

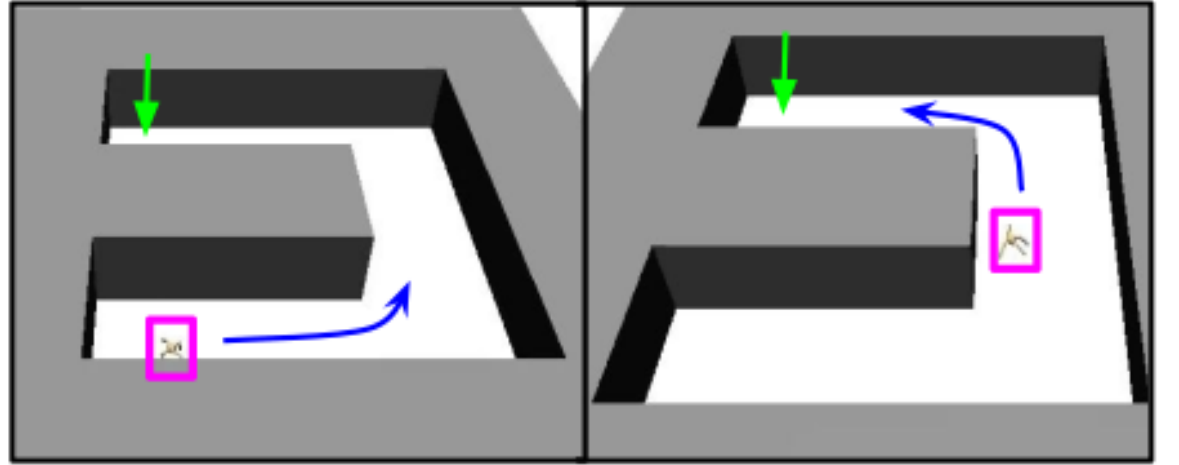
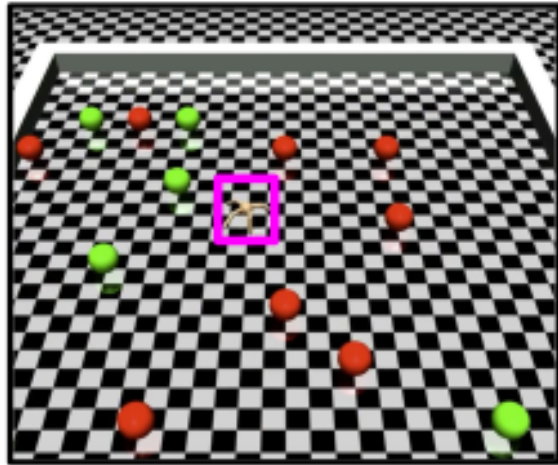
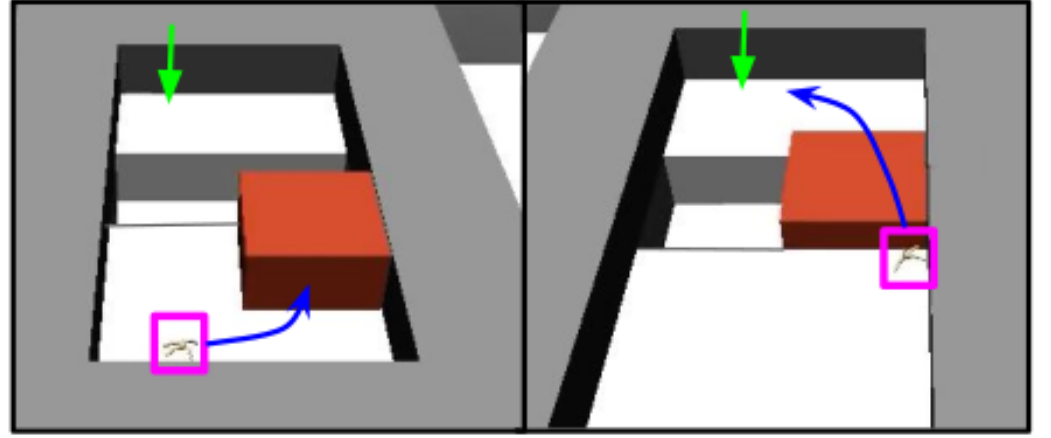
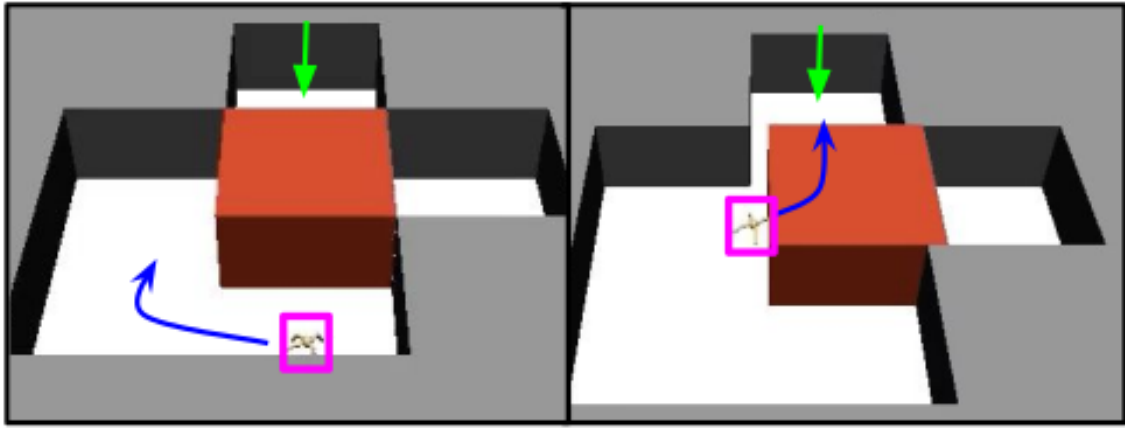
Hierarchical Deep Reinforcement Learning

