



Meta Reinforcement Learning and Imitation Learning

Prof. Joongheon Kim

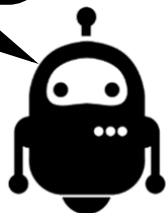
Korea University, School of Electrical Engineering

Artificial Intelligence and Mobility Laboratory

<https://joongheon.github.io>

joongheon@korea.ac.kr

I decide actions in the
direction of **maximizing**
the **cumulative reward**.



Agent

Markov Decision Process

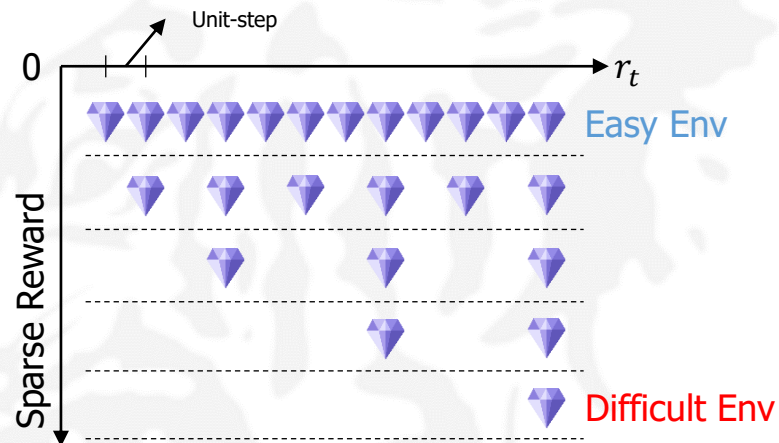
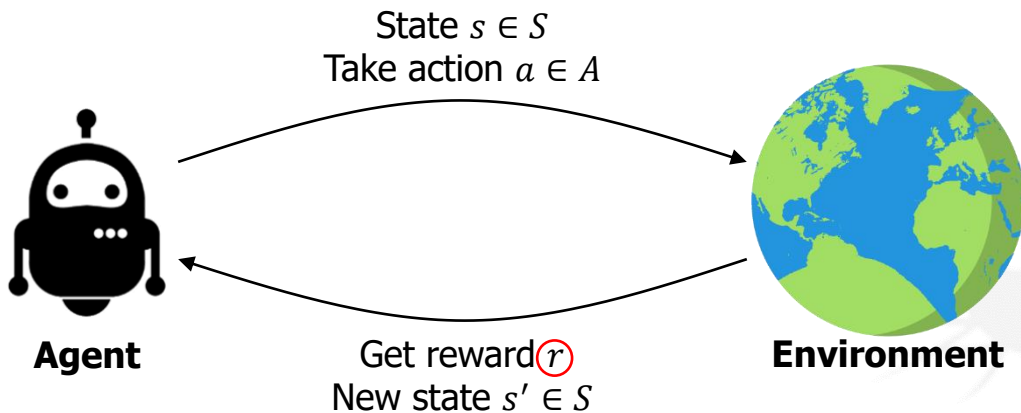
State $s \in \mathcal{S}$
Take action $a \in \mathcal{A}$



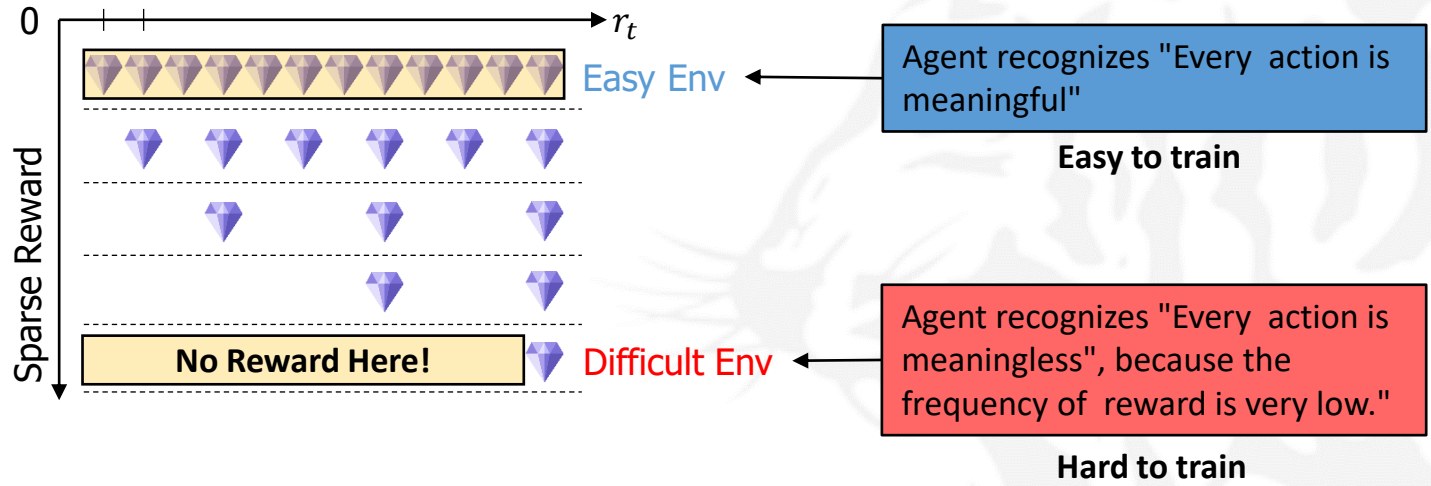
Environment

Get reward r
New state $s' \in \mathcal{S}$

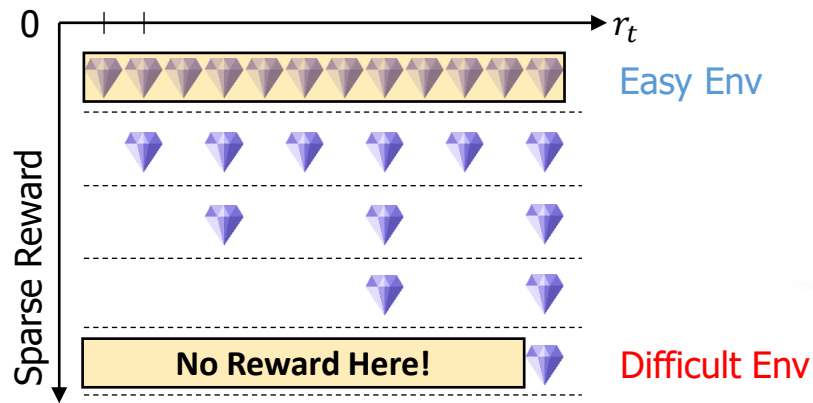
Reward Shaping is the
only way to achieve goal



