

META-LEARNING NEURAL PROCEDURAL BIASES

Anonymous authors

Paper under double-blind review

ABSTRACT

The goal of few-shot learning is to generalize and achieve high performance on new unseen learning tasks, where each task has only a limited number of examples available. Gradient-based meta-learning attempts to address this challenging task by learning how to learn new tasks by embedding inductive biases informed by prior learning experiences into the components of the learning algorithm. In this work, we build upon prior research and propose *Neural Procedural Bias Meta-Learning* (NPBML), a novel framework designed to meta-learn task-adaptive procedural biases. Our approach aims to consolidate recent advancements in meta-learned initializations, optimizers, and loss functions by learning them simultaneously and making them adapt to each individual task to maximize the strength of the learned inductive biases. This imbues each learning task with a unique set of procedural biases which is specifically designed and selected to attain strong learning performance in only a few gradient steps. The experimental results show that by meta-learning the procedural biases of a neural network, we can induce strong inductive biases towards a distribution of learning tasks, enabling robust learning performance across many well-established few-shot learning benchmarks.

1 INTRODUCTION

Humans have an exceptional ability to learn new tasks from only a few example instances. We can often quickly adapt to new domains effectively by building upon and utilizing past experiences of related tasks, leveraging only a small amount of information about the target domain. The field of meta-learning (Schmidhuber, 1987; Vanschoren, 2018; Peng, 2020; Hospedales et al., 2022) explores how deep learning techniques, which often require thousands or even millions of observations to achieve competitive performance, can acquire such a capability. In meta-learning, the learning process is often framed as a bilevel optimization problem (Bard, 2013; Maclaurin et al., 2015; Grefenstette et al., 2019; Lorraine et al., 2020). The outer optimization aims to learn the underlying regularities across a distribution of related tasks and embed them into the inductive biases of a learning algorithm. Consequently, in the inner optimization, the learning algorithm is utilized to quickly adapt to new learning tasks using only a few example instances.

Model-Agnostic Meta-Learning (MAML) (Finn et al., 2017) and its variants (Nichol & Schulman, 2018; Rajeswaran et al., 2019; Song et al., 2019; Triantafillou et al., 2020), are a popular approach to meta-learning. In MAML, the outer optimization aims to learn the underlying regularities across a set of related tasks and embed them into a shared parameter initialization. This initialization is then used in the inner optimization’s learning algorithm to encourage fast adaptation to new tasks. While successful, these methods resort to simple gradient descent using the cross-entropy loss for classification or squared loss for regression for the inner learning algorithm. Consequently, subsequent research has extended MAML to meta-learn additional components, such as the learning rate (Behl et al., 2019; Baik et al., 2020), gradient-based optimizer (Li et al., 2017; Lee & Choi, 2018; Simon et al., 2020; Flennerhag et al., 2020; Kang et al., 2023), loss function (Antoniou & Storkey, 2019), and more (Antoniou et al., 2019; Baik et al., 2023). This enables the meta-learning algorithm to induce stronger inductive biases on the learning algorithm, further enhancing performance.

In this paper, we propose *Neural Procedural Bias Meta-Learning* (NPBML), a novel gradient-based framework for meta-learning task-adaptive procedural biases for deep neural networks. Procedural biases are the subset of inductive biases that determine the order of traversal over the search space (Gordon & Desjardins, 1995), they play a central determining role in the convergence, sample efficiency, and generalization of a learning algorithm. As we will show, the procedural biases are

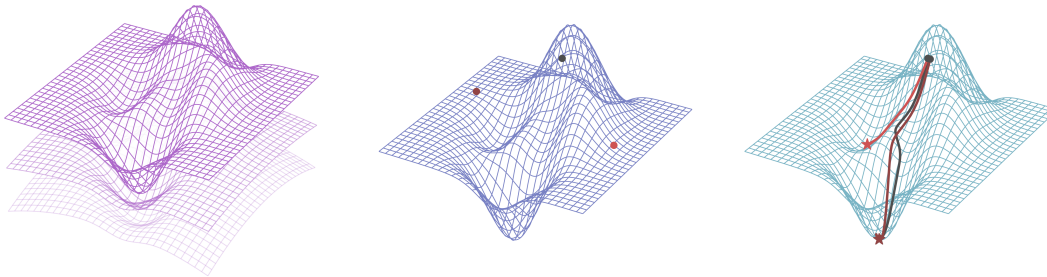


Figure 1: In NPBML, the procedural biases of a deep neural network are meta-learned. This involves meta-learning three key components: the loss function (left), the parameter initialization (center), and the optimizer (right). By meta-learning these components, a strong inductive bias towards fast adaptation can be induced into the learning algorithm.

primarily encoded into three fundamental components of a learning algorithm: the loss function, the optimizer, and the parameter initialization. These components define the geometry of the loss landscape, determine the starting point in this space, and guide the optimization process towards the optimum, respectively, as visualized in Figure 1. Therefore, we aim to meta-learn these three components to maximize the learning performance when using only a few gradient steps.

To achieve this ambitious goal, we first consolidate three related research areas into one unified end-to-end framework: MAML-based learned initializations (Finn et al., 2017), preconditioned gradient descent methods (Lee & Choi, 2018; Flennerhag et al., 2020; Kang et al., 2023), and meta-learned loss functions (Antoniou & Storkey, 2019; Baik et al., 2021; Bechtle et al., 2021; Raymond et al., 2023a;b). We then demonstrate how these meta-learned components can be made task-adaptive through feature-wise linear modulation (FiLM) (Perez et al., 2018) to facilitate downstream task-specific specialization towards each task. The proposed framework is highly flexible and general. As we will show, many existing gradient-based meta-learning approaches arise as special cases of NPBML. To validate the effectiveness of NPBML, we empirically evaluate our proposed algorithm on four well-established few-shot learning benchmarks. The results show that NPBML consistently outperforms many state-of-the-art gradient-based meta-learning algorithms.

2 BACKGROUND

In few-shot meta-learning, we are given access to a collection of tasks $\{\mathcal{T}_1, \mathcal{T}_2, \dots\}$, otherwise known as a meta-dataset, where each task is assumed to be drawn from a task distribution $p(\mathcal{T})$. Each task \mathcal{T}_i contains a support set \mathcal{D}^S , and a query set \mathcal{D}^Q (*i.e.* a training set and a testing set), where $\mathcal{D}^S \cap \mathcal{D}^Q = \emptyset$. Each of these sets contains a set of input-output pairs $\{(x_1, y_1), (x_2, y_2), \dots\}$. Let $x \in X$ and $y \in Y$ denote the inputs and outputs, respectively. In few-shot meta-learning, the goal is to learn a model of the form $f_\theta(x) : X \rightarrow Y$, where θ are the model parameters. The primary challenge of few-shot learning is that f_θ must be able to quickly adapt to any new task $\mathcal{T}_i \sim p(\mathcal{T})$ given only a very limited number of instances. For example, in an N -way K -shot few-shot classification task, f_θ is only given access to K labeled examples of N distinct classes.

2.1 MODEL AGNOSTIC META-LEARNING

MAML (Finn et al., 2017) is a highly influential and seminal method for gradient-based few-shot meta-learning. In MAML, the outer learning objective aims to meta-learn a shared parameter initialization θ over a distribution of related tasks $p(\mathcal{T})$. This shared initialization embeds prior knowledge learned from past learning experiences into the learning algorithm such that when a new unseen task is sampled fast adaptation can occur. The outer optimization prototypically occurs by minimizing the sum of the final losses $\sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}^{meta}$ on the query set at the end of each task’s learning trajectory using gradient descent as follows:

$$\theta_{new} = \theta - \eta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \left[\mathcal{L}^{meta}(\mathcal{D}_i^Q, \theta_{i,j}(\theta)) \right] \quad (1)$$

where $\theta_{i,j}(\boldsymbol{\theta})$ refers to the adapted model parameters on the i^{th} task at the j^{th} iterations of the inner update rule U^{MAML} , and η is the meta learning rate. The inner optimization for each task starts at the meta-learned parameter initialization $\boldsymbol{\theta}$ and applies U^{MAML} to the parameters J times:

$$\theta_{i,j}(\boldsymbol{\theta}) = \boldsymbol{\theta} - \alpha \sum_{j=0}^{J-1} \left[U^{MAML} \left(\mathcal{T}_i, \theta_{i,j}(\boldsymbol{\theta}) \right) \right]. \quad (2)$$

In the original version of MAML, the inner update rule U^{MAML} resorts to simple Stochastic Gradient Descent (SGD) minimizing the loss \mathcal{L}^{base} , typically set to the cross-entropy or squared loss, with a fixed learning rate α across all tasks

$$U^{MAML} \left(\mathcal{T}_i, \theta_{i,j}(\boldsymbol{\theta}) \right) := \theta_{i,j} - \alpha \nabla_{\theta_{i,j}} \left[\mathcal{L}^{base} \left(\mathcal{D}_i^S, \theta_{i,j} \right) \right]. \quad (3)$$

This approach assumes that all tasks in $p(\mathcal{T})$ should use the same fixed learning rule U in the inner optimization. Consequently, this greatly limits the performance that can be achieved when taking only a small number of gradient steps with few labeled samples available.

