

# NoRML:No-Reward Meta Learning

presented by Jason TOKO

# Outline

- Background & Motivation
- Review on MAML-RL
- No-Reward Meta Learning
  - Learned Advantage Function
  - Learned Offset
  - Algorithm
- Experiment
  - Point Agent
  - Continuous Control
- Conclusion

# Outline

- Background & Motivation
- Review on MAML-RL
- No-Reward Meta Learning
  - Learned Advantage Function
  - Learned Offset
  - Algorithm
- Experiment
  - Point Agent
  - Continuous Control
- Conclusion

# Background & Motivation

- MAML
  - MAML is successful at adapting policies to different tasks that are defined by **reward changes**.
  - However, it is **less effective** when adapting to other changes, such as **dynamics changes**, **sensor drifts**, or **missing reward signals**.
- An agent should recognize dynamics change from the state-action transitions alone to adapt its behavior.
- Goal:
  - To develop a **model-free meta-RL** algorithm that can learn to quickly adapt a policy to dynamics changes and sensor drifts **without external rewards**——  
NoRML

# Outline

- Background & Motivation
- Review on MAML-RL
- No-Reward Meta Learning
  - Learned Advantage Function
  - Learned Offset
  - Algorithm
- Experiment
  - Point Agent
  - Continuous Control
- Conclusion

# Review on MAML-RL

- Adaptation(Policy Gradient):

$$\begin{aligned}\theta_i &= \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta, D_i^{\text{train}}) \\ &= \theta + \alpha \sum_{(s_t, \mathbf{a}_t, r_t) \in D_i^{\text{train}}} A^{\pi}(s_t, \mathbf{a}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | s_t).\end{aligned}$$

- Meta-training:

$$\theta \leftarrow \theta - \beta \sum_{\mathcal{T}_i \sim \mathcal{T}} \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}}),$$

where

$$\mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}}) = \mathcal{L}_{\mathcal{T}_i}(\theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta, D_i^{\text{train}}), D_i^{\text{test}}).$$

# Review on MAML-RL

- Algorithm

---

**Algorithm 1** MAML Training [10]

---

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

**Require:**  $\alpha$ : adaptation learning rate

**Require:**  $\beta$ : meta learning rate

Randomly initialize  $\theta$

**while** not done **do**

Sample a batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$

**for all**  $\mathcal{T}_i$  in batch **do**

Sample  $K$  trajectories  $D_i^{\text{train}}$  using  $\pi_\theta$  on task  $\mathcal{T}_i$ .

Compute adapted parameters using  $D_i^{\text{train}}$ :

$$\theta_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta, D_i^{\text{train}}).$$

Sample  $K$  trajectories  $D_i^{\text{test}}$  using  $\pi_{\theta_i}$  on task  $\mathcal{T}_i$ .

**end for**

Update  $\theta$  using all  $D_i^{\text{train}}$ , and  $D_i^{\text{test}}$ :

$$\theta \leftarrow \theta - \beta \sum_{\mathcal{T}_i \sim \mathcal{T}} \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}})$$

with  $\mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}}) = \mathcal{L}_{\mathcal{T}_i}(\theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta, D_i^{\text{train}}), D_i^{\text{test}})$ .

**end while**

---

---

**Algorithm 2** MAML Fine-tuning [10]

---

**Require:**  $\mathcal{T}_i$ : a new test task

**Require:**  $\theta, \alpha$ : from MAML training

Sample  $K$  trajectories  $D$  using  $\pi_\theta$  on task  $\mathcal{T}_i$ .

Compute fine-tuned policy:

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(\theta, D).$$

---

# Outline

- Background & Motivation
- Review on MAML-RL
- No-Reward Meta Learning
  - Learned Advantage Function
  - Learned Offset
  - Algorithm
- Experiment
  - Point Agent
  - Continuous Control
- Conclusion



# No-Reward Meta Learning

- NoRML aims to address the challenge of MAML, consisting of two additional components:
  - A **learned advantage function** that internalizes the reward in a way that allows for **reward-free adaptation**.
  - A **learned parameter offset** that enables **better exploration**.

