

Prior-Knowledge and Attention based Meta-Learning for Few-Shot Learning

Yunxiao Qin^{1,2}, Weiguo Zhang¹, Chenxu Zhao², Zezheng Wang², Xiangyu Zhu³
Guojun Qi⁴, Jingping Shi¹, Zhen Lei³

¹Northwestern Polytechnical University of China ²AIBEE,

³Institute of Automation, Chinese Academy of Science, ⁴Huawei Cloud
{qyxqyx, zhangwg, shijingping}@mail.nwpu.edu.cn, guojunq@gmail.com,
{xiangyu.zhu, zlei}@nlpr.ia.ac.cn, {cxzhao, zezhengwang}@jd.com

Abstract

Recently, meta-learning has been shown as a promising way to solve few-shot learning. In this paper, inspired by the human cognition process which utilizes both prior-knowledge and vision attention in learning new knowledge, we present a novel paradigm of meta-learning approach with three developments to introduce attention mechanism and prior-knowledge for meta-learning. In our approach, prior-knowledge is responsible for helping meta-learner expressing the input data into high-level representation space, and attention mechanism enables meta-learner focusing on key features of the data in the representation space. Compared with existing meta-learning approaches which pay little attention to prior-knowledge and vision attention, our approach alleviates the meta-learner's few-shot cognition burden. Furthermore, a Task-Over-Fitting (TOF) problem¹, which indicates that the meta-learner has poor generalization on different K -shot learning tasks, is discovered and we propose a Cross Entropy across Tasks (CET) metric² to model and solve the TOF problem. Extensive experiments demonstrate that we improve the meta-learner with state-of-the-art performance on several few-shot learning benchmarks, and at the same time the TOF problem can also be released greatly.

1. Introduction

The development of deep learning makes remarkable progresses in many tasks [1–4]. To achieve all of them, large amounts of thousands and even millions of labeled data are required for the deep learning approach to obtain

¹When is tested on J -shot classification tasks, the meta-learner trained on K -shot tasks performs not as well as the one trained on J -shot tasks, where K and J are different unsigned integers denoting different numbers of shots for the meta-learner.

²A metric for quantizing how much a meta-learning method suffers from the TOF problem.

satisfactory performance. However, collecting and annotating abundant data is notoriously expensive. Therefore, few-shot learning [5–7] which requires the model to learn from a few data, has attracted researchers' attention in recent years.

Learning from few-data is challenging for Computer Vision. In comparison, we human beings can rapidly learn new categories from very few examples. Recently, meta-learning [8–19] has shown promising performance to improve the few-shot learning for Computer Vision. However, existing meta-learning methods commonly ignore prior-knowledge [20–24] and attention mechanism [25,26] which have been both demonstrated important for human cognitive and learning process. We illustrate a few-shot classification problem in Fig.1 for a better understanding of the role of prior-knowledge and attention mechanism in human few-shot learning. In Fig.1, we unconsciously leverage our learned knowledge about the world to understand and express these images into high-level compact representations, such as plant, animal, tree, and table *etc.* However, according to the four training images, we discover that only the feature of the tree and table are useful for us to recognize these two classes of images. Then, we quickly adjust ourselves to pay attention to the critical features and make the decision based on the focused features.

Evidently, we can summarize two main modules in human few-shot learning: **a stable Representation module that utilizes prior-knowledge to express the image into compact feature representations; and a smart attention-based decision logical module that adapts accurately and performs recognition based on the feature representations.** While existing meta-learning approaches commonly train meta-learners to learn adaptive networks directly based on the original input data with no attention mechanism and prior-knowledge.

In this paper, inspired by the human cognition process, we present a novel paradigm of meta-learning approach with three developments to introduce attention mechanism and prior-knowledge step-by-step for meta-learning. Here,



Figure 1: An Example of few-shot classification task. The six images come from two classes, where four labeled ones are training data with the two unlabeled for test. When predicting the two testing images, we utilize our prior-knowledge about the world to understand all components in these images and use our vision attention to pay attention to key components table and tree. Finally, we predict the image (c) belongs to class 1 that contains table, while the image (f) is associated with class 2 of tree.

we briefly introduce the proposed methods. **1)** The first method is **Attention based Meta-Learning (AML)** which leverages attention mechanism to enable the meta-learner paying more attention on essential feature. **2)** For the meta-learner enjoying not only attention but also prior-knowledge, we present another method **Representation and Attention based Meta-Learning (RAML)**. Its network contains a Representation module and an attention-based prediction (ABP) module. The Representation module is similar to the same module of human vision. It learns the prior-knowledge in a supervised fashion and is responsible for understanding and extracting stable compact feature representations from the input image. The ABP module plays the same role as the smart attention-based decision logic module of human vision. It enables the meta-learner to precisely adjusting first its attention to the most discriminative feature representations of input images and second the corresponding predictions. **3)** In the third method, to take full advantage of endless unlabeled data, we design a novel method where the Representation module learns the past knowledge in unsupervised fashion [27–32]. We call this method **Unsupervised Representation and Attention based Meta-Learning (URAML)**. With URAML, we show in our experiments that the growth of the number of unlabeled data and the development of unsupervised learning both improve the performance of URAML apparently.

In addition, we show a Task-Over-Fitting (TOF) problem for existing meta-learning and present a Cross-Entropy across Tasks (CET) metric to evaluate how much a meta-

learning method is troubled by the TOF problem. An example of the TOF problem is, the meta-learner trained on 5-way 1-shot tasks is not as capable as the one trained on 5-way 5-shot tasks when they are tested on 5-way 5-shot tasks, and vice versa. However, in practical applications, it is uncertain how much data and how many shot times are available to the meta-learner to learn. Therefore, we argue that the trained meta-learner should generalizes well to different K -shot tasks. The possible reason behind the TOF problem is that existing meta-learners are vulnerable to the features irrelevant to the presented tasks since they ignoring both priori knowledge and attention mechanism. Our experiment validates that by incorporating prior-knowledge and attention mechanism, our methods suffer less from the TOF problem than existing meta-learning methods.

We summarize the main contributions of our work as:

- We propose that both attention mechanism and prior-knowledge are crucial for meta-learner to reduce its cognition burden in few-shot learning, and we develop three methods AML, RAML, and URAML to step-by-step leverage attention mechanism and prior-knowledge in meta-learning.
- We discover the TOF problem for meta-learning, and design a novel metric Cross-Entropy across Tasks (CET) to measure how much meta-learning approaches suffer from the TOF problem.
- Through extensive experiments, we show that the proposed methods achieve state-of-the-art performance on several few-shot learning benchmarks and in the meantime, they are less sensitive to the TOF problem, especially the RAML and URAML.

2. Related Work

2.1. Meta-learning for Few-Shot Learning

An N -way K -shot learning task contains a support set and a query set. The support and query set contain K and L examples for each of the N classes, respectively. Existing meta-learning approaches usually solve the few-shot learning by training a meta-learner on the N -way K -shot learning tasks in the following way. Firstly, the meta-learner is required to inner-update itself on the support set. Secondly, after the inner-updating, meta-learner is evaluated on the query set. Finally, by minimizing the loss on the query set, the meta-learner learns a base learner which has easy-fine-tune weights [11, 14] or a skillful weight updater [13, 19] or both [12] or the ability to memorize the support set [15]. The methods train the meta-learner learning an easy-fine-tune base learner are also called as weight initialization based methods, as the meta-learner learns generalized initial weight for few-shot learning tasks. Recently, MAML,

which is a classical weight initialization based method, is popular and lots of MAML based methods have been proposed. For example, LLAML [33] uses a local Laplace approximation to model the task parameters, and MTL [34] trains a meta-transfer to adapt a pre-trained deep network to few-shot learning tasks. Besides, MetaGAN [18] shows that by coupling MAML with adversarial training, the meta-learner is trained to learn a better decision boundaries between different classes in few-shot learning. To reduce the computation and memory cost of MAML, iMAML [35] leverages implicit differentiation to remove the need of differentiation through the inner-update path.

Though existing meta-learning methods performs promising, they seldom consider the prior-knowledge and attention mechanism in meta-learning. In our paper, we improve meta-learning for few-shot learning by introducing prior-knowledge and attention mechanism to meta-learning.

2.2. Attention Mechanism

Recent years, attention mechanism [36–39] has been widely used in computer vision systems, machine translation and *etc.*. Several manners of the attention mechanism have been proposed, such as soft attention [36, 37], hard attention [38] and self attention [39] *etc.* Soft attention can be seen as simulating the attention mechanism by multiplying weight on the neural unit so that the network pays more attention on the neural unit which multiplies with larger weight. SENet [37] takes advantage of soft attention mechanism to win the champion on the image classification task of ILSVRC-2017 [40]. Hard attention [38] can be seen as a module that decides a block region of the input image where is visible to the network, and the other region is invisible. Self-attention [39] improves the performance of the machine translation system by training a network to find the inner dependency of the input and that of the output. In this paper, we use soft attention mechanism as the meta-learner’s attention mechanism.