3 NEURAL PROCEDURAL BIAS META-LEARNING

In this work, we propose *Neural Procedural Bias Meta-Learning* (NPBML), a novel framework that replaces the fixed inner update rule in MAML, *i.e.*, Equation (3), with a meta-learned task-adaptive learning rule. This modifies the inner update rule in three key ways:

1. An optimizer is meta-learned by leveraging the paradigm of preconditioned gradient descent. This involves meta-learning a parameterized preconditioning matrix P_ω with meta-parameters ω to warp the gradients of SGD.
2. The conventional loss function \mathcal{L}^{base} (*i.e.* the standard cross-entropy or squared loss) used in the inner optimization is replaced with a meta-learned loss function \mathcal{M}_ϕ , where ϕ are the meta-parameters.
3. The meta-learned initialization, optimizer, and loss function are adapted to each new task using Feature-Wise Linear Modulation $FiLM_\psi$, a general-purpose preconditioning method that has learnable meta-parameters ψ .

For clarity, we provide a high-level overview of NPBML and contrast it to MAML, before expanding on each of the new components in Sections 3.2, 3.3, and 3.4, respectively. Additionally, pseudocode for the outer and inner optimizations is provided in Algorithm 1 and 2, respectively, in Appendix A.

3.1 OVERVIEW

The central goal of NPBML is to meta-learn a task-adaptive parameter initialization, optimizer, and loss function. This changes the outer optimization previously seen in Equation (1) to the following, where $\Phi = \{\boldsymbol{\theta}, \omega, \phi, \psi\}$ refers to the set of meta-parameters:

$$\Phi_{new} = \Phi - \eta \nabla_{\Phi} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \left[\mathcal{L}^{meta} \left(\mathcal{D}_i^Q, \theta_{i,j}(\Phi) \right) \right]. \quad (4)$$

Unlike MAML which employs a fixed update rule U^{MAML} for all tasks, simple SGD using \mathcal{L}^{base} , NPBML uses a fully meta-learned update rule

$$\theta_{i,j}(\Phi) = \boldsymbol{\theta}(\psi) - \alpha \sum_{j=0}^{J-1} \left[U^{NPBML} \left(\mathcal{T}_i, \theta_{i,j}(\Phi) \right) \right] \quad (5)$$

which adjusts $\theta_{i,j}$ in the direction of the negative gradient of a meta-learned loss function \mathcal{M}_ϕ . Additionally, the gradient is warped via a meta-learned preconditioning matrix P_ω as follows:

$$U^{NPBML} \left(\mathcal{T}_i, \theta_{i,j}(\Phi) \right) := \theta_{i,j} - \alpha P_{(\omega, \psi)} \nabla_{\theta_{i,j}} \left[\mathcal{M}_{(\phi, \psi)} \left(\mathcal{D}_i^S, \theta_{i,j} \right) \right] \quad (6)$$

where both \mathcal{M}_ϕ and P_ω are adapted to each task using $FiLM_\psi$. This new task-adaptive learning rule empowers each task with a unique set of procedural biases enabling strong and robust learning performance on new unseen tasks in $p(\mathcal{T})$ using only a few gradient steps as depicted in Figure 2.

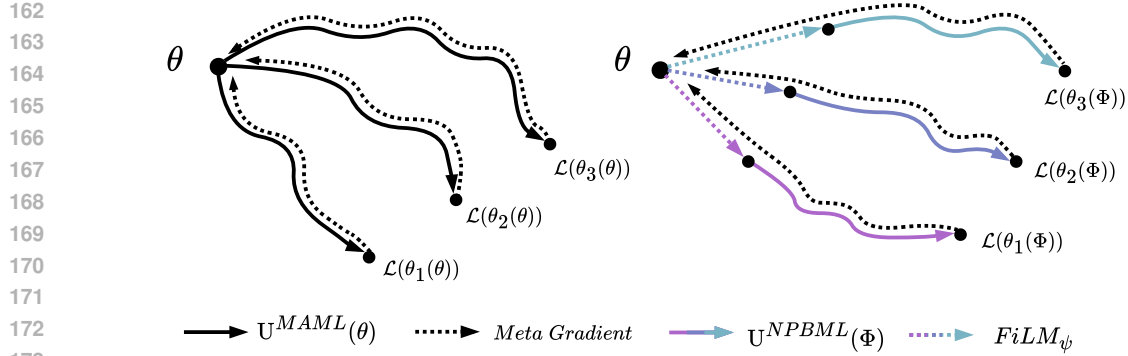


Figure 2: In MAML, the update rule U^{MAML} optimizes the base model parameters from a shared initialization using simple SGD minimizing \mathcal{L}^{base} . In contrast, NPBML adapts the model parameters from a task-adapted initialization using U^{NPBML} , a task-adaptive update rule employing a meta-learned preconditioning matrix P_ω and loss function \mathcal{M}_ϕ .

3.2 META-LEARNED OPTIMIZER

In NPBML, the paradigm of Preconditioned Gradient Descent (PGD) (Li et al., 2017; Lee & Choi, 2018; Park & Oliva, 2019; Flennerhag et al., 2020; Simon et al., 2020; Kang et al., 2023) is employed to meta-learn a gradient preconditioner P_ω that rescales the geometry of the parameter space by modifying the gradient descent update rule as follows:

$$\theta_{new} = \theta - \alpha P_\omega \nabla_\theta \mathcal{M}_{(\phi, \psi)}. \quad (7)$$

In this work, we take inspiration from T-Nets (Lee & Choi, 2018) and insert linear projection layers ω into our model $f_{(\theta, \omega, \psi)}$. The model, composed of an encoder z and a classification head h ,

$$f_{(\theta, \omega, \psi)} = z_{(\theta, \omega, \psi)} \circ h_\theta \quad (8)$$

is interleaved with linear projection layers ω between each layer in the L layer encoder (where σ refers to the non-linear activation functions):

$$z_{(\theta, \omega, \psi)}(x_i) = \sigma^{(L)}(\theta^{(L)} \omega^{(L)} (\dots \sigma^{(1)}(\omega^{(1)} \theta^{(1)} x) \dots)). \quad (9)$$

As described in Equations (4)–(6), ω is meta-learned in the outer loop and held fixed in the inner loop such that preconditioning of the gradients occurs. This form of precondition defines P_ω as a block-diagonal matrix, where each block is defined by the expression $(\omega \omega^\top)$, as shown in (Lee & Choi, 2018). For a simple model where $f_{(\theta, \omega)}(x) = \omega \theta x$, the update rule becomes:

$$\theta_{new} = \theta - \alpha (\omega \omega^\top) \nabla_\theta \mathcal{M}_{(\phi, \psi)}. \quad (10)$$

We leverage this style of parameterization for P_ω due to its relative simplicity and high expressive power. However, we emphasize that the NPBML framework is highly general, and other forms of preconditioning, such as those presented in (Lee & Choi, 2018; Park & Oliva, 2019; Flennerhag et al., 2020; Simon et al., 2020), could be used instead.

3.3 META-LEARNED LOSS FUNCTION

Unlike MAML, which defines the inner loss function \mathcal{L}^{base} to be the cross-entropy or squared loss for all tasks, NPBML uses a meta-learned loss function $\mathcal{M}_{(\phi, \psi)}$ that is learned in the outer optimization. In contrast to handcrafted loss functions, which typically consider only the ground truth label y and the model predictions $f_{(\theta, \omega, \psi)}(x)$, meta-learned loss functions can, in principle, be conditioned on any task-related information (Bechtle et al., 2021; Baik et al., 2021; Raymond, 2024). In NPBML, the meta-learned loss function is conditioned on three distinct sources of task-related information, which are subsequently processed by three small feed-forward neural networks:

$$\mathcal{M}_{(\phi, \psi)} = \mathcal{L}_{(\phi, \psi)}^S + \mathcal{L}_{(\phi, \psi)}^Q + \mathcal{R}_{(\phi, \psi)}. \quad (11)$$

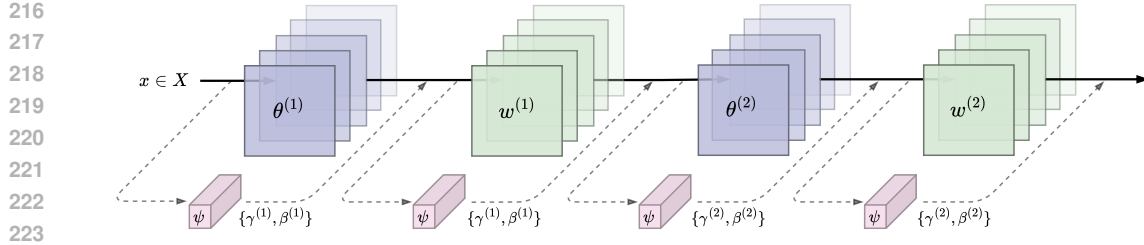


Figure 3: An example of a two-layer convolutional neural network in NPBML, where layers $\theta^{(1)}$ and $\theta^{(2)}$, are interleaved with warp preconditioning layers $\omega^{(1)}$ and $\omega^{(2)}$. Both types of layers are modulated in the inner loop using feature-wise linear modulation layers to induce task adaptation.

