# Few-shot Continual Learning: a Brain-inspired Approach

Liyuan Wang[1 2 3], Qian Li[1 2], Yi Zhong[1 2 *] and Jun Zhu[3 *]
[1] School of Life Sciences, IDG/McGovern Institute for Brain Research,
Tsinghua University, Beijing, China. [2] Tsinghua-Peking Center for Life Sciences,
Beijing, China. [3] Dept. of Comp. Sci. & Tech., Institute for AI, BNRist Center,
THBI Lab, Tsinghua University, Beijing, China. [*] *Corresponding authors*

## Abstract

*It is an important yet challenging setting to continually learn new tasks from a few examples. Although numerous efforts have been devoted to either continual learning or few-shot learning, little work has considered this new setting of few-shot continual learning (FSCL), which needs to minimize the catastrophic forgetting to the old tasks and gradually improve the ability of few-shot generalization. In this paper, we provide a first systematic study on FSCL and present an effective solution with deep neural networks. Our solution is based on the observation that continual learning of a task sequence inevitably interferes few-shot generalization, which makes it highly nontrivial to extend few-shot learning strategies to continual learning scenarios. We draw inspirations from the robust brain system and develop a method that (1) interdependently updates a pair of fast / slow weights for continual learning and few-shot learning to disentangle their divergent objectives, inspired by the biological model of meta-plasticity and fast / slow synapse; and (2) applies a brain-inspired two-step consolidation strategy to learn a task sequence without forgetting in the fast weights while improve generalization without overfitting in the slow weights. Extensive results on various benchmarks show that our method achieves a better performance than joint training of all the tasks ever seen. The ability of few-shot generalization is also substantially improved from incoming tasks and examples.*

## 1. Introduction

Few-shot learning (FSL) and continual learning (CL) are two important yet challenging tasks, which usually co-exist in scientific and engineering domains. For instance, an AI-assisted diagnosis system needs to continually learn from a few cases to provide diagnostic suggestions and incrementally improve its ability of diagnosis. However, FSL and CL are basically studied as two separate directions. FSL aims to learn a novel task from a few examples. A general strat-

egy of FSL is to learn transferable knowledge from large amounts of base tasks, and then generalize to a few examples of a novel task. Following this idea, *gradient-based FSL* learns to initialize parameters for fast adaptation to a few examples [9]; *metric-based FSL* learns a generalized embedding space to match the representations of target examples to support examples by distance metrics [31, 8]. On the other hand, CL attempts to learn a sequence of novel tasks by avoiding *catastrophic forgetting* [21], that a neural network catastrophically forgets the learned parameters for old tasks when accommodating for new tasks. To mitigate catastrophic forgetting, effective strategies include regularization of important parameters for the old tasks, and replay of the old data distributions [23].

Table 1. Comparison of settings that consider both continual learning and few-shot learning.

| | Objective of CL | Objective of FSL |
|---|---|---|
| OML [12] | ✓ | ✗ |
| CBCL [4] | ✓ | ✗ |
| TOPIC [32] | ✓ | ✗ |
| AAN [28] | ✗ | ✓ |
| IDA [18] | ✗ | ✓ |
| Ours | ✓ | ✓ |

A few very recent efforts start to consider the challenges of FSL and CL together. Although all of them try to continually learn from a few examples, they focus on either continually tackling a sequence of new tasks without forgetting (i.e. the objective of CL), or improving the ability to learn a novel task with a few examples (i.e. the objective of FSL), rather than achieving the two objectives together. Table 1 summarizes the aims of current settings. Specifically, OML [12] learns representations that are robust to catastrophic forgetting, but only works well in relatively simple datasets with online training. CBCL [4] stores the clustered feature embedding of incremental classes for joint prediction. TOPIC [32] applies a neural gas network to learn first several tasks from large amounts of data, followed by several new tasks using a few examples. AAN [27] uses an

attention attractor network to learn a few new classes without interfering the performance on the large amounts of base classes. IDA [18] improves FSL from incremental examples by indirect discriminant alignment.

Ideally, FSL and CL should benefit each other. The ability of FSL enables a model to learn incoming tasks from a few examples, and CL of new tasks and examples provides additional data sources to improve generalization. From analyzing the limitations of current efforts, we propose a more challenging but realistic setting as few-shot continual learning (FSCL): A model continually learns from a few examples and achieves the objectives of CL and FSL together, that: (1) tackle a sequence of novel tasks while avoiding catastrophic forgetting; and (2) improve FSL from incoming tasks and examples while avoiding overfitting. However, the proposed objectives of FSL and CL are not easy to achieve simultaneously, because CL tries to precisely memorize incoming tasks to mitigate forgetting, which inevitably "overfits" the task sequence and interferes FSL. If the model of FSCL is not properly designed, e.g., to simply combine the strategies of FSL and CL, the learned ability of FSL is not improved from CL but substantially decreased (Fig. 1, b), and thus interferes the performance on the entire task sequence (Fig. 1, a).
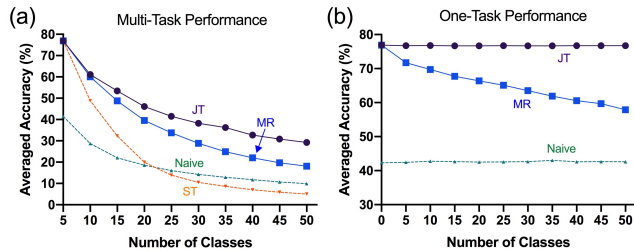


Figure 1. The ability of FSL is interfered by CL on miniImageNet. (a) is the performance to infer all the incremental classes ever seen. (b) is the performance on one 5-way 5-shot task during continual learning in (a). JT: Joint Training; ST: Sequential Training; MR: Memory replay of the training data ever seen; Naive: A randomly-initialized network.

Compared with artificial neural networks, the biological brain is able to continually learn from everyday occurrences and incrementally improve generalization [20, 22, 29]. Catastrophic interference can be effectively mitigated by replaying old experience together with new inputs [21, 20]. To fast generalize to an experience without overfitting, biological synapses interdependently update a fast and a slow components of plasticity, that a fast state quickly learns various inputs and then updates a slow state with a small step-size [29, 10]. Such synaptic plasticity across brain areas is modulated by replay activity [10]. Particularly, a replay strategy named *two-step consolidation* [22] addresses the challenges of CL and FSL simultaneously: The first step is a strong input of the current occurrences to tag the synap-

tic activities (*tagging*). While, the second step is to update the network with multiple weak replay inputs, where the cumulative activity of a synapse is applied to gate its stability (*capture*) and thus to selectively generalize memories.

Inspired by the model of fast / slow synapse, we interdependently update two sets of parameters, i.e. fast weights and slow weights, to disentangle the divergent objectives of CL and FSL. Inspired by the model of two-step consolidation, we establish a replay system to learn the fast / slow weights: The first step is to approximate the magnitude of expected gradients for incoming tasks in the slow weights, as the "activity tag". The second step is to update the fast weights for multiple iterations, with regularization on the gated cumulative activity to balance fitting and generalization. Then, the slow weights are updated for a small step by the fast weights to improve generalization. We sample task sequences and training examples from various benchmark datasets of FSL, including Omniglot, miniImageNet, tieredImageNet and CUB-200-2011, as the evaluation benchmarks for FSCL. Extensive evaluations on such benchmarks support the effectiveness of our method to achieve the proposed aims of FSCL.