2.3. Unsupervised Representation Learning

Supervised learning is a data-hungry manner to train deep network. Considering this, several unsupervised learning approach [27–32] have been proposed. A well-known way is training a neural network to reconstruct the original input through an Encoder-Decoder architecture, such as Auto-Encoder [27], Variational Auto-Encoder (VAE) [28] and *etc.* Given partial masked images, Context Auto-Encoder [29] trains a network to reconstruct not only the visible but also the masked region of the image. Colorization [30] uses *Lab* images to train a network to generate the unseen *ab* channels from the input *L* channel. Based on Colorization, Split-Brain [32] trains two separated networks to separately generate the *ab* channels from the *L* channel and generate the *L* channel from the *ab* channels. Different

from these methods, DeepCluster [31] couples deep learning with Cluster algorithm [41, 42]. However, in real world, many unlabeled images containing complex semantic information and are not suitable to be categorized into specific clusters. Therefore, we consider there might be a limitation for DeepCluster and we utilize Split-Brain as the unsupervised learning method in URAML.

3. Method

3.1. Problem of Learning from Few-Data

Learning from few-data is extremely difficult for the deep learning model. One reason is that the original input data is commonly represented in a large dimension space. Usually, tens or hundreds of thousands of dimension space is required. For example, for the image classification task, the original image is commonly stored in a large dimensional space (dimension of an 224x224 RGB image is 150528). In such a large dimension space, it is difficult for a few samples of one category to accurately reflect the character of this category.

Humans learn new categories efficiently because they utilize prior-knowledge and attention mechanism in cognition [20, 21, 23, 24, 43–47]. Prior-knowledge facilitates human to express perceptual images into high-level representations or descriptions, and attention mechanism helps human to focus on critical components of the representations. In this way, humans reduce the dimension of images and maintain the discriminative components of the images, which alleviates human cognition load and facilitate humans to efficiently learn new categories.

Existing meta-learning approaches improve deep learning a lot in few-shot learning. However, they train the meta-learner to quickly fit few-shot learning tasks directly on the few original high dimensional input data and pay little attention to the importance of prior-knowledge and attention mechanism, leading to unsatisfactory performance. Besides, as introduced before, we propose that ignoring prior-knowledge and attention mechanism is also the possible reason for existing meta-learning approaches to be vulnerable to suffer from the TOF problem.

In this paper, inspired by human cognition and for addressing the problem existing meta-learning approaches expose, we propose three methods step-by-step: Attention based Meta-Learning (AML), Representation and Attention based Meta-Learning (RAML), Unsupervised Representation and Attention based Meta-Learning (URAML).

3.2. AML

AML equips the meta-learner with the power of attention mechanism. We first introduce the network structure and then detail the training of AML.

AML Network

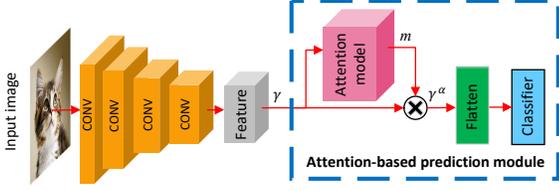


Figure 2: Network structure of the proposed method AML. There is an attention model inserted explicitly in the meta-learner’s network.

The network architecture of AML is shown in Fig.2. The network consists of a feature extractor and an attention-based prediction (ABP) module. The feature extractor is a CNN \mathcal{F} which is composed of four stacking convolutional layers. The ABP module contains a convolution-based attention model \mathcal{A} and a fully-connect layer based classifier \mathcal{C} . Eq.1 shows the inference of the network. θ_f , θ_a , and θ_c are weights of \mathcal{F} , \mathcal{A} , and \mathcal{C} , respectively. \mathcal{F} extracts features γ_i of the input image x_i and feed γ_i into the attention model \mathcal{A} . Then, \mathcal{A} calculates the soft attention mask m_i of the features γ_i . By channel-wise multiplication \odot between γ_i and m_i , the focused features γ_i^α is calculated. Finally, the classifier \mathcal{C} predicts the category of the input image, and \hat{y}_i is the corresponding prediction of x_i . We simplify and integrate the inference in Eq.1 as $\hat{y}_i = \mathbb{F}(x_i; \theta_f, \theta_a, \theta_c)$.

$$\begin{cases} \gamma_i = \mathcal{F}(x_i; \theta_f) \\ m_i = \mathcal{A}(\gamma_i; \theta_a) \\ \gamma_i^\alpha = \gamma_i \odot m_i \\ \hat{y}_i = \mathcal{C}(\gamma_i^\alpha; \theta_c) \end{cases} \quad (1)$$

In this paper, we use soft attention mechanism to build up the attention model. Although the soft attention mechanism is not exactly the same with the attention mechanism in human vision, it still plays a similar role with the human attention mechanism and helps the meta-learner to control its attention to key features. Fig.4(b) is used to better understand the soft attention processing of the meta-learner.

Fig.3 shows the attention model structure and Eq.2 shows the inference of the attention model. The input feature γ is firstly global-average-pooled to get feature γ' , and then a convolution layer coupled with a sigmoid activation layer are used to predict the attention mask m from the feature γ' .

$$\begin{cases} \gamma' = \mathcal{P}_a(\gamma), \\ m = \sigma(\mathcal{F}_a(\gamma'; \theta_a)) \end{cases} \quad (2)$$

\mathcal{P}_a is the global-average-pooling operation, and σ is the sigmoid activation, and \mathcal{F}_a is the convolution layer in the attention model.

AML Meta-Train Process

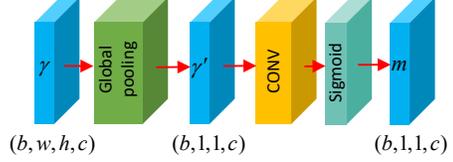


Figure 3: Inner network structure of the attention model. The shape of feature map γ is (b, w, h, c) which is shown at the left of the figure, where b , w , h , c are the batch size, width, height and umber of channels of the feature map γ , and the shape of γ' and m are both $(b, 1, 1, c)$.

Given a few-shot classification task τ , AML meta-trains the meta-learner to solve the task τ in the two steps. **First**, AML requires the meta-learner to inner-update itself on the the support set of τ , which can be formulated as Eq.3 and Eq.4.

$$\begin{cases} \hat{y}_i = \mathbb{F}(x_i; \theta_f, \theta_a, \theta_c), \\ \mathcal{L}_i(\theta_f, \theta_a, \theta_c) = l(\hat{y}_i, y_i), \\ \mathfrak{L}_s(\theta_f, \theta_a, \theta_c) = \frac{1}{N_s} \sum_{i=1}^{N_s} \mathcal{L}_i(\theta_f, \theta_a, \theta_c) \end{cases} \quad (3)$$

$$(\theta'_f, \theta'_a, \theta'_c) = (\theta_f, \theta_a, \theta_c) - \alpha \odot \nabla_{(\theta_f, \theta_a, \theta_c)} \mathfrak{L}_s(\theta_f, \theta_a, \theta_c) \quad (4)$$

In Eq.3, x_i is any image that belongs to the support set, l is the cross-entropy loss function, \mathcal{L}_i is the meta-learner’s loss on the image x_i , \mathfrak{L}_s is the meta-learner’s loss on the total support set, and N_s is the number of images in the support set. In Eq.4, inspired by Meta-SGD [12], we set α as a trainable vector which adjusts the inner-update direction and α has the same shape with the weights θ_f , θ_a , and θ_c . α can also be presented as $\alpha = [\alpha_f, \alpha_a, \alpha_c]$ and the Eq.4 can be split into three equations, *i.e.* $\theta'_f = \theta_f - \alpha_f \odot \nabla_{\theta_f} \mathfrak{L}_s(\theta_f, \theta_a, \theta_c)$ and *etc.*. For simplicity, we merge these three equations into one equation as Eq.4 shows. \odot is the element-wise multiplication. Supervised by the loss on the support set, the meta-learner inner-updates its weights $\theta_f, \theta_a, \theta_c$ to $\theta'_f, \theta'_a, \theta'_c$.