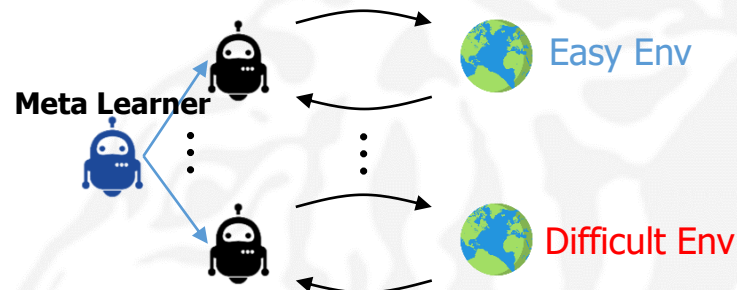
$$Q^*(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$$



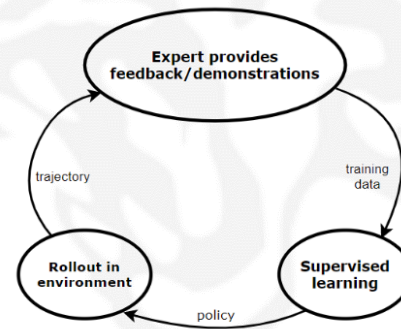
Two Approach to Solve



Meta Reinforcement Learning



Imitation Learning



**Meta
Reinforcement
Learning**

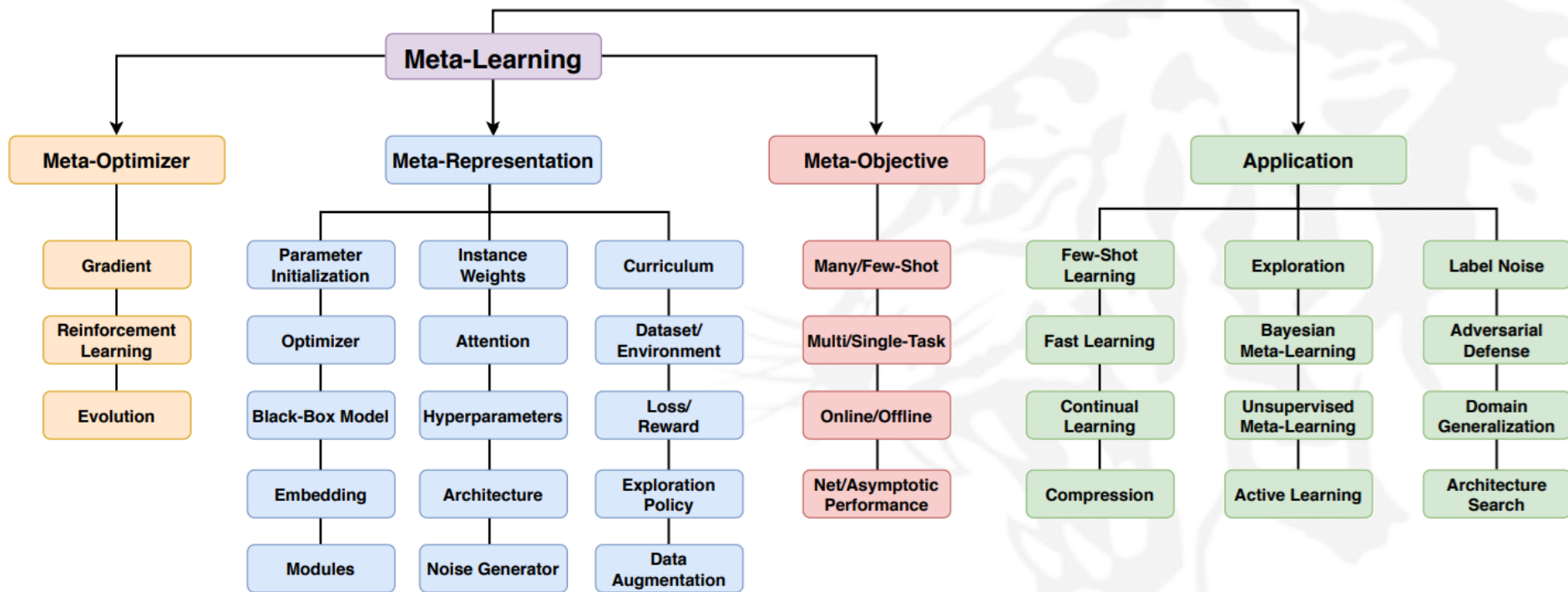
**Imitation
Learning**



Meta Reinforcement Learning



Meta Reinforcement Learning



[1] Timothy Hospedales, Antreas Antoniou, Paul Micaelli and Amos Storkey, "Meta-Learning in Neural Networks: A Survey," *arxiv*, 2020

Algorithm 1 Model-Agnostic Meta-Learning

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

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
- 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**

