

# Optimization-based Meta-Learning

# Outline

- Brief overview of optim-based meta-learning
  - ▶ Meta-learning models aim to learn **an global/initial params** and a **update rule** of how to quickly adapt to individual tasks (i.e., adjusting global initial params to individual tasks)
  - ▶ Meta-learning objective is to minimize the expected generalization loss of the meta-learner on the task space
- LSTM-based meta-learning (Meta-LSTM) (Ravi and Larochelle, 2016)
  - ▶ Meta-LSTM uses a sequential LSTM update rule
- Meta-SGD (Li et al., 2017)
  - ▶ Meta-SGD learns the global learning rate and update direction  $\alpha$  addition to the initial params
- Meta-Curvature (Park and Oliva, 2020)
  - ▶ Meta-Curvature further improves MAML by adding second-order information into the gradient
    - ★ the additional info is decomposed into 3 component matrices and is multiplied to gradient during inner gradient update steps
- IMO, the current trend is modifying and adding components to the update rules to learn more about individual tasks, but the components are controlled by some global learnable variables

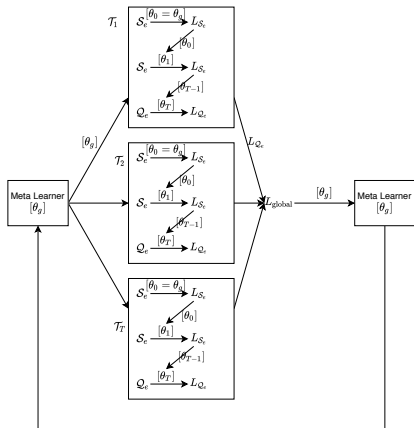
# Optim-based Meta-Learning: Intuition

- Meta-learning objective is to minimize the expected generalization loss of the meta-learner on the task space

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_e)] = \mathbb{E}_{\mathcal{T}_e} [\underbrace{L_{\mathcal{Q}_e}(\underbrace{\theta_g - \alpha \nabla_{\theta_g} L_{\mathcal{S}_e}(\theta_g))}_{\text{inner}}}_{\text{outer}}]$$

- ▶ **inner-loop:** given an individual task  $\mathcal{T}_e$ , the **learner** (params  $\theta_e$ ) is optimized based on the loss of task support set  $L_{\mathcal{S}_e}$
- ▶ **outer-loop:** the loss of task query set  $L_{\mathcal{Q}_e}$  is computed with the updated learner and is accumulated to the generalization loss  $L_{\text{global}} = \sum_e L_{\mathcal{Q}_e}$ . the **meta-learner** (params  $\theta_g$ ) is then optimized based the generalization loss.

# Model-Agnostic Meta-Learning (MAML): Model



## Algorithm : MAML's training procedure

**Require:** train meta-set  $\mathcal{T}^{\text{train}}$

**Return:** global params  $\theta_g$  that are adaptable to many tasks

- 1: initialize  $\theta_g$
- 2: **for** epoch in num epochs **do**
- 3:   **for** episode  $e$  in num episodes **do**
- 4:     sample a local task  $\mathcal{T}_e = (S_e, Q_e)$  from  $\mathcal{T}^{\text{train}}$
- 5:      $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set
- 6:     set  $\theta_0 = \theta_g$
- 7:     **for** step  $t$  in num steps **do**
- 8:       compute loss of the local task support set  

$$L_{S_e} = \sum_{(x,y) \in S_e} L_{\text{cls}}(f(x; \theta_{t-1}), y)$$
- 9:       compute gradient and update local params for the local task  

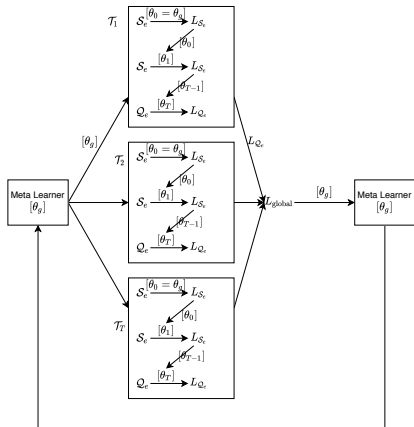
$$\theta_t = \theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{S_e}$$
- 10:     **end for**
- 11:      $\mapsto$  **outer loop:** learn global params via all the losses of all task query sets
- 12:     compute loss of the local task query set with local params  

$$L_{Q_e} = \sum_{(x,y) \in Q_e} L_{\text{cls}}(f(x; \theta_T), y)$$
- 13:     accumulate the loss of the local task query set to the global loss  

$$L_{\text{global}} += L_{Q_e}$$
- 14:   **end for**
- 15:   compute global gradient (w.r.t. sum of all local losses) and update global params  

$$\theta_g = \theta_g - \beta \nabla_{\theta_g} L_{\text{global}}$$
- 16: **end for**

