

Social Reinforcement Learning

Natasha Jaques

What abilities does an AI assistant need?



<https://clipart.me/free-vector/kitchen-robot>

Coordinate in
shared spaces



Learn from
human
interaction

MON	TUE	WED	THU	FRI	SAT
7	8	9	10	11	12
INS: NeurIPS					
	[Live stream] GWE EMEA Talks: Include Yourself, Deborah 9 - 10:30am	https://neu I/2020/pro 4.1816538	MOVE CAR STREET Cleaning		
		NeurIPS tu uncertainty estimation		Deep RL worksho p	Live Q1: Moderate session
		NeurIPS poster session 12 - 2pm https://neuri ps.gather.t	NeurIPS poster session 5 12 - 2pm https://neuri ps.gather.t	Deep RL worksho p	Cooperat 11:15a 11:15a
		NeurIPS poster session 12 - 2pm https://neuri ps.gather.t	[HOLD] [Option 1 - 2pm 1 - 2pm	Send in snippets for th multiagent SF-MOA 1 - 2pm	Offline RL works hop 12 - 9
NS / Doug/Natasha, 12 neuri 2 - 9:30 TensorFlow Office 1 - 2pm	DNS 12 - 9	Synthetic C: TF-Agents	Research Scientist - Ka 2pm, Video Link (above	Adversarial Surpris 8 - 9:30pm	
Learning S 1:30 - 2:30	NeurIPS at D/E/J Jour 2:30 - 2:30pm	DNS 2:30 - 9pm	Over/under on NeurIPS papers 3 - 5pm	Berkeley R 3 - 4pm	
Black / Host Black 2:30 - 2:30pm	NeurIPS at D/E/J Jour 3 - 5pm https://gather.town /app/86gkchhsk	WIML mentoring 3:15 - 5pm https://neuri ps.gather.t own/app/DoxKtP9W7 y8yX9f7e7m7worksho	Pay for neurips gather	Resistance 4:30pm	CooperativeAI 4pm, https://ne
Weekly Robots 3:30pm, MTV-Rush	Equity and Ethics in AI from the Perspective of Black Women in AI/STEM 5 - 7pm	Manfred / Natasha cha	Brain Research Virtual 7 - 8pm	Deep RL works 5pm, https://si Baby minds p 5:45pm, https	
	NeurIPS Q&A: Advers	NeurIPS invited talk: A Future of Work for the	neuri after party 9 - 11pm https://gather.town/app /86gkchhskXVdp7LD/n neuri after party	GO TO THIS: Deep RL GO TO THIS: deep RL 8 - 9pm, https	
EURIPS NeurIPS orals, reinforcement learning	NeurIPS orals/spotl lights: RL 5pm - 12a https://gather.town/app /86gkchhskXVdp7LD/n neuri after party	NeurIPS orals/spotl lights: RL 5pm - 12a https://gather.town/app /86gkchhskXVdp7LD/n neuri after party		Deep RL workshop post 9pm, https://www.goo	

Learn
complex
tasks



Generalize to new
environments

Hypothesis: social learning can help address all of these desiderata

Coordinate in
shared spaces

Learn from
**human
interaction**

Learn
complex
tasks

Generalize to new
environments

Social learning

= Learning from other intelligent agents in your environment

- Helps humans and animals...

(Laland, 2017;
Henrich, 2015)

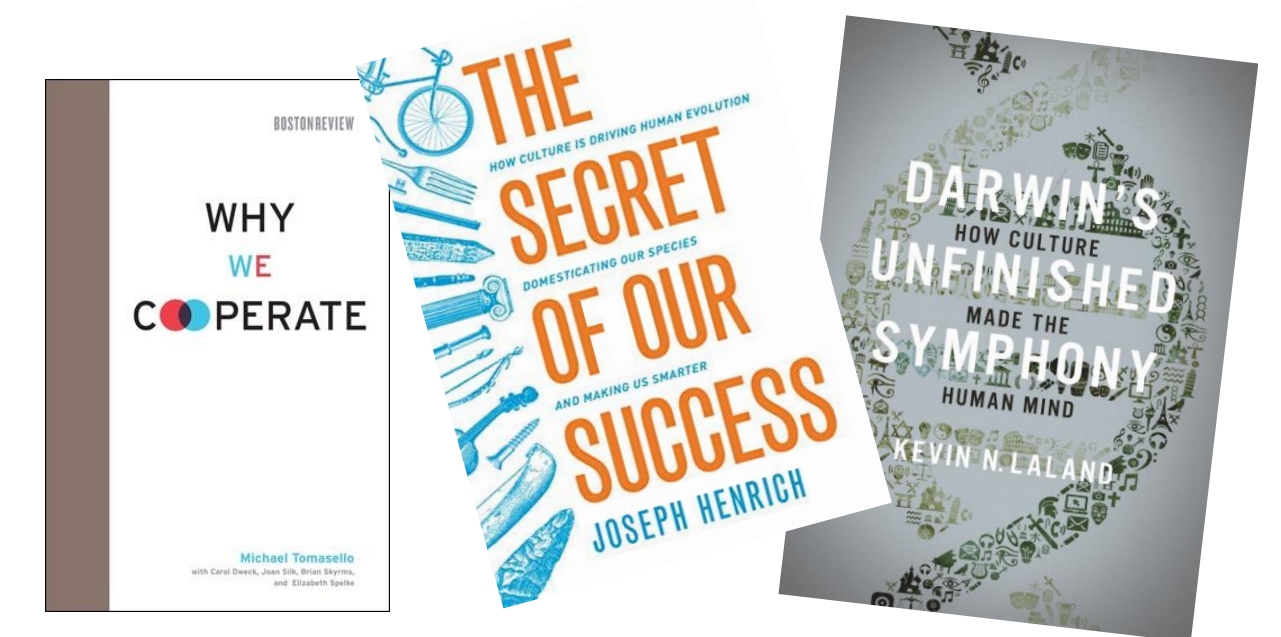


Generalize to new environments

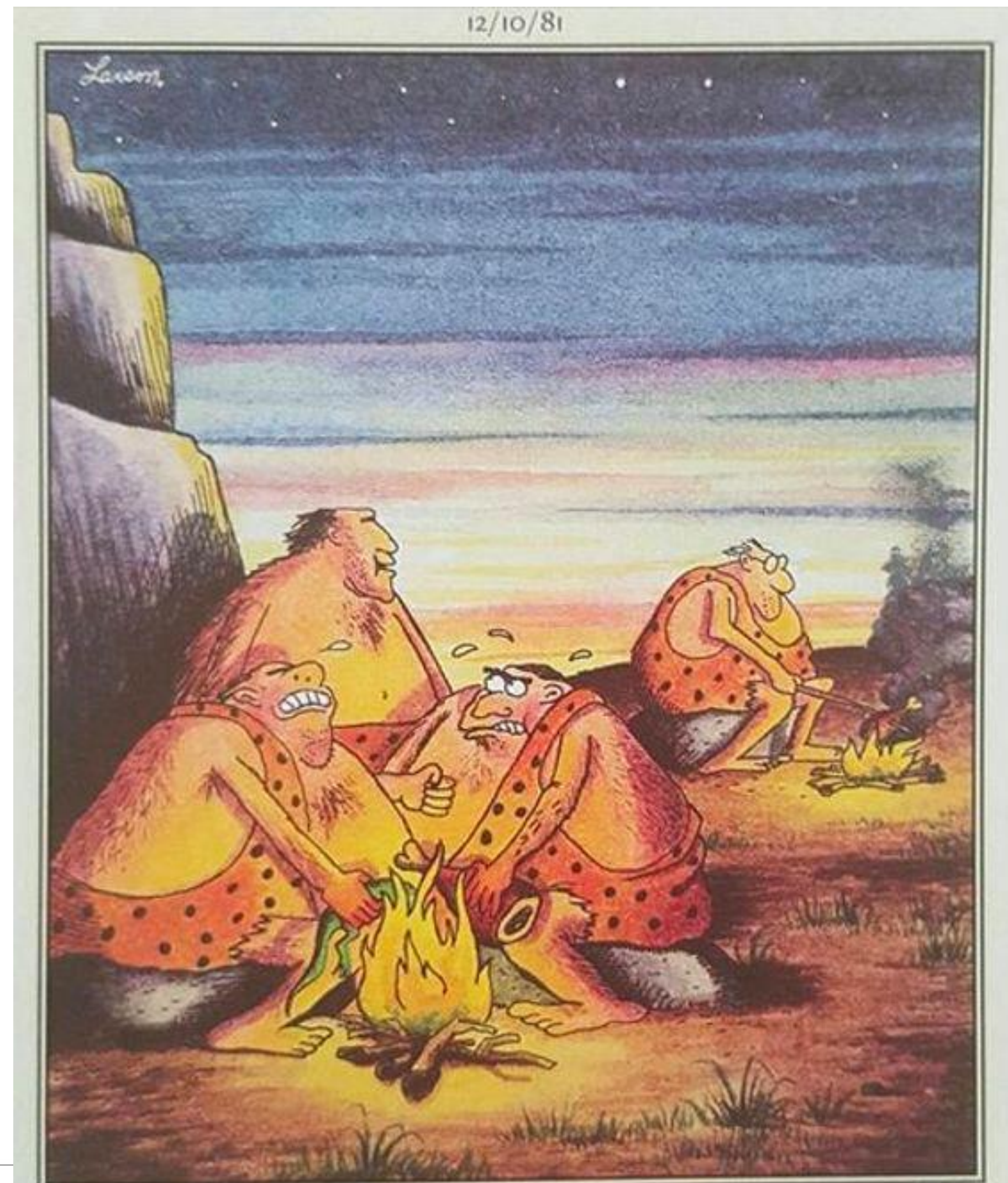
- Is key to **human** cognitive development, cultural and technological evolution (Henrich, 2015; Humphrey, 1976, Tomasello, 2009)



Learn complex behavior



Social learning can accelerate learning



“Hey! Look what Zog do!”

