

# Meta-learning

Yunqiu Xu

# Readings

1. Model-agnostic Meta-learning
2. MAML for One-shot Imitation Learning
3. Meta Learning Shared Hierarchies
4. MAML for Non-stationary Task

MAML : Model-agnostic Meta-learning

# MAML : Model-agnostic Meta-learning

- + Meta-learning: train a model with some learning tasks, then it can solve new tasks with only a few samples
- + MAML can be treated as an initialization method to get pretrained “base model” that is easy to fine-tune to new tasks
- + Compare with another recent work (Ravi & Larochelle, 2017) : do not need additional parameters, and “model agnostic”

# MAML: Problem Setup

General notion of task  $T = \{L(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H, \mathbf{a}_H), q(\mathbf{x}_1), q(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{a}_t), H\}$  :

- $L \rightarrow \mathbf{R}$  : loss function
- $q(\mathbf{x}_1)$  : distribution over initial observations
- $q(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{a}_t)$  : transition distribution
  - 对于监督学习, 不存在这个分布 :  $H = 1$
  - 对于强化学习,  $q(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{a}_t)$  代表某时间点观察值的分布, e.g. 初始观察值后观察值取自  $q(\mathbf{x}_2|\mathbf{x}_1, \mathbf{a}_1)$
- $H$  :
  - Episode length, model may generate samples of length  $H$  by choosing an output  $\mathbf{a}_t$  at each time  $t$
  - For supervised learning,  $H = 1$  and loss function  $L(\mathbf{x}_1, \mathbf{a}_1)$  could be MSE or cross entropy

# MAML : Algorithm

---

**Algorithm 1** Model-Agnostic Meta-Learning

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
  - 2: **while** not done **do**
  - 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
  - 4:   **for all**  $\mathcal{T}_i$  **do**
  - 5:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples
  - 6:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
  - 7:   **end for**
  - 8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
  - 9: **end while**
-

# MAML: K-shot Training and Testing

- Train model  $f$ :
  - Sample a new task  $T_i$  from  $p(T)$  (training taskset)
  - Learn  $T_i$  :
    - Train model with K samples drawn from  $q_i$
    - Get feedback  $L_{T_i}$  from  $T_i$
  - Test on new samples from  $T_i$  and get test error
  - Improve model  $f$  : treat the test error on sampled tasks  $T_i$  as the training error of meta-learning process
- Test meta-learning:
  - Sample new task from  $p(T)$  (testing taskset), try to adapt  $f$  to this new task
  - Learn the model with K samples
  - Treat the performance as "meta-performance"

# MAML : An Example of Inner-loop

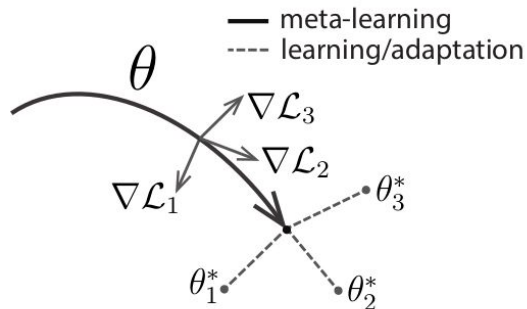


Figure 1. Diagram of our model-agnostic meta-learning algorithm (MAML), which optimizes for a representation  $\theta$  that can quickly adapt to new tasks.

- + 假定我们在当前内循环需要训练T2, 更新后的参数  $\theta_2'$  和原来相比会稍稍“上移”
- + 同理, 我们训练T1和T3后得到的  $\theta_1'$  和  $\theta_3'$  和原来相比方向也会有一定调整
- + 因此在meta-loop中我们用所有的  $\theta_i'$  来更新  $\theta$ , 得到一个比较 general 的方向



# MAML: SL and RL

---

## Algorithm 2 MAML for Few-Shot Supervised Learning

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
  - 2: **while** not done **do**
  - 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
  - 4:   **for all**  $\mathcal{T}_i$  **do**
  - 5:     Sample  $K$  datapoints  $\mathcal{D} = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$  from  $\mathcal{T}_i$
  - 6:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  using  $\mathcal{D}$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation (2) or (3)
  - 7:     Compute adapted parameters with gradient descent:  
       $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
  - 8:     Sample datapoints  $\mathcal{D}'_i = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$  from  $\mathcal{T}_i$  for the meta-update
  - 9:   **end for**
  - 10:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  using each  $\mathcal{D}'_i$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 2 or 3
  - 11: **end while**
- 

---

## Algorithm 3 MAML for Reinforcement Learning

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
  - 2: **while** not done **do**
  - 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
  - 4:   **for all**  $\mathcal{T}_i$  **do**
  - 5:     Sample  $K$  trajectories  $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_{\theta}$  in  $\mathcal{T}_i$
  - 6:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  using  $\mathcal{D}$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
  - 7:     Compute adapted parameters with gradient descent:  
       $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
  - 8:     Sample trajectories  $\mathcal{D}'_i = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_{\theta'_i}$  in  $\mathcal{T}_i$
  - 9:   **end for**
  - 10:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  using each  $\mathcal{D}'_i$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
  - 11: **end while**
-