Meta parameter

θ

— meta-learning

--- learning/adaptation

Optimal parameter of Env3

θ_3^*

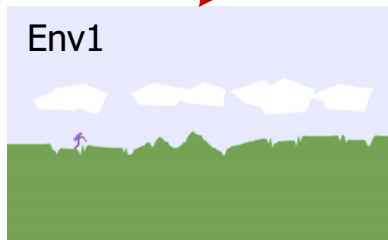
Optimal parameter of Env1

θ_1^*

Center of mass

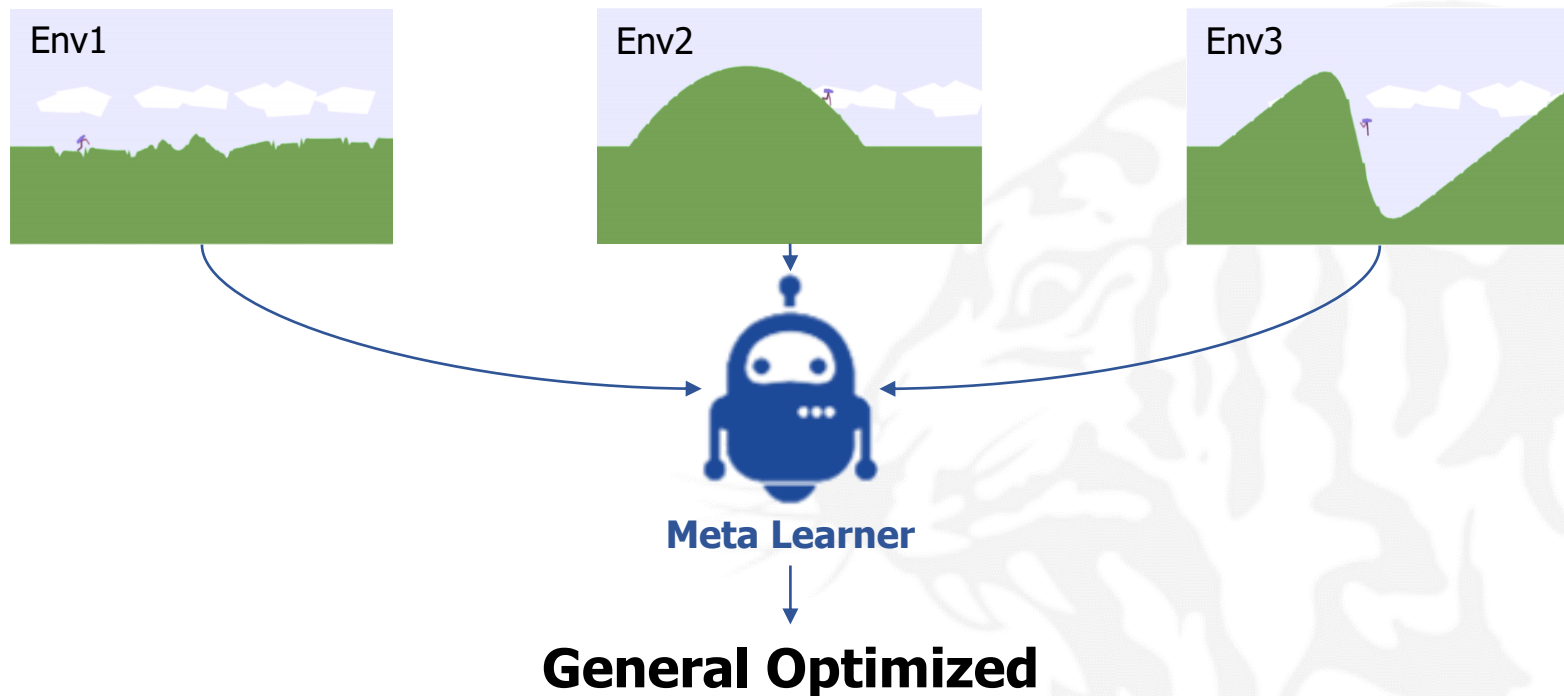
Optimal parameter of Env2

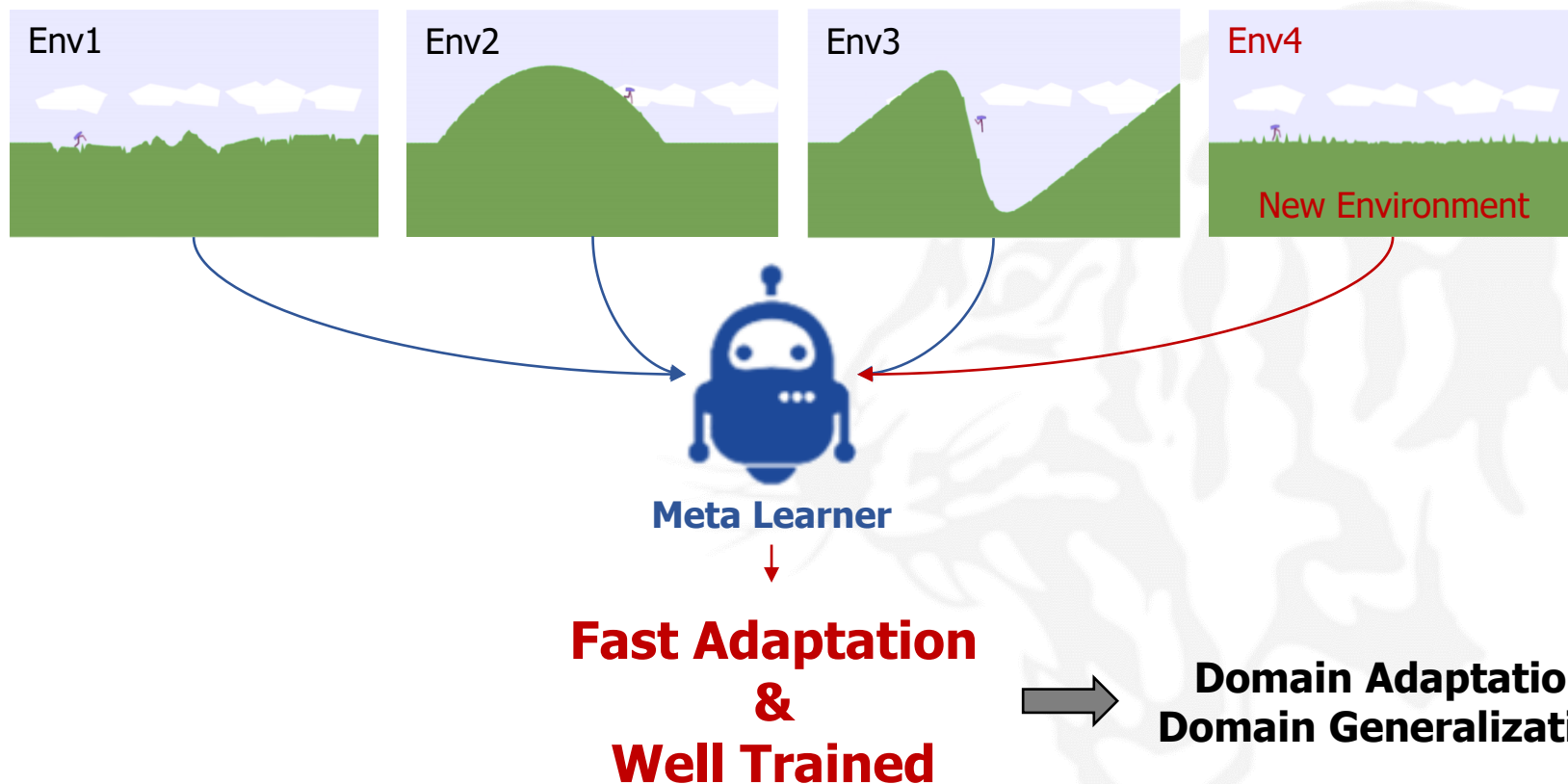
θ_2^*



[2] Chelsea Finn, Pieter Abbeel and Sergey Levine, "Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks," *ICML*, 2017