Second, as the inner-updated weight θ'_f, θ'_a , and θ'_c depend on not only the initial values of θ_f, θ_a , and θ_c , but also α , all $\theta_f, \theta_a, \theta_c$, and α can be meta-optimized. We formulate this process as Eq.5 and Eq.6.

$$\begin{cases} \hat{y}_i = \mathbb{F}(x_i; \theta'_f, \theta'_a, \theta'_c), \\ \mathcal{L}_i(\theta'_f, \theta'_a, \theta'_c) = l(\hat{y}_i, y_i), \\ \mathfrak{L}_q(\theta'_f, \theta'_a, \theta'_c) = \frac{1}{N_q} \sum_{i=1}^{N_q} \mathcal{L}_i(\theta'_f, \theta'_a, \theta'_c) \end{cases} \quad (5)$$

$$(\theta_f, \theta_a, \theta_c, \alpha) = (\theta_f, \theta_a, \theta_c, \alpha) - \beta \cdot \nabla_{(\theta_f, \theta_a, \theta_c, \alpha)} \mathfrak{L}_q(\theta'_f, \theta'_a, \theta'_c) \quad (6)$$

In Eq.5, x_i is an image belonging to the query set, and N_q denotes the number of images in the query set. \mathcal{L}_q is the inner-updated meta-learner’s loss on the query set. It should be noted that $\nabla_{(\theta_f, \theta_a, \theta_c, \alpha)} \mathcal{L}_q(\theta'_f, \theta'_a, \theta'_c)$ computes the gradient of \mathcal{L}_q towards $(\theta_f, \theta_a, \theta_c, \alpha)$ but not $(\theta'_f, \theta'_a, \theta'_c)$. By optimizing \mathcal{L}_q , the meta-learner is forced to learn not only the suitable initial weights $\theta_f, \theta_a, \theta_c$ but also α for task τ . With the learned initial weights and α , the meta-learner can inner-update itself precisely on the support set and then perform well on the query set.

In AML, the meta-learner is trained on lots of few-shot learning tasks with these two steps, which makes the meta-learner learn generalizable initial weights for not only the feature extractor \mathcal{F} and the classifier \mathcal{C} , but also the attention model \mathcal{A} . While existing initialization based meta-learning methods only train the meta-learner to learn initial weights for the feature extractor and the classifier. Therefore, compared with existing meta-learners, AML simplifies the few-shot problem and improves performance since its attention ability is meta-trained and can be easily adjusted to the crucial features for solving few-shot learning, which leads the classifier can make a precise prediction for the input. In our experiment, we show the positive effect of attention mechanism.

3.3. RAML

RAML assembles the meta-learner not only the attention mechanism but also the ability to well use the past learned knowledge.

Fig.4(a) shows the meta-learner’s network structure. Its network consists of a Representation module and an ABP module. The Representation module is different from the feature extractor in AML because the Representation module here is responsible for the meta-learner learning and leveraging prior-knowledge to understand the input image. While the feature extractor in AML is meta-trained for learning how to update itself for solving few-shot learning tasks. In our work, the Representation module is a ResNet-50 network. Similar to the ABP module in AML, the ABP module here also contains an attention model and a classifier. It is responsible for quickly adjusting the meta-learner’s attention and prediction based on the output feature from the Representation module. Besides, Fig.4(a) contains an Auxiliary module. The Auxiliary module does not belong to the meta-learner, and it is only used to assist the meta-learner learning prior-knowledge.

RAML Training Process

The training process of RAML can be separated into two stages: prior-knowledge learning and meta-training stage.

At the prior-knowledge learning stage, with the assist of the Auxiliary module, the Representation module is trained to learn prior-knowledge about image classification in a supervised manner. The training process can be formu-

lated as

$$\begin{cases} \gamma_i = \mathcal{F}_r(x_i; \theta_r) \\ \hat{y}_i = \mathcal{C}_{au}(\gamma_i; \theta_{au}) \\ L_{au} = \frac{1}{n} \sum_{i=1}^n l(\hat{y}_i, y_i) \\ \theta_r^*, \theta_{au}^* = \underset{\theta_r, \theta_{au}}{\operatorname{argmin}} L_{au} \end{cases} \quad (7)$$

\mathcal{F}_r and \mathcal{C}_{au} denote the Representation and Auxiliary modules, respectively, and θ_r and θ_{au} are their weights. x_i is an input image used for the representation model learning prior-knowledge, and n is the number of images. θ_r^* and θ_{au}^* are the learned values of θ_r and θ_{au} .

At the meta-training stage, for the meta-learner well using the learned knowledge to stably express the input image into high-level representations, the Representation module will not be meta-trained. Similar to AML, in RAML, we simplify the prediction of the meta-learner as $\hat{y}_i = \mathbb{F}(x_i; \theta_r^*, \theta_a, \theta_c)$, where all symbols denote the same meanings as those in AML. In RAML, the inner-update of the meta-learner on the support set can be formulated as Eq.8 and Eq.9. We can see that different from the inner-update of AML which update all weights of the network, the inner-update of RAML only update the weights θ_a and θ_c of the ABP module. The weight θ_r^* of the Representation module is fixed to keep the learned prior-knowledge.

$$\begin{cases} \hat{y}_i = \mathbb{F}(x_i; \theta_r^*, \theta_a, \theta_c), \\ \mathcal{L}_i(\theta_r^*, \theta_a, \theta_c) = l(\hat{y}_i, y_i), \\ \mathcal{L}_s(\theta_r^*, \theta_a, \theta_c) = \frac{1}{N_s} \sum_{i=1}^{N_s} \mathcal{L}_i(\theta_r^*, \theta_a, \theta_c) \end{cases} \quad (8)$$

$$(\theta'_a, \theta'_c) = (\theta_a, \theta_c) - \alpha \nabla_{(\theta_a, \theta_c)} \mathcal{L}_s(\theta_r^*, \theta_a, \theta_c) \quad (9)$$

The meta-optimizing in RAML can be formulated as Eq.10 and Eq.11.

$$\begin{cases} \hat{y}_i = \mathbb{F}(x_i; \theta_r^*, \theta'_a, \theta'_c), \\ \mathcal{L}_i(\theta_r^*, \theta'_a, \theta'_c) = l(\hat{y}_i, y_i), \\ \mathcal{L}_q(\theta_r^*, \theta'_a, \theta'_c) = \frac{1}{N_q} \sum_{i=1}^{N_q} \mathcal{L}_i(\theta_r^*, \theta'_a, \theta'_c) \end{cases} \quad (10)$$

$$(\theta_a, \theta_c, \alpha) = (\theta_a, \theta_c, \alpha) - \beta \cdot \nabla_{(\theta_a, \theta_c, \alpha)} \mathcal{L}_q(\theta_r^*, \theta'_a, \theta'_c) \quad (11)$$

The character of RAML is that the Representation module and the ABP module are trained separately. The Representation module is supervisorily trained to learn the prior-knowledge about image classification, and the ABP module is meta-trained to learn how to adjust itself quickly to solve few-shot learning tasks in the representation space provided by the Representation module. Compared with AML, which meta-trains the meta-learner not only adjusting the feature extractor but also the ABP module, RAML meta-trains the meta-learner simplify the few-shot learning

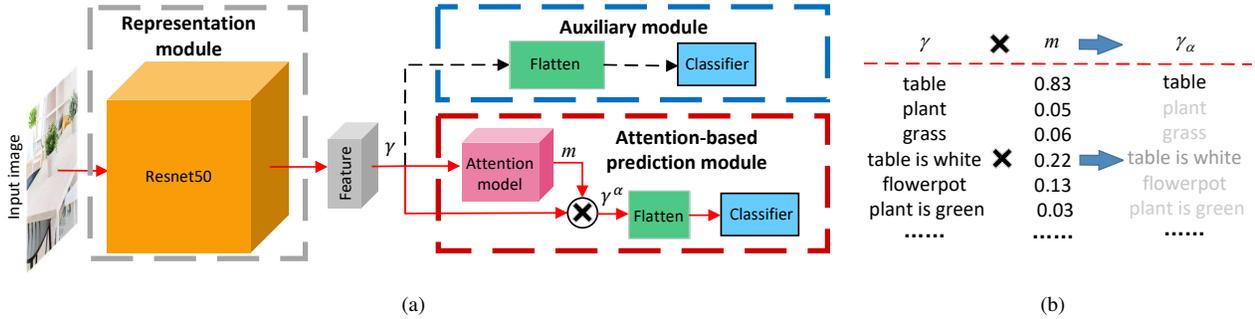


Figure 4: (a) Network structure of the proposed RAML. The meta-learner is composed of a Representation module and an ABP module. The Auxiliary module is used to assist the meta-learner to learn prior-knowledge. (b) An example that interprets the principle of soft attention mechanism for few-shot learning.

problem as the meta-learner only need to adjust its ABP module in the representation space. This is possibly the reason why RAML outperforms AML in our experiment.

3.4. URAML

The prior-knowledge can be learned on not only labeled data but also large-scale unlabelled data. Thus, we design the method URAML and show its network structure in Fig.5. Similar to RAML, the meta-learner is also composed of a Representation module and an ABP module, and the Auxiliary module does not belong to the meta-learner. The training process of URAML can be separated into two stages: prior-knowledge learning and meta-training stage.

