

A Relationship between iMAML and Prior Algorithms

The presented iMAML algorithm has close connections, as well as notable differences, to a number of related algorithms like MAML [15], first-order MAML, and Reptile [43]. Conventionally, these algorithms do not consider any explicit regularization in the inner-level and instead rely on early stopping, through only a few gradient descent steps. In our problem setting described in Eq. 4, we consider an explicitly regularized inner-level problem (refer to discussion in Section 2.2). We describe the connections between the algorithms in this explicitly regularized setting below.

MAML. The MAML algorithm first invokes an iterative algorithm to solve the inner optimization problem (see definition 1). Subsequently, it backpropagates through the path of the optimization algorithm to update the meta-parameters as:

$$\theta^{k+1} = \theta^k - \eta \frac{1}{M} \sum_{i=1}^M d_{\theta} \mathcal{L}_i(\text{Alg}_i(\theta^k)).$$

Since $\text{Alg}_i(\theta)$ approximates $\text{Alg}_i^*(\theta)$, it can be viewed that both MAML and iMAML intend to perform the same idealized update in Eq. 5. However, they perform the meta-gradient computation very differently. MAML backpropagates through the path of an iterative algorithm, while iMAML computes the meta-gradient through the implicit Jacobian approach outlined in Section 3.1 (see Figure 1 for a visual depiction). As a result, iMAML can be vastly more efficient in memory while having lesser or comparable computational requirements. It also allows for higher order optimization methods and non-differentiable components.

First-order MAML ignores the effect of meta-parameters θ on task parameters $\{\phi_i\}$ in the meta-gradient computation and updates the meta-parameters as:

$$\theta^{k+1} = \theta^k - \eta \frac{1}{M} \sum_{i=1}^M \nabla_{\phi} \mathcal{L}_i(\phi_i) \big|_{\phi_i = \text{Alg}_i(\theta^k)}$$

Note that iMAML strictly generalizes this, since first-order MAML is simply iMAML when the conjugate gradient procedure is not invoked (or corresponds to 0 steps of CG). Thus, iMAML allows for an easy way to interpolate from first-order MAML to the full MAML algorithm.

Reptile [43], similar to first-order MAML, ignores the dependence of task-parameters on meta-parameters. However, instead of following the gradients at $\phi_i = \text{Alg}_i(\theta^k)$, Reptile uses the task-parameters as targets and slowly moves meta-parameters towards them:

$$\theta^{k+1} = \theta^k - \eta \frac{1}{M} \sum_{i=1}^M (\theta^k - \phi_i).$$

From the proximal point equation in the proof of Lemma 1, we have $\phi_i = \theta^k - \frac{1}{\lambda} \nabla_{\phi} \mathcal{L}_i(\phi_i)$, using which we see that the Reptile equation becomes: $\theta^{k+1} = \theta^k - \frac{\eta}{\lambda M} \sum_{i=1}^M \nabla_{\phi} \mathcal{L}_i(\phi_i)$. Thus, Reptile and first-order MAML are identical in our problem formulation up to the choice of learning rate. Making the regularization explicit allows us to illustrate this equivalence.

B Optimization Preliminaries

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$. A function f is B Lipschitz (or B -bounded gradient norm) if for all $x \in \mathbb{R}^d$

$$\|\nabla f(x)\| \leq B.$$

Similarly, we say that a matrix valued function $M : \mathbb{R}^d \times \mathbb{R}^{d'} \rightarrow \mathbb{R}$ is ρ -Lipschitz if

$$\|M(x) - M(x')\| \leq \rho \|x - x'\|,$$

where $\|\cdot\|$ denotes the spectral norm.

We say that f is L -smooth if for all $x, x' \in \mathbb{R}^d$

$$\|\nabla f(x) - \nabla f(x')\| \leq L \|x - x'\|$$

and that f is μ -strongly convex if f is convex and if for all $x, x' \in \mathbb{R}^d$,

$$\|\nabla f(x) - \nabla f(x')\| \geq \mu \|x - x'\|.$$

We will make use of the following black-box complexity of first-order gradient methods for minimizing strongly convex and smooth functions.