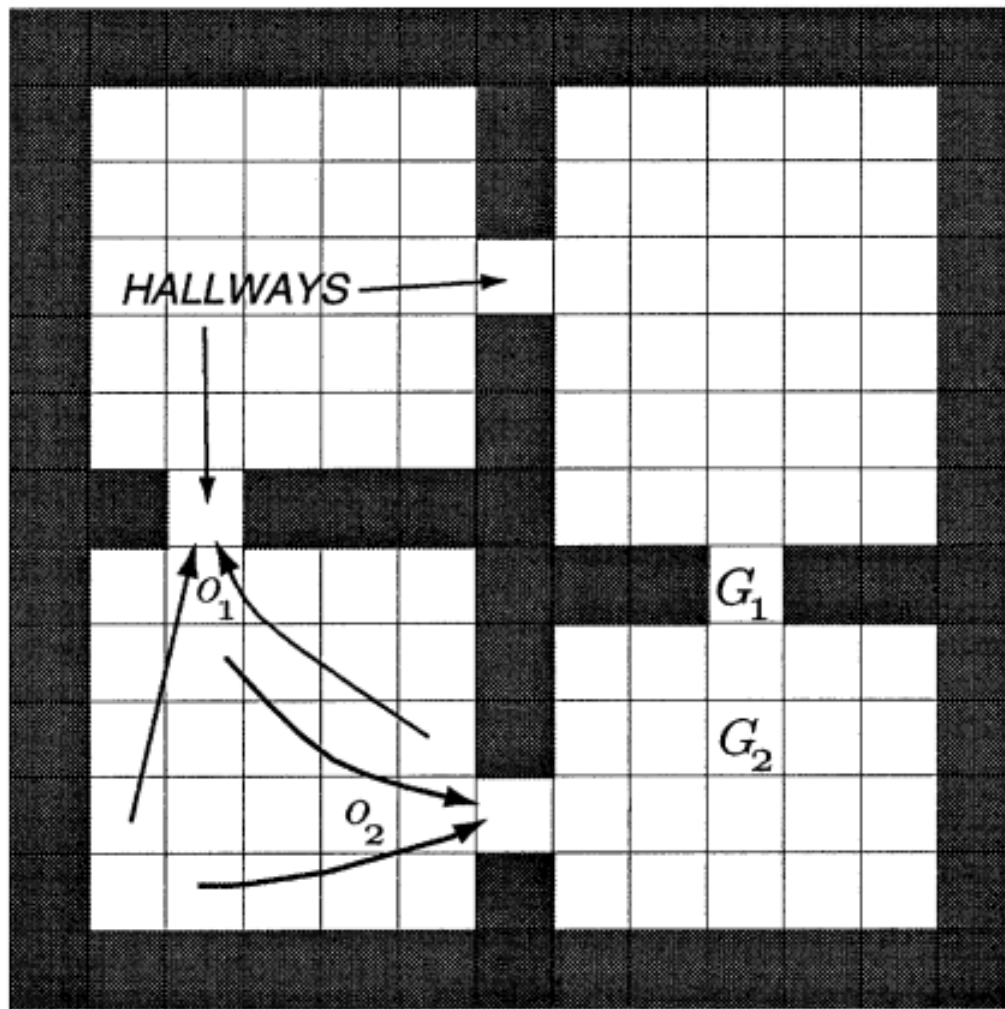
Seminar DRL

Lucas Brunner 17.03.2020

Goals

- What problem arise when doing HRL
- How can one solve them
- How are methods connected





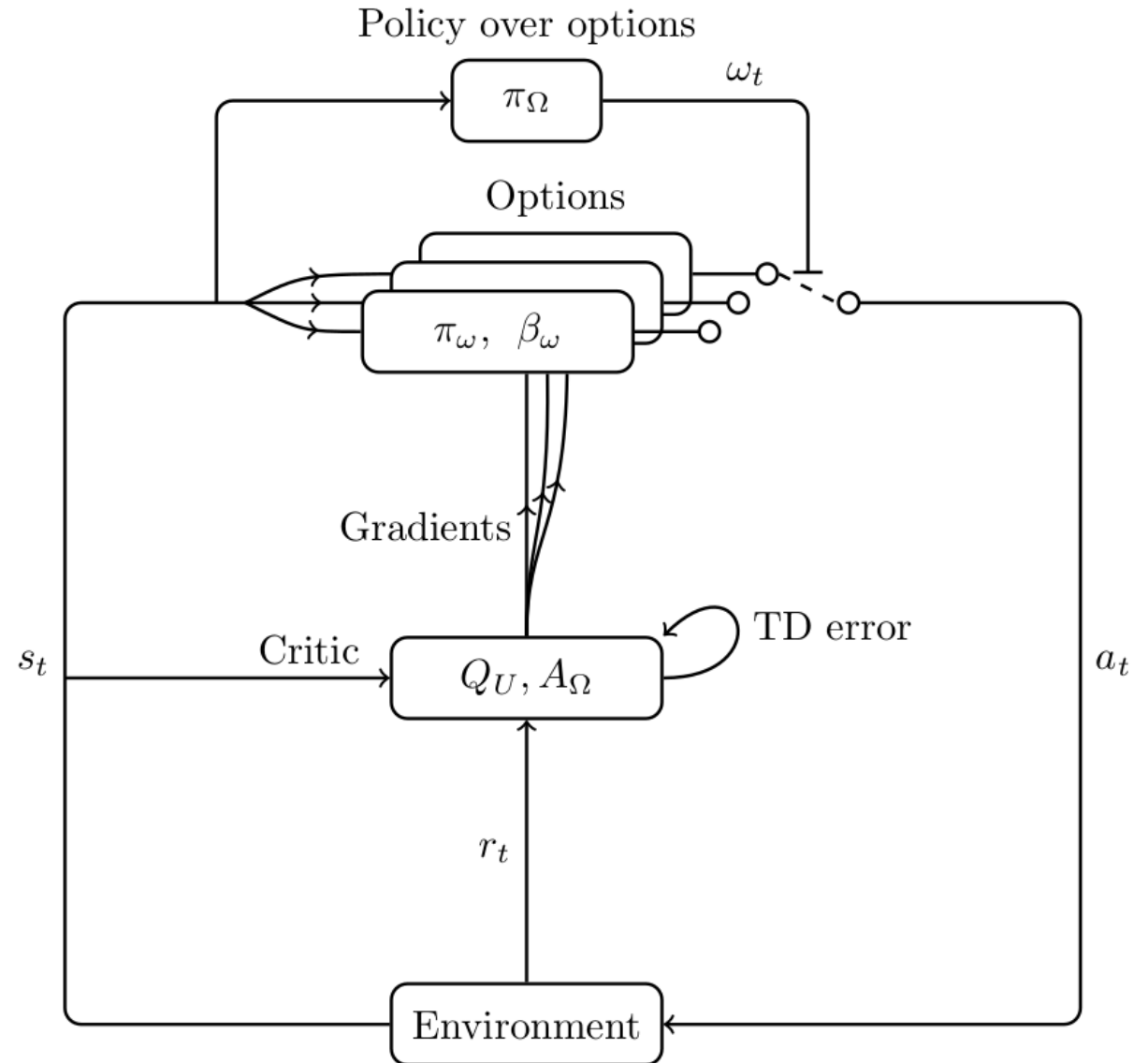
4 stochastic primitive actions



*8 multi-step options
(to each room's 2 hallways)*

(OC) The option-critic architecture

Bacon, Pierre-Luc, Jean Harb, and Doina Precup. 2017



Advantage Function

$$A(s,w) = Q(s,w) - V(s)$$

("measures how much better w is compared to alternative actions)