At the prior-knowledge learning stage, the Representation module learns the knowledge with an unsupervised learning algorithm: Split-Brain auto-encoder [32]. The Split-Brain auto-encoder simultaneously trains two auto-encoders with *Lab* images. In *Lab* color system, the *L* channel determines the brightness of the image, and the *ab* channels determine the color. One auto-encoder in Split-Brain is trained to predict the unseen *ab* channels of the input *Lab* image, given only the *L* channel. Another is trained to predict the unseen *L* channel, given the *ab* channels. As Fig.5 shows, the Representation module consists of two ResNet-50 based encoders and the Auxiliary module consists of two corresponding deconvolution [48] based decoders. We formulate the prior-knowledge learning process as Eq.12 and Eq.13.

$$\begin{cases} \gamma_i^l = \mathcal{F}_l(x_i^l; \theta_l) \\ \hat{x}_i^{ab} = \mathcal{D}_l(\gamma_i^l; \omega_l) \\ L_l(\theta_l, \omega_l) = \frac{1}{n} \sum_{i=1}^n l_2(x_i^{ab}, \hat{x}_i^{ab}) \\ \theta_l^*, \omega_l^* = \underset{\theta_l, \omega_l}{\operatorname{argmin}} L_l(\theta_l, \omega_l) \end{cases} \quad (12)$$

In Eq.12, x_i^l and x_i^{ab} are the *L* and *ab* channels of the

input *Lab* image x_i , respectively. \mathcal{F}_l and \mathcal{D}_l are the encoder and decoder that predict x_i^{ab} based on x_i^l , respectively, and \hat{x}_i^{ab} is the prediction. θ_l and ω_l are the weights of \mathcal{F}_l and \mathcal{D}_l , respectively, and θ_l^* and ω_l^* are the optimized values of θ_l and ω_l . γ_i^l is the squeezed feature of x_i^l by the encoder \mathcal{F}_l . L_l is the loss of \mathcal{F}_l and \mathcal{D}_l , and l_2 is the *MSE* loss function. n is the number of *Lab* images that trains \mathcal{F}_l and \mathcal{D}_l . In Eq.13, all symbols are defined in the same way with those in Eq.12.

$$\begin{cases} \gamma_i^{ab} = \mathcal{F}_{ab}(x_i^{ab}; \theta_{ab}) \\ \hat{x}_i^l = \mathcal{D}_i^{ab}(\gamma_i^{ab}; \omega_{ab}) \\ L_{ab}(\theta_{ab}, \omega_{ab}) = \frac{1}{n} \sum_{i=1}^n l_2(x_i^l, \hat{x}_i^l) \\ \theta_{ab}^*, \omega_{ab}^* = \underset{\theta_{ab}, \omega_{ab}}{\operatorname{argmin}} L_{ab}(\theta_{ab}, \omega_{ab}) \end{cases} \quad (13)$$

After unsupervised learning, the representations γ_i of an *Lab* image x_i can be calculated by first concatenating γ_i^l with γ_i^{ab} and second average-pooling, which is shown as $\gamma_i = \mathcal{P}_a(\gamma_i^l, \gamma_i^{ab})$, where \mathcal{P}_a is an average-pooling layer.

At the meta-training stage, the ABP module is trained in the same way with that in RAML. Note that, the learned weight of the Representation module in URAML is $\theta_r^* = [\theta_l^*, \theta_{ab}^*]$.

At the end of our methodology, we summarize our three methods briefly. Inspired by human cognition which makes full use of attention mechanism and prior-knowledge to efficiently learn new knowledge, we design a novel paradigm with three methods to step-by-step utilize attention mechanism and prior-knowledge in meta-learning. Firstly, the method AML is designed to leverage attention mechanism in meta-learning. Secondly, the method RAML is designed to use not only the attention mechanism but also prior-knowledge in meta-learning. Compared with RAML, the method URAML learns the prior-knowledge with unsupervised learning, which brings URAML the advantage that with the growth of available unlabeled images used in the

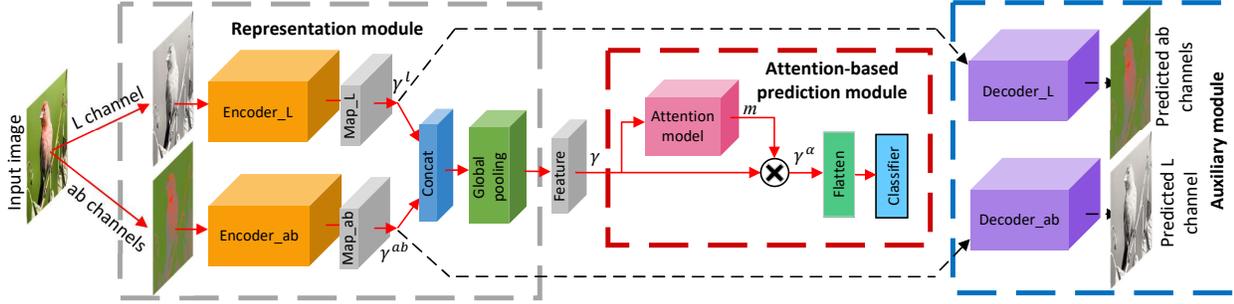


Figure 5: The network structure of URAML. The meta-learner is composed of a Representation module and an ABP module. The Auxiliary module is used to assist the meta-learner to learn prior-knowledge.

prior-knowledge learning stage and the progress of unsupervised learning algorithm, the performance of the meta-learner will be boosted up.

4. Experiments

In this section, we firstly present the datasets we used in our experiments, and then the details and results of our experiments.

4.1. Dataset

We use several datasets in all our experiments: MiniImagenet [13], Omniglot [49], MiniImagenet-900, Places2 [50], COCO [51], and OpenImages-300. Note that, we resize all the images in Omniglot into 28x28 resolution, and all the other images into 84x84.

4.1.1 MiniImagenet

MiniImagenet [13] is popularly used for evaluating few-shot learning and meta-learning. It contains 100 image classes, including 64 training classes, 16 validation classes, and 20 testing classes. Each image class with 600 images are sampled from the ImageNet dataset [52].

4.1.2 Omniglot

Omniglot [49] is another widely used dataset for meta-learning. It contains 50 different alphabets and 1623 characters from these alphabets, and each character has 20 images that hand-drawn by 20 different people.

4.1.3 MiniImagenet-900

MiniImagenet-900 dataset is designed for the Representation modules in RAML and URAML learning prior-knowledge, and it is composed of 900 image classes. Each image class with 1300 images are collected from the original ImageNet dataset. It is worth noting that there is no im-

age class in MiniImageNet-900 coincides with the classes in the MiniImagenet dataset.

4.1.4 Other Datasets

As the Representation module of URAML is trained by unsupervised learning, we take full advantage of this characteristic by training the Representation module of URAML on not only MiniImagenet-900 but also Places2 [50], COCO2017 [51], and OpenImages-300. The dataset OpenImages-300 is a subset of the OpenImages-V4 dataset [53]. The total OpenImages-V4 dataset contains 9 million images, and we randomly downloaded 3 million images from the OpenImages-V4 website to form the OpenImages-300 dataset.

4.2. Experiments on MiniImagenet

On MiniImagenet, we test all our methods on 5-way 1-shot and 5-way 5-shot classification tasks. The testing accuracy is averaged by the accuracies on 600 tasks, with 95% confidence intervals, and all these 600 tasks are randomly generated on the test set of MiniImagenet. The support and query set of each N -way K -shot task contains NK and $15 * N$ images, respectively.

In AML, the network structure of the meta-learner is shown in Fig.2. The feature extractor is composed of 4 Convolution layers and the classifier is a fully-connect layer, and the attention model structure is shown in Fig.3. Each Convolution layer consists of 64 channels and is followed with a ReLU and batch-normalization layer. We train the meta-learner on 200000 randomly generated tasks for 60000 iterations, and set the learning rate to 0.001, and decay the learning rate to 0.0001 after 30000 iterations. Moreover, Dropout with dropout-rate 0.2, L1 and L2 normalization with 0.001 and 0.00001, respectively, are used to prevent the meta-learner from over-fitting.

The experimental result of the method AML on MiniImagenet shows in Tab.2. Note that in Tab.2, the method whose name is printed as black uses a shallow network con-

Table 1: Few-shot learning performance on Omniglot. The method which is colored with blue uses deep network (ResNet) to extract image features, while the other use shallow network (4 cascading convolution layers). The accuracy is tested as the same way as MAML [11]

Method	Venue	5-way Accuracy		20-way Accuracy	
		1-shot	5-shot	1-shot	5-shot
MAML [11]	ICML-17	98.70±0.40%	99.90±0.10%	95.80±0.30%	98.90±0.20%
Prototypical Nets [6]	NIPS-17	98.80%	99.70%	96.00%	98.90%
Meta-SGD [12]	/	99.53±0.26%	99.93±0.09%	95.93±0.38%	98.97±0.19%
Relation Net [54]	CVPR-18	99.60±0.20%	99.80±0.10%	97.60±0.20%	99.10±0.10%
GNN [55]	ICLR-18	99.20%	99.70%	97.40%	99.00%
Spot-Learn [56]	CVPR-19	97.56±0.31%	99.65±0.06%	/	/
iMAML HF [35]	NIPS-19	99.50±0.26%	99.74±0.11%	96.18±0.36%	99.14±0.10%
SNAIL [15]	ICLR-18	99.07±0.16%	99.78±0.09%	97.64±0.30%	99.36±0.18%
MetaGAN+RN [18]	NIPS-18	99.67±0.18%	99.86±0.11%	97.64±0.17%	99.21±0.10%
AML(ours)	/	99.65±0.10%	99.85±0.04%	98.48±0.09%	99.55±0.06%

Table 2: Few-shot learning performance on MiniImagenet. The method which is colored with blue uses deep network to extract image features, while the other use shallow network. We separately highlight the best result of the methods using shallow network and that of the methods using deep network, for each task.