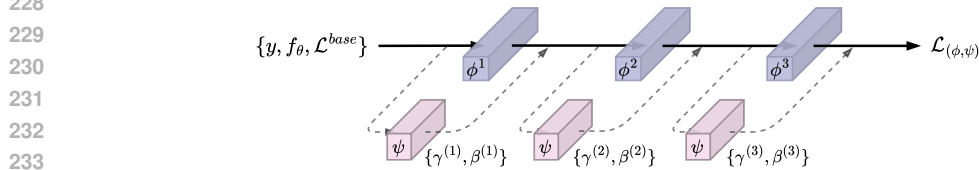


Figure 4: An example of a meta-learned loss function in NPBML, represented as a composition of feed-forward (linear) layers. These layers are modulated using feature-wise linear modulation, resulting in a task-adaptive meta-learned loss function.

First, $\mathcal{L}^S : \mathbb{R}^{2N+1} \rightarrow \mathbb{R}^1$ is an inductive loss function conditioned on task-related information derived from the support set; namely, the one-hot encoded ground truth target and model predictions, and the corresponding loss calculated using \mathcal{L}^{base} . Next, $\mathcal{L}^Q : \mathbb{R}^{2N+1} \rightarrow \mathbb{R}^1$ is a transductive loss function conditioned on task-related information derived from the query set. Here, we give \mathcal{L}^Q access to the model predictions on the query set, embeddings (*i.e.*, relation scores) from a pre-trained relation network (Sung et al., 2018), and the corresponding loss between the model predictions and embeddings using \mathcal{L}^{base} . Note that similar embedding functions have previously been used in (Rusu et al., 2019; Antoniou & Storkey, 2019). Finally, the adapted model parameters $\theta_{i,j}$ are used as inputs to meta-learn a weight regularizer $\mathcal{R} : \mathbb{R}^{4L} \rightarrow \mathbb{R}^1$. To improve efficiency, we condition \mathcal{R} on the mean, standard deviation, L1, and L2 norm of each layer’s weights, as opposed to $\theta_{i,j}$ directly.

3.4 TASK-ADAPTIVE MODULATION

Although all tasks in few-shot learning are assumed to be sampled from the same task distribution $p(\mathcal{T})$, the optimal parameter initialization, optimizer, and loss function may differ between tasks. Therefore, in NPBML the meta-learned components Φ are modulated for each new task, providing each task with a unique set of task-adaptive procedural biases. To achieve this, Feature-wise Linear Modulation (FiLM) layers (Perez et al., 2018; Dumoulin et al., 2018) are inserted into both the encoder $z_{(\theta, \omega, \psi)}$ and the meta-learned loss function $M_{(\phi, \psi)}$ as shown in Figures 3 and 4. $FiLM_{\psi}$ is defined as follows, where γ and β are the scaling and shifting vectors, respectively, and ψ are the meta-learnable FiLM parameters:

$$FiLM_{\psi}(x) = (\gamma_{\psi}(x) + \mathbf{1}) \odot x + \beta_{\psi}(x). \quad (12)$$

Affine transformations conditioned on some input have become increasingly popular, being used by several few-shot learning works to make the learned component adaptive (Oreshkin et al., 2018; Jiang et al., 2018; Vuorio et al., 2019; Zintgraf et al., 2019; Baik et al., 2021; 2023). Furthermore, FiLM layers can help alleviate issues related to batch normalization (Ioffe & Szegedy, 2015), which have been empirically observed to cause training instability due to different distributions of features being passed through the same model in few-shot learning (De Vries et al., 2017; Antoniou et al., 2019). In our work, we have found that conditioning the FiLM on the output activations of the previous layers is an effective way to achieve task adaptability. This form of conditioning is in essence a simplified version of that used in CNAPs (Requeima et al., 2019); however, we have omitted the use of global embeddings as we found it was not necessary for our method.

3.5 INITIALIZATION

Due to the large number of learnable meta-parameters, initialization becomes an important and necessary aspect to consider. Here we detail how to initialize each of the meta-learned components, *i.e.*, $\Phi_0 = \{\theta_0, \omega_0, \phi_0, \psi_0\}$, in NPBML. Firstly, we pre-train the encoder weights θ_0 prior to meta-learning, following many recent methods in few-shot learning (Rusu et al., 2019; Qiao et al., 2018; Requeima et al., 2019; Ye et al., 2020; Ye & Chao, 2021), see Appendix A.3 for more details. For the linear projection layers ω , we leverage the fact that in PGD, setting P_ω to the identity \mathbf{I} recovers SGD. Therefore, we set $\forall l \in \{1, \dots, L\} : \omega^{(l)} = \mathbf{I}$; note for convolutional layers this corresponds to Dirac initialization. Regarding the meta-learned loss function $\mathcal{M}_{(\phi, \psi)}$, the weights ϕ at the start of meta-training are randomly initialized $\phi_0 \sim \mathcal{N}(0, 1e^{-2})$; therefore, the $\mathbb{E}[\mathcal{M}_{(\phi_0, \psi_0)}] = 0$, assuming an identity output activation. Consequently, the definition of the meta-learned loss function in Equation (11) can be modified to

$$\mathcal{M}_{(\phi, \psi)} = \mathcal{L}^{base} + \mathcal{L}_{(\phi, \psi)}^S + \mathcal{L}_{(\phi, \psi)}^Q + \mathcal{R}_{(\phi, \psi)} \quad (13)$$

such that the meta-learned loss function approximately recovers the base loss function at the start of meta-training, *i.e.*, $\mathcal{M}_{(\phi_0, \psi_0)} \approx \mathcal{L}^{base}$. Finally, the FiLM layers in NPBML are initialized using a similar strategy taking advantage of the fact that when $\psi_0 \sim \mathcal{N}(0, 1e - 2)$ the $\mathbb{E}[\gamma_{\psi_0}(x)] = \mathbb{E}[\beta_{\psi_0}(x)] = 0$; consequently, $FiLM_{\psi_0}(x) \approx x$. When initialized in this manner, the update rule for NPBML at the start of meta-training closely approximates the update rule of MAML.

$$\mathbb{U}^{MAML}(\mathcal{T}_i, \theta_{i,j}(\theta_0)) \approx \mathbb{E} \left[\mathbb{U}^{NPBML}(\mathcal{T}_i, \theta_{i,j}(\Phi_0)) \right] \quad (14)$$

4 IMPLICIT META-LEARNING

In NPBML, the parameter initialization, gradient-based optimizer, and loss function are explicitly meta-learned in the outer loop. Critically, we make a novel observation that many other key procedural biases are also implicitly learned by meta-learning these three fundamental components (hence the name given to our algorithm). For example, consider the scalar learning rate α , which is implicitly meta-learned since the following equality holds:

$$\exists \alpha \exists \phi : \theta_{i,j} - \alpha \nabla_{\theta_{i,j}} \mathcal{L}^{base} \approx \theta_{i,j} - \nabla_{\theta_{i,j}} \mathcal{M}_\phi, \quad (15)$$

and if ϕ is made to adapt on each inner step as done in (Baik et al., 2021; Raymond et al., 2023b), then by extension NPBML also learns a learning rate schedule. Another straightforward related observation is that NPBML implicitly learns a layer-wise learning rate $\{\alpha^{(1)}, \dots, \alpha^{(L)}\}$, since for each block $\{\omega^{(1)}, \dots, \omega^{(L)}\}$ in the block diagonal preconditioning matrix P_ω the following holds:

$$\forall l (\exists \omega^{(l)} \exists \alpha^{(l)}) : (\omega^{(l)} (\omega^{(l)})^\top) \nabla_{\theta^{(l)}} \mathcal{M}_\phi \approx \alpha^{(l)} \nabla_{\theta^{(l)}} \mathcal{M}_\phi. \quad (16)$$

There are also less obvious connections that can be drawn. For example, through a transitive relationship, NPBML implicitly learns early stopping when the implicitly learned learning rate approaches zero, as discussed in (Baydin et al., 2018). Furthermore, since there is a linear scaling rule between the batch size and the learning rate (Smith et al., 2017; Smith & Le, 2017; Goyal et al., 2017), NPBML implicitly learns the regularization behavior of the batch size hyperparameter. Another non-trivial example is label smoothing regularization (Müller et al., 2019), which, as proven in (Gonzalez & Miikkulainen, 2020), can be implicitly induced when meta-learning a loss function.

5 RELATED WORK

Meta-learning approaches to few-shot learning aim to equip models with the ability to quickly adapt to new tasks given only a limited number of examples by leveraging prior learning experiences across a distribution of related tasks. These approaches are commonly partitioned into three categories: (1) metric-based methods, which aim to learn a similarity metric for efficient class differentiation (Koch et al., 2015; Vinyals et al., 2016; Sung et al., 2018; Snell et al., 2017); (2) memory-based methods, which utilize architectures that store training examples in memory or directly encode fast adaptation algorithms in the model weights (Santoro et al., 2016; Ravi & Larochelle, 2017); and (3) optimization-based methods, which aim to learn an optimization algorithm specifically designed for

fast adaptation and few-shot learning (Finn et al., 2017). This work explored the latter approach by meta-learning a gradient update rule.