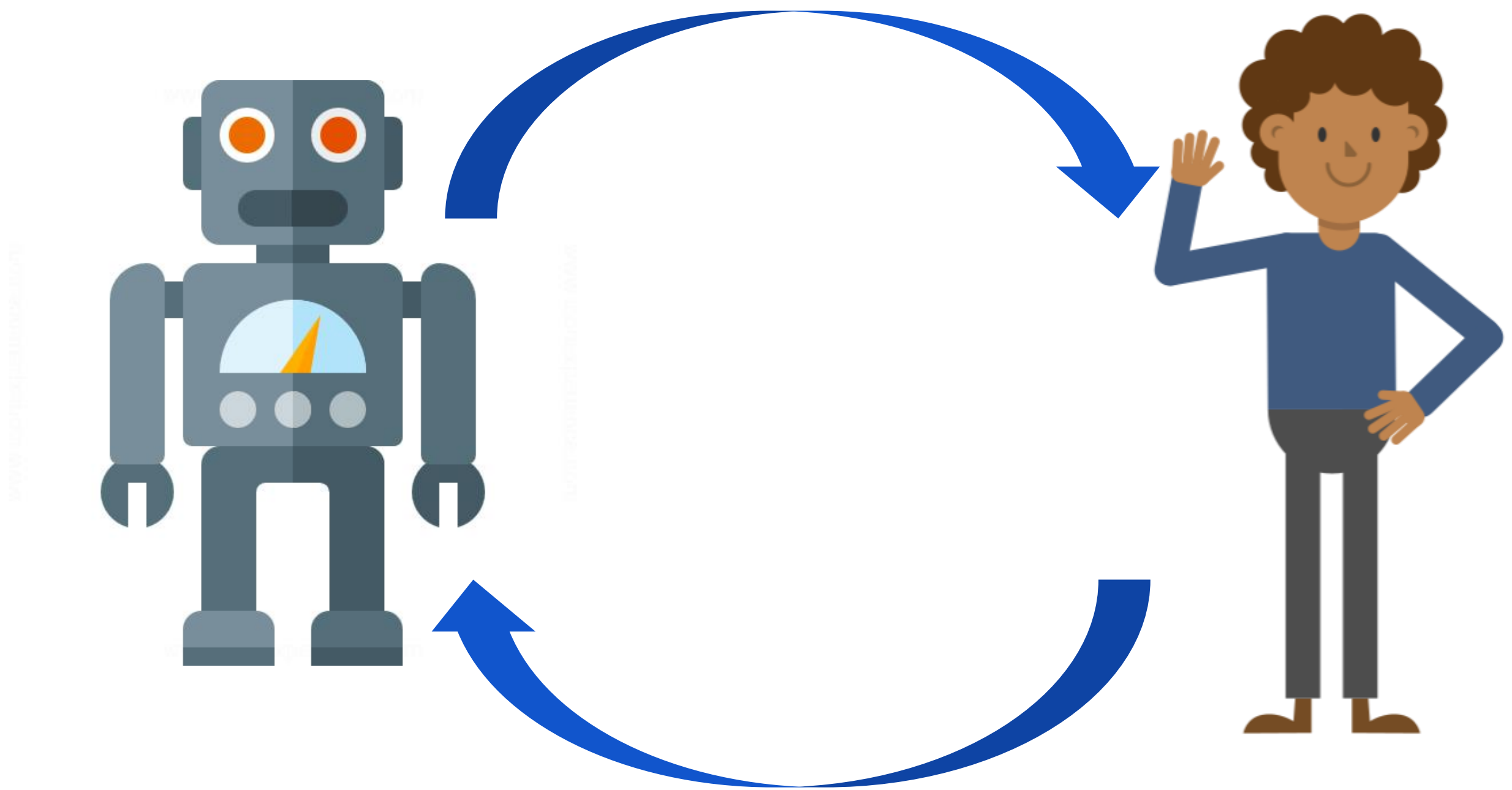
Individual learning can be unsafe,
error prone, time consuming

Social learning can enable you to
“stand on the shoulders of giants”

