Meta-Learning with Self-Improving Momentum Target

Anonymous Author(s) Affiliation Address email

Abstract

The idea of using a separately trained target model (or *teacher*) to improve the 1 performance of the student model has been increasingly popular in various machine 2 learning domains, and meta-learning is no exception; a recent discovery shows that 3 utilizing task-wise target models can significantly boost the generalization perfor-4 mance. However, obtaining a target model for each task can be highly expensive, 5 especially when the number of tasks for meta-learning is large. To tackle this issue, 6 we propose a simple yet effective method, coined Self-improving Momentum Target 7 (SiMT). SiMT generates the target model by adapting from the temporal ensemble 8 of the meta-learner, i.e., the momentum network. This momentum network and its 9 task-specific adaptations enjoy a favorable generalization performance, enabling 10 self-improving of the meta-learner through knowledge distillation. Moreover, we 11 found that perturbing parameters of the meta-learner, e.g., dropout, further stabilize 12 this self-improving process by preventing fast convergence of the distillation loss 13 during meta-training. Our experimental results demonstrate that SiMT brings a 14 significant performance gain when combined with a wide range of meta-learning 15 methods under various applications, including few-shot regression, few-shot classi-16 17 fication, and meta-reinforcement learning.

18 1 Introduction

19 Meta-learning [48] is the art of extracting and utilizing the knowledge from the distribution of tasks 20 to better solve a relevant task. This problem is typically approached by training a meta-model that can transfer its knowledge to a task-specific solver, where the performance of the meta-model is evaluated 21 on the basis of how well each solver performs on the corresponding task. To learn such meta-model, 22 one should be able to (a) train an appropriate solver for each task utilizing the knowledge transferred 23 from the meta-model, and (b) accurately evaluate the performance of the solver. A standard way to do 24 this is the so-called S/Q (support/query) protocol [52, 32]: for (a), use a set of support set samples to 25 train the solver; for (b), use another set of samples, called *query set* samples to evaluate the solver¹. 26

Recently, however, an alternative paradigm—called S/T (support/target) protocol—has received much attention [55, 59, 30]. The approach assumes that the meta-learner has an access to task-specific *target models*, i.e., an expert model for each given task, and uses these models to evaluate task-specific solvers by measuring the discrepancy of the solvers from the target models. Intriguingly, it has been observed that such knowledge distillation procedure [40, 20] helps to improve the meta-generalization performance [59], in a similar way that such teacher-student framework helps to avoid overfitting under non-meta-learning contexts [28, 22].

³⁴ Despite such advantage, the S/T protocol is difficult to be used in practice, as training target models ³⁵ for each task usually requires an excessive amount of compute, especially when the number of tasks

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

¹We give an overview of terminologies used in the paper to guide readers new to this field (see the Appendix).



Figure 1: An overview of the proposed *Self-improving Momentum Target (SiMT)*: the momentum network efficiently generates the target model, and by distilling knowledge to the task-specific solver, it forms a self-improving process. S and Q denote the support and query datasets, respectively.

is large. Prior works aim to alleviate this issue by proposing methods to generate target models in

³⁷ a compute-efficient manner. For instance, Lu et al. [30] consider the case where the learner has an

access to a model pre-trained on a global data domain that covers most tasks (to be meta-trained

³⁹ upon), and propose to generate task-wise target models by simply fine-tuning the model for each task.

40 However, the method still requires the compute for fine-tuning on a large number of tasks, and more 41 importantly, is hard to be used when there is no effective pre-trained model available at hand. For

⁴¹ importantly, is hard to be used when there is no effective pre-trained model available at hand. For ⁴² example, a globally pre-trained model is usually not available in reinforcement learning, as collecting

43 "global" data is a nontrivial task [8].

In this paper, we ask whether we can generate the task-specific target models by (somewhat ironically) using meta-learning. We draw inspiration from recent observations in semi/self-supervised learning literature [47, 15, 5] that the temporal ensemble of a model, i.e., the momentum network [25], can be an effective teacher of the original meta-model. It turns out that a similar phenomenon happens in the meta-learning scenario: one can construct a momentum network of the meta-model, whose task-specific adaptation is an effective target model from which the task-specific knowledge can be distilled to train the original meta-model.

⁵¹ **Contribution.** We establish a novel framework, coined *Meta-Learning with Self-improving Momen-*⁵² *tum Target (SiMT)*, which brings the benefit of the S/T protocol to the S/Q-like scenario where ⁵³ task-specific target models are not available (but have access to query data). In a nutshell, SiMT is ⁵⁴ comprised of two (iterative) steps:

- *Momentum target*: We generate the target model by adapting from the momentum network, which shows better adaptation performance than the meta-model itself. In this regard, generating the target model becomes highly efficient, e.g., one single forward is required when obtaining the momentum target for ProtoNet [42].
- Self-improving process: The meta-model enables to improve through the knowledge distillation from the momentum target, and this recursively improves the momentum network by the temporal ensemble. Furthermore, we find that perturbing parameters of the task-specific solver of the meta-model, e.g., dropout [44], further stabilizes this self-improving process by preventing fast
- ⁶³ convergence of the distillation loss during meta-training.

We verify the effectiveness of SiMT under various applications of meta-learning, including few-64 shot regression, few-shot classification, and meta-reinforcement learning (meta-RL). Overall, our 65 experimental results show that incorporating the proposed method can consistently and significantly 66 improve the baseline meta-learning methods [9, 29, 34, 42]. In particular, our method improves the 67 few-shot classification accuracy of Conv4 [52] trained with MAML [9] on mini-ImageNet [52] from 68 $47.33\% \rightarrow 51.49\%$ for 1-shot, and from $63.27\% \rightarrow 68.74\%$ for 5-shot, respectively. Moreover, we 69 show that our framework could even notably improve on the few-shot regression and meta-RL tasks, 70 which supports that our proposed method is indeed domain-agnostic. 71

