# Not All Tasks are Equal - Task Attended Meta-learning for Few-shot Learning

**Anonymous authors**
**Paper under double-blind review**

## Abstract

Meta-learning (ML) has emerged as a promising direction in learning models under constrained resource settings like few-shot learning. The popular approaches for ML either learn a generalizable initial model or a generic parametric optimizer through batch episodic training. In this work, we study the importance of tasks in a batch for ML. We hypothesize that the common assumption in batch episodic training where each task in a batch has an equal contribution to learning an optimal meta-model need not be true. We propose to weight the tasks in a batch according to their "importance" in improving the meta-model's learning. To this end, we introduce a training curriculum called task attended meta-training to learn a meta-model from weighted tasks in a batch. The task attention module is a standalone unit and can be integrated with any batch episodic training regimen. Comparison of task-attended ML models with their non-task-attended counterparts on complex datasets, performance improvement of proposed curriculum over state-of-the-art task scheduling algorithms on noisy datasets, and cross-domain few shot learning setup validate its effectiveness.

## 1 Introduction

The ability to infer knowledge and discover complex representations from data has made deep learning models widely popular in the machine learning community. However, these models are data-hungry, often requiring large volumes of labeled data for training. Collection and annotation of such large amounts of training data may not be feasible for many real life applications, especially in domains that are inherently data constrained, like medical and satellite image classification, drug toxicity estimation, etc. Meta-learning (ML) has emerged as a promising direction for learning models in such settings, where only a limited amount (few-shots) of labeled training data is available. A typical ML algorithm employs an episodic training regimen that differs from the training procedure of conventional learning tasks. This episodic meta-training regimen is backed by the assumption that a machine learning model quickly generalizes to novel unseen data with minimal fine-tuning when trained and tested under similar circumstances (Vinyals et al., 2016). To facilitate such a generalization capacity, a meta-training phase is undertaken, where the model is trained to optimize its performance on several homogeneous tasks/episodes randomly sampled from a dataset. Each episode or task is a learning problem in itself. In the few-shot setting each task is a classification problem, a collection of $K$ support (train) and $Q$ query (test) samples corresponding to each of the $N$ classes. Task-specific knowledge is learned using the support data, and meta-knowledge across the tasks is learned using query samples, which essentially encodes "how to learn a new task effectively."

The learned meta-knowledge is generic and agnostic to tasks from the same distribution. It is typically characterized in two different forms - either as an optimal initialization for the machine learning model or a learned parametric optimizer. Under the optimal initialization view, the learned meta-knowledge represents an optimal prior over the model parameters, that is equidistant, but close to the optimal parameters for all individual tasks. This enables the model to rapidly adapt to unseen tasks from the same distribution (Finn et al., 2017; Li et al., 2017; Jamal & Qi, 2019). Under the parametric optimizer view, meta-knowledge pertaining to the traversal of the loss surface of tasks is learned by the meta-optimizer. Through learning task specific and task agnostic characteristics of the loss surface, a parametric optimizer can thus effectively

guide the base model to traverse the loss surface and achieve superior performance on unseen tasks from the same distribution (Ravi & Larochelle, 2017).

Initialization based ML approaches accumulate the meta-knowledge by simultaneously optimizing over a batch of tasks. On the other hand, a parametric optimizer sequentially accumulates meta-knowledge across individual tasks. The sequential accumulation process leads to a long oscillatory optimization trajectory and a bias towards the last task, limiting the parametric optimizer's task agnostic potential. However, recently meta-knowledge has been accumulated in a batch mode even for the parametric optimizer (Aimen et al., 2021). Further, under such batch episodic training (for both initialization and optimization views), a common assumption in ML that the randomly sampled episodes of a batch contribute equally to improving the learned meta-knowledge need not hold good. Due to the latent properties of the sampled tasks in a batch and the model configuration, some tasks may be better aligned with the optimal meta-knowledge than others. We hypothesize that proportioning the contribution of a task as per its alignment towards the optimal meta-knowledge can improve the meta-model's learning. This is analogous to classical machine learning algorithms like sample re-weighting, which however, operate at sample granularity. In re-weighting, samples leading to false positives are prioritized and therefore replayed. Hence, the latent properties due to which a sample is prioritized are explicitly defined. For complex task distributions, explicitly handcrafting the notion of "importance" of a task would be hard.

To this end, we propose a task attended meta-training curriculum that employs an attention module that learns to assign weights to the tasks of a batch with experience. The attention module is parametrized as a neural network that takes meta-information in terms of the model's performance on the tasks in a batch as input and learns to associate weights to each of the tasks according to their contribution in improving the meta-model. Overall, we make the following contributions,

- We propose a task attended meta-training strategy wherein different tasks of a batch are weighted according to their "importance" defined by the attention module. This attention module is a standalone unit that can be integrated into any batch episodic training regimen.

- We extend the empirical investigation of the batch-mode parametric optimizer (MetaLSTM++) to complex datasets like miniImagenet, FC100, and tieredImagenet and validate its efficiency over its sequential counter-part (MetaLSTM).

- We conduct extensive experiments on miniImagenet, FC100, and tieredImagenet datasets and compare ML algorithms like MAML, MetaSGD, ANIL, and MetaLSTM++ with their non-task-attended counterparts to validate the effectiveness of the task attention module and its coupling with any batch episodic training regimen.

- We compare task-attended curriculum with state-of-the-art task scheduling approaches and also show the merit of the proposed approach on the miniImagenet-noisy dataset and cross-domain few shot learning (CDFSL) setup.

- We also perform exhaustive empirical analysis and visual inspections to decipher the working of the task attention module.

## 2 Related Work

ML literature is profoundly diverse and may broadly be classified into *initialization* (Finn et al., 2017; Li et al., 2017; Jamal & Qi, 2019; Raghu et al., 2020; Rusu et al., 2019; Sun et al., 2019) and *optimization approaches* (Ravi & Larochelle, 2017) depending on the metaknowledge. However, these approaches assume uniform contribution of tasks in learning a meta-model. In supervised learning, assigning non-uniform priorities to the samples is not new (Kahn & Marshall, 1953; Shrivastava et al., 2016). Self-paced learning (Kumar et al., 2010) and hard example mining (Shrivastava et al., 2016) have popularly been used to reweight the samples and various attributes like losses, gradients, and uncertainty have been used to assign priorities to samples (Lin et al., 2017; Zhao & Zhang, 2015; Chang et al., 2017). Zhao & Zhang (2015) introduce importance sampling to reduce variance and improve the convergence rate of stochastic optimization algorithms over

uniform sampling. They theoretically prove that the reduction in the variance is possible if the sampling distribution depends on the norm of the gradients of the loss function. Chang et al. (2017) conclude that mini-batch SGD for classification is improved by emphasizing the uncertain examples. Lin et al. (2017) propose reshaped cross-entropy loss (focal loss) that down-weights the loss of confidently classified samples. Nevertheless, assigning non-uniform priorities to tasks in meta-learning is under-explored and has recently drawn attention (Kaddour et al., 2020; Gutierrez & Leonetti, 2020; Liu et al., 2020; Yao et al., 2021; Arnold et al., 2021). Gutierrez & Leonetti (2020) propose Information-Theoretic Task Selection (ITTS) algorithm to filter training tasks that are distinct from each other and close to the tasks of the target distribution. This algorithm results in a smaller pool of training tasks. A model trained on the smaller subset learns better than the one trained on the original set. On the other hand, Kaddour et al. (2020) propose probabilistic active meta-learning (PAML) that learns probabilistic task embeddings. Scores are assigned to these embeddings to select the next task presented to the model. These algorithms are, however, specific to meta-reinforcement learning (meta-RL). On the contrary, our focus is on the few shot classification problem. Liu et al. (2020) propose a greedy class-pair potential-based adaptive task sampling strategy wherein task selection depends on the difficulty of all class-pairs in a task. This sampling technique is static and operates at a class granularity. On the other hand, our approach is dynamic and operates at a task granularity. Assigning non-uniform weights to samples prevents overfitting on corrupt data points (Ren et al., 2018b; Jiang et al., 2018). Ren et al. (2018b) used gradient directions to re-weight the data points, and Jiang et al. (2018) learned a curriculum on examples using a mentor network. However, these approaches assume availability of abundant labeled data. Yao et al. (2021) extended (Jiang et al., 2018) to few-shot learning setup. They propose a neural schedular to predict the sampling probability of tasks in a candidate pool. Parallel to (Jiang et al., 2018), they consider noisy and imbalanced task distributions. Our work is different from these approaches as we do not propose a task sampling strategy but a dynamic task-batch re-weighting mechanism for the meta-model update in a few-shot learning setup. Also, (Yao et al., 2021) is more expensive than the proposed approach as it performs an additional warm start to the scheduler, utilizes more task batches in a run, and uses REINFORCE for reward estimation. Arnold et al. (2021) hypothesize and empirically validate that task difficulty approximately follows a normal distribution. They find the sampling uniformly over episode difficulty outperforms other sampling schemes like curriculum, easy and hard-mining. Our approach differs from Uniform Sampling as we do not explicitly handicraft the notion of task difficulty and do not assume the normal distribution over task difficulty. Instead, we let an attention network learn the suitable weights for the tasks in a batch. Contrary to our idea is TAML (Jamal & Qi, 2019) - a meta-training curriculum that enforces equity across the tasks in a batch. We show that weighting the tasks according to their "importance" and hence utilizing the diversity present in a batch given the meta-model's current configuration offers better performance than enforcing equity in a batch of tasks.

## 3 Preliminary

In a typical ML setting, the principal dataset $\mathcal{D}$ is divided into disjoint meta-sets $\mathcal{M}$ (meta-train set), $\mathcal{M}_v$ (meta-validation set) and $\mathcal{M}_t$ (meta-test set) for training the model, tuning its hyperparameters and evaluating its performance, respectively. Every meta-set is a collection of tasks $\mathcal{T}$ drawn from the joint task distribution $P(\mathcal{T})$ where each task $\mathcal{T}_i$ consists of support set $D_i = \{(x_k^c, y_k^c)_{k=1}^K\}_{c=1}^N$ and query set $D_i^* = \{(x_q^{*c}, y_q^{*c})_{q=1}^Q\}_{c=1}^N$. Here $(x, y)$ represents a (sample, label) pair and $N$ is the number of classes, $K$ and $Q$ are the number of samples belonging to each class in the support and query set, respectively. According to support-query characterization $\mathcal{M}$, $\mathcal{M}_v$ and $\mathcal{M}_t$ could be represented as $\{(D_i, D_i^*)\}_{i=1}^M$, $\{(D_i, D_i^*)\}_{i=1}^R$, $\{(D_i, D_i^*)\}_{i=1}^S$ where $M$, $R$ and $S$ are the total number of tasks in $\mathcal{M}$, $\mathcal{M}_v$ and $\mathcal{M}_t$ respectively. During meta-training on $\mathcal{M}$, meta-model $\theta$ is adapted on $D_i$ of each $\mathcal{T}_i$ to $\phi_i$. The adapted model $\phi_i$ is then evaluated on $D_i^*$ to update $\theta$. The output of this episodic training is either an optimal prior or a parametric optimizer, both aiming to facilitate the rapid adaptation of the model on unseen tasks from $\mathcal{M}_t$.