0 if w is optimal action

<0 if w is suboptimal

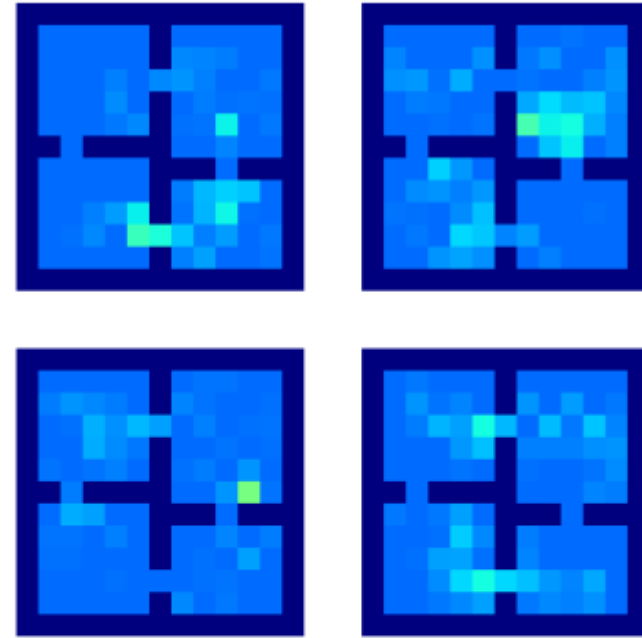
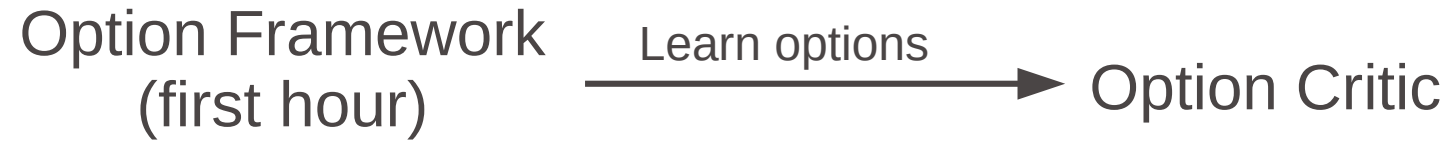
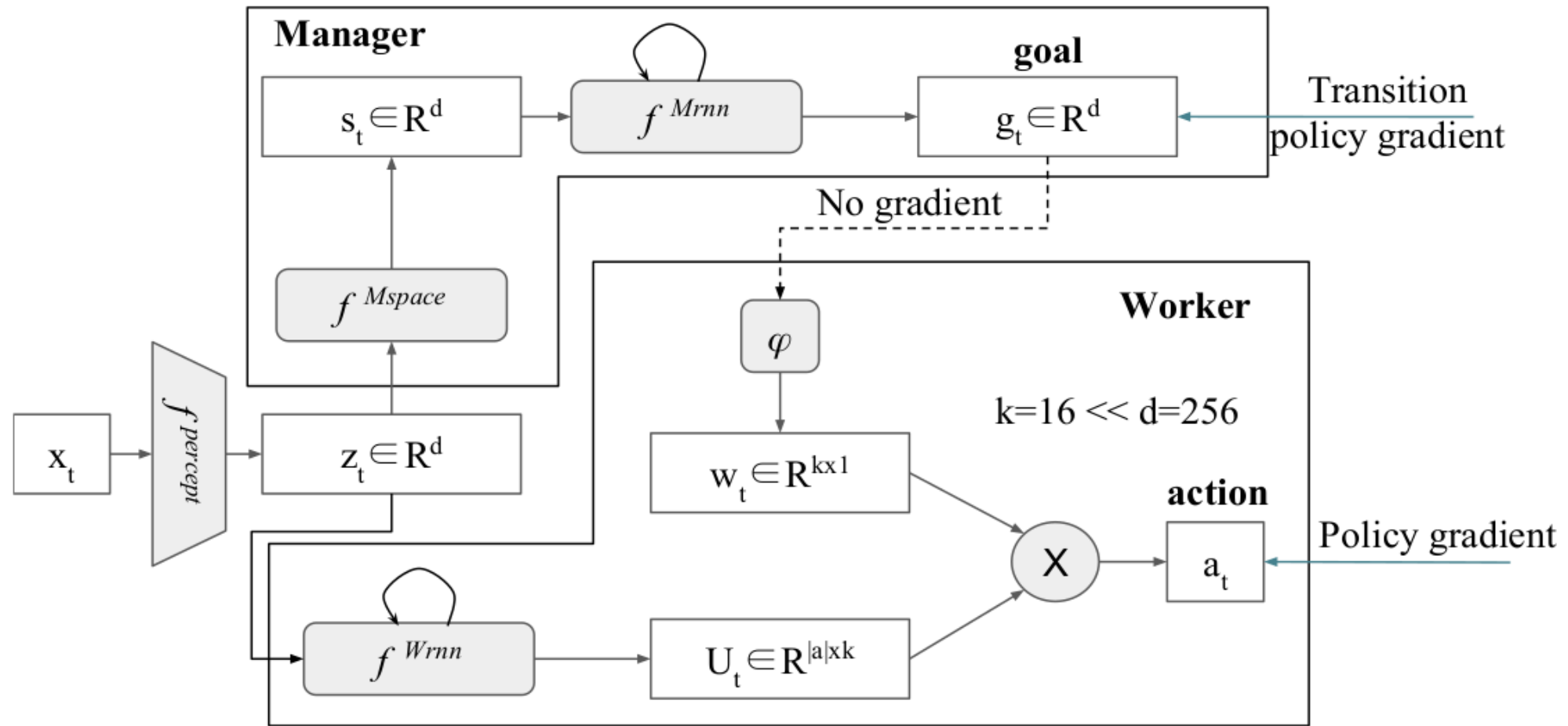


Figure 3: Termination probabilities for the option-critic agent learning with 4 options. The darkest color represents the *walls* in the environment while lighter colors encode higher termination probabilities.



FeUdal Networks for Hierarchical Reinforcement Learning

Vezhnevets, Alexander Sasha, et al. 2017



Option Framework
(first hour)

