



Abstract

Model Agnostic Meta Learning (MAML) is widely used to find a good initialization for a family of tasks. Despite its success, a critical challenge in MAML is to calculate the gradient w.r.t. the initialization of a long training trajectory for the sampled tasks, because the computation graph can rapidly explode and the computational cost is very expensive. To address this problem, we propose Adjoint MAML (A-MAML). We view gradient descent in the inner optimization as the evolution of an Ordinary Differential Equation (ODE). To efficiently compute the gradient of the validation loss w.r.t. the initialization, we use the adjoint method to construct a companion, backward ODE. To obtain the gradient w.r.t. the initialization, we only need to run the standard ODE solver twice --- one is forward in time that evolves a long trajectory of gradient flow for the sampled task; the other is backward and solves the adjoint ODE. We need not create or expand any intermediate computational graphs, adopt aggressive approximations, or impose proximal regularizers in the training loss. Our approach is cheap, accurate, and adaptable to different trajectory lengths. We demonstrate the advantage of our approach in both synthetic and real-world meta-learning tasks. The code is available at <https://github.com/shibo0li/Adjoint-MAML>.

Motivation

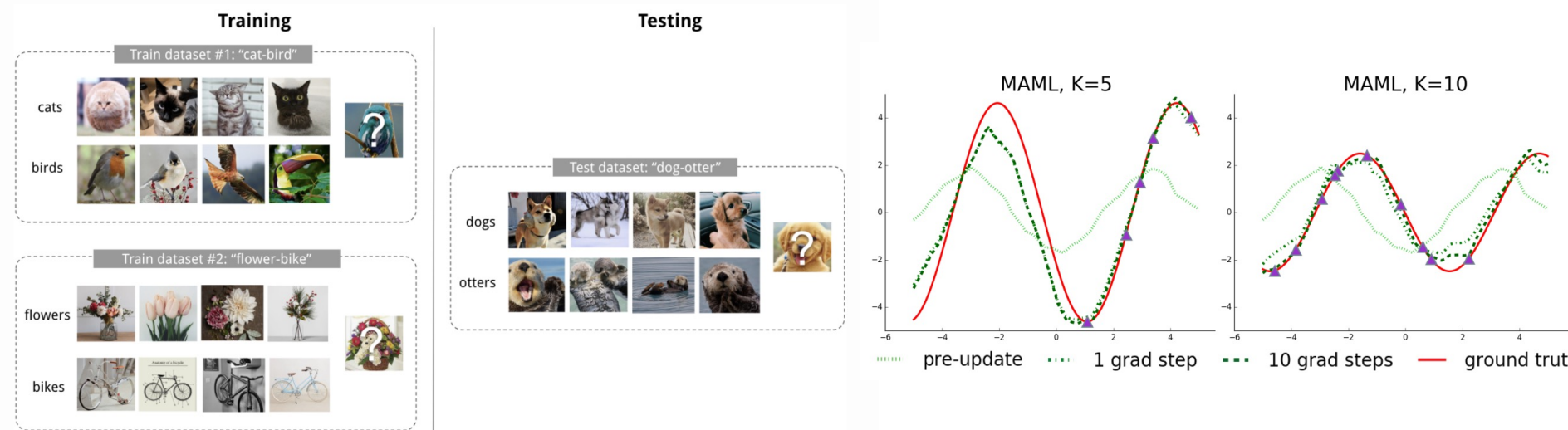
Meta-Learning:

- Standard (Deep) Machine Learning: **cheap, safe, easy** to collect large amount of data



- Data **Costly, Sensitive** Applications: Robotics, User personalization, etc.
- Meta Learning (Learning to Learn Fast)**: Enable **efficient** learning on **new tasks** with encoding **adaptable representations**

Examples: *Few-shot Classification; Few-shot Regression*



Meta-Learning Methods: *Metric based; Model based; Optimization based*

Mode-Agnostic Meta Learning

- A fairly *general* optimization algorithm
- Compatible with any model that learns through gradient descent