# MAML: Experiment

Regression task: 模拟sine曲线

Classification task: few-shot image recognition

RL task: 2D Navigation, MuJoCo Simulation

# MAML: Regression Result

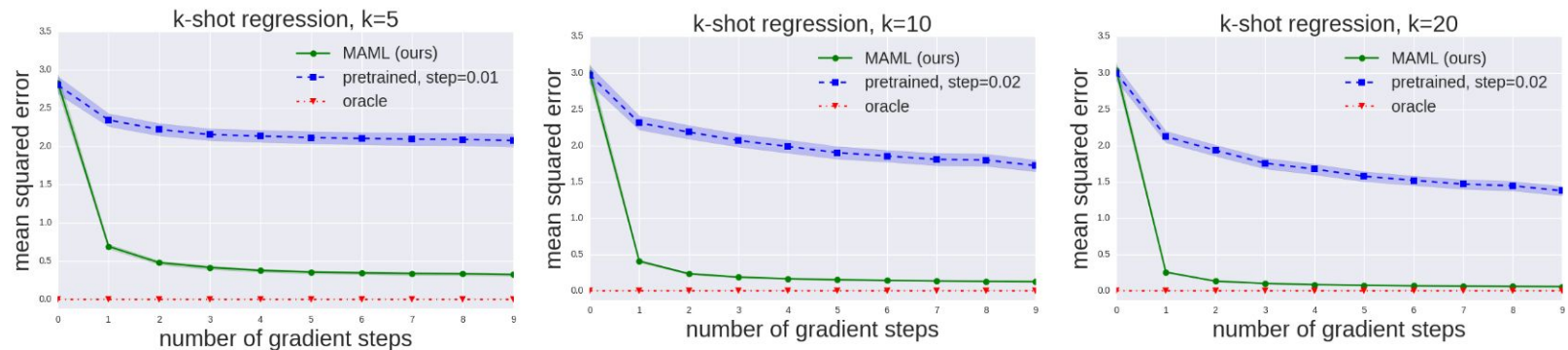


Figure 6. Quantitative sinusoid regression results showing test-time learning curves with varying numbers of  $K$  test-time samples. Each gradient step is computed using the same  $K$  examples. Note that MAML continues to improve with additional gradient steps without overfitting to the extremely small dataset during meta-testing, and achieves a loss that is substantially lower than the baseline fine-tuning approach.

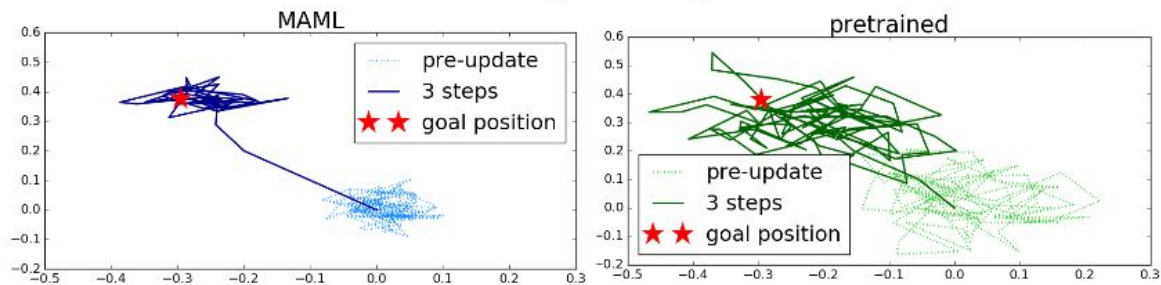
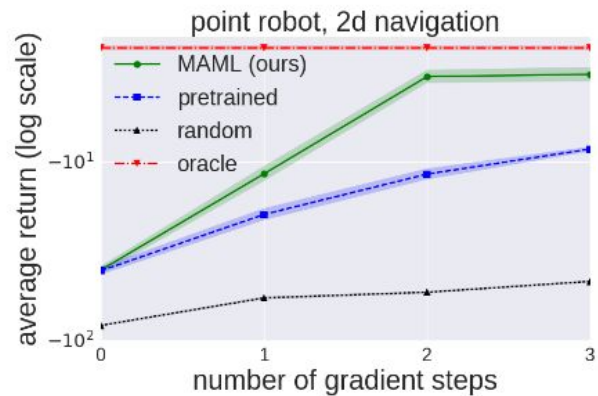
- When all the points are in one half, the model can still infer the shape of the other half → model the periodic nature
- Regression可以用很少的样本/循环进行finetune, 不会overfitting
- 因为效果本身就很不错了, 增加迭代次数并没有明显提升

# MAML: Classification Result

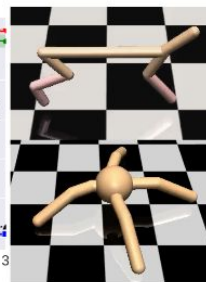
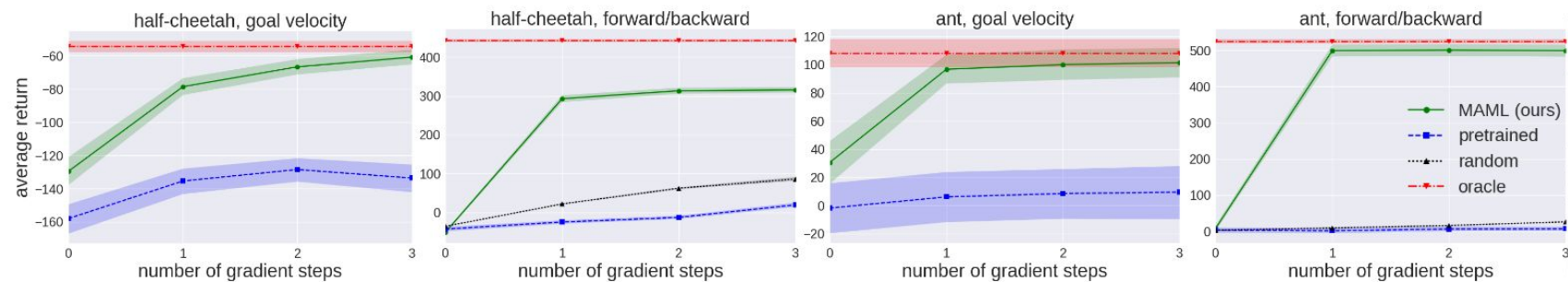
	5-way Accuracy		20-way Accuracy	
	1-shot	5-shot	1-shot	5-shot
Omniglot (Lake et al., 2011)				
MANN, no conv (Santoro et al., 2016)	82.8%	94.9%	–	–
<b>MAML, no conv (ours)</b>	<b><math>89.7 \pm 1.1\%</math></b>	<b><math>97.5 \pm 0.6\%</math></b>	–	–
Siamese nets (Koch, 2015)	97.3%	98.4%	88.2%	97.0%
matching nets (Vinyals et al., 2016)	98.1%	98.9%	93.8%	98.5%
neural statistician (Edwards & Storkey, 2017)	98.1%	99.5%	93.2%	98.1%
memory mod. (Kaiser et al., 2017)	98.4%	99.6%	95.0%	98.6%
<b>MAML (ours)</b>	<b><math>98.7 \pm 0.4\%</math></b>	<b><math>99.9 \pm 0.1\%</math></b>	<b><math>95.8 \pm 0.3\%</math></b>	<b><math>98.9 \pm 0.2\%</math></b>