Learn options

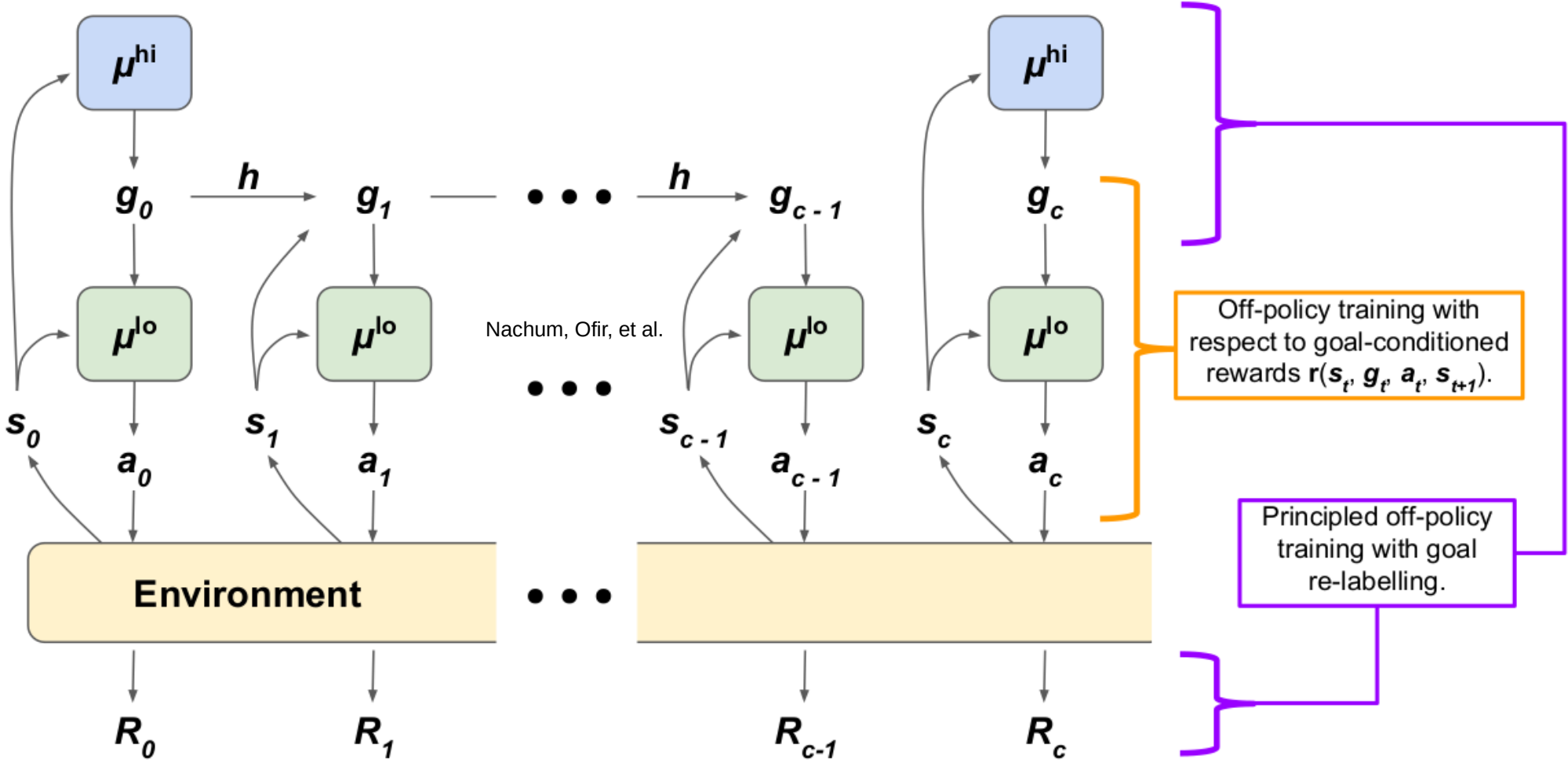
Option Critic

Decoupling Manager and Worker

FeUdal Net



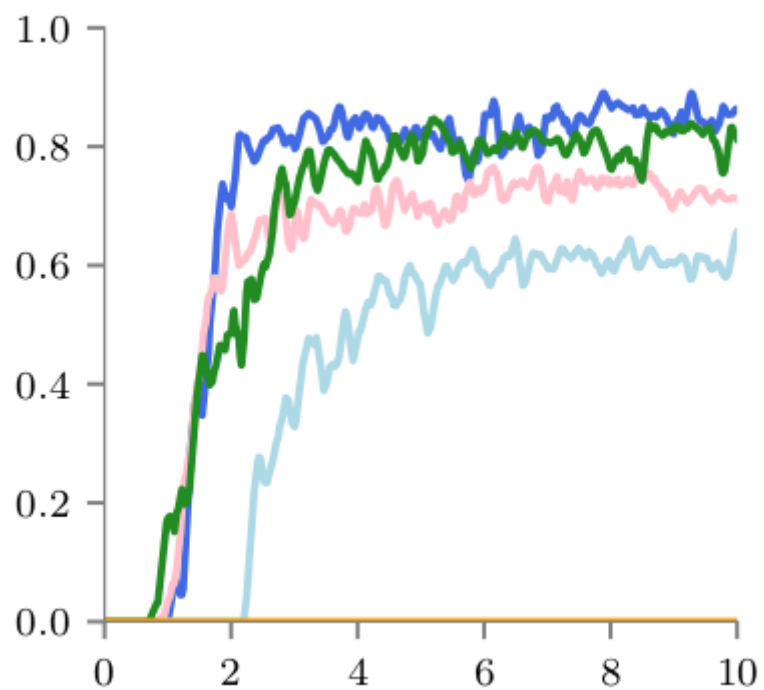
Data-Efficient Hierarchical Reinforcement Learning (HIRO) Nachum, Ofir, et al. 2017



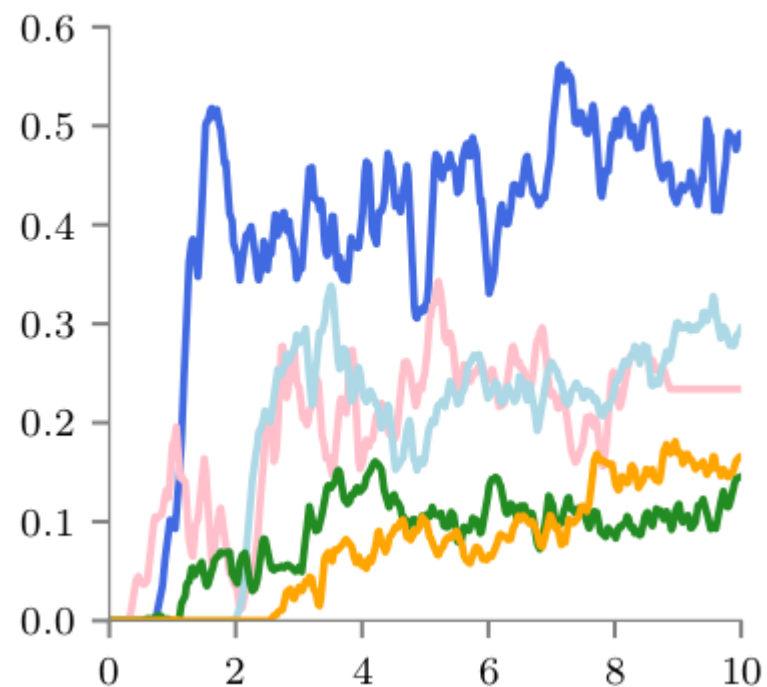
Does it help?