72 2 Related work

Learning from target models. Learning from an expert model, i.e., the target model, has shown its 73 74 effectiveness across various domains [28, 33, 62, 49]. As a follow-up, recent papers demonstrate that meta-learning can also be the case [55, 59]. However, training independent task-specific target models 75 is highly expensive due to the large space of task distribution in meta-learning. To this end, recent 76 work suggests pre-training a global encoder on the whole meta-training set and finetune target models 77 on each task [30]; however, they are limited to specific domains and still require some computations, 78 e.g., they take more than 6.5 GPU hours to pre-train only 10% of target models while ours require 2 79 80 GPU hours for the entire meta-learning process (ProtoNet [42] of ResNet-12 [32]) on the same GPU. Another recent relevant work is bootstrapped meta-learning [10], which generates the target model 81 from the meta-model by further updating the parameters of the task-specific solver for some number 82 of steps with the query dataset. While the bootstrapped target models can be obtained efficiently, their 83 approach is specialized in gradient-based meta-learning schemes, e.g., MAML [9]. In this paper, we 84 suggest an efficient and more generic way to generate the target model during the meta-training. 85 Learning with momentum networks. The idea of temporal ensembling, i.e., the momentum network, 86 has become an essential component of the recent semi/self-supervised learning algorithms [3, 5]. For 87 example, Mean Teacher [47] first showed that the momentum network improves the performance of 88 semi-supervised image classification, and recent advanced approaches [2, 43] adopted this idea for 89 achieving state-of-the-art performances. Also, in self-supervised learning methods which enforce 90 invariance to data augmentation, such momentum networks are widely utilized as a target network 91

[18, 15] to prevent collapse by providing smoother changes in the representations. In this paper, we
 empirically demonstrate that the momentum network shows better adaptation performance compare

⁹⁴ to the original meta-model, which motivates us to utilize it for generating the target model in a ⁹⁵ compute-efficient manner.

36 3 Problem setup and evaluation protocols

⁹⁷ In this section, we formally describe the meta-learning setup under consideration, and S/Q and ⁹⁸ S/T protocols studied in prior works.

Problem setup: Meta-learning. Let $p(\tau)$ be a distribution of tasks. The goal of meta-learning is 99 to train a meta-model f_{θ} , parameterized by the meta-model parameter θ , which can transfer its 100 knowledge to help to train a *solver* for a new task. More formally, we consider some *adaptation* 101 subroutine Adapt(\cdot, \cdot) which uses both information transferred from θ and the task-specific dataset 102 (which we call support set) S^{τ} to output a task-specific solver as $\phi^{\tau} = \text{Adapt}(\theta, S^{\tau})$. For example, 103 the model-agnostic meta-learning algorithm (MAML; [9]) uses the adaptation subroutine of taking 104 a fixed number of SGD on \mathcal{S}^{τ} , starting from the initial parameter θ . In this paper, we aim to give a 105 general meta-learning framework that can be used in conjunction with any adaptation subroutine, 106 instead of designing a method specialized for a specific one. 107

The objective is to learn a nice meta-model parameter θ from a set of tasks sampled from $p(\tau)$ (or sometimes the task distribution itself), such that the expected loss of the task-specific adaptations is small, i.e., $\min_{\theta} \mathbb{E}_{\tau \sim p(\tau)} [\ell^{\tau} (\operatorname{Adapt}(\theta, S^{\tau}))]$, where $\ell^{\tau}(\cdot)$ denotes the test loss on task τ . To train such meta-model, we need a mechanism to evaluate and optimize θ (e.g., via gradient descent). For this purpose, existing approaches take one of two approaches: the S/Q protocol or the S/T protocol.

¹¹³ *S*/*Q* **protocol.** The majority of the existing meta-learning frameworks (e.g., [52, 32]) splits the ¹¹⁴ task-specific training data into two, and use them for different purposes. One is the support set S^{τ} ¹¹⁵ which is used to perform the adaptation subroutine. Another is the *query set* Q^{τ} which is used for ¹¹⁶ evaluating the performance of the adapted parameter and compute the gradient with respect to θ . In ¹¹⁷ other words, given the task datasets $(S_1, Q_1), (S_2, Q_2), \ldots, (S_N, Q_N)$,² the *S*/*Q* protocol solves

...

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}(\texttt{Adapt}(\theta, \mathcal{S}^{\tau_i}), \mathcal{Q}^{\tau_i}),$$
(1)

where $\mathcal{L}(\phi, Q)$ denotes the empirical loss of a solver ϕ on the dataset Q.

²Here, while we assumed a static batch of tasks for notational simplicity, the expression is readily extendible to the case of a *stream of tasks* drawn from $p(\tau)$.

¹¹⁹ S/T **protocol.** Another line of work considers the scenario where the meta-learner additionally has ¹²⁰ an access to a set of *target models* ϕ_{target} for each training task [55, 30]. In such case, one can ¹²¹ use a teacher-student framework to regularize the adapted solver to behave similarly (or have low ¹²² *prediction discrepancy*, equivalently) to the target model. Here, a typical practice is to *not split* each ¹²³ task dataset and measure the discrepancy using the support dataset that is used for the adaptation ¹²⁴ [30]. In other words, given the task datasets S_1, S_2, \ldots, S_N and the corresponding target models

125 $\phi_{\text{target}}^{\tau_1}, \phi_{\text{target}}^{\tau_2}, \dots, \phi_{\text{target}}^{\tau_N}$, the S/T protocol updates the meta-model by solving

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_{\texttt{teach}} \big(\texttt{Adapt}(\theta, \mathcal{S}^{\tau_i}), \phi_{\texttt{target}}^{\tau_i}, \mathcal{S}^{\tau_i} \big), \tag{2}$$

where $\mathcal{L}_{\text{teach}}(\phi, \phi_{\text{target}}, S)$ denotes a discrepancy measure between the adapted model ϕ and the target model ϕ_{target} , measured using the dataset S.