	5-way Accuracy	
	1-shot	5-shot
MiniImagenet (Ravi & Larochelle, 2017)		
fine-tuning baseline	$28.86 \pm 0.54\%$	$49.79 \pm 0.79\%$
nearest neighbor baseline	$41.08 \pm 0.70\%$	$51.04 \pm 0.65\%$
matching nets (Vinyals et al., 2016)	$43.56 \pm 0.84\%$	$55.31 \pm 0.73\%$
meta-learner LSTM (Ravi & Larochelle, 2017)	$43.44 \pm 0.77\%$	$60.60 \pm 0.71\%$
<b>MAML, first order approx. (ours)</b>	<b><math>48.07 \pm 1.75\%</math></b>	<b><math>63.15 \pm 0.91\%</math></b>
<b>MAML (ours)</b>	<b><math>48.70 \pm 1.84\%</math></b>	<b><math>63.11 \pm 0.92\%</math></b>

# MAML: RL Result



# MAML: RL Result



MIL : MAML for One-shot Imitation Learning

# MIL : MAML for One-shot Imitation Learning

- + Generalize meta-learning technique to apply to imitation learning
- + Raw visual inputs
- + Reuse past experience to train the "base model", then adapt it to new task with only a single demonstration



# MIL: Problem Setup

- + Goal: learn a policy that can quickly adapt to new tasks from a single demonstration of that task
- + Task:

Each imitation task  $T_i = \{\tau = \{o_1, a_1, \dots, o_T, a_T\} \sim \pi_i^*, L(a_{1:T}, \hat{a}_{1:T}), T\}$

- $\tau$  : a demonstration generated by policy  $\pi_i^*$
- $L(a_1, \dots, a_T, \hat{a}_1, \dots, \hat{a}_T) \rightarrow R$  : loss function to give feedback

# MIL: Algorithm

- $o_t$  is the observation at time  $t$ , i.e. an image, while  $a_t$  is the action
- For demonstration trajectory  $\tau$ , we use MSE to compute loss:

$$L_{T_i}(f_\phi) = \sum_{\tau_j \sim T_i} \sum_t \|f_\phi(o_t^{(j)}) - a_t^{(j)}\|_2^2 \quad (2)$$

---

## Algorithm 1 Meta-Imitation Learning with MAML

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
  - 2: **while** not done **do**
  - 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
  - 4:   **for all**  $\mathcal{T}_i$  **do**
  - 5:     Sample demonstration  $\tau = \{\mathbf{o}_1, \mathbf{a}_1, \dots, \mathbf{o}_T, \mathbf{a}_T\}$  from  $\mathcal{T}_i$
  - 6:     Evaluate  $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$  using  $\tau$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation (2)
  - 7:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
  - 8:     Sample demonstration  $\tau'_i = \{\mathbf{o}'_1, \mathbf{a}'_1, \dots, \mathbf{o}'_T, \mathbf{a}'_T\}$  from  $\mathcal{T}_i$  for the meta-update
  - 9:   **end for**
  - 10:   Update  $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\tau_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$  using each  $\tau'_i$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation (2)
  - 11: **end while**
  - 12: **return** parameters  $\theta$  that can be quickly adapted to new tasks through imitation.
-

# MIL: Algorithm

- Meta-training:

- Assume each training task has at least 2 demonstrations, thus we can sample a set of tasks with **two demonstrations per task**
- For each task  $T_i$ , train  $\theta'_i$  with its one demonstration  $\tau_i \rightarrow$  inner loop of meta-learning
- Use another demonstration  $\tau'_i$  to "test"  $\theta'_i$ , i.e. check the mse of predicted actions and demonstration actions
- Then we can update  $\theta$  according to the gradient of meta-objective
- As we get a series of  $\theta'_i$ s and their testing error, we can update  $\theta$
- Finally we can get trained parameters  $\theta$  for meta-learner

- Meta-testing:

- Sample a new task  $T$  and its one demonstration
- This task can involve new goals or manipulating new, previously unseen objects
- Then we can adapt  $\theta$  to this task

# MIL: Network Architecture

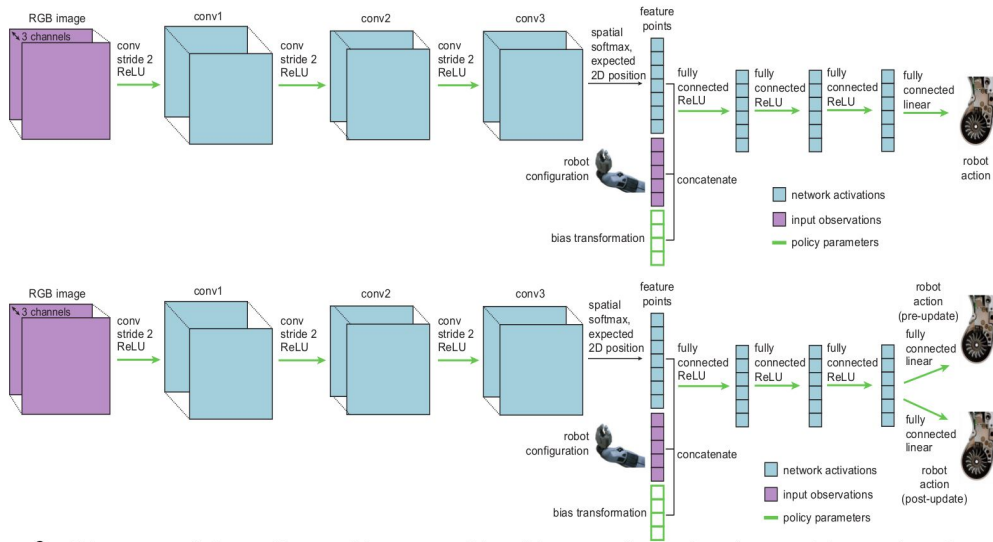


Figure 2: Diagrams of the policy architecture with a bias transformation (top and bottom) and two heads (bottom). The green arrows and boxes indicate weights that are part of the meta-learned policy parameters  $\theta$ .