# Learned Advantage Function

- Learned Advantage Function(LAF):  $A_\psi(s_t, a_t, s_{t+1})$
- Difference from  $A^\pi(s_t, a_t) = \sum_{t'=t}^H \gamma^{t'-t} r_{t'} - V^\pi(s_t)$  :
  - $A_\psi$  is a **feed-forward neural network** that takes in consecutive states and action  $(s_t, a_t, s_{t+1})$  and it's trained **end-to-end**.
  - $A_\psi$  is **only used during fine-tuning**, while the  $A^\pi$  is still used to compute the outer gradient during meta training.
  - $A_\psi$  is **not a true advantage function**, it is optimized to **transform or reshape** the policy gradient  $\nabla_\theta \log \pi_\theta(a_t | s_t)$  that achieves more **effective adaptation** in a single fine-tune step.

# Learned Advantage Function

- The benefit of LAF:
  - Since  $A_\psi$  takes in  $(s_t, a_t, s_{t+1})$  as input, this allows  $A_\psi$  to **detect changes** in the dynamics, and provide a **more informed evaluation** of the actions, compared to using only  $(s_t, a_t)$ .
  - We **eliminate the need to estimate the value function** to calculate the observed advantage (during meta-test fine-tuning)
  - $A_\psi$  **directly transforms** the policy gradient and can provide **accurate information** to the fine-tune step.
  - $A_\psi$  **keeps fixed** during meta-test fine-tuning.