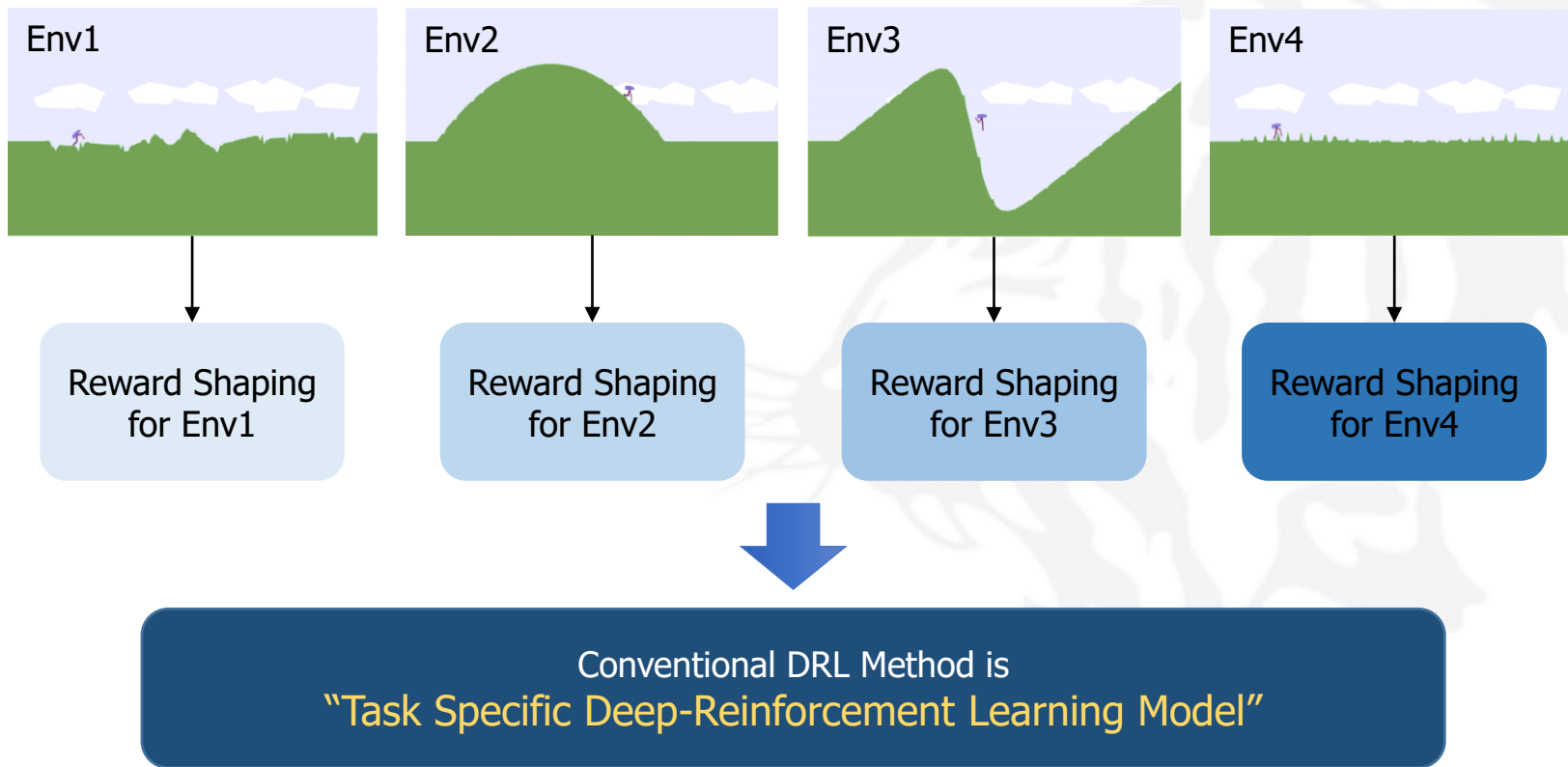
[3] Rui Wang et al., "Enhanced POET: Open-Ended Reinforcement Learning through Unbounded Invention of Learning Challenges and their Solutions," *arxiv*, 2020





[4] Code Available: <https://github.com/uber-research/poet>, https://github.com/cbfinn/maml_rl

To optimize for all environments with conventional DRL method,



MAML in Details

Algorithm 3 MAML for Reinforcement Learning

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

Require: α, β : step size hyperparameters

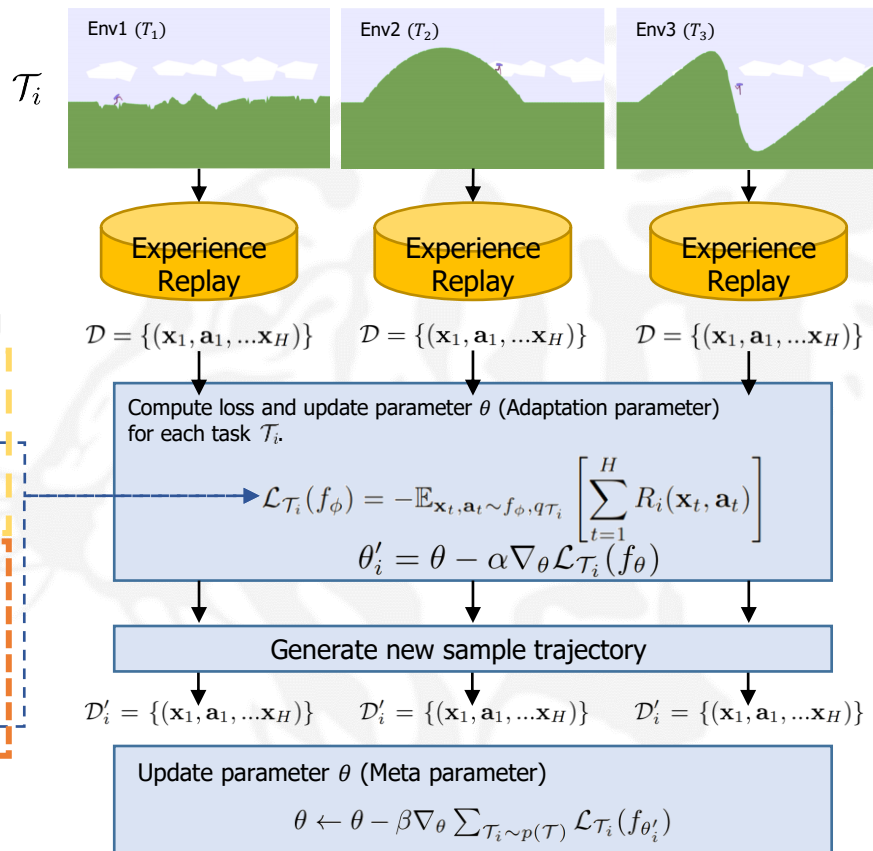
```

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  $f_\theta$ 
6:     in  $\mathcal{T}_i$ 
7:     Evaluate  $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$  using  $\mathcal{D}$  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
8:     Compute adapted parameters with gradient descent:
9:      $\theta'_i = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ 
10:    Sample trajectories  $\mathcal{D}'_i = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$  using  $f_{\theta'_i}$ 
11:    in  $\mathcal{T}_i$ 
12:  end for
13:  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  $\mathcal{D}'_i$ 
14:  and  $\mathcal{L}_{\mathcal{T}_i}$  in Equation 4
15: end while

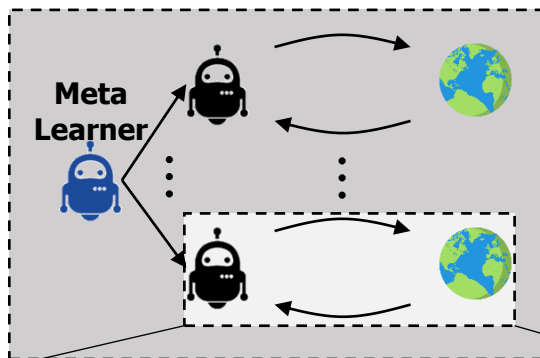
```