Two-head structure : more flexibility during adapting

- Modification : parameters of final layers are not shared, forming two heads
  - Change loss function as:

$$L_{T_i}(f_\phi) = \sum_{\tau_j \sim T_i} \sum_t \|W y_t^{(j)} + b - a_t^{(j)}\|_2^2 \quad (3)$$

- $y_t^{(j)}$  : post-synaptic activations of the last hidden layer
- $W, b$  : weights and bias for last layer
- Then the meta-objective is about  $\theta, W, b$

$$\min_{\theta, W, b} \sum_{T_i \sim \mathcal{P}(T)} L_{T_i}(f_{\theta_i}) = \sum_{T_i \sim \mathcal{P}(T)} L_{T_i}(f_{\theta - \alpha \nabla_{\theta} L_{T_i}(f_{\theta})}) \quad (4)$$

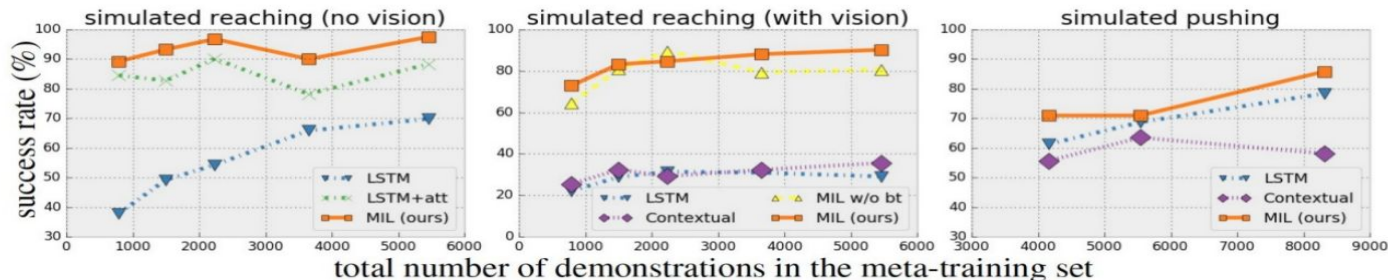
# MIL: Network Architecture

- **Layer normalization** after each layer
  - Data within a demonstration trajectory is highly correlated across time
  - Thus BN was not effective
- Bias transformation  $\rightarrow$  improve the performance of meta-learning
  - Concatenate a vector of parameters to a hidden layer of post-synaptic activations
  - Thus vector is treated as same as other parameters during meta-learning and final testing

$$y = Wx + b \rightarrow y = W_1x + W_2z + b$$

# MIL: Experiments

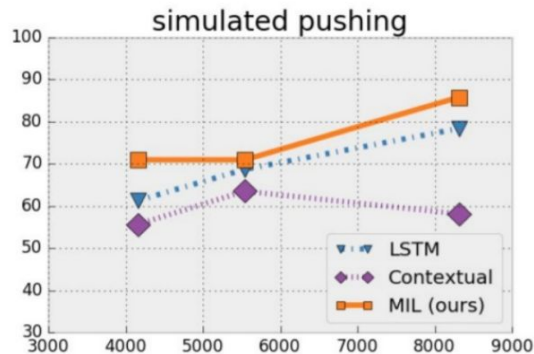
## Simulated Reaching



## Simulated Pushing

method		video+state +action	video +state	video
LSTM	1-shot	78.38%	37.61%	34.23%
contextual		n/a	58.11%	56.98%
MIL (ours)	1-shot	<b>85.81%</b>	<b>72.52%</b>	<b>66.44%</b>
LSTM	5-shot	83.11%	39.64%	31.98%
contextual		n/a	64.64%	59.01%
MIL (ours)	5-shot	<b>88.75%</b>	<b>78.15%</b>	<b>70.50%</b>

Table 1: One-shot and 5-shot simulating pushing success rate with varying demonstration information provided at test-time. MIL can more successfully learn from a demonstration without actions and without robot state and actions than LSTM and contextual policies.



MLSH : Meta-learning Shared Hierarchies

# MLSH : Meta-learning Shared Hierarchies

- Hierarchical model similar to "options framework"
- Contain a set of shared sub-policies (primitives) → these primitives are shared within a distribution of tasks
- How to switch these sub-tasks : by using a task-specific master policy
- For new tasks, we can just **learn master policy only** about how to switch the sub-policies correctly



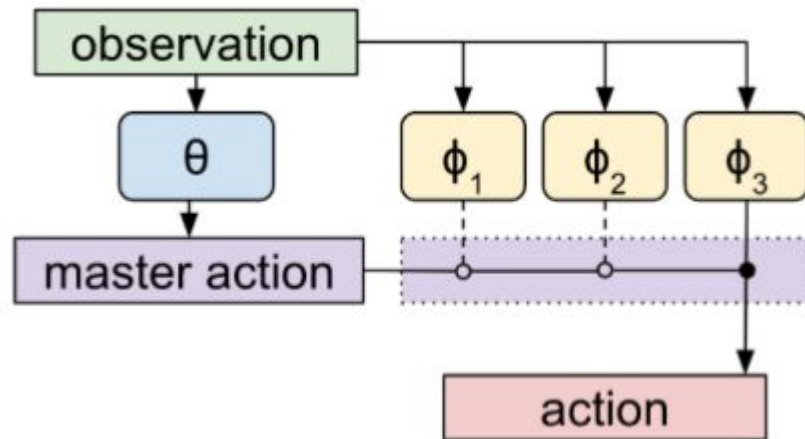
# MLSH : Problem Setup

- Define a policy  $\pi_{\phi, \theta}(a|s)$ 
  - $\phi$ :
    - A set of parameters **shared between all tasks**
    - $\phi = \{\phi_1, \phi_2, \dots, \phi_K\}$
    - Each  $\phi_k \rightarrow$  the parameters of a sub-policy  $\pi_{\phi_k}(a|s)$
  - $\theta$ :
    - The parameters of master policy
    - Task-specific  $\rightarrow$  zero or random initialized at the beginning
    - Choose a sub-task to activate for given timestep

