

Shared Experience Actor-Critic for Multi-Agent Reinforcement Learning

Filippos Christianos, Lukas Schäfer, Stefano V. Albrecht



In 34th Conference on Neural Information Processing Systems 2020.

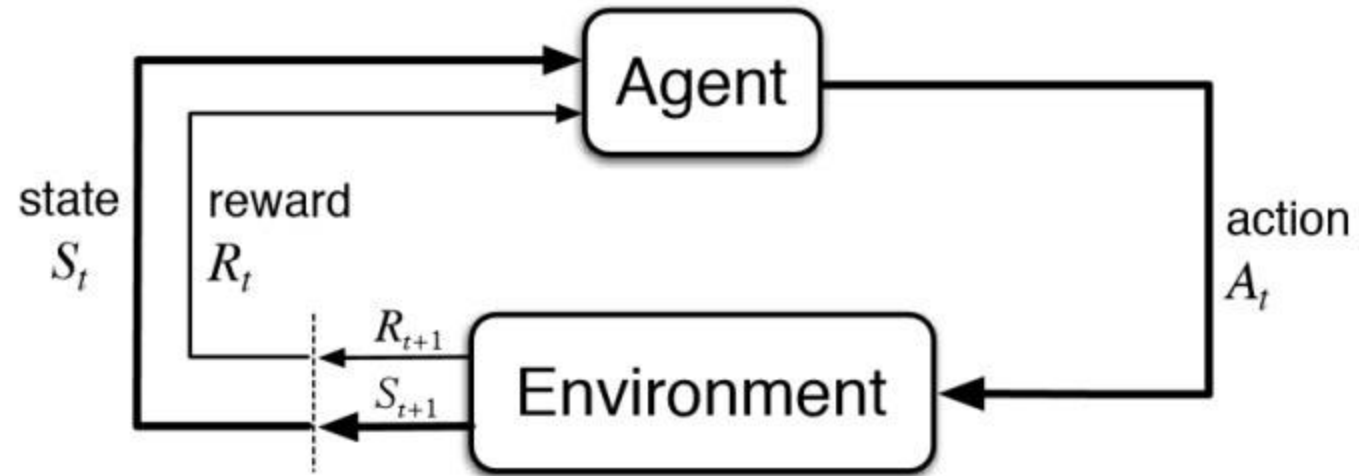


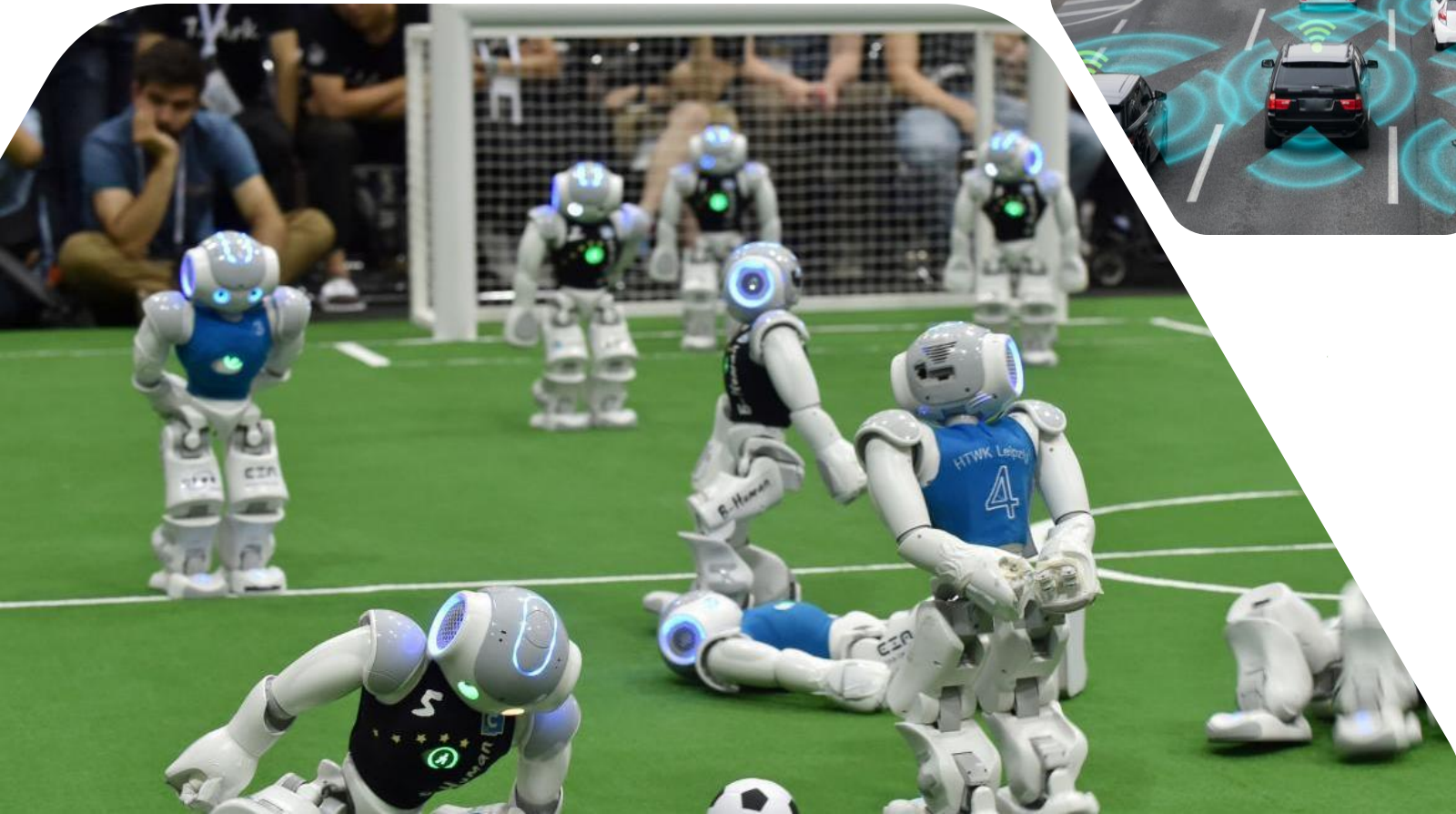
THE UNIVERSITY of EDINBURGH
informatics



AUTONOMOUS AGENTS
RESEARCH GROUP

