

Second-order Optimization in MAML

□ Model-Agnostic Meta-Learning (MAML) [1]

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\tau_i} [\mathcal{L}_{\text{val}}^{\tau_i}(\theta - \alpha \nabla \mathcal{L}_{\text{tr}}^{\tau_i}(\theta))]$$

□ Meta-SGD [2]

$$\operatorname{argmin}_{\theta, \alpha} \mathbb{E}_{\tau_i} [\mathcal{L}_{\text{val}}^{\tau_i}(\theta - \alpha \nabla \mathcal{L}_{\text{tr}}^{\tau_i}(\theta))]$$

□ Meta-Curvature

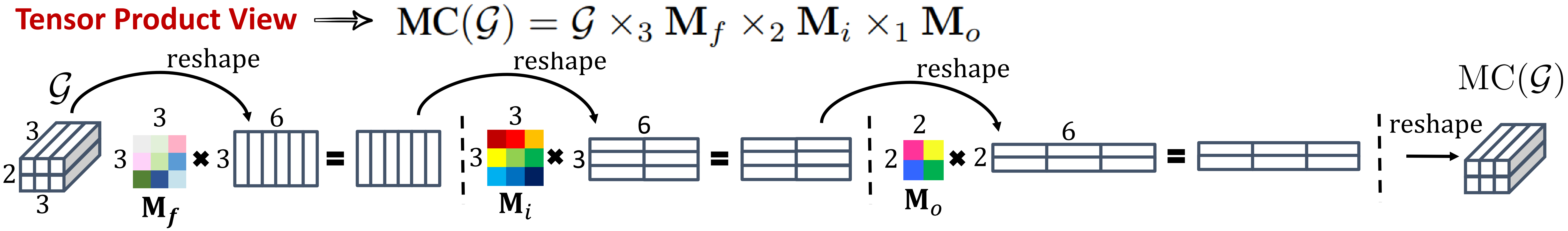
$$\operatorname{argmin}_{\theta, \mathbf{M}} \mathbb{E}_{\tau_i} [\mathcal{L}_{\text{val}}^{\tau_i}(\theta - \mathbf{M} \nabla \mathcal{L}_{\text{tr}}^{\tau_i}(\theta))]$$

Learning a curvature matrix for better generalization and fast model adaptation in MAML framework.

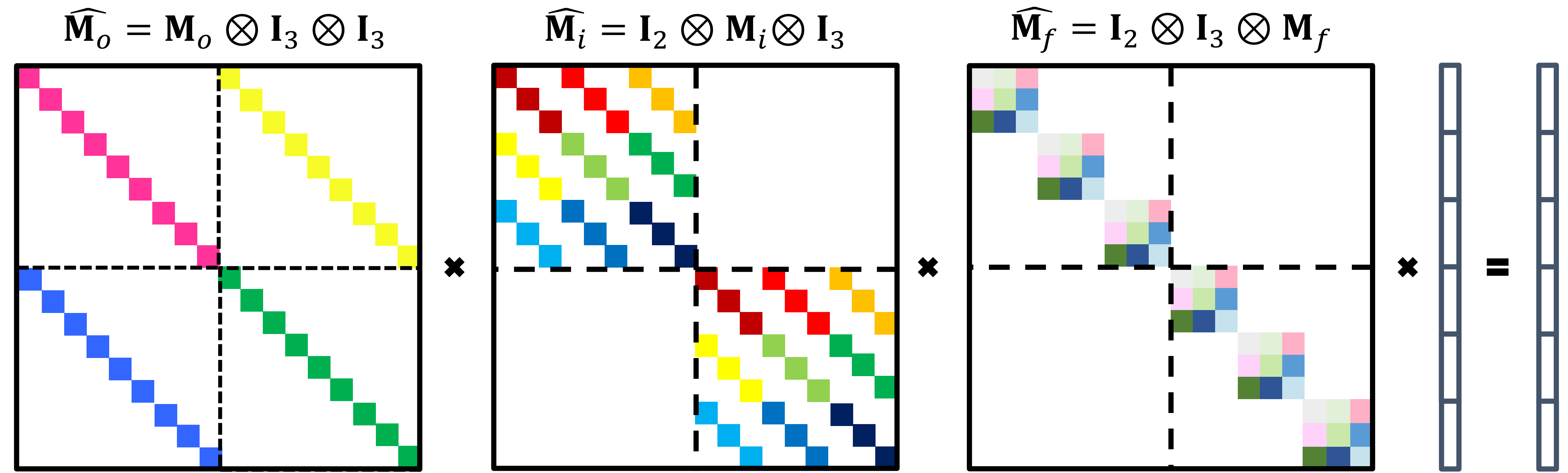
Modeling the second-order dependencies \Rightarrow Second-order optimization in the inner optimization

Reducing space and computational complexity \Rightarrow Decomposition of the curvature matrix