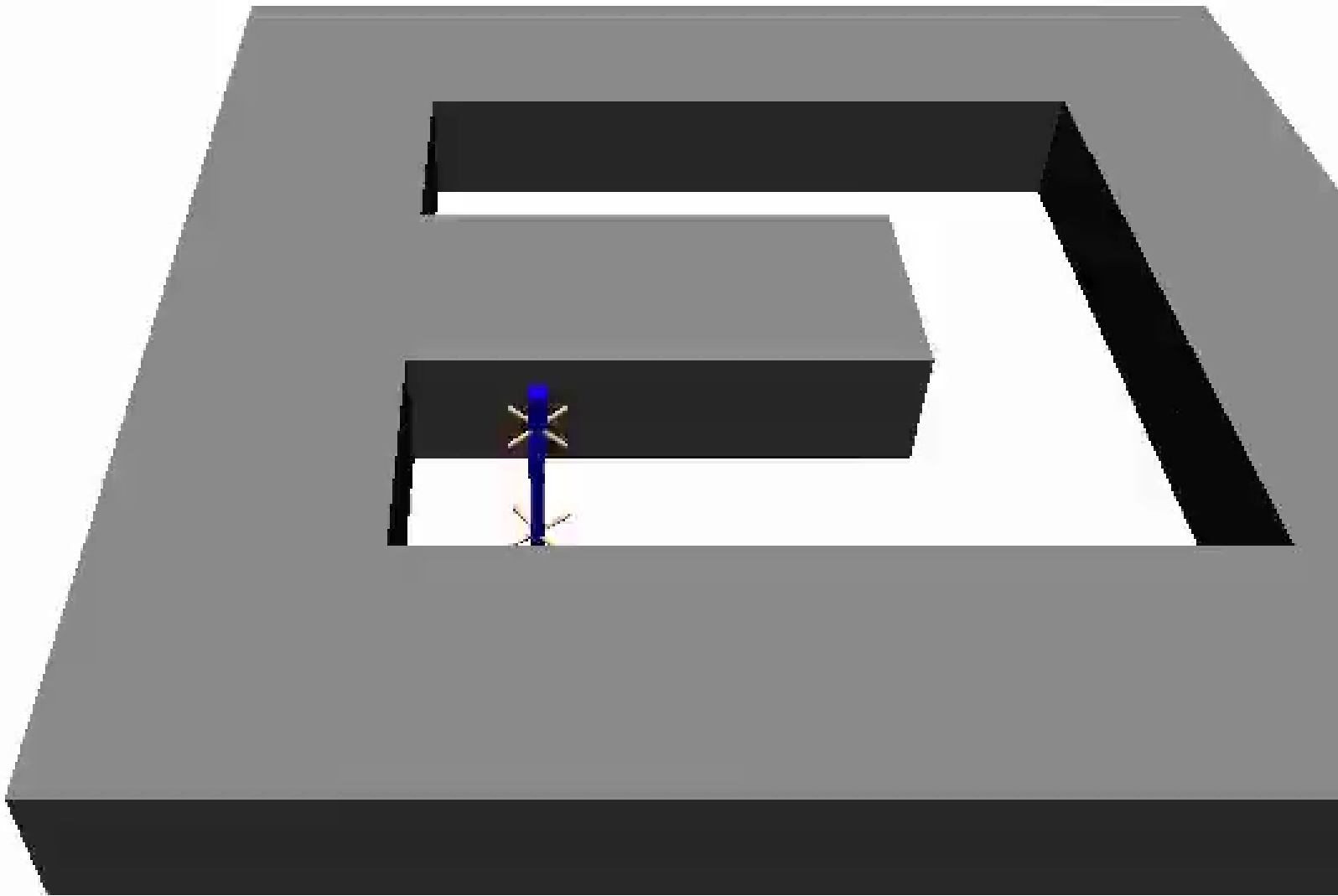


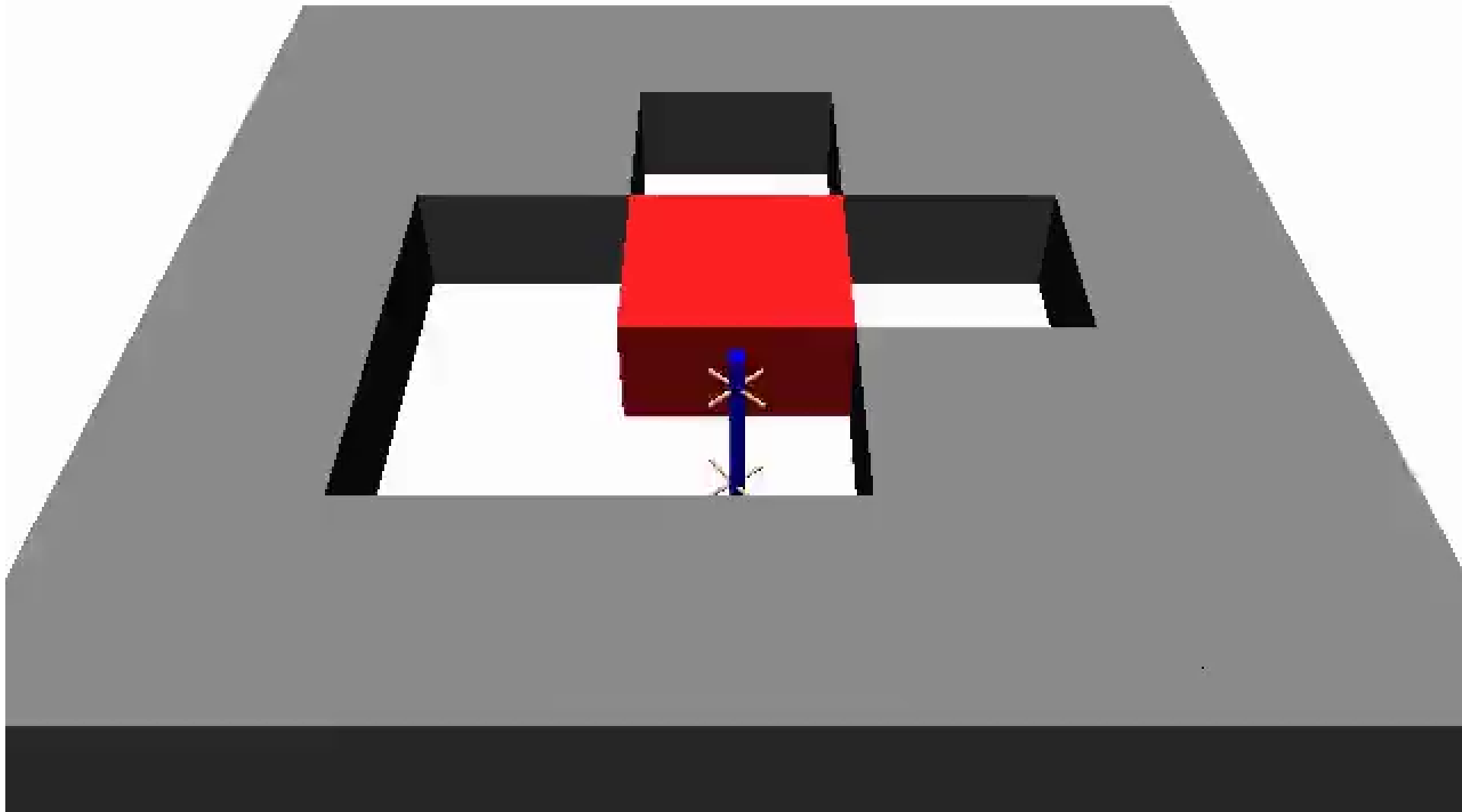
Ant Maze



Ant Push



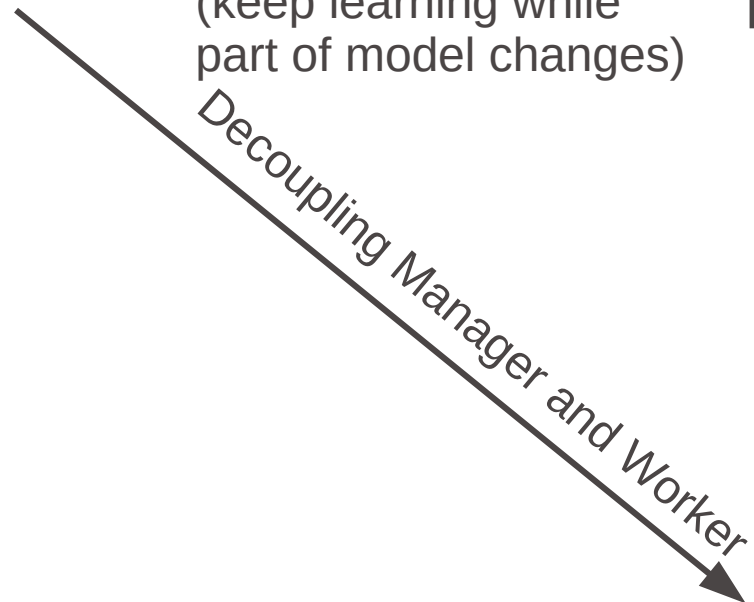
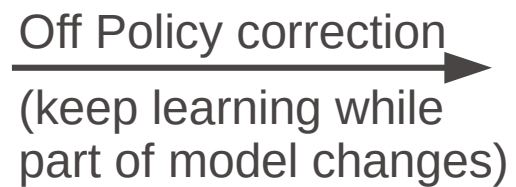