Our contributions include: 1. We propose a challenging but realistic setting as FSCL, in which the objectives of FSL and CL should be achieved simultaneously to improve each other, and provide a systematic analysis to disclose the relation between FSL and CL in FSCL; 2. We draw inspirations from the biological models of synaptic plasticity and memory replay to propose a novel approach for FSCL, named two-step consolidation (TSC), to interdependently update parameters for FSL and CL, respectively; 3. Extensive evaluations on benchmarks of FSCL validate the effectiveness of our strategy to mitigate catastrophic forgetting and incrementally improve FSL from CL.

## 2. Problem Formulation

We start by introducing the notations and setup of few-shot continual learning (FSCL). Then we analyze the non-trivial interference between CL and FSL in FSCL.

### 2.1. Setup and Problem Formulation

Few-shot continual learning (FSCL) is a setting where a learner continually learns $T$ tasks and each task has a few training examples. For a typical classification task, let $N$ be the number of classes and $K$ be the number of training examples for each class.[1] We call such a task $N$-way $K$-shot classification. Let $D_t = \{(x_i^t, y_i^t)\}_{i=1}^{N \times K}$ be the fully-labeled training set of task $t$. Then, the collection $D_T = \bigcup_{t=1}^{T} D_t$ is called a query dataset for learning on the task sequence. In FSCL, $D_T$ is sequentially provided to the learner, that $D_t$ is introduced when training on task $t$.

---

[1]It is straightforward to have different $K$ for various classes.

An ideal model for FSCL should achieve the objectives of both few-short learning (FSL) and continual learning (CL) simultaneously — it can quickly learn a new task from a few examples, mitigate catastrophic forgetting of the old tasks, and incrementally improve FSL from incoming tasks and examples. Correspondingly, the evaluation should reflect the two objectives — it should perform single-head evaluation [7], i.e. all the learned tasks are evaluated without access to the task label, to test the ability of CL; and also evaluate on a new task with a few training samples to test the ability of FSL. Table 1 summarizes the difference of FSCL from the existing studies, which simply improve one of them.

## 2.2. Interference between Continual Learning and Few-Shot Generalization in FSCL

We now analyze the nontrivial interference between continual learning and few-shot generalization if the FSCL model is not properly designed.

To solve FSCL, a straightforward idea would be to directly combine the strategies of FSL and CL to build a model. Specifically, we extend the state-of-the-art FSL models to a continual learning setting. Here, we choose gradient-based FSL (e.g., MAML [9]) as the base model for FSCL. [2] The entire training process includes two stages: The model first learns transferable knowledge from large amounts of base tasks in a support set $S$, called a meta-training (*MT*) stage. $S$ should be a much larger dataset with transferable domains to $D_T$ while without overlapped classes and examples ($S \bigcap D_T = \emptyset$): e.g., if $D_T$ is a dataset of birds, $S$ can be a larger dataset of other animals. Then, the model continually learns a sequence of novel tasks from $D_T$ without access to $S$, called a *FSCL* stage. CL strategies, like memory replay and weight regularization, are implemented to mitigate catastrophic forgetting. After training of each task, the performances of both CL and FSL are evaluated.

However, to precisely memorize the incoming tasks inevitably "overfits" the tasks ever seen in the FSCL stage, which interferes the ability of FSL on following tasks. Fig. 1 shows the results of FSCL on miniImageNet. Joint training (JT) of all the tasks ever seen is the upper-bound performance for the objective of CL. Compared with sequential training (ST) of each task, memory replay (MR) of the training data ever seen can completely avoid catastrophic forgetting, but significantly underperforms JT (Fig. 1, a). The ability of FSL also gradually decreases from learning the task sequence (Fig. 1, b). Either, a simple combination of regularization-based CL methods cannot solve this issue (Table 2).

To explicitly show the interference of CL to FSL, we

first visualize the feature embedding in CL of several few-shot classification tasks on miniImageNet with MR to avoid catastrophic forgetting. As shown in Fig. 2 (a, b), the learned feature embedding is only slightly adjusted to learn a new 5-way 5-shot (5W5S) task. However, after learning several additional tasks and then training on the 5W5S task, the learned feature embedding is incrementally deformed, which interferes the performance of FSL. Next, when the model moves from the MT stage to the FSCL stage, we observe a significant collapse of Fisher Information of the feature extractor parameters (Fig. 2, c), which indicates forgetting of the important parameters for FSL [15]. Further, the confusion matrix of the above experiments indicates that CL of new task interferes generalization to the entire task sequence (Fig. 2, d), which cannot be addressed by a simple combination of MR and weight regularization (Fig. 5).

We provides a slightly more formal illustration of FSCL in Fig. 3 (a, b), where $\theta_0^*$ and $\theta_0$ are the parameters of a network with or without training on the support set $S$, respectively. After training on large amounts of base tasks sampled from $S$, $\theta_0^*$ moves to the center of the solution manifolds of these tasks, and thus can quickly generalize to a novel task from a few examples. We assume a CL approach can completely avoid catastrophic forgetting. However, after continual learning of several tasks, $\theta_0^*$ finally converges to a solution manifold close to $\theta_0$ with poor performance, i.e. bad solution. By contrast, if the distributions of $S$ and $D_T$ are the same and all $T$ tasks are *i.i.d.*, good solution of FSCL should be close to the center. If the distributions of $S$ and $D_T$ are (slightly) different, the optimal solution will (slightly) deviate the center. Thus, the network should be carefully adjusted on the sequentially arrived $D_t$ to improve generalization while avoiding overfitting.

## 3. Our Proposal

We now present our solution to address the above issues that draws inspirations from biological brain — an ideal system for FSCL compared with artificial neural networks.

### 3.1. Two-Step Consolidation

Biological brain can continually learn new tasks from a few examples and apply the learned information for generalization [22, 20, 17]. Replay of old memory with new inputs is able to mitigate catastrophic forgetting [20]. Particularly, a replay strategy named *two-step consolidation* (TSC) has been shown effective to integrate a few examples to a network and selectively generalize memories [22] (Fig. 4, a). The first step is a strong input of the current occurrence to tag the synaptic activity. A *synaptic tagging and capture* (STC) process selectively gates the stability of synapses dependent on the cumulative tags. Then, multiple weak replay inputs in the second step can selectively generalize the learned memories for inference [22]. Correspondingly, the
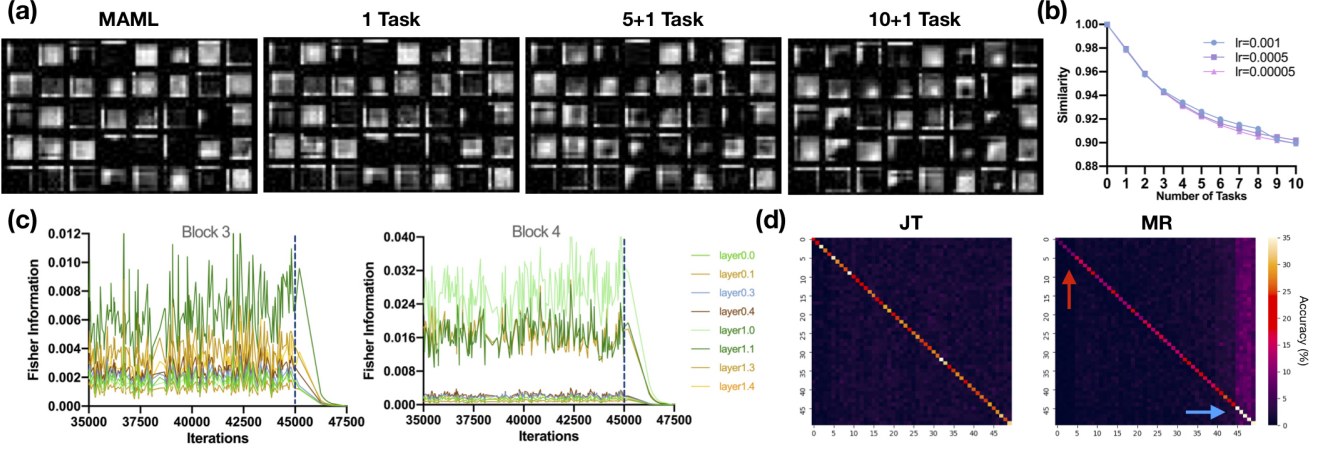
---