# MLSH : Problem Setup

- For a task  $M$  sampled from  $P_M$ 
  - Randomly initialized  $\theta$  and shared  $\phi$
  - Goal: learn  $\theta$ , note that this is just the **objective for current task**
- The objective of meta-learning:
  - By learning training tasks, try to **find shared parameter  $\psi$**  which can be generalize to a new MDP
  - Then for a new task, only learn  $\theta$

$$\text{maximize}_{\phi} E_{M \sim P_M, t=0, \dots, T-1} [R]$$



# MLSH : Algorithm

---

**Algorithm 1** Meta Learning Shared Hierarchies

---

```
Initialize  $\phi$ 
repeat
  Initialize  $\theta$ 
  Sample task  $M \sim P_M$ 
  for  $w = 0, 1, \dots, W$  (warmup period) do
    Collect  $D$  timesteps of experience using  $\pi_{\phi, \theta}$ 
    Update  $\theta$  to maximize expected return from  $1/N$  timescale viewpoint
  end for
  for  $u = 0, 1, \dots, U$  (joint update period) do
    Collect  $D$  timesteps of experience using  $\pi_{\phi, \theta}$ 
    Update  $\theta$  to maximize expected return from  $1/N$  timescale viewpoint
    Update  $\phi$  to maximize expected return from full timescale viewpoint
  end for
until convergence
```

---

- The goal of meta-training (policy update) is to learn  $\phi$  which can be shared for all tasks
- The goal of learning a single task is to learn  $\theta \rightarrow$  choose  $\phi$  correctly
- At the very beginning we random initialize  $\phi$ , and for each new task, we random initialize  $\theta \rightarrow$  对每个新任务都要重设 $\theta$

# MLSH : Warm-up Period

```
for  $w = 0, 1, \dots, W$  (warmup period) do  
  Collect  $D$  timesteps of experience using  $\pi_{\phi, \theta}$   
  Update  $\theta$  to maximize expected return from  $1/N$  timescale viewpoint  
end for
```

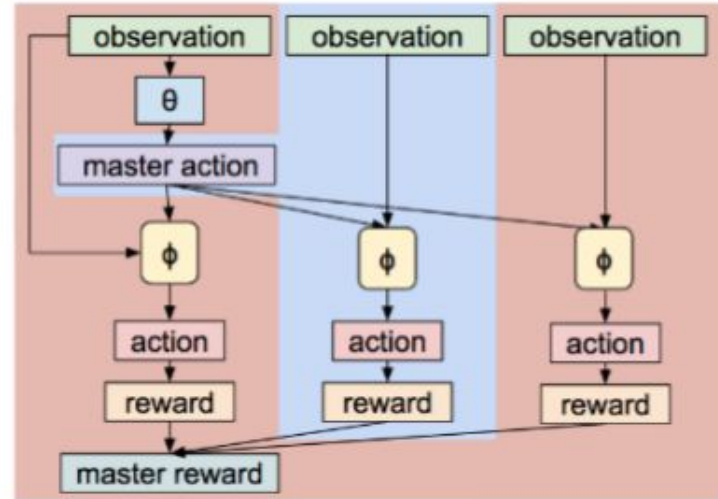
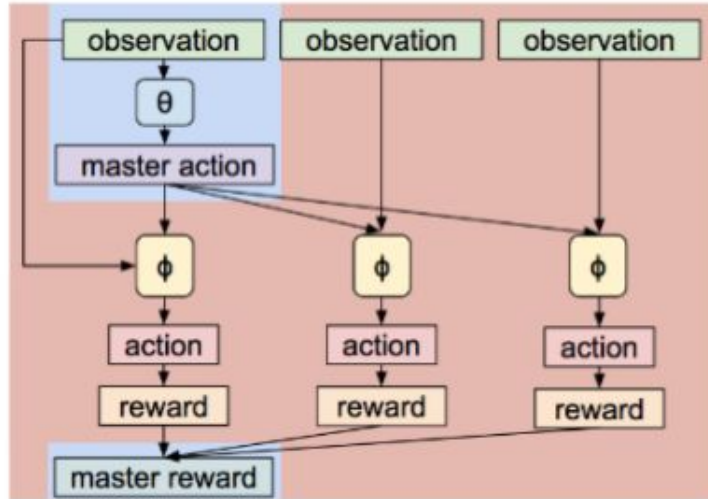
- **Goal : try to optimize  $\theta$  to nearly optimal**
- In this period, we hold  ***$\phi$***  fixed
- For each iteration sample  $D$  timesteps of experience
- For each  $1/N$  timescale, consider a sub-policy as an "action"
- 注意这里  $1/N$  timescale 的意思就是每间隔  $\frac{1}{N} * total\_time$  选一个动作

# MLSH : Joint-update Period

```
for  $u = 0, 1, \dots, U$  (joint update period) do  
  Collect  $D$  timesteps of experience using  $\pi_{\phi, \theta}$   
  Update  $\theta$  to maximize expected return from  $1/N$  timescale viewpoint  
  Update  $\phi$  to maximize expected return from full timescale viewpoint  
end for
```

- **Both  $\theta$  and  $\phi$  are updated**
- For each iteration, collect experience and optimize  $\theta \rightarrow$  **same as warm-up**
- Update  $\phi$  : reuse these  $D$  samples, but viewed via sub-policy
- Treat the master policy as an extension of the environment  $\rightarrow$  a discrete portion of observation
- For each  $N$ -timestep slice of experience, we only update the parameters of the sub-policy that had been activated by master policy

# MLSH : An Example

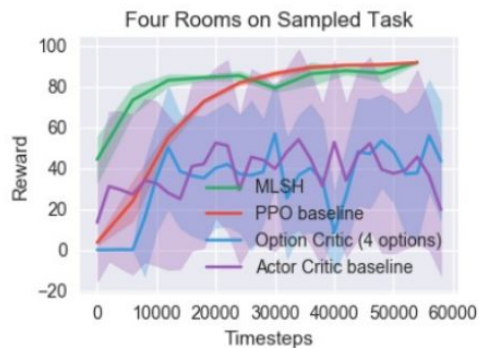
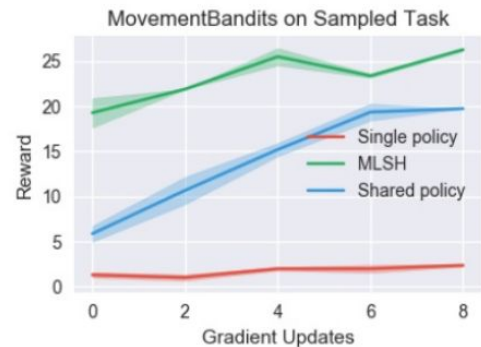
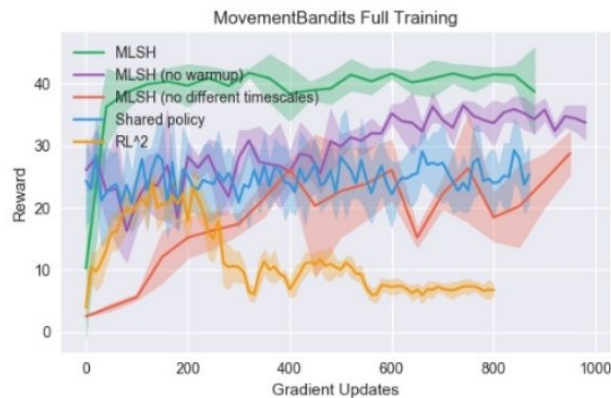