Option Framework
(first hour)



Option Critic



HIRO
Data-Efficient



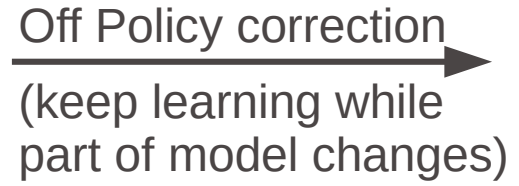
Full-state to
worker

FeUdal Net

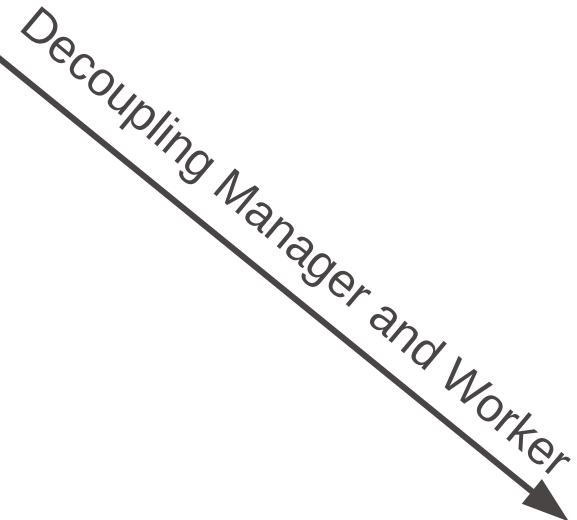
Option Framework
(first hour)