128 4 Meta-learning with self-improving momentum target

In this section, we develop a compute-efficient framework which bring the benefits of S/T protocol 129 to the settings where we do not have access to target-specific tasks or a general pretrained model, as 130 in general S/Q-like setups. In a nutshell, our framework iteratively generates a *meta-target model* 131 which generalizes well when adapted to the target tasks, by constructing a momentum network [47] 132 133 of the meta-model itself. The meta-model is then trained, using both the knowledge transferred from the momentum target and the knowledge freshly learned from the query sets. We first briefly describe 134 our meta-model update protocol (Section 4.1), and then the core component, coined Self-Improving 135 *Momentum Target (SiMT)*, which efficiently generates the target model for each task (Section 4.2). 136

137 4.1 Meta-model update with a S/Q-S/T hybrid loss

To update the meta-model, we use a hybrid loss function of the S/Q protocol (1) and the S/T protocol (2). Formally, let $(S_1, Q_1), (S_2, Q_2), \dots, (S_N, Q_N)$ be given task datasets with support-query split, and let $\phi_{\text{target}}^{\tau_1}, \phi_{\text{target}}^{\tau_2}, \dots, \phi_{\text{target}}^{\tau_N}$ be task-specific target models generated by our target generation procedure (which will be explained with more detail in Section 4.2). We train the meta-model as

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^{N} \left((1-\lambda) \cdot \mathcal{L}(\texttt{Adapt}(\theta, \mathcal{S}^{\tau_i}), \mathcal{Q}^{\tau_i}) + \lambda \cdot \mathcal{L}_{\texttt{teach}}(\texttt{Adapt}(\theta, \mathcal{S}^{\tau_i}), \phi_{\texttt{target}}^{\tau_i}, \mathcal{Q}^{\tau_i}) \right), \quad (3)$$

where $\lambda \in [0, 1)$ is the weight hyperparameter. We note two things about Eq. 3. First, while we are 142 training using the target model, we also use a S/Q loss term. This is because our method trains the 143 meta-target model and the meta-model simultaneously from scratch, instead of requiring fully-trained 144 target models. Second, unlike in the S/T protocol, we evaluate the discrepancy \mathcal{L}_{teach} using the 145 query set Q^{τ_i} instead of the support set, to improve the generalization performance of the student 146 model. In particular, the predictions of adapted models on query set samples are softer (i.e., having 147 less confidence) than on support set samples, and such soft predictions are known to be beneficial on 148 the generalization performance of the student model in the knowledge distillation literature [61, 46]. 149

150 4.2 SiMT: Self-improving momentum target

We now describe the algorithm we propose, SiMT (Algorithm 1), to generate the target model in a compute-efficient manner. In a nutshell, SiMT is comprised of two iterative steps: *momentum target* and *self-improving process*. To efficiently generate a target model, SiMT utilizes the temporal ensemble of the network, i.e., the momentum network, then distills the knowledge of the generated target model into the task-specific solver of the meta-model to form a self-improving process.

Momentum target. For the compute-efficient generation of target models, we utilize the momentum network θ_{moment} of the meta-model. Specifically, after every meta-model training iteration, we compute the exponential moving average of the meta-model parameter θ as

$$\theta_{\text{moment}} \leftarrow \eta \cdot \theta_{\text{moment}} + (1 - \eta) \cdot \theta,$$
 (4)

where $\eta \in [0, 1)$ is the momentum coefficient. We find that θ_{moment} can adapt better than the metamodel θ itself and observe that the loss landscape has flatter minima (see Section 5.4), which can Algorithm 1 SiMT: Self-Improving Momentum Target

Require: Distribution over tasks $p(\tau)$, adaptation subroutine Adapt(·), momentum coefficient η , weight hyperparameter λ , dropout probability p, task batch size N, learning rate β .

- 1: Initialize θ using the standard initialization scheme.
- 2: Initialize the momentum network with the meta-model parameter, $\theta_{moment} \leftarrow \theta$.
- 3: while not done do
- Sample N tasks $\{\tau_i\}_{i=1}^N$ from $p(\tau)$ 4:
- for i = 1 to N do 5:
- Sample support set S^{τ_i} and query set Q^{τ_i} from τ_i 6:
- 7: $\phi_{\texttt{moment}}^{\tau_i} = \texttt{Adapt}(\theta_{\texttt{moment}}, \mathcal{S}^{\tau_i}).$ ▷ Generate a momentum target. $\phi^{\tau_i} = \operatorname{Adapt}(\theta, \mathcal{S}^{\tau_i}).$ ▷ Adapt a task-specific solver. 8:
 $$\begin{split} \phi_{\text{drop}}^{\tau_i} &= \text{Dropout}(\phi^{\tau_i}, p). & \triangleright \text{Perturb the solver.} \\ \mathcal{L}_{\text{total}}^{\tau_i}(\theta) &= (1-\lambda) \cdot \mathcal{L}(\phi_{\text{drop}}^{\tau_i}, \mathcal{Q}^{\tau_i}) + \lambda \cdot \mathcal{L}_{\text{teach}}(\phi_{\text{drop}}^{\tau_i}, \phi_{\text{moment}}^{\tau_i}, \mathcal{Q}^{\tau_i}) & \triangleright \text{Compute loss.} \end{split}$$
 9: 10: end for 11: $\begin{array}{l} \theta \leftarrow \theta - \frac{\beta}{N} \cdot \nabla_{\theta} \sum_{i=1}^{N} \mathcal{L}_{\texttt{total}}^{\tau_{i}}(\theta). \\ \theta_{\texttt{moment}} \leftarrow \eta \cdot \theta_{\texttt{moment}} + (1 - \eta) \cdot \theta. \end{array}$ \triangleright Train the meta-model. 12: 13: ▷ Update the momentum network. 14: end while
- be a hint for understanding the generalization improvement [27, 11]. Based on this, we propose 161 to generate the task-specific target model, i.e., the momentum target ϕ_{moment} , by adapting from the 162 momentum network θ_{moment} . For a given support set S, we generate the target model for each task as 163
 - $\phi_{\texttt{moment}}^{\tau_i} = \texttt{Adapt}(\theta_{\texttt{moment}}, \mathcal{S}^{\tau_i}), \qquad \forall i \in \{1, 2, \dots, N\}.$ (5)

We remark that generating momentum targets does not require an excessive amount of compute (see 164 Section 5.4), e.g., ProtoNet [42] requires a single forward of a support set, and MAML [9] requires 165

few-gradient steps without second-order gradient computation for the adaptation. 166