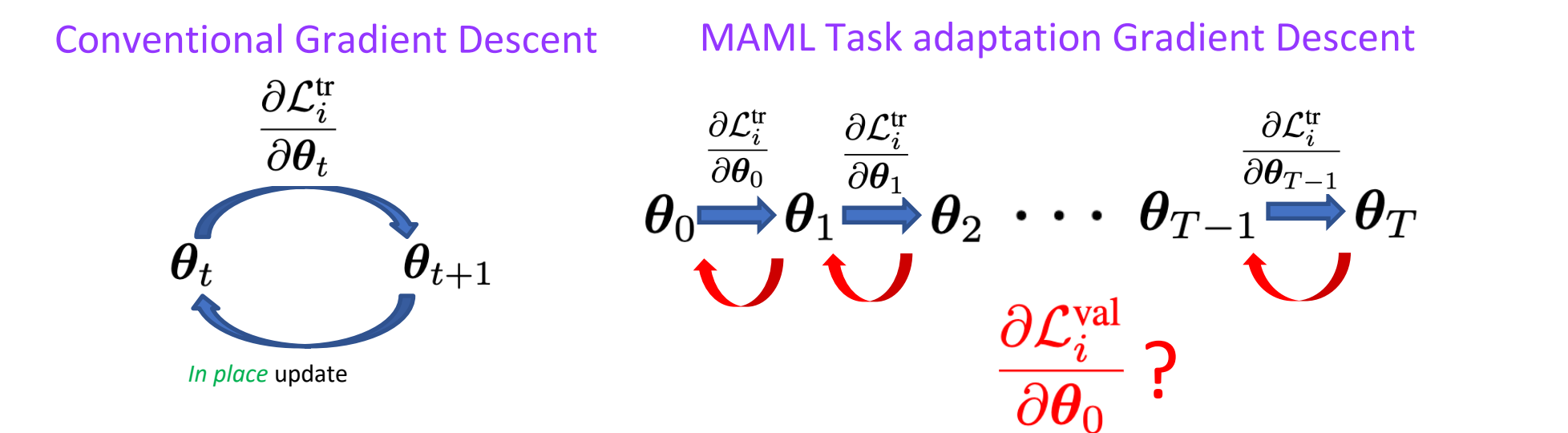
$$\theta^* = \arg \min_{\theta} \sum_{\tau_i \sim p(\tau)} \mathcal{L}_{\tau_i}^{(1)}(f_{\theta_i}) = \arg \min_{\theta} \sum_{\tau_i \sim p(\tau)} \mathcal{L}_{\tau_i}^{(1)}(f_{\theta - \alpha \nabla_{\theta} \mathcal{L}_{\tau_i}^{(0)}(f_{\theta})})$$

Task adaptation

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\tau_i \sim p(\tau)} \mathcal{L}_{\tau_i}^{(1)}(f_{\theta - \alpha \nabla_{\theta} \mathcal{L}_{\tau_i}^{(0)}(f_{\theta})})$$