Lemma 2. (*δ -approximate solver; see [9]*) Suppose f is a function that is L -smooth and μ strongly convex. Define $\kappa := L/\mu$, and let $x^* = \operatorname{argmin} f(x)$. Nesterov's accelerated gradient descent can be used to find a point x such that:

$$\|x - x^*\| \leq \delta$$

using a number of gradient computations of f that is bounded as follows:

$$\# \text{ gradient computations of } f(\cdot) \leq 2\sqrt{\kappa} \log \left(2\kappa \frac{\|x^*\|}{\delta} \right).$$

C Review: Time and Space Complexity of Hessian-Vector Products

We briefly discuss the time and space complexity of Hessian-vector product computation using the reverse mode of automatic differentiation. The reverse mode of automatic differentiation [6, 22] is the widely used method for automatic differentiation in modern software packages like TensorFlow and PyTorch [7]. Recall that for a differentiable function $f(x)$, the reverse mode of automatic differentiation computes $\nabla f(x)$ in time that is no more than a factor of 5 of the time it takes to compute $f(x)$ itself (see [22] for review). As our algorithm makes use of Hessian vector products, we will make use of the following assumption as to how Hessian vector products will be computed when executing Algorithm 2.

Assumption 1. (*Complexity of Hessian-vector product*) We assume that the time to compute the Hessian-vector product $\nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi) \mathbf{v}$ is no more than a (universal) constant over the time used to compute $\nabla \hat{\mathcal{L}}_i(\phi)$ (typically, this constant is 5). Furthermore, we assume that the memory used to compute the Hessian-vector product $\nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi) \mathbf{v}$ is no more than twice the memory used when computing $\nabla \hat{\mathcal{L}}_i(\phi)$. This assumption is valid if the reverse mode of automatic differentiation is used to compute Hessian vector products (see [21]).

A few remarks about this assumption are in order. With regards to computation, first observe that the gradient of the scalar function $\nabla_{\phi} \hat{\mathcal{L}}_i(\phi)^{\top} \mathbf{v}$ is the desired Hessian vector product $\nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi) \mathbf{v}$. Thus computing the Hessian vector product using the reverse mode is within a constant factor of computing the function itself, which is simply the cost of computing $\nabla \hat{\mathcal{L}}_i(\phi)^{\top} \mathbf{v}$. The issue of memory is more subtle (see [21]), which we now discuss. The memory used to compute the gradient of a scalar cost function $f(x)$ using the reverse mode of auto-differentiation is proportional to the size of the computation graph; precisely, the memory required to compute the gradient is equal to the total space required to store all the intermediate variables used when computing $f(x)$. In practice, this is often much larger than the memory required to compute $f(x)$ itself, due to that all intermediate variables need not be simultaneously stored in memory when computing $f(x)$. However, for the special case of computing the gradient of the function $f(\phi) = \nabla_{\phi} \hat{\mathcal{L}}_i(\phi)^{\top} \mathbf{v}$, the factor of 2 in the memory bound is a consequence of the following reason: first, using the reverse mode to compute $f(\phi)$ means we already have stored the computation graph of $\hat{\mathcal{L}}_i(\phi)$ itself. Furthermore, the size of the computation graph for computing $f(\phi) = \nabla_{\phi} \hat{\mathcal{L}}_i(\phi)^{\top} \mathbf{v}$ is essentially the same size as the computation graph of $\hat{\mathcal{L}}_i(\phi)$. This leads to the factor of 2 memory bound; see Griewank [21] for further discussion.

D Additional Discussion About Compute and Memory Complexity

Our main complexity results are summarized in Table 1. For these results, we consider two notions of error that are subtly different, which we explicitly define below. Let \mathbf{g}_i be the computed meta-gradient for task \mathcal{T}_i . Then, the errors we consider are:

Definition 3. *Exact-solve error (our notion of error):* Our goal is to accurately compute the gradient of $F(\theta)$ as defined in Equation 4, where $\text{Alg}_i^*(\theta)$ is an exact algorithm. Specifically, we seek to compute a \mathbf{g}_i such that:

$$\|\mathbf{g}_i - \mathbf{d}_{\theta} \mathcal{L}_i(\text{Alg}_i^*(\theta))\| \leq \epsilon$$

where ϵ is the error in the gradient computation.