Figure 2. The objectives of FSL and CL interfere with each other in FSCL on miniImageNet. (a, b) The learned feature embedding for FSL is incrementally deformed in CL on various learning rates (lr). (c) Fisher Information of feature extractor parameters on ResNet-18 in meta-training (before 45k iterations) and training on a specific task (after 45k iterations). (Best viewed in color) (d) Confusion matrix of CL of ten 5W5S tasks. The blue and the red arrows indicate overfitting to the current task and forgetting of the earlier tasks, respectively. x-axis: predicted labels; y-axis: ground-truth labels.
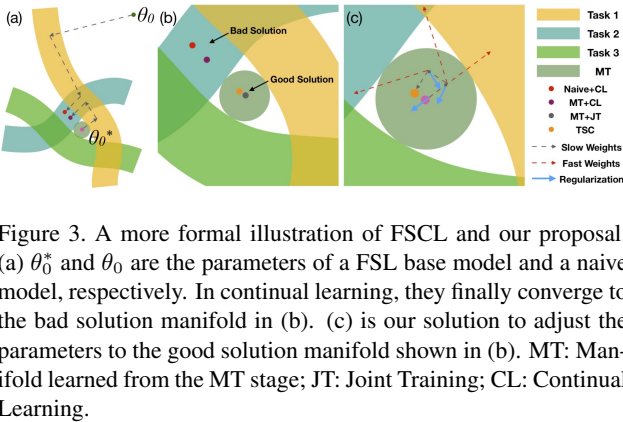


Figure 3. A more formal illustration of FSCL and our proposal. (a) $\theta_0^*$ and $\theta_0$ are the parameters of a FSL base model and a naive model, respectively. In continual learning, they finally converge to the bad solution manifold in (b). (c) is our solution to adjust the parameters to the good solution manifold shown in (b). MT: Manifold learned from the MT stage; JT: Joint Training; CL: Continual Learning.

biological synaptic strength is modulated by replay activity through two components of plasticity, which depend on the history of synaptic modifications (called *meta-plasticity*): The state of synaptic strengths is copied from the previous state. A fast component quickly encodes various replay inputs, and then updates a slow component by a small step as the current state, which will be copied to the downstream states [10].

## 3.2. Our Solution to FSCL

Inspired by such brain strategies to continually generalize memories to incoming occurrences and learn from them, we propose a novel solution to FSCL named TSC, which enables a network to simultaneously achieve the objectives of FSL and CL (Fig. 4, b and Algorithm 1). To disentangle the divergent objectives of FSL and CL, we use two sets of parameters for a network, i.e. *fast weights* $\theta_f$ and *slow weights* $\theta_s$, and interdependently update them for CL
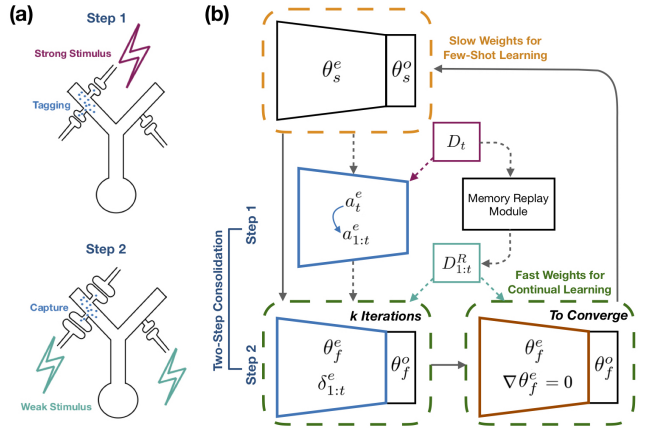


Figure 4. The biological and computational models of two-step consolidation (TSC). (a) Synapses are tagged by a strong input in the first step and then captured by multiple weak inputs in the second step. (b) Architecture of TSC to interdependently update the fast / slow weights for CL and FSL, respectively.

and FSL, respectively. We consider a typical network for FSCL that includes a feature embedding layer ($e$) and an output layer ($o$). Correspondingly, both $\theta_f$ and $\theta_s$ consist of the parameters of $e$ and $o$ (if applicable) that $\theta_f = \theta_f^e \cup \theta_f^o$, $\theta_s = \theta_s^e \cup \theta_s^o$. For gradient-based few-shot classification, the output layer is a linear classifier (or fully connected layer) in general. Our method draws key inspirations from how biological synapses are modulated by replay inputs to precisely memorize the old experience and better generalize to a new experience. As shown in Fig. 4 (b), we first implement a memory replay module, and then apply a brain-inspired two-step consolidation (TSC) strategy to interdependently

4

update $\theta_f$ and $\theta_s$ on a replay dataset, as detailed below.

---

**Algorithm 1** Algorithm of TSC for FSCL
___
1: **Require:** $\theta_s$: parameters learned from the MT stage; $D_t$: training dataset of task $t$; $\beta$, $\lambda$, $m$: hyperparameters; $k$: iterations to learn $\theta_f^e$.
2: **for** task $t = 1, 2, 3, ...$ **do**
3:    Establish the replay dataset $D_{1:t}^R = D_t \bigcup D_{\leq t-1}^R$
4:    Step 1: Compute $a_t^e$ and update $a_{1:t}^e$, $\delta_{1:t}^e$
5:    $\theta_f \leftarrow \theta_s$
6:    **for** $1, 2, ..., k$ **do**
7:       Step 2: Update $L_e\left(\theta_f^e, \theta_f^o, D_{1:t}^R\right)$
8:    **end for**
9:    **while** not done **do**
10:      Update $L_o\left(\theta_f^o, D_{1:t}^R\right)$
11:    **end while**
12:    $\theta_s \leftarrow \theta_s - \beta\left(\theta_s - \theta_f\right)$
13: **end for**
___

To mitigate catastrophic forgetting of the old experience, we apply a memory replay module to provides a replay dataset of the old tasks as $D_{\leq t-1}^R$ when learning a new task $t$. This module can be a memory buffer to replay old training data, or a generative model to learn the old data distributions and replay generated data.