Social Reinforcement Learning



Multi-agent interaction



Human-AI interaction

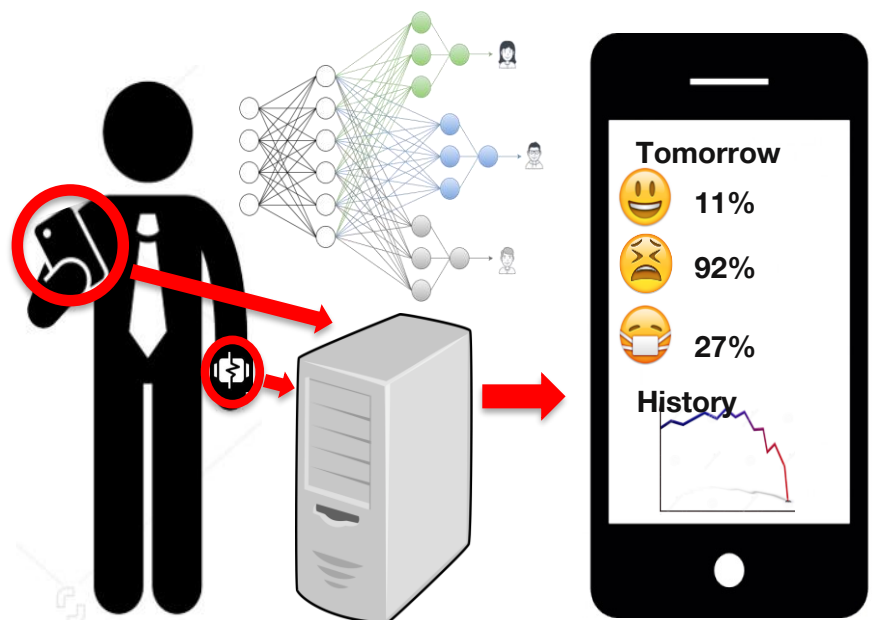
Broad overview of Social RL

Human-AI

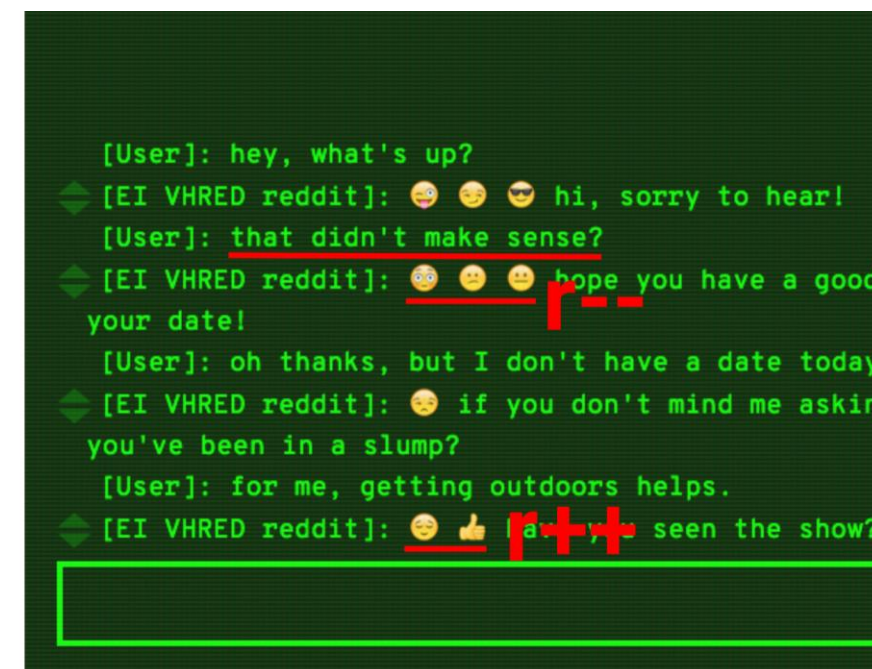


Multi-agent

Detecting social
& affective cues



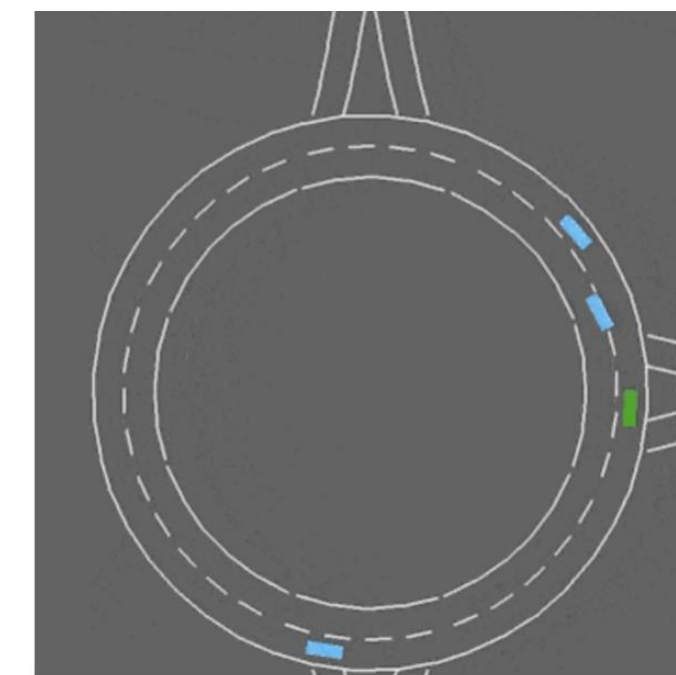
Learning from
human social cues



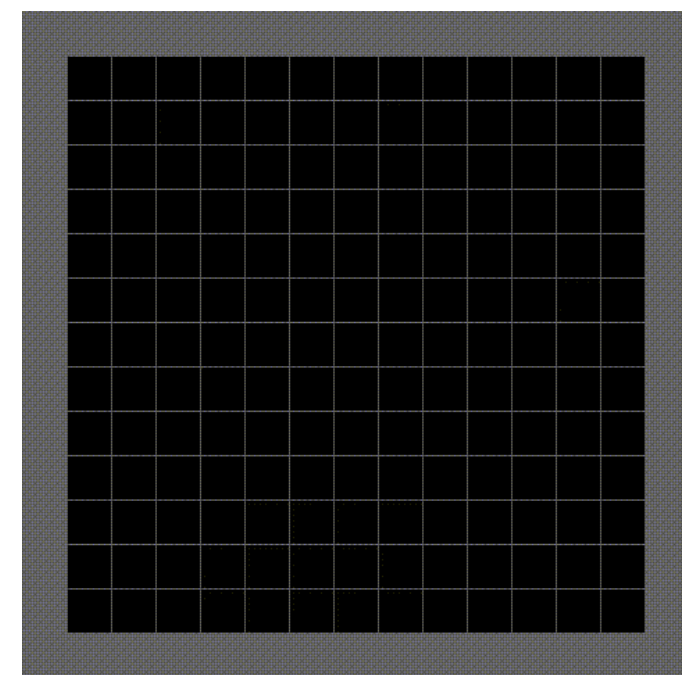
Coordination



Multi-agent
social learning



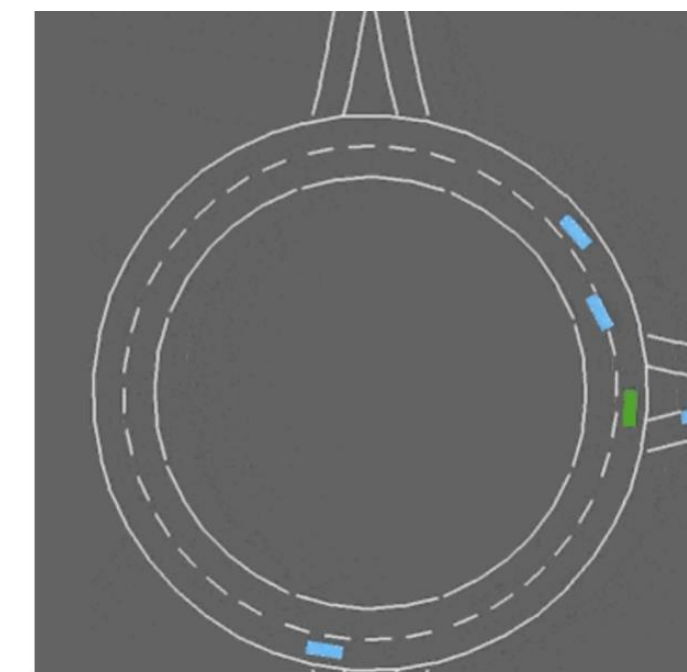
Emergent
Complexity



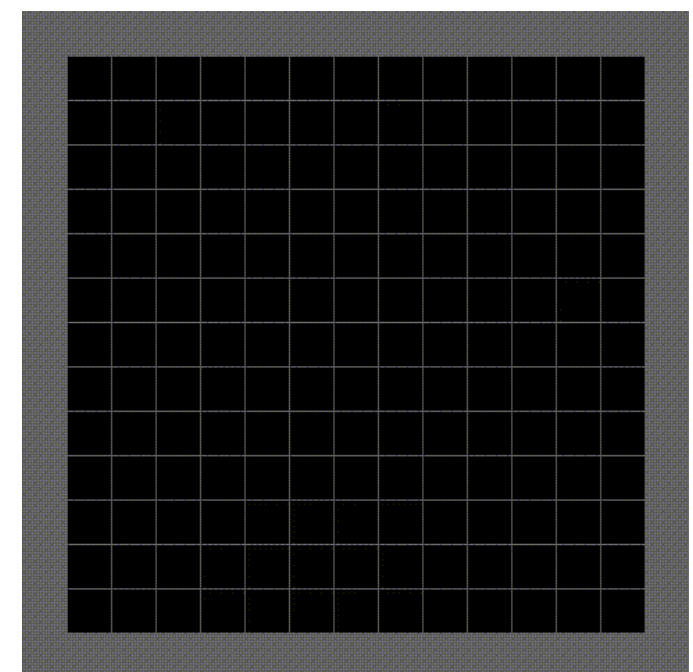
Broad overview of Social RL

Human-AI ←————→ Multi-agent

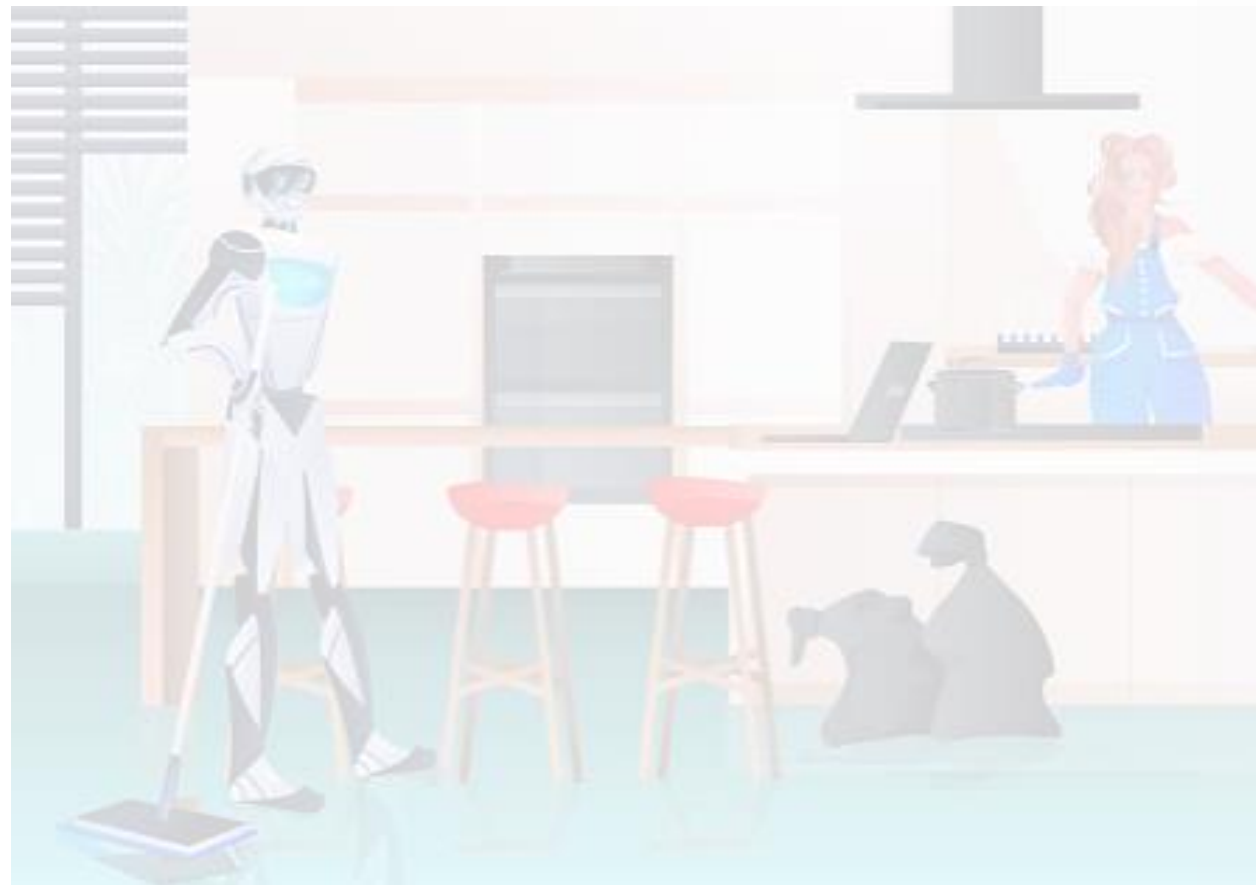
**Multi-agent
social learning**



**Emergent
Complexity**

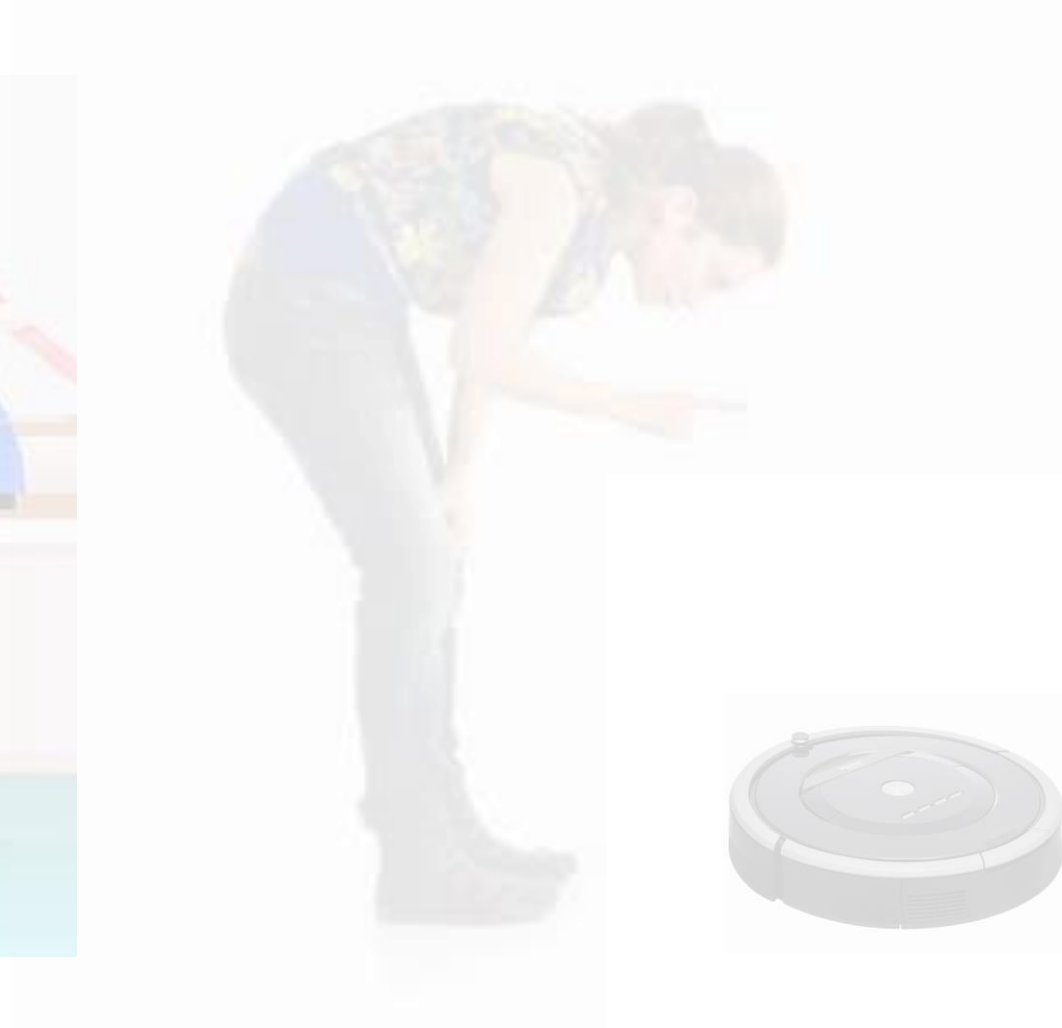


What abilities does an AI assistant need?