# Model-Agnostic Meta-Learning (MAML): Model




---

## Algorithm : MAML's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  

$$L_{\mathcal{S}_e} = \sum_{(\mathbf{x}, y) \in \mathcal{S}_e} L_{\text{cls}}(f(\mathbf{x}; \theta_{t-1}), y)$$
  - 5:   compute gradient and update local params for the local task  

$$\theta_t = \theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e}$$
  - 6: **end for**
- 

## Algorithm : MAML's training procedure, outer-loop

---

**Require:** task query set  $\mathcal{Q}_e$ , individual task params  $\theta_T$

**Return:** loss of task query set  $L_{\mathcal{Q}_e}$

- 1:  $\mapsto$  **outer loop:** learn global params via all the losses of all task query sets
  - 2: compute loss of the local task query set with local params  

$$L_{\mathcal{Q}_e} = \sum_{(\mathbf{x}, y) \in \mathcal{Q}_e} L_{\text{cls}}(f(\mathbf{x}; \theta_T), y)$$
  - 3: accumulate the loss of the local task query set to the global loss  

$$L_{\text{global}} += L_{\mathcal{Q}_e}$$
-

# Meta-LSTM vs MAML: Model

---

**Algorithm :** MAML's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  
 $L_{\mathcal{S}_e} = \sum_{(x,y) \in \mathcal{S}_e} L_{\text{cls}}(f(x; \theta_{t-1}), y)$
  - 5:   compute gradient and update local params for the local task  
 $\theta_t = \theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e}$
  - 6: **end for**
- 

---

**Algorithm :** Meta-LSTM's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set sequentially with the LSTM mechanism
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  
 $L_{\mathcal{S}_e} = \sum_{(x,y) \in \mathcal{S}_e} L_{\text{cls}}(f(x; \theta_{t-1}), y)$
  - 5:   compute input gate  $i_t$  and forget gate  $f_t$   
 $i_t = \sigma(W^I \cdot [\nabla_{\theta_{t-1}} L_{\mathcal{S}_e}, L_{\mathcal{S}_e}, \theta_{t-1}, i_{t-1}] + b^I)$   
 $f_t = \sigma(W^F \cdot [\nabla_{\theta_{t-1}} L_{\mathcal{S}_e}, L_{\mathcal{S}_e}, \theta_{t-1}, f_{t-1}] + b^F)$
  - 6:   update local params for the local task according to LSTM rule  
 $\theta_t = f_t \odot \theta_{t-1} - i_t \odot \nabla_{\theta_{t-1}} L_{\mathcal{S}_e}$
  - 7: **end for**
- 

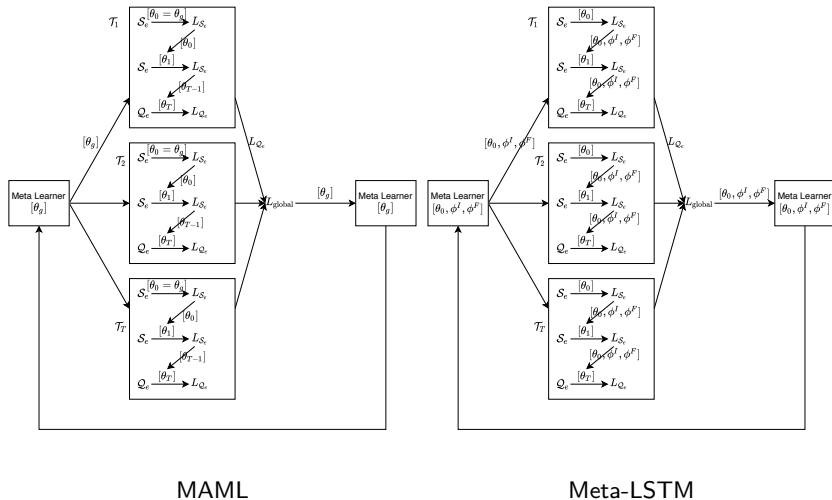
- MAML's objective

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_e)] = \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e})]$$