MAML (Finn et al., 2017), a highly flexible task and model-agnostic method for meta-learning a parameter initialization from which fast adaptation can occur. Many follow-up works have sought to enhance MAML’s performance by addressing limitations in the outer-optimization algorithm, such as the memory and compute efficiency (Nichol & Schulman, 2018; Rajeswaran et al., 2019; Raghu et al., 2019; Oh et al., 2020), or the meta-level overfitting (Rusu et al., 2019; Flennerhag et al., 2018). Relatively fewer works have focused on enhancing the inner-optimization update rule, as is done in our work. Most MAML variants continue to use an inner update rule consisting of SGD with a fixed learning rate minimizing a loss function such as the cross-entropy or squared loss.

Of the works that have explored improving the inner update rule, the vast majority focus on improving the optimizer. For example, early MAML-based methods such as (Behl et al., 2019; Antoniou et al., 2019; Li et al., 2017) explored meta-learning the scalar, layer-wise, and parameter-wise learning rates, respectively. More recent methods have explored more powerful parameterization for the meta-learned optimizer through the utilization of preconditioned gradient descent methods (Lee & Choi, 2018; Park & Oliva, 2019; Flennerhag et al., 2020; Simon et al., 2020; Kang et al., 2023), which rescale the geometry of the parameter space by modifying the update rule with a learned preconditioning matrix. While these methods have advanced MAML-based few-shot learning, they often lack a task-adaptive property, falsely assuming that all tasks should use the same optimizer.

A small number of recent works have also investigated replacing the inner loss function (e.g., cross-entropy loss) with a meta-learned loss function. In (Antoniou & Storkey, 2019), a fully transductive loss function represented as a dilated convolutional neural network is meta-learned. Meanwhile, in (Baik et al., 2021), a set of loss functions and task adapters are meta-learned for each step taken in the inner optimization. Although these approaches have shown a lot of promise, their potential has yet to be fully realized, as they have not yet been meta-learned in tandem with the optimizer as we have done in this work.

6 EXPERIMENTAL EVALUATION

In this section, we evaluate the performance of the proposed method on a set of well-established few-shot learning benchmarks. The experimental evaluation aims to answer the following key questions: (1) Can NPBML perform well across a diverse range of few-shot learning tasks? (2) Do the novel components meta-learned in NPBML individually enhance performance? (3) To what extent does each component synergistically contribute to the overall performance of the proposed algorithm?

6.1 RESULTS AND ANALYSIS

To evaluate the performance of NPBML, experiments are performed on four well-established few-shot learning datasets: *mini*-Imagenet (Ravi & Larochelle, 2017), *tiered*-ImageNet (Ren et al., 2018), CIFAR-FS (Bertinetto et al., 2018), and FC-100 (Oreshkin et al., 2018). For each dataset, experiments are performed using both 5-way 1-shot and 5-way 5-shot configurations. Results are also reported on both the 4-CONV (Finn et al., 2017; Zintgraf et al., 2019; Flennerhag et al., 2020; Kang et al., 2023) and ResNet-12 (He et al., 2016; Baik et al., 2020; 2021) network architectures. The full details of all experiments, including a comprehensive description of all datasets, models, and training hyperparameters, can be found in Appendix A. The code for our experiments can be found at https://github.com/*redacted*

6.1.1 MINI-IMAGENET AND TIERED-IMAGENET

We first assess the performance of NPBML and compare it to a range of MAML-based few-shot learning methods on two popular ImageNet derivatives (Deng et al., 2009): *mini*-ImageNet (Ravi & Larochelle, 2017) and *tiered*-ImageNet (Ren et al., 2018). The results, presented in Table 1, demonstrate that the proposed method NPBML, which uses a fully meta-learned update rule in the inner optimization, significantly improves upon the performance of MAML-based few-shot learning methods. The proposed method achieves higher meta-testing accuracy in the 1-shot and 5-shot settings using both low-capacity (4-CONV) and high-capacity (ResNet-12) models.

Table 1: Few-shot classification meta-testing accuracy on 5-way 1-shot and 5-way 5-shot *mini-ImageNet* and *tiered-ImageNet* where \pm represents the 95% confidence intervals.

Method	Base Learner	<i>mini-ImageNet</i> (5-way)		<i>tiered-ImageNet</i> (5-way)	
		1-shot	5-shot	1-shot	5-shot
MAML ¹	4-CONV	48.70 \pm 1.84%	63.11 \pm 0.92%	50.98 \pm 0.26%	66.25 \pm 0.19%
MetaSGD ²	4-CONV	50.47 \pm 1.87%	64.03 \pm 0.94%	-	-
T-Net ³	4-CONV	50.86 \pm 1.82%	-	-	-
MAML++ ⁴	4-CONV	52.15 \pm 0.26%	68.32 \pm 0.44%	-	-
SCA ⁵	4-CONV	54.84 \pm 0.99%	71.85 \pm 0.53%	-	-
WarpGrad ⁷	4-CONV	52.30 \pm 0.80%	68.40 \pm 0.60%	57.20 \pm 0.90%	74.10 \pm 0.70%
ModGrad ⁸	4-CONV	53.20 \pm 0.86%	69.17 \pm 0.69%	-	-
MeTAL ⁹	4-CONV	52.63 \pm 0.37%	70.52 \pm 0.29%	54.34 \pm 0.31%	70.40 \pm 0.21%
ALFA ¹⁰	4-CONV	50.58 \pm 0.51%	69.12 \pm 0.47%	53.16 \pm 0.49%	70.54 \pm 0.46%
GAP ¹¹	4-CONV	54.86 \pm 0.85%	71.55 \pm 0.61%	57.60 \pm 0.93%	74.90 \pm 0.68%
NPBML	4-CONV	57.49\pm0.83%	75.01\pm0.64%	64.24\pm0.97%	79.17\pm0.71%
MAML ¹	ResNet-12	58.60 \pm 0.42%	69.54 \pm 0.38%	59.82 \pm 0.41%	73.17 \pm 0.32%
MC ⁶	WRN-28-10	-	-	64.40 \pm 0.10%	80.21 \pm 0.10%
ModGrad ⁸	WRN-28-10	-	-	65.72 \pm 0.21%	81.17 \pm 0.20%
MeTAL ⁹	ResNet-12	59.64 \pm 0.38%	76.20 \pm 0.19%	63.89 \pm 0.43%	80.14 \pm 0.40%
ALFA ¹⁰	ResNet-12	59.74 \pm 0.49%	77.96 \pm 0.41%	64.62 \pm 0.49%	82.48 \pm 0.38%
NPBML	ResNet-12	61.59\pm0.80%	78.18\pm0.60%	72.22\pm0.96%	85.41\pm0.61%

¹(Finn et al., 2017) ²(Li et al., 2017) ³(Lee & Choi, 2018) ⁴(Antoniou et al., 2019) ⁵(Antoniou & Storkey, 2019)
⁶(Park & Oliva, 2019) ⁷(Flennerhag et al., 2020) ⁸(Simon et al., 2020) ⁹(Baik et al., 2021) ¹⁰(Baik et al., 2023)
¹¹(Kang et al., 2023)

In contrast to PGD methods Meta-SGD, T-Net, WarpGrad, ModGrad, ALFA, and GAP, which meta-learn an optimizer alongside the parameter initialization, NPBML shows clear gains in generalization performance. This improvement is also evident when compared to SCA and MeTAL, which replace the inner optimization’s loss function with a meta-learned loss function. These results empirically demonstrate that meta-learning an optimizer and loss function are complementary and orthogonal approaches to improving MAML-based few-shot learning methods.

On *tiered-ImageNet*, the larger of the two datasets, we find that the difference between NPBML and its competitors is even more pronounced than on *mini-ImageNet*. This result suggests that when given enough data, NPBML can learn highly expressive inner update rules that significantly enhances few-shot learning performance. However, meta-overfitting can occur on smaller datasets, necessitating regularization techniques as discussed in Appendix A. Alternatively, we conjecture that less expressive representations for P_ω would also reduce meta-overfitting.

6.1.2 CIFAR-FS AND FC-100

Next, we further validate the effectiveness of NPBML on two popular CIFAR-100 derivatives (Krizhevsky & Hinton, 2009): CIFAR-FS (Ravi & Larochelle, 2017) and FC-100 (Ren et al., 2018). The results, presented in Table 2, show that NPBML continues to achieve strong and robust generalization performance across all settings and models. These results are particularly impressive, given that both MeTAL and ALFA ensemble the top 5 performing models from the same run, which significantly increases the model size and capacity. These experimental results reinforce our claim that meta-learning a task-adaptive update rule is an effective approach to improving the performance of MAML-based few-shot learning algorithms.

Table 2: Few-shot classification meta-testing accuracy on 5-way 1-shot and 5-way 5-shot *CIFAR-FS* and *FC-100* where \pm represents the 95% confidence intervals.