### 3.1 Meta-knowledge as an Optimal Initialization

When meta-knowledge is a generic initialization on the model parameters learned through the experience over various tasks, it is enforced to be close to each individual training tasks' optimal parameters. A model initialized with such an optimal prior quickly adapts to unseen tasks from the same distribution during

¹³⁶ meta-testing. **MAML** (Finn et al., 2017) employs a nested iterative process to learn the task-agnostic
¹³⁷ optimal prior $\theta$. In the inner iterations representing the task adaptation steps, $\theta$ is separately fine-tuned for
¹³⁸ each meta-training task $\mathcal{T}_i$ of a batch using $D_i$ to obtain $\phi_i$ through gradient descent on the train loss $L$
¹³⁹ using learning rate $\alpha$. Specifically, $\phi_i$ is initialized as $\theta$ and updated using $\phi_i \leftarrow \phi_i - \alpha\nabla_{\phi_i}L(\phi_i)$, $T$ times
¹⁴⁰ resulting in the adapted model $\phi_i^T$. In the outer loop, meta-knowledge is gathered by optimizing $\theta$ over
¹⁴¹ loss $L^*$ computed with the task adapted model parameters $\phi_i^T$ on query dataset $D_i^*$. Specifically, during
¹⁴² meta-optimization $\theta \leftarrow \theta - \beta\nabla_\theta \sum_{i=1}^{B} L^*(\phi_i^T)$ using a task batch of size $B$ and learning rate $\beta$. **MetaSGD**
¹⁴³ (Li et al., 2017) improves upon MAML by learning parameter-specific learning rates $\boldsymbol{\alpha}$ in addition to the
¹⁴⁴ optimal initialization in a similar nested iterative procedure. Meta-knowledge is gathered by optimizing $\theta$
¹⁴⁵ and $\boldsymbol{\alpha}$ in the outer loop using the loss $L^*$ computed on query set $D_i^*$. Specifically, during meta-optimization
¹⁴⁶ $(\theta, \boldsymbol{\alpha}) \leftarrow (\theta, \boldsymbol{\alpha}) - \beta\nabla_{(\theta,\boldsymbol{\alpha})} \sum_{i=1}^{B} L^*(\phi_i^T)$. Learning dynamic learning rates for each parameter of a model
¹⁴⁷ makes MetaSGD faster and more generalizable than MAML. A single adaptation step is sufficient to adjust
¹⁴⁸ the model towards a new task. The performance of MAML is attributed to the reuse of the features
¹⁴⁹ across tasks rather than the rapid learning of new tasks (Raghu et al., 2020). Exploiting this characteristic,
¹⁵⁰ **ANIL** freezes the feature backbone layers $(1, \ldots, l-1)$ and only adapts classifier layer $(l)$ in the inner
¹⁵¹ loop $T$ times. Specifically during adaptation $\phi_i^l \leftarrow \phi_i^l - \alpha\nabla_{\phi_i^l}L(\phi_i^l)$. During meta-optimization $\theta^{1,\ldots,l} \leftarrow$
¹⁵² $\theta^{1,\ldots,l} - \beta\nabla_{\theta^{1,\ldots,l}} \sum_{i=1}^{B} L^*(\phi_i^{lT})$ i.e., all layers are learned in the outer loop. Freezing the feature backbone
¹⁵³ during adaptation reduces the overhead of computing gradient through the gradient (differentiating through
¹⁵⁴ the inner loop), and thereby heavier backbones could be used for the feature extraction. **TAML** (Jamal
¹⁵⁵ & Qi, 2019) suggests that the optimal prior learned by MAML may still be biased towards some tasks.
¹⁵⁶ They propose to reduce this bias and enforce equity among the tasks by explicitly minimizing the inequality
¹⁵⁷ among the performances of tasks in a batch. The inequality defined using statistical measures such as Theil
¹⁵⁸ Index, Atkinson Index, Generalized Entropy Index, and Gini Coefficient among the performances of tasks
¹⁵⁹ in a batch is used as a regularizer while gathering the meta-knowledge. For the baseline comparison, in
¹⁶⁰ our experiments, we use the Theil index for TAML owing to its average best results. Specifically during
¹⁶¹ meta-optimization $\theta \leftarrow \theta - \beta\nabla_\theta \left[ \sum_{i=1}^{B} L^*(\phi_i^T) + \lambda \left\{ \frac{L^*(\phi_i^0)}{\bar{L}^*(\phi_i^0)} \ln \frac{L^*(\phi_i^0)}{\bar{L}^*(\phi_i^0)} \right\} \right]$ (for TAML-Theil Index) where $B$
¹⁶² is the number of tasks in a batch, $L^*(\phi_i^0)$ is the loss incurred by initial model $\phi_i^0$ on the query set $D_i^*$ of
¹⁶³ task $\mathcal{T}_i$ and $\bar{L}^*(\phi_i^0)$ is the average query loss of initial model on a batch of tasks. As TAML enforces equity
¹⁶⁴ of the optimal prior towards meta-train tasks, it counters the adaptation, which leads to slow and unstable
¹⁶⁵ training largely dependent on $\lambda$.

## 3.2 Meta-knowledge as a Parametric Optimizer

¹⁶⁷ A regulated gradient-based optimizer gathers the task-specific and task-agnostic meta-knowledge to traverse
¹⁶⁸ the loss surfaces of tasks in the meta-train set during meta-training. A base model guided by such a
¹⁶⁹ learned parametric optimizer quickly finds the way to minima even for unseen tasks sampled from the
¹⁷⁰ same distribution during meta-testing. **MetaLSTM** (Ravi & Larochelle, 2017) is a recurrent parametric
¹⁷¹ optimizer $\theta$ that mimics the gradient-based optimization of a base model $\phi$. This recurrent optimizer is an
¹⁷² LSTM (Hochreiter & Schmidhuber, 1997) and is inherently capable of performing two-level learning due to its
¹⁷³ architecture. During adaptation of $\phi_i$ on $D_i$, $\theta$ takes meta information of $\phi_i$ characterized by its current loss
¹⁷⁴ $L$ and gradients $\nabla_{\phi_i}(L)$ as input and outputs the next set of parameters for $\phi_i$. This adaptation procedure
¹⁷⁵ is repeated $T$ times resulting in the adapted base-model $\phi_i^T$. Internally, the cell state of $\theta$ corresponds to $\phi_i$,
¹⁷⁶ and the cell state update for $\theta$ resembles a learned and controlled gradient update. The emphasis on previous
¹⁷⁷ parameters and the current update is regulated by the learned forget and input gates respectively. While
¹⁷⁸ adapting $\phi_i$ to $D_i$, information about the trajectory on the loss surface across the adaptation steps is captured
¹⁷⁹ in the hidden states of $\theta$, representing the task-specific knowledge. During meta-optimization, $\theta$ is updated
¹⁸⁰ based on the loss of the adapted model $L^*(\phi_i^T)$ computed on the query set $D_i^*$ to garner the meta-knowledge
¹⁸¹ across tasks. Specifically, during meta-optimization, $\theta \leftarrow \theta - \beta\nabla_\theta L^*(\phi_i^T)$. MetaLSTM updates parametric
¹⁸² optimizer $\theta$ after adapting the base model $\phi$ to each task. This causes $\theta$ to follow optima's of all adapted
¹⁸³ base models leading to its elongated and fluctuating optimization trajectory, which is biased towards the last
¹⁸⁴ task. **MetaLSTM++** (Aimen et al., 2021) circumvents these issues as $\theta$ is updated by an aggregate query
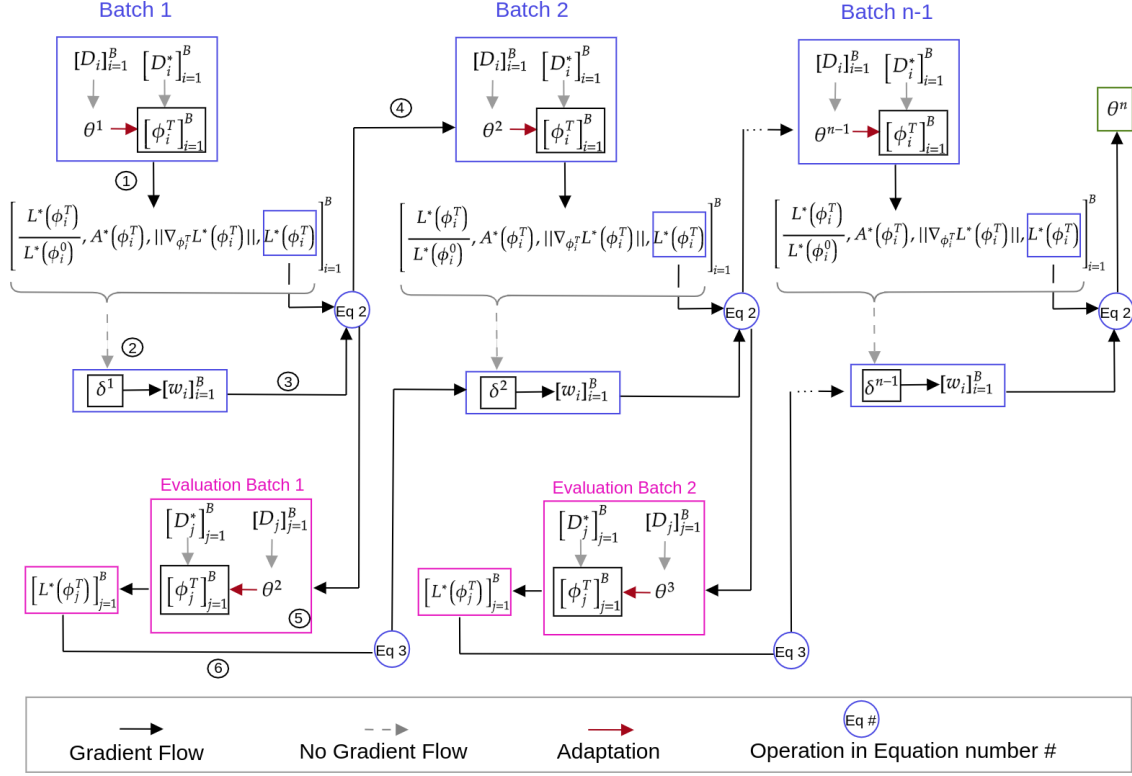
Figure 1: Computational Graph of the forward pass of the meta-model using task attended meta-training curriculum. The output of this procedure is a meta-model $\theta^n$. Gradients are propagated through solid lines and restricted through dashed lines.

loss of the adapted models on a batch of tasks. Batch updates smoothen the optimization trajectory of $\theta$ and eliminate its bias towards the last task. Specifically, during meta-optimization $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{i=1}^{B} L^*(\phi_i^T)$.

## 4 Task Attention in Meta-learning

A common assumption under the batch-wise episodic training regimen adopted by ML is that each task in a batch has an equal contribution in improving the learned meta-knowledge. However, this need not always be true. It is likely that given the current configuration of the meta-model, some tasks may be more important for the meta-model's learning. A contributing factor to this difference is that tasks sampled from complex data distributions can be profoundly diverse. The diversity and latent properties of the tasks coupled with the model configuration may induce some tasks to be better aligned with the optimal meta-knowledge than others. The challenging aspect in the meta-learning setting is to define the "importance" and associate weights to the tasks of a batch proportional to their contribution to improving the meta-knowledge. As human beings, we *learn* to associate importance to events subjective to meta-information about the events and prior experience. This motivates us to define a learnable module that can map the meta-information of tasks to their importance weights.