Method	Venue	5-way Accuracy	
		1-shot	5-shot
MAML [11]	ICML-17	48.70±1.84%	63.11±0.92%
Prototypical Nets [6]	NIPS-17	49.42±0.78%	68.20±0.66%
Meta-SGD [12]	/	50.47±1.87%	64.03±0.94%
LLAMA [33]	ICLR-18	49.40±1.83%	/
Relation Net [54]	CVPR-18	51.38±0.82%	67.07±0.69%
GNN [55]	ICLR-18	50.33±0.36%	66.41±0.63%
Spot-Learn [56]	CVPR-19	51.03±0.78%	67.96±0.71%
iMAML HF [35]	NIPS-19	49.30±1.88%	/
Meta-MinibatchProx [57]	NIPS-19	50.77±0.90%	67.43±0.89
AML(ours)	/	52.25±0.85%	69.46±0.68%
SNAIL [15]	ICLR-18	55.71±0.99%	68.88±0.92%
TADAM [58]	NIPS-18	58.50±0.30%	76.70±0.30%
MetaGAN+RN [18]	NIPS-18	52.71±0.64%	68.63±0.67%
AM3-TADAM [59]	ICLR-19	65.30±0.49%	78.10±0.36%
Incremental [60]	NIPS-19	54.95±0.30%	63.04±0.30%
RAML(ours)	/	63.66±0.85%	80.49±0.45%
URAML(ours)	/	49.56±0.79%	63.42±0.76%

sists of 4 or 5 Convolution layers and one or two fully-connect layers, and the method whose name is printed as blue uses a deep ResNet-based network. Among all the methods using shallow network, AML attained the state-of-the-art on both the 5-way 1-shot and 5-way 5-shot image classification tasks.

In RAML, the Representation module is a ResNet-50 [61] network, and the Auxiliary module is a fully-connect layer. The attention model is the same as that in AML, and the classifier is composed of two fully-connect layers.

At the prior-knowledge learning stage, we set the batch size to 256, and the learning rate to 0.001, and decay the learning rate to 0.0001 after 30000 iterations, and use L2

normalization with 0.00001 and Dropout with 0.2 to prevent the Representation module from over-fitting. At the meta-training stage, the ABP module is meta-trained with the same setting as AML. The experiment result of RAML is shown in Tab.2. Compared to method AML, RAML improves the meta-learner’s performance more significantly. It rises the accuracy on 5-way 1-shot tasks from 52.25% to 63.66%, and the accuracy on 5-way 5-shot tasks from 69.46% to 80.49%.

The most likely reason why RAML performs well is: before the meta-training stage, the Representation module has learned old knowledge to help the meta-learner understanding new input image and provides high-level meaningful representations and features of the input image. In the meta-training stage, the meta-learner’s work becomes more comfortable because it only needs to learn how to quickly adjust its ABP module according to the compact features the Representation module provided, and do not need to take care of the original high dimensional input data. While the meta-learner of AML works harder than the meta-learner of RAML, as it has to adjust its total network to fit new few-shot learning tasks according to the original input data.

In URAML, the Representation module learns the prior-knowledge with an unsupervised learning algorithm: Split-Brain. As Fig.5 shows, two independent ResNet-50 network-based encoders compose the Representation module, and we halve all the filters in each encoder so that the Representation module outputs feature vector with a dimension of 2048, which is the same with that in RAML. The Auxiliary module is composed of two deconvolution-based decoders, and Tab.4 shows the detail of the decoder network structure. The last Conv-layer’s number of filters is 1 or 2 according to that the decoder is recovering the L channel or the ab channels of the Lab image.

At both the prior-knowledge learning and meta-training stage, we set all hyperparameters the same with those in the

Table 3: Ablation experimental results about the attention mechanism on Omniglot.

Method	5-way Accuracy		20-way Accuracy	
	1-shot	5-shot	1-shot	5-shot
MAML*	97.40±0.27%	99.71±0.05%	93.37±0.23%	97.46±0.11%
MAML+attention	97.41±0.28%	99.48±0.12%	92.99±0.25%	97.94±0.10%
Meta-SGD*	98.94±0.17%	99.51±0.07%	95.82±0.21%	98.40±0.09%
Meta-SGD+attention	99.26±0.15%	99.79±0.04%	97.94±0.14%	98.99±0.10%

Table 4: Detailed structure of the decoder module in URAML.

Layers	Number of filters	Kernel
CONV	1024	5
DeCONV	512	3
DeCONV	256	3
CONV	1 or 2	1

Table 5: Ablation experimental results about the attention mechanism on MiniImagenet

Method	5-way Accuracy	
	1-shot	5-shot
MAML*	48.03±0.83%	64.11±0.73%
MAML+attention	48.52±0.85%	64.94±0.69%
Reptile*	48.23±0.43%	63.69±0.49%
Reptile+attention	48.30±0.45%	64.22±0.39%
Meta-SGD*	48.15±0.93%	63.73±0.85%
Meta-SGD+attention	49.11±0.94%	65.54±0.84%

RAML experiment. Noted that for saving the training computation cost, the decoders in the Auxiliary module recover the ab and L channels into 11×11 resolution, but not the original 84×84 . When calculating the MSE losses $L_l(\theta_l, \omega_l)$ and $L_{ab}(\theta_{ab}, \omega_{ab})$ shown in Eq.12 and Eq.13, we first resize ab and L channels of the input Lab image into 11×11 resolution and then calculate $L_l(\theta_l, \omega_l)$ and $L_{ab}(\theta_{ab}, \omega_{ab})$. The experiment result of URAML is shown in Tab.2. We also highlight the result of URAML in Tab.2, even though its result is not state-of-the-art. In our viewpoint, the reason why URAML lags behind RAML is that the Representation module in URAML learns the prior-knowledge with unsupervised learning while the Representation module in RAML learns with supervised learning.

4.3. Experiments on Omniglot

As Omniglot is a much easier dataset than MiniImagenet that existing meta-learners can easily achieve more than 95% accuracy on most testing tasks generated on Omniglot, we only test method AML on Omniglot.

Same to the experiments on Miniimagenet, we also train the meta-learner on 200000 randomly generated tasks for 60000 iterations and set the learning rate to 0.001. The ex-

periment results are shown in Tab.1

It is clear that the proposed method AML attains state-of-the-art performance on 2 of all 4 kinds of few-shot image classification tasks. On the 5-way 1-shot task, though the method MetaGAN+RN performs slightly better than AML, we still highlight AML as MetaGAN+RN uses a deeper ResNet-based network while AML uses a shallower network. On the 20-way 1-shot task, our method AML surpasses other methods by a large margin. For example, compared to IMAML HF, AML improves the meta-learner’s performance from 96.18% to 98.48%.

4.4. Ablation Study

4.4.1 Ablation Study about the Attention Mechanism

To confirm the promotion effect of the attention mechanism for meta-learning, we conduct experiments to compare the performance of the meta-learner which is equipped with the attention model and its counterpart which is not. The experimental results show in Tab.5 and Tab.3. The compared meta-learner which is marked with * is the meta-learner re-implemented by ourselves. The performances of our re-implemented meta-learners differ slightly from those reported in their original papers. This is probably caused by different hyper-parameters or experiment settings (all methods in this experiment use convolution layers with 32 filters). The comparisons in Tab.5 and Tab.3 revealing that in most cases, the attention mechanism improves the meta-learner significantly, which demonstrates the reason-ability of our idea.

As attention mechanism brings the meta-learner more weights and computation cost, we do another experiment to validate that the improvement of AML is the contribution of the attention mechanism but not the growth of the number of weights and computation cost. The experiment detail is: since the attention model in AML is a convolution layer with the kernel size of 1×1 , we remove the attention model, and stack a convolution layer with the same kernel size on the top of the CNN feature extractor. We name the meta-learner with this network as AML-attention, and its number of weight is the same as that of AML. The corresponding experimental result is shown in Tab.6, and it is clear that AML outperforms AML-attention, which further shows the improvement effect of attention mechanism for meta-learning.

Table 6: Results of several ablation experiments.