<https://clipart.me/free-vector/kitchen-robot>

Coordinate in shared spaces



Learn from **social cues**

MON 7	TUE 8	WED 9	THU 10	FRI 11	SAT 12
INS: NeurIPS					
			MOVE CAR Street Cleaning		
	[Live stream] GWE EMEA Talks: Include Yourself, Deborah 9 - 10:30am	https://neurips.cc/virtual/2020/protected/papers.html			
		NeurIPS tutorial uncertainty estimation			
NS Doug/Natasha, 12 neurips 2-9:30	DNS 12-9	NeurIPS poster session 12-2pm	NeurIPS poster session 5 12-2pm	Deep RL workshop 1-2pm	Live Q1 Moderate session Cooperative AI 11:15a-11:55a
TensorFlow Office 1-2pm	Attending 1:30-2:30	https://neurips.cc/virtual/2020/protected/papers.html	[HOLD] Option 1-2pm	Send in snippets for the multiagent SF-MOA 1-2pm	Offline RL workshop 12-9
Learning S Black Host Black 2:30-2:30pm	NeurIPS at D/E/J Jour party: Europe edition 3-5pm	Synthetic C TF-Agents 2:30-9pm	Research Scientist - Ka 2pm, Video Link (above)	Adversarial Surpris 8-9:30pm	
Weekly Robots 3:30pm, MTV-Rush	https://gather.town/join/8gkchHak	WIML mentoring 3:15-5pm	Over/under on NeurIPS papers 3-5pm	Berkeley R 3-4pm	
	Equity and Ethics in AI from the Perspective of Black Women in AI/STEM 5-7pm	Manfred / Natasha cha	Pay for neurips gather	Resistance 4:30pm	CooperativeAI 4pm, https://ne
		NeurIPS Q&A: Advers	Brain Research Virtual 7-8pm	GO TO THIS: Deep RL v GO TO THIS: deep RL f	
NeurIPS orals, reinforcement learning	NeurIPS orals/spotlights: RL 5pm-12a	NeurIPS orals/spotlights: RL 5pm-12a	NeurIPS invited talk: A Future of Work for the	Deep RL workshop 9pm, https://www.goog	
		NeurIPS orals/spotlights: RL 5pm-12a	neurips after party 9-11pm		

Learning complex tasks



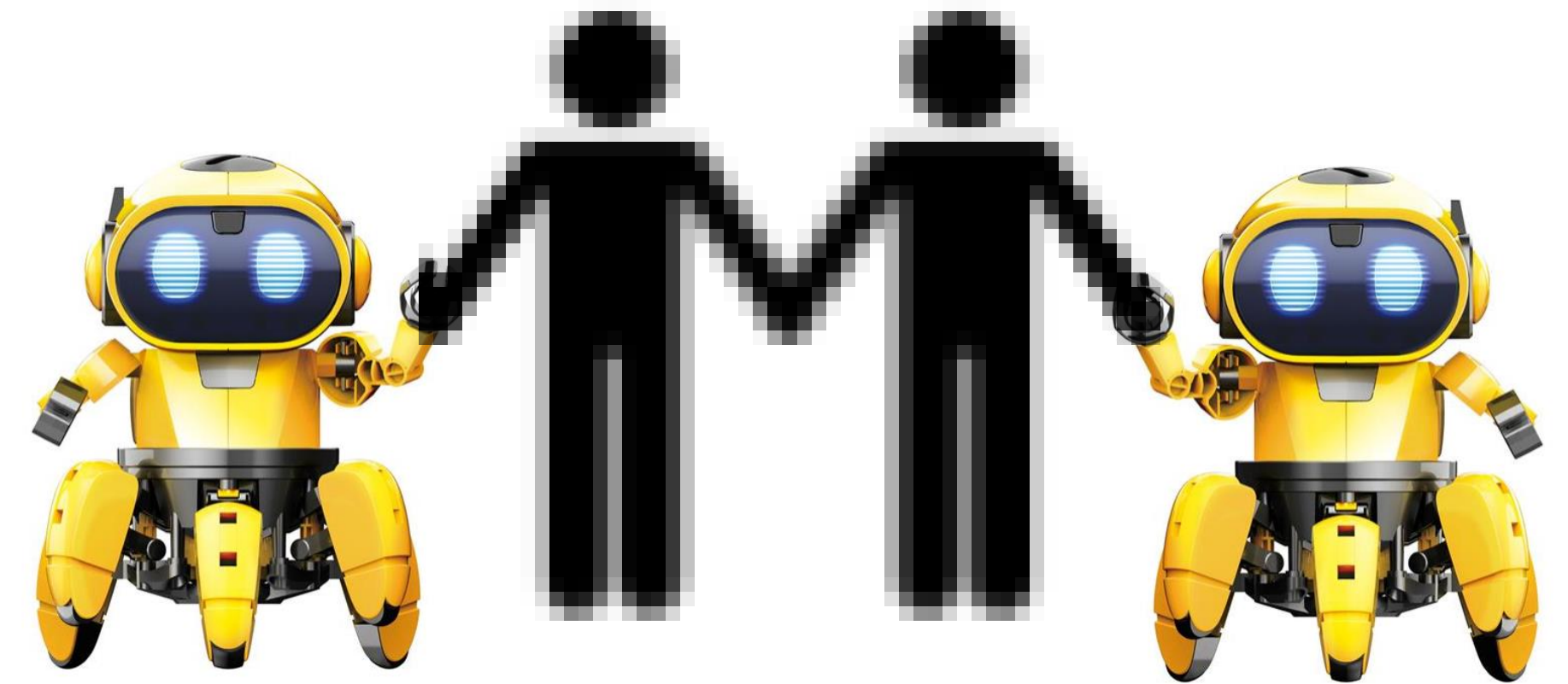
Generalize to new environments

Outline

Multi-agent...