Fine Tuning

Update Meta-policy



Summary of MAML



Model-free Reinforcement Learning

Markov Decision Process: $\langle S, A, R, P, \gamma \rangle$

With **trajectory** $D = \{s_1, a_1, r_1, \dots, s_{T+1}\}$

$$\text{min: } \mathcal{L}(\theta, D) = -\mathbb{E}_{(s_t, a_t) \sim \pi_\theta} \left[\sum_{t=1}^H R(s_t, a_t) \right] = -\mathbb{E}_{(s_t, a_t) \sim \pi_\theta} \left[\sum_{t=1}^H r_t \right]$$

With **Gradient Descent Method**

$$\nabla_\theta \mathcal{L}(\theta, D) = -\mathbb{E}_{(s_t, a_t) \sim \pi_\theta} [A^\pi(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t | s_t)]$$

$$\text{s.t. } A^\pi(s_t, a_t) = \sum_{t'=t}^H \gamma^{t'-t} r_{t'} - V^\pi(s_t)$$

Model-Agnostic Reinforcement Learning

Tasks: $T = \{T_1, \dots, T_i, \dots, T_I\}$

Meta-Policy: π_θ

Fine-tuned Policy: $\{\pi_{\theta_1}, \dots, \pi_{\theta_i}, \dots, \pi_{\theta_I}\}$

With **Trajectory** $D_i^{\text{train}} = \{(s_1, a_1, r_1, \dots, s_{T+1})_i\}$:

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

With **Trajectory** $D_i^{\text{test}}, \theta_i$:

$$\begin{aligned} \theta &\leftarrow \theta - \beta \sum_{\mathcal{T}_i \sim \mathcal{T}} \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}}) \\ \text{s.t. } \mathcal{L}_{\mathcal{T}_i}(\theta_i, D_i^{\text{test}}) &= \mathcal{L}_{\mathcal{T}_i}(\theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(\theta, D_i^{\text{train}}), D_i^{\text{test}}) \\ &\rightarrow \text{Update meta-policy} \end{aligned}$$

NoRML: No-Reward Meta Learning

Model-Agnostic Reinforcement Learning

1) Tasks: $T = \{T_1, \dots, T_i, \dots, T_I\}$

2) Meta-Policy: π_θ

3) Fine-tuned Policy: $\{\pi_{\theta_1}, \dots, \pi_{\theta_i}, \dots, \pi_{\theta_I}\}$

With **Trajectory** $D_i^{\text{train}} = \{(s_1, \mathbf{a}_1, \mathbf{r}_1, \dots, s_{T+1})_i\}$:

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

$$= \theta + \alpha \sum_{(s_t, \mathbf{a}_t, \mathbf{r}_t) \in D_i^{\text{train}}} \boxed{A^\pi(s_t, \mathbf{a}_t)} \nabla_\theta \log \pi_\theta(\mathbf{a}_t | s_t)$$

→ fine-tuning policy

With **Trajectory** $D_i^{\text{test}}, \theta_i$:

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

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

→ Update meta-policy

$$A^\pi(s_t, \mathbf{a}_t) = \sum_{t'=t}^H \gamma^{t'-t} r_{t'} - V^\pi(s_t)$$

No Reward Meta Learning

Goal: develop model-free meta-RL algorithm that can learn to quickly adapt a policy to dynamics changes and sensor drifts *w/o external reward*.

Learned advantage function

$$A_\psi(s_t, \mathbf{a}_t, s_{t+1})$$

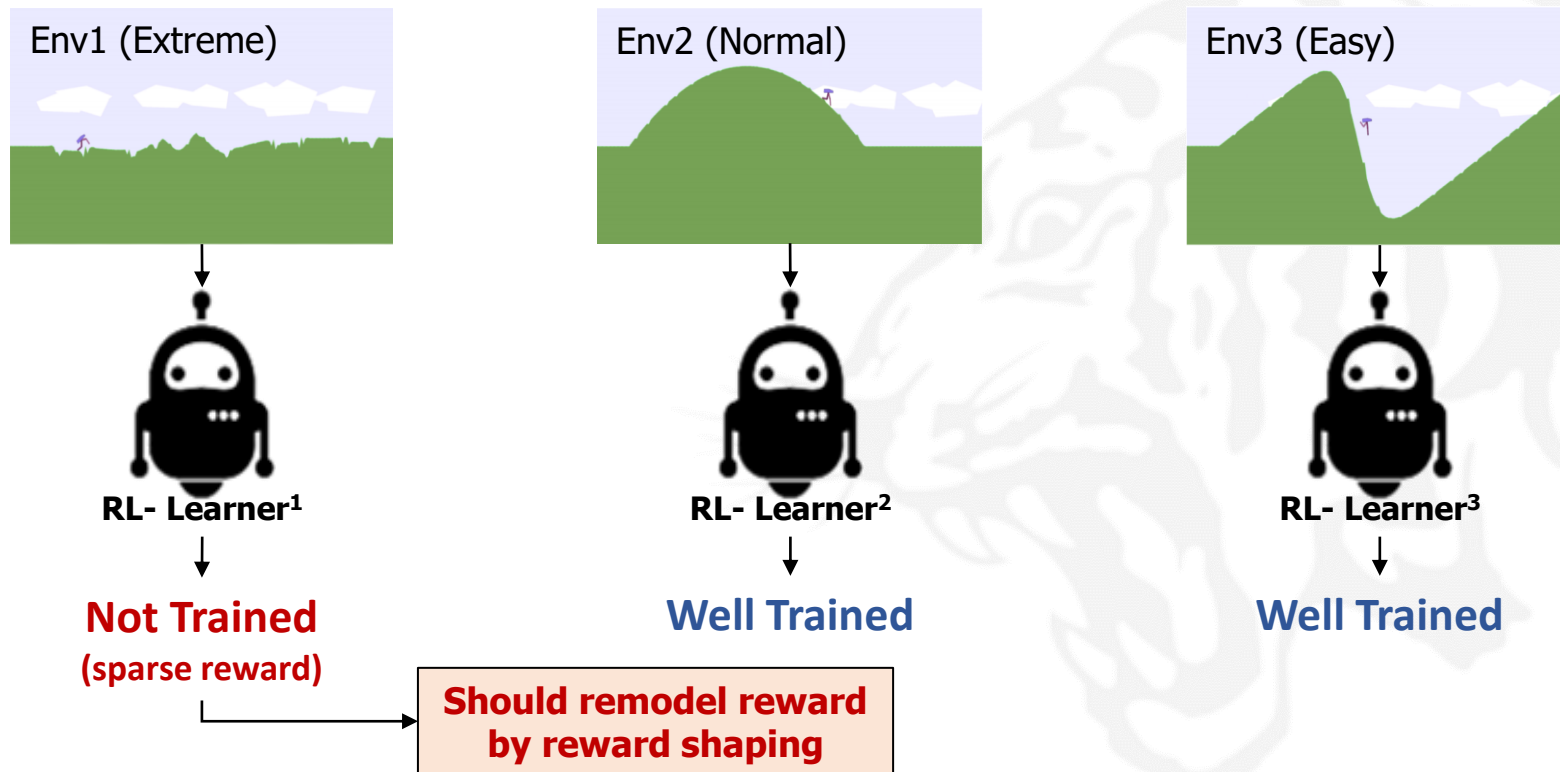
Offset learning for better exploration

[5] Yuxiang Yang, Ken Caluwaerts, Atil Iscen, Jie Tan, Chelsea Finn, "NoRML: No-Reward Meta Learning," *AAMAS*, 2019