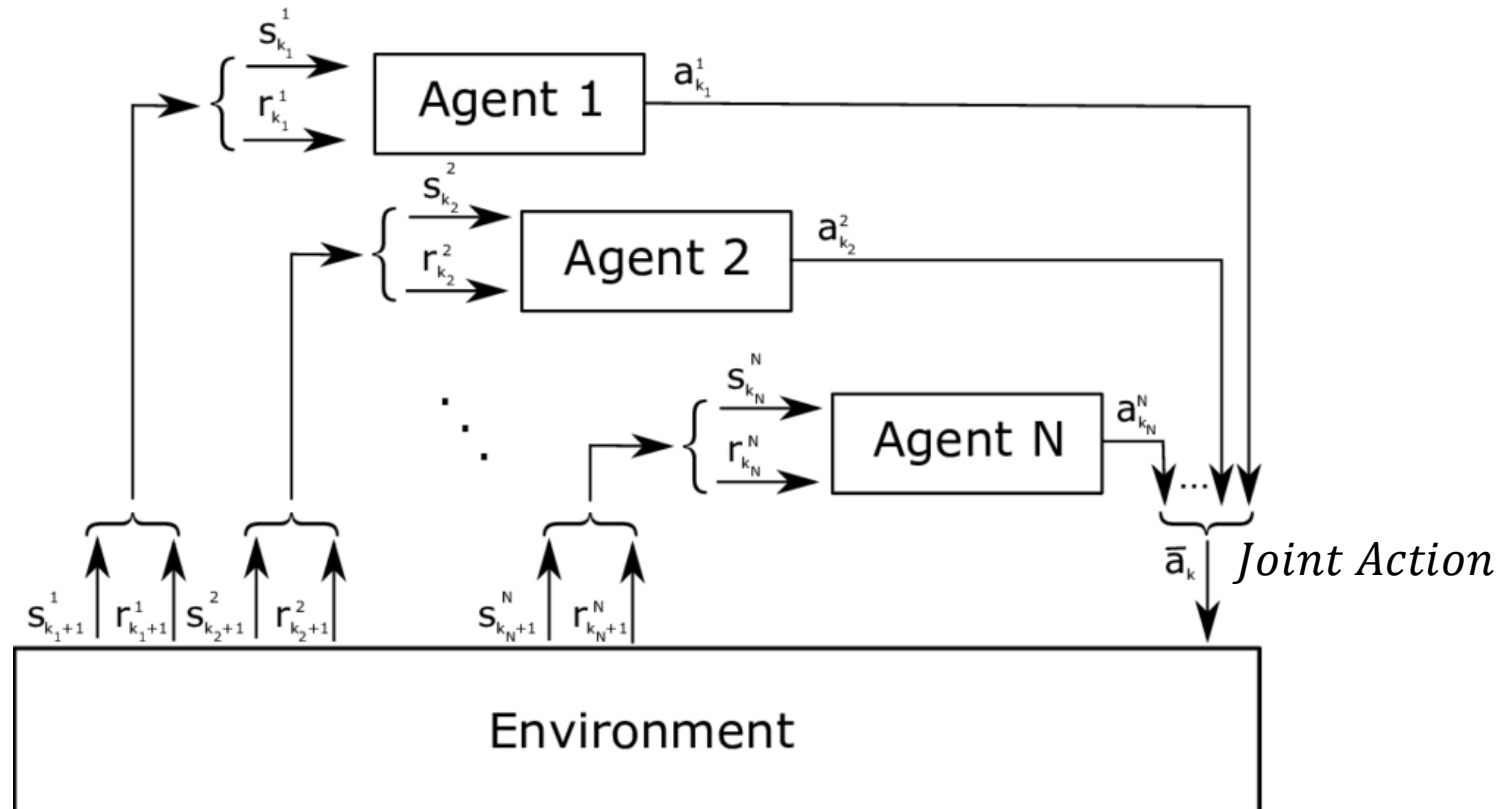
Reinforcement Learning



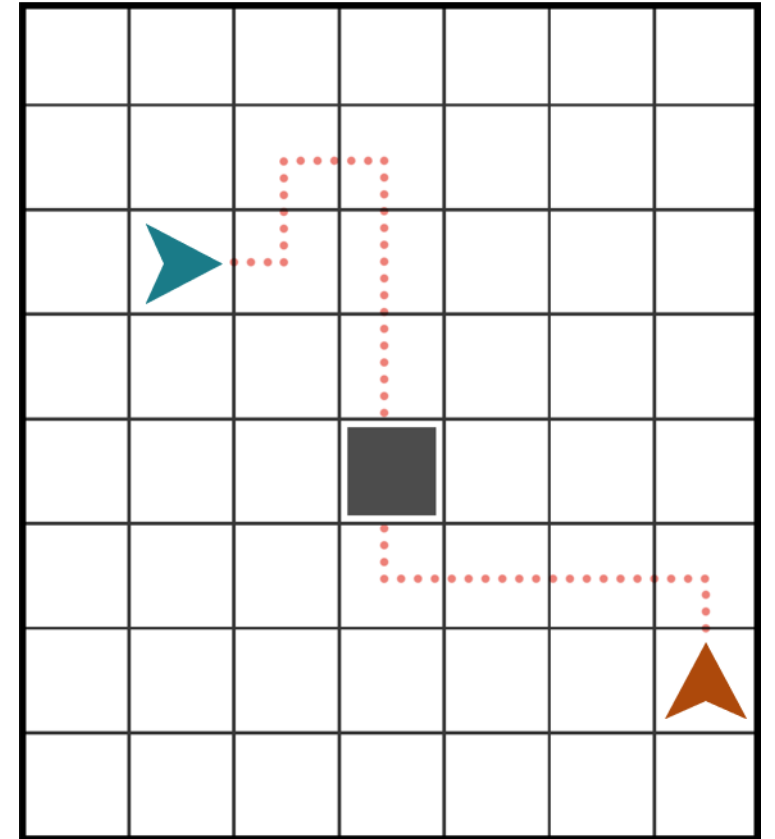


Multi-Agent Systems

Multi-Agent Reinforcement Learning

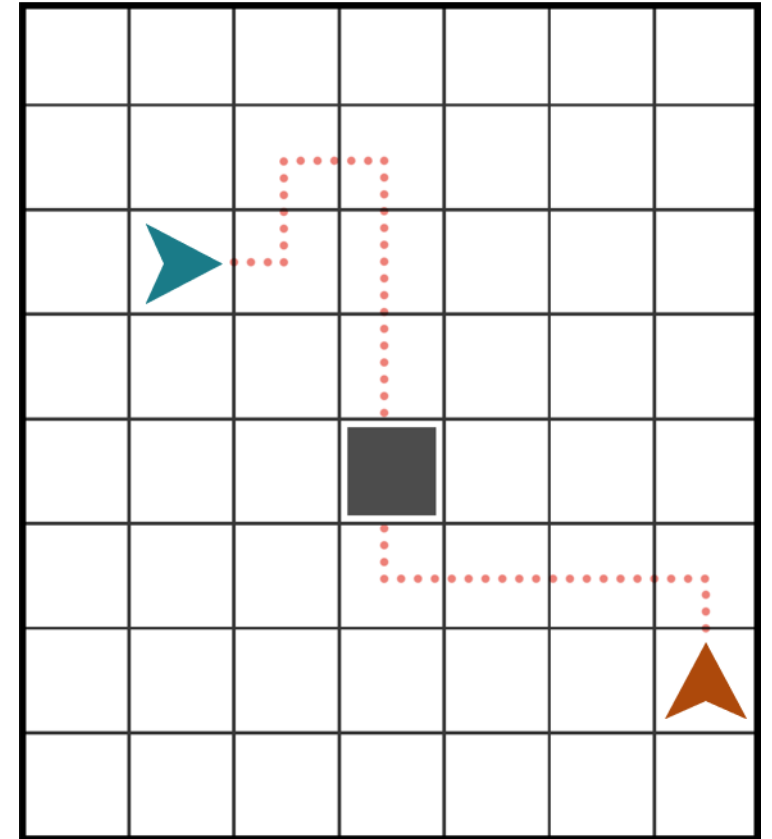


Motivational Example



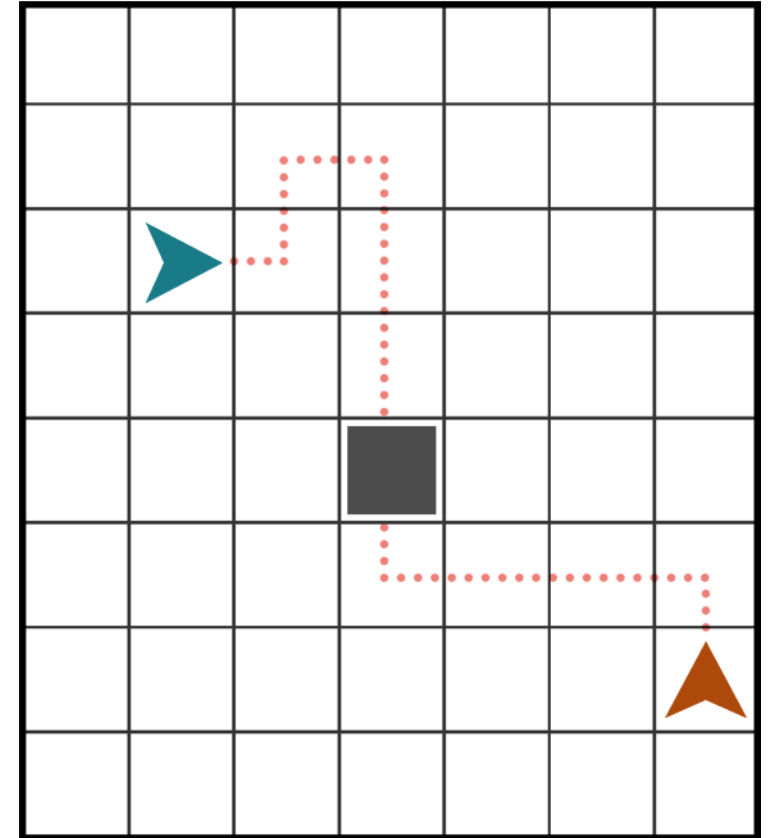
Motivational Example

- Both agents must reach goal simultaneously
 - Sparse reward signal

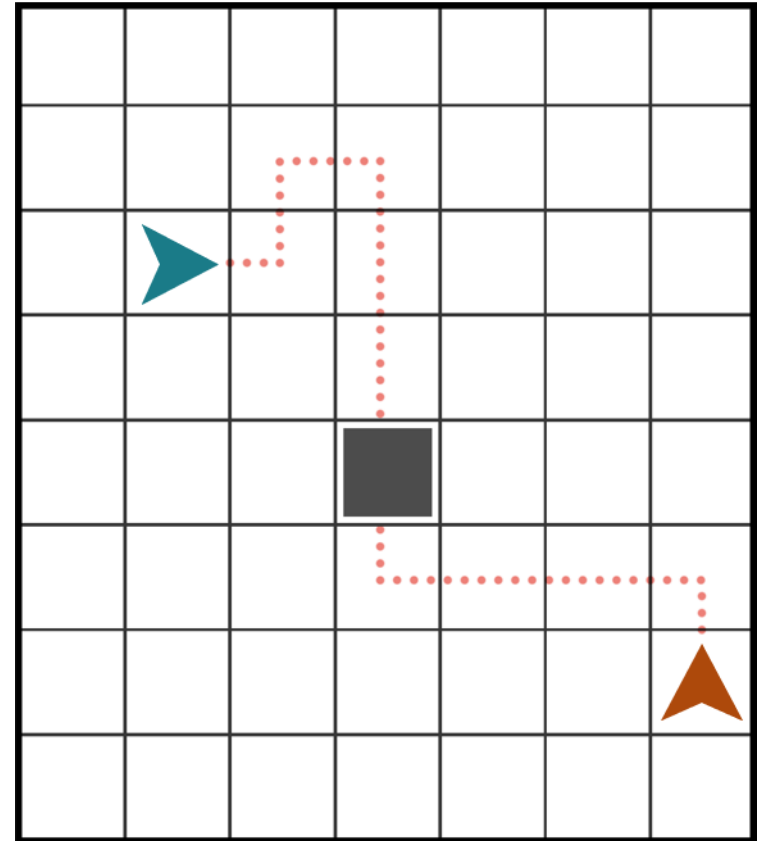