Meta-Curvature Computation



Matrix Product View $\Rightarrow \text{vec}(\text{MC}(\mathcal{G})) = \widehat{\mathbf{M}}_o \widehat{\mathbf{M}}_i \widehat{\mathbf{M}}_f \text{vec}(\mathcal{G})$



Experimental Results

[Few-shot classification - Omniglot]

	5w-1s	5w-5s	20w-1s	20w-5s
MAML [1]	98.7 \pm 0.4	99.9 \pm 0.1	95.8 \pm 0.3	98.9 \pm 0.2
Meta-SGD [2]	99.53 \pm 0.26	99.93 \pm 0.09	95.93 \pm 0.38	98.97 \pm 0.19
MAML++ [¶] [3]	99.47	99.93	97.65	99.33
MC	99.77 \pm 0.17	99.79 \pm 0.10	97.86 \pm 0.26	99.24 \pm 0.07
MC [¶]	99.97 \pm 0.06	99.89 \pm 0.06	99.12 \pm 0.16	99.65 \pm 0.05

5w-1s: 5way 1 shot [¶]: 3 model ensemble

[Few-shot classification – Imagenet]

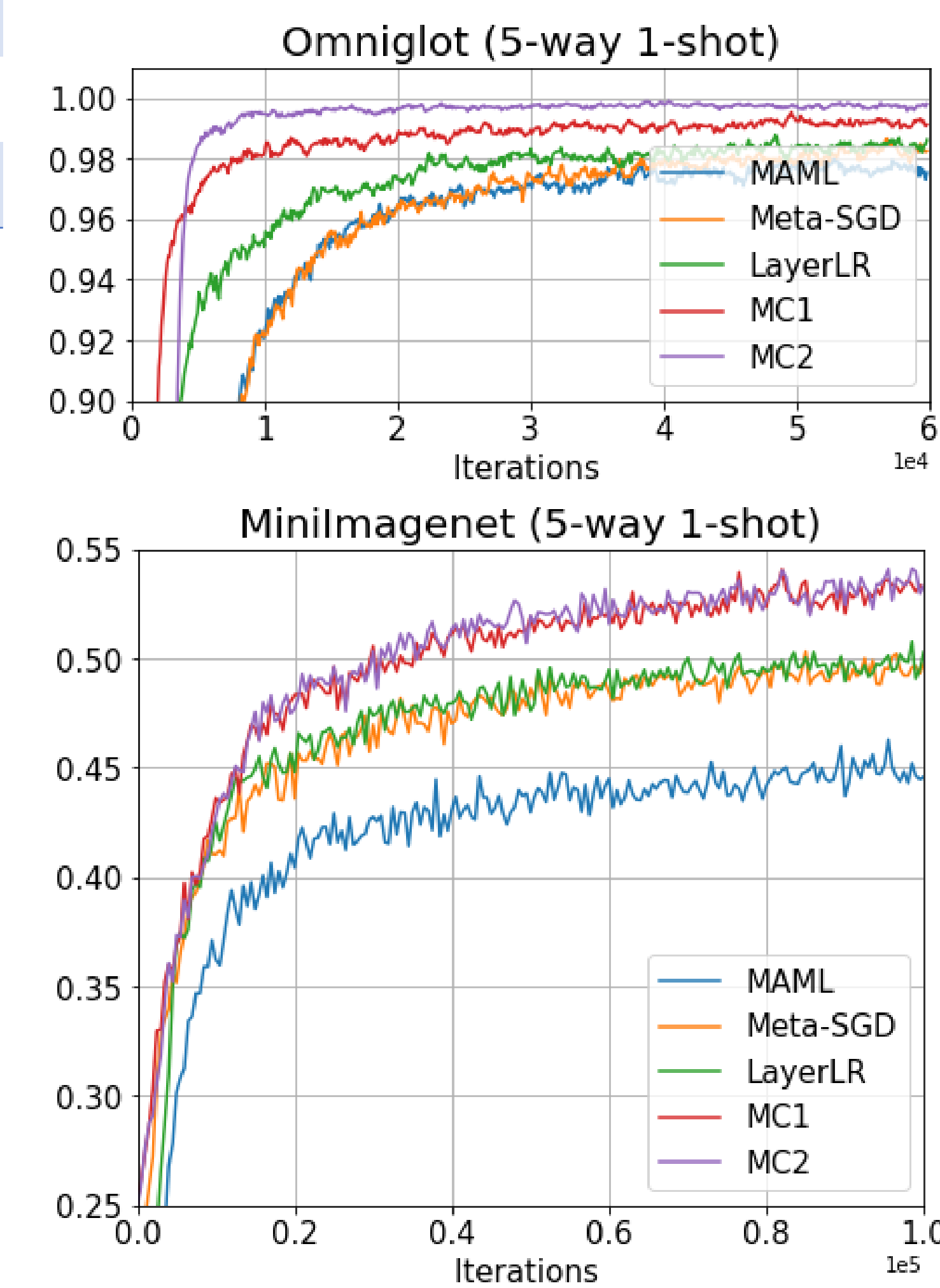
	mini-Imagenet		tiered-Imagenet	
	5w-1s	5w-1s	5w-5s	5w-5s
LEO – center [4]	61.76 \pm 0.08	77.59 \pm 0.12	66.33 \pm 0.05	81.44 \pm 0.09
LEO – multiview [4]	63.97 \pm 0.20	79.49 \pm 0.70		
MetaOptNet [◊] [5]	64.09 \pm 0.62	80.00 \pm 0.45	65.81 \pm 0.74	81.75 \pm 0.53
Meta-SGD [2]	56.58 \pm 0.21	68.84 \pm 0.19	59.75 \pm 0.25	69.04 \pm 0.22
MC – center	61.85 \pm 0.10	77.02 \pm 0.11	67.21 \pm 0.10	82.61 \pm 0.08
MC – multiview	64.40 \pm 0.10	80.21 \pm 0.10		