Method	5-way Accuracy	
	1-shot	5-shot
AML	52.25±0.85%	69.46±0.68%
AML-attention	51.27±0.78%	67.73±0.65%
RAML	63.66±0.85%	80.49±0.45%
RAML-Places2	58.82±0.89%	74.09±0.76%

4.4.2 Prior-Knowledge Learning Dataset

We do experiments to test how does the prior-knowledge learning dataset affects RAML and URAML.

a) affects to RAML: In RAML, the default prior-knowledge learning dataset is our reorganized *Miniimagenet-900* dataset. In this experiment, the Representation module learns the prior-knowledge on Places2 [50] instead of *Miniimagenet-900*, and all the other experiment settings and hyper-parameters are constant with the primordial RAML. We denote this meta-learner as RAML-Places2. Corresponding experimental result shows in Tab.6. It is clear that prior-knowledge learning dataset affects the meta-learner. The reason is that different prior-knowledge learning dataset leads the Representation module learning different knowledge and expressing image features differently. Places2 is a dataset commonly used for scene classification, which results in that the Representation module learning the knowledge about scene understanding rather than object classification.

b) affects to URAML: In this experiment, we test how the quantity of unlabeled *Lab* images in the prior-knowledge learning dataset affect URAML. We design two new versions of URAML: URAML-V1 and URAML-V2. The Representation module of URAML-V1 learns prior-knowledge only on MiniImagenet-900, and that of URAML-V2 learns prior-knowledge on not only MiniImagenet-900, but also the Places2 and COCO2017. Compared with URAML-V1 and URAML-V2, the quantity of unlabeled *Lab* used in the primordial URAML is the largest, as MiniImagenet-900, places365, COCO2017, and OpenImages-300 are all used in the primordial URAML. Tab.7 shows the prior-knowledge learning dataset and the performances of URAML-V1, URAML-V2, and the primordial URAML. It is clear that the primordial URAML performs the best, and the more the unlabeled *Lab* images used for the meta-learner to learn prior-knowledge, the better the meta-learner performs. Besides, there remains a large performance progress space as we can use more unlabeled data in URAML.

4.4.3 Unsupervised Learning for URAML

The development of unsupervised learning also affects URAML a lot. To verify this viewpoint, we do an experi-

ment that the Representation module in URAML learns the prior-knowledge with a basic unsupervised learning method Auto-Encoder [27], and we name this version of URAML as URAML-AE. The experimental result of URAML-AE shown in Tab.7 revealing that the unsupervised learning algorithm affects the meta-learner significantly. Maybe the most promising way to improve the performance of URAML is to develop the unsupervised learning algorithm and collect more unlabeled data.

4.5. Cross-Testing Experiment

We find that existing meta-learning methods generally suffer from a Task-Over-Fitting (TOF) problem, and this problem has seldom been studied. An example of the TOF problem is that the meta-learner to be tested on 5-way 1-shot classification tasks should be trained on 5-way 1-shot tasks rather than on other tasks, and similarly, the meta-learner to be tested on 5-way 5-shot tasks should be trained on 5-way 5-shot tasks. This is because the meta-learner trained on 5-shot tasks over-fits to 5-shot tasks, and when testing it on 1-shot tasks, it will perform obviously worse than the meta-learner trained on 1-shot tasks.

We do lots of cross-testing experiments to test how much does MAML, Meta-SGD, AML, RAML, and URAML suffer from the TOF problem, and the experimental results show that compared with the other methods, our methods suffer less from this problem, especially RAML and URAML.

For each tested meta-learning method, we do the cross-testing experiments in the following way: 1) train the meta-learner on 5-way K -shot image classification tasks, where $K \in \{1,3,5,7,9\}$, 2) test the meta-learner on 5-way J -shot tasks, where $J \in \{1,3,5,7,9\}$. For example, we train a meta-learner with MAML on 5-way 3-shot tasks and test its performance on all 5-way K -shot tasks, $K \in \{1,3,5,7,9\}$. The experimental results are shown in Fig.6.

Obviously, Fig.6 shows that MAML suffers seriously from the TOF problem, because its meta-learner which performs best on K -shot tasks probably performs not well on J -shot tasks, where $K \neq J$. For example, in MAML, the meta-learner trained on 1-shot tasks performs best on the 1-shot tasks, but it can not perform as well as the other meta-learners on 3-, 5-, 7-, and 9-shot tasks, which means the meta-learner trained on 1-shot tasks over-fits to 1-shot tasks. The meta-learner trained by URAML troubled little by the TOF problem because the meta-learner which performs best on K -shot tasks probably performs best on J -shot tasks, where $K, J \in \{1,5,7,9\}$. For example, in URAML, the meta-learner trained on 1-shot tasks performs best not only on the 1-shot tasks but also on 5-, 7-, and 9-shot tasks, which means the meta-learner trained on 1-shot tasks generalizes well to the other J -shot tasks.

We design a metric Cross-Entropy across Tasks (CET),

Table 7: Ablation experimental results about URAML.

Method	Dataset	Number of images	5-way Accuracy	
			1-shot	5-shot
URAML-V1	MiniImagenet-900	1.15million	45.91±0.79%	61.04±0.71%
URAML-V2	MiniImagenet-900, places365, COCO2017	4.10million	48.82±0.79%	62.84±0.78%
URAML-AE	MiniImagenet-900, places365, COCO2017, OpenImages-300	7.10million	33.29±0.71%	43.60±0.66%
URAML	MiniImagenet-900, places365, COCO2017, OpenImages-300	7.10million	49.56±0.79%	63.42±0.76%

testing task \ training task	Meta-SGD					MAML					AML					RAML					URAML				
	1-shot	3-shot	5-shot	7-shot	9-shot	1-shot	3-shot	5-shot	7-shot	9-shot	1-shot	3-shot	5-shot	7-shot	9-shot	1-shot	3-shot	5-shot	7-shot	9-shot	1-shot	3-shot	5-shot	7-shot	9-shot
1-shot	48.15	47.71	46.44	47.67	46.38	48.03	45.56	40.54	41.07	39.69	52.25	51.58	51.79	51.66	51.03	63.66	63.54	63.14	63.46	63.48	49.56	48.04	47.89	47.43	46.26
3-shot	58.24	59.18	58.90	58.75	59.15	58.49	60.12	59.12	59.99	59.43	63.23	64.97	65.19	65.02	64.56	74.73	76.60	76.18	76.58	76.78	59.03	58.48	58.90	59.32	58.21
5-shot	62.85	63.56	63.73	63.93	63.79	62.06	64.18	64.11	64.32	64.31	67.64	68.82	69.46	69.32	69.08	78.10	80.15	80.49	80.17	80.29	63.64	62.96	63.42	62.04	62.22
7-shot	65.10	65.79	66.07	65.95	65.64	64.52	66.85	67.42	67.33	67.19	69.51	71.35	71.60	72.57	72.68	79.93	82.19	82.62	82.28	82.31	65.52	64.54	63.77	64.60	64.06
9-shot	66.25	67.04	67.36	67.53	67.73	65.06	68.44	69.01	69.49	69.31	71.32	73.16	73.43	73.76	73.12	80.83	83.69	83.46	83.51	83.72	66.82	66.37	65.87	65.85	63.30

Figure 6: Results of cross-testing experiments among MAML, Meta-SGD, AML, RAML and URAML. Each column presents a meta-learner trained on specific K -shot training tasks and each row presents specific J -shot testing tasks, where $K, J \in \{1, 3, 5, 7, 9\}$. For each method, the value at J -shot row K -shot column presents the J -shot testing accuracy of the meta-learner trained on K -shot training tasks. For example, the value 59.99 at 3-shot row 7-shot column of MAML presents the 3-shot testing accuracy of the MAML meta-learner trained on 7-shot training tasks. The value 80.83 at 9-shot row 1-shot column of RAML presents the 9-shot testing accuracy of the RAML meta-learner trained on 1-shot training tasks.

to quantize how much does a meta-learning approach be vulnerable to the TOF problem. The evaluation process is shown as Eq.14, where $i, j \in \{1, 3, 5, 7, 9\}$ and overstriking variables are vector. \mathcal{S} and \mathcal{D} are the softmax and cross-entropy operation. \mathbf{a}_i is the testing accuracies of five meta-learners trained on 1-, 3-, 5-, 7-, 9-shot tasks when they are tested on i -shot tasks. \mathbf{d}_i is the meta-learners' accuracy distribution on i -shot tasks. $l_{i,j}$ presents the similarity between accuracy distribution vector \mathbf{d}_i and \mathbf{d}_j , where $i, j \in \{1, 3, 5, 7, 9\}$. L presents the overall similarities of $l_{i,j}$ for a specific approach.

$$\begin{cases} \mathbf{d}_i = \mathcal{S}(\mathbf{a}_i / \max(\mathbf{a}_i)) \\ l_{ij} = \mathcal{D}(\mathbf{d}_i, \mathbf{d}_j) \\ L = \sum_{i,j \in \{1,3,5,7,9\}} l_{ij} \end{cases} \quad (14)$$

For example, the testing accuracies \mathbf{a}_3 of Meta-SGD [58.24%, 59.18%, 58.90%, 58.75%, 59.15%] is the five trained meta-learners of Meta-SGD when they are tested on 3-shot tasks. So, $\mathbf{a}_3 / \max(\mathbf{a}_3) = [58.24\%, 59.18\%, 58.90\%, 58.75\%, 59.15\%] / 59.18\%$, and $\mathbf{d}_3 = \mathcal{S}(\mathbf{a}_3 / \max(\mathbf{a}_3)) = [0.116, 0.255, 0.202, 0.178, 0.249]$. Similarly, $\mathbf{d}_7 = [0.122, 0.206, 0.255, 0.233, 0.184]$. Then, $l_{3,7} = 1.603$, and $L = 34.22$.