Self-improving process via knowledge distillation. After generating the momentum target, we 167 utilize its knowledge to improve the generalization performance of the meta-model. To this end, we 168 choose the knowledge distillation scheme [20], which is simple yet effective across various domains, 169 including meta-learning [30]. Here, our key concept is that the momentum target self-improves 170 during the training due to the knowledge transfer. To be specific, the knowledge distillation from 171 the momentum target improves the meta-model itself, which recursively improves the momentum 172 through the temporal ensemble. Formally, for a given query set \mathcal{Q} , we distill the knowledge of the 173 174

momentum target ϕ_{moment} to the task-specific solver of the meta-model ϕ as

$$\mathcal{L}_{\texttt{teach}}(\phi, \phi_{\texttt{moment}}, \mathcal{Q}) := \frac{1}{|\mathcal{Q}|} \sum_{(x,y) \in \mathcal{Q}} l_{\texttt{KD}} \big(f_{\phi_{\texttt{moment}}}(x), f_{\phi}(x) \big), \tag{6}$$

where l_{KD} is the distillation loss and $|\cdot|$ is the cardinality of the set. For regression tasks, we use the 175 MSE loss, i.e., $l_{\text{KD}}(z_1, z_2) := ||z_1 - z_2||_2^2$, and for classification tasks, we use the KL divergence with temperature scaling [16], i.e., $l_{\text{KD}}(z_1, z_2) := T^2 \cdot \text{KL}(\sigma(z_1/T) || \sigma(z_2/T))$, where T is the temperature temperature scaling [16], i.e., $l_{\text{KD}}(z_1, z_2) := T^2 \cdot \text{KL}(\sigma(z_1/T) || \sigma(z_2/T))$, where T is the temperature tempera 176 177 hyperparameter, σ is the softmax function and z_1, z_2 are logits of the classifier, respectively. We 178 present the detailed distillation objective of reinforcement learning tasks in the Appendix. Also, 179 note that optimizing the distillation loss only propagates gradients to the meta-model θ , not to the 180 momentum network θ_{moment} , i.e., known as the stop-gradient operator [5, 6]. 181

Furthermore, we find that the distillation loss (6) sometimes converges too fast during the meta-182 training, which can stop the self-improving process. To prevent this, we suggest perturbing the 183 parameter space of ϕ . Intuitively, injecting noise to the parameter space of ϕ forces an asymmetricity 184 between the momentum target's prediction, hence, preventing f_{ϕ} and $f_{\phi_{moment}}$ from reaching a similar 185 prediction. To this end, we choose the standard dropout regularization [44] due to its simplicity and 186 generality across architectures and also have shown its effectiveness under distillation research [57]: 187 $\phi_{drop} := Dropout(\phi, p)$ where p is the probability of dropping activations. In the end, we use the 188 perturbed task-specific solver ϕ_{drop} and the momentum target ϕ_{moment} for our evaluation protocol (3). 189

Table 1: Few-shot regression results on ShapeNet and Pascal datasets. We report the angular error for ShapeNet, and MSE for Pascal. SiMT utilizes the momentum network for the adaptation. The reported results are averaged over three trials, subscripts denote the standard deviation, and bold denotes the best result of each group.

	ShapeNet		Pascal	
Method	10-shot	15-shot	10-shot	15-shot
MAML [9] MAML [9] + SiMT	$\begin{array}{c} 29.555 {\pm} 0.600 \\ \textbf{18.913} {\pm} \textbf{2.655} \end{array}$	$\begin{array}{c} 22.286 {\pm} 3.369 \\ \textbf{16.100} {\pm} \textbf{1.318} \end{array}$	2.612±0.280 1.462±0.230	2.513±0.250 1.229±0.074
ANIL [34] ANIL [34] + SiMT	39.915±0.665 37.424±0.951	38.202±1.388 29.478±0.212	$\begin{array}{c} 6.600{\pm}0.360\\ \textbf{5.339}{\pm}\textbf{0.321}\end{array}$	6.517±0.420 5.007±0.145
MetaSGD [29] MetaSGD [29] + SiMT	$\begin{array}{c} 17.353 {\pm} 1.110 \\ \textbf{16.121} {\pm} \textbf{1.322} \end{array}$	$\begin{array}{c} 15.768 {\pm} 1.266 \\ \textbf{14.377} {\pm} \textbf{0.358} \end{array}$	3.532 ± 0.381 2.300 ± 0.871	$\begin{array}{c} 2.833 {\pm} 0.216 \\ \textbf{1.879} {\pm} \textbf{0.134} \end{array}$

Table 2: Few-shot in-domain adaptation accuracy (%) on mini-ImageNet, and tiered-ImageNet. SiMT utilizes the momentum network for the adaptation. The reported results are averaged over three trials, subscripts denote the standard deviation, and bold denotes the best result of each group.

		mini-ImageNet		tiered-ImageNet	
Model	Method	1-shot	5-shot	1-shot	5-shot
Conv4 [52]	MAML [9] MAML [9] + SiMT	$\begin{array}{c} 47.33 {\pm} 0.45 \\ \textbf{51.49} {\pm} \textbf{0.18} \end{array}$	63.27±0.14 68.74±0.12	$\begin{array}{c} 50.19{\pm}0.21\\ \textbf{52.51}{\pm}\textbf{0.21}\end{array}$	66.05±0.19 69.58±0.11
	ANIL [34] ANIL [34] + SiMT	47.71±0.47 50.81±0.56	63.13±0.43 67.99±0.19	49.57±0.04 51.66±0.26	$\begin{array}{c} 66.34{\pm}0.28\\ \textbf{68.88}{\pm}\textbf{0.08}\end{array}$
	MetaSGD [29] MetaSGD [29] + SiMT	50.66±0.18 51.70±0.80	65.55±0.54 69.13±1.40	$52.48{\scriptstyle\pm1.22}\\52.98{\scriptstyle\pm0.07}$	71.06±0.20 71.46±0.12
	ProtoNet [42] ProtoNet [42] + SiMT	$\begin{array}{c} 47.97 {\pm} 0.29 \\ \textbf{51.25} {\pm} \textbf{0.55} \end{array}$	65.16±0.67 68.71±0.35	51.90±0.55 53.25±0.27	71.51±0.25 72.69±0.27
ResNet-12 [32]	MAML [9] MAML [9] + SiMT	$52.66{\scriptstyle\pm0.60} \\ \textbf{56.28}{\scriptstyle\pm0.63}$	68.69±0.33 72.01±0.26	57.32±0.59 59.72±0.22	73.78±0.27 74.40±0.90
	ANIL [34] ANIL [34] + SiMT	51.80±0.59 54.44 ± 0.27	$\begin{array}{c} 68.38 {\pm} 0.20 \\ \textbf{69.98} {\pm} \textbf{0.66} \end{array}$	57.52±0.68 58.18±0.31	$\begin{array}{c} 73.50 {\pm} 0.35 \\ \textbf{75.59} {\pm} \textbf{0.50} \end{array}$
	MetaSGD [29] MetaSGD [29] + SiMT	54.95±0.11 55.72±0.96	$70.65{\scriptstyle\pm0.43} \\ \textbf{74.01}{\scriptstyle\pm0.79}$	58.97±0.89 61.03±0.05	$76.37{\scriptstyle\pm 0.11} \\ \textbf{78.04}{\scriptstyle\pm \textbf{0.48}}$
	ProtoNet [42] ProtoNet [42] + SiMT	52.84±0.21 55.84±0.57	$\begin{array}{c} 68.35 {\pm} 0.29 \\ \textbf{72.45} {\pm} \textbf{0.32} \end{array}$	61.16±0.17 62.01±0.42	$\begin{array}{c} 79.94{\scriptstyle\pm0.20}\\ \textbf{81.82}{\scriptstyle\pm0.12}\end{array}$