[6] Code Available: <https://github.com/google-research/google-research/tree/master/norml>

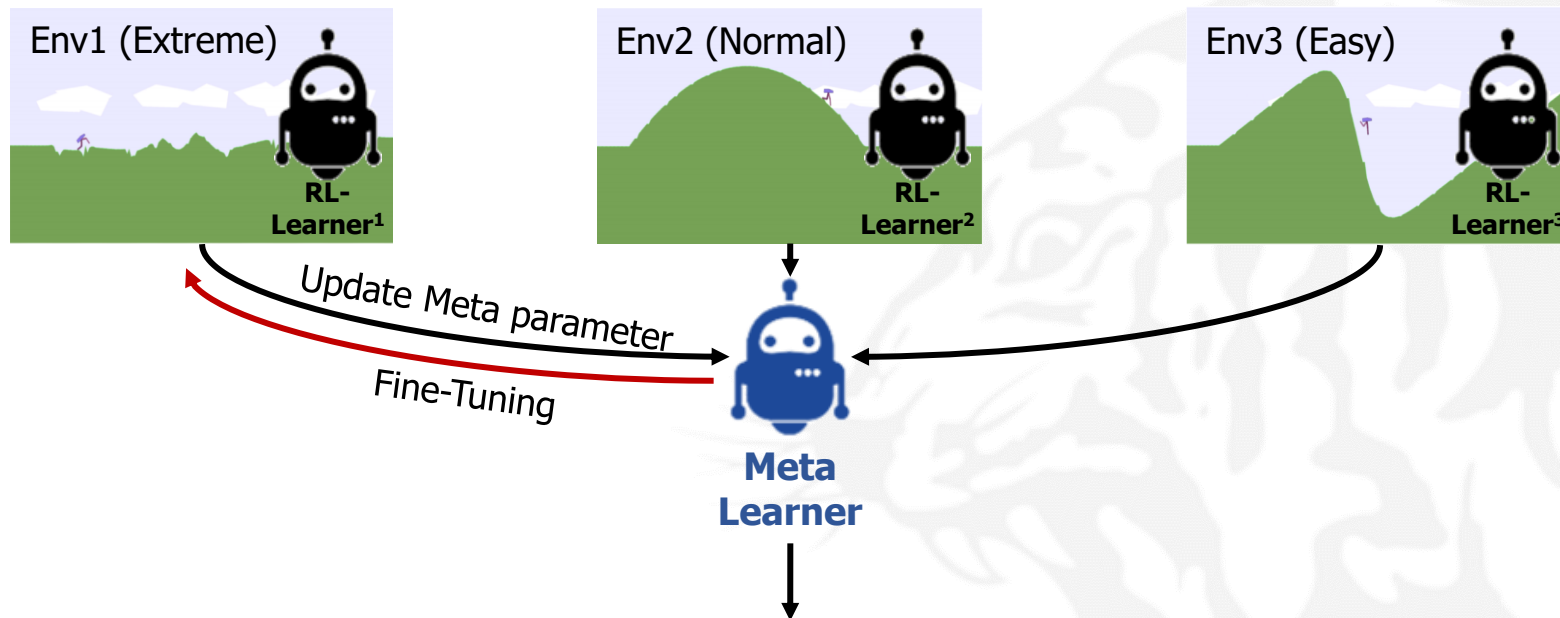
Summary of Meta-RL

In Previous RL (Task Specific RL Method)



Summary of Meta-RL

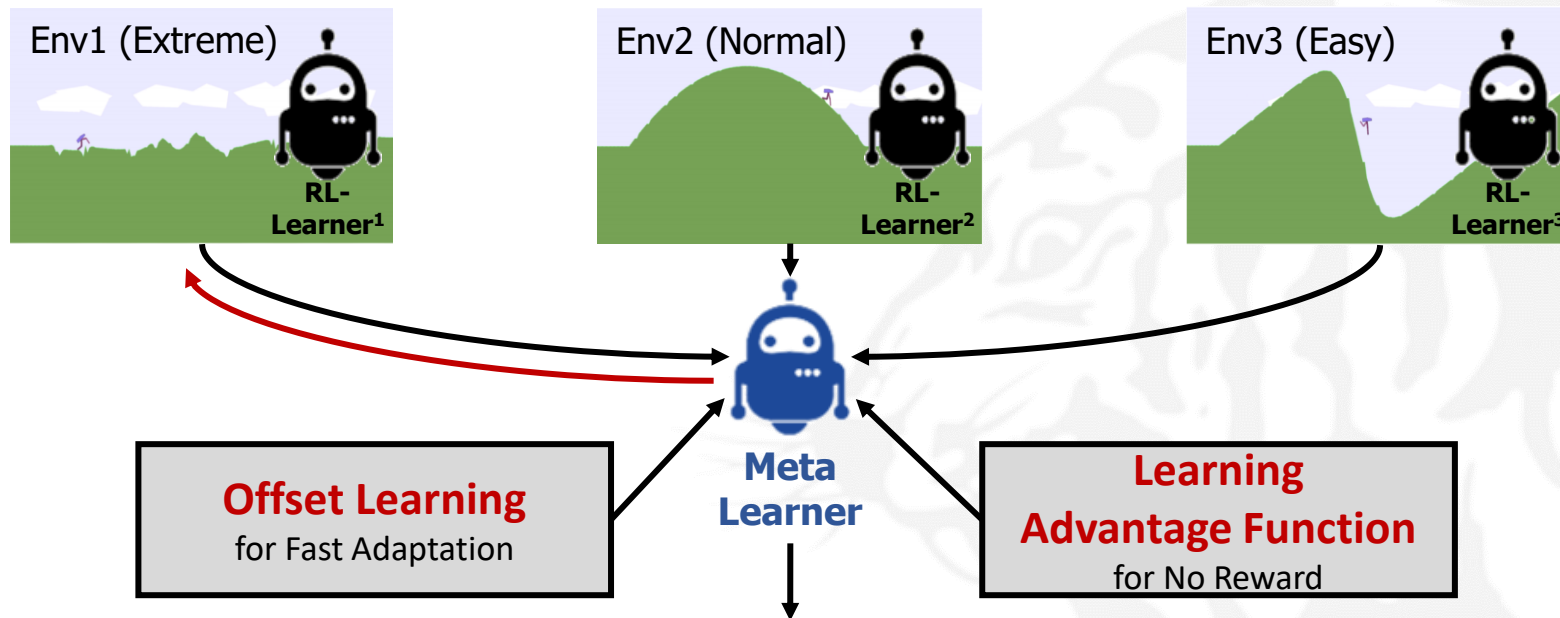
In MAML-RL (Task General RL Method)