- Meta-SGD uses a sequential LSTM update rule, its objective becomes

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_e)] = \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(f_{t-1} \odot \theta_{t-1} - i_{t-1} \odot \nabla_{\theta_{t-1}} L_{\mathcal{S}_e})]$$

# Meta-LSTM vs MAML: Model



# Meta-SGD vs MAML: Model

---

**Algorithm :** MAML's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  
     $L_{\mathcal{S}_e} = \sum_{(\mathbf{x}, y) \in \mathcal{S}_e} L_{\text{cls}}(f(\mathbf{x}; \theta_{t-1}), y)$
  - 5:   compute gradient and update local params for the local task  
     $\theta_t = \theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e}$
  - 6: **end for**
- 

---

**Algorithm :** Meta-SGD's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$ , global lr  $\alpha_g$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set using global params  $\theta_g$  and global lr  $\alpha_g$
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  
     $L_{\mathcal{S}_e} = \sum_{(\mathbf{x}, y) \in \mathcal{S}_e} L_{\text{cls}}(f(\mathbf{x}; \theta_{t-1}), y)$
  - 5:   compute gradient and update local params for the local task  
     $\theta_t = \theta_{t-1} - \alpha_g \nabla_{\theta_{t-1}} L_{\mathcal{S}_e}$
  - 6: **end for**
- 

- MAML's objective

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e(\theta_e)}] = \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e(\theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e})}]$$

- In addition to learning global initial params as in MAML, Meta-SGD also learns the global learning rate (and update direction) alpha

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e(\theta_e)}] = \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e(\theta_{t-1} - \alpha_g \nabla_{\theta_{t-1}} L_{\mathcal{S}_e})}]$$

# Meta-Curvature vs MAML: Model

---

**Algorithm :** MAML's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  
     $L_{\mathcal{S}_e} = \sum_{(x,y) \in \mathcal{S}_e} L_{\text{cls}}(f(\mathbf{x}; \theta_{t-1}), y)$
  - 5:   compute gradient and update local params for the local task  
     $\theta_t = \theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e}$
  - 6: **end for**
- 

---

**Algorithm :** Meta-Curvature's training procedure, inner-loop

---

**Require:** task support set  $\mathcal{S}_e$ , global params  $\theta_g$ , curvature params  $\mathbf{M}_f, \mathbf{M}_i, \mathbf{M}_o$

**Return:** individual task params  $\theta_T$  after  $T$  update steps

- 1:  $\mapsto$  **inner loop:** learn local params for the local task via the loss of task support set
  - 2: set  $\theta_0 = \theta_g$
  - 3: **for** step  $t$  in num steps **do**
  - 4:   compute loss of the local task support set  
     $L_{\mathcal{S}_e} = \sum_{(x,y) \in \mathcal{S}_e} L_{\text{cls}}(f(\mathbf{x}; \theta_{t-1}), y)$
  - 5:   compute gradient and update local params for the local task  
     $\theta_t = \theta_{t-1} - \alpha \text{MC}(\nabla_{\theta_{t-1}} L_{\mathcal{S}_e}; \mathbf{M}_f, \mathbf{M}_i, \mathbf{M}_o)$
  - 6: **end for**
- 

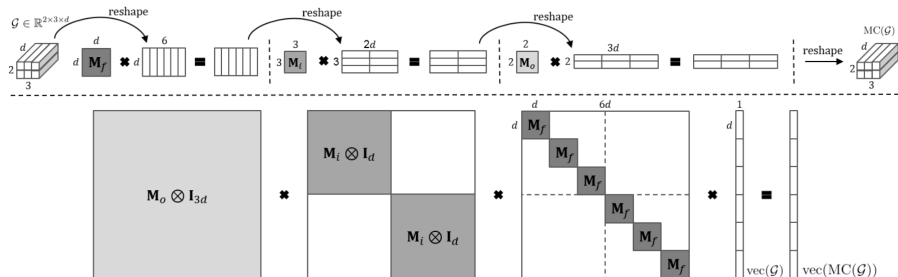
## ● MAML's objective

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_e)] = \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_{t-1} - \alpha \nabla_{\theta_{t-1}} L_{\mathcal{S}_e})]$$