### 4.1 Characteristics of Meta-Information

Given a task-batch $\{\mathcal{T}_i\}_{i=1}^{B}$, the task attention module takes as input meta-information about each task $(\mathcal{T}_i)$ in the batch, defined as the four tuple below:

$$\mathcal{I} = \left\{ \left( ||\nabla_{\phi_i^T} L^*(\phi_i^T)||, L^*(\phi_i^T), A^*(\phi_i^T), \frac{L^*(\phi_i^T)}{L^*(\phi_i^0)} \right) \right\}_{i=1}^{B} \tag{1}$$

where corresponding to each task $i$ in the batch $||\nabla_{\phi_i^T} L^*(\phi_i^T)||$ denotes the norm of gradient, $L^*(\phi_i^T)$ and $A^*(\phi_i^T)$ are the test loss and accuracy of the adapted model respectively, and $\frac{L^*(\phi_i^T)}{L^*(\phi_i^0)}$ is the ratio of the model's test loss post and prior adaptation.

### 4.1.1 Gradient Norm

Let $P = \left\{\phi_i^T\right\}_{i=1}^B$ be the parameters of the models obtained after adapting the initial model (for $T$ iterations) on the support data $\{D_i\}_{i=1}^B$ of tasks $\{\mathcal{T}_i\}_{i=1}^B$. Also, let $G = \left\{\nabla_{\phi_i^T} L^*(\phi_i^T)\right\}_{i=1}^B$ be the gradients of the adapted model parameters w.r.t the query losses $\{L^*(\phi_i^T)\}_{i=1}^B$. The gradient norm $\left\{||\nabla_{\phi_i^T} L^*(\phi_i^T)||\right\}_{i=1}^B$ is the $L_2$ norm of the gradients and quantifies the magnitude of the consolidated displacement of the adapted model parameters during a gradient descent update on query data. Larger gradient norm on query dataset could indicate that the model has either not learned the support set or has overfitted. Hence the model is not generalizable on query set compared to the models with low gradient norm. Gradient norm, therefore, carries information about the convergence and generalizability of the adapted models which has been theoretically studied in (Li et al., 2019).

### 4.1.2 Test Loss

$\{L^*(\phi_i^T)\}_{i=1}^B$ represents the empirical error (cross entropy loss) of the adapted base models on unseen query instances and hence characterizes their generalizability. Unlike gradient norm, which characterizes the generalizability in parameter space, query loss quantifies generalizability in the output space as the divergence between the real and predicted probability distributions. As $\{L^*(\phi_i^T)\}_{i=1}^B$ is a key component in the meta-update equation, it is an important factor influencing the meta-model's learning. Further, test errors of classes have been widely used to determine their "easy or hardness" (Bengio et al., 2009; Liu et al., 2021; Arnold et al., 2021). Thus $\{L^*(\phi_i^T)\}_{i=1}^B$ acquaints the attention module with the generalizability aspect of task models and their influence in updating the meta-model.

### 4.1.3 Test Accuracy

$\{A^*(\phi_i^T)\}_{i=1}^B$ corresponds to the accuracies of $\{\phi_i^T\}_{i=1}^B$ on $\{D_i^*\}_{i=1}^B$ scaled in the range [0,1]. $A^*(\phi_i^T)$ evaluates the thresholded predictions (predicted labels) unlike $L^*(\phi_i^T)$, which evaluates the confidence of the model's predictions on the true class labels. Two task models may predict the same class labels but differ in the confidence of the predictions. In such scenarios, neither loss nor accuracy is individually sufficient to comprehend this relationship among the tasks. So, the combination of these two entities is more reflective of the nature of the learned task models.

### 4.1.4 Loss-ratio

Let $L^*(\phi_i^0)$ be the loss of $\theta$ on the $D_i^*$, and $L^*(\phi_i^T)$ be the loss of the adapted model $\phi_i^T$ on $D_i^*$. The loss-ratio $\frac{L^*(\phi_i^T)}{L^*(\phi_i^0)}$ is representative of the relative progress of a meta-model on each task. Higher values ($> 1$) of the loss-ratio suggests adapting $\theta$ to $D_i$ has an adverse effect on generalizing it to $D_i^*$ (negative impact), while lower values ($< 1$) of the loss-ratio indicates the benefit of adaptation of $\theta$ on $D_i$ (positive impact). Loss-ratio of exactly one signifies adaptation attributes to no additional benefit (neutral impact). Therefore, loss-ratio provides information regarding the impact of adaptation on each task for a given meta-model.

## 4.2 Task Attention Module

We learn a task attention module parameterized by $\delta$, which attends to the tasks that contribute more to the model's learning i.e., the objective of the task attention module is to learn the relative importance of each task in the batch for the meta-model's learning. Thus the output of the module is a $B-$dimensional vector $\mathbf{w} = [w_1, \ldots, w_B]$, $(\sum_{i=1}^B w_i = 1 \text{ and } \forall \mathcal{T}_i, \ w_i \geq 0)$ quantifying the attention-score (weight - $w_i$) for each task.

**Algorithm 1:** Task Attended Meta-Training

---

**Input:**

*Dataset:* $\mathcal{M} = \{D_i, D_i^*\}_{i=1}^M$

*Models:* Meta-model $\theta$, Base-model $\phi$, Att-module $\delta$

*Learning-rates:* $\alpha$, $\beta$, $\gamma$

*Parameters:* Iterations $n_{iter}$, Batch-size $B$,
                     Adaptation-steps $T$

**Output:** Meta-model $\theta$

**1 Initialization:** $\theta, \delta \leftarrow$ Random Initialization

**2 for** *iteration in $n_{iter}$* **do**

**3**    $\{\mathcal{T}_i\}_{i=1}^B = \{D_i, D_i^*\}_{i=1}^B \leftarrow$ Sample task-batch($\mathcal{M}$)

**4**    **for** *all $\mathcal{T}_i$* **do**

**5**      $\phi_i^0 \leftarrow \theta$

**6**      $L^*(\phi_i^0), \_\_ \leftarrow evaluate(\phi_i^0, D_i^*)$   $\triangleright$ Compute loss
         and accuracy of input model on given dataset.

**7**      $\phi_i^T = adapt(\phi_i^0, D_i)$

**8**      $L^*(\phi_i^T), A^*(\phi_i^T) \leftarrow evaluate(\phi_i^T, D_i^*)$

**9**    **end**

**10**   $[w_i]_{i=1}^B \leftarrow$ Att_module

       $\left( \left[ \dfrac{L^*(\phi_i^T)}{L^*(\phi_i^0)}, A^*(\phi_i^T), ||\nabla_{\phi_i^T} L^*(\phi_i^T)||, L^*(\phi_i^T) \right]_{i=1}^B \right)$

**11**   $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{i=1}^B w_i L^*(\phi_i^T)$

**12**   $\{D_j, D_j^*\}_{j=1}^B \leftarrow$ Sample task-batch($\mathcal{M}$)

**13**   **for** *all $\mathcal{T}_j$* **do**

**14**      $\phi_j^0 \leftarrow \theta$

**15**      $\phi_j^T = adapt(\phi_j^0, D_j)$

**16**   **end**

**17**   $\delta \leftarrow \delta - \gamma \nabla_\delta \sum_{j=1}^B L^*(\phi_j^T)$

**18 end**

**19 Return** $\theta$

**20 Function** `adapt`($\phi_i^t, D_i$)**:**

**21**   $\theta \leftarrow \phi_i^t$

**22**   **if** *$\theta$ is optimal-initialization* **then**

**23**      **for** *t=1 to T* **do**

**24**        $\phi_i^{t+1} \leftarrow \phi_i^t - \alpha \nabla_{\phi_i^t} L(\phi_i^t)$

**25**      **end**

**26**   **end**

**27**   **else if** *$\theta$ is parametric-optimizer* **then**

**28**      **for** *t=1 to T* **do**

**29**        $\phi_i^{t+1} \leftarrow \theta \left( L(\phi_i^t), \nabla_{\phi_i^t} L(\phi_i^t) \right)$   $\triangleright$ Parameter
         updates given by cell state of $\theta$.

**30**      **end**

**31**   **end**

**32**   ***Return*** $\phi_i^T$

The attention vector $\mathbf{w}$ is multiplied with the corresponding task losses of the adapted models $L^*(\phi_i^T)$ on the held-out datasets $D_i^*$ to update the meta-model $\theta$:

$$\theta^{t+1} \leftarrow \theta^t - \beta \nabla_{\theta^t} \sum_{i=1}^B w_i L^*(\phi_i^T) \qquad (2)$$

After the meta-model is updated using the weighted task losses, we evaluate the goodness of the generated attention weights. We sample a new batch of tasks $\{D_j, D_j^*\}_{j=1}^B$ and adapt a base-model $\phi_j$ using the updated meta-model $\theta^{t+1}$ on the train data $\{D_j\}$ of each task. The mean test-loss of the adapted models $\{\phi_j^T\}_{j=1}^B$ reflect the goodness of the weights assigned by the attention-module in the previous iteration. The attention module $\delta$ is thus updated using the gradients flowing back into it w.r.t to this mean test-loss. The attention network is trained simultaneously with the meta-model in an end to end fashion using the update rule:

$$\delta^{t+1} \leftarrow \delta^t - \gamma \nabla_{\delta^t} \sum_{j=1}^B L^*(\phi_j^T) \qquad (3)$$

where $\phi_j^T$ is adapted from $\theta^{t+1}$ and $\gamma$ is the learning rate .

### 4.3 Task Attended Meta-Training Algorithm

We demonstrate the meta-training curriculum using the proposed task attention in Figure 1 and formally summarize it in Algorithm 1. As with the classical meta-training process, we first sample a batch of tasks from the task distribution. For each task $\mathcal{T}_i$, we adapt the base-model $\phi_i$ using the train data $D_i$ for $T$ time-steps (line 7 and lines 20-32 in Algorithm 1). Specifically, for initialization approaches, adaptation is performed by gradient descent on train loss $L$ (lines 22-26 in Algorithm 1). However, for optimization approaches, current loss and gradients are inputted to the meta-model $\theta$, which outputs the updated base-model parameters (lines 27-31 in Algorithm 1). Then we compute the meta-information about the adapted model corresponding to each task. It comprises of the loss $L^*(\phi_i^T)$, accuracy $A^*(\phi_i^T)$, loss-ratio $\dfrac{L^*(\phi_i^T)}{L^*(\phi_i^0)}$ and gradient norm $||\nabla_{\phi_i^T} L^*(\phi_i^T)||$ on the test data $D_i^*$. This meta-information corresponding to each task in a batch is given as input to the task attention module (Figure 1 - Label: ②) which outputs the attention vector (line 10 in Algorithm 1). The attention vector and test losses $\{L^*(\phi_i^T)\}_{i=1}^B$ are used to update meta-model parameters $\theta$ according to equation 2 (line 11 in Algorithm 1, Figure 1 - Label: ④). We sample a new batch of tasks $\{D_j, D_j^*\}_{j=1}^B$ and adapt the base-models $\{\phi_j^T\}_{j=1}^B$ using the updated meta-model (lines 12-16 in Algorithm

1, Figure 1 - Label: ⑤). We compute the mean test loss over the adapted base-models $\{L^*(\phi_j^T)\}_{j=1}^B$, which is then used to update the parameters of the task attention module $\delta$ according to equation 3 (line 17 in Algorithm 1, Figure 1 - Label: ⑥).