1. Emergent complexity

2. Social Learning

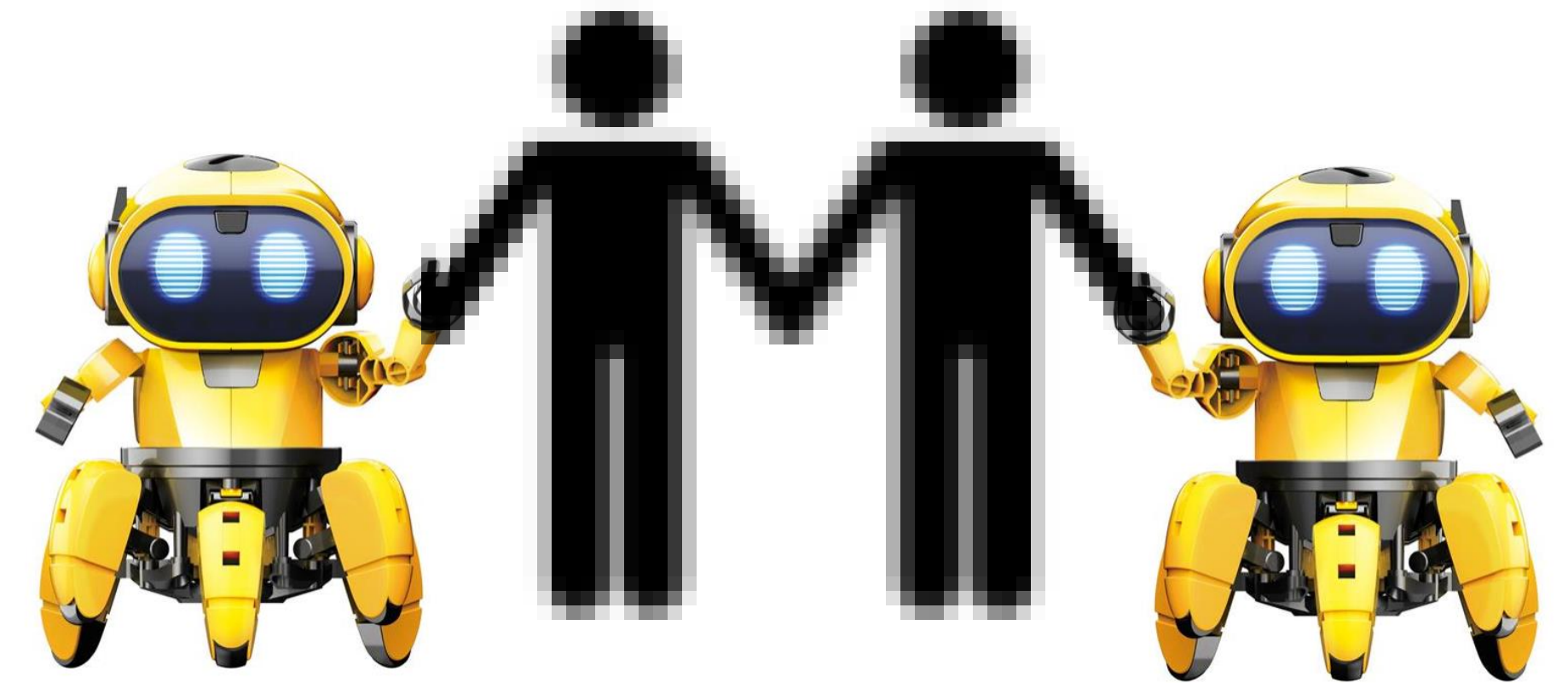


Outline

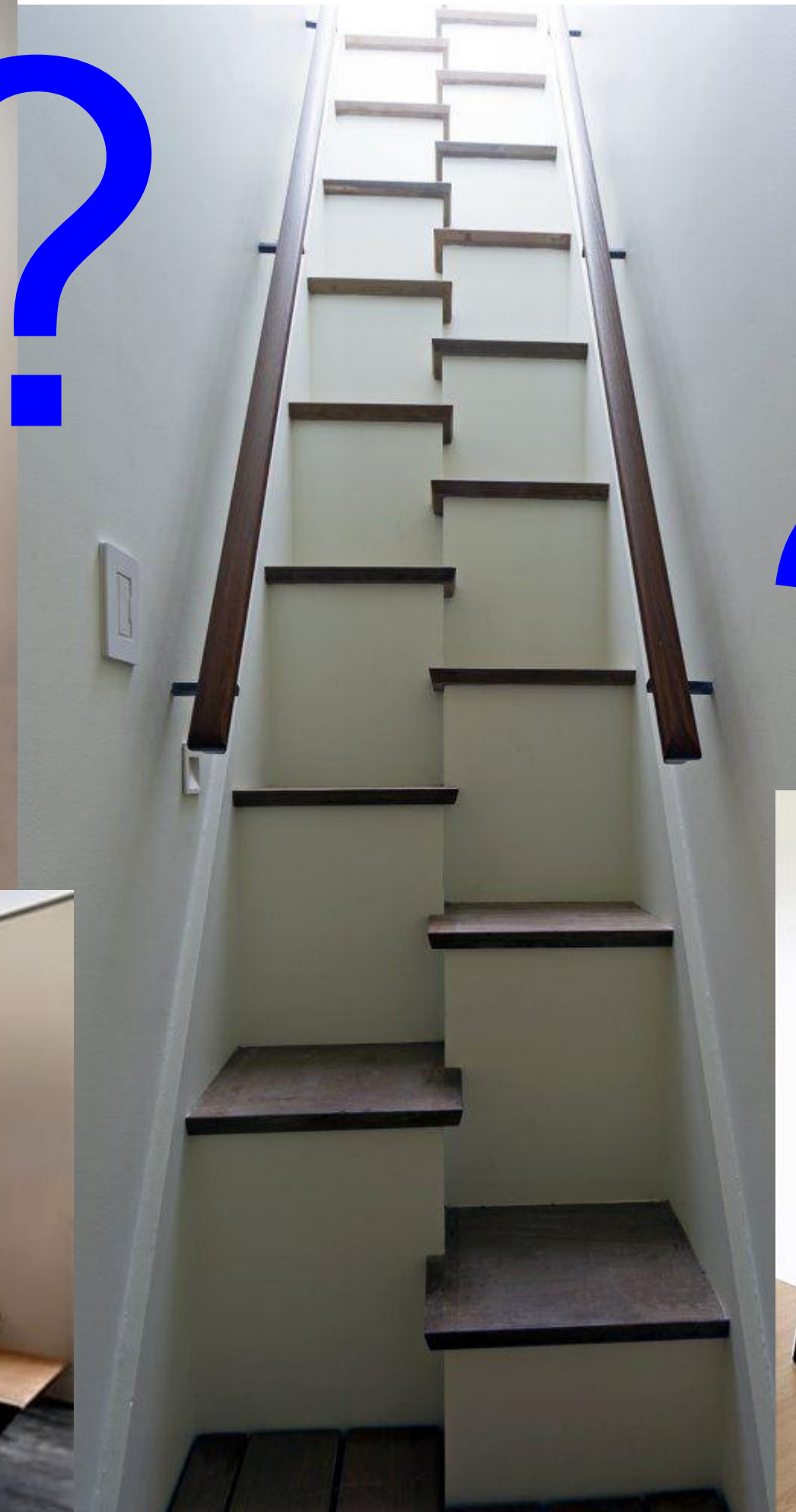
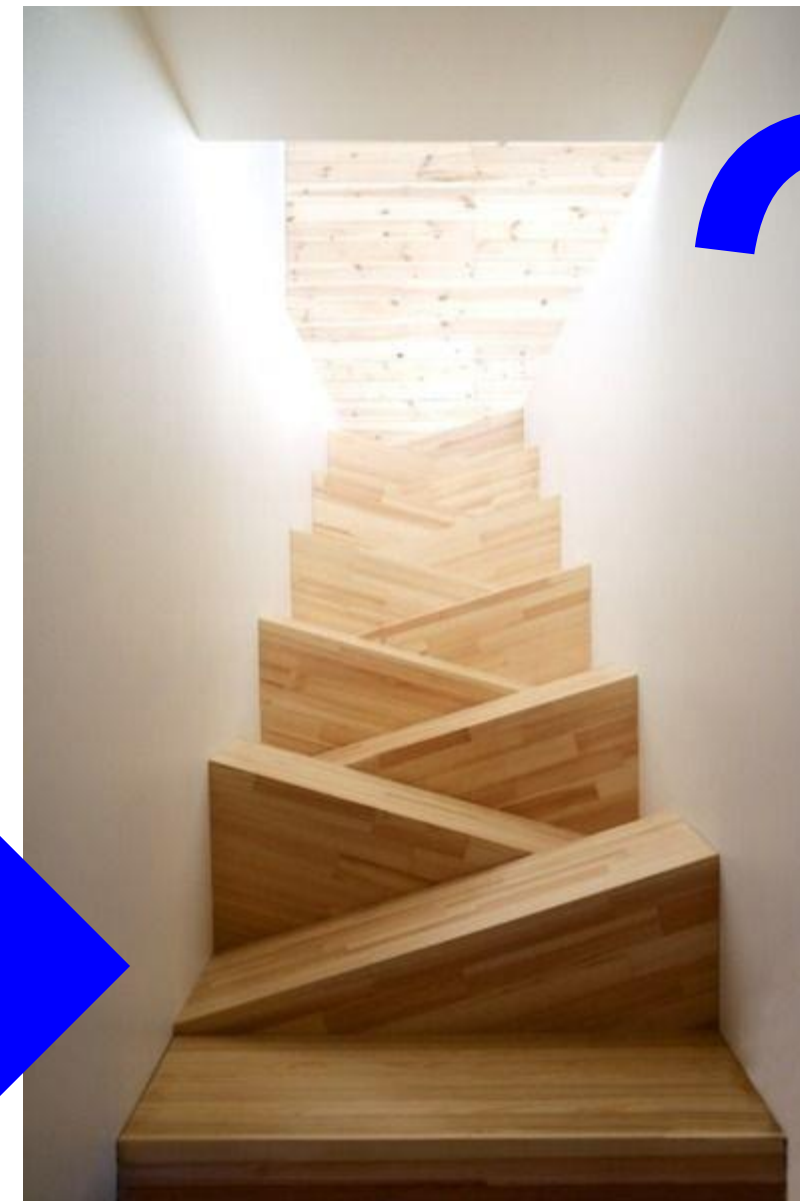
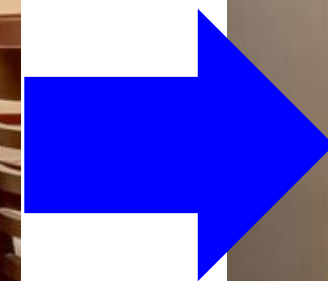
Multi-agent...

1. Emergent complexity

2. Social Learning



Generalization and transfer in RL



Domain Randomization:

[1] Fereshteh Sadeghi and Sergey Levine. Cad2rl: Real single-image flight without a single real image. arXiv preprint arXiv:1611.04201, 2016.

[2] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS), page 23–30. IEEE, 2017.

How can we get:

- **complex** training environments
- that cover **unknown, real world** challenges
- without having to program them ourselves **by hand**?