Obviously, the smaller the total distance L appears, the less the meta-learning approach suffers from the TOF problem. We show different meta-learning approaches' performance on the CET metric in Tab.8. This experiment shows that the proposed AML, RAML, and URAML performs better than MAML and Meta-SGD on the CET metric, and RAML and URAML performs best. The possible reason for this is that prior-knowledge and attention mechanism are both helpful for the meta-learner to reduce its few-shot cognitive load and to avoid itself be affected by redundant

Table 8: Performance of different meta-learning methods on the CET metric.

Method	MAML	Meta-SGD	AML	RAML	URAML
CET	57.19	34.22	33.35	32.13	32.16

useless information.

We can see an interesting phenomenon in Fig.6, that the meta-learner trained by RAML on 5-way 9shot tasks performs best in most of the test tasks, while the meta-learner trained by URAML on 5-way 1-shot tasks performs best. The possible reason behind this phenomenon is that the Representation module of RAML learns knowledge by supervised learning, while the Representation module of URAML learns knowledge by unsupervised learning, which results in the output features between these two kinds of Representation module be different.

4.6. Feature Analysis

To understand the effect of attention mechanism, we visualize the distributions of feature γ and γ^α (shown in Fig.2, Fig.4(a) and Fig.5) in Fig.7 with t-SNE [62]. In Fig.7, 500 feature points of each picture represent 500 γ or γ^α of the query set images of a 5-way 1 or 5 shot task that randomly generated on the test set of MiniImagenet.

The average distribution inner-class distance D1 of γ^α is smaller than that of γ , and the average inter-class distance D2 of γ^α is larger than that of γ . This result indicates that among different image classes, the distribution of γ^α is more distinguishable than that of γ . The reason for this is that the attention mechanism makes the meta-learner be able to adjust its attention quickly to critical image features and makes γ^α more distinguishable than γ to differentiate

images of different classes.

4.7. Heat-Map of γ and γ^α

To further analyze how the attention mechanism affects the meta-learner, we visualize the heat-maps of γ and γ^α in Fig.8. To get the heat-map of γ , we first inner-update the RAML meta-learner on the support set of a randomly generated 5-way 1-shot testing task on MiniImagenet. Then, we feed the meta-learner with the query set images and average the feature maps γ across the channel axis to get the heat-maps of γ . Similarly, the heat-maps of γ^α can be got.

From the heat-maps shown in Fig.8, we can see that compared with γ , γ^α is more sensitive to the distinguishable part of the input image, revealing that the meta-learner changes its attention to the most discriminative image feature. For example, the first column of Fig.8 is a fish. Besides the fish body, γ is also sensitive to some background region of the image. However, the meta-learner discovers that only the fish body is the crucial feature to category this image and shrinks its attention region so that γ^α sensitive only to the fish body.

Through the visualization and analysis of the heat-map of γ and γ^α , we can see that the attention mechanism helps the meta-learner to focus on the most distinguishable image feature, and further helps the meta-learner to do a better few-shot learning task.

5. Conclusion and the Future Work

In this paper, be inspired by human cognition and learning process, we find the importance of attention mechanism and the prior-knowledge for meta-learning based few-shot learning. To solve a few-shot learning task, the meta-learner should first well use stable prior-knowledge to understand images and extract compact feature representations of images so that it can solve the task in the compact representation space rather than the original image space. Then, the meta-learner should adjust its attention to the crucial feature of the extracted feature representations, and make the final decision based on its attention. Therefore, we step-by-step propose three methods AML, RAML, and URAML to introduce attention mechanism and the prior-knowledge to meta-learning. All of them work successfully with state-of-the-art performance on several few-shot learning benchmarks, which indicating the rationality of our viewpoints and methods.

Besides, we find existing meta-learning approaches suffer from the TOF problem, which is unfriendly to practical applications. We design a novel Cross-Entropy across Tasks (CET) metric to evaluate how much does a meta-learning suffers from TOF. The experiment shows that compared to existing meta-learning methods, the proposed methods suffer less from the TOF problem, especially the RAML and URAML methods.

Among all the proposed methods, though URAML performs not the best, we think it is the most promising method yet because there is a large development space for the performance of URAML method which will also be the direction of our future work. From the ablation study, two manners seem can improve the performance of URAML significantly. One is to develop the unsupervised learning algorithm or self-supervised learning. RAML performs better than URAML revealing that the current unsupervised learning algorithm falls behind supervised learning. Bridging the gap between unsupervised learning and supervised learning algorithms will boost up the performance of URAML in a substantial probability. The other manner is to use more unlabeled data for URAML to learn prior-knowledge. Although 7.1 million unlabeled images are used in URAML, it still dramatically falls behind the images that humans have ever seen in terms of both quantity and quality. As for the quantity, we assume that, if a person watches 1 image per second and keep watching 15 hours per day, he/she can see 100 million images in 5 years. As for quality, humans see the world in a multimodal way, that is, the human can not only see the object but also touch and move around the object, which helps humans understand the world more accurately than Computer Vision. In a word, developing the unsupervised or self-supervised learning algorithm and collecting more unlabeled images will both help URAML to perform well.

Acknowledgements

This work is supported by the National Natural Science Foundation of China (No. 61573286).

References

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [2] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning." in *AAAI*, vol. 4, 2017, p. 12.
- [3] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks." in *CVPR*, vol. 1, no. 2, 2017, p. 3.
- [4] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [5] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," in

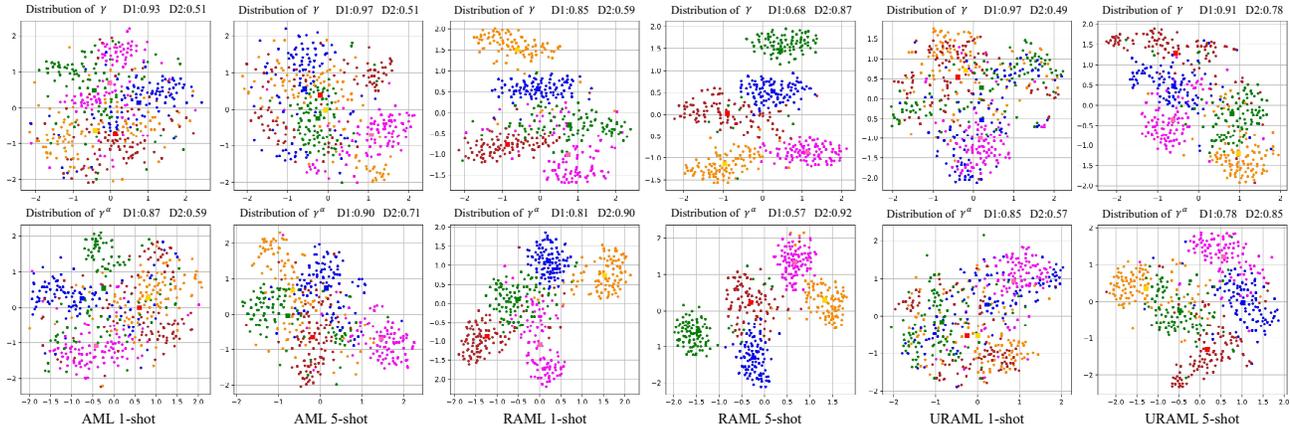


Figure 7: Visualization of the distributions of features γ and γ^α of all our three methods. For each method, we randomly generate a 5-way 1-shot and a 5-way 5-shot testing task on Miniimagenet, and the query set of each task contains 100 images for each image class. For each testing task, after the meta-learner inner-updating on the support set, we use t-SNE to visualize the distributions of the meta-learner’s γ and γ^α of the query set images. For each picture, five colors are used to represents 5 image classes in the testing task and each point denotes the feature γ or γ^α . We also show the average distribution inner-class distance D1 and inter-class distance D2 above each picture to better understand the distributions.

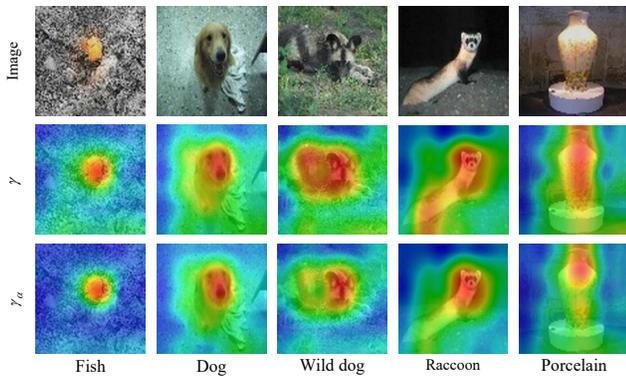


Figure 8: We show some images which are sampled from the query set of a 5-way 1-shot classification task, and the corresponding heat-maps of γ and γ^α .