The attention network is designed as a stand-alone module to learn the mapping from the meta-information space to the importance of tasks in a batch. The meta-model is learned according to equation 2 and aims to minimize the weighted loss. It is important to decouple the learning of the attention network from that of the meta-model. If there is information flow from the task attention module to the meta-model, the latter may reduce its weighted loss by learning an initialization that is suboptimal, but for which the task attention network assigns lower weights. This would introduce an undesirable bias to the learning process. To circumvent this bias, we restrict the flow of gradients to the meta-model $\theta$ through the task attention module $\delta$ by enforcing $\nabla_\theta w_i L^*(\phi_i^T) = w_i \nabla_\theta L^*(\phi_i^T)$ i.e., $\nabla_\theta w_i$ is not computed. Also, gradients flowing through the attention network to the meta-model create additional computational overhead. Specifically, the term $\nabla_\theta \sum_i w_i L^*(\phi_i^T)$ from equation 2 can be expanded as follows -

$$\nabla_\theta \sum_i w_i L^*(\phi_i^T) = \sum_i \nabla_\theta w_i L^*(\phi_i^T) = \underbrace{\sum_i w_i \nabla_\theta L^*(\phi_i^T)}_{\text{Term 1}} + \underbrace{\sum_i L^*(\phi_i^T) \nabla_\theta w_i}_{\text{Term 2}}$$

The $\nabla_\theta w_i$ in Term 2 is computationally expensive as $\nabla_\theta w_i = \nabla_\delta w_i . \nabla_I \delta . \nabla_\phi I . \nabla_\theta \phi$. Restricting the gradient flow avoids these additional computations. We also note that the meta-model and attention network are updated only once during each training iteration, although on different batches of tasks.

## 5 Experiments and Results

We consider different few-shot learning settings on the benchmark datasets - miniImagenet, miniImagenet-noisy, Fewshot Cifar 100 (FC100) and tieredImagenet to test the effectiveness of the proposed attention module. All the experimental results and comparisons correspond to our re-implementation of the ML algorithms integrated into learn2learn library (Arnold et al., 2020) to ensure fairness and uniformity. We believe that integrating the proposed attention module and additional ML algorithms into the learn2learn library will benefit the ML community. We perform individual hyperparameter tuning for all the models over the same hyperparameter space to ensure a fair comparison. The source code is publicly available.[1]



Figure 2: Architecture of Task-attention module.

### 5.1 Dataset and Implementation Details

In line with the state-of-the-art literature (Sun et al., 2020; Arnold et al., 2021), we use miniImagenet, FC100, and tieredImagenet for evaluating the effectiveness of the proposed attention module as they are more challenging datasets comprising of highly diverse tasks. We also test the efficacy of the proposed approach on noisy datasets like miniImagenet-noisy, and for CDFSL, we use miniImagenet → CUB-200 and miniImagenet → FGVC-Aircrafts datsets. The details of the datasets are presented in the supplementary material.

We use a 4-layer CNN from (Finn et al., 2017) as a base model and a two-layer LSTM (Ravi & Larochelle, 2017) for the parametric optimizer. The architecture of the task-attention module is illustrated in Figure 2 and described as follows. The task attention module is implemented as a 4-layer neural network. The first layer performs a 1×1 convolution over the input (meta-information) of size B×4 where B denotes the meta-batch size, producing a vector of size B×1 as output. This vector is then passed through two fully connected layers with 32 hidden nodes, each followed by a ReLU activation. This output is then passed through a fully connected layer with B nodes, followed by a softmax activation to produce the normalized attention weights.
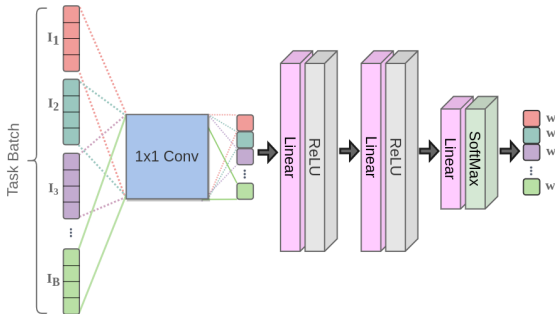
[1]https://github.com/taskattention/task-attended-metalearning.git

We perform a grid search over 30 different configurations for 5000 iterations to find the optimal hyper-parameters for each setting. The search space is shared across all meta-training algorithms and datasets. The meta, base and attention model learning rates are sampled from a log uniform distribution in the ranges $[1e^{-4}, 1e^{-2}]$, $[1e^{-2}, 5e^{-1}]$ and $[1e^{-4}, 1e^{-2}]$ respectively (see appendix for more details). The hyperparameter $\lambda$ for TAML (Theil) is sampled from a log uniform distribution over the range of $[1e^{-2}, 1]$. The number of adaptation steps is fixed to 5 for all settings except for 10-way 5-shot setting, where we use 2 adaptation steps owing to the computational expenses. The meta-batch size is set to 4 for all settings (Finn et al., 2017; Jamal & Qi, 2019). However, we study its impact in Table 1. All models were trained for 55000 iterations (early stopping was employed for tieredImagenet) using the optimal set of hyper-parameters using an Adam optimizer (Kingma & Ba, 2015).

Table 1: Comparison of few-shot classification performance of MAML and TA-MAML on miniImagenet dataset with meta-batch size 4 and 6 and 8 for 5 and 10-way (1 and 5-shot) settings. The ± represents the 95% confidence intervals over 300 tasks. Algorithms denoted by * are rerun on their optimal hyper-parameters. We observe that TA-MAML consistently performs better than MAML, and an increase in the tasks in a batch improves the performance of both MAML and TA-MAML.

| | Test Accuracy (%) on miniImagenet | | | |
| | 5-Way | | 10-Way | |
| Model | 1 Shot | 5 Shot | 1 Shot | 5 Shot |
|---|---|---|---|---|
| Batch Size 4 | | | | |
| MAML* | 46.10 ± 0.19 | 60.16 ± 0.17 | 29.42 ± 0.11 | 41.98 ± 0.10 |
| **TA-MAML** | **48.36 ± 0.23** | **62.48 ± 0.18** | **31.15± 0.11** | **43.70 ± 0.09** |
| Batch Size 6 | | | | |
| MAML* | 47.72 ± 1.041 | 63.45 ± 1.083 | 31.55 ± 0.626 | 46.27 ± 0.64 |
| **TA-MAML** | **49.14 ± 1.211** | **65.26 ± 0.956** | **32.62± 0.635** | **46.67 ± 0.63** |
| Batch Size 8 | | | | |
| MAML* | 47.68±1.20 | 63.81±0.98 | 31.54±0.66 | 46.15±0.58 |
| **TA-MAML** | **50.35±1.22** | **65.69±1.08** | **32.00±0.68** | **48.33±0.63** |

## 5.2 Influence of Task Attention on Meta-Training

As task-attention (TA) is a standalone module, it can be integrated with any batch episodic training regimen. We, therefore, use MetaLSTM++ (batch mode of MetaLSTM) for our experiments. In (Aimen et al., 2021), authors demonstrated the merit of MetaLSTM++ on MetaLSTM only on Omniglot dataset. We extend upon this empirical investigation by comparing the performance of MetaLSTM and MetaLSTM++ on complex datasets like miniImagenet, FC100, and tieredImagenet (Table 2). It is evident from the results that batch-wise episodic training is more effective than sequential episodic training. We also investigate the performance of the models trained with the TA meta-training regimen with their non-TA counterparts. Specifically, we compare MAML, MetaSGD, MetaLSTM++ and ANIL with TA-MAML, TA-MetaSGD, TA-MetaLSTM++ and TA-ANIL respectively over 5 and 10-way (1 and 5-shot) settings on miniImagenet, FC100 and tieredImagenet datasets and report the results in Table 2. For ANIL and TA-ANIL, we consider 1000 testing tasks. We observe that models trained with TA regimen generalize better to the unseen meta-test tasks than their non-task-attended versions across all the settings in all datasets. We also observe that the TA mechanism performs better than uniform sampling (Arnold et al., 2021) on the miniImagenet dataset on 1 and 5 shot settings for MAML and 1 shot setting on ANIL. Sampling episodes uniformly for ANIL in 5 way 5 shot setting is, however, better than attending to tasks in a batch. Note that the proposed task attention mechanism aims not to surpass the state-of-the-art meta-learning algorithms but provides new insight into the batch episodic meta-training regimen, which as per our knowledge, is common to all meta-learning algorithms.

We also compare the performance of TA-MAML against TAML - a meta-training regimen that forces the meta-model to be equally close to all the tasks. The results, as presented in Table 2, suggest that TA-MAML performs better than TAML on all benchmarks across all settings. Note that both TAML and TA-MAML are approaches that built upon MAML to address the inequality/diversity of tasks in a batch. Our aim is thus to compare TAML and TA-MAML and not to assess the efficacy of TAML when meta-trained using task attention. We investigate the influence of the TA meta-training regimen on the model's convergence by analyzing the trend of the model's validation accuracy over iterations. Figure 3 depicts the mean validation accuracy over 300 tasks on miniImagenet and tieredImagenet datasets for a 5-way 1-shot setting across

training iterations. We observe that the models meta-trained with TA regimen tend to achieve higher/at-par performance in fewer iterations than the corresponding models meta-trained with the non-TA regimen.