To effectively learn the parameters, we need to consider the trade-off between fitting and generalization. We first propose a strategy to interdependently update $\theta_f$ and $\theta_s$, inspired by the biological fast / slow synapse. In the MT stage, $\theta_s$ are first trained on a support set $S$. In the FSCL stage, $\theta_s$ transfer parameters to $\theta_f$ before training on each task. $\theta_f$ learn incoming tasks and examples for inference on the task sequence. $\theta_s$ are updated by $\theta_f$ for a small step $\beta$ to improve the ability of generalization from the new tasks and examples. $\theta_s$ then transfer the updated parameters to $\theta_f$ to learn the next task. In such a framework, $\theta_f$ tend to precisely learn and memorize the task sequence for better inference, but should only feedback generalized information to $\theta_s$.

To achieve this aim, we learn $\theta_f^e$ on $D_{1:t}^R = D_t \bigcup D_{\leq t-1}^R$ for limited $k$ iterations and regularize the training: large changes of the parameters should be penalized to avoid overfitting, while a too strong limitation generally results in underfitting. We draw inspirations from the biological synaptic tagging and capture (STC) in two-step consolidation: Synaptic activities in a network are first tagged by a new input of the current experience. Then the network is updated by multiple replay inputs, where the synaptic stability is gated by the cumulative activity tags (Fig. 4, a). Specifically, we propose to learn $\theta_f^e$ with a STC regularization strategy, which penalizes changes of $\theta_f^e$ from $\theta_s^e$ dependent on the gated cumulative activity of learning history:

$$L_e(\theta_f^e, \theta_f^o, D_{1:t}^R) = L_{(x,y)\sim D_{1:t}^R}(\theta_f^e, \theta_f^o)$$
$$+\lambda \sum_l \sum_{i=1}^{N_l} \delta_{1:t,i}^{e,l}(\theta_{f,i}^{e,l} - \theta_{s,i}^{e,l})^2, \quad (1)$$

$$\delta_{1:t,i}^{e,l} = \sigma[m(a_{1:t,i}^{e,l} - \frac{1}{t}\frac{1}{N_l}\sum_{i=1}^{N_l} a_{1:t,i}^{e,l})], \quad (2)$$

where $N_l$ is the number of parameter $i$ in layer $l$, $a_{1:t,i}^{e,l}$ is the cumulative activity in task $1:t$ of the parameter, and $m$ is a positive scale factor. For task $t$, we calculate the magnitude of the expected gradients as the "activity" of a parameter in $\theta_s^e$, that $a_{t,i}^{e,l} = \| \mathbb{E}_{(x,y)\sim D_t}[\nabla L(\theta_{s,i}^{e,l})] \|$. The cumulative activity is updated as $a_{1:t,i}^{e,l} = a_{1:t-1,i}^{e,l} + a_{t,i}^{e,l}$. For each parameter $i$ in layer $l$, we average $a_{1:t,i}^{e,l}$ in that layer of $t$ tasks as the "threshold", and use a sigmoid function $\sigma$ as the "gate" to project the threshold-centered cumulative activity to the values in $(0, 1)$. Therefore, the changes of $\theta_f^e$ are limited by their cumulative magnitude of the expected gradients to alleviate overfitting, while the sigmoid function clips the penalty to avoid underfitting. $L_{(x,y)\sim D_{1:t}^R}$ is the loss function for the target tasks trained on the replay dataset $D_{1:t}^R$. We use a cross entropy loss for classification.

Then we update $\theta_f^o$ with fixed $\theta_f^e$ for better inference on the task sequence, which can be multiple iterations till convergence:

$$L_o(\theta_f^o, D_{1:t}^R) = L_{(x,y)\sim D_{1:t}^R}(\theta_f^o). \quad (3)$$

The overall objective of our method is defined as:

$$L_e(\theta_f^e, \theta_f^o, D_{1:t}^R) + L_o(\theta_f^o, D_{1:t}^R), \ t = 1, 2, ..., T. \quad (4)$$

As shown in Fig. 3 (c), the brain-inspired designs in TSC provide an effective solution to FSCL: The fast / slow weight strategy enables a network to fit a task sequence by fast weights, but only update slow weights by a small step to avoid overfitting. Then, the regularization strategy further regularizes the changes of parameters to the direction of good solution manifold.

## 4. Experiment

### 4.1. Experiment Setup

Most methods of gradient-based FSL, e.g., MAML [9], MAML++[2], with linear classifier [8] can be easily adapted to TSC. We also consider metric-based FSL as the base model by implementing ProtoNet [31] in Appendix A. We evaluate TSC for FSCL of *New Class* and *New Instance* [19, 3], first on three benchmark datasets of FSL, i.e. Omniglot, miniImageNet and tieredImageNet. Then we consider huge domain differences, i.e. $S$ is the training set of

tieredImageNet while $D_T$ is sampled from CUB-200-2011 dataset (referred to as the CUB). We apply the averaged accuracy ($A_t$) on the test set of $t$ tasks trained so far to evaluate the ability of CL, and evaluate the ability of FSL on an additional novel task. All the experiment results are averaged by 20 randomly sampled tasks or task sequences, if not specified.

**Dataset:** Omniglot [16] is a dataset of handwritten characters, consisting of 963 classes for training and 660 classes for testing with 20 images per class . miniImageNet [33] is a dataset of color images derived from the 1000-class ILSVRC-12 dataset [30], including 100 classes with 600 images per class. We follow the same split of miniImageNet from ILSVRC-12 as [33], but apply all 100 classes for FSCL and the remained 900 classes as the support set. tieredImageNet [28] is also derived from ILSVRC-12, including 608 classes from 34 super-classes. Each super-class consists of 20, 6, and 8 classes for training, validation and testing, respectively. CUB-200-2011 [34] is a dataset including 200 classes and 11,788 colored images of birds. Omniglot and other three datasets are resized to $32 \times 32$ and $64 \times 64$ respectively before experiments [3]. Since the hierarchy in tieredImageNet splits the testing classes sufficiently distinct from the training classes [28], domain differences between $S$ and $D_T$ are generally larger and larger from miniImageNet, tieredImageNet to tieredImageNet→CUB.

**Baseline:** First, we compare TSC with several baselines of representative CL methods, using MAML as the base model of FSL. Since replay-based methods show great advantages in single-head evaluation on incremental classes [13], we consider the baselines as below: (1) Joint training (JT) of all the tasks ever seen, i.e. the upper bound performance of CL. (2) Memory replay (MR): We store the old training data ever seen in a memory buffer to theoretically avoid forgetting of the learned tasks. Although many existing work [26, 6] use a fixed memory buffer (with a pre-reserved size of 20 images per class in general) selected from much larger amounts of training data, in FSCL the training data (typically 1- or 5-shot per class) is insufficient to perform selection. Also, the memory buffer of all the old data is much smaller than a fixed size of 20 images per class and is indeed competitive to commonly-used methods of selection as analyzed in [7]. (3) MR with regularization-based methods: We combine MR with representative regularization-based methods (appended with '-M'), e.g., SI [35], EWC [15] and MAS [1], which penalize changes of the important parameters for the old tasks. (4) Generative replay (GR): To evaluate the strategy of generative replay and augment the limited training data, we train a generative model to recover the old data distributions.