Motivational Example

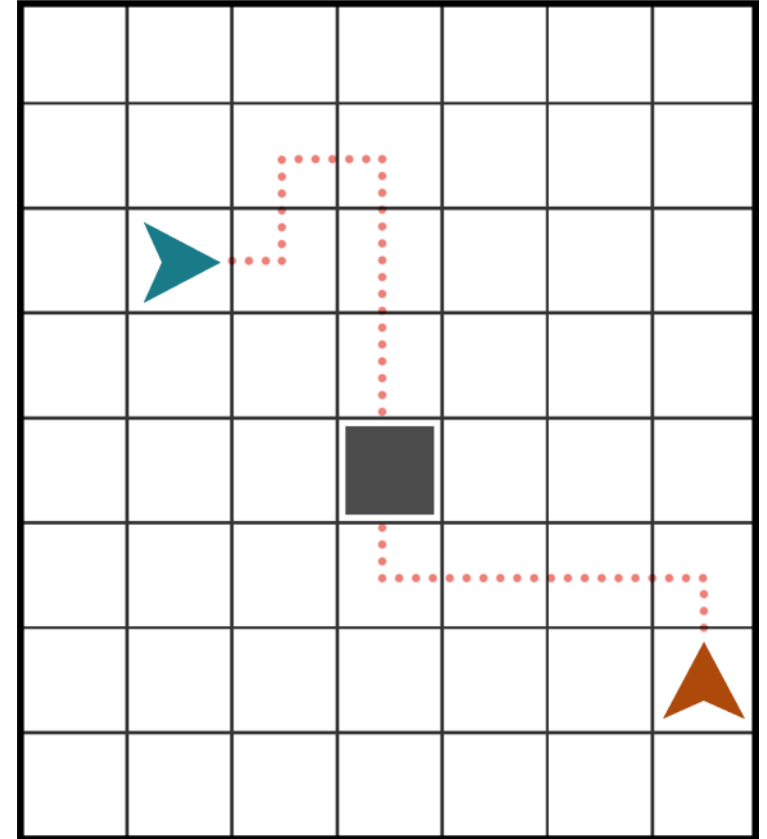
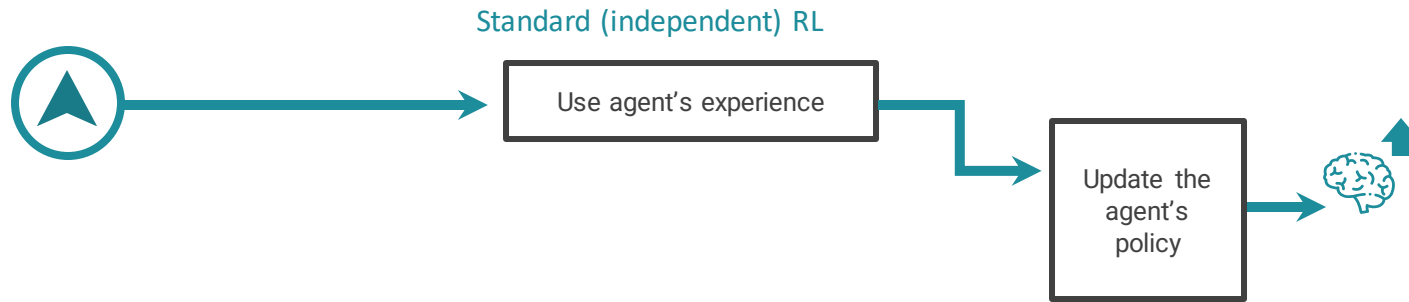
- Both agents must reach goal simultaneously
 - Sparse reward signal
- **Idea:** Make use of both agents' exploration
 - Share experience of agents



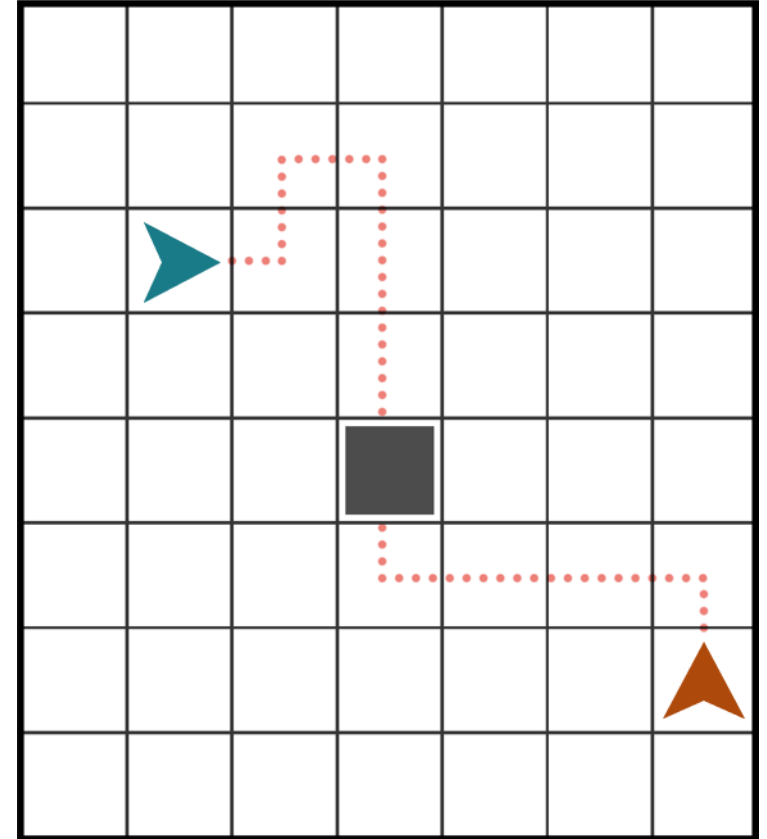
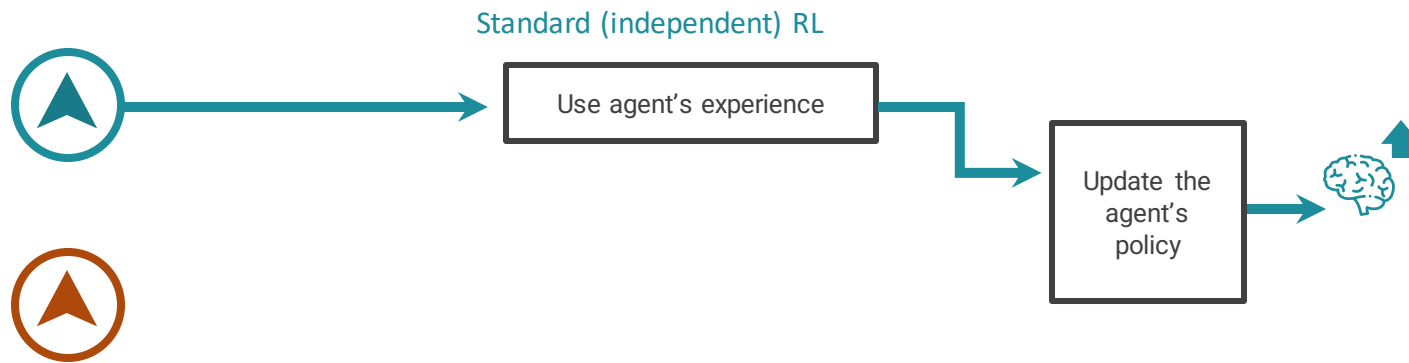
SHARED EXPERIENCE ACTOR-CRITIC