Definition 4. *Approx-solve error:* Here we suppose that Alg_i computes a δ -accurate solution to the inner optimization problem over G_i in Eq. 4, i.e. that Alg_i satisfies $\|\text{Alg}_i(\theta) - \text{Alg}_i^*(\theta)\| \leq \delta$, as per definition 1. Then the objective is to compute a \mathbf{g} such that:

$$\|\mathbf{g} - \mathbf{d}_{\theta} \mathcal{L}_i(\text{Alg}_i(\theta))\| \leq \epsilon$$

where ϵ is the error in the gradient computation of $\mathbf{d}_{\theta} \mathcal{L}_i(\text{Alg}_i(\theta))$. Subtly, note that the gradient is with respect to the δ -approximate algorithm, as opposed to using Alg_i^* .

For the complexity results, we assume that MAML invokes Alg_i to get a δ -approximate solution for inner problem (recall definition 1). The exact-solve error for MAML is not known in the literature; in particular, even as $\delta \rightarrow 0$ it is not evident if the approx-solve solution tends to the exact-solve solution, unless further regularity conditions are imposed. The approx-solve error for MAML is 0, ignoring finite-precision and numerical issues, since it backpropagates through the path. Truncated backprop [53] also invokes Alg_i to obtain a δ -approximate solution but instead performs a truncated or partial back-propagation so that it uses a smaller number of iterations when computing the gradient through the path of $\text{Alg}_i(\theta)$. Exact-solve error for truncated backprop is also not known, but a small approx-solve error can be obtained with less memory than full back-prop. We use Prop 3.1 of Shaban et al. [53] to provide a guarantee that leads to an ϵ -accurate approximation of the full-backprop (i.e. MAML) gradient. It is not evident how accurate the truncated procedure is when an accelerated method is used instead. Finally, our iMAML algorithm also invokes an approximate solver Alg_i rather than Alg_i^* . However, importantly, we guarantee a small exact-solve error even though we do not require access to Alg_i^* . Furthermore, the iMAML algorithm also requires substantially less memory. Up to small constant factors, it only utilizes the memory required for computing a single gradient of $\hat{\mathcal{L}}_i(\cdot)$.

E Theoretical Results and Proofs

Lemma 1, restated. Consider $\text{Alg}_i^*(\theta)$ as defined in Eq. 4 for task \mathcal{T}_i . Let $\phi_i = \text{Alg}_i^*(\theta)$ be the result of $\text{Alg}_i^*(\theta)$. If $\left(\mathbf{I} + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i)\right)$ is invertible, then the derivative Jacobian is

$$\frac{d\text{Alg}_i^*(\theta)}{d\theta} = \left(\mathbf{I} + \frac{1}{\lambda} \nabla_{\phi}^2 \hat{\mathcal{L}}_i(\phi_i)\right)^{-1}.$$

Proof. We drop the task i subscripts in the proof for convenience. Since $\phi = \text{Alg}^*(\theta)$ is the minimizer of $G(\phi', \theta)$ in Eq. 4, the stationary point conditions imply that

$$\nabla_{\phi'} G(\phi', \theta) |_{\phi'=\phi} = 0 \implies \nabla \hat{\mathcal{L}}(\phi) + \lambda(\phi - \theta) = 0 \implies \phi = \theta - \frac{1}{\lambda} \nabla \hat{\mathcal{L}}(\phi),$$

which is an implicit equation that often arises in proximal point methods. When the derivative exists, we can differentiate the above equation to obtain:

$$\frac{d\phi}{d\theta} = \mathbf{I} - \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \frac{d\phi}{d\theta} \implies \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi)\right) \frac{d\phi}{d\theta} = \mathbf{I}.$$

which completes the proof. \square

Recall that:

$$G_i(\phi', \theta) := \hat{\mathcal{L}}_i(\phi') + \frac{\lambda}{2} \|\phi' - \theta\|^2.$$

Assumption 2. (Regularity conditions) Suppose the following holds for all tasks i :

1. $\mathcal{L}_i(\cdot)$ is B Lipshitz and L smooth.

2. For all θ , $G_i(\cdot, \theta)$ is both a β -smooth function and a μ -strongly convex function. Define:

$$\kappa := \frac{\beta}{\mu}.$$

3. $\hat{\mathcal{L}}_i(\cdot)$ is ρ -Lipshitz Hessian, i.e. $\nabla^2 \hat{\mathcal{L}}_i(\cdot)$ is ρ -Lipshitz.