Next, we compare with several recent meta-learning approaches that improve the ability of CL through a better

meta-training, including OML javed2019meta, ANML [5] and APL [25]. OML [12] is a meta-learning approach to learn sparse representations that are more robust to catastrophic forgetting. ANML [5] meta-learns an activation-gating function for context-dependent selective activation to improve OML. APL [25] learns a memory module from the support set for FSL, which also benefits CL since the memory module can be used for memory replay. In contrast to the methods above, TSC aims to extend a base model that achieves considerable ability of FSL to FSCL. We implement MAML++ [2], a more recent meta-learning approach for FSL, as the base model of TSC to show this advantage. For all experiments, we extensively evaluate hyperparameters and report the best performance.

**Architecture and Training:** The base model of MAML applies a ResNet-18 feature extractor with similar implementations of MAML as [8], validated in Appendix B. For generative replay, we apply a 4-layer DCGAN architecture [24] with WGAN gradient penalty to stabilize generation [11], but replace BatchNorm layers by InstanceNorm to learn a few examples. The implementations and hyperparameters are detailed in Appendix C.

## 4.2. FSCL of New Class

In Table 2, we evaluate TSC on a sequence of 5-way 5-shot (5W5S) classification tasks. Because of the interference between CL and FSL, the performances of MR and GR are much lower than JT, although GR slightly outperforms MR due to the augmentation effects of generated data. The regularization-based methods penalize changes of the parameters and thus outperform MR, but still significantly underperform JT. In contrast, through interdependently updating the fast / slow weights (FSW) with the STC regularization (Reg), TSC substantially outperforms joint training (JT) and other baselines. The performance of TSC can be further improved by using GR as the memory replay module (Table 2).

Then, we compare with the meta-learning approaches that improve the ability of CL through a better meta-training. As shown in Table 3, although OML javed2019meta and ANML [5] learn a sparse representation in meta-training to alleviate catastrophic forgetting in CL, the learned representation sacrifices the ability of FSL and thus decreases the performance on the entire task sequence. APL [25] requires a larger number of training data, e.g., 20-shot per class, to select prototypes in the memory modules, and significantly underperforms other methods. By contrast, TSC extends a base model that acquires a promising ability of FSL to FSCL. Therefore, TSC with MAML [9] as the base model substantially outperforms other methods. The performance of FSCL can be further improved by using more recent meta-learning approaches for FSL as the base model, e.g., MAML++ [2].

Table 2. Averaged accuracy (%) on a sequence of 5W5S classification tasks. All the baselines and TSC are initialized by MAML [9]. The results are shown as top-1 accuracy or top-1 / top-5 accuracy.

| | Method | Omniglot $A_{30}$ | Omniglot $A_{50}$ | miniImageNet $A_{30}$ | miniImageNet $A_{50}$ | tieredImageNet $A_{30}$ | tieredImageNet $A_{50}$ | tieredImageNet→CUB $A_{30}$ | tieredImageNet→CUB $A_{50}$ |
|---|---|---|---|---|---|---|---|---|---|
| Baseline | JT | 96.20 | 94.97 | 36.83 / 70.46 | 28.75 / 59.60 | 31.51 / 65.34 | 23.04 / 53.76 | 29.83 / 64.58 | 23.43 / 53.86 |
| | MR | 94.38 | 92.13 | 28.85 / 54.39 | 18.07 / 38.42 | 24.21 / 50.57 | 15.10 / 34.80 | 20.86 / 45.79 | 13.17 / 31.87 |
| | GR | 94.93 | 93.49 | 28.49 / 55.02 | 19.81 / 41.45 | 24.72 / 50.40 | 16.30 / 37.42 | 21.26 / 48.14 | 15.68 / 36.90 |
| | SI-M [31] | 94.97 | 94.51 | 28.08 / 54.13 | 19.09 / 39.67 | 24.66 / 51.60 | 16.57 / 36.39 | 20.11 / 45.92 | 14.64 / 33.71 |
| | EWC-M [15] | 95.52 | 93.95 | 27.46 / 54.88 | 25.29 / 53.19 | 25.76 / 53.44 | 21.99 / 50.01 | 23.97 / 57.47 | 15.87 / 41.97 |
| | MAS-M [1] | 95.81 | 94.47 | 28.94 / 60.48 | 27.01 / 58.98 | 26.48 / 58.25 | 22.11 / 53.84 | 22.33 / 54.41 | 18.60 / 47.84 |
| TSC (Ours) | Reg | 95.84 | 95.06 | 35.02 / 72.02 | 29.56 / 60.44 | 30.35 / 67.95 | 22.93 / 55.72 | 27.24 / 61.48 | 20.48 / 48.82 |
| | FSW | 96.44 | 95.57 | 37.13 / 70.51 | 30.32 / 61.51 | 30.66 / 64.95 | 23.29 / 54.27 | 29.75 / 63.89 | 23.76 / 54.29 |
| | Reg+FSW | **96.67** | **96.04** | **38.49 / 71.87** | **30.70 / 62.38** | **32.73 / 67.61** | **24.92 / 56.73** | **30.52 / 65.00** | **24.38 / 55.20** |
| | Reg+FSW+GR | 96.89 | 96.63 | 39.23 / 72.56 | 33.66 / 61.55 | 33.33 / 67.99 | 25.48 / 56.43 | 31.63 / 64.74 | 24.86 / 54.63 |

Table 3. Averaged accuracy (%) of FSCL on a sequence of 5-way tasks on Omniglot. *APL uses 20-shot per class while other methods use 5-shot per class.

| Method | $A_{30}$ | $A_{50}$ |
|---|---|---|
| OML [12] | 91.47 | 88.55 |
| ANML [5] | 96.71 | 95.47 |
| APL [25] | 94.88* | 93.42* |
| MAML[9]+JT | 96.20 | 94.97 |
| MAML[9]+MR | 94.38 | 92.13 |
| MAML[9]+TSC | 96.67 | 96.04 |
| MAML++[2]+JT | 97.33 | 95.67 |
| MAML++[2]+MR | 95.53 | 93.26 |
| MAML++[2]+TSC | 97.78 | 96.63 |

## 4.3. Ablation Study and Analysis

We provide the ablation study of TSC in Table 2 that both the brain-inspired designs of the STC regularization (Reg) and the fast / slow weight (FSW) significantly improve FSCL. Here we provide a more extensive analysis of how the two strategies balance fitting and generalization in FSCL. First, we examine two hyperparameters to control FSW in Table 4, including learning rate of the slow weight $\theta_s$ ($\beta$) and iterations to learn the fast weight $\theta_f$ ($k$). $\beta = 0$ is equal to JT, while $\beta = 1$ is a batch version of MR. Smaller $\beta$ mitigates the decreased performance of CL through more slowly updating $\theta_s$, while the performances on $\beta = 0.01$ and $\beta = 0.001$ are comparable. Then we evaluate $k$ on $\beta = 0.01$. Under the extreme condition of $k = 0$, i.e. the feature embedding layer is fixed, the model still achieves considerable performance and significantly outperforms MR (Table 2), consistent with our analysis that the good solution of FSCL is located in the MT manifold if the generalization error is trivial. The network can be fully updated by a larger $k$, but generally suffers from more overfitting. Compared with $k = 0$ and $k = 500$, the perfor-

mance is improved by updating the feature extractor with limited iterations, i.e. $k = 100$.