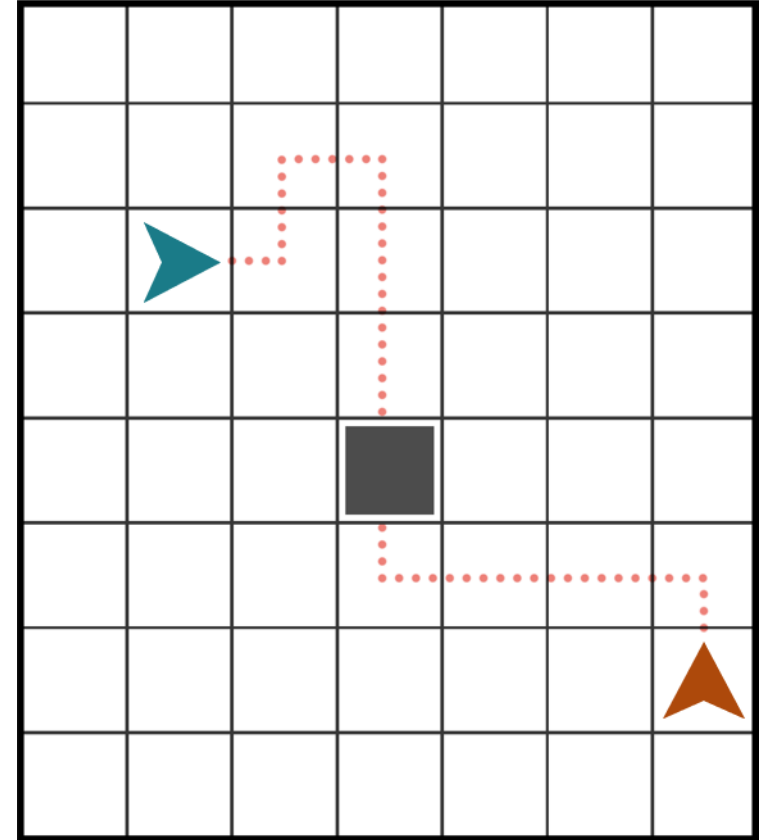
SHARED EXPERIENCE ACTOR-CRITIC



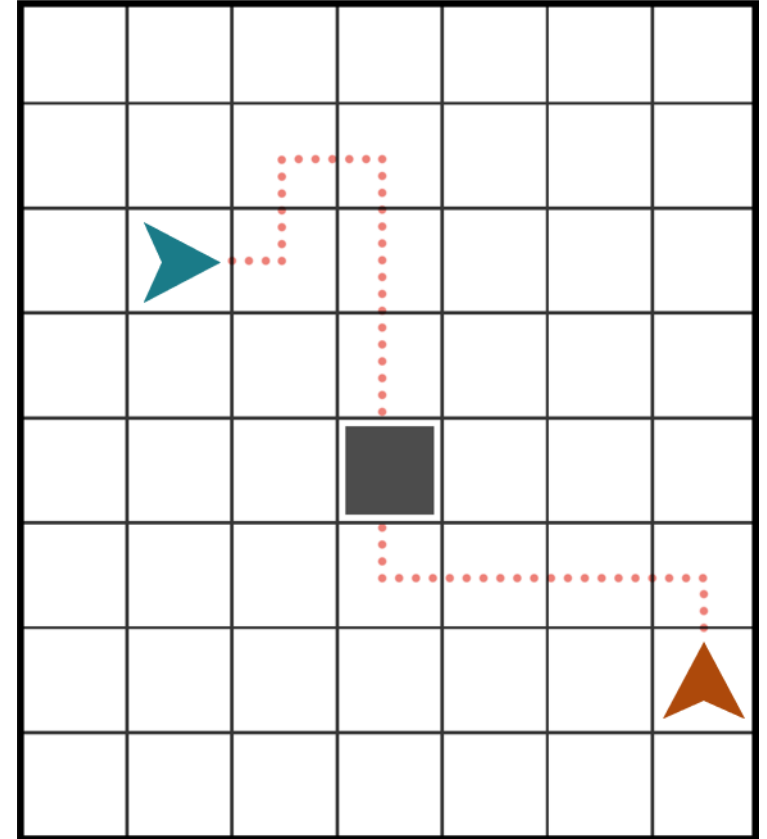
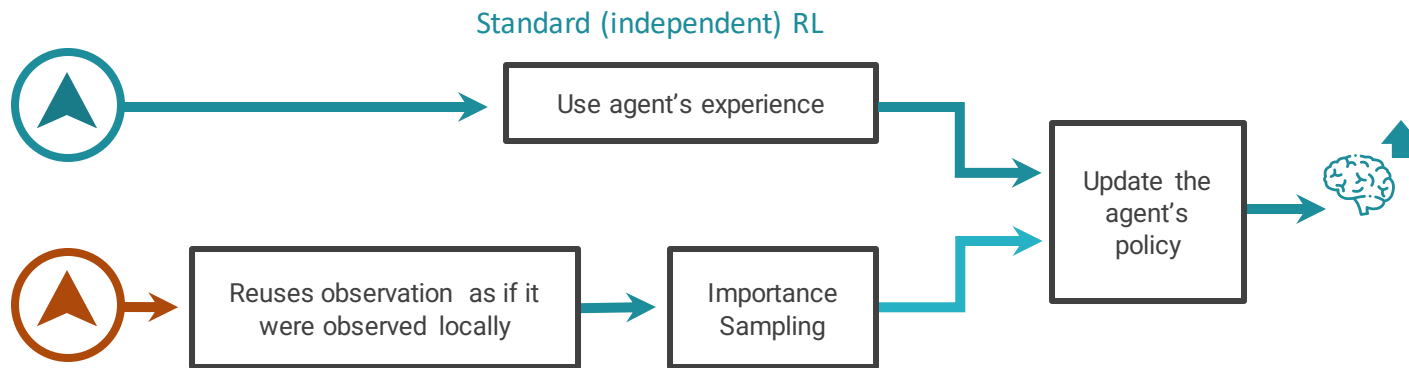
SHARED EXPERIENCE ACTOR-CRITIC



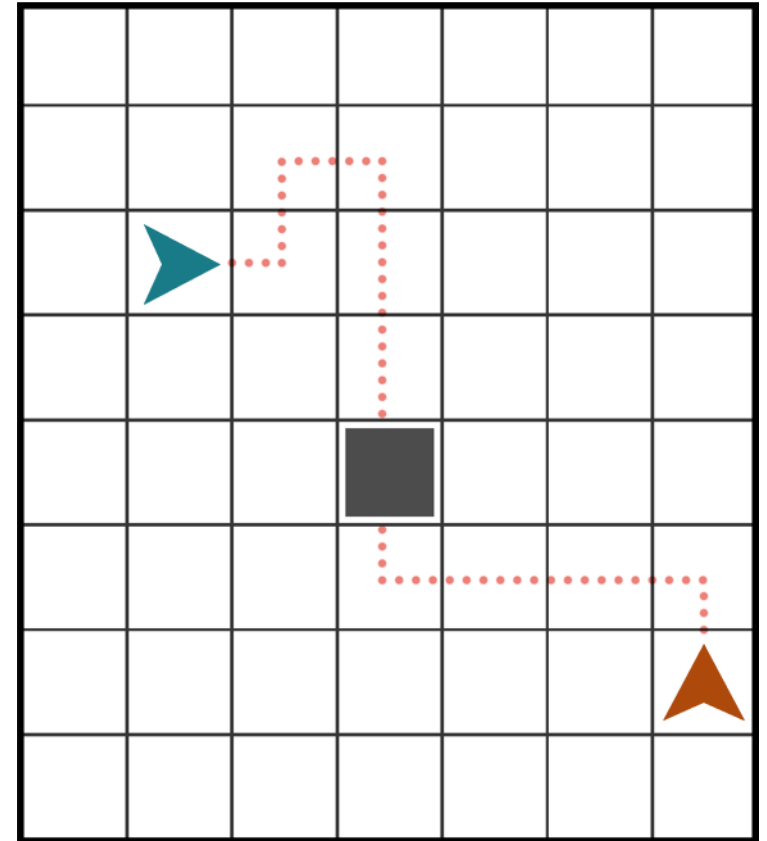
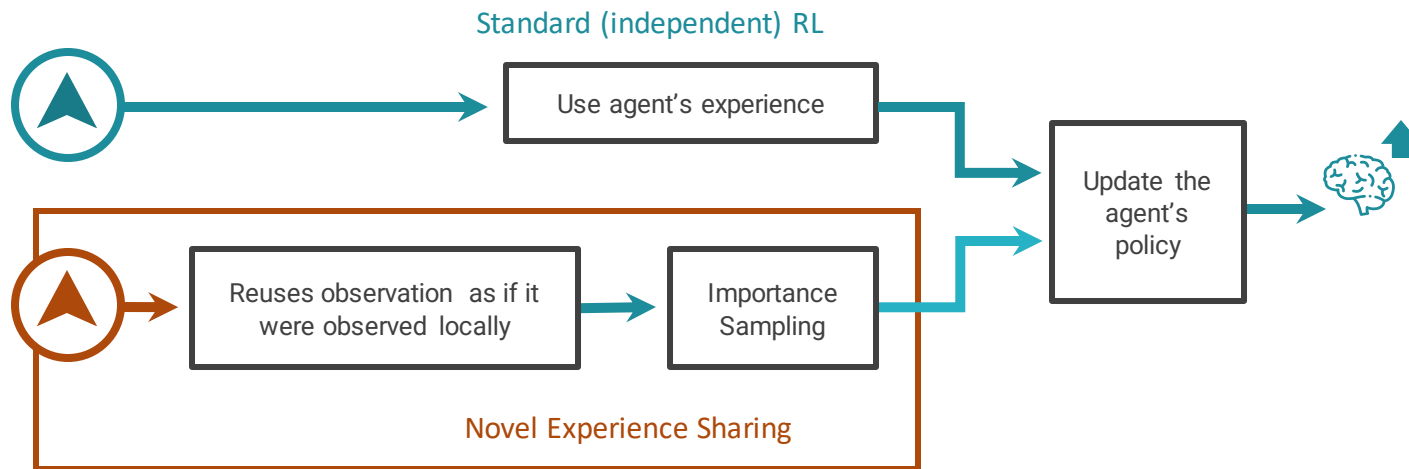
SHARED EXPERIENCE ACTOR-CRITIC



SHARED EXPERIENCE ACTOR-CRITIC

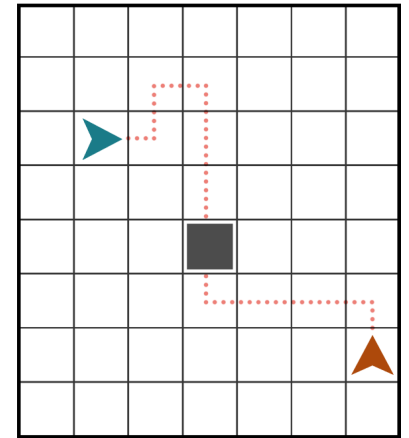


SHARED EXPERIENCE ACTOR-CRITIC



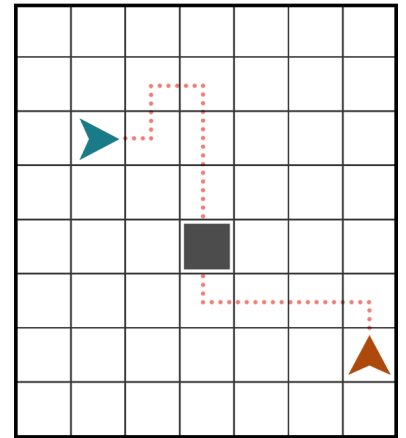
SHARED EXPERIENCE ACTOR-CRITIC

$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) (r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))$$



SHARED EXPERIENCE ACTOR-CRITIC

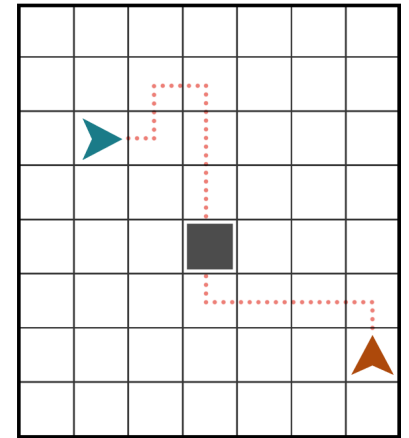
$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) \underbrace{(r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))}_{\text{Estimate of Returns}}$$



SHARED EXPERIENCE ACTOR-CRITIC

$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) (r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))$$

↑
Baseline

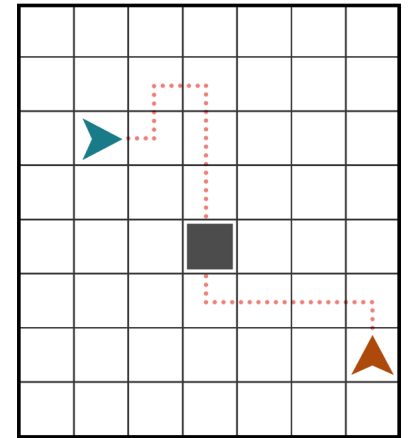


SHARED EXPERIENCE ACTOR-CRITIC

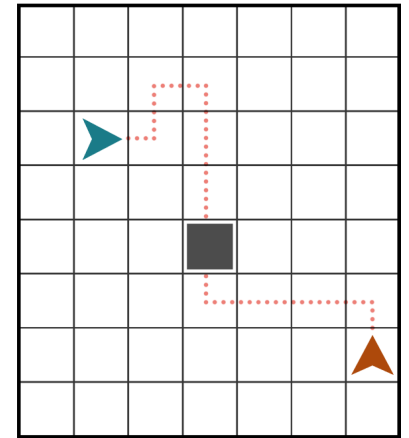
$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) \underbrace{(r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))}_{\text{Advantage: high/low}}$$