Method	Base Learner	CIFAR-FS (5-way)		FC-100 (5-way)	
		1-shot	5-shot	1-shot	5-shot
MAML ¹	4-CONV	57.63±0.73%	73.95±0.84%	35.89±0.72%	49.31±0.47%
BOIL ²	4-CONV	58.03±0.43%	73.61±0.32%	38.93±0.45%	51.66±0.32%
MeTAL ³	4-CONV	59.16±0.56%	74.62±0.42%	37.46±0.39%	51.34±0.25%
ALFA ⁴	4-CONV	59.96±0.49%	76.79±0.42%	37.99±0.48%	53.01±0.49%
NPBML	4-CONV	64.90±0.94%	79.24±0.69%	40.56±0.76%	53.48±0.68%
MAML ¹	ResNet-12	63.81±0.54%	77.07±0.42%	37.29±0.40%	50.70±0.35%
MeTAL ³	ResNet-12	67.97±0.47%	82.17±0.38%	39.98±0.39%	53.85±0.36%
ALFA ⁴	ResNet-12	66.79±0.47%	83.62±0.37%	41.46±0.49%	55.82±0.50%
NPBML	ResNet-12	69.30±0.91%	83.72±0.64%	43.63±0.71%	59.85±0.70%

¹(Finn et al., 2017) ²(Oh et al., 2020) ³(Baik et al., 2021) ⁴(Baik et al., 2023)

6.2 ABLATION STUDIES

To further investigate the performance of the proposed method, we conduct two sets of ablation studies to analyze the effectiveness of each component. All ablation experiments are performed using the 4-CONV network architecture in a 5-way 5-shot setting on the *mini-ImageNet* dataset.

6.2.1 META-LEARNED COMPONENTS

First, we examine the importance of the meta-learned optimizer P_ω , loss function \mathcal{M}_ϕ , and task-adaptive conditioning method $FILM_\psi$. The results are presented in Table 3, and they demonstrate that each of the proposed components clearly and significantly contributes to the performance of NPBML. In (2) MAML is modified to include gradient preconditioning, which increases accuracy by 2.09%. Conversely in (3) we modify MAML with our meta-learned loss function, resulting in a 6.37% performance increase. Interestingly, the meta-learned loss function enhances performance by a larger margin; however, this may be due to the relatively simple T-Net style optimizer used in NPBML. This suggests that a more powerful parameterization, such as (Flennerhag et al., 2018) or (Kang et al., 2023), may further improve performance. In (4), MAML is modified to include both the optimizer and loss function, resulting in a 7.41% performance increase. This further supports our claim that meta-learning both an optimizer and a loss function are complementary and orthogonal approaches to improving MAML. Finally, in (5), we add our task-adaptive conditioning method, increasing performance by 2.22% over the prior experiment and 9.63% over MAML.

6.2.2 META-LEARNED LOSS FUNCTION

The prior ablation study shows that the meta-learned loss function \mathcal{M}_ϕ is a crucial component in NPBML. Therefore, we further investigate each of the components; namely, the meta-learned inductive and transductive loss functions, and weight regularizer. The results are presented in Table 4, and surprisingly, they show that each of the components in isolation (7), (8), and (9), improves performance by approximately 5%. However, when combined in (10), the total performance increase is 6.37%. We hypothesize that this result is a consequence of the implicit meta-learning of the learning rate identified in Equation (15), which not only holds for \mathcal{M}_ϕ , but also for each of its components, *i.e.*, the equality is also true when \mathcal{M}_ϕ is replaced with \mathcal{L}_ϕ^S , \mathcal{L}_ϕ^Q , or \mathcal{R}_ϕ . Since all components share implicit learning rate tuning, the performance gains from this behavior do not accumulate; however, the improvement in (10) is better than each component in isolation indicating that each component provides additional unique benefits to the meta-learning process.

Table 3: Ablation study of the meta-learned components in NPBML, reporting the meta-testing accuracy on *mini-ImageNet* 5-way 5-shot. A \checkmark denotes that the component is meta-learned, with variant (1) reducing to MAML, while variant (5) represents our final proposed algorithm.

	Initialization	Optimizer	Loss Function	Task-Adaptive	Accuracy
(1)	\checkmark				65.38 \pm 0.67%
(2)	\checkmark	\checkmark			67.47 \pm 0.68%
(3)	\checkmark		\checkmark		71.75 \pm 0.69%
(4)	\checkmark	\checkmark	\checkmark		72.79 \pm 0.67%
(5)	\checkmark	\checkmark	\checkmark	\checkmark	75.01\pm0.64%

Table 4: Ablation study of the meta-learned loss function \mathcal{M}_ϕ in NPBML, reporting the meta-testing accuracy on *mini-ImageNet* 5-way 5-shot. Note, variants (6) and (10) correspond to variants (1) and (3), respectively in Table 3.

	Base Loss	Inductive Loss	Transductive Loss	Weight Regularizer	Accuracy
(6)	\checkmark				65.38 \pm 0.67%
(7)	\checkmark	\checkmark			70.68 \pm 0.66%
(8)	\checkmark		\checkmark		70.92 \pm 0.68%
(9)	\checkmark			\checkmark	70.04 \pm 0.65%
(10)	\checkmark	\checkmark	\checkmark	\checkmark	71.75\pm0.69%

7 CONCLUSION

In this work, we propose a novel meta-learning framework for learning the procedural biases of a deep neural network. The proposed technique, *Neural Procedural Bias Meta-Learning* (NPBML), consolidates recent advancements in MAML-based few-shot learning methods by replacing the fixed inner update rule with a fully meta-learned update rule. This is achieved by meta-learning a task-adaptive loss function, optimizer, and parameter initialization. The experimental results confirm the effectiveness and scalability of the proposed approach, demonstrating strong few-shot learning performance across a range of popular benchmarks. We believe NPBML provides a principled framework for advancing general-purpose meta-learning in deep neural networks. Looking ahead, numerous compelling future research directions exist, such as developing more powerful parameterizations for the meta-learned optimizer or loss function. We expect that further investigation of this topic will result in more expressive inner update rules, resulting in increased robustness and efficiency within the context of optimization-based meta-learning. Finally, broadening the scope of the proposed framework to encompass the related domains of cross-domain few-shot learning and continual learning may be a promising avenue for future exploration.

REFERENCES

- 540
541
542 Antreas Antoniou and Amos J Storkey. Learning to learn by self-critique. *Advances in Neural*
543 *Information Processing Systems*, 2019.
- 544 Antreas Antoniou, Harrison Edwards, and Amos Storkey. How to train your maml. *International*
545 *Conference on Learning Representations*, 2019.
- 546
547 Sungyong Baik, Myungsub Choi, Janghoon Choi, Heewon Kim, and Kyoung Mu Lee. Meta-
548 learning with adaptive hyperparameters. *Advances in Neural Information Processing Systems*,
549 2020.
- 550 Sungyong Baik, Janghoon Choi, Heewon Kim, Dohee Cho, Jaesik Min, and Kyoung Mu Lee. Meta-
551 learning with task-adaptive loss function for few-shot learning. In *International Conference on*
552 *Computer Vision*, 2021.
- 553 Sungyong Baik, Myungsub Choi, Janghoon Choi, Heewon Kim, and Kyoung Mu Lee. Learning to
554 learn task-adaptive hyperparameters for few-shot learning. *Transactions on Pattern Analysis and*
555 *Machine Intelligence*, 2023.
- 556
557 Jonathan F Bard. *Practical Bilevel Optimization: Algorithms and Applications*. Springer Science &
558 Business Media, 2013.
- 559 Atılım Güneş Baydin, Robert Cornish, David Martínez Rubio, Mark Schmidt, and Frank Wood. On-
560 line learning rate adaptation with hypergradient descent. In *International Conference on Learning*
561 *Representations*, 2018.
- 562 Sarah Bechtle, Artem Molchanov, Yevgen Chebotar, Edward Grefenstette, Ludovic Righetti, Gaurav
563 Sukhatme, and Franziska Meier. Meta learning via learned loss. In *International Conference on*
564 *Pattern Recognition (ICPR)*, 2021.
- 565
566 Harkirat Singh Behl, Atılım Güneş Baydin, and Philip HS Torr. Alpha maml: Adaptive model-
567 agnostic meta-learning. *arXiv preprint arXiv:1905.07435*, 2019.
- 568 Luca Bertinetto, Joao F Henriques, Philip HS Torr, and Andrea Vedaldi. Meta-learning with differ-
569 entiable closed-form solvers. *arXiv preprint arXiv:1805.08136*, 2018.
- 570
571 Harm De Vries, Florian Strub, Jérémie Mary, Hugo Larochelle, Olivier Pietquin, and Aaron C
572 Courville. Modulating early visual processing by language. *Advances in Neural Information*
573 *Processing Systems*, 2017.
- 574
575 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
576 hierarchical image database. In *Conference on Computer Vision and Pattern Recognition*, 2009.
- 577 Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville,
578 and Yoshua Bengio. Feature-wise transformations. *Distill*, 2018.
- 579
580 Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation
581 of deep networks. In *International Conference on Machine Learning*, 2017.
- 582 Sebastian Flennerhag, Pablo G Moreno, Neil D Lawrence, and Andreas Damianou. Transferring
583 knowledge across learning processes. *International Conference on Learning Representations*,
584 2018.
- 585
586 Sebastian Flennerhag, Andrei A Rusu, Razvan Pascanu, Francesco Visin, Hujun Yin, and Raia
587 Hadsell. Meta-learning with warped gradient descent. *International Conference on Learning*
588 *Representations*, 2020.
- 589 Santiago Gonzalez and Risto Miikkulainen. Effective regularization through loss function meta-
590 learning. *arXiv preprint arXiv:2010.00788*, 2020.
- 591
592 Diana F Gordon and Marie Desjardins. Evaluation and selection of biases in machine learning.
593 *Machine Learning*, 1995.