The STC regularization further addresses this issue by selectively stabilizing parameters, that $k = 100$+Reg significantly outperforms JT, $k = 100$ and $k = 500$ (Table 4). While adding a constant penalty (CP) cannot improve $k = 100$. On the other hand, the strategy of FSW with Reg further improves FSCL, compared with learning one-set of parameters in Table 2. Next, we analyze how the STC regularization is more effective in FSCL than other regularization-based methods for CL. As shown in 5, STC incrementally stabilizes parameters with a relatively small but constant penalty to balance fitting and generalization. While, other methods estimate a strong penalty on very few parameters, which cannot fit the entire task sequence well without interfering generalization (Fig. 5, a, b).

Table 4. Ablation study of TSC. We report averaged accuracy (%) (top-1 / top-5) of 5W5S class-incremental learning. JT: joint training; MR: memory replay; Reg: the STC regularizer; CP: a constant penalty.

| | Hyperparameters | miniImageNet $A_{30}$ | miniImageNet $A_{50}$ | tieredImageNet $A_{30}$ | tieredImageNet $A_{50}$ |
|---|---|---|---|---|---|
| $k = 500$ | $\beta = 0$ (JT) | 36.83 / 70.46 | 28.75 / 59.60 | 31.51 / 65.34 | 23.04 / 53.76 |
| | $\beta = 1$ (MR) | 28.54 / 54.85 | 18.17 / 38.58 | 24.65 / 50.73 | 15.28 / 35.55 |
| | $\beta = 0.1$ | 35.25 / 66.64 | 24.17 / 50.96 | 28.75 / 61.03 | 19.28 / 44.47 |
| | $\beta = 0.01$ | 36.17 / 69.57 | 28.51 / 59.13 | 30.98 / 65.07 | 22.78 / 52.72 |
| | $\beta = 0.001$ | 35.85 / 69.17 | 28.52 / 59.69 | 31.17 / 65.41 | 22.84 / 52.70 |
| $\beta = 0.01$ | $k = 0$ | 33.97 / 72.00 | 26.77 / 60.73 | 29.56 / 64.29 | 21.91 / 52.95 |
| | $k = 100$ | 37.13 / 70.51 | 30.32 / 61.51 | 30.66 / 64.95 | 23.29 / 54.27 |
| | $k = 100$ +CP | 37.06 / 70.93 | 30.11 / 60.89 | 30.93 / 65.56 | 23.30 / 54.11 |
| | $k = 100$ +Reg | **38.49 / 71.87** | **30.70 / 62.38** | **32.73 / 67.61** | **24.92 / 56.73** |

## 4.4. FSL is Incrementally Improved from CL

To evaluate the ability to transfer knowledge in FSCL, we evaluate TSC on miniImageNet, tieredImageNet and tieredImageNet→CUB during incremental learning of ten 5W5S tasks. We first consider the backward transfer (BWT), which indicate the averaged influence that learning a new task has on the performance on the previous tasks. Then we examine the ability of the model to learn an extra 5W5S classification task consisting of novel class
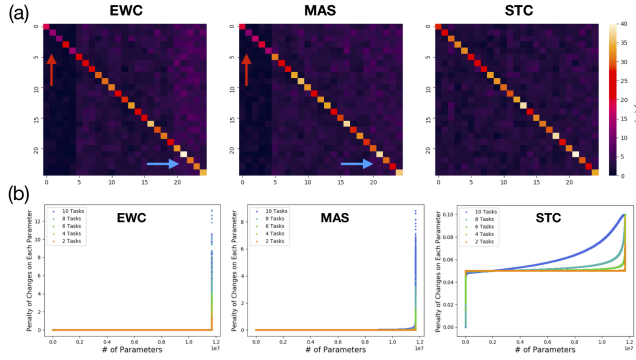
Figure 5. (a) Confusion matrix of various weight regularizations in FSCL of five 5W5S classification tasks on tieredImageNet. The blue and the red arrows indicate overfitting to the current task and forgetting of the earlier tasks, respectively. x-axis: predicted labels; y-axis: ground-truth labels. The distributions of penalty on the feature extractor parameters are shown in (b) (The scales of y-axis are different).
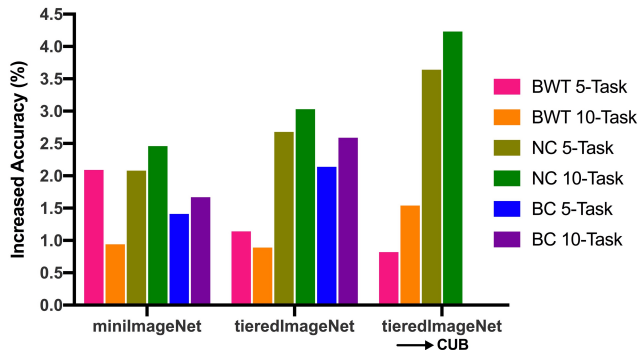


Figure 6. TSC improves the backward transfer (BWT, %) and the 5W5S performance on both novel class (NC, %) and base class (BC, %) in FSCL, averaged by more than 20 randomly sampled task sequences.

(NC, sampled from the query set) or base class (BC, sampled from the support set) in FSCL. As shown in Fig. 6, TSC substantially improves the ability of FSL of the learned tasks, the query set and the support set from incremental tasks and examples in FSCL. Particularly, TSC can flexibly adapt to domain difference, since larger domain differences generally result in a more significant improvement of NC and BC.

### 4.5. FSCL of New Instance

We now evaluate our framework on *New Instance* [19], i.e. FSCL on incremental batches of a few new instances. Since all instances belong to the same task, the model sequentially learns new batch of instances *without* replay of old batches. In Fig. 7, we randomly sample a task (20-class for Omniglot and 5-class for miniImageNet, tiered-
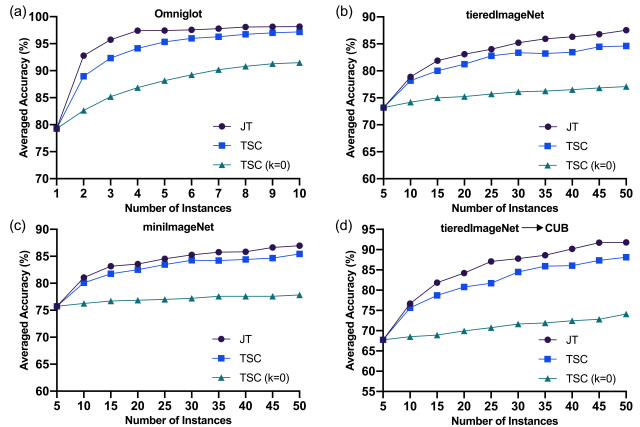


Figure 7. FSCL on new instance. FSL on one task is improved from sequentially seeing new instances by TSC.