↑
Action

Advantage: high/low



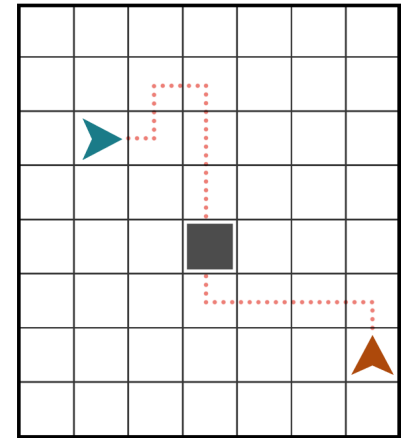
SHARED EXPERIENCE ACTOR-CRITIC



$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) (r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))$$

$$- \lambda \sum_{k \neq i} \log \pi(a_t^k | o_t^k; \phi_i) (r_t^k + \gamma V(o_{t+1}^k; \theta_i) - V(o_t^k; \theta_i))$$

SHARED EXPERIENCE ACTOR-CRITIC



$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) (r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))$$

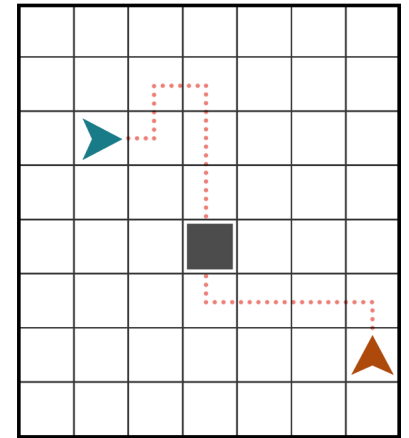
$$- \lambda \sum_{k \neq i} \frac{\pi(a_t^k | o_t^k; \phi_i)}{\pi(a_t^k | o_t^k; \phi_k)} \log \pi(a_t^k | o_t^k; \phi_i) (r_t^k + \gamma V(o_{t+1}^k; \theta_i) - V(o_t^k; \theta_i))$$

SHARED EXPERIENCE ACTOR-CRITIC

Policy Gradient Actor Loss:

$$\mathcal{L}(\phi_i) = -\log \pi(a_t^i | o_t^i; \phi_i) (r_t^i + \gamma V(o_{t+1}^i; \theta_i) - V(o_t^i; \theta_i))$$

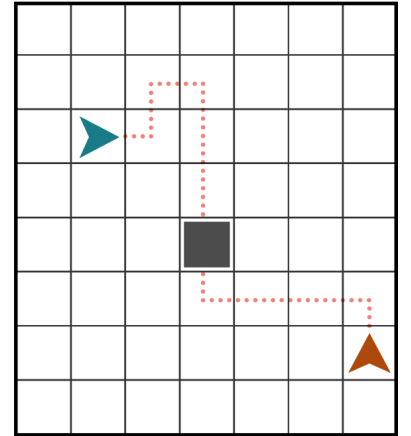
$$- \lambda \sum_{k \neq i} \frac{\pi(a_t^k | o_t^k; \phi_i)}{\pi(a_t^k | o_t^k; \phi_k)} \log \pi(a_t^k | o_t^k; \phi_i) (r_t^k + \gamma V(o_{t+1}^k; \theta_i) - V(o_t^k; \theta_i))$$



SHARED EXPERIENCE ACTOR-CRITIC

$$\mathcal{L}(\theta_i) = \left\| V(o_t^i; \theta_i) - y_i^i \right\|^2$$

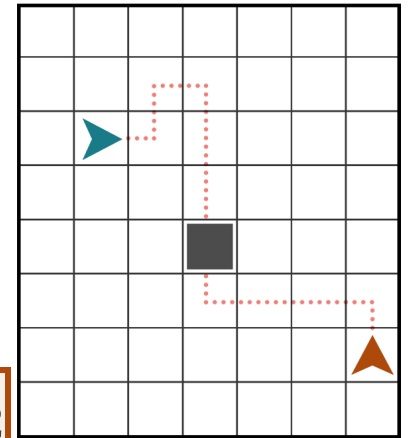
$$y_k^i = r_t^k + \gamma V(o_{t+1}^k; \theta_i)$$



SHARED EXPERIENCE ACTOR-CRITIC

$$\mathcal{L}(\theta_i) = \left\| V(o_t^i; \theta_i) - y_i^i \right\|^2 + \lambda \sum_{k \neq i} \left\| V(o_t^k; \theta_i) - y_k^i \right\|^2$$

$$y_k^i = r_t^k + \gamma V(o_{t+1}^k; \theta_i)$$

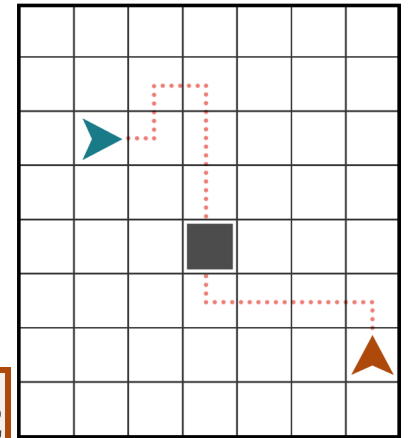


SHARED EXPERIENCE ACTOR-CRITIC

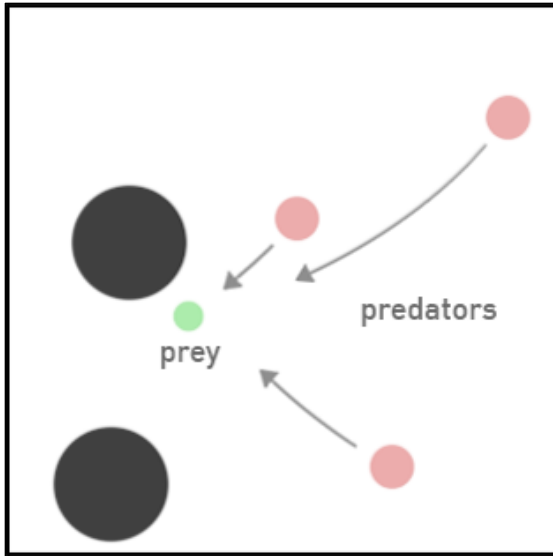
Value Function Critic Loss:

$$\mathcal{L}(\theta_i) = \left\| V(o_t^i; \theta_i) - y_i^i \right\|^2 + \lambda \sum_{k \neq i} \frac{\pi(a_t^k | o_t^k; \phi_i)}{\pi(a_t^k | o_t^k; \phi_k)} \left\| V(o_t^k; \theta_i) - y_k^i \right\|^2$$

$$y_k^i = r_t^k + \gamma V(o_{t+1}^k; \theta_i)$$



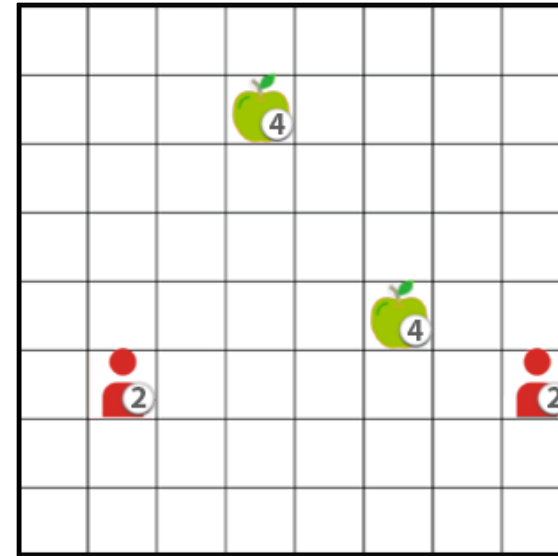
Evaluation - Domains



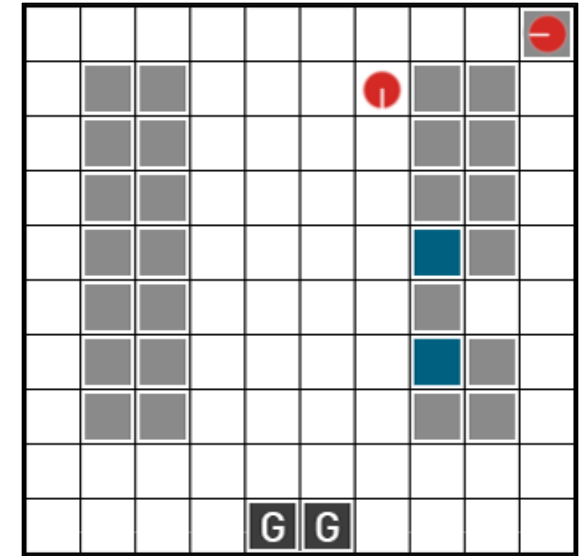
Predator Prey (sparse)



SMAC - 3m (sparse)



Level-Based Foraging (LBF)



Multi-Robot Warehouse (RWARE)