- 594 Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, An-
595 drew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet
596 in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- 597 Edward Grefenstette, Brandon Amos, Denis Yarats, Phu Mon Htut, Artem Molchanov, Franziska
598 Meier, Douwe Kiela, Kyunghyun Cho, and Soumith Chintala. Generalized inner loop meta-
599 learning. *arXiv preprint arXiv:1910.01727*, 2019.
- 600 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
601 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- 602 Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural
603 networks: A survey. *Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 604 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
605 reducing internal covariate shift. In *International Conference on Machine Learning*, 2015.
- 606 Xiang Jiang, Mohammad Havaei, Farshid Varno, Gabriel Chartrand, Nicolas Chapados, and Stan
607 Matwin. Learning to learn with conditional class dependencies. In *International Conference on*
608 *Learning Representations*, 2018.
- 609 Suhyun Kang, Duhun Hwang, Moonjung Eo, Taesup Kim, and Wonjong Rhee. Meta-learning with
610 a geometry-adaptive preconditioner. In *Conference on Computer Vision and Pattern Recognition*,
611 2023.
- 612 Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. Siamese neural networks for one-shot
613 image recognition. In *International Conference on Machine Learning - Deep Learning Workshop*,
614 2015.
- 615 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- 616 Kwonjoon Lee, Subhransu Maji, Avinash Ravichandran, and Stefano Soatto. Meta-learning with
617 differentiable convex optimization. In *Conference on Computer Vision and Pattern Recognition*,
618 2019.
- 619 Yoonho Lee and Seungjin Choi. Gradient-based meta-learning with learned layerwise metric and
620 subspace. In *International Conference on Machine Learning*, 2018.
- 621 Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-sgd: Learning to learn quickly for few-
622 shot learning. *arXiv preprint arXiv:1707.09835*, 2017.
- 623 Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters by
624 implicit differentiation. In *International Conference on Artificial Intelligence and Statistics*, 2020.
- 625 Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimiza-
626 tion through reversible learning. In *International Conference on Machine Learning*, 2015.
- 627 Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? *Ad-*
628 *vances in Neural Information Processing Systems*, 2019.
- 629 Alex Nichol and John Schulman. Reptile: A scalable meta-learning algorithm. *arXiv preprint*
630 *arXiv:1803.02999*, 2018.
- 631 Jaehoon Oh, Hyungjun Yoo, ChangHwan Kim, and Se-Young Yun. Boil: Towards representation
632 change for few-shot learning. *arXiv preprint arXiv:2008.08882*, 2020.
- 633 Boris Oreshkin, Pau Rodríguez López, and Alexandre Lacoste. Tadam: Task dependent adaptive
634 metric for improved few-shot learning. *Advances in Neural Information Processing Systems*,
635 2018.
- 636 Eunbyung Park and Junier B Oliva. Meta-curvature. *Advances in Neural Information Processing*
637 *Systems*, 2019.
- 638 Huimin Peng. A comprehensive overview and survey of recent advances in meta-learning. *arXiv*
639 *preprint arXiv:2004.11149*, 2020.

- 648 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual
649 reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial*
650 *Intelligence*, 2018.
- 651 Apostolos F Psaros, Kenji Kawaguchi, and George Em Karniadakis. Meta-learning pinn loss func-
652 tions. *Journal of computational physics*, 2022.
- 653 Siyuan Qiao, Chenxi Liu, Wei Shen, and Alan L Yuille. Few-shot image recognition by predicting
654 parameters from activations. In *Conference on Computer Vision and Pattern Recognition*, 2018.
- 655 Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. Rapid learning or feature reuse?
656 towards understanding the effectiveness of maml. *arXiv preprint arXiv:1909.09157*, 2019.
- 657 Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with im-
658 plicit gradients. *Advances in Neural Information Processing Systems*, 2019.
- 659 Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International*
660 *Conference on Learning Representations*, 2017.
- 661 Christian Raymond. Meta-learning loss functions for deep neural networks. *arXiv preprint*
662 *arXiv:2406.09713*, 2024.
- 663 Christian Raymond, Qi Chen, Bing Xue, and Mengjie Zhang. Learning symbolic model-agnostic
664 loss functions via meta-learning. *Transactions on Pattern Analysis and Machine Intelligence*,
665 2023a.
- 666 Christian Raymond, Qi Chen, Bing Xue, and Mengjie Zhang. Online loss function learning. *arXiv*
667 *preprint arXiv:2301.13247*, 2023b.
- 668 Mengye Ren, Eleni Triantafillou, Sachin Ravi, Jake Snell, Kevin Swersky, Joshua B. Tenenbaum,
669 Hugo Larochelle, and Richard S. Zemel. Meta-learning for semi-supervised few-shot classifica-
670 tion. In *International Conference on Learning Representations*, 2018.
- 671 James Requeima, Jonathan Gordon, John Bronskill, Sebastian Nowozin, and Richard E Turner. Fast
672 and flexible multi-task classification using conditional neural adaptive processes. *Advances in*
673 *Neural Information Processing Systems*, 2019.
- 674 Andrei A. Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero,
675 and Raia Hadsell. Meta-learning with latent embedding optimization. In *International Conference*
676 *on Learning Representations*, 2019.
- 677 Adam Santoro, Sergey Bartunov, Matthew Botvinick, Daan Wierstra, and Timothy Lillicrap. Meta-
678 learning with memory-augmented neural networks. In *International Conference on Machine*
679 *Learning*, 2016.
- 680 Jürgen Schmidhuber. *Evolutionary Principles in Self-Referential Learning*. PhD thesis, Technische
681 Universität München, 1987.
- 682 Christian Simon, Piotr Koniusz, Richard Nock, and Mehrtash Harandi. On modulating the gradient
683 for meta-learning. In *European Conference on Computer Vision*, 2020.
- 684 Samuel L Smith and Quoc V Le. A bayesian perspective on generalization and stochastic gradient
685 descent. *arXiv preprint arXiv:1710.06451*, 2017.
- 686 Samuel L Smith, Pieter-Jan Kindermans, Chris Ying, and Quoc V Le. Don’t decay the learning rate,
687 increase the batch size. *arXiv preprint arXiv:1711.00489*, 2017.
- 688 Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Ad-*
689 *vances in Neural Information Processing Systems*, 2017.
- 690 Xingyou Song, Wenbo Gao, Yuxiang Yang, Krzysztof Choromanski, Aldo Pacchiano, and Yunhao
691 Tang. Es-maml: Simple hessian-free meta learning. *arXiv preprint arXiv:1910.01215*, 2019.
- 692
693
694
695
696
697
698
699
700
701

702 Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang, Philip HS Torr, and Timothy M Hospedales.
703 Learning to compare: Relation network for few-shot learning. In *Conference on Computer Vision*
704 *and Pattern Recognition*, 2018.

705 Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross
706 Goroshin, Carles Gelada, Kevin Jordan Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle.
707 Meta-dataset: A dataset of datasets for learning to learn from few examples. In *International*
708 *Conference on Learning Representations*, 2020.

709

710 Joaquin Vanschoren. Meta-learning: A survey. *arXiv preprint arXiv:1810.03548*, 2018.

711

712 Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for
713 one-shot learning. *Advances in Neural Information Processing Systems*, 2016.

714 Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J Lim. Multimodal model-agnostic meta-
715 learning via task-aware modulation. *Advances in Neural Information Processing Systems*, 2019.

716

717 Han-Jia Ye and Wei-Lun Chao. How to train your maml to excel in few-shot classification. In
718 *International Conference on Learning Representations*, 2021.

719 Han-Jia Ye, Hexiang Hu, De-Chuan Zhan, and Fei Sha. Few-shot learning via embedding adaptation
720 with set-to-set functions. In *Conference on Computer Vision and Pattern Recognition*, 2020.

721

722 Luisa Zintgraf, Kyriacos Shiarli, Vitaly Kurin, Katja Hofmann, and Shimon Whiteson. Fast context
723 adaptation via meta-learning. In *International Conference on Machine Learning*, 2019.