190 5 Experiments

In this section, we experimentally validate the effectiveness of the proposed SiMT by measuring its performance on various meta-learning applications, including few-shot regression (Section 5.1), few-shot classification (Section 5.2), and meta-reinforcement learning (meta-RL; Section 5.3).

Common setup. By following the prior works, we chose the checkpoints and the hyperparameters on the meta-validation set for the few-shot learning tasks [31, 53]. For RL, we chose it based on the best average return during the training [9]. We find that the hyperparameters, e.g., momentum coefficient η or the weight hyperparameter λ , are not sensitive across datasets and architectures but can vary on the type of the meta-learning scheme or tasks. We provide further details in the Appendix. Moreover, we report the adaptation performance of the *momentum network* for SiMT.

200 5.1 Few-shot regression

For regression tasks, we demonstrate our experiments on ShapeNet [12] and Pascal [60] datasets, where they aim to predict the object pose of a gray-scale image relative to the canonical orientation. To this end, we use the following empirical loss \mathcal{L} to train the meta-model: the angular loss for the ShapeNet $(\sum_{(x,y)\in Q} \|\cos(f_{\phi}(x)) - \cos(y)\|^2 + \|\sin(f_{\phi}(x)) - \sin(y)\|^2)$ and the MSE loss for the

Table 3: Few-shot cross-domain adaptation accuracy (%) on ResNet-12 trained with mini-ImageNet and tiered-ImageNet. We consider CUB and Cars as cross-domain datasets. SiMT utilizes the momentum network for the adaptation. The reported results are averaged over three trials, subscripts denote the standard deviation, and bold denotes the best result of each group.

		mini-ImageNet $ ightarrow$		tiered-ImageNet $ ightarrow$	
Problem	Method	CUB	Cars	CUB	Cars
1-shot	MAML [9] MAML [9] + SiMT	$\begin{array}{c} 39.50 {\pm} 0.91 \\ \textbf{42.32} {\pm} \textbf{0.62} \end{array}$	32.87±0.20 33.73±0.63	$\begin{array}{c} 42.32{\pm}0.69\\ \textbf{44.33}{\pm}\textbf{0.43}\end{array}$	36.62±0.12 37.21±0.35
	ANIL [34] ANIL [34] + SiMT	$\begin{array}{c} 37.30 {\pm} 0.89 \\ \textbf{38.86} {\pm} \textbf{0.98} \end{array}$	31.28±1.03 32.34±0.95	$\begin{array}{c} 42.29 {\pm} 0.33 \\ \textbf{44.53} {\pm} \textbf{1.21} \end{array}$	36.27±0.58 36.92±0.56
	MetaSGD [29] MetaSGD [29] + SiMT	41.98±0.18 43.50±0.89	34.52±0.56 33.92±0.30	46.48±2.10 46.62±0.41	38.09±1.21 38.69±0.26
	ProtoNet [42] ProtoNet [42] + SiMT	$\begin{array}{c} 41.22{\pm}0.81\\ \textbf{44.13}{\pm}\textbf{0.30} \end{array}$	32.79±0.61 34.53±0.40	$\begin{array}{c} 47.75 {\pm} 0.56 \\ \textbf{48.89} {\pm} \textbf{0.65} \end{array}$	37.59 ± 0.80 38.07 \pm 0.42
5-shot	MAML [9] MAML [9] + SiMT	56.17±0.92 59.22±0.39	44.56±0.79 46.59 ± 0.21	65.00±0.89 67.58±0.61	51.08±0.28 51.88±0.52
	ANIL [34] ANIL [34] + SiMT	53.42±0.97 56.03±1.40	41.65±0.67 45.88 ± 0.82	62.48±0.85 66.30±0.99	50.50±1.18 54.60±0.91
	MetaSGD [29] MetaSGD [29] + SiMT	$\begin{array}{c} 58.90{\pm}1.30\\ \textbf{65.07}{\pm}\textbf{1.89} \end{array}$	47.44±1.55 49.86 ± 0.84	$70.38{\scriptstyle\pm0.27} \\ \textbf{73.93}{\scriptstyle\pm0.42}$	$\begin{array}{c} 56.28 {\pm} 0.07 \\ \textbf{57.97} {\pm} \textbf{1.34} \end{array}$
	ProtoNet [42] ProtoNet [42] + SiMT	57.87±0.77 63.85±0.76	48.06±1.10 51.67±0.29	$\begin{array}{c} 74.35 {\pm} 0.93 \\ \textbf{75.97} {\pm} \textbf{0.09} \end{array}$	$57.23{\pm}0.25\\\textbf{59.01}{\pm}\textbf{0.50}$

Pascal $(\sum_{(x,y)\in Q} || f_{\phi}(x) - y ||^2)$, following the prior works [60, 12]. For the backbone meta-learning schemes we use gradient-based approaches, including MAML [9], ANIL [34], and MetaSGD [29]. For all methods, we train the convolutional neural network with 7 layers [60] and apply dropout regularization [44] before the max-pooling layer for SiMT. Table 1 summarizes the results, showing that SiMT significantly improves the overall meta-learning schemes in all tested cases.