Evaluation - Baselines

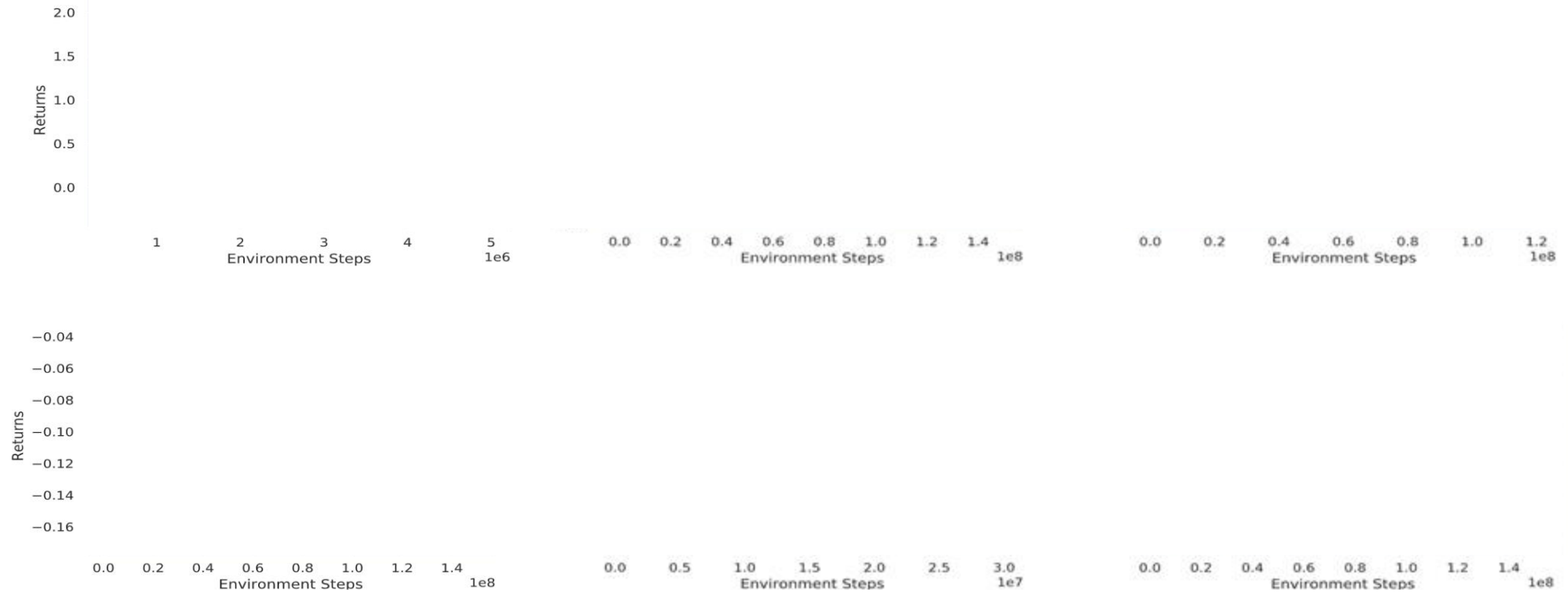
Baselines:

- (1) Independent Actor-Critic (IAC)
- (2) Shared Network Actor-Critic (SNAC)

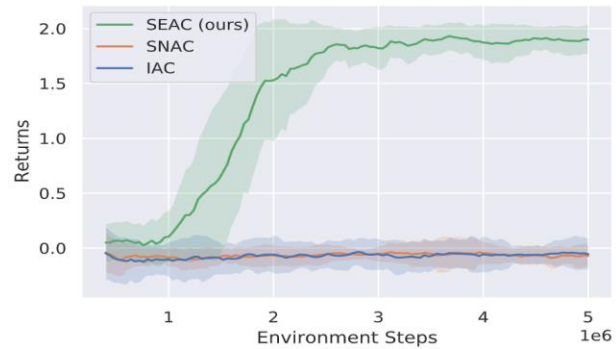
State-of-the-art MARL:

- (1) MADDPG
- (2) QMIX
- (3) ROMA

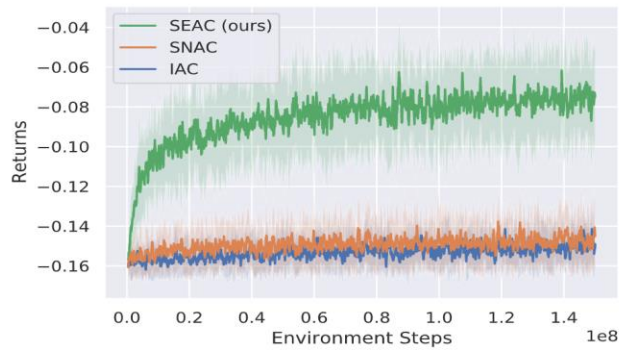
Results



Results

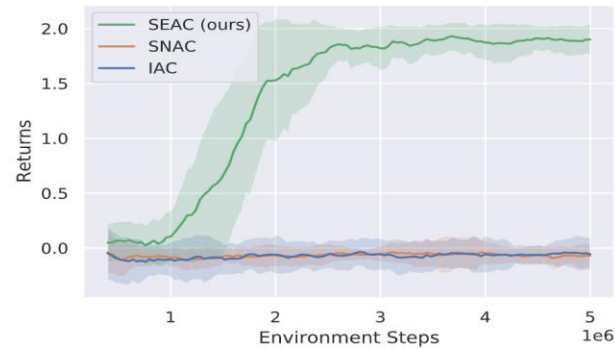


(a) PP, sparse rewards

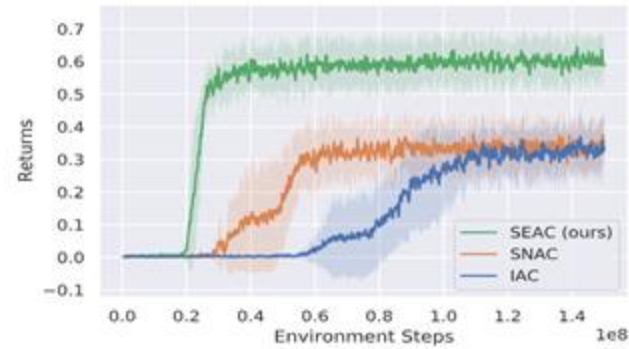


(b) SMAC with three marines, sparse rewards

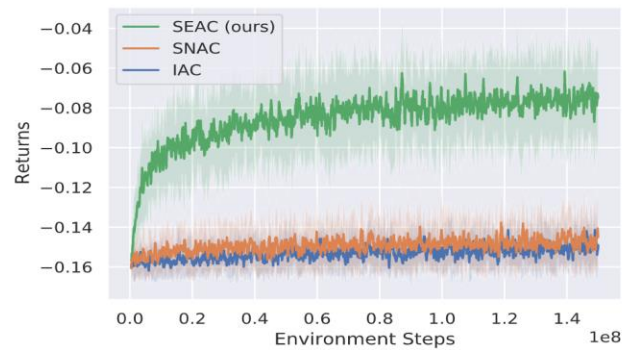
Results



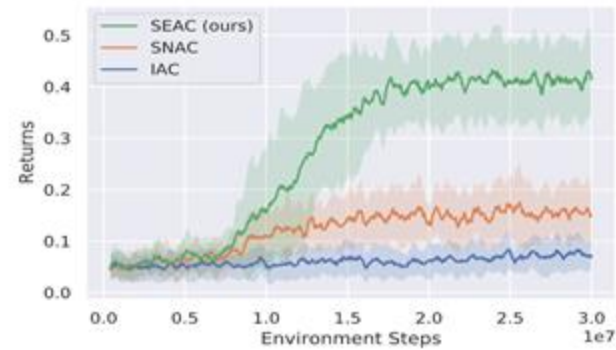
(a) PP, sparse rewards



(d) LBF: (8 x 8), two agents, two foods, cooperative



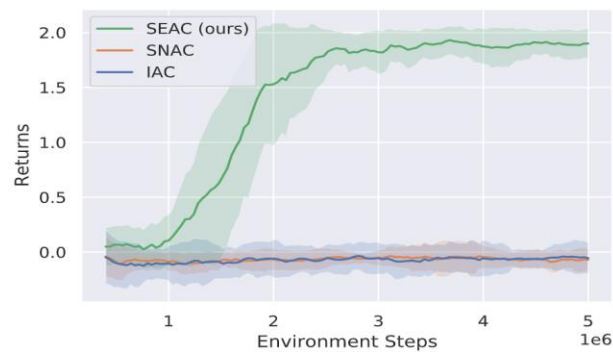
(b) SMAC with three marines, sparse rewards



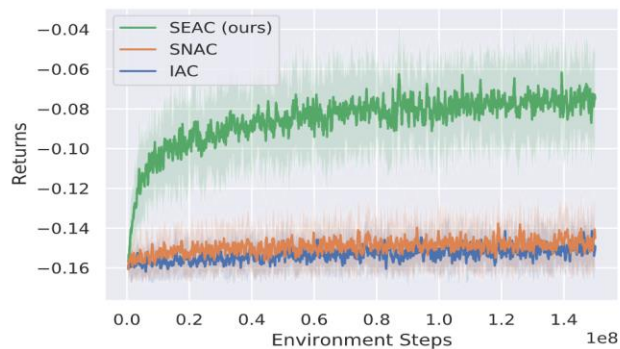
(c) LBF: (15 x 15), three agents, four food



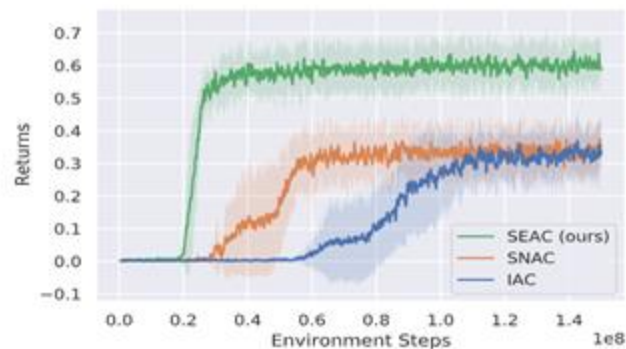
Results



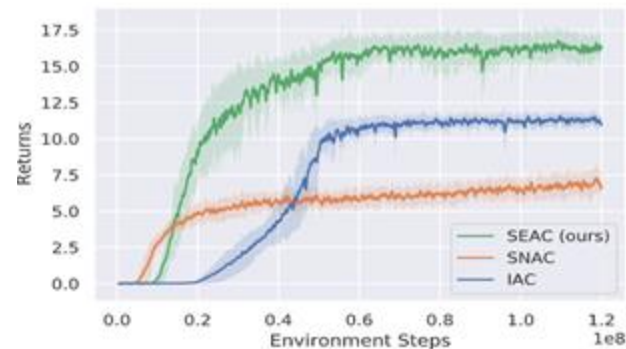
(a) PP, sparse rewards



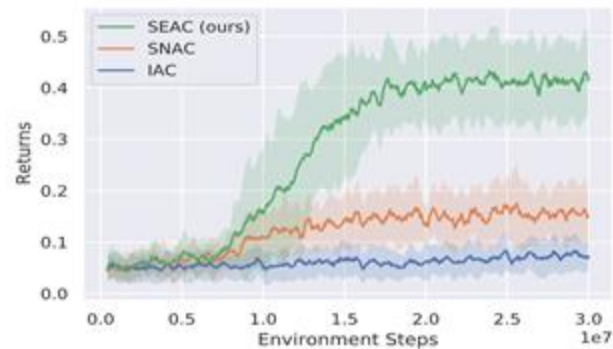
(b) SMAC with three marines, sparse rewards