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

Algorithm 1 Meta-Learning (Outer Loop)

Input: $\mathcal{L}^{meta} \leftarrow$ Meta loss function
Input: $p(\mathcal{T}) \leftarrow$ Task distribution
Input: $\eta \leftarrow$ Meta learning rate

- 1: $\Phi_0 \leftarrow$ Initialize meta-parameters $\{\theta, \omega, \phi, \psi\}$
- 2: **for** $t \in \{0, \dots, S^{meta}\}$ **do**
- 3: $\mathcal{T}_0, \mathcal{T}_1, \dots, \mathcal{T}_B \leftarrow$ Sample tasks from $p(\mathcal{T})$
- 4: **for** $i \in \{0, \dots, B\}$ **do**
- 5: $\mathcal{D}_i^S = \{(x_i^s, y_i^s)\}_{s=0}^S \leftarrow$ Sample support from \mathcal{T}_i
- 6: $\mathcal{D}_i^Q = \{(x_i^q, y_i^q)\}_{q=0}^Q \leftarrow$ Sample query from \mathcal{T}_i
- 7: $\theta_{i,j} \leftarrow$ Base-Learning using Algorithm (2)
- 8: $\Phi_{t+1} \leftarrow \Phi_t - \eta \frac{1}{B} \nabla_{\Phi_t} \sum_i \mathcal{L}_i^{meta}(\mathcal{D}_i^Q, \theta_{i,j})$
- 9: **return** Φ_t

Algorithm 2 Base-Learning (Inner Loop)

Input: $\mathcal{L}^{base} \leftarrow$ Base loss function
Input: $\mathcal{D}_i^S, \mathcal{D}_i^Q \leftarrow$ Support and query sets
Input: $\Phi \leftarrow$ Meta parameters $\{\theta, \omega, \phi, \psi\}$
Input: $\alpha \leftarrow$ Base learning rate
Input: $g \leftarrow$ Relation network

- 1: $\theta_{i,0} \leftarrow$ Initialize base weights with θ
- 2: **for** $j \in \{0, \dots, S^{base}\}$ **do**
- 3: $\hat{y}_i^S, \hat{y}_i^Q \leftarrow f_{(\theta_{i,j}, \omega, \psi)}(x_i^S \cup x_i^Q)$
- 4: $\mathcal{L}_{i,j}^{base} \leftarrow \frac{1}{|\mathcal{D}^S|} \sum \mathcal{L}^{base}(y_i^S, \hat{y}_i^S)$
- 5: $\mathcal{L}_{i,j}^S \leftarrow \frac{1}{|\mathcal{D}^S|} \sum \mathcal{L}_{(\phi, \psi)}^S(y_i^S, \hat{y}_i^S)$
- 6: $\mathcal{L}_{i,j}^Q \leftarrow \frac{1}{|\mathcal{D}^Q|} \sum \mathcal{L}_{(\phi, \psi)}^Q(g(x_i^Q), \hat{y}_i^Q)$
- 7: $\mathcal{R}_{i,j} \leftarrow \mathcal{R}_{(\phi, \psi)}(\theta_{i,j})$
- 8: $\mathcal{M}_{(\phi, \psi)} \leftarrow \mathcal{L}_{i,j}^{base} + \mathcal{L}_{i,j}^S + \mathcal{L}_{i,j}^Q + \mathcal{R}_{i,j}$
- 9: $\theta_{i,j+1} \leftarrow \theta_{i,j} - \alpha P \nabla_{\theta_{i,j}} \mathcal{M}_{(\phi, \psi)}$
- 10: **return** $\theta_{i,j}$

A EXPERIMENTAL SETUP

In this section, we provide a comprehensive description of the experimental settings used in this work. Section A.1 offers an overview of the datasets utilized, Section A.2 discusses the network architectures, and Section A.3 details the hyperparameters specific to our proposed method, NPBML. For further details, please refer to Algorithms 1 and 2, as well as our code made available at https://github.com/*redacted*

A.1 DATASET CONFIGURATIONS

The few-shot learning experiments follow a prototypical episodic learning setup (Ravi & Larochelle, 2017), where each dataset contains a set of non-overlapping meta-training, meta-validation, and meta-testing tasks. Each task \mathcal{T}_i has a support \mathcal{D}^S and query \mathcal{D}^Q set (*i.e.*, training and testing set, respectively). For an N -way K -shot classification task, the support set contains $N \times K$ instances, and the query set contains $N \times 15$ instances, where N refers to the number of randomly sampled classes, and K refers to the number of instances (*i.e.*, shots) available from each of those classes. For example, in a 5-way 5-shot classification task, the support set contains $5 \times 5 = |\mathcal{D}^S|$ instances to train on, and the query set contains $5 \times 15 = |\mathcal{D}^Q|$ instances to validate the performance on.

810 A.1.1 IMAGENET DATASETS

811
812 Two commonly used datasets for few-shot learning are two ImageNet derivatives (Deng et al., 2009):
813 *mini-ImageNet* (Ravi & Larochelle, 2017) and *tiered-ImageNet* (Ren et al., 2018). Both datasets are
814 composed of three subsets (training, validation, and testing), each of which consists of images with
815 a size of 84×84 . The two datasets differ in regard to how the classes are partitioned into mutually
816 exclusive subsets.

817 **mini-ImageNet** randomly samples 100 classes from the 1000 base classes in ImageNet. Follow-
818 ing (Vinyals et al., 2016) the sampled classes are partitioned such that 64 classes are allocated for
819 meta-training, 16 for meta-validation, and 20 for meta-testing, where for each class 600 images are
820 available.

821 **tiered-ImageNet** as described in (Ren et al., 2018) alternatively stratifies 608 classes into 34 higher-
822 level categories in the ImageNet human-curated hierarchy. The 34 classes are partitioned into 20
823 categories for meta-training, 6 for meta-validation, and 8 for meta-testing. Each class in *tiered-*
824 *ImageNet* has a minimum of 732 instances and a maximum of 1300.

826 A.1.2 CIFAR-100 DATASETS

827 Additional experiments are also conducted on CIFAR-FS (Bertinetto et al., 2018) and FC-100 (Ore-
828 shkin et al., 2018), which are few-shot learning datasets derived from the popular CIFAR-100 dataset
829 (Krizhevsky & Hinton, 2009). The CIFAR100 dataset contains 100 classes containing 600 images
830 each, each with a resolution of 32×32 .

831 **CIFAR-FS** uses a similar sampling procedure to *mini-ImageNet* (Ravi & Larochelle, 2017), CIFAR-
832 FS is derived by randomly sampling 100 classes from the 100 base classes in CIFAR100. The
833 sampled classes are partitioned such that 60 classes are allocated for meta-training, 20 for meta-
834 validation, and 20 for meta-testing.

835 **FC-100** is obtained by using a dataset construction process similar to *tiered-ImageNet* (Ren et al.,
836 2018), in which class hierarchies are used to partition the original dataset to simulate more challeng-
837 ing few-shot learning scenarios. In FC100 there are a total of 20 high-level classes, where 12 classes
838 are allocated for meta-training, 4 for meta-validation, and 4 for meta-testing.

841 A.2 NETWORK ARCHITECTURES

842 **4-CONV:** Following the encoder architecture settings in (Zintgraf et al., 2019; Flennerhag et al.,
843 2020; Kang et al., 2023), the network consists of four modules. Each module contains a 3×3
844 convolutional layer with 128 filters, followed by a batch normalization layer, a ReLU non-linearity,
845 and a 2×2 max-pooling downsampling layer. Following the four modules, we apply average pooling
846 and flatten the embedding to a size of 128.

847 **ResNet-12:** The encoder z_θ follows the standard network architecture settings used in (Baik et al.,
848 2020; 2021). The network, similar to 4-CONV, consists of four modules. Each module contains
849 a stack of three 3×3 convolutional layers, where each layer is followed by batch normalization
850 and a leaky ReLU non-linearity. A skip convolutional over the convolutional stack is used in each
851 module. Each skip connection contains a 1×1 convolutional layer followed by batch normalization.
852 Following the convolutional stack and skip connection, a leaky ReLU non-linearity and 2×2 max-
853 pooling downsampling layer are placed at the end of each module. The number of filters used in each
854 module is [64, 128, 256, 512], respectively. Following the four modules, we apply average pooling
855 and flatten the embedding to a size of 512.

856 **Classifier:** As MAML and its variants are sensitive to label permutation during meta-testing, we opt
857 to use the permutation invariant classification head proposed by (Ye & Chao, 2021). This replaces
858 the typical classification head $h_\theta \in \mathbb{R}^{in \times N}$ with a *single-weight vector* $\theta_{head} \in \mathbb{R}^{in \times 1}$ which is
859 meta-learned at meta-training time. During meta-testing the weight vector is duplicated into each
860 output, *i.e.*, $h_\theta = \{\theta_c = \theta_{head}\}_{c=1}^N$, and adapted in an identical fashion to traditional MAML.
861 This reduces the number of parameters in the classification head and makes NPBML permutation
862 invariant similar to MetaOptNet (Lee et al., 2019), CNAPs (Requeima et al., 2019), and ProtoMAML
863 (Triantafillou et al., 2020).

864 **Loss Network:** The meta-learned loss function \mathcal{M} is composed of three separate networks \mathcal{L}^S , \mathcal{L}^Q ,
 865 and \mathcal{R} whose scalar outputs are summed to produce the final output as shown in Equation (11).
 866 The network architecture for all three networks is identical except for the input dimension of the
 867 first layer as discussed in Section 3.3. The network architecture is taken from (Bechtle et al., 2021;
 868 Psaros et al., 2022; Raymond et al., 2023b) and is a small feedforward rectified linear neural network
 869 with two hidden layers and 40 hidden units in each layer. Note \mathcal{L}^S and \mathcal{L}^Q are applied instance-wise
 870 and reduced using a mean reduction, such that they are both invariant to the number of instances
 871 made available at meta-testing time in the support and query set.