# MLSH : Why Faster

- Traditional meta-learning : optimize reward (master policy  $\theta$ ) over an entire inner loop
- MLSH :
  - **注意这里引入假设 1 : 经过warm-up可以学习到 optimal  $\theta$**
  - $\theta$  is updated per N-timesteps, only a much smaller task  $\phi$  to update over entire inner loop
  - $\theta$  is learned to nearly optimal in warm-up, so it will not take much time to reach optimal in joint-update
- 为什么限制joint-update的次数:
  - **这里引入假设 2: 在 joint-update 时因为  $\theta$  已经是优化的了, 即使更新也和原来区别不大**
  - 因此 joint-update 的主要作用是更新  $\phi$ , 而子任务相对容易学习, 在参数上微调就好
  - 因此我们不需要在joint-update上花费太多时间, 只需要固定训练循环数就好
  - 这里我的理解是因为一开始  $\phi$  并不够robust, 因此在warm-up过程中得到的  $\theta$  只能是近似优化的, 然后在 joint-update 过程中进一步优化, 并优化  $\phi$

# MLSH : Experiments

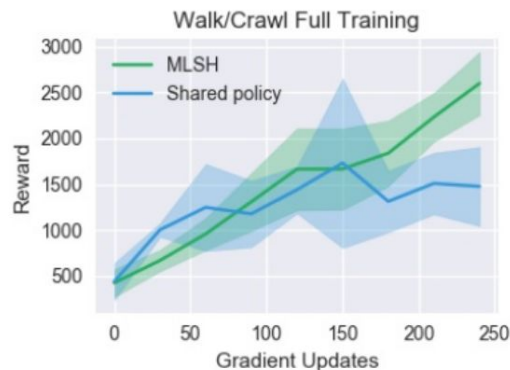
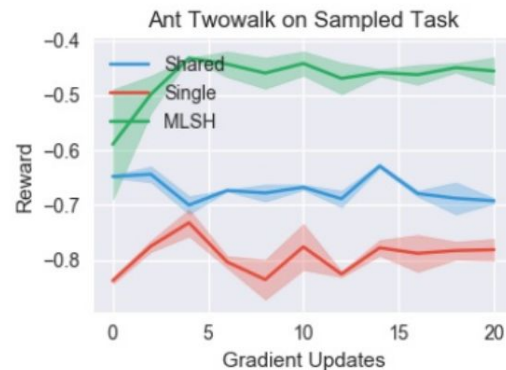
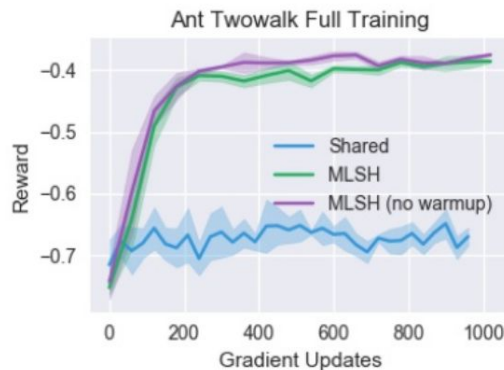
## 2D moving bandits





# MLSH : Experiments

## Simulated Walk



# MAML for Non-stationary Task

- + Challenge: non-stationary  $\rightarrow$  changing environment, multiple actors
- + Insight: non-stationary task can be seen as a sequence of stationary tasks  $\rightarrow$  multi-task problem
- + Our work: gradient-based meta-learning algorithm for adaptation in dynamically changing and adversarial scenarios
- + Goal: handle continuous adaptation with few-shot

# Probabilistic View of MAML

- Task  $T$ , trajectories  $\tau$  and policies  $\pi_\theta$  are random variables,  $\phi$  is generated from conditional distribution  $P_T(\phi|\theta, \tau_{1:k})$
- Inner loop update: equivalent to assuming the delta distribution

$$P_T(\phi|\theta, \tau_{1:k}) := \delta(\theta - \alpha \nabla_\theta \frac{1}{K} \sum_{k=1}^K L(\tau_k))$$

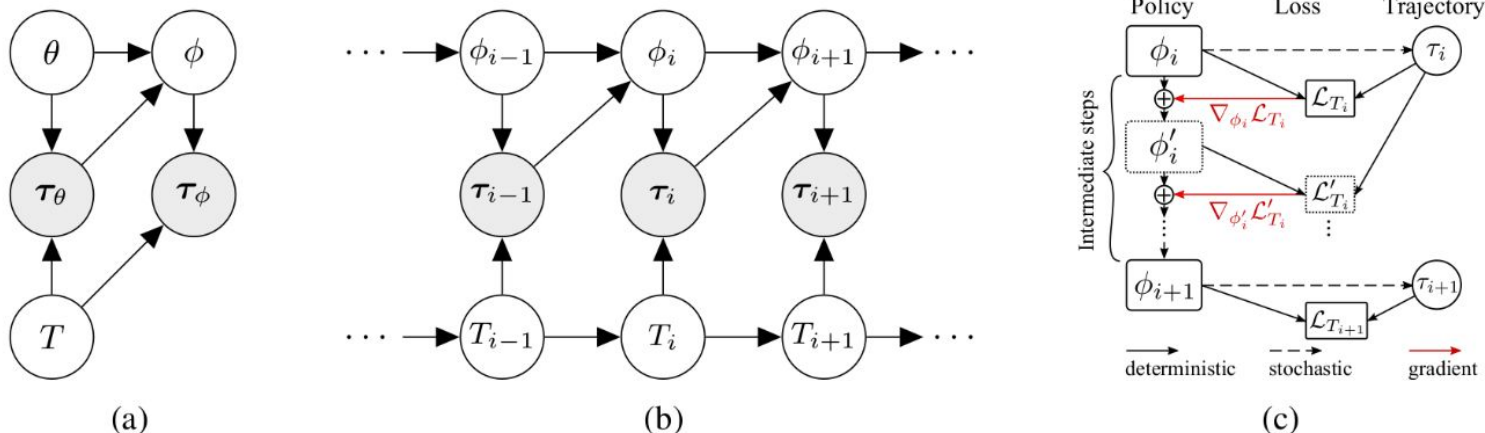
- Optimize meta-objective: PG where the gradient of  $R_T(\theta)$