210 5.2 Few-shot classification

For few-shot classification tasks, we use the cross-entropy loss for the empirical loss term $\mathcal L$ to 211 train the meta-model θ , i.e., $\sum_{(x,y)\in\mathcal{Q}} l_{ce}(f_{\phi}(x), y)$ where l_{ce} is the cross-entropy loss. We train the 212 meta-model on mini-ImageNet $\begin{bmatrix} 52 \end{bmatrix}$ and tiered-ImageNet $\begin{bmatrix} 36 \end{bmatrix}$ datasets, following the prior works 213 [30, 53]. Here, we consider the following gradient-based and metric-based meta-learning approaches 214 as our backbone algorithm to show the wide usability of our method: MAML, ANIL, MetaSGD, and 215 ProtoNet [42]. We train each method on Conv4 [52] and ResNet-12 [32], and apply dropout before 216 the max-pooling layer for SiMT. For the training details, we mainly follow the setups from each 217 backbone algorithm paper. See the Appendix for more details. 218

In-domain adaptation. In this setup, we evaluate the adaptation performance on different classes 219 of the same dataset used in meta-training. As shown in Table 2, incorporating SiMT into existing 220 meta-learning methods consistently and significantly improves the in-domain adaptation performance. 221 In particular, SiMT achieves higher accuracy gains on the mini-ImageNet dataset, e.g., 5-shot 222 performance improves from $63.27\% \rightarrow 68.74\%$ on Conv4. We find that this is due to the overfitting 223 issue of backbone algorithms on the mini-ImageNet dataset, where SiMT is more robust to such 224 issues. For instance, when training mini-ImageNet 5-shot classification on Conv4, MAML starts to 225 overfit after the first 40% of the training process, while SiMT does not overfit during the training. 226

Cross-domain adaptation. We also consider the cross-domain adaptation scenarios. Here, we adapt the meta-model on different datasets from the meta-training: we use CUB [54] and Cars [24] datasets. Such tasks are known to be challenging, as there exists a large distribution shift between training and testing domains [17]. Table 3 shows the results. Somewhat interestingly, SiMT also improves the cross-domain adaptation performance of the base meta-learning methods across the considered datasets. These results indicate that SiMT successfully learns the ability to generalize to unseen tasks even for the distributions that highly differ from the training.

Table 4: Ablation study on each component of SiMT. We report the few-shot in-domain adaptation accuracy (%) on Conv4 trained with mini-ImageNet. Here, we use the learned momentum network at meta-test time, except for the first experiment of the table. The reported results are averaged over three trials, subscripts denote the standard deviation, and bold denotes the best result.

Momentum	Distillation	Dropout	1-shot	5-shot
-	-	-	$47.33{\scriptstyle\pm0.45}$	$63.27 {\pm} 0.14$
\checkmark	-	-	$48.98{\scriptstyle\pm0.32}$	66.12 ± 0.21
\checkmark	\checkmark	-	$49.23{\scriptstyle\pm0.24}$	66.52 ± 0.15
\checkmark	-	\checkmark	$49.25{\scriptstyle\pm0.41}$	$65.25{\pm}0.15$
\checkmark	\checkmark	\checkmark	$51.49{\scriptstyle\pm0.18}$	$68.74{\scriptstyle\pm0.12}$

234 5.3 Reinforcement learning

235 The goal of meta-RL is training an agent to quickly adapt a policy to maximize the expected return for unseen tasks using only a limited number of sample trajectories. Since the expected return is usually 236 not differentiable, we use policy gradient methods to update the policy. Specifically, we use vanilla 237 policy gradient [56], and trust-region policy optimization (TRPO; [37]) for the task-specific solver 238 and meta-model, respectively, following MAML [9]. The overall training objective of meta-RL is in 239 the Appendix, including the empirical loss \mathcal{L} , and the knowledge distillation loss \mathcal{L}_{teach} . We evaluate 240 SiMT on continuous control tasks based on OpenAI Gym [4] environments. In these experiments, we 241 242 choose MAML as our backbone algorithm, and train a multi-layer perceptron policy network with two hidden layers of size 100 by following the prior setup [9]. We find that the distillation loss is 243 already quite effective even without the dropout regularization, and applying it does not improve 244 more. We conjecture that dropout on such a small network may not be effective as it is designed to 245 reduce the overfitting of large networks [44]. We provide more experimental details in the Appendix. 246

2D Navigation. We first evaluate SiMT on 247 a 2D Navigation task, where a point agent 248 moves to different goal positions which are 249 randomly chosen within a 2D unit square. 250 Figure 2 shows the adaptation performance 251 of learned models with up to 3 gradient up-252 dates. These results demonstrate that SiMT 253 (red) could consistently improve the adap-254 tation performance of MAML (blue). Also, 255 SiMT makes faster performance improve-256 ments than vanilla MAML with additional 257 gradient updates. 258

Locomotion. To further demonstrate our



Figure 2: Meta-RL results for (a) 2D navigation and (b) Half-cheetah locomotion tasks. The solid line and shaded regions represent the truncated mean and standard deviation, respectively, across five runs.

method, we also study high-dimensional, complex locomotion tasks based on the MuJoCo [50] simulator. We choose a set of goal direction 261 tasks with a planar cheetah ("Half-cheetah"), following previous works [9, 34]. In the goal direction 262 tasks, the reward is the magnitude of the velocity in either the forward or backward direction, 263 randomly chosen for each task. Figure 2b shows that SiMT significantly improves the adaptation 264 performance of MAML even with a single gradient step. 265

5.4 Ablation study 266

259

260

267 Throughout this section, unless otherwise specified, we perform the experiments in 5-shot in-domain adaptation on mini-ImageNet with Conv4, where MAML is the backbone meta-learning scheme. 268

Component analysis. We perform an analysis on each component of our method in both 1-shot 269 and 5-shot classification on mini-ImageNet: namely, the use of (a) the momentum network θ_{moment} , 270 (b) the distillation loss \mathcal{L}_{teach} (6), and (c) the dropout regularization $Dropout(\cdot)$, by comparing the 271 accuracies. The results in Table 4 show each component is indeed important for the improvement. 272 We find that a naïve combination of the distillation loss and the momentum network does not show 273 significant improvements. But, by additionally applying the dropout, the distillation loss becomes 274 more effective and further improves the performance. Note that this improvement does not fully come 275 from the dropout itself, as only using dropout slightly degrades the performance in some cases. 276



(a) Computation efficiency

(b) Network choice for the adaptation

Figure 3: Validation accuracy curves of 5-shot mini-ImageNet on Conv4: we compare the adaptation performance of (a) MAML and SiMT, under the same training wall-clock time, and (b) the meta-model and the momentum network of SiMT, under the same number of training steps. The solid line and shaded regions represent the mean and standard deviation, respectively, across three runs.