Advances in Neural Information Processing Systems, 2016, pp. 3630–3638.

- [6] J. Snell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” in *Advances in Neural Information Processing Systems*, 2017, pp. 4077–4087.
- [7] W. Li, L. Wang, J. Xu, J. Huo, Y. Gao, and J. Luo, “Revisiting local descriptor based image-to-class measure for few-shot learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7260–7268.

- [8] Y. Bengio, S. Bengio, and J. Cloutier, *Learning a synaptic learning rule*. Université de Montréal, Département d’informatique et de recherche opérationnelle, 1990.
- [9] S. Bengio, Y. Bengio, J. Cloutier, and J. Gecsei, “On the optimization of a synaptic learning rule,” in *Preprints Conf. Optimality in Artificial and Biological Neural Networks*. Univ. of Texas, 1992, pp. 6–8.
- [10] J. Schmidhuber, “Learning to control fast-weight memories: An alternative to dynamic recurrent networks,” *Neural Computation*, vol. 4, no. 1, pp. 131–139, 1992.
- [11] C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep networks,” *international conference on machine learning*, 2017.
- [12] 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.
- [13] S. Ravi and H. Larochelle, “Optimization as a model for few-shot learning,” 2016.
- [14] A. Nichol, J. Achiam, and J. Schulman, “On first-order meta-learning algorithms,” 2018.
- [15] N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, “A simple neural attentive meta-learner,” 2018.
- [16] S. Gidaris and N. Komodakis, “Dynamic few-shot visual learning without forgetting,” in *Proceedings of*

- the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4367–4375.
- [17] Y. Duan, J. Schulman, X. Chen, P. L. Bartlett, I. Sutskever, and P. Abbeel, “RI²: Fast reinforcement learning via slow reinforcement learning,” *arXiv preprint arXiv:1611.02779*, 2016.
- [18] R. Zhang, T. Che, Z. Ghahramani, Y. Bengio, and Y. Song, “Metagan: An adversarial approach to few-shot learning,” in *Advances in Neural Information Processing Systems*, 2018, pp. 2367–2376.
- [19] T. Munkhdalai and H. Yu, “Meta networks,” *arXiv preprint arXiv:1703.00837*, 2017.
- [20] J. A. Langer and M. Nicolich, “Prior knowledge and its relationship to comprehension,” *Journal of Reading Behavior*, vol. 13, no. 4, pp. 373–379, 1981.
- [21] F. Dochy, “Instructional implications of recent research and empirically-based theories on the effect of prior knowledge on learning,” in *Learning Environments*. Springer, 1990, pp. 339–355.
- [22] A. M. Shapiro, “How including prior knowledge as a subject variable may change outcomes of learning research,” *American Educational Research Journal*, vol. 41, no. 1, pp. 159–189, 2004.
- [23] J. Wylie and C. McGuinness, “The interactive effects of prior knowledge and text structure on memory for cognitive psychology texts,” *British Journal of Educational Psychology*, vol. 74, no. 4, pp. 497–514, 2004.
- [24] I. Hsin and F. Paas, “Effects of computer-based visual representation on mathematics learning and cognitive load,” *Journal of Educational Technology & Society*, vol. 18, no. 4, pp. 70–77, 2015.
- [25] G. Logan, D. Dagenbach, and T. Carr, “Inhibitory processes in attention, memory and language,” *Academic Press, San Diego*, pp. 189–239, 1994.
- [26] S. A. Hillyard, E. K. Vogel, and S. J. Luck, “Sensory gain control (amplification) as a mechanism of selective attention: electrophysiological and neuroimaging evidence,” *Philosophical Transactions of the Royal Society of London B: Biological Sciences*, vol. 353, no. 1373, pp. 1257–1270, 1998.
- [27] G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [28] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [29] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, “Context encoders: Feature learning by inpainting,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2536–2544.
- [30] R. Zhang, P. Isola, and A. A. Efros, “Colorful image colorization,” in *European Conference on Computer Vision*. Springer, 2016, pp. 649–666.
- [31] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, “Deep clustering for unsupervised learning of visual features,” *arXiv preprint arXiv:1807.05520*, vol. 3, 2018.
- [32] R. Zhang, P. Isola, and A. A. Efros, “Split-brain autoencoders: Unsupervised learning by cross-channel prediction,” in *CVPR*, vol. 1, no. 2, 2017, p. 5.
- [33] E. Grant, C. Finn, S. Levine, T. Darrell, and T. Griffiths, “Recasting gradient-based meta-learning as hierarchical bayes,” *arXiv preprint arXiv:1801.08930*, 2018.
- [34] Q. Sun, Y. Liu, T. Chua, and B. Schiele, “Meta-transfer learning for few-shot learning,” pp. 403–412, 2019.
- [35] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine, “Meta-learning with implicit gradients,” in *Advances in Neural Information Processing Systems*, 2019, pp. 113–124.
- [36] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, “Residual attention network for image classification,” *arXiv preprint arXiv:1704.06904*, 2017.
- [37] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” *arXiv preprint arXiv:1709.01507*, vol. 7, 2017.
- [38] V. Mnih, N. Heess, A. Graves *et al.*, “Recurrent models of visual attention,” in *Advances in neural information processing systems*, 2014, pp. 2204–2212.
- [39] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [40] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

- [41] L. Xu, J. Neufeld, B. Larson, and D. Schuurmans, “Maximum margin clustering,” in *Advances in neural information processing systems*, 2005, pp. 1537–1544.
- [42] F. Lin and W. W. Cohen, “Power iteration clustering,” in *ICML*, vol. 10. Citeseer, 2010, pp. 655–662.
- [43] T. Ormerod, “Human cognition and programming,” in *Psychology of programming*. Elsevier, 1990, pp. 63–82.
- [44] A. Oliva, A. Torralba, M. S. Castelhana, and J. M. Henderson, “Top-down control of visual attention in object detection,” in *Proceedings 2003 International Conference on Image Processing (Cat. No. 03CH37429)*, vol. 1. IEEE, 2003, pp. I–253.
- [45] A. H. van der Heijden, *Selective attention in vision*. Routledge, 2003.
- [46] M. I. Posner and S. E. Petersen, “The attention system of the human brain,” *Annual review of neuroscience*, vol. 13, no. 1, pp. 25–42, 1990.
- [47] S. K. Ungerleider and L. G., “Mechanisms of visual attention in the human cortex,” *Annual review of neuroscience*, vol. 23, no. 1, pp. 315–341, 2000.
- [48] W. Shi, J. Caballero, L. Theis, F. Huszar, A. Aitken, C. Ledig, and Z. Wang, “Is the deconvolution layer the same as a convolutional layer?” *arXiv preprint arXiv:1609.07009*, 2016.
- [49] B. Lake, R. Salakhutdinov, J. Gross, and J. Tenenbaum, “One shot learning of simple visual concepts,” in *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 33, no. 33, 2011.
- [50] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, “Places: A 10 million image database for scene recognition,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [51] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [52] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. Ieee, 2009, pp. 248–255.
- [53] I. Krasin, T. Duerig, N. Alldrin, V. Ferrari, S. Abu-El-Haija, A. Kuznetsova, H. Rom, J. Uijlings, S. Popov, S. Kamali, M. Mallocci, J. Pont-Tuset, A. Veit, S. Belongie, V. Gomes, A. Gupta, C. Sun, G. Chechik, D. Cai, Z. Feng, D. Narayanan, and K. Murphy, “Openimages: A public dataset for large-scale multi-label and multi-class image classification.” *Dataset available from <https://storage.googleapis.com/openimages/web/index.html>*, 2017.
- [54] F. S. Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, “Learning to compare: Relation network for few-shot learning,” in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Salt Lake City, UT, USA*, 2018.
- [55] V. Garcia and J. Bruna, “Few-shot learning with graph neural networks,” *arXiv preprint arXiv:1711.04043*, 2017.
- [56] W.-H. Chu, Y.-J. Li, J.-C. Chang, and Y.-C. F. Wang, “Spot and learn: A maximum-entropy patch sampler for few-shot image classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6251–6260.
- [57] P. Zhou, X. Yuan, H. Xu, S. Yan, and J. Feng, “Efficient meta learning via minibatch proximal update,” in *Advances in Neural Information Processing Systems*, 2019, pp. 1532–1542.
- [58] B. N. Oreshkin, A. Lacoste, and P. Rodriguez, “Tadam: Task dependent adaptive metric for improved few-shot learning,” *arXiv preprint arXiv:1805.10123*, 2018.
- [59] C. Xing, N. Rostamzadeh, B. N. Oreshkin, and P. O. Pinheiro, “Adaptive cross-modal few-shot learning,” *arXiv preprint arXiv:1902.07104*, 2019.
- [60] M. Ren, R. Liao, E. Fetaya, and R. Zemel, “Incremental few-shot learning with attention attractor networks,” in *Advances in Neural Information Processing Systems*, 2019, pp. 5276–5286.
- [61] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [62] L. v. d. Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.