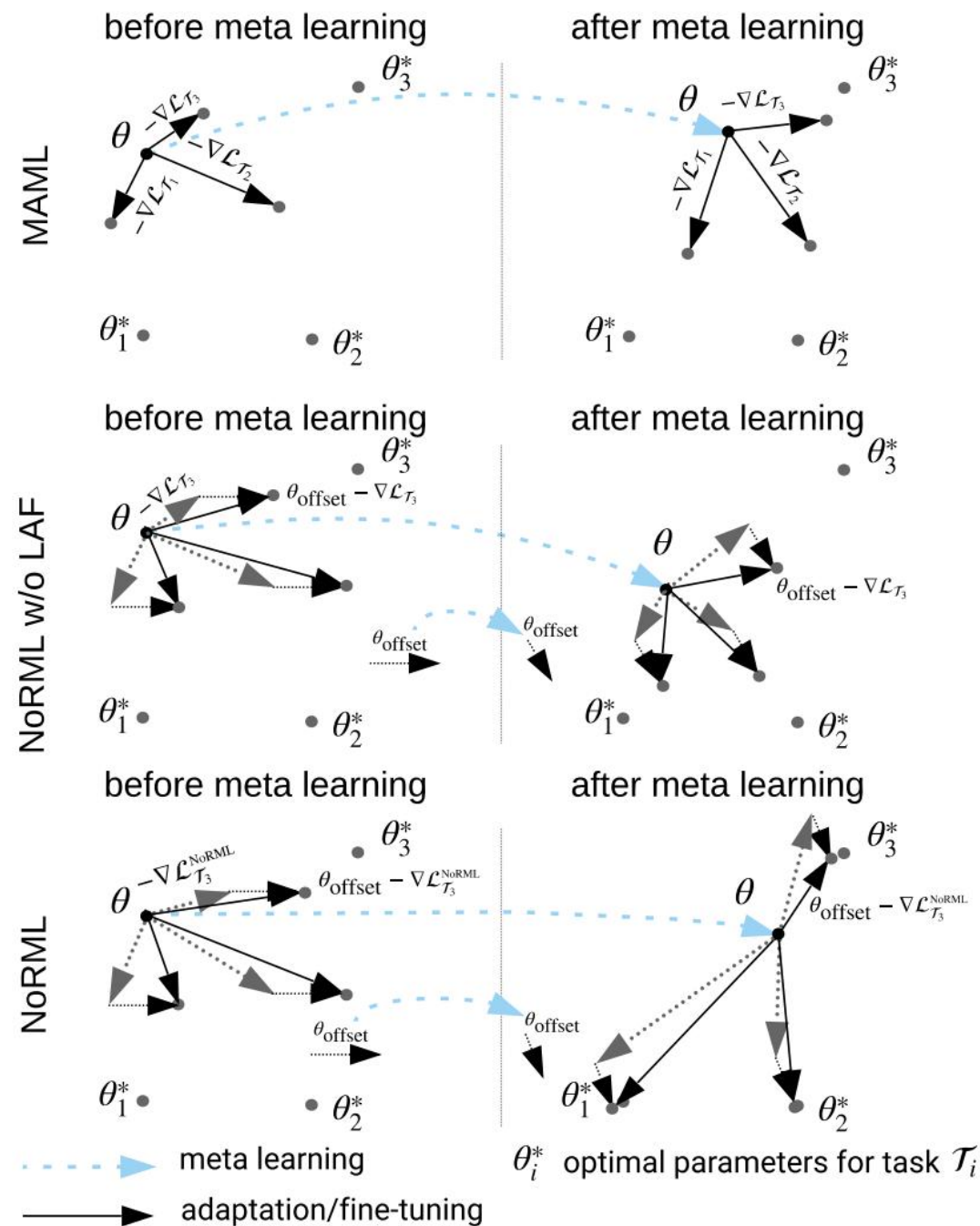
# Learned Offset

- One policy gradient step may be **insufficient to adapt** an exploratory meta-policy into a policy for the new task.
- A learned offset  $\vartheta_{\text{offset}}$  that is added to the policy parameters  $\vartheta$  when calculating an adapted policy:

$$\theta_i = \theta + \theta_{\text{offset}} - \alpha \sum_{D_i^{\text{train}}} A_{\psi}(s_t, a_t, s_{t+1}) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t).$$

# Learned Offset

Geometric interpretations of MAML, NoRML w/o Learned Advantage Function (LAF) and NoRML.



# Algorithm

---

## Algorithm 3 NoRML Training

---

Differences between MAML and NoRML are highlighted.

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

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

**Require:**  $\theta_\sigma$ : initial log standard deviation of meta-policy

Randomly initialize  $\theta_\mu$  and set  $\theta = \begin{bmatrix} \theta_\mu^T & \theta_\sigma^T \end{bmatrix}^T$

Randomly initialize  $\psi$

Initialize  $\theta_{\text{offset}}$  to  $[0 \dots 0]^T$

**while** not done **do**

    Sample a batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$

**for all**  $\mathcal{T}_i$  **do**

        Sample  $K$  trajectories without rewards using  $\pi_\theta$  on task  $\mathcal{T}_i$  and store all state transitions as a set  $D_i^{\text{train}} = \{(s_t, a_t, s_{t+1}) : \forall k < K, \forall t \leq H\}$ .

        Compute adapted parameters using  $D_i^{\text{train}}$ :

$\theta_i = \theta + \theta_{\text{offset}} - \alpha \sum_{D_i^{\text{train}}} A_\psi(s_t, a_t, s_{t+1}) \nabla_\theta \log \pi_\theta(a_t | s_t)$ .

        Sample  $K$  trajectories  $D_i^{\text{test}}$  using  $\pi_{\theta_i}$  on task  $\mathcal{T}_i$ .

**end for**

    Update  $\theta, \theta_{\text{offset}}$ , and  $\psi$  using all  $D_i^{\text{train}}$ , and  $D_i^{\text{test}}$ :

$\begin{bmatrix} \theta \\ \theta_{\text{offset}} \\ \psi \end{bmatrix} \leftarrow \begin{bmatrix} \theta \\ \theta_{\text{offset}} \\ \psi \end{bmatrix} - \beta \sum_{\mathcal{T}_i \sim \mathcal{T}} \nabla_{[\theta^T \theta_{\text{offset}}^T \psi^T]^T} \mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}})$

        with  $\nabla \mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}}) = \sum_{D_i^{\text{test}}} A^\pi(s_t, a_t) \nabla \log \pi_{\theta_i}(a_t | s_t)$ .