277 Loss landscape of the momentum network.

We visualize the loss landscape of the momen-278 tum network θ_{moment} and the meta-model θ , to 279 give insights into the generalization improve-280 ment. To do this, we train MAML with a momen-281 tum network (without distillation and dropout) 282 and visualize the loss by perturbing each pa-283 rameter space [27] (See the Appendix for the 284 detail of the visualization method). As shown in 285 Figure 4, the momentum network forms a flat-286 ter loss landscape than the meta-model, where 287 recent studies demonstrate that such a flat land-288 scape is effective under various domains [11]. 289





Figure 4: Loss landscape visualization of Conv4 trained under 5-shot mini-ImageNet with MAML.

Computational efficiency. Our method may be seemingly compute-inefficient when incorporating meta-learning methods (due to the momentum target generation); however, we show that it is not. Although SiMT increases the total training time of MAML by roughly 1.2 times, we have observed that it is 3 times faster to achieve the best performance of MAML: in Figure 3a, we compare the accuracy under the same training wall-clock time with MAML.

Comparison of the momentum network and meta-model. To understand how the momentum network improves the performance of the meta-model, we compare the adaptation performance of the momentum network and the meta-model during training SiMT. As shown in Figure 3b, we observe that the performance of the momentum network is consistently better than the meta-model, which implies that the proposed momentum target is a nice target model in our self-improving mechanism.

300 6 Discussion and conclusion

In this paper, we propose a simple yet effective method, SiMT, for improving meta-learning. Our key idea is to efficiently generate target models using a momentum network and utilize its knowledge to self-improve the meta-learner. Our experiments demonstrate that SiMT significantly improves the performance of meta-learning methods on various applications.

Limitations and future work. While SiMT is a compute-efficient way to use target models in meta-learning, it is still built on top of existing meta-model update techniques. Since existing metalearning methods have limited scalability (to large-scale scenarios) [41], SiMT is no exception. Hence, improving the scalability of meta-learning schemes is an intriguing future research direction, where we believe incorporating SiMT into such scenarios is worthwhile.

Potential negative impacts. Meta-learning often requires a large computation due to the numerous
 task adaptation during meta-training, therefore raising environmental concerns, e.g., carbon generation
 [39]. As SiMT is built upon the meta-learning methods, practitioners may need to consider some
 computation for successful training. To address this issue, efficient meta-learning schemes [53] or
 lightweight methods for meta-learning [26] would be required for the applications.

315 **References**

- [1] S. M. R. Arnold, P. Mahajan, D. Datta, I. Bunner, and K. S. Zarkias. learn2learn: A library for
 Meta-Learning research. *arXiv preprint arXiv:2008.12284*, 2020.
- [2] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel. MixMatch:
 A holistic approach to semi-supervised learning. *Advances in Neural Information Processing Systems*, 2019.
- [3] D. Berthelot, N. Carlini, E. D. Cubuk, A. Kurakin, K. Sohn, H. Zhang, and C. Raffel. ReMix-Match: Semi-supervised learning with distribution alignment and augmentation anchoring. In *International Conference on Learning Representations*, 2020.
- [4] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. OpenAI gym. *arXiv preprint arXiv:1606.01540*, 2016.
- [5] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in self-supervised vision transformers. In *IEEE International Conference on Computer Vision*, 2021.
- [6] X. Chen and K. He. Exploring simple Siamese representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [7] Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement learning for continuous control. In *International Conference on Machine Learning*, 2016.
- [8] G. Dulac-Arnold, D. Mankowitz, and T. Hester. Challenges of real-world reinforcement learning.
 In *International Conference on Machine Learning*, 2019.
- [9] C. Finn, P. Abbeel, and S. Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, 2017.
- [10] S. Flennerhag, Y. Schroecker, T. Zahavy, H. van Hasselt, D. Silver, and S. Singh. Bootstrapped
 meta-learning. In *International Conference on Learning Representations*, 2022.
- [11] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for efficiently
 improving generalization. In *International Conference on Learning Representations*, 2021.
- [12] N. Gao, H. Ziesche, N. A. Vien, M. Volpp, and G. Neumann. What matters for meta-learning
 vision regression tasks? In *IEEE Conference on Computer Vision and Pattern Recognition*,
 2022.
- [13] V. Garcia and J. Bruna. Few-shot learning with graph neural networks. In *International Conference on Learning Representations*, 2018.
- [14] G. Ghiasi, T.-Y. Lin, and Q. V. Le. Dropblock: A regularization method for convolutional
 networks. In *Advances in Neural Information Processing Systems*, 2018.
- I.5] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch,
 B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, et al. Bootstrap your own latent a new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, 2020.
- [16] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In International Conference on Machine Learning, 2017.
- Y. Guo, N. C. Codella, L. Karlinsky, J. V. Codella, J. R. Smith, K. Saenko, T. Rosing, and
 R. Feris. A broader study of cross-domain few-shot learning. In *European Conference on Computer Vision*, 2020.
- [18] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.