Table 2: Comparison of few-shot classification performance of vanilla ML algorithms with their task attended versions on miniImagenet, FC100 and tieredImagenet datasets for 5 and 10-way (1 and 5-shot) settings. The ± represents the 95% confidence intervals over 300 tasks. Algorithms denoted by * are rerun on their optimal hyper-parameters for a fair comparison. Attention-based ML algorithms perform better than their corresponding vanilla approaches across all the settings. Further, MetaLSTM++ and TA-MAML perform better than MetaLSTM and TAML, respectively, across all settings and datasets.

| | Test Accuracy (%) | | | |
| | 5-Way | | 10-Way | |
| Model | 1 Shot | 5 Shot | 1 Shot | 5 Shot |
|---|---|---|---|---|
| **miniImagenet** | | | | |
| MAML* | 46.10 ± 0.19 | 60.16 ± 0.17 | 29.42 ± 0.11 | 41.98 ± 0.10 |
| TAML* | 46.26 ± 0.21 | 53.40 ± 0.14 | 29.76 ± 0.11 | 36.88 ± 0.10 |
| MAML+UNIFORM (Offline) | 46.67 ± 0.63 | 62.09 ± 0.55 | - | - |
| MAML+UNIFORM (Online) | 46.70 ± 0.61 | 61.62 ± 0.54 | - | - |
| **TA-MAML** | **48.36 ± 0.23** | **62.48 ± 0.18** | **31.15± 0.11** | **43.70 ± 0.09** |
| MetaSGD* | 47.65± 0.21 | 61.60 ± 0.17 | 30.09± 0.10 | 42.22 ± 0.11 |
| **TA-MetaSGD** | **49.28 ± 0.20** | **63.37 ± 0.16** | **31.50± 0.11** | **44.06 ± 0.10** |
| MetaLSTM* | 41.48 ± 1.02 | 58.87 ± 0.94 | 28.62 ± 0.64 | 44.03 ± 0.69 |
| MetaLSTM++ | 48.00 ± 0.19 | 62.73 ± 0.17 | 31.16 ± 0.09 | 45.46 ± 0.10 |
| **TA-MetaLSTM++** | **49.18 ± 0.17** | **64.89 ± 0.16** | **32.07± 0.11** | **46.66 ± 0.09** |
| ANIL* | 46.92 ± 0.62 | 58.68 ± 0.54 | 28.84 ± 0.34 | 40.95 ± 0.32 |
| ANIL+UNIFORM (Offline) | 46.93 ± 0.62 | **62.75 ± 0.60** | - | - |
| ANIL+UNIFORM (Online) | 46.82 ± 0.63 | 62.63 ± 0.59 | - | - |
| **TA-ANIL** | **48.84 ± 0.62** | 60.80± 0.55 | **31.14± 0.34** | **42.52 ± 0.34** |
| **FC100** | | | | |
| MAML* | 36.40 ± 0.38 | 46.76±0.21 | 23.93±0.14 | 31.14 ± 0.07 |
| TAML* | 38.00 ± 0.26 | 48.05± 0.13 | 21.60± 0.14 | 33.19± 0.07 |
| **TA-MAML** | **39.86± 0.25** | **49.56 ± 0.13** | **25.46± 0.15** | **36.06± 0.08** |
| MetaSGD* | 33.46 ± 0.23 | 43.96± 0.13 | 21.40±0.15 | 30.59± 0.07 |
| **TA-MetaSGD** | **35.66±0.25** | **49.49± 0.12** | **23.80±0.15** | **32.08±0.07** |
| MetaLSTM* | 37.20 ± 0.26 | 47.89 ± 0.13 | 21.70 ± 0.14 | 32.11 ± 0.07 |
| MetaLSTM++ | 38.60 ±0.23 | 49.82 ± 0.12 | 22.80 ± 0.14 | 33.46 ± 0.08 |
| **TA-MetaLSTM++** | **41.53 ±0.28** | **51.17 ±0.13** | **25.33 ±0.15** | **34.18 ±0.08** |
| ANIL* | 34.08 ± 1.29 | 44.74 ± 0.68 | 20.65 ± 0.77 | 27.93 ± 0.42 |
| **TA-ANIL** | **38.06 ± 1.26** | **46.94± 0.69** | **23.27± 0.79** | **28.29 ± 0.40** |
| **tieredImagenet** | | | | |
| MAML* | 44.40 ± 0.49 | 57.07 ± 0.22 | 27.40 ± 0.25 | 34.30 ± 0.14 |
| TAML* | 46.40 ± 0.40 | 56.80 ± 0.23 | 26.40 ± 0.25 | 34.40 ± 0.15 |
| **TA-MAML** | **48.40 ± 0.46** | **60.40 ± 0.25** | **31.00± 0.26** | **37.60± 0.15** |
| MetaSGD* | 52.80 ± 0.44 | 62.35 ± 0.26 | 31.90 ± 0.27 | 44.16 ± 0.15 |
| **TA-MetaSGD** | **56.20 ± 0.45** | **64.56 ± 0.24** | **33.20± 0.29** | **47.12 ± 0.16** |
| MetaLSTM* | 37.00 ± 0.44 | 59.83 ± 0.25 | 29.80 ± 0.28 | 39.28 ± 0.13 |
| MetaLSTM++ | 47.60 ± 0.49 | 63.24 ± 0.25 | 30.70 ± 0.27 | 47.97 ± 0.16 |
| **TA-MetaLSTM++** | **49.00 ± 0.44** | **66.15 ± 0.23** | **32.10± 0.27** | **51.35 ± 0.17** |
| ANIL* | 45.08 ± 1.37 | 59.71 ±0.77 | 29.32 ± 0.83 | 42.76 ± 0.50 |
| **TA-ANIL** | **45.96 ± 1.32** | **60.96± 0.72** | **32.68± 0.92** | **47.56 ± 0.51** |

## 5.3 Comparison with Baselines

Yao et al. (2021) proposed Adaptive Task Scheduler (ATS) and ascertained the merit of ATS over Greedy class-pair (GCP) technique (Liu et al., 2020) on miniImagenet dataset. We extend this comparison and show in Table 3 that the proposed approach performs better than state-of-the-art ATS and GCP only in 1 shot setting. ATS has been designed for noisy and imbalanced task distributions. So, we compare the proposed approach with GCP, ATS, and other sampling techniques on the miniImagenet-noisy dataset (Yao et al., 2021) and report the results in Table 3.
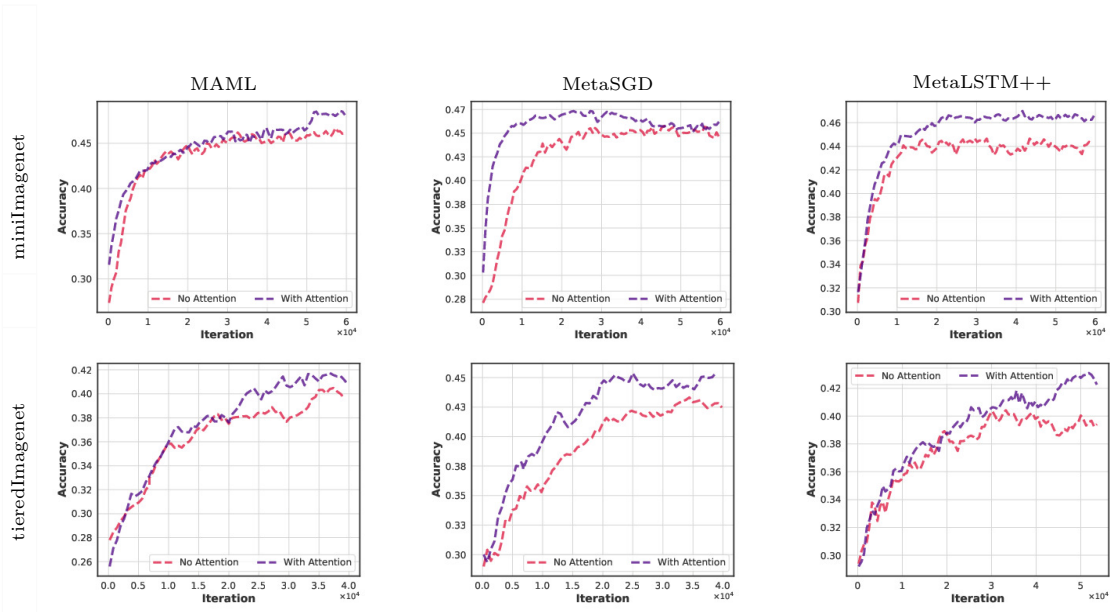
Figure 3: Mean validation accuracies of MAML (Col-1), MetaSGD (Col-2) and MetaLSTM++ (Col-3) across 300 tasks with/without attention on 5-way 1-shot setting on miniImagenet (Row-1) and tieredImagenet (Row-2) datasets.

We observe that task attention outperforms all scheduling algorithms on the miniImagenet-noisy dataset. As ATS is the most competitive baseline for the proposed method on the miniImagenet-noisy dataset, we compare the TA-ANIL and ATS on varying noise ratios for the miniImagenet dataset on 5 way 1 shot setting (Table 4). We observe that the proposed method outperforms ATS on all noise ratios except 0.8. Note that the algorithm used for all sampling approaches is ANIL.

## 5.4 Effectiveness of Task Attention in CDFSL setup

Classical meta-learning approaches assume meta-train and meta-test data belong to the same distribution such that the meta-trained model extends its knowledge to the meta-test set. This is, however, not always the case. The difference in the data acquisition techniques, or evolution of data with time, may cause a discrepancy between the meta-train and meta-test distributions. This realistic setting is popularly termed as cross-domain few-shot learning (CDFSL) (Guo et al., 2020). We conducted experiments to show the merit of the proposed approach in CDFSL setup. Specifically, we train a model using TA meta-training regimen on the miniImagenet dataset and meta-test it on CUB-200 and FGVC-Aircraft datasets. The results reported for 5 way 1 and 5 shot settings in Table 5 indicate that the proposed approach outperforms the state-of-the-art task scheduling approach (Uniform Sampling (Arnold et al., 2021)) on CDFSL setup by a large margin.

## 5.5 Ablation Studies

To examine the significance of each input given to the task attention model, we conduct an ablation study on 5-way 1 and 5 shot TA-MAML on miniImagenet dataset and report the results in Table 6. We observe that all the components of meta-information contribute to the learning of a more generalizable meta-model.

Table 3: Comparison (Test Accuracy (%)) of task attention with GCP and ATS for MAML and MetaSGD on miniImagenet dataset and various sampling techniques for ANIL on the miniImagenet-noisy dataset for 5 way 1 and 5 shot settings.

| | 5-Way | |
|---|---|---|
| **Model** | 1 Shot | 5 Shot |
| **miniImagenet** | | |
| MAML with GCP | $46.92 \pm 0.83$ | $63.28 \pm 0.66$ |
| MAML with ATS | $47.89 \pm 0.77$ | **$64.07 \pm 0.70$** |
| **TA-MAML (Ours)** | **$48.36 \pm 0.23$** | $62.48 \pm 0.18$ |
| MetaSGD with GCP | $47.77 \pm 0.75$ | $63.50 \pm 0.71$ |
| MetaSGD with ATS | $48.59 \pm 0.79$ | **$64.79 \pm 0.74$** |
| **TA-MetaSGD (Ours)** | **$49.28 \pm 0.20$** | $63.37 \pm 0.16$ |
| **miniImagenet-noisy** | | |
| Uniform | $41.67 \pm 0.80$ | $55.80 \pm 0.71$ |
| SPL | $42.13 \pm 0.79$ | $56.19 \pm 0.70$ |
| Focal Loss | $41.91 \pm 0.78$ | $53.58 \pm 0.75$ |
| GCP | $41.86 \pm 0.75$ | $54.63 \pm 0.72$ |
| PAML | $41.49 \pm 0.74$ | $52.45 \pm 0.69$ |
| DAML | $41.26 \pm 0.73$ | $55.46 \pm 0.70$ |
| ATS | $44.21 \pm 0.76$ | $59.50 \pm 0.71$ |
| **TA-ANIL (Ours)** | **$45.17 \pm 0.23$** | **$62.15 \pm 1.01$** |

To further support this observation, we investigate the relationship between the meta-information and weights assigned by the task attention module by analyzing the mean Pearson correlation of each of the components (four tuple) of the meta-information with the attention vector across the training iterations. This is depicted in Figure 4 for TA-MAML on 5-way 1 and 5 shot settings for miniImagenet dataset. We observe that the loss ratio and loss are positively correlated with the attention vector, while accuracy and gradient norm are negatively correlated.