## ● Meta-Curvature further improves MAML by adding second-order information into the gradient

$$\min_{\theta_g} \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_e)] = \mathbb{E}_{\mathcal{T}_e} [L_{\mathcal{Q}_e}(\theta_{t-1} - \alpha \text{MC}(\nabla_{\theta_{t-1}} L_{\mathcal{S}_e}; \mathbf{M}_f, \mathbf{M}_i, \mathbf{M}_o))]$$

# Meta-Curvature vs MAML: Model



- The second-order information is characterized/decomposed by into 3 component matrices, the info is multiplied to the gradient during updates

$$MC(\nabla_{\theta_g} L_{S_e}(\theta_g); M_f, M_i, M_o) = \nabla_{\theta_g} L_{S_e}(\theta_g) \times_3 M_f \times_2 M_i \times_1 M_o$$

- all component matrices are updated globally similar to global initial params

# FSL Benchmark

Model	<i>miniImageNet</i> test accuracy	
	1-shot	5-shot
Matching networks (Vinyals et al., 2016)	43.56 $\pm$ 0.84%	55.31 $\pm$ 0.73%
Meta-learner LSTM (Ravi & Larochelle, 2017)	43.44 $\pm$ 0.77%	60.60 $\pm$ 0.71%
MAML (Finn et al., 2017)	48.70 $\pm$ 1.84%	63.11 $\pm$ 0.92%
LLAMA (Grant et al., 2018)	49.40 $\pm$ 1.83%	-
REPTILE (Nichol & Schulman, 2018)	49.97 $\pm$ 0.32%	65.99 $\pm$ 0.58%
PLATIPUS (Finn et al., 2018)	50.13 $\pm$ 1.86%	-
Meta-SGD (our features)	54.24 $\pm$ 0.03%	70.86 $\pm$ 0.04%
SNAIL (Mishra et al., 2018)	55.71 $\pm$ 0.99%	68.88 $\pm$ 0.92%
Gidaris & Komodakis, 2018	56.20 $\pm$ 0.86%	73.00 $\pm$ 0.64%
(Bauer et al., 2017)	56.30 $\pm$ 0.40%	73.90 $\pm$ 0.30%
(Munkhdalai et al., 2017)	57.10 $\pm$ 0.70%	70.04 $\pm$ 0.63%
DEML+Meta-SGD (Zhou et al., 2018) <sup>4</sup>	58.49 $\pm$ 0.91%	71.28 $\pm$ 0.69%
TADAM (Oreshkin et al., 2018)	58.50 $\pm$ 0.30%	76.70 $\pm$ 0.30%
(Qiao et al., 2017)	59.60 $\pm$ 0.41%	73.74 $\pm$ 0.19%
<b>LEO (ours)</b>	<b>61.76 <math>\pm</math> 0.08%</b>	<b>77.59 <math>\pm</math> 0.12%</b>
Model	<i>tieredImageNet</i> test accuracy	
	1-shot	5-shot
MAML (deeper net, evaluated in (Liu et al., 2018))	51.67 $\pm$ 1.81%	70.30 $\pm$ 0.08%
Prototypical Nets (Ren et al., 2018)	53.31 $\pm$ 0.89%	72.69 $\pm$ 0.74%
Relation Net (evaluated in (Liu et al., 2018))	54.48 $\pm$ 0.93%	71.32 $\pm$ 0.78%
Transductive Prop. Nets (Liu et al., 2018)	57.41 $\pm$ 0.94%	71.55 $\pm$ 0.74%
Meta-SGD (our features)	62.95 $\pm$ 0.03%	79.34 $\pm$ 0.06%
<b>LEO (ours)</b>	<b>66.33 <math>\pm</math> 0.05%</b>	<b>81.44 <math>\pm</math> 0.09%</b>

Table 1: Test accuracies on *miniImageNet* and *tieredImageNet*. For each dataset, the first set of results use convolutional networks, while the second use much deeper residual networks, predominantly in conjunction with pre-training.

*Thank you !*

- Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-SGD: Learning to Learn Quickly for Few-Shot Learning. *arXiv:1707.09835 [cs]*, 2017.
- Eunbyung Park and Junier B. Oliva. Meta-Curvature. In *Proceedings of the NeurIPS 2019*, 2020.
- Sachin Ravi and Hugo Larochelle. Optimization as a Model for Few-Shot Learning. In *Proceedings of the ICLR 2016*, 2016.