ImageNet and CUB) and sequentially train the model with batches of new instances (1-instance / class for Omniglot, 5-instance / class for miniImageNet, tieredImageNet and CUB). After CL of several batches, we evaluate the performance of FSL on a new batch. The performance of TSC on one batch is continually improved from sequentially seeing new instances, only slightly lower than the joint training (JT) of all the instances ever seen. Note that TSC significantly outperforms TSC ($k = 0$), which applies a fixed feature embedding layer and only updates the output layer. The results above demonstrate the effectiveness of TSC to continually learn new instances, incrementally improve FSL from CL, and adapt to domain differences.

## 5. Conclusions

In this work, we present a first systematic study to a challenging but realistic setting of few-shot continual learning (FSCL), in which the objectives of few-shot learning (FSL) and continual learning (CL) should be achieved simultaneously to improve each other. We present a solution to FSCL by drawing inspirations from the biological models of synaptic plasticity and replay. Extensive results on various datasets demonstrate the effectiveness of our approach for FSCL. Further work will focus on extending our method to other FSL approaches and CL scenarios.

## References

[1] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 139–154, 2018. 6, 7

[2] Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *arXiv preprint arXiv:1810.09502*, 2019. 5, 6, 7

[3] Antreas Antoniou, Massimiliano Patacchiola, Mateusz Ochal, and Amos Storkey. Defining benchmarks for continual few-shot learning. *arXiv preprint arXiv:2004.11967*, 2020. 5, 6

[4] Ali Ayub and Alan Wagner. Brain-inspired model for incremental learning using a few examples. *arXiv preprint arXiv:2002.12411*, 2020. 1

[5] Shawn Beaulieu, Lapo Frati, Thomas Miconi, Joel Lehman, Kenneth O Stanley, Jeff Clune, and Nick Cheney. Learning to continually learn. *arXiv preprint arXiv:2002.09571*, 2020. 6, 7

[6] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 233–248, 2018. 6

[7] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 532–547, 2018. 3, 6

[8] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. *arXiv preprint arXiv:1904.04232*, 2019. 1, 5, 6, 11, 12

[9] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1126–1135. JMLR. org, 2017. 1, 3, 5, 6, 7, 11, 12

[10] Stefano Fusi. Computational models of long term plasticity and memory. *arXiv preprint arXiv:1706.04946*, 2017. 2, 4

[11] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in neural information processing systems*, pages 5767–5777, 2017. 6, 12

[12] Khurram Javed and Martha White. Meta-learning representations for continual learning. In *Advances in Neural Information Processing Systems*, pages 1818–1828, 2019. 1, 6, 7

[13] Ronald Kemker, Marc McClure, Angelina Abitino, Tyler L Hayes, and Christopher Kanan. Measuring catastrophic forgetting in neural networks. In *Thirty-second AAAI conference on artificial intelligence*, 2018. 6

[14] Jongmin Kim, Taesup Kim, Sungwoong Kim, and Chang D Yoo. Edge-labeling graph neural network for few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11–20, 2019. 11, 12

[15] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017. 3, 6, 7

[16] Brenden Lake, Ruslan Salakhutdinov, Jason Gross, and Joshua Tenenbaum. One shot learning of simple visual concepts. In *Proceedings of the annual meeting of the cognitive science society*, volume 33, 2011. 6

[17] Penelope A Lewis and Simon J Durrant. Overlapping memory replay during sleep builds cognitive schemata. *Trends in cognitive sciences*, 15(8):343–351, 2011. 3

[18] Qing Liu, Orchid Majumder, Alessandro Achille, Avinash Ravichandran, Rahul Bhotika, and Stefano Soatto. Incremental meta-learning via indirect discriminant alignment. *arXiv preprint arXiv:2002.04162*, 2020. 1, 2

[19] Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. *arXiv preprint arXiv:1705.03550*, 2017. 5, 8

[20] James L McClelland. Incorporating rapid neocortical learning of new schema-consistent information into complementary learning systems theory. *Journal of Experimental Psychology: General*, 142(4):1190, 2013. 2, 3

[21] James L McClelland, Bruce L McNaughton, and Randall C O'Reilly. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological review*, 102(3):419, 1995. 1, 2

[22] Cian O'Donnell and Terrence J Sejnowski. Selective memory generalization by spatial patterning of protein synthesis. *Neuron*, 82(2):398–412, 2014. 2, 3

[23] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 2019. 1

[24] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 6, 12

[25] Tiago Ramalho and Marta Garnelo. Adaptive posterior learning: few-shot learning with a surprise-based memory module. *arXiv preprint arXiv:1902.02527*, 2019. 6, 7

[26] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017. 6

[27] Mengye Ren, Renjie Liao, Ethan Fetaya, and Richard Zemel. Incremental few-shot learning with attention attractor networks. In *Advances in Neural Information Processing Systems*, pages 5276–5286, 2019. 1, 11

[28] Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B Tenenbaum, Hugo Larochelle, and Richard S Zemel. Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676*, 2018. 1, 6

[29] Alex Roxin and Stefano Fusi. Efficient partitioning of memory systems and its importance for memory consolidation. *PLoS computational biology*, 9(7), 2013. 2

[30] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015. 6

[31] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*, pages 4077–4087, 2017. 1, 5, 7, 11

[32] Xiaoyu Tao, Xiaopeng Hong, Xinyuan Chang, Songlin Dong, Xing Wei, and Yihong Gong. Few-shot class-incremental learning. *arXiv preprint arXiv:2004.10956*, 2020. 1, 11

[33] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016. 6, 11

[34] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 6

[35] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3987–3995. JMLR. org, 2017. 6

## A. Few-Shot Continual Learning (FSCL) with Metric-Based Method

In our setting of FSCL, a model aims to continually learn from a few examples, tackle a sequence of novel tasks, and incrementally improve the ability of few-shot learning (FSL). Although metric-based method is also a representative strategy for few-shot classification, continual learning (CL) of novel classes from a few examples is not easy. To continually improve few-shot generalization from incoming tasks and examples, metric-based FSL generally requires additional query data to update the feature embedding layer.

We notice that several recent works in this direction attempt to learn a few novel tasks / classes (sampled from a small query set) from their relations to the base tasks / classes (sampled from a large support set) and improve the joint prediction of both [32, 27]. We stress it is different from our setting of FSCL: The information of base tasks / classes and the support set might be unavailable in continual learning of novel tasks / classes, which is a practical assumption in real world scenarios. Also, prediction of only novel tasks / classes can better evaluate the ability of few-shot generalization in continual learning, compared with evaluation on the tasks / classes seen in the support set.

To achieve a considerable performance on a sequence of novel tasks / classes, another straightforward idea is to store all the training data and use a fixed feature embedding layer for joint prediction. However, such "continual inference" strategy cannot improve few-shot generalization from novel tasks and examples, and thus cannot adapt to domain differences or changing domains in continual learning.
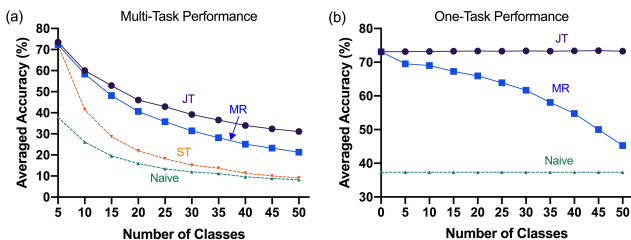


Figure 8. FSCL on miniImageNet with ProtoNet *Linear Classifier*. (a) is the performance to infer all the incremental classes ever seen. (b) is the performance on one 5-way 5-shot task during continual learning in (a). JT: Joint Training; ST: Sequential Training; MR: ST with memory replay of all the data ever seen; Naive: A randomly-initialized network without pretraining (PT).