872 **Relation Network:** The transductive loss L^Q takes an embedding from a pre-trained relation net-
 873 work (Sung et al., 2018) as one of its inputs, similar to (Rusu et al., 2019; Antoniou & Storkey,
 874 2019). The relation network used in our experiments employs the previously mentioned ResNet-12
 875 encoder to generate a set of embeddings for each instance in $D^S \cup D^Q$. These embeddings are
 876 then concatenated and processed using a relation module. The relation module in our experiments
 877 consists of two ResNet blocks with 2×512 and 512 filters, respectively. The output is then flattened
 878 and averaged pooled to an embedding size of 512, which is passed through a final linear layer of
 879 size 512×1 which outputs each relation score. In regard to training, the relation network is trained
 880 using the settings described in (Sung et al., 2018).

881 A.3 PROPOSED METHOD SETTINGS

882 A.3.1 PRE-TRAINING SETTINGS

883 In NPBML, the encoder portion of the network z_θ is pre-trained prior to meta-learning, following
 884 many recent methods in few-shot learning (Rusu et al., 2019; Qiao et al., 2018; Requeima et al.,
 885 2019; Ye et al., 2020; Ye & Chao, 2021). For both the 4-CONV and ResNet-12 models, we fol-
 886 low the pre-training pipeline used in (Ye et al., 2020; Ye & Chao, 2021). During pre-training, we
 887 append the feature encoder z_θ with a temporary classification head f_θ and train it to classify all
 888 classes in the D^{train} (e.g., over all 64 classes in the *mini*-ImageNet) using the cross-entropy loss.
 889 During pre-training, the model is trained over 200000 gradient steps using SGD with Nesterov mo-
 890 mentum and weight decay with a batch size of 128. The learning rate, momentum coefficient, and
 891 weight decay penalty are set to 0.01, 0.9, and 0.0005, respectively. Additionally, we use a multistep
 892 learning rate schedule, decaying the learning rate by a factor of 0.1 at the following gradient steps
 893 $\{100000, 150000, 175000, 190000\}$. Note that in our ResNet-12 experiments on *mini*-ImageNet,
 894 CIFAR-FS, and FC-100, we observed some meta-level overfitting, which we hypothesize is due
 895 to the expressive network architecture coupled with the relatively small datasets. We found that
 896 increasing the pre-training weight decay to 0.01 in these experiments resolved the issue.
 897

898 A.3.2 META-LEARNING SETTINGS

899 In both the meta-training and meta-testing phases, we adhere to the standard hyperparameter values
 900 from the literature (Finn et al., 2017). In the outer loop, our algorithm is trained over $S^{meta} =$
 901 30,000 meta-gradient steps using Adam with a meta learning rate of $\eta = 0.00001$. As with previous
 902 works, our models are trained with a meta-batch size of 2 for 5-shot classification and 4 for 1-shot
 903 classification. In the inner loop, the models are adapted using $S^{base} = 5$ base-gradient steps using
 904 SGD with Nesterov momentum and weight decay. The learning rate, momentum coefficient, and
 905 weight decay penalty are set to 0.01, 0.9, and 0.0005 respectively. During meta-testing, the same
 906 inner loop hyperparameter settings are employed, and the final model is evaluated over 600 tasks
 907 sampled from $D^{testing}$. Notably, unlike prior work (Antoniou et al., 2019; Baik et al., 2020; 2021;
 908 2023), we do not ensemble the top 3 or 5 performing models from the same run, as this significantly
 909 increases the number of parameters and expressive power of the final models.
 910

911 A.3.3 NPBML SETTINGS

912 In our experiments, the backbone encoders z_θ (*i.e.*, 4-CONV and ResNet-12) are modified in two key
 913 ways during meta-learning: (1) with linear projection (warp) layers ω for preconditioning gradients,
 914 and (2) feature-wise linear modulation layers $FiLM_\psi$ for task adaptation. In all our experiments,
 915 we modify only the last (*i.e.*, fourth) module with warp and FiLM layers based on the findings from
 916 (Raghu et al., 2019), which showed that features near the start of the network are primarily reused
 917 and do not require adaptation. Additionally, we freeze all preceding modules in both the inner and

918 outer loops, which significantly reduces the storage and memory footprint of the proposed method,
 919 with no noticeable effect on the performance.

920 Following the recommendations of Flennerhag et al. (2020), warp layers are inserted after the main
 921 convolutional stack and before the residual connection ends in the ResNet modules. Regarding
 922 the FiLM layers, they are inserted after each batch normalization layer in a manner similar to the
 923 backbone encoders used in (Requeima et al., 2019). Where for convolutional layers, we first pool
 924 in the spatial dimension to obtain the average activation scores per map, and then flatten in the
 925 depth/filter dimension before processing $FiLM_\psi$ to obtain γ and β . As discussed in Section 3.4,
 926 the loss networks \mathcal{L}^S , \mathcal{L}^Q , and \mathcal{R} also use FiLM layers, these are interleaved between each of the
 927 linear layers as shown in Figure 4.

928 B FURTHER DISCUSSION

929 The proposed NPBML framework is a highly flexible and general approach to gradient-based meta-
 930 learning. As a result of this generality, many existing meta-learning algorithms are closely related
 931 to NPBML, with some potentially arising as special cases of NPBML. In what follows, we discuss
 932 some salient examples:

- 933 • **MAML++**: In (Antoniou et al., 2019), the inner optimization of MAML is modified by
 934 learning per-step batch normalization running statistics, weights and biases, and layer-wise
 935 learning rates. This is closely related to the FiLM layers used in NPBML, as they also
 936 output meta-learned scale and shift parameters, *i.e.*, γ_ψ and β_ψ , which improve robustness
 937 to varying feature distributions. Furthermore, as shown by Equation (16), NPBML implicitly
 938 learns a layer-wise learning rate through P_ω .
- 939 • **SCA**: In (Antoniou & Storkey, 2019), the authors propose replacing \mathcal{L}^{base} with a meta-
 940 learned loss function \mathcal{M}_ϕ similar to NPBML. The primary difference between their param-
 941 eterization and ours is that their meta-learned loss function is represented as a 1D dilated
 942 convolutional neural network with DenseNet-style residual connectivity. In our experiments,
 943 we found that a simple feedforward ReLU network is sufficiently expressive to obtain high
 944 performance as shown in the ablation study in Table 3.
- 945 • **MeTAL**: In (Baik et al., 2021), a set of task-adaptive loss functions, represented as small
 946 feedforward networks are meta-learned for each step in the inner optimization. Task-
 947 adaptivity is achieved by using a separate set of networks to generate FiLM scale and shift pa-
 948 rameters which are applied directly to ϕ instead of the network’s activations. Unlike MeTAL,
 949 NPBML does not maintain separate networks for each step; thus, it can be applied to more
 950 gradient steps at meta-testing time than it was initially trained to do.
- 951 • **ALFA**: In (Baik et al., 2023), MAML’s inner optimization is modified with meta-learned
 952 layer-wise learning rate values and weight decay coefficients. As shown in Equation (16),
 953 NPBML implicitly learns layer-wise learning rate values. Additionally, since the learned
 954 loss function \mathcal{M}_ϕ in NPBML is conditioned on $\theta_{i,j}$, our method also implicitly learns a
 955 scalar weight decay values, since $\exists \lambda \exists \phi : \theta_{i,j} - \nabla_{\theta_{i,j}} (\mathcal{L}^{base} + \lambda \|\theta\|^2) \approx \theta_{i,j} - \nabla_{\theta_{i,j}} \mathcal{M}_\phi$.
- 956 • **MetaSGD**: In (Li et al., 2017), a parameter-wise learning rate is meta-learned in tandem with
 957 the parameter initialization. This corresponds to using an identical inner update rule to Equa-
 958 tion (6) where $\mathcal{M}_{(\phi,\psi)} = \mathcal{L}^{base}$ and $P_{(\omega,\psi)} = \text{diag}(\omega)$, *i.e.*, no meta-learned loss function,
 959 and a diagonal matrix used instead of a block diagonal matrix for gradient preconditioning.
- 960 • **WarpGrad**: In (Flennerhag et al., 2020), a full gradient preconditioning matrix is induced by
 961 modifying the meta-learned linear projection layers employed in T-Nets (Lee & Choi, 2018)
 962 with non-linear activations. Furthermore, a trajectory-agnostic meta-objective is used. As
 963 NPBML utilizes T-Net style gradient preconditioning, simply introducing non-linear activa-
 964 tion functions and replacing the outer objective \mathcal{L}^{meta} in Equation (4) would be sufficient to
 965 recover WarpGrad’s gradient preconditioning.
- 966 • **ModGrad**: In (Simon et al., 2020), a low-rank gradient preconditioning matrix is meta-
 967 learned. This is, in essence, an identical inner update rule to Equation (6) where
 968 $\mathcal{M}_{(\phi,\psi)} = \mathcal{L}^{base}$ and $P_{(\omega,\psi)} = \omega_1 \otimes \omega_2$, where $\omega_1 \in \mathbb{R}^{in \times r}$, $\omega_2 \in \mathbb{R}^{r \times out}$, and \otimes is the
 969 outer product between the two low-rank matrices with rank r .