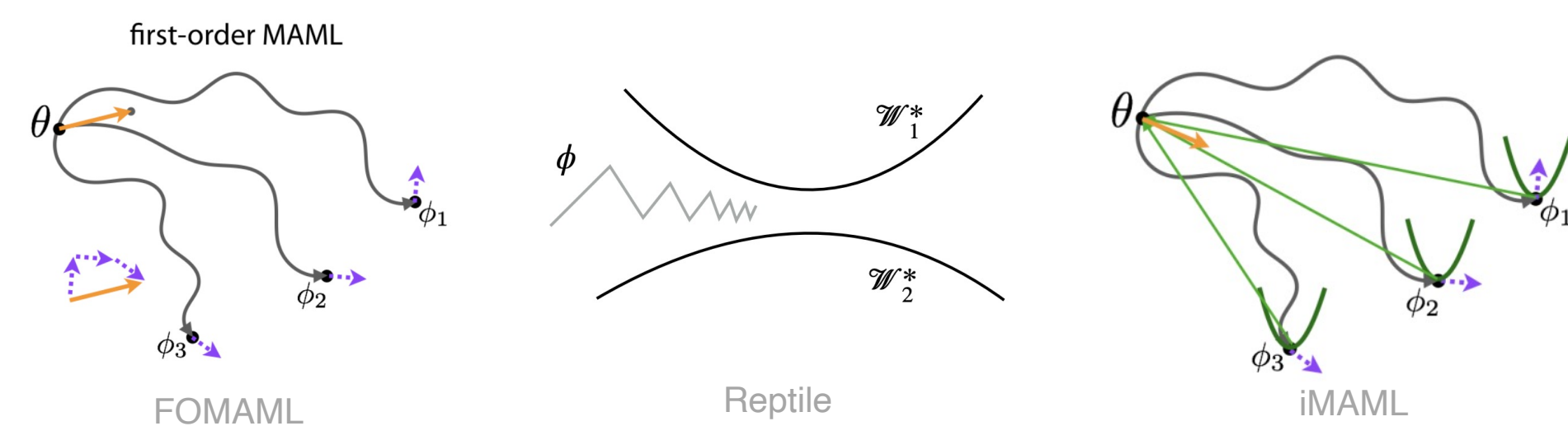
Meta update/optimization

- Exhaustive memory usage** due to saving the entire gradient descent trajectory.



Approximations of MAML

- First-order MAML; Reptile; Meta-Learning with Implicit Gradient (iMAML)



- One or few** inner steps is/are sufficient to approximate the meta gradients. (original MAML)
- Don't look back at all. Ignores the **outer level optimization**. (FOMAML, Reptile, iMAML)
- Strong (local) curvature** heuristic about the meta gradients (iMAML)

Our Contribution: Adjoint-MAML

- Forward ODE: Task adaptation (gradient flow)
- (Backward) Adjoint ODE: Dynamics of the meta-gradient
- Resurgence** the "ideal" MAML (trace back the **exact** task adaptation)
- Memory efficiency** for long adaptation trajectory while yet accurate meta gradients on validation

Methods

Adjoint Method:

$$\begin{cases} \frac{du}{dt} = f(u, t; \theta) \\ u(t_0) = u_0 \end{cases} \quad u(t) \in \mathbb{R}^N; f \in \mathbb{R}^N; \theta \in \mathbb{R}^P$$

$$\min_{\theta} J(u, \theta) = \min_{\theta} \int_0^T g(u, \theta) dt$$

For example: Quadratic Loss: $u^T Q u$

Goal: To Compute the total derivative

$$\frac{dJ}{d\theta} = \int_0^T \left(\frac{\partial g}{\partial \theta} + \frac{\partial g}{\partial u} \cdot \frac{\partial u}{\partial \theta} \right) dt$$

Auto-differentiation

$$\frac{dJ}{d\theta} = \int_0^T \left(\frac{dg}{d\theta} + \lambda(t) \frac{df}{d\theta} \right) dt + \lambda(0) \frac{du_0}{d\theta}$$

Solve P+1 ODE systems 😞

Solve **only 2** ODE systems, scale constantly 😊

Adjoint Model Agnostic Meta-Learning (AMAML):

Step 1: Task adaptation, Forward ODE (Gradient Flow)

$$\text{I.V.P} \begin{cases} u_n(0) = \theta, \\ \frac{du_n}{dt} = -\frac{\partial \mathcal{L}(u_n, \mathcal{D}_n^{tr})}{\partial u} \end{cases} \quad u_n(t + \alpha) \leftarrow u_n(t) - \alpha \frac{\partial \mathcal{L}(u_n, \mathcal{D}_n^{tr})}{\partial u}$$

Adjoint ODE

$$\text{Validation loss } J(\theta) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}(u_n(T), \mathcal{D}_n^{val})$$

Forward ODE

Step 2: Derived the Adjoint ODE

Target Constrained Optimization problem

$$\min_{\theta} \mathbb{E}[J(\theta)] = \min_{\theta} \frac{1}{N} \sum_n \mathcal{L}(u_n(T), \mathcal{D}_n^{val})$$

s.t. $\forall T_n \sim p(T)$

Lagrangian Relaxation

$$\hat{J}_n = J_n(u_n(T)) + \int_0^T \lambda(t)^T \left(f(u_n, \mathcal{D}_n^{tr}) - \frac{du_n}{dt} \right) dt$$