To adapt metric-based FSL to FSCL, we propose a simple strategy: to use a linear classifier with softmax as the output layer in the FSCL stage, which has been shown as an effective strategy to improve generalization from $S$ to $D_T$, particularly in large domain differences [8]. Then we examine if the ability of FSL is interfered by CL in FSCL. We choose ProtoNet [31] as the base model of metric-based FSL and train such a classifier on $S$. Next, we use a lin-

ear classifier with softmax as the output layer and continually train the classifier on a sequence of novel classification tasks. Similar to the observations in gradient-based FSL, the performance of FSCL with replay of all the old training data (ST+Replay) is significantly lower than joint training (JT) of all the tasks (Appendix Fig. 1, a). The ability of FSL also gradually decreases in CL (Appendix Fig. 1, b).

The extension of ProtoNet *Linear Classifier* can be easily implemented as the base model in TSC. We conduct the same FSCL experiments on *New Class* as the main text. As shown in Appendix Table 1 and 2, TSC on ProtoNet *Linear Classifier* achieves a better performance than joint training of all the tasks ever seen, and the ability of few-shot generalization is incrementally improved from continual learning.

Table 5. Averaged accuracy (%) (top-1 / top-5) of class-incremental learning, averaged by 20 randomly sampled task sequences.

| Methods | miniImageNet* | | tieredImageNet | | tieredImageNet → CUB | |
|---|---|---|---|---|---|---|
| | $A_{30}$ | $A_{50}$ | $A_{30}$ | $A_{50}$ | $A_{30}$ | $A_{50}$ |
| Joint Training | 36.65 / 69.01 | 27.86 / 56.03 | 37.95 / 70.93 | 29.25 / 60.47 | 34.79 / 68.72 | 27.53 / 57.75 |
| Memory Replay | 27.36 / 54.69 | 17.05 / 36.81 | 26.78 / 53.71 | 16.08 / 35.85 | 21.96 / 48.76 | 14.47 / 33.53 |
| STC-M | 36.77 / 69.60 | 28.88 / 57.23 | 35.94 / 72.33 | 27.31 / 58.69 | 31.24 / 68.48 | 23.56 / 54.06 |
| TSC | **37.26 / 71.05** | **29.81 / 58.94** | 37.39 / **73.24** | **30.33 / 62.77** | **35.80 / 69.49** | **27.79 / 58.04** |

Table 6. Averaged absolute accuracy (%) / increased accuracy (%) of 5W5S task on various benchmark datasets in FSCL, averaged by more than 100 randomly sampled tasks.

| | | Task 0 | Task 5 | Task 10 |
|---|---|---|---|---|
| miniImageNet* | $k = 20$ (Novel Class) | 75.88 | 77.92 / +2.04 | 78.47 / +2.59 |
| | $k = 100$ (Novel Class) | 75.88 | 78.85 / +2.97 | **79.28 / +3.40** |
| tieredImageNet | $k = 20$ (Novel Class) | 70.31 | 73.01 / +2.70 | 73.51 / +3.20 |
| | $k = 100$ (Novel Class) | 70.31 | 73.83 / +3.52 | **74.27 / +3.96** |
| tieredImageNet → CUB | $k = 20$ (Novel Class) | 60.77 | 63.46 / +2.69 | 64.09 / +3.32 |
| | $k = 100$ (Novel Class) | 60.77 | 64.96 / +4.19 | **65.71 / +4.94** |

## B. Implementations of Base Model

Our implementation of MAML [9] is validated in Appendix Table 3. miniImageNet is the original split proposed by [33], which applies 64, 16 and 20 classes for training, validation and testing, respectively. To evaluate FSCL on large amounts of classes, we applies the entire 100 classes of miniImageNet (refer to as miniImageNet*) as the query set and the remained 900 classes in ImageNet as the support set. We follow similar implementations as [8, 9] and reproduce the reported performance on various datasets and backbones [9, 8, 14]. Note that we use 5-way classification for meta-training, because the ability of FSL might be improved from the incremental classes (ways) in FSCL. Then we apply the same implementation, but use ResNet-18 as the feature extractor, on miniImageNet* and tieredImageNet.

## C. Model Hyperparameters

We use the Adam optimizer with $beta1 = 0.5$ and $beta2 = 0.999$ to learn the classifier in all experiments.

Table 7. Averaged 5W5S accuracy (%) of MAML on various datasets. The performance of our results is averaged by more than 100 randomly sampled tasks.

|  | Backbone | Omniglot | miniImageNet | tieredImageNet |
|---|---|---|---|---|
| MAML (Reported) | Conv-4 | 99.9 [9] | $63.15 \pm 0.91$ [9] | 70.30 [14] |
|  |  |  | $62.71 \pm 0.71$ [8] |  |
| MAML (Ours) | Conv-4 | $99.80 \pm 0.11$ | $63.22 \pm 0.74$ | $70.12 \pm 0.69$ |

|  | Backbone | miniImageNet | miniImageNet* | tieredImageNet |
|---|---|---|---|---|
| MAML (Reported) | ResNet-18 | $65.72 \pm 0.77$ [8] | – | – |
| MAML (Ours) | ResNet-18 | $67.21 \pm 0.73$ | $77.90 \pm 0.67$ | $73.22 \pm 0.60$ |

The hyperparameters to update fast /slow weights include the learning rate of $\theta_s$ ($\beta$) and the iterations to learn $\theta_f^e$ ($k$), both of which have been extensively evaluated in the main text. The model hyperparameters are summarized in Table 8.

Table 8. Hyperparameters of our methods.

| Hyperparameters | Omniglot | miniImageNet* | tieredImageNet | tieredImageNet $\rightarrow$ CUB |
|---|---|---|---|---|
| Learning Rate | 0.0005 | 0.0005 | 0.0005 | 0.0005 |
| $\beta$ | 0.01 | 0.01 | 0.01 | 0.01 |
| $k$ | 100 | 100 | 100 | 100 |
| $\lambda$ (STC-M) | 0.1 | 0.1 | 0.1 | 0.1 |
| $\lambda$ (TSC) | $10^{-10}$ | $10^{-10}$ | $10^{-10}$ | $10^{-10}$ |
| $m$ | 1 | 1 | 0.1 | 0.1 |
| # of Epoch | 500 | 500 | 500 | 500 |

For the experiments of generative replay, we applies a 4-layer DCGAN architecture [24] with WGAN-GP[11] as the generative model to learn the old data distributions and replay generated data. However, a typical DCGAN architecture shows severe mode collapse to learn only a few training examples. This issue is caused by the BatchNorm layers in the generator and the discriminator, which can be replaced by InstanceNorm. The generator used in our implementation takes 3.80M parameters. The memory cost is equal to around 1853 images of the size $64 \times 64$ (one image takes around 0.0082MB). To avoid the large memory cost of saving a generator, we sample 20 generated images after training on each class and store them in the memory buffer for replay.