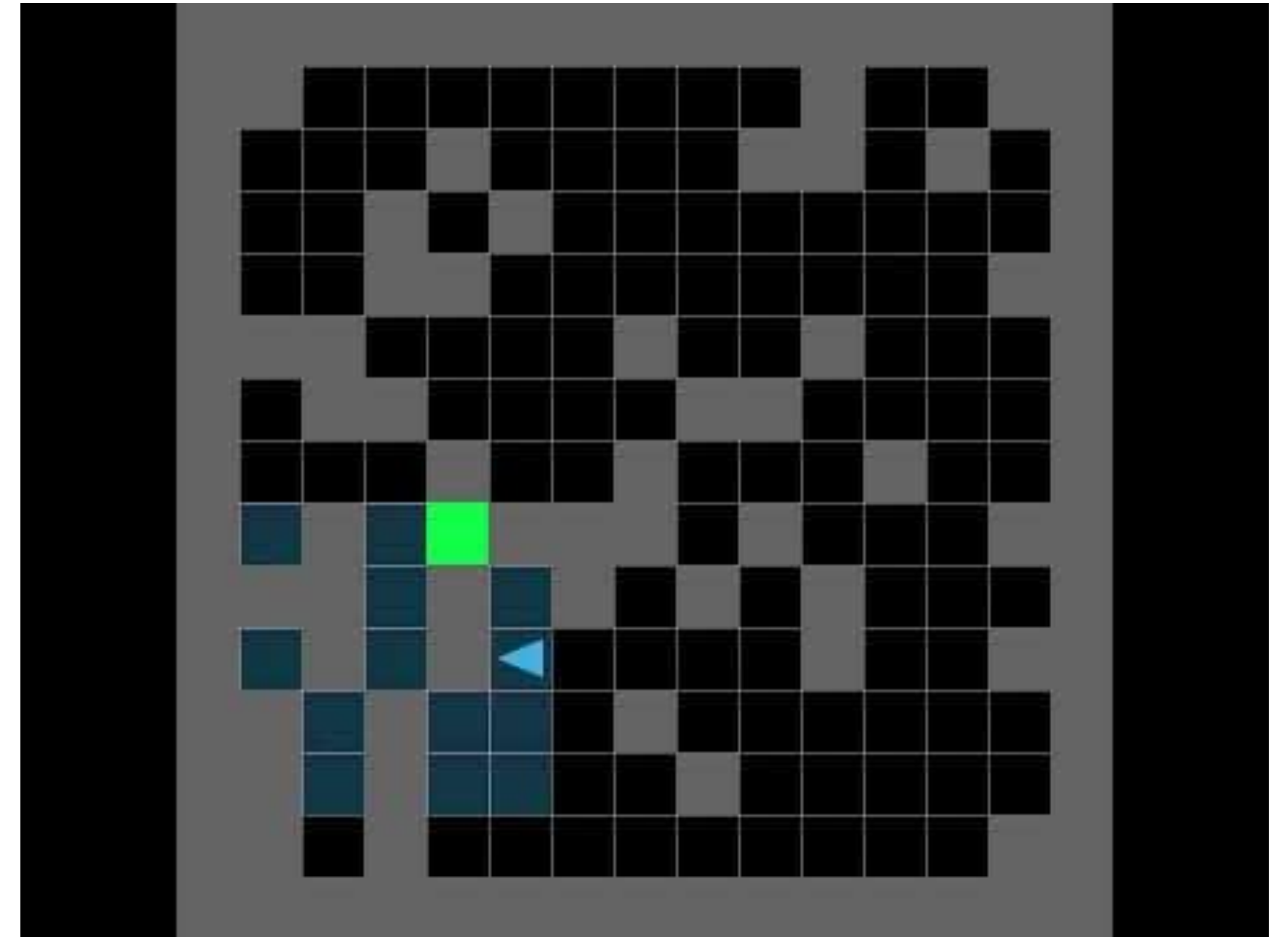


Adversarial environment generation

Environment-generating
adversary policy

$$\min_{\pi^E} \max_{\pi^P} R(\tau)$$

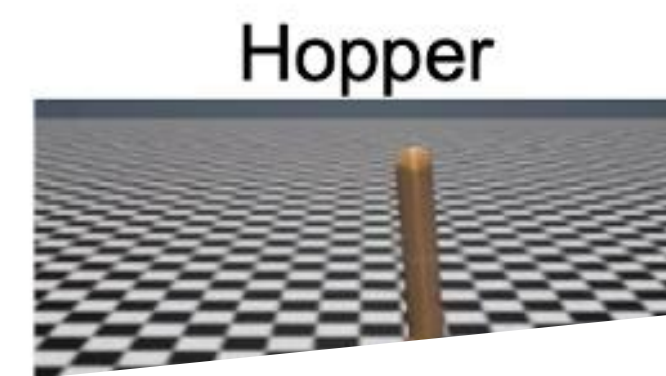
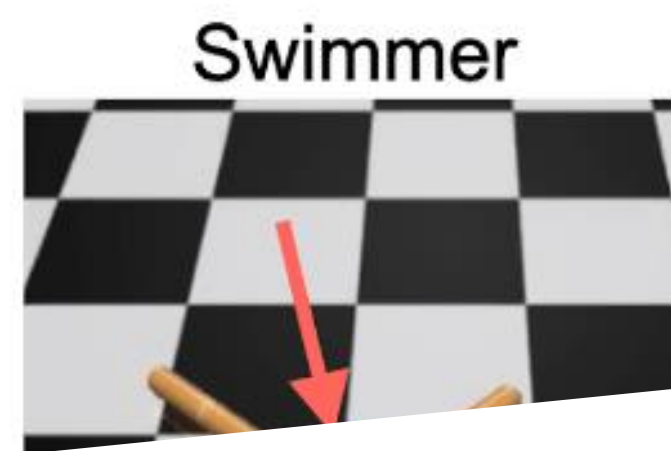
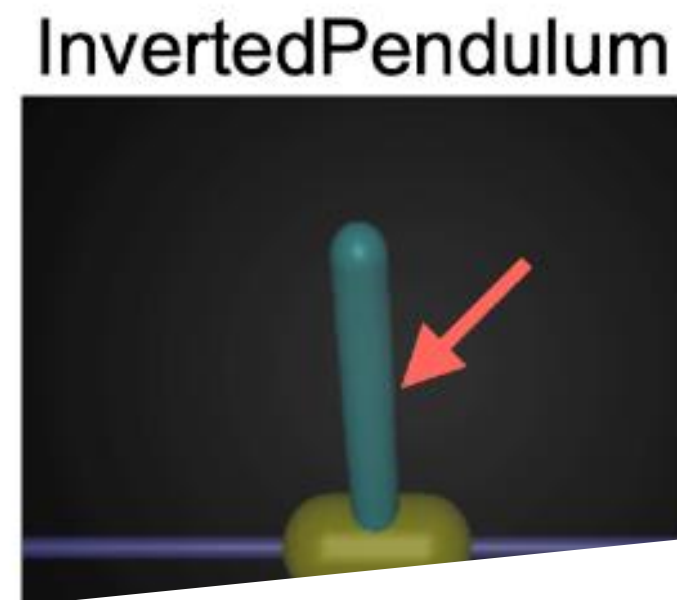
Player agent policy



Minimax adversary

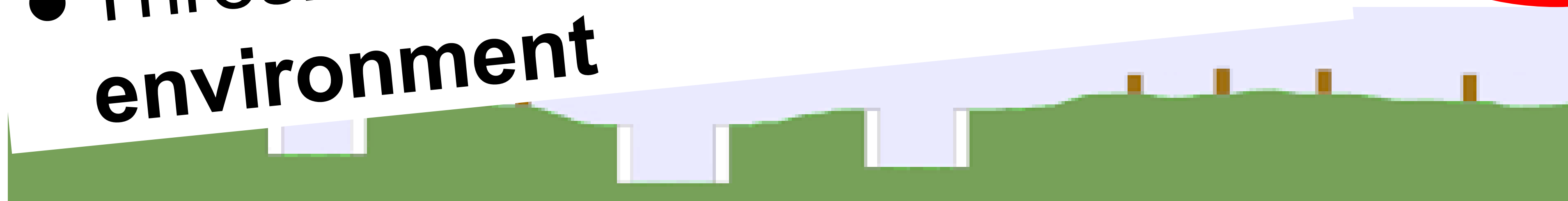
Prior work with deep RL adversaries

RARL



- Hand-tuned constraints on adversary
- Thresholds specific to environment

POET



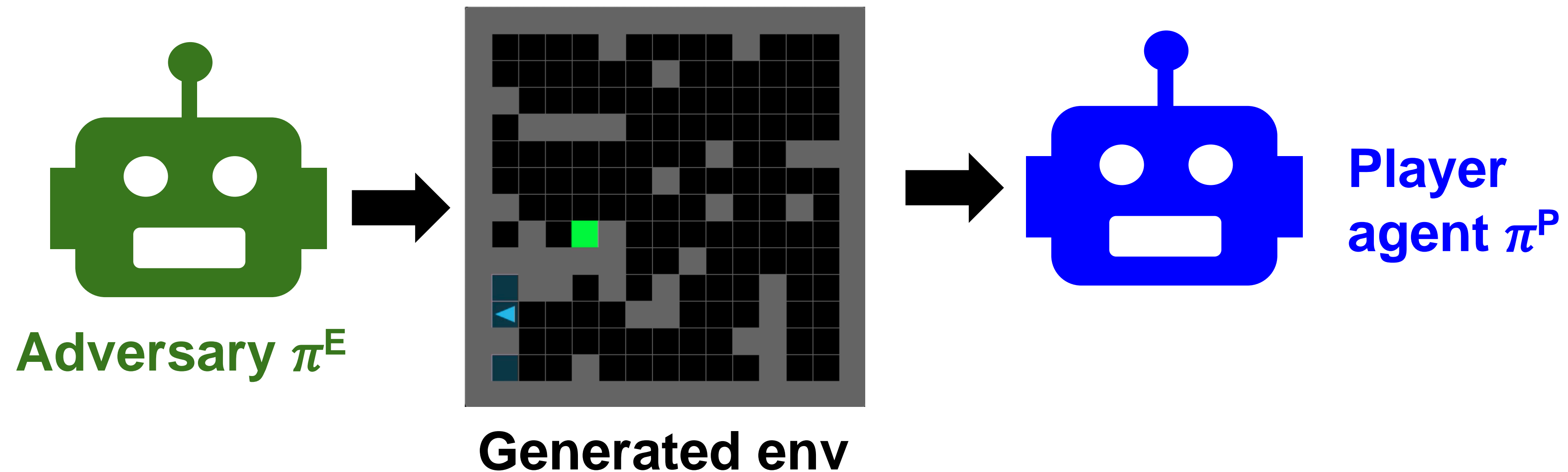
Is there a more elegant way to ensure the adversary does not create **impossible** environments?

Or even environments **tailored to current skill** of the agent? (**automatic curriculum**)

→ **Let's make this more multi-agent**

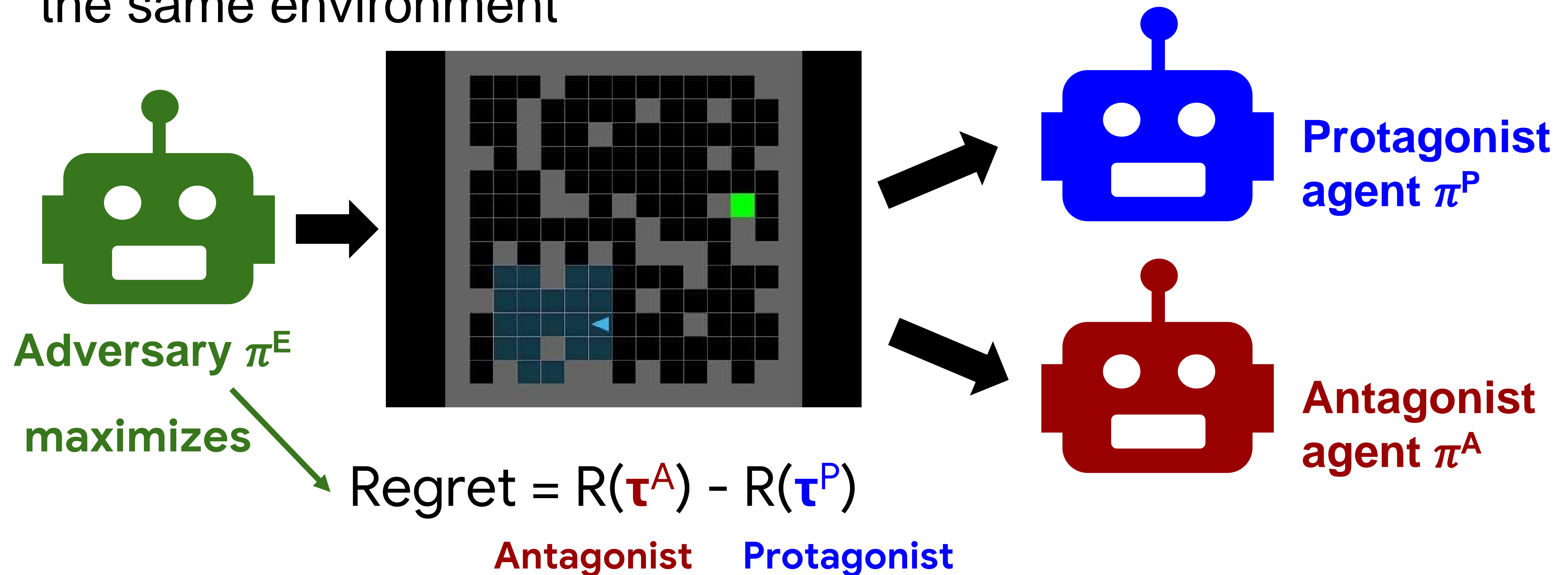
PAIRED

- Constrain the adversary using the performance of a **second agent** in the same environment