**Fast Adaptation & Well Trained
but, still not guarantee sparse reward environment**

Summary of Meta-RL

In NoRML-RL (Task General RL Method)



Fast Adaptation & Well Trained
Good performance on sparse reward environment

Imitation Learning



- ICML 2018 Tutorial
 - <https://sites.google.com/view/icml2018-imitation-learning/>



Imitation Learning Tutorial ICML 2018

- ICML 2019 Tutorial
 - <https://slideslive.com/38917941/imitation-prediction-and-modelbased-reinforcement-learning-for-autonomous-driving>



Imitation, Prediction, and Model-Based Reinforcement Learning for Autonomous Driving

Sergey Levine

15th June 2019 - 10:50am



- Gameplay

Pro-Gamer



Trained Agent



The goal of Imitation Learning is to train a policy to mimic
the expert's demonstrations

Introduction to Imitation Learning

- Problems of RL



1. Reward Shaping



2. Safe Learning

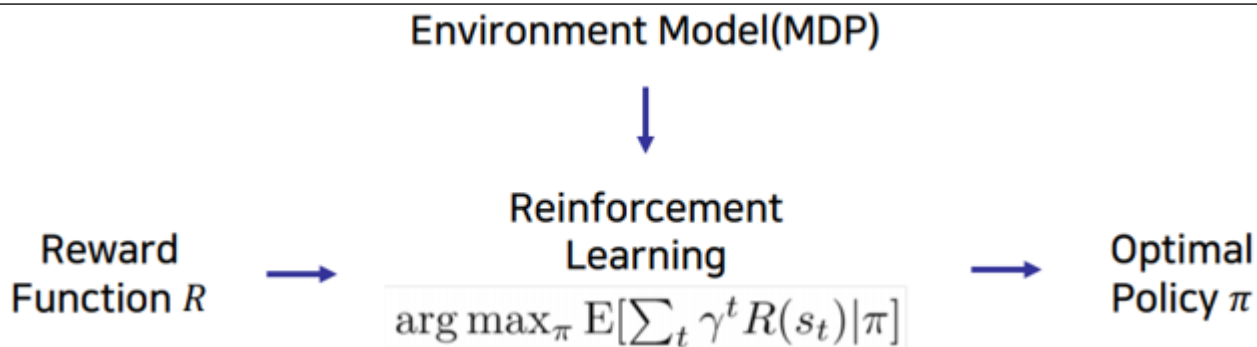


3. Exploration process

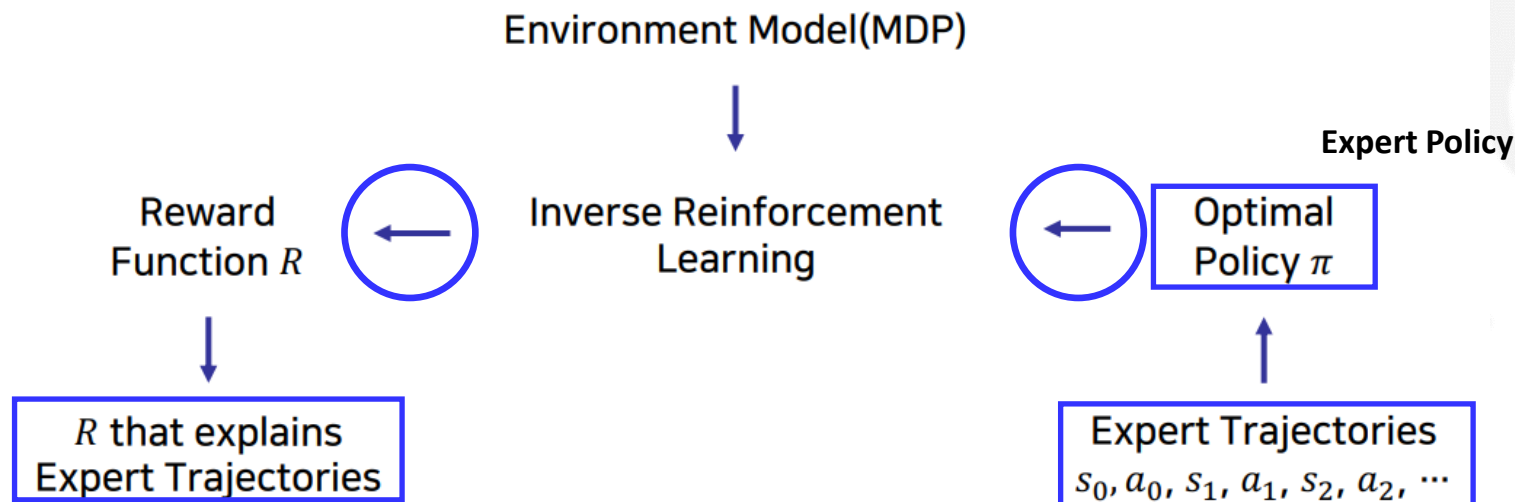
Imitation Learning **handles with** these problems
through the demonstration of the experts.

Inverse Reinforcement Learning (IRL)

RL



IRL



• Behavior Cloning

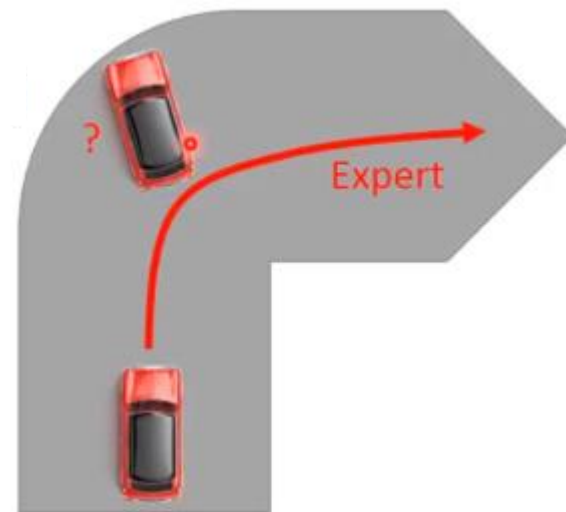
- Define $P^* = P(s|\pi^*)$ (distribution of states visited by **expert**)
- **Learning objective**

$$\operatorname{argmin}_{\theta} E_{(s,a_E) \sim P^*} L(a_E, \pi_{\theta}(s))$$

$$L(a_E, \pi_{\theta}(s)) = (a_E - \pi_{\theta}(s))^2$$

• Discussion

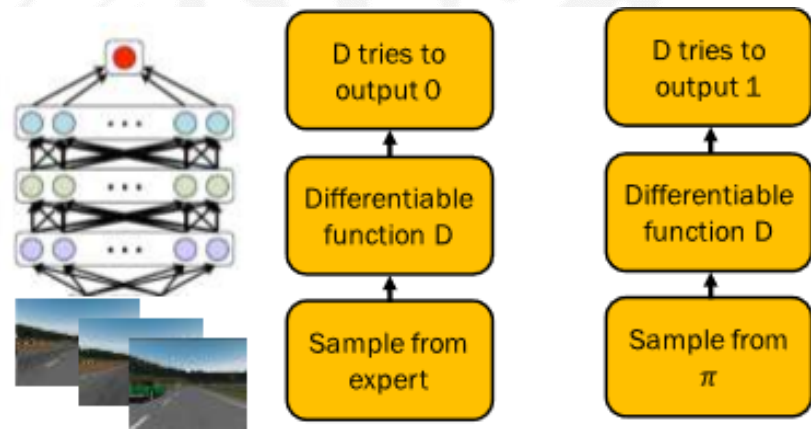
- Works well when P^* close to the distribution of states visited by π_{θ}
- **Minimize 1-step deviation error** along the expert trajectories