WRN-28-10 pretrained features + MLP [◊]: ResNet-12 and 15-shot meta-training
multiview: features averaged over four corners, central crops, and horizontal mirrored

[Few-shot regression]

	5-shot	10-shot
MAML [1]	0.69 \pm 0.07	0.44 \pm 0.04
Meta-SGD [2]	0.48 \pm 0.06	0.26 \pm 0.03
MC	0.41 \pm 0.05	0.20 \pm 0.02

Sinusoidal : amp [0.1, 5.0], phase [0, π]



Analysis

Meta-gradient w.r.t \mathbf{M} : $\theta^{\tau_i}(\mathbf{M}) = \theta - \alpha \mathbf{M} \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_i}(\theta)$

$$\nabla_{\mathbf{M}} \mathcal{L}_{\text{val}}^{\tau_i}(\theta^{\tau_i}(\mathbf{M})) = -\alpha \nabla_{\theta^{\tau_i}} \mathcal{L}_{\text{val}}^{\tau_i}(\theta^{\tau_i}) \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_i}(\theta)^{\top}$$

Given a fixed point θ and a meta training set $T = \{\tau_i\}$, gradient descents from an initial \mathbf{M}_0 :

$$\mathbf{M}_T = \mathbf{M}_0 - \beta \sum_{i=1}^{|\mathcal{T}|} \nabla_{\mathbf{M}_{i-1}} \mathcal{L}_{\text{val}}^{\tau_i}(\theta^{\tau_i}) = \mathbf{M}_0 + \alpha \beta \sum_{i=1}^{|\mathcal{T}|} \nabla_{\theta^{\tau_i}} \mathcal{L}_{\text{val}}^{\tau_i}(\theta^{\tau_i}) \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_i}(\theta)^{\top}$$

Applying \mathbf{M}_T to the gradients of a new task τ_{new} :

$$\begin{aligned} \mathbf{M}_T \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_{\text{new}}}(\theta) &= (\mathbf{M}_0 + \alpha \beta \sum_{i=1}^{|\mathcal{T}|} \nabla_{\theta^{\tau_i}} \mathcal{L}_{\text{val}}^{\tau_i}(\theta^{\tau_i}) \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_i}(\theta)^{\top}) \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_{\text{new}}}(\theta) \\ &= \mathbf{M}_0 \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_{\text{new}}}(\theta) + \beta \sum_{i=1}^{|\mathcal{T}|} (\nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_i}(\theta)^{\top} \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_{\text{new}}}(\theta)) \alpha \nabla_{\theta^{\tau_i}} \mathcal{L}_{\text{val}}^{\tau_i}(\theta^{\tau_i}) \\ &= \mathbf{M}_0 \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_{\text{new}}}(\theta) + \beta \sum_{i=1}^{|\mathcal{T}|} (\nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_i}(\theta)^{\top} \nabla_{\theta} \mathcal{L}_{\text{tr}}^{\tau_{\text{new}}}(\theta)) (\underbrace{\alpha \nabla_{\theta} \mathcal{L}_{\text{val}}^{\tau_i}(\theta)}_{\text{Gradient similarity}} + \underbrace{\mathcal{O}(\alpha^2)}_{\text{Taylor expansion}}). \end{aligned}$$

References

- [1] Meta-SGD: Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks, Finn et al, ICML 2017
[2] Meta-SGD: Learning to Learn Quickly for Few-Shot Learning, Li et al, arXiv 2017 [3] How to train your MAML, Antoniou et al, ICLR 2019
[4] Meta-Learning with Latent Embedding Optimization, Rusu et al, ICLR 2019 [5] Meta-Learning with Differentiable Convex Optimization, Lee et al, CVPR 2019