4. For all θ , suppose the arg-minimizer of $G_i(\cdot, \theta)$ is unique and bounded in a ball of radius D , i.e. for all θ ,

$$\|\mathcal{A}lg_i^*(\theta)\| \leq D.$$

Lemma 3. (Implicit Gradient Accuracy) Suppose Assumption 2 holds. Fix a task i . Suppose that ϕ_i satisfies:

$$\|\phi_i - \mathcal{A}lg_i^*(\theta)\| \leq \delta$$

and that \mathbf{g}_i satisfies:

$$\|\mathbf{g}_i - \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}_i(\phi) \right)^{-1} \nabla_{\phi} \mathcal{L}_i(\phi)\| \leq \delta'.$$

Assuming that $\delta < \mu/(2\rho)$, we have that:

$$\|\mathbf{g}_i - \mathbf{d}_{\theta} \mathcal{L}_i(\mathcal{A}lg_i^*(\theta))\| \leq \left(2 \frac{\lambda \rho}{\mu^2} B + \frac{\lambda L}{\mu} \right) \delta + \delta'$$

Proof. First, observe that:

$$\mathbf{d}_{\theta} \mathcal{L}_i(\mathcal{A}lg_i^*(\theta)) = \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}_i(\mathcal{A}lg_i^*(\theta)) \right)^{-1} \nabla_{\phi} \mathcal{L}_i(\mathcal{A}lg_i^*(\theta))$$

For notational convenience, we drop the i subscripts within the proof. We have:

$$\begin{aligned} & \|\mathbf{d}_{\theta} \mathcal{L}(\mathcal{A}lg^*(\theta)) - \mathbf{g}\| \\ & \leq \|\mathbf{d}_{\theta} \mathcal{L}(\mathcal{A}lg^*(\theta)) - \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} \nabla_{\phi} \mathcal{L}(\phi)\| + \delta' \\ & \leq \|\mathbf{d}_{\theta} \mathcal{L}(\mathcal{A}lg^*(\theta)) - \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} \nabla_{\phi} \mathcal{L}(\mathcal{A}lg^*(\theta))\| + \\ & \quad \left\| \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} (\nabla_{\phi} \mathcal{L}(\mathcal{A}lg^*(\theta)) - \nabla_{\phi} \mathcal{L}(\phi)) \right\| + \delta' \end{aligned}$$

where the first inequality uses the triangle inequality.

We now bound each of these terms. For the second term,

$$\begin{aligned} & \left\| \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} (\nabla_{\phi} \mathcal{L}(\mathcal{A}lg^*(\theta)) - \nabla_{\phi} \mathcal{L}(\phi)) \right\| \\ & \leq \left\| \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} \right\| \|\nabla_{\phi} \mathcal{L}(\mathcal{A}lg^*(\theta)) - \nabla_{\phi} \mathcal{L}(\phi)\| \\ & \leq \lambda L \left\| \left(\lambda \mathbf{I} + \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} \right\| \|\mathcal{A}lg^*(\theta) - \phi\| \\ & = \lambda L \|\nabla_{\phi}^2 G(\phi, \theta)^{-1}\| \|\mathcal{A}lg^*(\theta) - \phi\| \\ & \leq \frac{\lambda L}{\mu} \delta \end{aligned}$$

where we the second inequality uses that $\nabla_{\phi} \mathcal{L}$ is L -smooth and the final inequality uses that G is μ strongly convex.

For the first term, we have:

$$\begin{aligned}
& \|d_{\theta}\mathcal{L}(\mathcal{A}lg^*(\theta)) - \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi)\right)^{-1} \nabla_{\phi}\mathcal{L}(\mathcal{A}lg^*(\theta))\| \\
&= \left\| \left(\left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\mathcal{A}lg^*(\theta))\right)^{-1} - \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}(\phi)\right)^{-1} \right) \nabla_{\phi}\mathcal{L}(\mathcal{A}lg^*(\theta)) \right\| \\
&\leq \lambda \left\| \left(\lambda \mathbf{I} + \nabla^2 \hat{\mathcal{L}}(\mathcal{A}lg^*(\theta)) \right)^{-1} - \left(\lambda \mathbf{I} + \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} \right\| B,
\end{aligned}$$

using that $\nabla_{\phi}\mathcal{L}$ is B Lipshitz. Now let

$$\Delta := \nabla^2 \hat{\mathcal{L}}(\mathcal{A}lg^*(\theta)) - \nabla^2 \hat{\mathcal{L}}(\phi), \quad M := \nabla_{\phi}^2 G(\phi, \theta) = \lambda \mathbf{I} + \nabla^2 \hat{\mathcal{L}}(\phi)$$