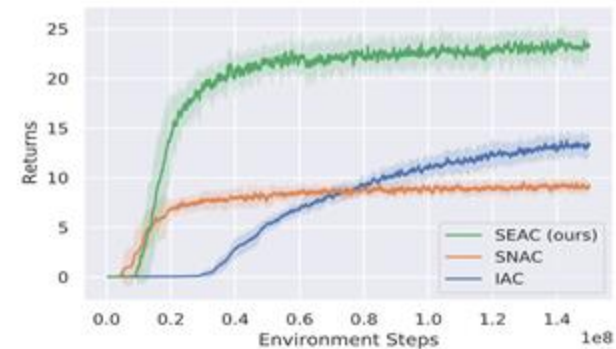
(d) LBF: (8 × 8), two agents, two foods, cooperative



(c) RWARE: (10 × 11), two agents, hard



(c) LBF: (15 × 15), three agents, four food

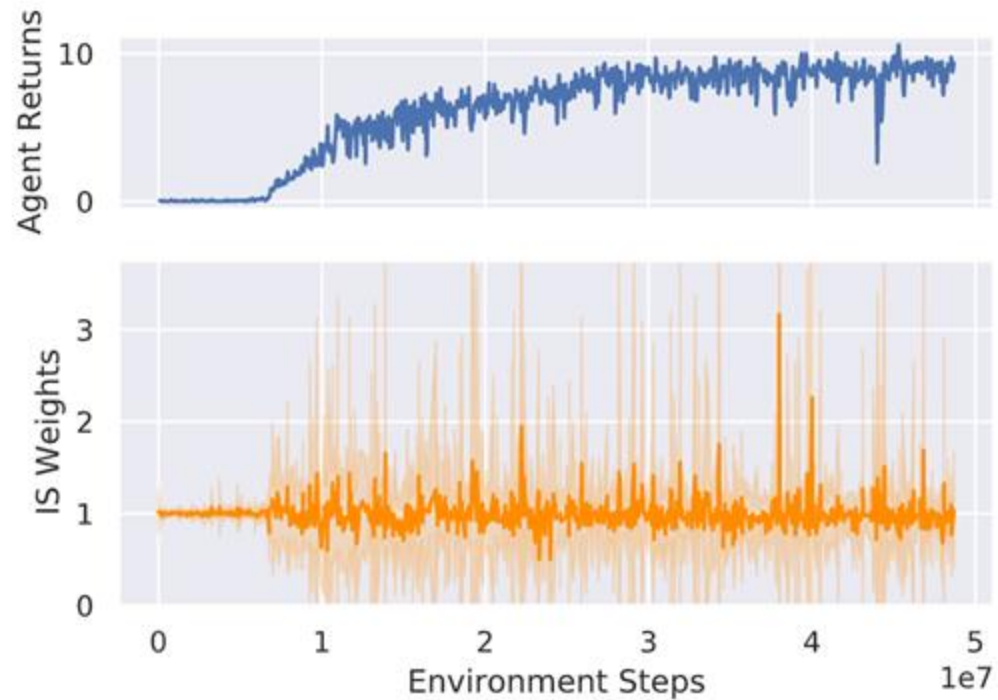


(d) RWARE: (10 × 20), four agents

Results

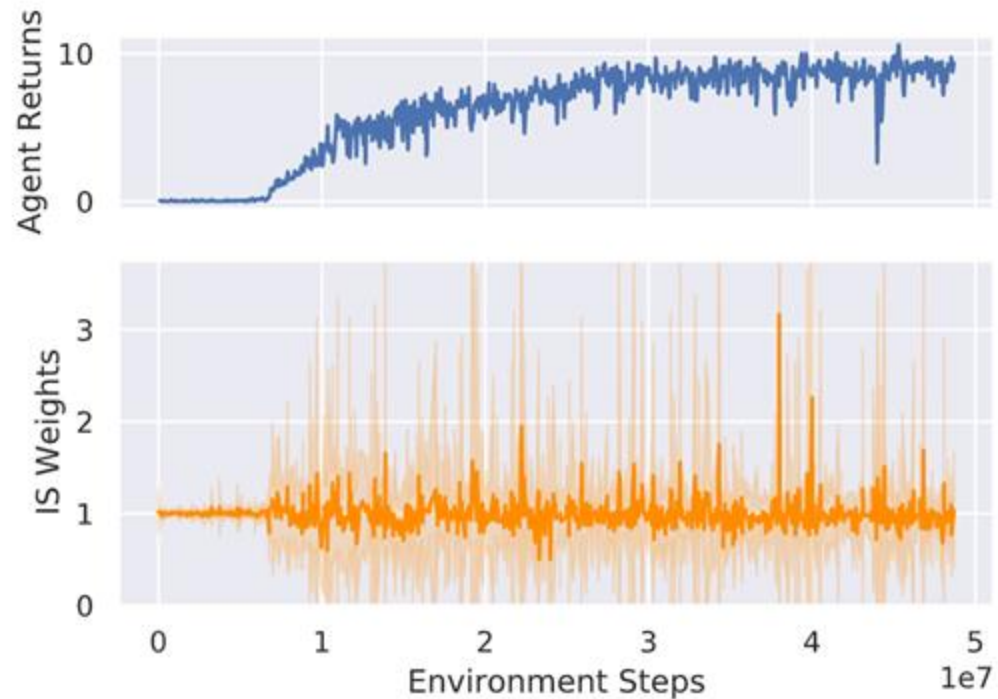
	IAC	SNAC	SEAC (ours)	QMIX	MADDPG	ROMA
PP (sparse)	-0.04 \pm 0.13	-0.04 \pm 0.1	1.93 \pm0.13	0.05 \pm 0.07	2.04 \pm0.08	0.04 \pm 0.07
SMAC-3m (sparse)	-0.13 \pm 0.01	-0.14 \pm 0.02	-0.03 \pm0.03	0.00 \pm0.00	-0.01 \pm0.01	0.00 \pm0.00
LBF-(15x15)-3ag-4f	0.13 \pm 0.04	0.18 \pm 0.08	0.43 \pm0.09	0.03 \pm 0.01	0.01 \pm 0.02	0.03 \pm 0.02
LBF-(8x8)-2ag-2f-coop	0.37 \pm 0.10	0.38 \pm 0.10	0.64 \pm0.08	0.79 \pm0.31	0.01 \pm 0.02	0.01 \pm 0.02
RWARE-(10x20)-4ag	13.75 \pm 1.26	9.53 \pm 0.83	23.96 \pm1.92	0.00 \pm 0.00	0.00 \pm 0.00	0.00 \pm 0.00
RWARE-(10x11)-4ag	40.10 \pm5.60	36.79 \pm 2.36	45.11 \pm2.90	0.00 \pm 0.00	0.00 \pm 0.00	0.01 \pm 0.01

Analysis (1)



Importance weights of one SEAC agent in RWARE,
(10x11), two agents, hard

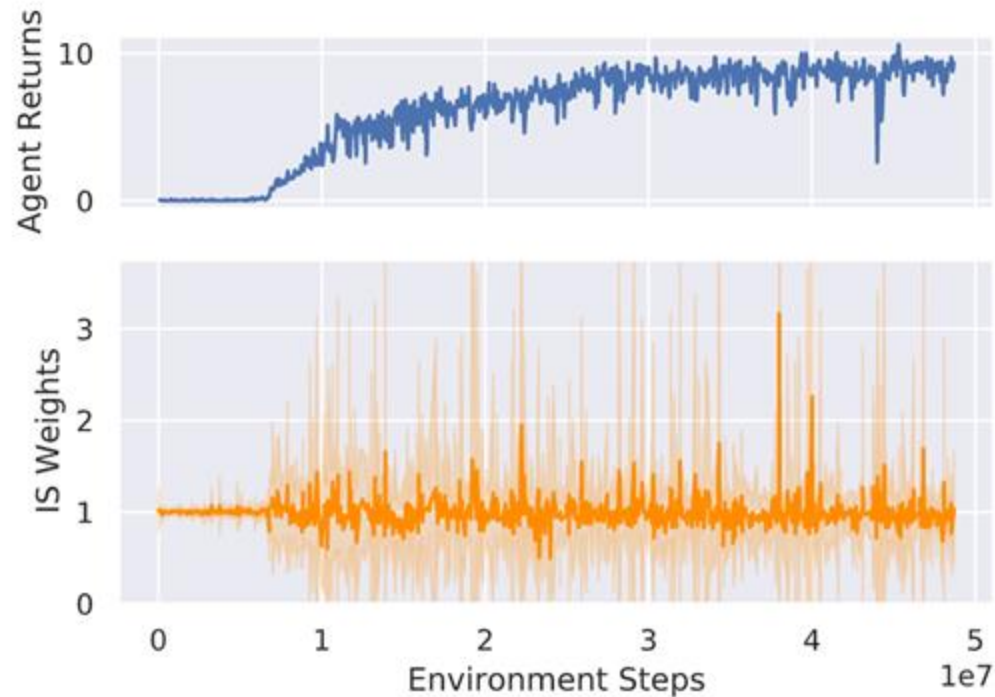
Analysis (1)