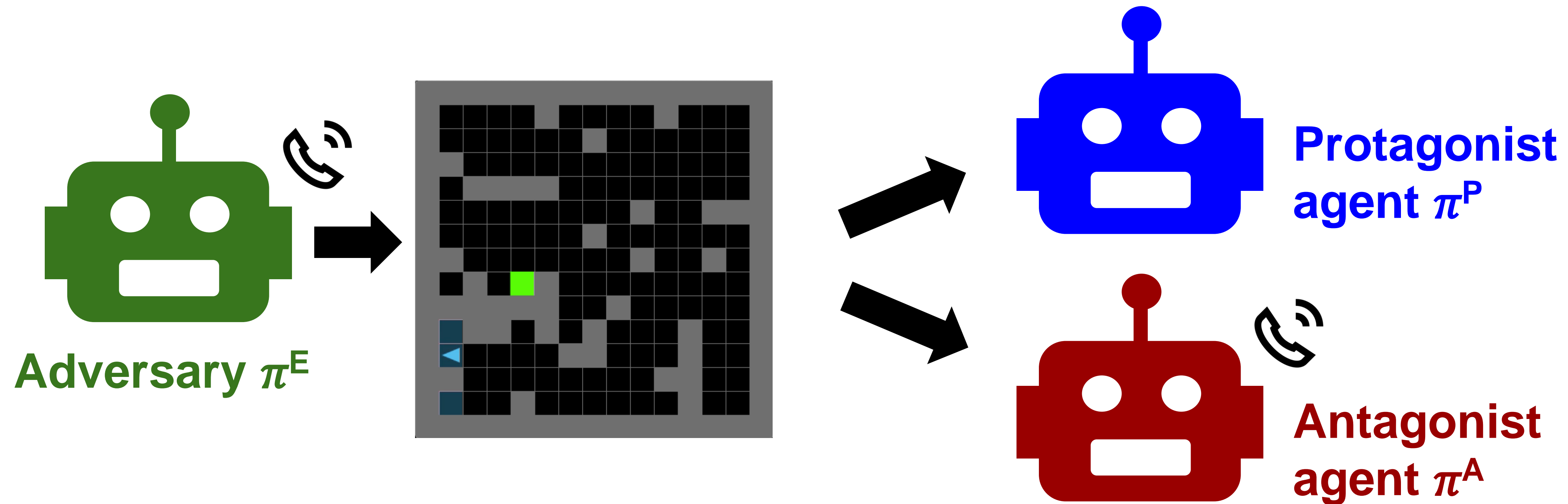
PAIRED *(Protagonist Antagonist Induced Regret Environment Design)*

- Constrain the adversary using the performance of a **second agent** in the same environment

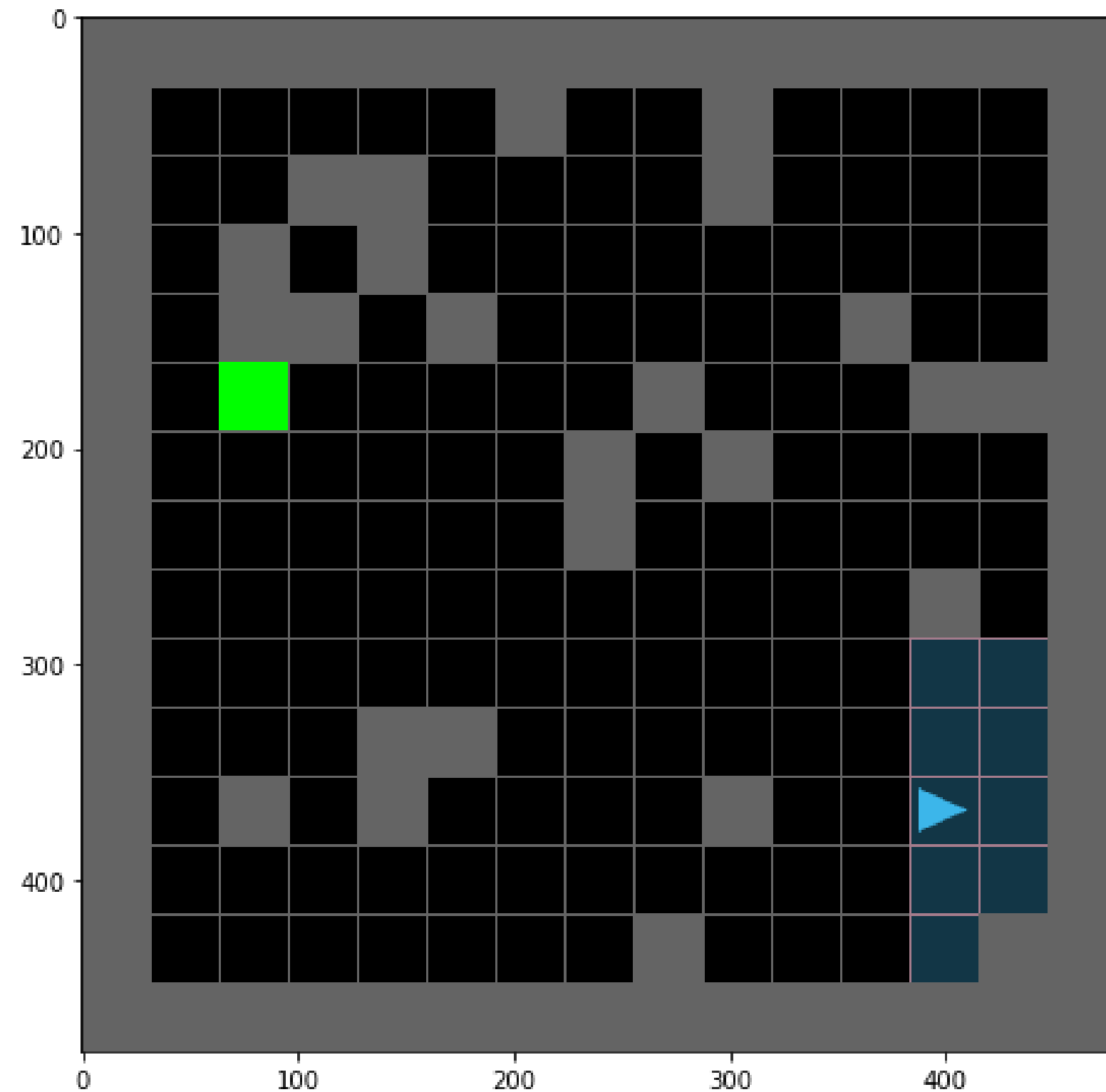


PAIRED as Minimax Regret

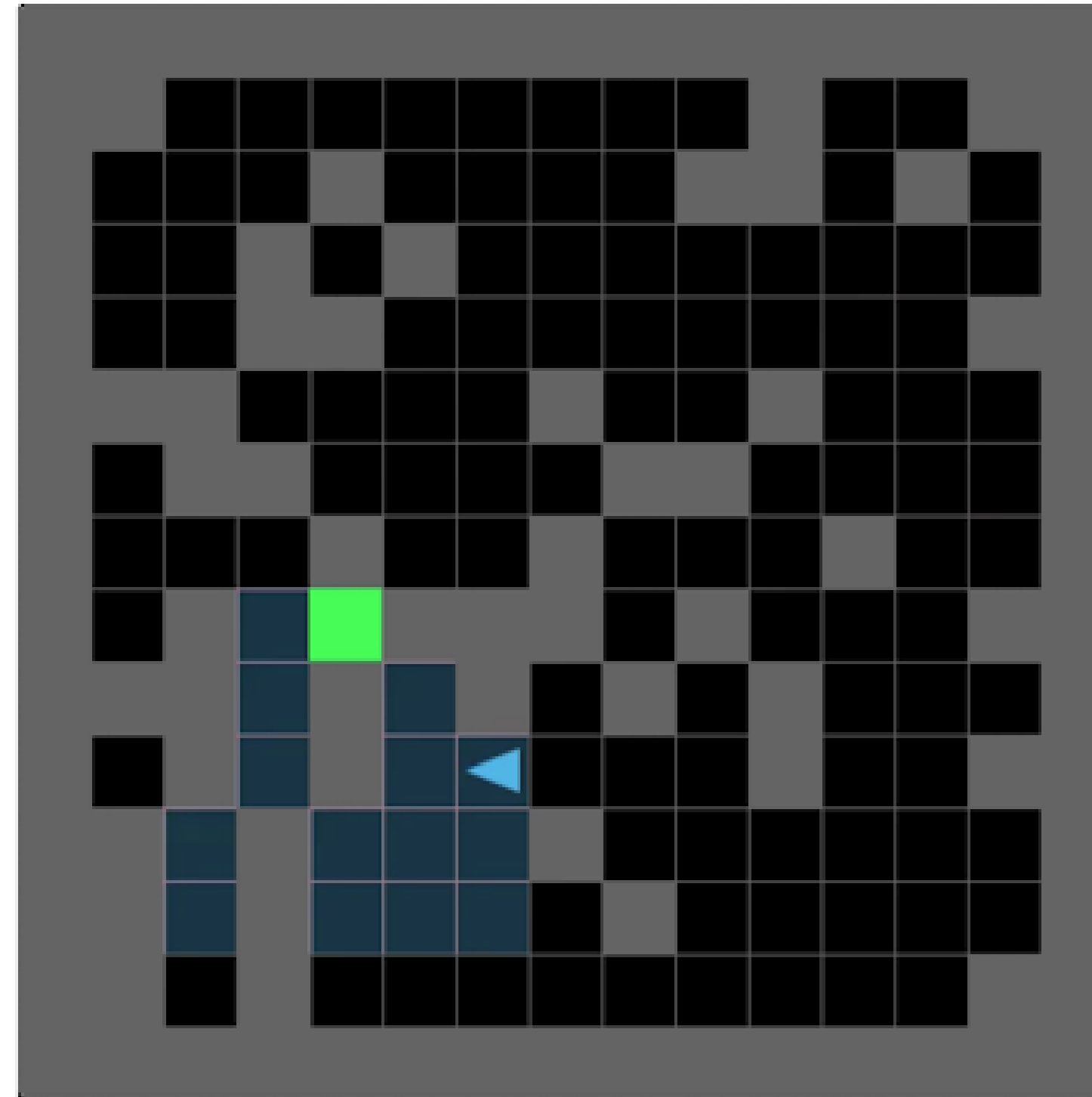
Theorem 2. *Let $(\pi^P, \pi^A, \vec{\theta})$ be in Nash equilibrium and the pair $(\pi^A, \vec{\theta})$ be jointly a best response to π^P . Then $\pi^P \in \operatorname{argmin}_{\pi^P \in \Pi^P} \{ \operatorname{argmax}_{\pi^A, \vec{\theta} \in \Pi^A \times \Theta^T} \{ \operatorname{REGRET}^{\vec{\theta}}(\pi^P, \pi^A) \} \}$.*