**end while**

---



---

## Algorithm 4 NoRML Fine-tuning

---

**Require:**  $\mathcal{T}_i$ : a task

**Require:**  $\theta, \theta_{\text{offset}}, \psi, \alpha$ : from NoRML training

    Sample  $K$  trajectories without rewards using  $\pi_\theta$  and store all state transitions as a set  $D = \{(s_t, a_t, s_{t+1}) : \forall k < K, \forall t < H\}$ .

    Compute fine-tuned policy:

$\theta \leftarrow \theta + \theta_{\text{offset}} - \alpha \sum_D A_\psi(s_t, a_t, s_{t+1}) \nabla_\theta \log \pi_\theta(a_t | s_t)$ .

---

Note  $r_t$  is only used to train  $A^\pi$ .  
During meta-test fine-tuning, we only apply learned offset and LAF without the requirement of reward.

# Outline

- Background & Motivation
- Review on MAML-RL
- No-Reward Meta Learning
  - Learned Advantage Function
  - Learned Offset
  - Algorithm
- Experiment
  - Point Agent
  - Continuous Control
- Conclusion

# Experiment

- Comparison objective:
  - vanilla MAML
  - Domain Randomization(DR): Implement DR by setting  $\alpha$  and  $\theta_{\text{offset}}$  to zero
  - NoRML w/o offset
  - NoRML w/o LAF
- Some details:
  - Policy:  $\pi_{\theta}(a_t|s_t) = \mathcal{N}(f(s_t|\theta_{\mu}), \text{diag}(e^{\theta_{\sigma}})^2)$  ( $\theta = [\theta_{\mu}^T \ \theta_{\sigma}^T]^T$ )
  - Apply **Meta-SGD** extension to **MAML**: replace the fixed inner learning rate of MAML with a learned vector of the same dimension as the policy parameter.
  - Apply **PPO** for **MAML**'s **adaptation and meta objective** to improve performance. For **NoRML** PPO is only used for the **meta objective**.
  - we use **polynomial regression** to fit the value function



# Point Agent

- Purpose: Study the adaptation to **dynamics change** and **reward change**.
- Task: move a point agent from (0,0) to (1,0);
- State:  $(x_t, y_t)$ ; Action:  $(dx_t, dy_t)$ ;
- Dynamics:  $\begin{bmatrix} x_{t+1} \\ y_{t+1} \end{bmatrix} = \begin{bmatrix} \cos \phi & -\sin \phi \\ \sin \phi & \cos \phi \end{bmatrix} \begin{bmatrix} dx_t \\ dy_t \end{bmatrix} + \begin{bmatrix} x_t \\ y_t \end{bmatrix}$  (one  $\phi$  per task);
- Reward setting:
  - **Shaped reward** case: negative Euclidean distance to (1, 0)
  - **Sparse reward** case: -1 for each step