Sensitivity Cancellation

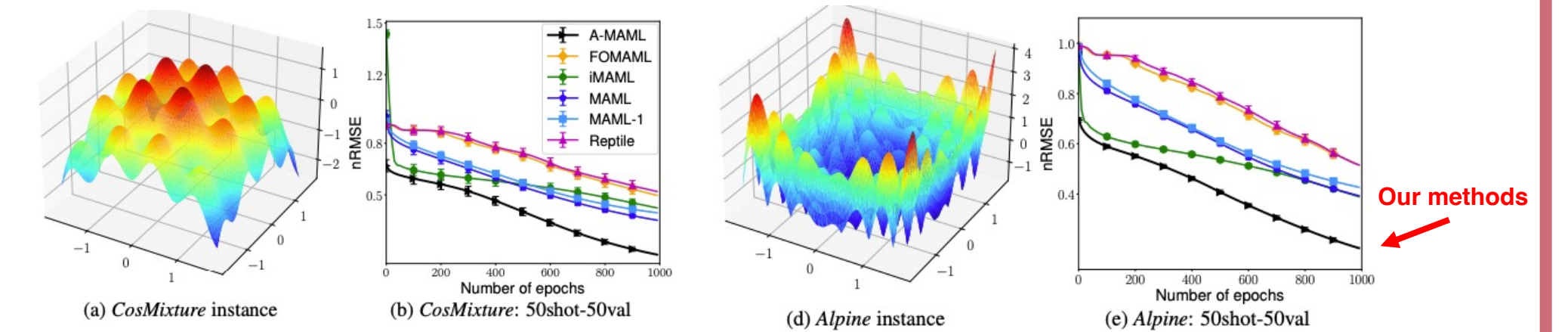
$$\frac{dJ_n}{d\theta} = \frac{\partial J_n}{\partial u_n(T)} \frac{du_n}{d\theta}(T) - \lambda(T)^T \frac{du_n}{d\theta}(T) + \lambda(0)^T \frac{du_n}{d\theta}(0) + \int_0^T \left\{ \lambda^T \frac{\partial f}{\partial u_n} \frac{du_n}{d\theta} + \left(\frac{d\lambda}{dt} \right)^T \frac{du_n}{d\theta} \right\} dt$$

$$\text{T.V.P} \begin{cases} \lambda(T) = \left(\frac{\partial J_n}{\partial u_n(T)} \right)^T \\ \left(\frac{d\lambda}{dt} \right)^T = -\lambda(t)^T \frac{\partial f}{\partial u_n} \end{cases}$$

$$\frac{\partial f}{\partial u_n} = H(u_n) = -\frac{\partial^2 \mathcal{L}(u_n, \mathcal{D}_n^{tr})}{\partial u_n^2}$$

Experiment

Evaluation 1: Regression Tasks



Evaluation 2: Collaborative Filtering

Method	Jester-1		MovieLens100K		MovieLens1M	
	10shot-15val	20shot-30val	10shot-15val	20shot-30val	10shot-15val	20shot-30val
A-MAML	0.074±0.005	0.027±0.002	0.053±0.005	0.023±0.003	0.094±0.008	0.035±0.004
iMAML	0.114±0.007	0.050±0.003	0.082±0.004	0.033±0.002	0.138±0.010	0.052±0.004
MAML	0.120±0.001	0.036±0.000	0.123±0.001	0.050±0.003	0.140±0.002	0.059±0.001
FOMAML	0.292±0.012	0.115±0.004	0.174±0.008	0.068±0.004	0.270±0.011	0.104±0.006
Reptile	0.270±0.012	0.106±0.004	0.166±0.008	0.063±0.003	0.266±0.011	0.101±0.006

Table 1: Meta-test error (nRMSE) with 50 inner gradient descent steps (MAML used 5 GD steps) The results were averaged over 100 tasks.