- **Generative Adversarial Imitation Learning (GAIL), NIPS 2016**

- Generative adversarial imitation learning (GAIL) learns a policy that can imitate expert demonstration using **the adversarial network** from generative adversarial network (GAN).
- **Learning Objective**

$$\operatorname{argmin}_{\theta} \operatorname{argmax}_{\phi} E[\log(D_{\phi}(s, a))] + E[\log(1 - D_{\phi}(s, a))]$$



Imitation Learning Applications: Starcraft2

• Starcraft2

States: $s = \text{minimap, screen}$

Action: $a = \text{select, drag}$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$

States: s

Action: a

Policy: π_{θ}

- Policy maps states to actions : $\pi_{\theta}(s) \rightarrow a$
- Distributions over actions : $\pi_{\theta}(s) \rightarrow P(a)$

State Dynamics: $P(s'|s,a)$

- Typically not known to policy
- Essentially the simulator/environment

Rollout: sequentially execute $\pi_{\theta}(s_0)$ on initial state

- Produce trajectories τ

$P(\tau|\pi)$: distribution of trajectories induced by a policy

$P(s|\pi)$: distribution of states induced by a policy



Imitation Learning Applications: Autonomous Driving

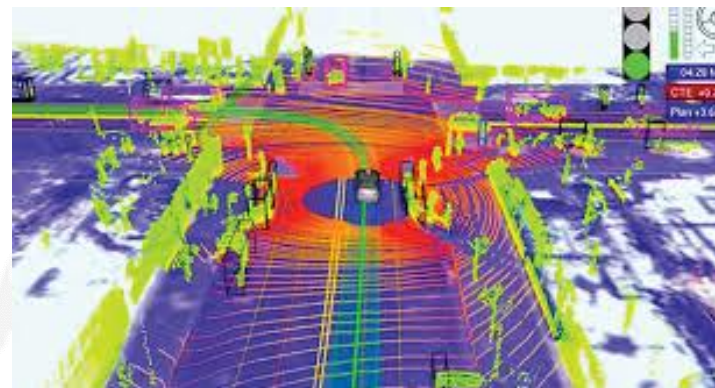
- Autonomous Driving Control

States: s = **sensors**

Action: a = **steering wheel, brake, ...**

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$



- Smartphone Security

States: $s = \text{apps}, \dots$

Action: $a = \text{use patterns}, \dots$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$



- PPF/RFTN Injection Control in Medicine

States: $s = \text{BIS}, \text{BP}, \dots$

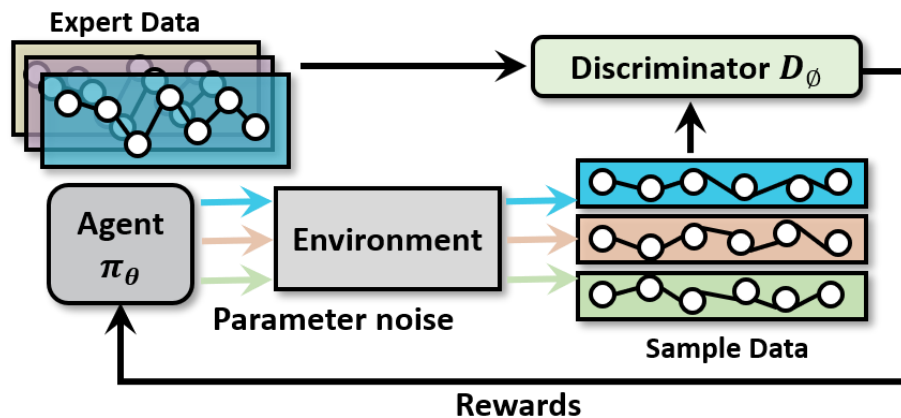
Action: $a = \text{PPF}, \text{RFTN}, \dots$

Training set: $D = \{\tau := (s, a)\}$ from expert

Goal: learn $\pi_{\theta}(s) \rightarrow a$

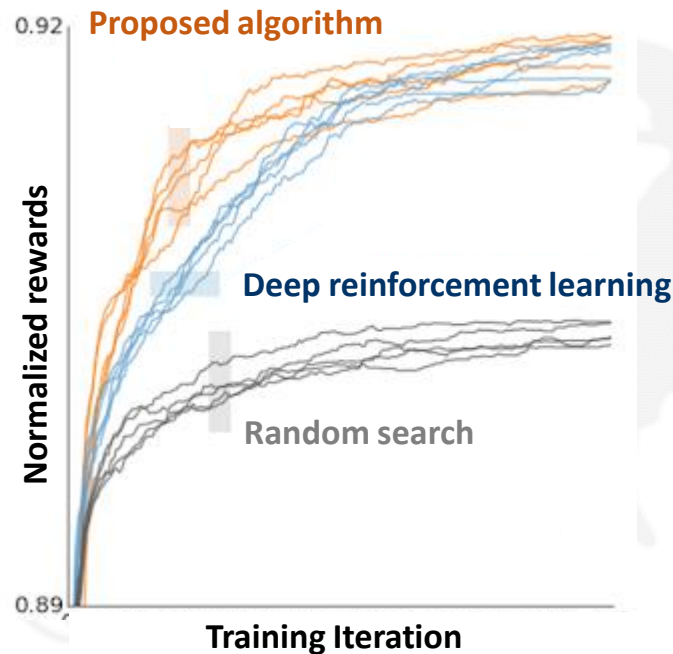


Autonomous Driving with Imitation Learning



M. Shin and J. Kim, "Adversarial Imitation Learning via Random Search in Lane Change Decision-Making," *ICML 2019 Workshop on AI for Autonomous Driving*, 2019.

M. Shin and J. Kim, "Randomized Adversarial Imitation Learning for Autonomous Driving," *IJCAI*, 2019., (Acceptance Rate: 850/4752=17.89%)



Generative Adversarial Network (GAN) + Random Search for Autonomous Driving

Thank you for your attention!

- More questions?
 - joongheon@korea.ac.kr