$$\nabla_\theta \mathcal{R}_T(\theta) = \mathbb{E}_{\substack{\tau_\theta^{1:K} \sim P_T(\tau|\theta) \\ \tau_\phi \sim P_T(\tau|\phi)}} \left[ \mathcal{L}_T(\tau_\phi) \left[ \nabla_\theta \log \pi_\phi(\tau_\phi) + \nabla_\theta \sum_{k=1}^K \log \pi_\theta(\tau_\theta^k) \right] \right] \quad (4)$$

# Continuous adaptation via meta-learning

- $D(T)$  is defined by the environment changes, and the tasks become sequentially dependent
- Our goal:
  - Find the dependence between consecutive tasks
  - Meta-learn a rule to minimize total expected loss during interacting

# Continuous adaptation via meta-learning



- Here is probabilistic graph of MAML
  - (a) MAML in a multi-task RL setting
  - (b) Extend to continuous adaptation
    - Policy and trajectories at a previous step are used to construct a new policy for the current step, i.e.  $\phi_i, \tau_i \rightarrow \phi_{i+1}$
  - (c) Computation graph for the meta-update from  $\phi_i$  to  $\phi_{i+1}$ 
    - The model is optimized via truncated backpropagation through time starting from  $L_{T_{i+1}}$

# Continuous adaptation via meta-learning

- From (b) we represent nonstationary as a sequence of stationary tasks, then the goal is try to minimize the expected loss of  $L$  tasks, where  $L$  is a given number

$$\min_{\theta} \mathbb{E}_{\mathcal{P}(T_0), \mathcal{P}(T_{i+1}|T_i)} \left[ \sum_{i=1}^L \mathcal{R}_{T_i, T_{i+1}}(\theta) \right] \quad (5)$$

- Expected loss on a pair of consecutive tasks

$$\mathcal{R}_{T_i, T_{i+1}}(\theta) := \mathbb{E}_{\tau_{i,\theta}^{1:K} \sim P_{T_i}(\tau|\theta)} \left[ \mathbb{E}_{\tau_{i+1,\phi} \sim P_{T_{i+1}}(\tau|\phi)} [\mathcal{L}_{T_{i+1}}(\tau_{i+1,\phi}) \mid \tau_{i,\theta}^{1:K}, \theta] \right] \quad (6)$$

The principal difference between the loss in (3) and (6) is that trajectories  $\tau_{i,\theta}^{1:K}$  come from the current task,  $T_i$ , and are used to construct a policy,  $\pi_{\phi}$ , that is good for the upcoming task,  $T_{i+1}$ .

# Continuous adaptation via meta-learning

- To construct the parameters of policy for task  $\mathbf{T}_{i+1}$ , we start from  $\theta$ , note that here we can also change inner step size  $\alpha$

$$\begin{aligned}
 \phi_i^0 &:= \theta, \quad \tau_{\theta}^{1:K} \sim P_{T_i}(\tau \mid \theta), \\
 \phi_i^m &:= \phi_i^{m-1} - \alpha_m \nabla_{\phi_i^{m-1}} \mathcal{L}_{T_i} \left( \tau_{i, \phi_i^{m-1}}^{1:K} \right), \quad m = 1, \dots, M-1, \\
 \phi_{i+1} &:= \phi_i^{M-1} - \alpha_M \nabla_{\phi_i^{M-1}} \mathcal{L}_{T_i} \left( \tau_{i, \phi_i^{M-1}}^{1:K} \right)
 \end{aligned} \tag{7}$$

- PG for (c), as expectation is taken to both  $\mathbf{T}_i$  and  $\mathbf{T}_{i+1}$ , we change (4) to

$$\begin{aligned}
 &\nabla_{\theta, \alpha} \mathcal{R}_{T_i, T_{i+1}}(\theta, \alpha) = \\
 &\mathbb{E}_{\substack{\tau_{i, \theta}^{1:K} \sim P_{T_i}(\tau \mid \theta) \\ \tau_{i+1, \phi} \sim P_{T_{i+1}}(\tau \mid \phi)}} \left[ \mathcal{L}_{T_{i+1}}(\tau_{i+1, \phi}) \left[ \nabla_{\theta, \alpha} \log \pi_{\phi}(\tau_{i+1, \phi}) + \nabla_{\theta} \sum_{k=1}^K \log \pi_{\theta}(\tau_{i, \theta}^k) \right] \right]
 \end{aligned} \tag{8}$$

# MAML for Nonstationary Task : Training

---

## Algorithm 1 Meta-learning at training time.

---

**input** Distribution over pairs of tasks,  $\mathcal{P}(T_i, T_{i+1})$ ,  
learning rate,  $\beta$ .

1: Randomly initialize  $\theta$  and  $\alpha$ .

2: **repeat**

3:   Sample a batch of task pairs,  $\{(T_i, T_{i+1})\}_{i=1}^n$ .

4:   **for all** task pairs  $(T_i, T_{i+1})$  in the batch **do**

5:     Sample traj.  $\tau_{1:K}$  from  $T_i$  using  $\pi_\theta$ .

6:     Compute  $\phi = \phi(\tau_{1:K}, \theta, \alpha)$  as given in (7).

7:     Sample traj.  $\tau$  from  $T_{i+1}$  using  $\pi_\phi$ .

8:   **end for**

9:   Construct  $\nabla_\theta \mathcal{R}_T(\theta, \alpha)$  and  $\nabla_\alpha \mathcal{R}_T(\theta, \alpha)$  using  $\tau_{1:K}$  and  $\tau$  as given in (8).