Evaluation 3: Image Classification

Method	Ominiglot	Mini-ImageNet
MAML	95.8 ± 0.3%	48.70 ± 1.84%
FOMAML	89.4 ± 0.5%	48.07 ± 1.75%
Reptile	89.43 ± 0.14%	49.97 ± 0.32%
iMAML-GD	94.46 ± 0.42%	48.96 ± 1.84%
iMAML-HF	96.18 ± 0.36%	49.30 ± 1.88%
A-MAML(T = 0.5)	96.36 ± 0.39%	49.43 ± 1.64%
A-MAML(T = 1.0)	96.79 ± 0.34%	49.47 ± 1.77%

Table 3: Meta-test accuracy for 20-way 1-shot on *Ominiglot* and 5-way 1-shot on *Mini-ImageNet*.

Evaluation 4: Efficiency Analysis

Evaluation 5: Trade-off of the Meta-Gradients

(Appendix A of the paper)

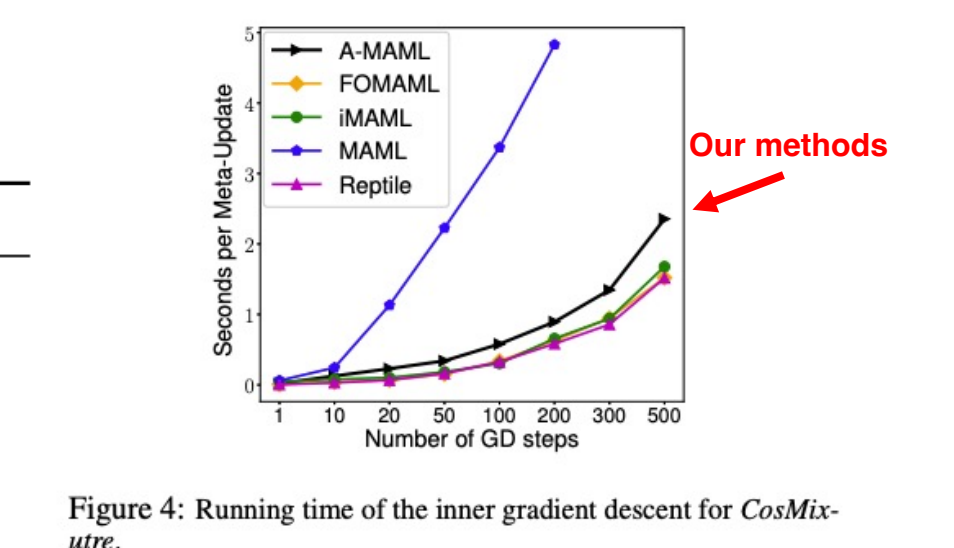


Figure 4: Running time of the inner gradient descent for *CosMixture*.

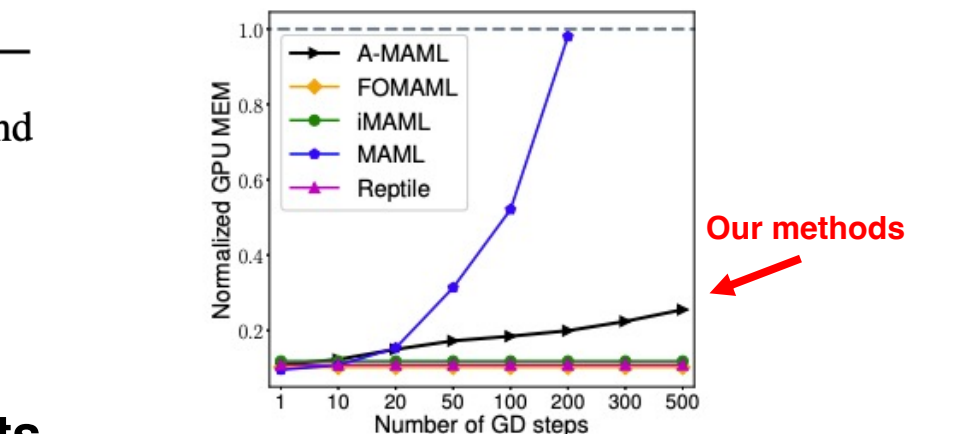


Figure 3: Normalized GPU usage in meta learning of *CosMixture* with 100shot-100validation. The dashed line indicates the capacity of available GPU memory.