- Agents learn similar, but not identical policies which improves coordination

Importance weights of one SEAC agent in RWARE, (10x11), two agents, hard

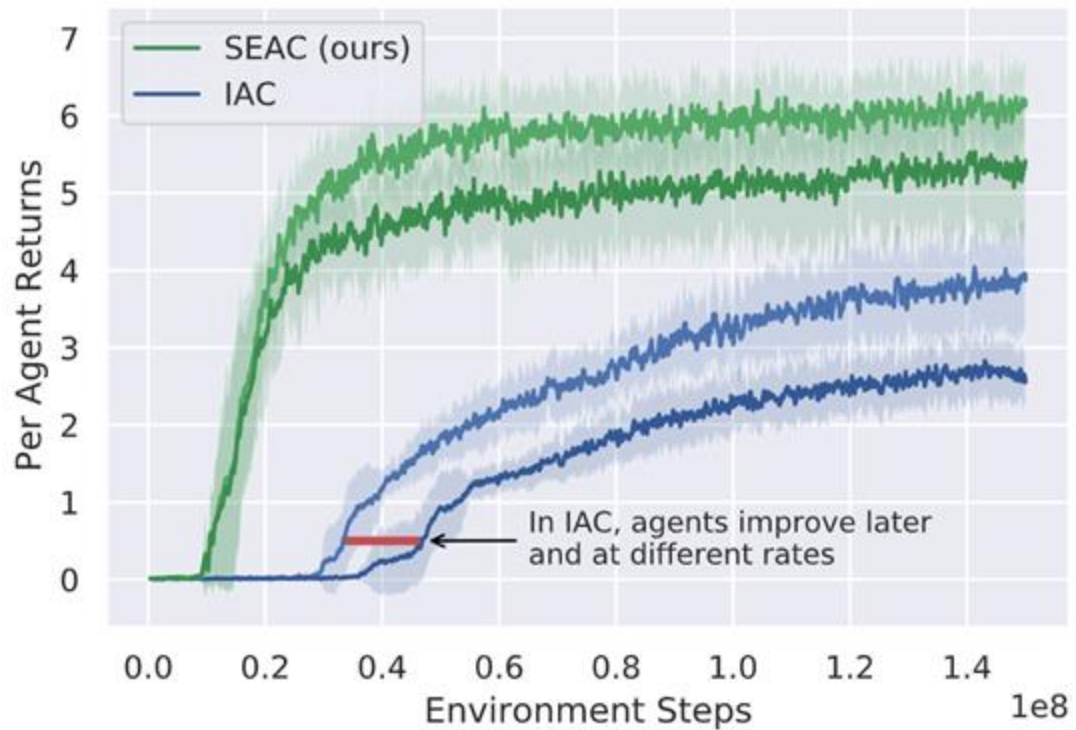
Analysis (1)



Importance weights of one SEAC agent in RWARE,
(10x11), two agents, hard

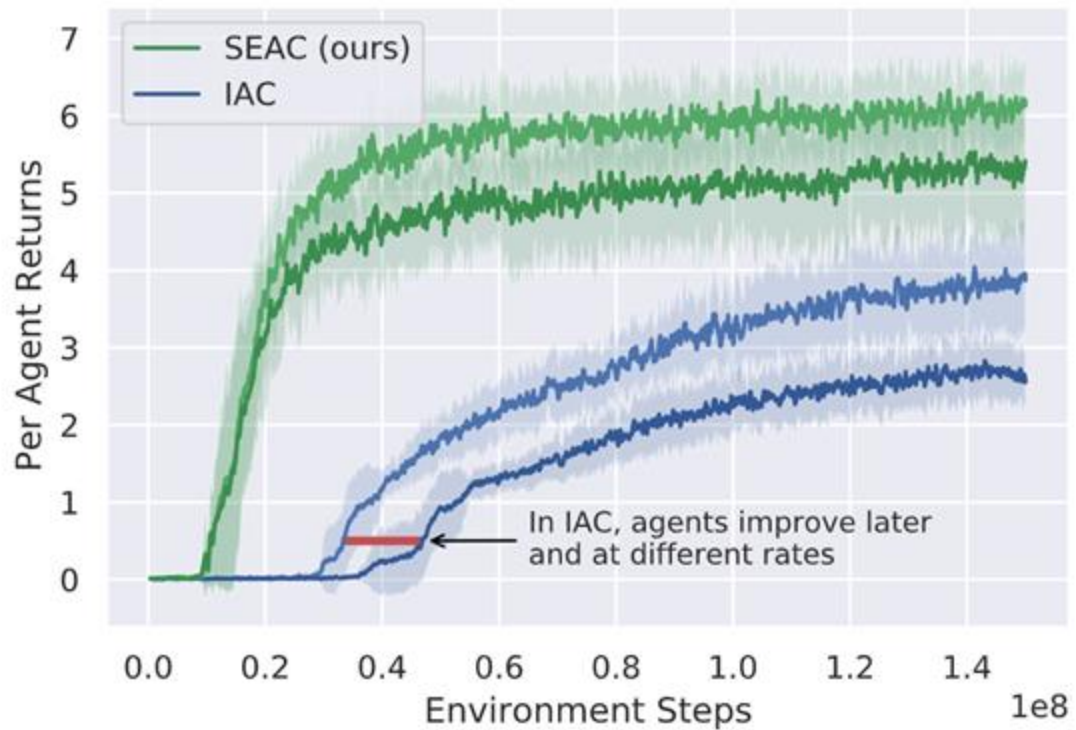
- Agents learn similar, but not identical policies which improves coordination
- Policies diverge because of ...
 1. Random network initialization
 2. Entropy regularization term in final policy loss (based on own policy)

Analysis (2)



Best vs. Worst performing agents on RWARE,
(10x20), four agents

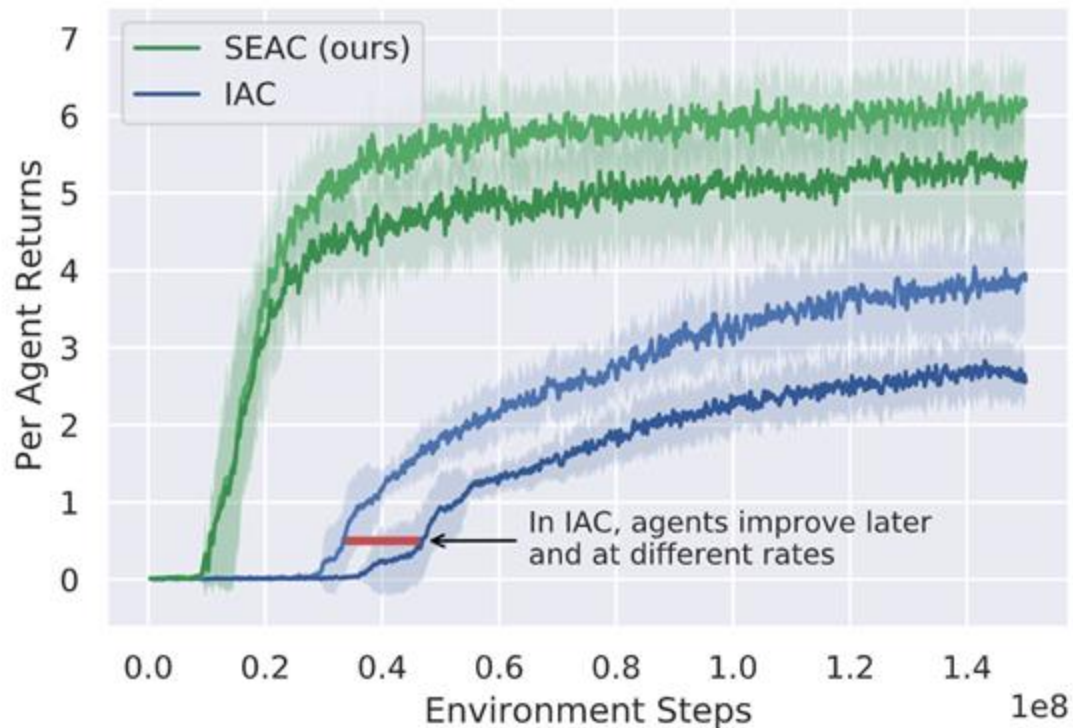
Analysis (2)



- Agents learn simultaneously which helps in exploring promising joint actions more

Best vs. Worst performing agents on RWARE, (10x20), four agents

Analysis (2)



Best vs. Worst performing agents on RWARE, (10x20), four agents

- Agents learn simultaneously which helps in exploring promising joint actions more
- Synchronise training progress of agents

Contributions

We propose a novel experience sharing method (Shared Experience Actor-Critic or SEAC) that combines gradients of multiple agents to share experience between agents.

- Evaluated in four sparse-reward multi-agent environments
- Consistently outperforms baselines and three state-of-the-art MARL algorithms (MADDPG, QMIX, ROMA)
- SEAC learns in fewer steps and converges to higher returns
- In harder tasks, sharing experience makes the difference between not learning at all and learning

Conclusion

- Using our method, agents learn similar – but not identical policies.
 - Facilitates coordination between agents
- Exploration is improved:
 - Agents tend to pick-up behaviors concurrently: more promising joint actions are explored more
- Simple and general method (can be used to extend any on- and even off-policy algorithms)

Future Work

- Relax assumptions about the task required to share experience
- Learn λ for agents (which agents share experience with whom)

Shared Experience Actor-Critic for Multi-Agent Reinforcement Learning

Filippos Christianos, Lukas Schäfer, Stefano V. Albrecht

Links to code:

SEAC:	https://github.com/uoel-agents/seac
RWARE:	https://github.com/uoel-agents/robotic-warehouse
LBF:	https://github.com/uoel-agents/lb-foraging

Shared Experience Actor-Critic for Multi-Agent Reinforcement Learning

Filippos Christianos, Lukas Schäfer, Stefano V. Albrecht

Arxiv: <https://arxiv.org/pdf/2006.07169.pdf>

Contact: f.christianos@ed.ac.uk

Filippos Christianos, Lukas Schäfer, and Stefano V. Albrecht. "Shared Experience Actor-Critic for Multi-Agent Reinforcement Learning." In *34th Conference on Neural Information Processing Systems 2020*.

Thank you for your Attention! Any Questions?

Filippos Christianos, Lukas Schäfer, Stefano V. Albrecht

NeurIPS

Poster: **Poster Session 5**
on Thu, Dec 10th, 2020 @ 18:00 – 20:00 CET