Due to that $\nabla^2 \hat{\mathcal{L}}(\cdot)$ is Lipshitz Hessian, $\|\Delta\| \leq \rho\delta$. Also, by our assumption on δ , we have that:

$$\|M^{-1}\Delta\| \leq \|\Delta\|/\mu \leq \rho\delta/\mu \leq 1/2,$$

which implies that $\|(\mathbf{I} + M^{-1}\Delta)^{-1}\| \leq 2$. Hence,

$$\begin{aligned}
& \left\| \left(\lambda \mathbf{I} + \nabla^2 \hat{\mathcal{L}}(\mathcal{A}lg^*(\theta)) \right)^{-1} - \left(\lambda \mathbf{I} + \nabla^2 \hat{\mathcal{L}}(\phi) \right)^{-1} \right\| \\
&= \left\| (M + \Delta)^{-1} - M^{-1} \right\| \\
&\leq \|M^{-1}\| \left\| (\mathbf{I} + M^{-1}\Delta)^{-1} - \mathbf{I} \right\| \\
&= \|M^{-1}\| \left\| (\mathbf{I} + M^{-1}\Delta)^{-1} (\mathbf{I} - (\mathbf{I} + M^{-1}\Delta)) \right\| \\
&\leq \|M^{-1}\| \left\| (\mathbf{I} + M^{-1}\Delta)^{-1} \right\| \|M^{-1}\Delta\| \\
&\leq \frac{1}{\mu} \cdot 2 \cdot \frac{\rho\delta}{\mu} = 2 \frac{\rho}{\mu^2} \delta.
\end{aligned}$$

The proof is completed by substitution. \square

Theorem 2. (Approximate Implicit Gradient Computation) Suppose Assumption 2 holds. Fix a task i . Let

$$B_1 := 2 \frac{\lambda\rho}{\mu^2} B + \frac{\lambda L}{\mu}$$

Suppose Nesterov's accelerated gradient descent algorithm is used to compute ϕ (as desired in Algorithm 2), using a number of iterations that is:

$$2\sqrt{\kappa} \log \left(8\kappa D \left(\frac{B_1}{\epsilon} + \frac{\rho}{\mu} \right) \right)$$

and suppose Nesterov's accelerated gradient descent algorithm (or the conjugate gradient algorithm²) is used to compute \mathbf{g}_i using a number of iterations that is:

$$2\sqrt{\kappa} \log \left(4\kappa \frac{(\lambda/\mu)B}{\epsilon} \right).$$

We have that:

$$\|\mathbf{g}_i - d_{\theta}\mathcal{L}_i(\mathcal{A}lg_i^*(\theta))\| \leq \epsilon.$$

Proof. The result will follow from the guarantees in Lemma 2. Specifically, let us set $\delta = \min\{\epsilon/(2B_1), \mu/(2\rho)\}$ and $\delta' = \epsilon/2$. To ensure the bound of δ , by Lemma 3, it suffices to use a number of iterations that is bounded by:

$$2 \log \left(2\kappa \frac{\|D\|}{\delta} \right) \leq 2\sqrt{\kappa} \log \left(8\kappa D \left(\frac{B_1}{\epsilon} + \frac{\rho}{\mu} \right) \right)$$

²The conjugate gradient descent algorithm also suffices and give a slightly improved iteration complexity in terms of log factors.

To ensure the bound of δ' , the algorithm will be solving the sub-problem in Equation 7. First observe that in the context of in Lemma 2, note that $\|x^*\| = \left\| \left(\mathbf{I} + \frac{1}{\lambda} \nabla^2 \hat{\mathcal{L}}_i(\phi) \right)^{-1} \nabla \mathcal{L}_i(\phi) \right\| \leq (\lambda/\mu)B$, and so it suffices to use a number of iterations that is bounded by:

$$2 \log \left(2\kappa \frac{\|x^*\|}{\delta} \right) \leq 2 \log \left(4\kappa \frac{(\lambda/\mu)B}{\epsilon} \right),$$

which completes the proof. \square

Finally, we present a corollary of previous theorem that shows that iMAML finds approximate stationary points due to controllable error in gradient computation.