Option Critic



HIRO
Data-Efficient



HAAR

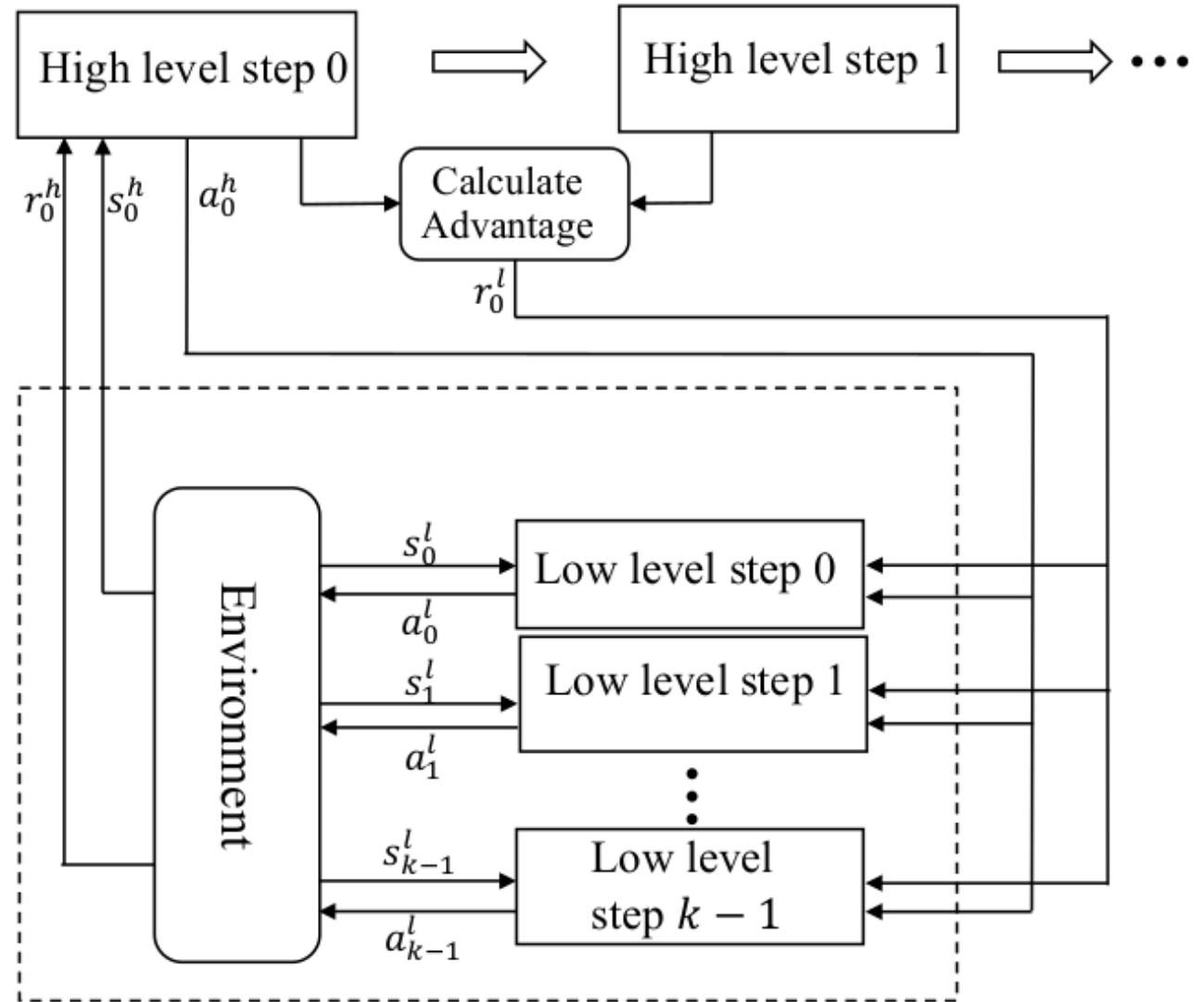


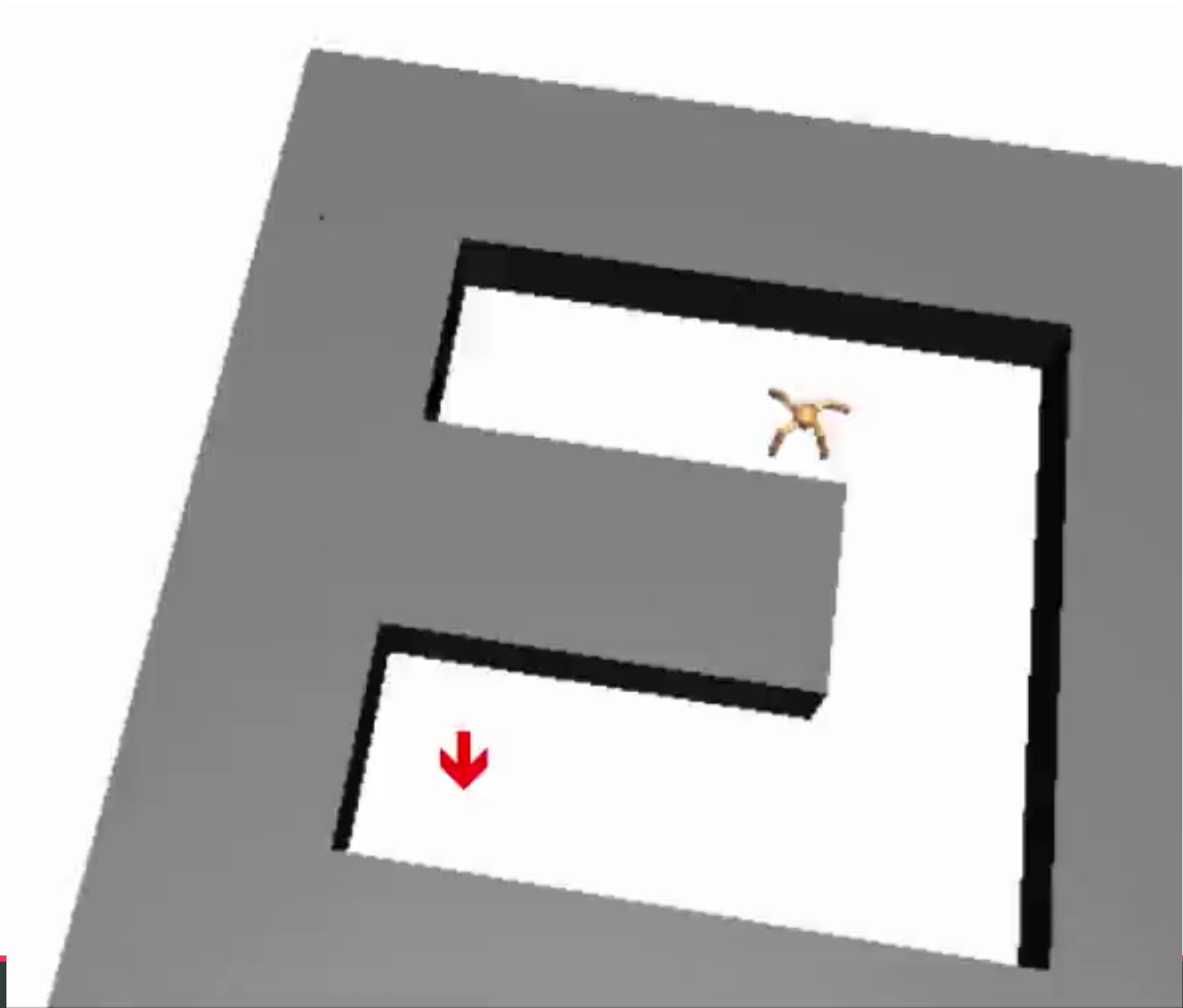
FeUdal Net



Hierarchical Reinforcement Learning with Advantage-Based Auxiliary Rewards (HAAR)

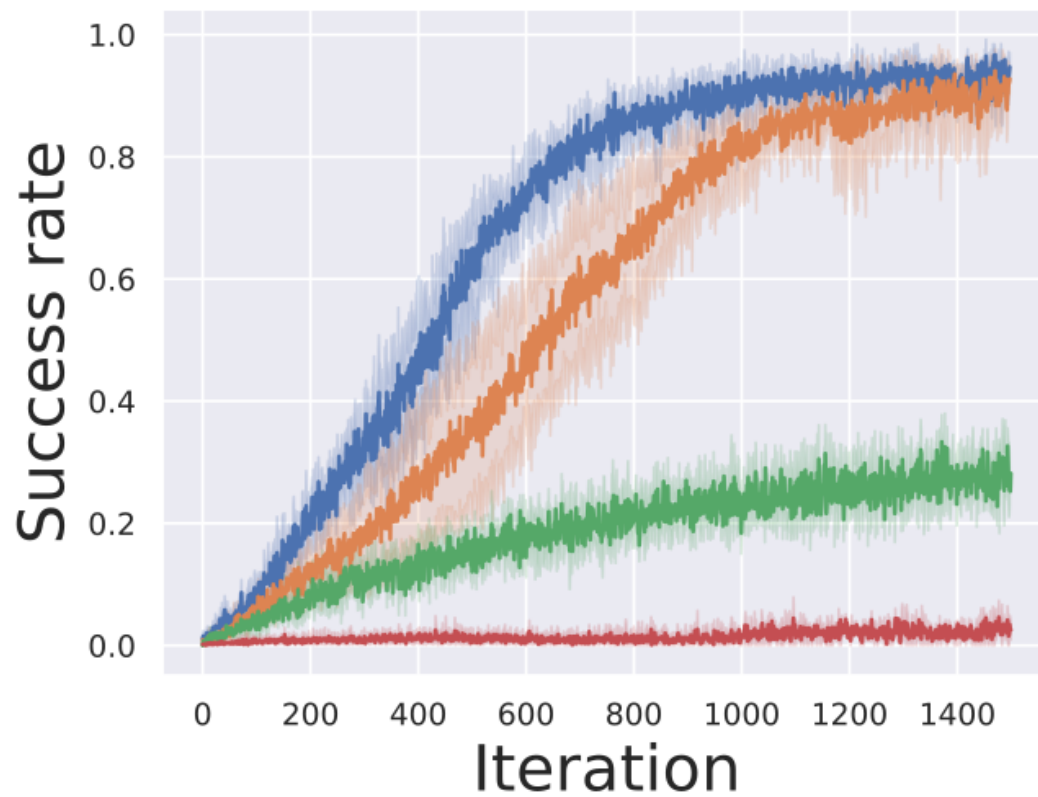
Li, Siyuan, et al. 2019



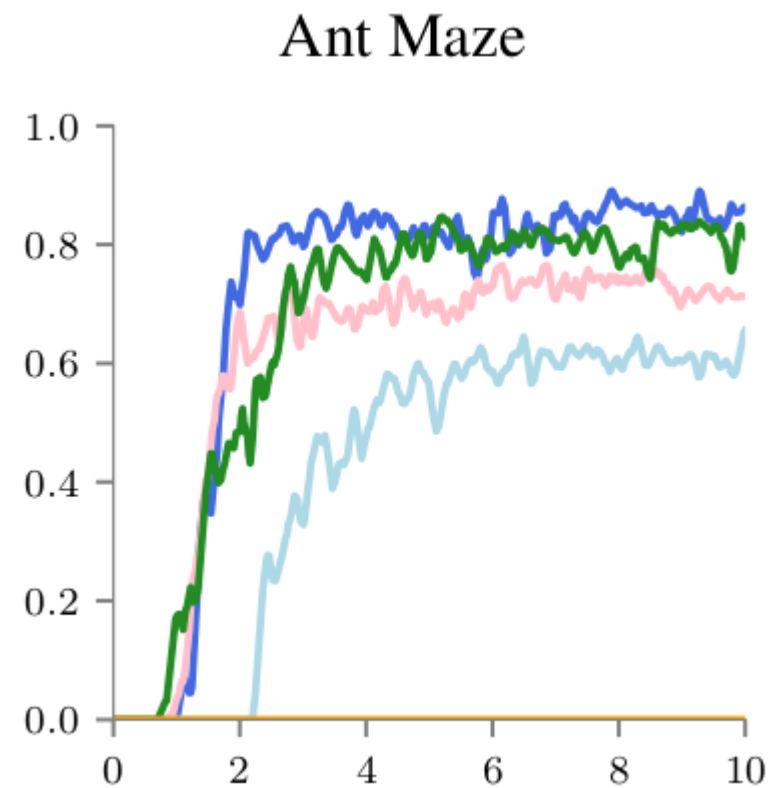


Does it help?

