single network to multiple
422 tasks by learning to mask weights. In *Proceedings of the European Conference on Computer
423 Vision*.

424 Mallya, A. and Lazebnik, S. (2018). Packnet: Adding multiple tasks to a single network by iterative
425 pruning. In *Computer Vision and Pattern Recognition*.

426 Maltoni, D. and Lomonaco, V. (2019). Continuous learning in single-incremental-task scenarios.
427 *Neural Networks*.

428 McCloskey, M. and Cohen, N. J. (1989). Catastrophic interference in connectionist networks: The
429 sequential learning problem. Elsevier.

430 Mitchell, T., Cohen, W., Hruschka, E., Talukdar, P., Yang, B., Betteridge, J., Carlson, A., Dalvi, B.,
431 Gardner, M., Kisiel, B., et al. (2018). Never-ending learning. *Communications of the ACM*.

432 Parisi, G. I., Kemker, R., Part, J. L., Kanan, C., and Wermter, S. (2019). Continual lifelong learning
433 with neural networks: A review. *Neural Networks*.

434 Patacchiola, M., Turner, J., Crowley, E. J., and Storkey, A. (2019). Deep kernel transfer in gaussian
435 processes for few-shot learning. *arXiv preprint arXiv:1910.05199*.

436 Ravi, S. and Larochelle, H. (2016). Optimization as a model for Few-Shot Learning. In *International
437 Conference On Learning Representations*.

438 Rebuffi, S.-A., Kolesnikov, A., Sperl, G., and Lampert, C. H. (2017). ICARL: Incremental Classifier
439 and Representation Learning. In *Computer Vision and Pattern Recognition*.

440 Ren, M., Liao, R., Fetaya, E., and Zemel, R. S. (2018a). Incremental Few-Shot Learning with
441 Attention Attractor Networks. *arXiv preprint arXiv:1810.07218*.

442 Ren, M., Triantafillou, E., Ravi, S., Snell, J., Swersky, K., Tenenbaum, J. B., Larochelle, H., and
443 Zemel, R. S. (2018b). Meta-Learning for Semi-Supervised Few-Shot Classification. *arXiv preprint
444 arXiv:1803.00676*.

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

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

450 Rusu, A. A., Rao, D., Sygnowski, J., Vinyals, O., Pascanu, R., Osindero, S., and Hadsell, R. (2019).
451 Meta-Learning with Latent Embedding Optimization. *International Conference On Learning
452 Representations*.

453 Santoro, A., Raposo, D., Barrett, D. G., Malinowski, M., Pascanu, R., Battaglia, P., and Lillicrap,
454 T. (2017). A simple neural network module for Relational Reasoning. In *Neural Information
455 Processing Systems*.

456 Santurkar, S., Tsipras, D., Ilyas, A., and Madry, A. (2018). How Does Batch Normalization Help
457 Optimization? *Neural Information Processing Systems*.

458 Snell, J., Swersky, K., and Zemel, R. (2017). Prototypical Networks for Few-Shot Learning. In
459 *Neural Information Processing Systems*.

460 Spigler, G. (2019). Meta-learned priors slow down catastrophic forgetting in neural networks. *arXiv
461 preprint arXiv:1909.04170*.

462 van de Ven, G. M. and Tolias, A. S. (2018). Generative replay with feedback connections as a general
463 strategy for continual learning. *arXiv preprint arXiv:1809.10635*.

464 Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al. (2016). Matching Networks for One Shot
465 Learning. In *Neural Information Processing Systems*.

466 Vuorio, R., Cho, D.-Y., Kim, D., and Kim, J. (2018). Meta Continual Learning. *arXiv preprint
467 arXiv:1806.06928*.

468 Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., and Perona, P. (2010).
469 Caltech-UCSD birds 200.

- 470 Xiao, H., Rasul, K., and Vollgraf, R. (2017). Fashion-MNIST: a Novel Image Dataset for Bench-
471 marking Machine Learning Algorithms.
- 472 Zenke, F., Poole, B., and Ganguli, S. (2017). Continual Learning through Synaptic Intelligence. In
473 *International Conference on Machine Learning*.
- 474 Zeno, C., Golan, I., Hoffer, E., and Soudry, D. (2018). Task Agnostic Continual Learning Using
475 Online Variational Bayes. *arXiv preprint arXiv:1803.10123*.