# Point Agent

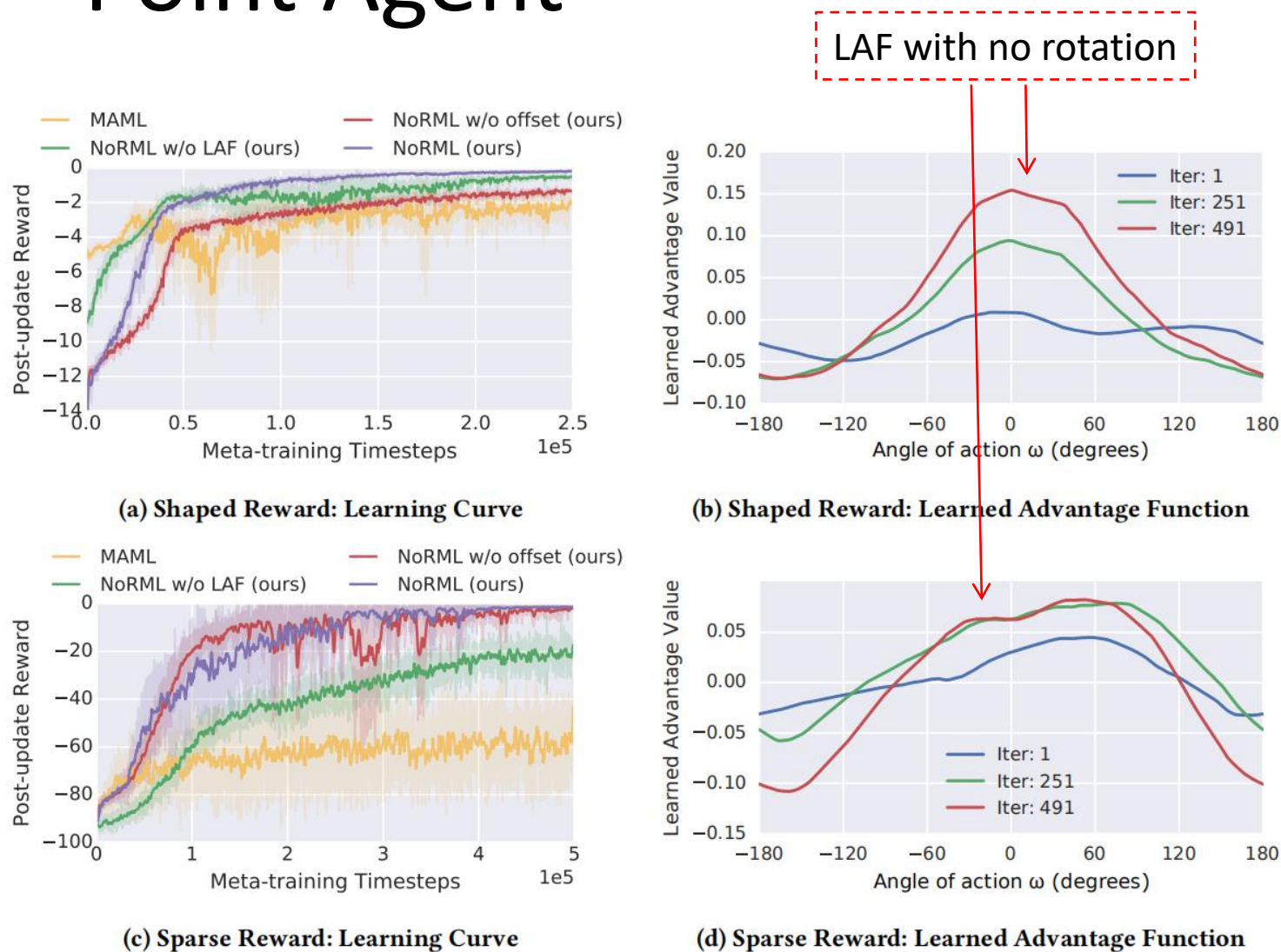


Figure 2

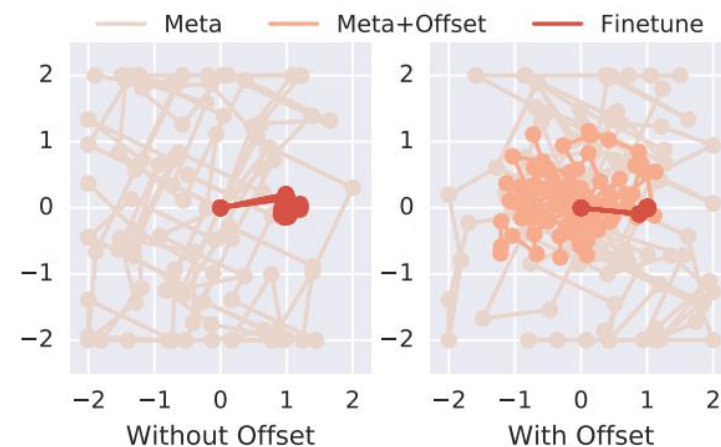


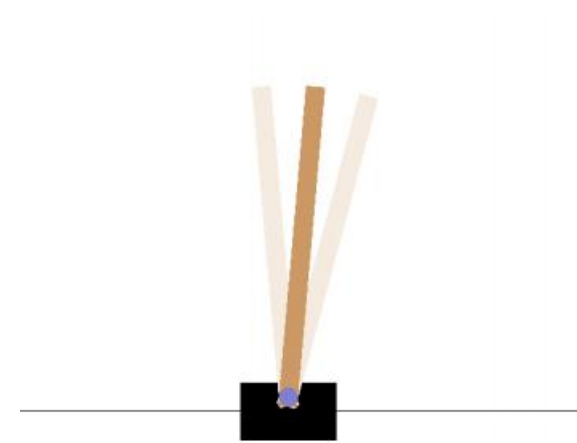
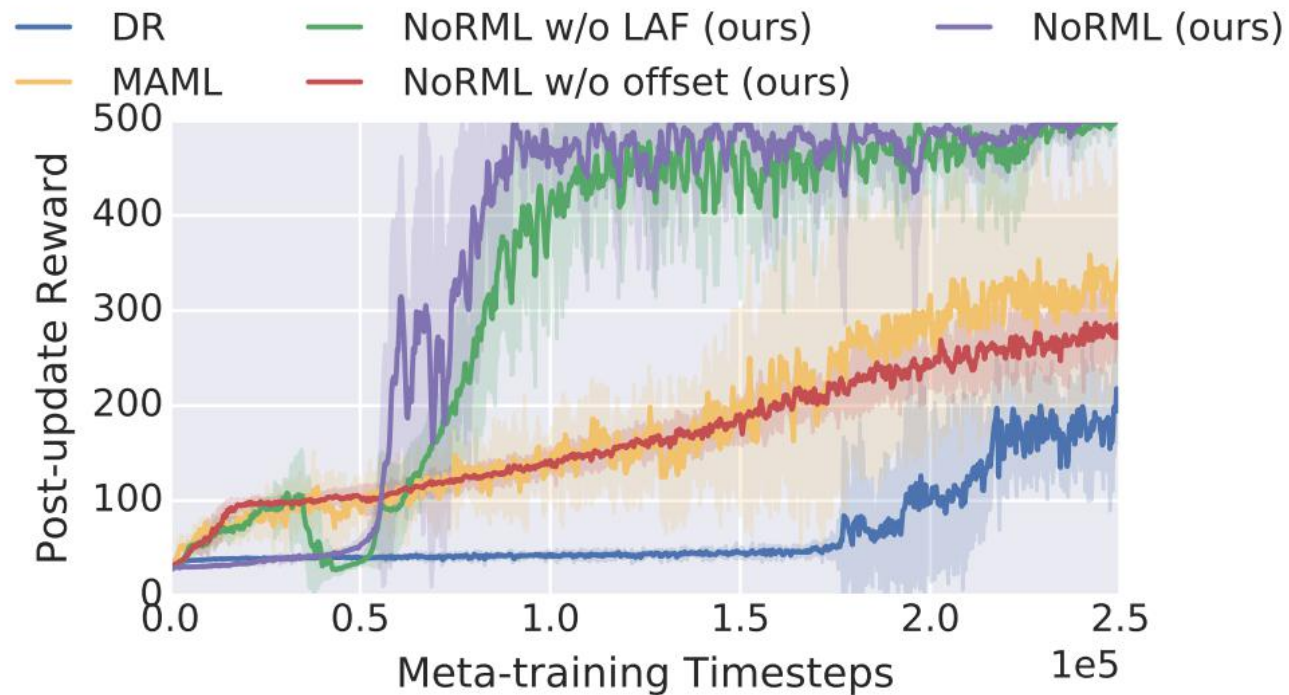
Figure 3: 10 rollout trajectories for NoRML policies trained with and without the offset. "Meta+Offset" means we only add the learned offset  $\theta_{\text{offset}}$  to the meta-policy parameters  $\theta$  without a gradient step and evaluate policy  $\pi_{\theta+\theta_{\text{offset}}}$ .