In 5-way 5-shot setting, we observe that the correlation pattern is comparable to 5-way 1-shot setting, but the mean correlation value of grad norm across iterations is less than that of the 5-way 1-shot setting. This could be because the 5-way 5-shot setting is richer in data than the 5-way 1-shot setting, which allows better learning and therefore has low average values of grad norm (Section 4.1.1). The critical observation, however, is that the meta-information components have a weak correlation with the attention weights, indicating that the TA module does not trivially follow any single component of meta-information. We also analyze the ranks of the tasks for maximum and minimum values of : loss, loss ratio, accuracy, and grad norm in a batch, as per the weights across training iterations, and describe results in the supplementary material. The rank analysis also reinforces the same observation. We ascertain the decreasing trend of mean weighted loss across iterations in the supplementary material.

Table 4: Comparative analysis of ANIL integrated with ATS and proposed method on miniImagenet dataset with varying noise ratios for 5 way 1 shot setting. BNS is the best non-adaptive scheduler.

| | Test Accuracy (%) on miniImagenet-noisy | | | |
|---|---|---|---|---|
| Noise ratio | 0.2 | 0.4 | 0.6 | 0.8 |
| ANIL with Uniform | $43.46 \pm 0.82$ | $42.92 \pm 0.78$ | $41.67 \pm 0.80$ | $36.53 \pm 0.73$ |
| ANIL with BNS | $44.04 \pm 0.81$ | $43.36 \pm 0.75$ | $42.13 \pm 0.79$ | $38.21 \pm 0.75$ |
| ANIL with ATS | $45.55 \pm 0.80$ | $44.50 \pm 0.86$ | $44.21 \pm 0.76$ | $\mathbf{42.18 \pm 0.73}$ |
| **TA-ANIL (Ours)** | $\mathbf{47.98 \pm 0.26}$ | $\mathbf{46.69 \pm 0.22}$ | $\mathbf{45.17 \pm 0.23}$ | $40.35 \pm 1.14$ |

Table 5: Comparative analysis of proposed approach and uniform sampling (Arnold et al., 2021) in a CDFSL setting after training on miniImagenet dataset and tested on CUB-200 and FGVC-Aircraft datasets for 5 way 1 and 5 shot settings.

| | 5-Way | |
|---|---|---|
| Model | 1 Shot | 5 Shot |
| | **CUB-200** | |
| MAML+ UNIFORM (Online) | $35.84 \pm 0.54$ | $46.67 \pm 0.55$ |
| **TA-MAML (Ours)** | $\mathbf{42.87 \pm 1.18}$ | $\mathbf{57.49 \pm 0.99}$ |
| | **FGVC-Aircraft** | |
| MAML+ UNIFORM (Online) | $26.62 \pm 0.39$ | $34.41 \pm 0.44$ |
| **TA-MAML (Ours)** | $\mathbf{29.42 \pm 0.78}$ | $\mathbf{36.34 \pm 0.86}$ |

## 5.6 Analysis of Attention Network

To gain further insights into the operation of the attention module, we also examine the trend of the attention-vector (Figure 5) while meta-training TA-MAML for 5 way 1 and 5 shot settings on the miniImagenet dataset. We plot the maximum and the minimum attention score assigned to the tasks of a batch across iterations together with a few weighted task batches in 5-way 1-shot setting for illustration. We note that the weighted task batches are only intended to demonstrate the change in the tasks' attention scores across iterations. The next experiment presents a more rigorous analysis studying the relationship among classes in a task and attention scores assigned. We note that the mean attention score is always 0.25 as we follow a meta-batch size of 4. We observe that the TA module's output follows an interesting trend. Initially, the TA module assigns almost uniform weights to all the tasks of a batch; however, as the iterations increase, it assigns unequal scores to the tasks in a batch, preferring some over the other. This suggests that during the initial phases of the meta-model's training, all tasks have equal contribution towards learning a *generic structure* of the meta-knowledge. As the meta-model's learning proceeds, learning the further *fine-grained meta-knowledge structure* requires

Table 6: Effect of ablating components of meta-information in TA-MAML for 5 way 1 and 5 shot settings on miniImagenet dataset.

| Ablation on inputs | | | | | |
|---|---|---|---|---|---|
| Grad norm | Loss | Loss-ratio | Accuracy | Test Accuracy | |
| | | | | 5-way 1-shot | 5-way 5-shot |
| ✗ | ✗ | ✗ | ✗ | $46.10 \pm 0.19$ | $60.16 \pm 0.17$ |
| ✓ | ✓ | ✓ | ✗ | $47.30 \pm 0.16$ | $60.48 \pm 0.16$ |
| ✓ | ✓ | ✗ | ✓ | $47.62 \pm 0.17$ | $62.17 \pm 0.17$ |
| ✓ | ✗ | ✓ | ✓ | $48.10 \pm 0.18$ | $60.90 \pm 0.20$ |
| ✗ | ✓ | ✓ | ✓ | $47.30 \pm 0.18$ | $61.52 \pm 0.16$ |
| ✓ | ✓ | ✓ | ✓ | $\mathbf{48.36 \pm 0.23}$ | $\mathbf{62.48 \pm 0.18}$ |

prioritizing some tasks in a batch over the others, which are potentially better aligned with learning the optimal meta-knowledge. We study the computational feasibility of TA regimen in the appendix.

We further decipher the functioning of the black box attention network by analyzing the qualitative relation among weights and the classes of task batches (Figure 10 is presented in appendix due to space constraints). In Figure 10 left column (col-1) corresponds to the cases where the assignment of attention scores to the tasks is human interpretable. In contrast, the right column (col-2) refers to the uninterpretable attention scores. From the human perspective, tasks con-



Figure 4: Mean Pearson correlation of TA-MAML on 5-way 1-shot (left) and 5-shot (right) setting on miniImagenet.

taining images from similar classes are hard to distinguish and are assigned higher attention scores indicated by red bounding boxes (Figure 10 col-1). Specifically, (col-1, row-1) task 2 is regarded as most important, possibly because it includes three breeds of dogs followed by task 4, which comprises two species of fish. However, the aforementioned is not a hard constraint, as there are some task batches (Figure 10 col-2) in which the distribution of weights cannot be explained qualitatively.

## 6 Conclusion

In this work we have shown that the batch wise episodic training regimen adopted by ML strategies can benefit from leveraging knowledge about the importance of tasks within a batch. Unlike prior approaches that assume uniform importance for each task in a batch, we propose task attention as a way to learn the relevance of each task according to its alignment with the optimal meta-knowledge. We have validated the effectiveness of task attention by augmenting it to popular initialization and optimization based ML strategies. We have demonstrated through experiments on miniImagenet, FC100 and tieredImagenet datasets that augmenting task at-



Figure 5: Trend of an attention vector in 5-way 1-shot (left) and 5-shot (right) settings on miniImagenet dataset for TA-MAML.

tention helps attain better generalization to unseen tasks from the same distribution while requiring fewer iterations to converge. We also show that the task attention is meritorious over existing task scheduling algorithms, even on noisy and CDFSL setups. We also conduct an exhaustive empirical analysis on the distribution of attention weights to study the nature of the meta-knowledge and task attention module. We leave the theoretical motivation of the meta-information components and the proof of convergence of the proposed curriculum as part of our future work. We believe that this end-to-end attention-based meta training paves the way towards efficient and automated meta-training.

## References

Aroof Aimen, Sahil Sidheekh, Vineet Madan, and Narayanan C Krishnan. Stress Testing of Meta-learning Approaches for Few-shot Learning. In *AAAI Workshop on Meta-Learning and MetaDL Challenge*, 2021.

Sébastien Arnold, Guneet Dhillon, Avinash Ravichandran, and Stefano Soatto. Uniform sampling over episode difficulty. *Advances in Neural Information Processing Systems*, 34:1481–1493, 2021.

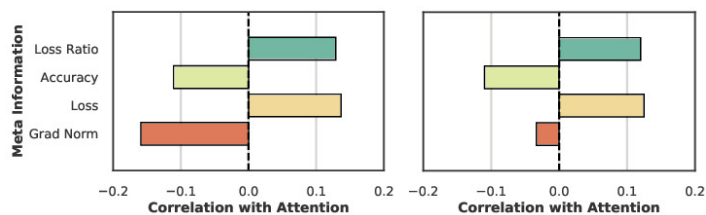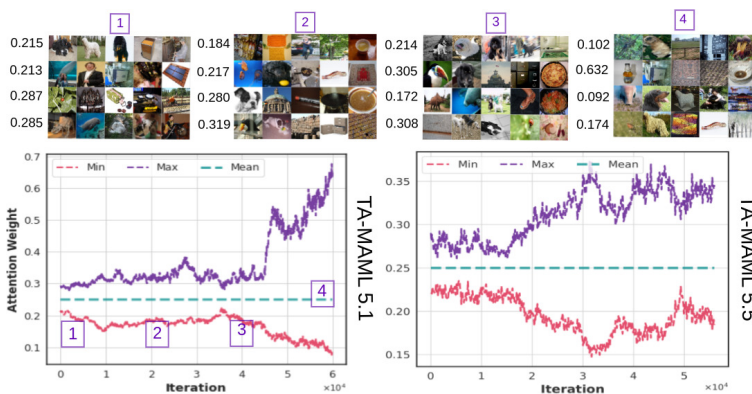Sébastien MR Arnold, Praateek Mahajan, Debajyoti Datta, Ian Bunner, and Konstantinos Saitas Zarkias. learn2learn: A library for meta-learning research. *CoRR*, 2020.

Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009*, volume 382 of *ACM International Conference Proceeding Series*, pp. 41–48. ACM, 2009.

Haw-Shiuan Chang, Erik G. Learned-Miller, and Andrew McCallum. Active bias: Training more accurate neural networks by emphasizing high variance samples. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 1002–1012, 2017.

Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.

Yunhui Guo, Noel C Codella, Leonid Karlinsky, James V Codella, John R Smith, Kate Saenko, Tajana Rosing, and Rogerio Feris. A broader study of cross-domain few-shot learning. In *European conference on computer vision*, pp. 124–141. Springer, 2020.

Ricardo Luna Gutierrez and Matteo Leonetti. Information-theoretic task selection for meta-reinforcement learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.

Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 1997.

Muhammad Abdullah Jamal and Guo-Jun Qi. Task agnostic meta-learning for few-shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 11719–11727. Computer Vision Foundation / IEEE, 2019.

Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2309–2318. PMLR, 2018.

Jean Kaddour, Steindór Sæmundsson, and Marc Peter Deisenroth. Probabilistic active meta-learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.

H. Kahn and A. W. Marshall. Methods of reducing sample size in monte carlo computations. *Oper. Res.*, 1953.

Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

M. Pawan Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada*, pp. 1189–1197. Curran Associates, Inc., 2010.

Jian Li, Xuanyuan Luo, and Mingda Qiao. On generalization error bounds of noisy gradient methods for non-convex learning. In *International Conference on Learning Representations*, 2019.

Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few-shot learning, 2017.

Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pp. 2999–3007. IEEE Computer Society, 2017.

Chenghao Liu, Zhihao Wang, Doyen Sahoo, Yuan Fang, Kun Zhang, and Steven C. H. Hoi. Adaptive task sampling for meta-learning. In *ECCV*, 2020.

Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *ICML*, 2021.

Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.

Boris N. Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. TADAM: task dependent adaptive metric for improved few-shot learning. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 719–729, 2018.

Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse? towards understanding the effectiveness of MAML. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017.

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. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018a.

Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4331–4340. PMLR, 2018b.

Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

Abhinav Shrivastava, Abhinav Gupta, and Ross B. Girshick. Training region-based object detectors with online hard example mining. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 761–769. IEEE Computer Society, 2016.

Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 403–412. Computer Vision Foundation / IEEE, 2019.

Qianru Sun, Yaoyao Liu, Zhaozheng Chen, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning through hard tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 3630–3638, 2016.

Peter Welinder, Steve Branson, Takeshi Mita, Catherine Wah, Florian Schroff, Serge Belongie, and Pietro Perona. Caltech-ucsd birds 200. 2010.

Huaxiu Yao, Yu Wang, Ying Wei, Peilin Zhao, Mehrdad Mahdavi, Defu Lian, and Chelsea Finn. Meta-learning with an adaptive task scheduler. *Advances in Neural Information Processing Systems*, 2021.

Peilin Zhao and Tong Zhang. Stochastic optimization with importance sampling for regularized loss minimization. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 1–9. JMLR.org, 2015.

# 7 Appendix

## 7.1 Experiments

### 7.1.1 Datasets Details

**miniImagenet** dataset (Vinyals et al., 2016) comprises 600 color images of size $84 \times 84$ from each of 100 classes sampled from the Imagenet dataset. The 100 classes are split into 64, 16 and 20 classes for meta-training, meta-validation and meta-testing respectively. **miniImagenet-noisy** (Yao et al., 2021) is constructed from the miniImagenet dataset with the additional constraint that tasks have noisy support labels and clean query labels. The noise in support labels is introduced by symmetry flipping, and the default noise ratio is 0.6. **Fewshot Cifar 100 (FC100)** dataset (Oreshkin et al., 2018) has been created from Cifar 100 object classification dataset. It contains 600 color images of size $32 \times 32$ corresponding to each of 100 classes grouped into 20 super-classes. Among 100 classes, 60 classes belonging to 12 super-classes correspond to the meta-train set, 20 classes from 4 super-classes to the meta-validation set, and the rest to the meta-test set. **tieredImagenet** (Ren et al., 2018a) is a more challenging benchmark for few-shot image classification. It contains 779,165 color images sampled from 608 classes of Imagenet and are grouped into 34 super-classes. These super-classes are divided into 20, 6, and 8 disjoint sets for meta-training, meta-validation, and meta-testing. **CUB-200** (Welinder et al., 2010) comprises of 6033 bird images corresponding to 200 species. We use its modified version (Arnold et al., 2021), wherein the images overlapping with Imagenet dataset have been removed. This avoids bias during CDFSL from miniImagenet $\rightarrow$ CUB-200. The meta-test set contains images from 30 classes. **FGVC Aircrafts** (Maji et al., 2013) contains 10200 aircraft images from 102 classes, among which 15 classes are present in the test split. Each class contains 100 examples.

### 7.1.2 Ablation Studies

We analyze the ranks of the tasks for maximum and minimum values of : loss, loss ratio, accuracy, and grad norm in a batch wrt attention weights throughout meta-training of TA-MAML on a 5-way 1 and 5 shot settings on miniImagenet dataset (Figure 6 and 7). Specifically, the highest weighted task is given rank one, and the least weighted task in a batch is given the last rank. We observe that the TA module does not assign maximum weight to the tasks with maximum or minimum values of : test loss, loss ratio, grad norm or accuracy throughout meta-training. Thus, the TA module does not trivially learn to assign weights to the tasks based on some component of meta-information but learns useful latent information from all the components to assign importance for the tasks in a batch.

## 7.2 Relation of Weights with Meta-Information

In Figure 8, we illustrate the trend of mean weighted loss across iterations for TA-MAML on 5-way 1 and 5 shot settings on miniImagenet dataset. The trend indicates that the average weighted loss decreases over the meta-training iterations. The shaded region represents a 95% confidence interval over 100 tasks.

**5-way 1-shot setting**



Figure 6: Rank Analysis of tasks for maximum and minimum values of : loss, loss-ratio, accuracy and grad norm throughout the training of TA-MAML for 5-way 1 shot setting on miniImagenet dataset.

**5-way 5-shot setting**
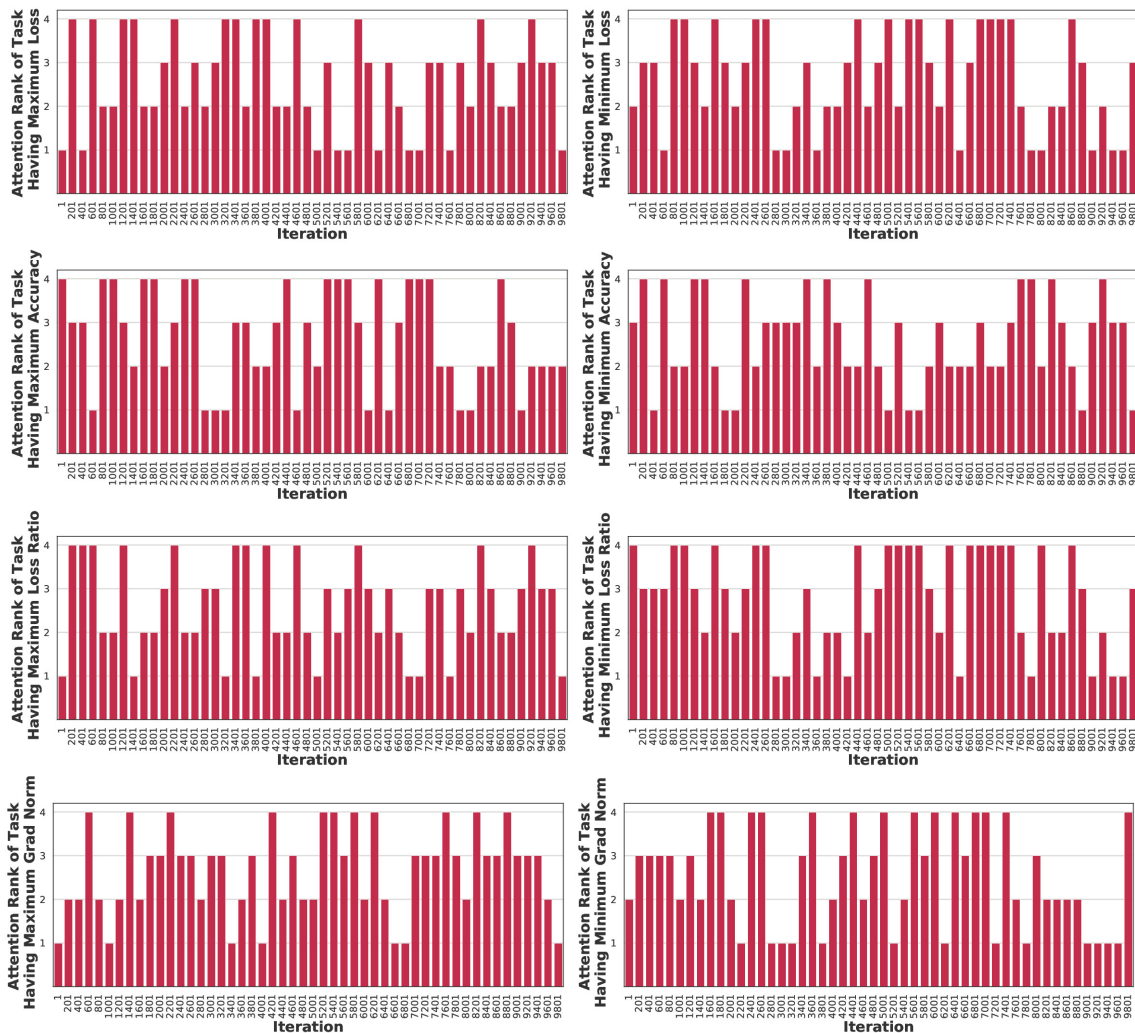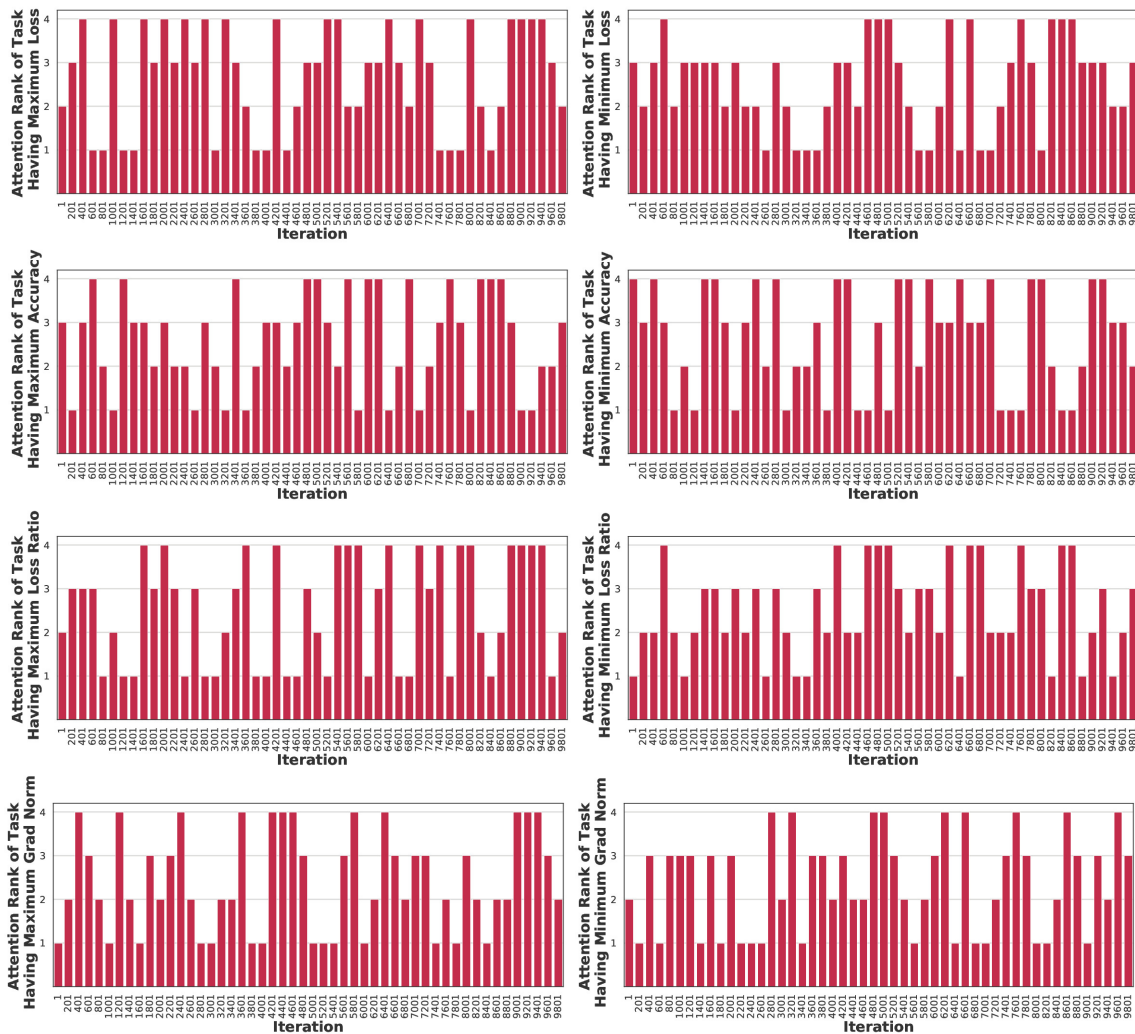


Figure 7: Rank Analysis of tasks for maximum and minimum values of : loss, loss-ratio, accuracy and grad norm throughout the training of TA-MAML for 5-way 5 shot setting on miniImagenet dataset.
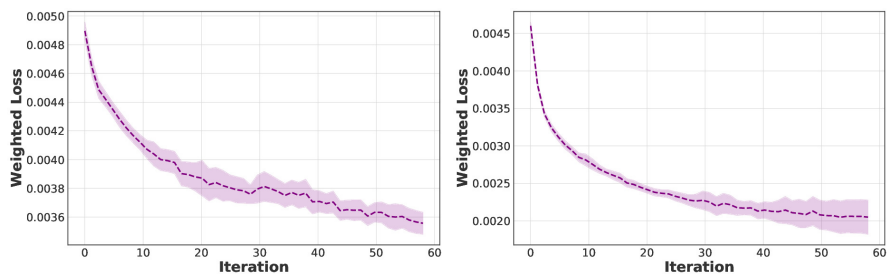


Figure 8: Trend analysis of weighted loss across meta-training iterations for TA-MAML on 5-way 1-shot (left) and 5-shot (right) settings on miniImagenet dataset. Iterations are in thousands.