10:   Update  $\theta \leftarrow \theta + \beta \nabla_\theta \mathcal{R}_T(\theta, \alpha)$ .

11:   Update  $\alpha \leftarrow \alpha + \beta \nabla_\alpha \mathcal{R}_T(\theta, \alpha)$ .

12: **until** Convergence

**output** Optimal  $\theta^*$  and  $\alpha^*$ .

---

• **Assumption: trajectories of  $T_i$  contain some information about  $T_{i+1}$**

• Difference from original MAML:

◦ Goal: 不仅仅优化  $\theta$  还有一系列步长  $\alpha$

◦ Use a pair of consecutive tasks for training

◦ During inner loop:

▪ 首先基于原策略  $\pi_\theta$  从  $T_i$  中获取一系列 trajectories  $\tau_{1:K}$

▪ 此处存疑, 原文中说获取trajectories的过程中还需要与  $T_{i+1}$  互动

▪ 根据这些  $T_i$  的 trajectories 计算  $\phi$

▪  $\pi_\phi$  用于解决  $T_{i+1}$ , 基于该策略从  $T_{i+1}$  中获取一个 trajectory  $\tau$

◦ Meta-update:

▪ 通过一系列子任务获取  $\tau_{1:K}$  以及  $\tau$  后, 计算  $\nabla_\theta \mathcal{R}_T(\theta, \alpha)$ ,  $\nabla_\alpha \mathcal{R}_T(\theta, \alpha)$

▪ 更新  $\theta$  与  $\alpha$



# MAML for Nonstationary Task : Testing

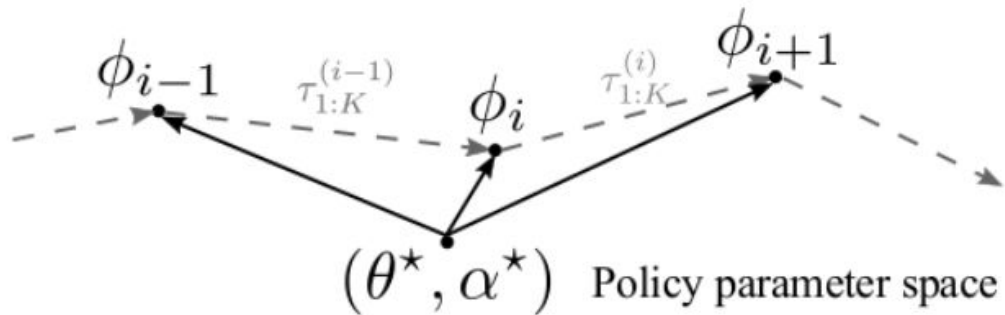
---

## Algorithm 2 Adaptation at execution time.

---

**input** A stream of tasks,  $T_1, T_2, T_3, \dots$

- 1: Initialize  $\phi = \theta$ .
- 2: **while** there are new incoming tasks **do**
- 3:   Get a new task,  $T_i$ , from the stream.
- 4:   Solve  $T_i$  using  $\pi_\phi$  policy.
- 5:   While solving  $T_i$ , collect trajectories,  $\tau_{1:K}^{(i)}$ .
- 6:   Update  $\phi \leftarrow \phi(\tau_{1:K}^{(i)}, \theta^*, \alpha^*)$  using importance-corrected meta-update as in (9).
- 7: **end while**



- Nonstationary  $\rightarrow$  can not access to same task multiple times
- How to handle : keep acting according to  $\pi_\phi$  and re-use past experience to for computing updates of  $\phi$  for each new incoming task
- $\pi_\phi$  获取的past experience与 $\pi_\theta$ 会有不同  $\rightarrow$  使用 importance weight correction 进行调整

$$\phi_i := \theta - \alpha \frac{1}{K} \sum_{k=1}^K \left( \frac{\pi_\theta(\tau^k)}{\pi_{\phi_{i-1}}(\tau^k)} \right) \nabla_{\theta} \mathcal{L}(\tau^k), \quad \tau^{1:K} \sim P_{T_{i-1}}(\tau \mid \phi_{i-1}), \quad (9)$$

# Experiment : RoboSumo



Figure 2: (a) The three types of agents used in experiments. The robots differ in the anatomy: the number of legs, their positions, and constraints on the thigh and knee joints. (b) The nonstationary locomotion environment. The torques applied to red-colored legs are scaled by a dynamically changing factor. (c) The RoboSumo environment.

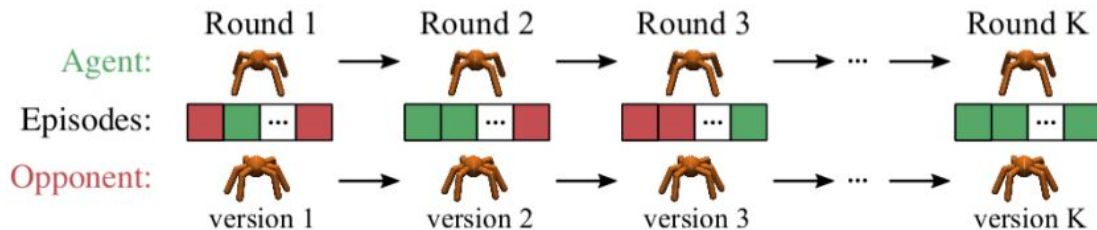


Figure 3: An agent competes with an opponent in an iterated adaptation games that consist of multi-episode rounds. The agent wins a round if it wins the majority of episodes (wins and losses illustrated with color). Both the agent and its opponent may update their policies from round to round (denoted by the version number).

# Conclusion and Future Work

- Regard nonstationarity as a sequence of stationary tasks
- Train agents to exploit the dependencies between consecutive tasks
- During testing, the model can handle similar nonstationarities + Problem : 该工作把动态环境看作一系列静态环境的序列, 但真实环境会更复杂
- Future work

Current	Future
One-step-ahead update of the policy	Extend to fully recurrent meta-updates → take into full history of interaction
(c) requires second order derivatives at training time	Utilize information provided by the loss but avoid explicit backpropagation through the gradients
Can't handle sparse reward → the meta-updates use policy gradients and heavily rely on the reward signal	Try to introduce auxiliary dense rewards designed to enable meta-learning