# Point Agent

- Impact of LAF:
  - In the shaped-reward case, MAML and NoRML both perform well.(Fig.2a)
  - In the sparse-reward case, MAML's performance degrades dramatically but LAF in NoRML enables the agent to adapt in one policy gradient step.(Fig.2c)
- Impact of the learned offset:
  - Better fine-tuning performance(Fig2a&c)
  - Reduce variance(Fig.3)

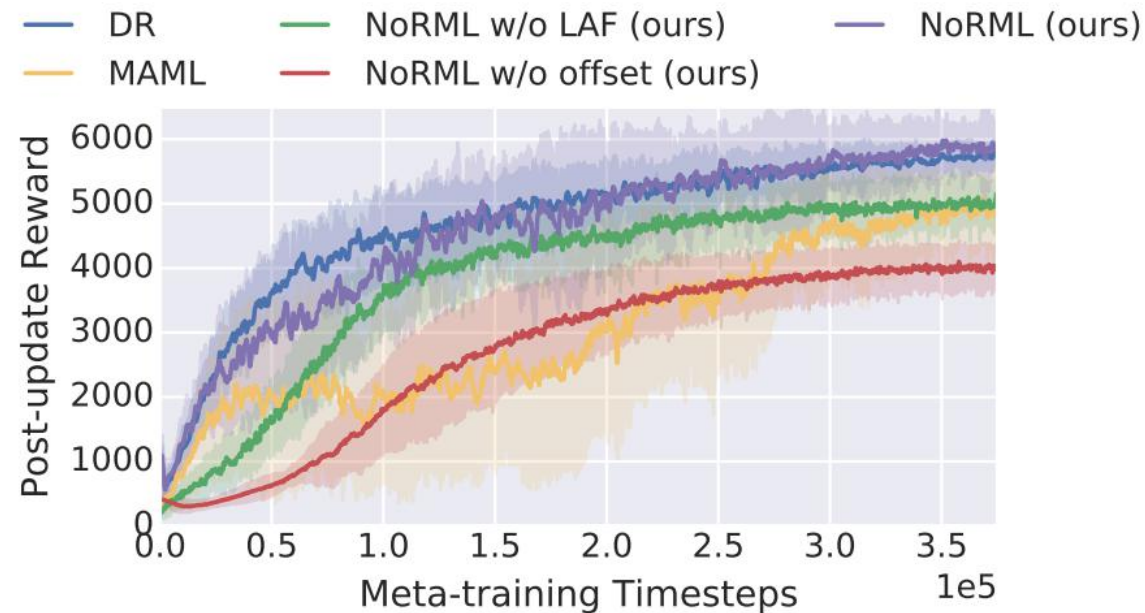
# Continuous Control

- Experiment 1: Cartpole with **sensor bias**.



# Continuous Control

- Experiment 2: Half Cheetah with **Swapped Actions**
  - Allow the torque outputs of the two hip joints to be swapped
  - Remove the **position** and **linear velocity** from half cheetah's observation space. (makes it **difficult** to compute the **distance-based reward function**)



# Continuous Control

- Exp1:
  - NoRML converges **faster** despite having **fewer assumptions**—it does not require an external reward signal for adaptation.
  - Without the offset, NoRML could not converge to a **high final reward**, and without the learned advantage function, convergence is **slower**.
- Exp2:
  - NoRML outperforms MAML both in **convergence speed and final return**.

# Outline

- Background & Motivation
- Review on MAML-RL
- No-Reward Meta Learning
  - Learned Advantage Function
  - Learned Offset
  - Algorithm
- Experiment
  - Point Agent
  - Continuous Control
- Conclusion

# Conclusion

- MAML cannot address dynamics change, sensor drift and missing reward signal well.
- By incorporating both a learned advantage function and a learned offset, NoRML can adapt to all these types of changes.
- LAF mainly **transforms** the policy gradient with more informed evaluation, leading **faster and robust adaptation**.
- Offset provides an **additional adjustment** to get **better fine-tuning performance**.