Corollary 1. (*iMAML finds stationary points*) Suppose the conditions of Theorem 1 hold and that $F(\cdot)$ is an L_F smooth function. Then the implicit MAML algorithm (Algorithm 1), when the batch size is M (so that we are doing gradient descent), will find a point θ such that:

$$\|\nabla F(\theta)\| \leq \epsilon$$

in a number of calls to *Implicit-Meta-Gradient* that is at most $\frac{4ML_f(F(0) - \min_{\theta} F(\theta))}{\epsilon^2}$. Furthermore, the total number of gradient computations (of $\nabla \hat{\mathcal{L}}_i$) is at most $\tilde{O} \left(M \sqrt{\kappa \frac{L_f(F(0) - \min_{\theta} F(\theta))}{\epsilon^2}} \log \left(\frac{\text{poly}(\kappa, D, B, L, \rho, \mu, \lambda)}{\epsilon} \right) \right)$, and only $\tilde{O}(\text{Mem}(\nabla \hat{\mathcal{L}}_i))$ memory is required throughout.

F Experiment Details

Here, we provide additional details of the experimental set-up for the experiments in Section 4. All training runs were conducted on a single NVIDIA (Titan Xp) GPU.

F.1 Synthetic Experiments

For the synthetic experiments, we consider a linear regression problem. We consider parametric models of the form $h_{\phi}(\mathbf{x}) = \phi^T \mathbf{x}$, where \mathbf{x} can either be the raw inputs or features (e.g. Fourier features) of the input. For task \mathcal{T}_i , we can equivalently write a quadratic objective that represents the task loss as:

$$\hat{\mathcal{L}}_i(\phi) = \frac{1}{2} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_i^{\text{tr}}} [\|h_{\phi}(\mathbf{x}) - \mathbf{y}\|^2] = \frac{1}{2} \phi^T A_i \phi + \phi^T b_i,$$

where $A_i = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_i^{\text{tr}}} [\mathbf{x}\mathbf{x}^T]$ and $b_i = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_i^{\text{tr}}} [\mathbf{x}^T \mathbf{y}]$. Thus, the inner level objective and corresponding minimizer can be written as:

$$G_i(\phi', \theta) = \frac{1}{2} \phi'^T A_i \phi' + \phi'^T b_i + \frac{\lambda}{2} (\phi' - \theta)^T (\phi' - \theta)$$

$$\text{Alg}_i^*(\theta) = (A_i + \lambda \mathbf{I})^{-1} (\lambda \theta - b_i)$$

Thus, the exact meta-gradient can be written as

$$\mathbf{d}_{\theta} \mathcal{L}_i(\text{Alg}_i^*(\theta)) = \lambda (A_i + \lambda \mathbf{I})^{-1} \nabla_{\phi} \mathcal{L}_i(\theta) \big|_{\phi = \text{Alg}_i^*(\theta)}.$$

We compare this gradient with the gradients computed by the iMAML and MAML algorithms. We considered the case of $\mathbf{x} \in \mathbb{R}^{50}$, $\mathbf{y} \in \mathbb{R}$, $\lambda = 5.0$, and $\kappa = 50$, for the presented results.

F.2 Omniglot and Mini-ImageNet experiments

We follow the standard training and evaluation protocol as in prior works [50, 57, 15].

Omniglot Experiments The GD version of iMAML uses 16 gradient steps for 5-way 1-shot and 5-way 5-shot settings, and 25 gradient steps for 20-way 1-shot and 20-way 5-shot settings. A regularization strength of $\lambda = 2.0$ was used for both. 5 steps of conjugate gradient was used to compute the meta-gradient for each task in the mini-batch, and the meta-gradients were averaged before taking a step with the default parameters of Adam in the outer loop.

The Hessian-free version of MAML proceeds by using Hessian-free or Newton-CG method for solving the inner optimization problem (with respect to ϕ) with objective $G_i(\phi, \theta)$. This method proceeds by constructing a local quadratic approximation to the objective and approximately computing the Newton direction with conjugate gradient. 5 CG steps are used for this process in our experiments. This allows us to compute the search direction, following which a step size has to be picked. We pick the step size through line-search. This procedure of computing the approximate Newton direction and linesearch is repeated 3 times in our experiments to solve the inner optimization problem well.

Mini-ImageNet For the GD version of iMAML, 10 GD steps were used with regularization strength of $\lambda = 0.5$. Again, 5 CG steps are used to compute the meta-gradient. Similarly, in the Hessian-Free variant, we again use 5 CG steps to compute the search direction followed by line search. This process is repeated 3 times to solve the inner level optimization. Again, to compute the meta-gradient, 5 steps of CG are used.