— HAAR w/ annealing — HAAR w/o annealing — TRPO



(a) Ant Maze



Option Framework
(first hour)

Learn options →

Option Critic

Off Policy correction
(keep learning while
part of model changes) →

HIRO
Data-Efficient

Fixed length
option ↓

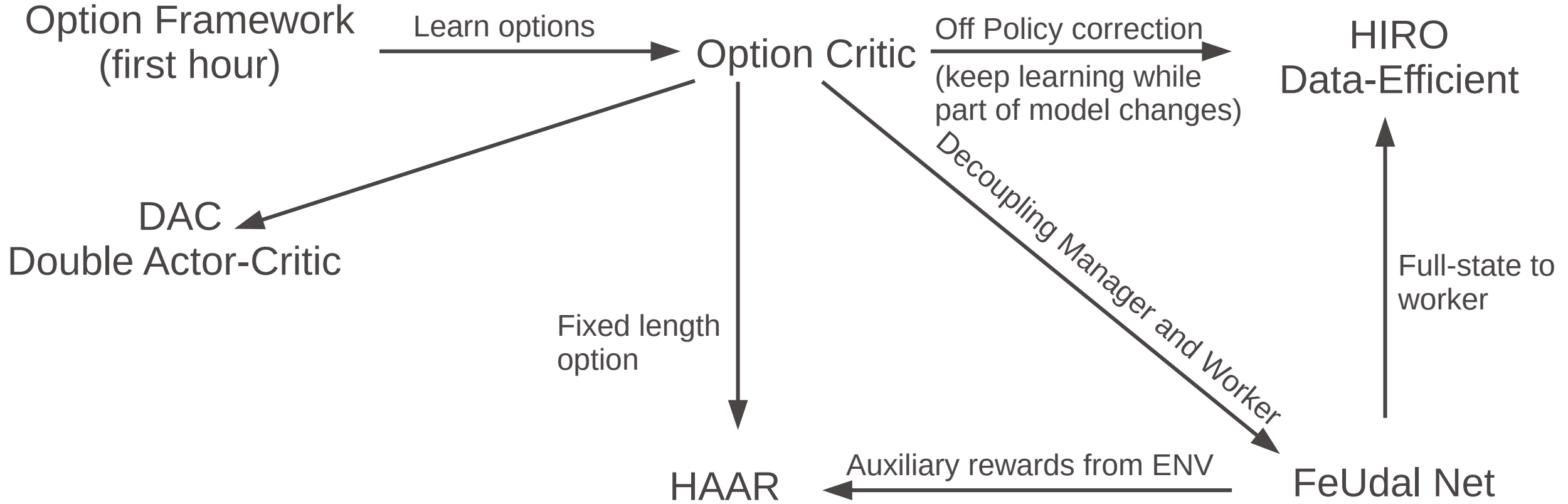
HAAR

Decoupling Manager and Worker ↘

Auxiliary rewards from ENV ←

FeUdal Net

Full-state to
worker ↑



Double Actor-Critic (DAC)

Shangdong Zhang, Shimon Whiteson 2019

Algorithm 1: Pseudocode of DAC

Input:

Parameterized $\pi, \{\pi_o, \beta_o\}_{o \in \mathcal{O}}$

Policy optimization algorithms $\mathbb{A}_1, \mathbb{A}_2$

Get an initial state S_0

$t \leftarrow 0$

while *True* do

Sample O_t from $\pi^{\mathcal{H}}(\cdot | (O_{t-1}, S_t))$

Sample A_t from $\pi^{\mathcal{L}}(\cdot | (S_t, O_t))$

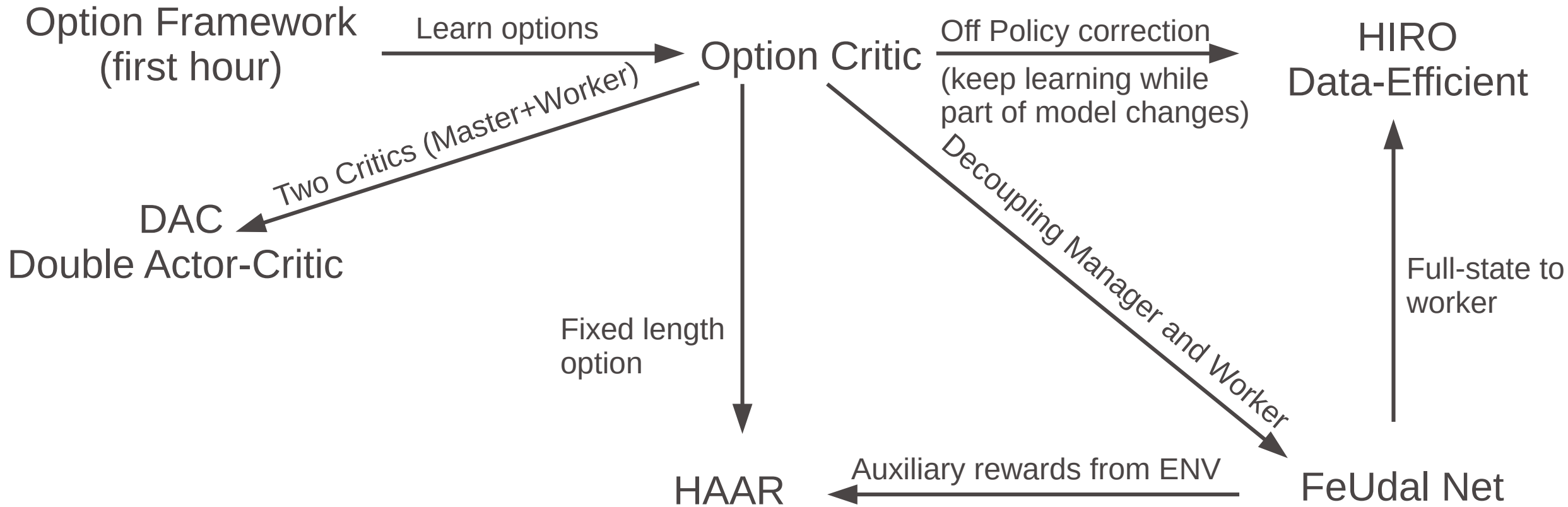
Execute A_t , get R_{t+1}, S_{t+1}

Optimize $\pi^{\mathcal{H}}$ with $(S_t^{\mathcal{H}}, A_t^{\mathcal{H}}, R_{t+1}, S_{t+1}^{\mathcal{H}})$ and \mathbb{A}_1

Optimize $\pi^{\mathcal{L}}$ with $(S_t^{\mathcal{L}}, A_t^{\mathcal{L}}, R_{t+1}, S_{t+1}^{\mathcal{L}})$ and \mathbb{A}_2

$t \leftarrow t + 1$

end



Thank You

Advantage Function

- Do we really need to compute both