Figure 9: Trend of an attention vector for TA-MAML when attention module is frozen after 15000 iterations in 5-way 1-shot setting on miniImagenet dataset.

### 7.2.1 Analysis of Attention Network

To reduce the computational burden, we freeze the weights of the attention module after 15000 iterations, i.e., only inputs of the attention module vary beyond 15000 iterations. We obtained a similar performance as when the attention module was trained throughout the meta-train phase ($\approx 48\%$ for 5-way 1-shot setting on miniImagenet dataset). From Figure 9, we observe that the attention vector still follows a similar trend as when trained end-to-end, indicating 15000 iterations are sufficient for the attention module's training. Thus, we note that proposed approach is computationally feasible.

Due to space constraints in the main paper, we illustrate the qualitative relation among weights and the classes of task batches in Figure 10.

### 7.2.2 Hyperparameter Details

Figure 10: Explanations of TA module in TA-MAML on miniImagenet. **Left Col)** Higher weights accredited to tasks with comparable classes marked by red bounding boxes. **Right Col)** Association of weights and task data is qualitatively uninterpretable. Rows correspond to the batches.

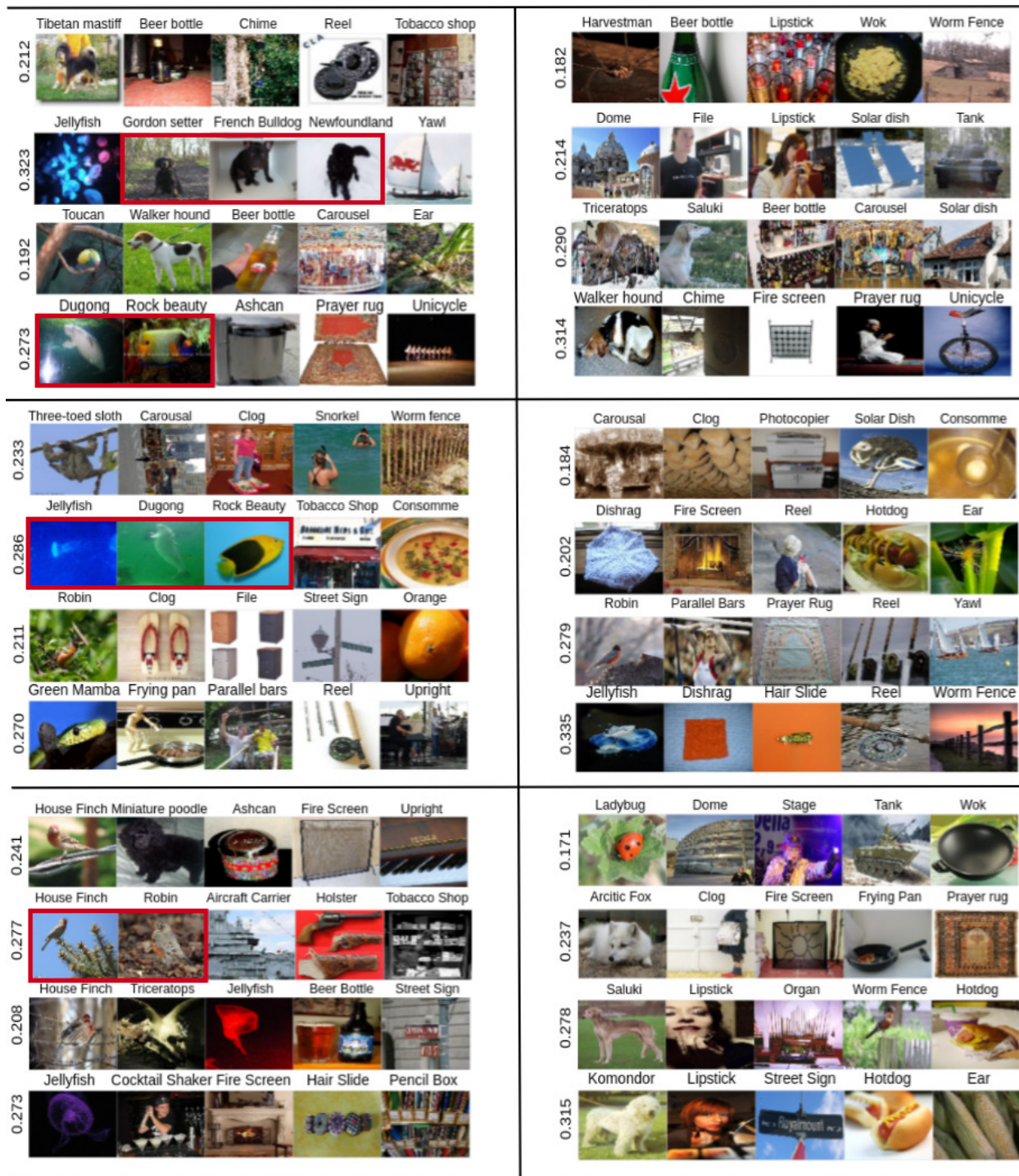| Setting | Model | base lr | meta lr | attention lr | lambda |
|---------|-------|---------|---------|--------------|--------|
| | | **miniImagenet** | | | |
| 5.1 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.0748 |
| | TA-MAML | 0.0763 | 0.0005 | 0.0004 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0529 | 0.0011 | 0.0004 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0012 | - | - |
| | TA-MetaLSTM++ | - | 0.0012 | 0.0031 | - |
| | ANIL | 0.3000 | 0.0006 | - | - |
| | TA-ANIL | 0.0763 | 0.0005 | 0.0004 | - |
| 5.5 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.7916 |
| | TA-MAML | 0.0763 | 0.0005 | 0.0004 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0529 | 0.0011 | 0.0004 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0012 | - | - |
| | TA-MetaLSTM++ | - | 0.0004 | 0.0001 | - |
| | ANIL | 0.3000 | 0.0006 | - | - |
| | TA-ANIL | 0.0763 | 0.0005 | 0.0004 | - |
| 10.1 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.2631 |
| | TA-MAML | 0.2551 | 0.0015 | 0.0001 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0627 | 0.0008 | 0.0013 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0015 | - | - |
| | TA-MetaLSTM++ | - | 0.0009 | 0.0015 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.2551 | 0.0015 | 0.0001 | - |
| 10.5 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.0741 |
| | TA-MAML | 0.2551 | 0.0015 | 0.0001 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0627 | 0.0008 | 0.0013 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0036 | - | - |
| | TA-MetaLSTM++ | - | 0.0024 | 0.0002 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.2551 | 0.0015 | 0.0001 | - |

| Setting | Model | base lr | meta lr | attention lr | lambda |
|---------|-------|---------|---------|--------------|--------|
| | | | **FC100** | | |
| 5.1 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.0164 |
| | TA-MAML | 0.2826 | 0.0003 | 0.0024 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0349 | 0.0008 | 0.0001 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0010 | - | - |
| | TA-MetaLSTM++ | - | 0.0002 | 0.0074 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.2826 | 0.0003 | 0.0024 | - |
| 5.5 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.0153 |
| | TA-MAML | 0.2826 | 0.0003 | 0.0024 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0349 | 0.0008 | 0.0001 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0002 | - | - |
| | TA-MetaLSTM++ | - | 0.0007 | 0.0003 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.2826 | 0.0003 | 0.0024 | - |
| 10.1 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.0794 |
| | TA-MAML | 0.2353 | 0.0002 | 0.0001 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.2583 | 0.0029 | 0.0007 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0021 | - | - |
| | TA-MetaLSTM++ | - | 0.0005 | 0.0014 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.2826 | 0.0003 | 0.0024 | - |
| 10.5 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.0193 |
| | TA-MAML | 0.2353 | 0.0002 | 0.0001 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.2583 | 0.0029 | 0.0007 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0004 | - | - |
| | TA-MetaLSTM++ | - | 0.0004 | 0.0090 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.2826 | 0.0003 | 0.0024 | - |

| Setting | Model | base lr | meta lr | attention lr | lambda |
|---------|-------|---------|---------|--------------|--------|
| | | | **tieredImagenet** | | |
| 5.1 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.3978 |
| | TA-MAML | 0.0261 | 0.0005 | 0.0015 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0944 | 0.0003 | 0.0002 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0002 | - | - |
| | TA-MetaLSTM++ | - | 0.0010 | 0.0006 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.0261 | 0.0005 | 0.0015 | - |
| 5.5 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.7733 |
| | TA-MAML | 0.0261 | 0.0005 | 0.0015 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0944 | 0.0003 | 0.0002 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0009 | - | - |
| | TA-MetaLSTM++ | - | 0.0012 | 0.0001 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.0261 | 0.0005 | 0.0015 | - |
| 10.1 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.4752 |
| | TA-MAML | 0.0821 | 0.0002 | 0.0006 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0512 | 0.0007 | 0.0018 | - |
| | MetaLSTM | - | 0.005 | - | - |
| | MetaLSTM++ | - | 0.0011 | - | - |
| | TA-MetaLSTM++ | - | 0.0018 | 0.0002 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.0821 | 0.0002 | 0.0006 | - |
| 10.5 | MAML | 0.5000 | 0.0030 | - | - |
| | TAML | 0.5000 | 0.0030 | - | 0.2501 |
| | TA-MAML | 0.0821 | 0.0002 | 0.0006 | - |
| | MetaSGD | 0.5000 | 0.0030 | - | - |
| | TA-MetaSGD | 0.0512 | 0.0007 | 0.0018 | - |
| | MetaLSTM | - | 0.0050 | - | - |
| | MetaLSTM++ | - | 0.0024 | - | - |
| | TA-MetaLSTM++ | - | 0.0015 | 0.0019 | - |
| | ANIL | 0.5000 | 0.0030 | - | - |
| | TA-ANIL | 0.0821 | 0.0002 | 0.0006 | - |