- [19] N. Hilliard, L. Phillips, S. Howland, A. Yankov, C. D. Corley, and N. O. Hodas. Few-shot
 learning with metric-agnostic conditional embeddings. *arXiv preprint arXiv:1802.04376*, 2018.
- [20] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [21] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing
 internal covariate shift. In *International Conference on Machine Learning*, 2015.
- [22] Y. Jang, H. Lee, S. J. Hwang, and J. Shin. Learning what and where to transfer. In *International Conference on Machine Learning*, 2019.
- [23] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- J. Krause, M. Stark, J. Deng, and L. Fei-Fei. 3D object representations for fine-grained
 categorization. In *4th International IEEE Workshop on 3D Representation and Recognition* (*3dRR-13*), 2013.
- [25] S. Laine and T. Aila. Temporal ensembling for semi-supervised learning. In *International Conference on Learning Representations*, 2017.
- [26] J. Lee, J. Tack, N. Lee, and J. Shin. Meta-learning sparse implicit neural representations. In
 Advances in Neural Information Processing Systems, 2021.
- [27] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems*, 2018.
- [28] Z. Li and D. Hoiem. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [29] Z. Li, F. Zhou, F. Chen, and H. Li. Meta-sgd: Learning to learn quickly for few-shot learning.
 arXiv preprint arXiv:1707.09835, 2017.
- [30] S. Lu, H.-J. Ye, L. Gan, and D.-C. Zhan. Towards enabling meta-learning from target models.
 In Advances in Neural Information Processing Systems, 2021.
- [31] J. Oh, H. Yoo, C. Kim, and S.-Y. Yun. Boil: Towards representation change for few-shot learning.
 In *International Conference on Learning Representations*, 2021.
- [32] B. Oreshkin, P. Rodríguez López, and A. Lacoste. Tadam: Task dependent adaptive metric for
 improved few-shot learning. In *Advances in Neural Information Processing Systems*, 2018.
- [33] W. Park, D. Kim, Y. Lu, and M. Cho. Relational knowledge distillation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [34] A. Raghu, M. Raghu, S. Bengio, and O. Vinyals. Rapid learning or feature reuse? towards under standing the effectiveness of maml. In *International Conference on Learning Representations*,
 2020.
- [35] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations*, 2017.
- [36] M. Ren, E. Triantafillou, S. Ravi, J. Snell, K. Swersky, J. B. Tenenbaum, H. Larochelle, and R. S.
 Zemel. Meta-learning for semi-supervised few-shot classification. In *International Conference* on *Learning Representations*, 2018.
- [37] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization.
 In *International Conference on Machine Learning*, 2015.
- [38] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous
 control using generalized advantage estimation. In *International Conference on Learning Representations*, 2016.

- [39] R. Schwartz, J. Dodge, N. A. Smith, and O. Etzioni. Green ai. *arXiv preprint arXiv:1907.10597*, 2019.
- [40] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. Cnn features off-the-shelf: an
 astounding baseline for recognition. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2014.
- ⁴⁰⁹ [41] J. Shin, H. B. Lee, B. Gong, and S. J. Hwang. Large-scale meta-learning with continual ⁴¹⁰ trajectory shifting. In *International Conference on Machine Learning*, 2021.
- [42] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, 2017.
- [43] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. A. Raffel, E. D. Cubuk, A. Kurakin,
 and C.-L. Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence.
 In Advances in Neural Information Processing Systems, 2020.
- [44] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple
 way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*,
 2014.
- 419 [45] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT Press, 2018.
- [46] J. Tang, R. Shivanna, Z. Zhao, D. Lin, A. Singh, E. H. Chi, and S. Jain. Understanding and
 improving knowledge distillation. *arXiv preprint arXiv:2002.03532*, 2020.
- [47] A. Tarvainen and H. Valpola. Mean teachers are better role models: Weight-averaged consistency
 targets improve semi-supervised deep learning results. In *Advances in Neural Information Processing Systems*, 2017.
- 425 [48] S. Thrun and L. Pratt. *Learning to Learn*. Springer, 1998.
- [49] Y. Tian, D. Krishnan, and P. Isola. Contrastive representation distillation. In *International Conference on Learning Representations*, 2020.
- [50] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2012.
- [51] H.-Y. Tseng, H.-Y. Lee, J.-B. Huang, and M.-H. Yang. Cross-domain few-shot classification via
 learned feature-wise transformation. In *International Conference on Learning Representations*, 2020.
- [52] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning.
 In Advances in Neural Information Processing Systems, 2016.
- I. Von Oswald, D. Zhao, S. Kobayashi, S. Schug, M. Caccia, N. Zucchet, and J. Sacramento.
 Learning where to learn: Gradient sparsity in meta and continual learning. In *Advances in Neural Information Processing Systems*, 2021.
- [54] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011
 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- [55] Y.-X. Wang and M. Hebert. Learning to learn: Model regression networks for easy small sample
 learning. In *European Conference on Computer Vision*, 2016.
- [56] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
 learning. *Machine learning*, 1992.
- 444 [57] Q. Xie, M.-T. Luong, E. Hovy, and Q. V. Le. Self-training with noisy student improves imagenet 445 classification. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [58] H. Yao, L.-K. Huang, L. Zhang, Y. Wei, L. Tian, J. Zou, J. Huang, et al. Improving generalization in meta-learning via task augmentation. In *International Conference on Machine Learning*, 2021.

- [59] H. Ye, L. Ming, D. Zhan, and W. Chao. Few-shot learning with a strong teacher. *CoRR*, abs/2107.00197, 2021.
- [60] M. Yin, G. Tucker, M. Zhou, S. Levine, and C. Finn. Meta-learning without memorization. In *International Conference on Learning Representations*, 2020.
- [61] L. Yuan, F. E. Tay, G. Li, T. Wang, and J. Feng. Revisiting knowledge distillation via label
 smoothing regularization. In *IEEE Conference on Computer Vision and Pattern Recognition*,
 2020.
- [62] S. Yun, J. Park, K. Lee, and J. Shin. Regularizing class-wise predictions via self-knowledge
 distillation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.

458 Checklist

459	1. For all authors
460 461	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
462	(b) Did you describe the limitations of your work? [Yes] See Section 6.
463 464	(c) Did you discuss any potential negative societal impacts of your work? [Yes] See Section 6.
465 466	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
467	2. If you are including theoretical results
468	(a) Did you state the full set of assumptions of all theoretical results? [N/A]
469	(b) Did you include complete proofs of all theoretical results? [N/A]
470	3. If you ran experiments
471	(a) Did you include the code, data, and instructions needed to reproduce the main experi-
472	mental results (either in the supplemental material or as a URL)? [Yes] The code and
473	instructions are in the supplementary material.
474	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
475	(a) Did you report array here (a g, with respect to the renders coad ofter running details.
476 477	ments multiple times)? [Yes] See Section 5.
478 479	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Appendix for details.
480	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
481	(a) If your work uses existing assets, did you cite the creators? [Yes]
482	(b) Did you mention the license of the assets? [Yes]
483	(c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
484	
485	(d) Did you discuss whether and how consent was obtained from people whose data you're
486	using/curating? [N/A]
487	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? $[N/A]$
400	5. If you used anoudcoursing or conducted response with human subjects
489	5. If you used crowdsourcing or conducted research with numan subjects
490 491	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
492 493	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
494 495	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]