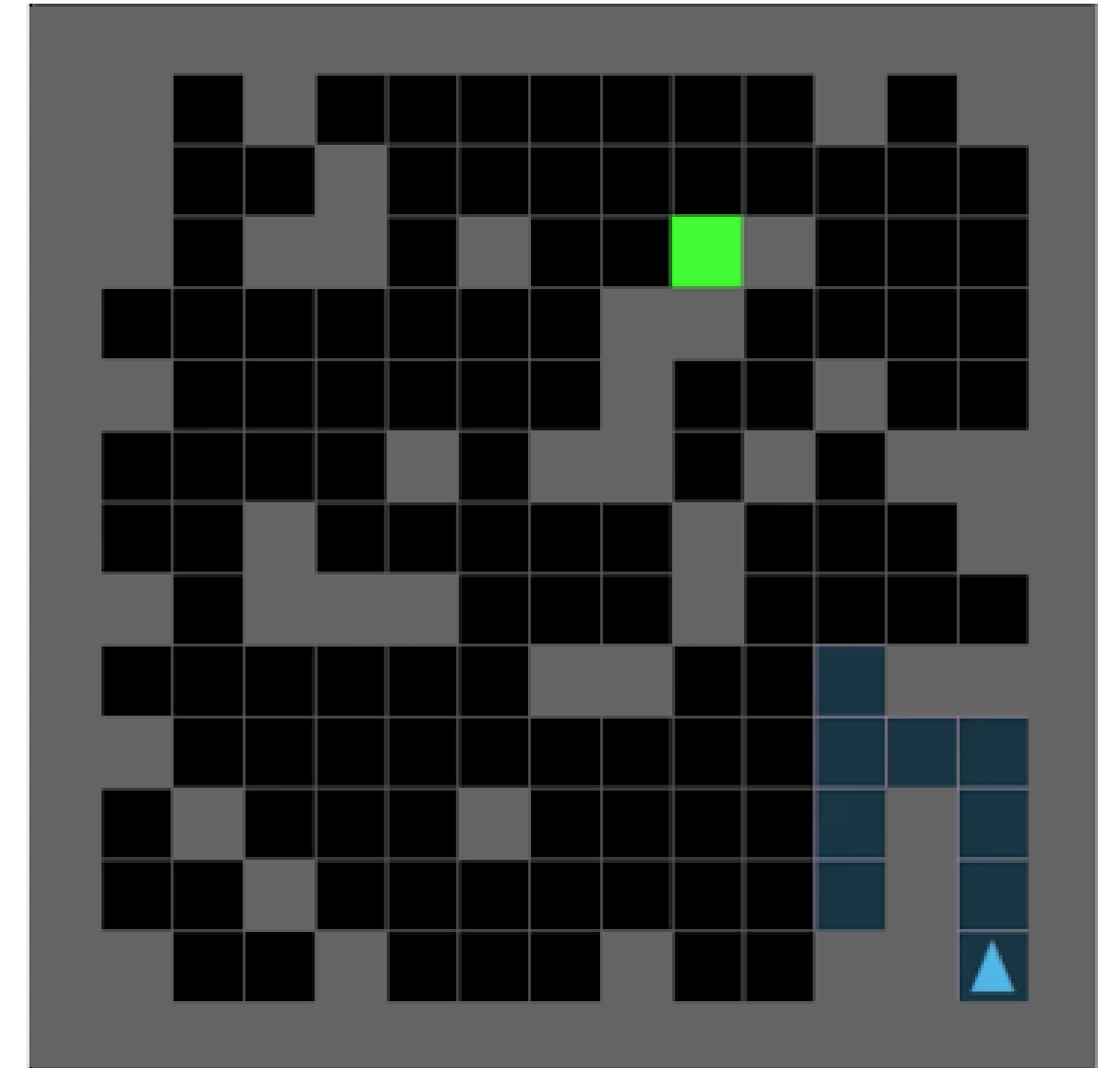
Generated environments



(a) Domain Randomization



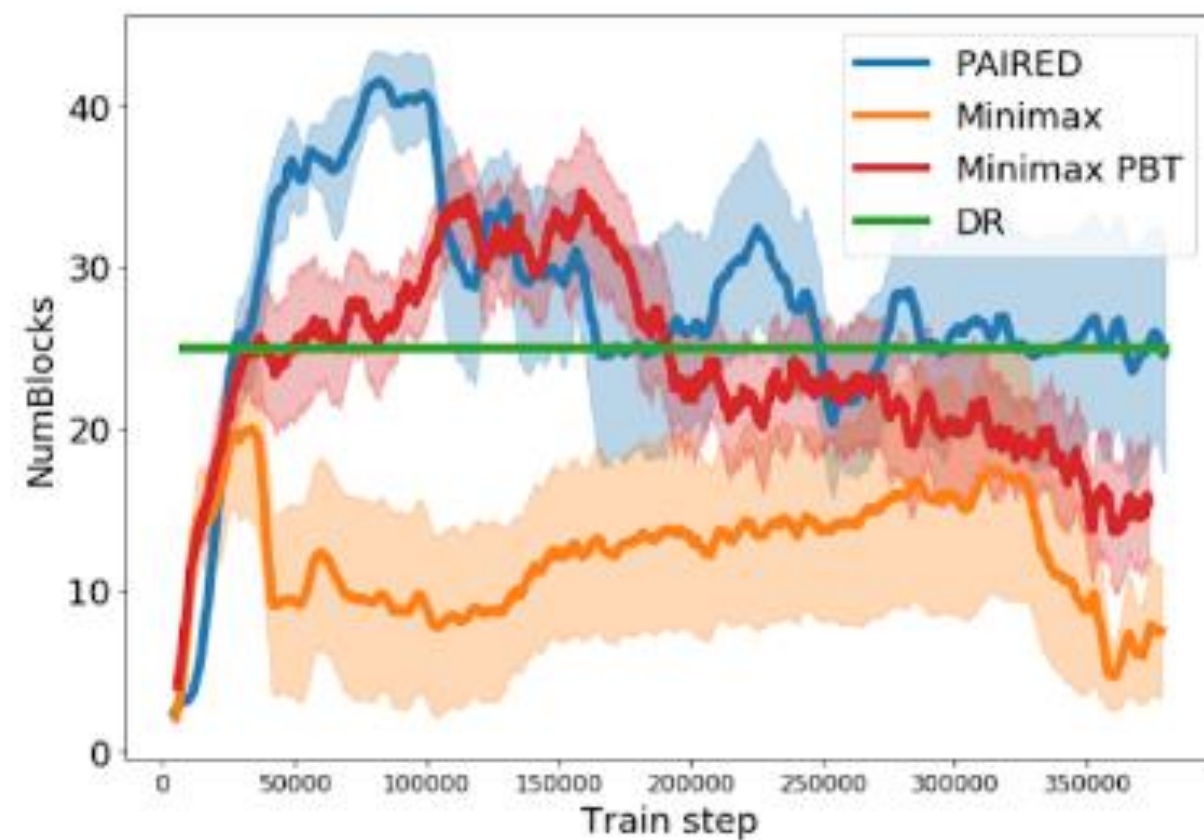
(b) Minimax Adversarial



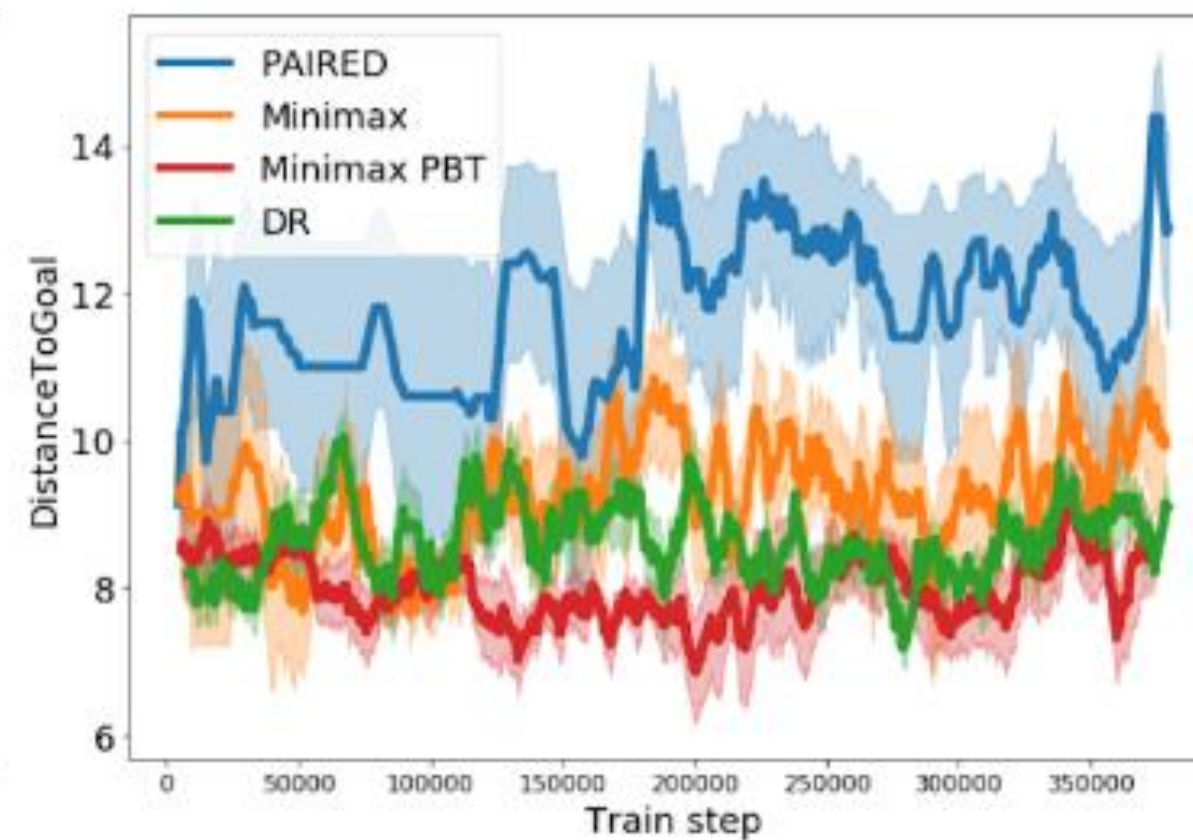
(c) PAIRED (ours)

Emergent complexity and curriculum

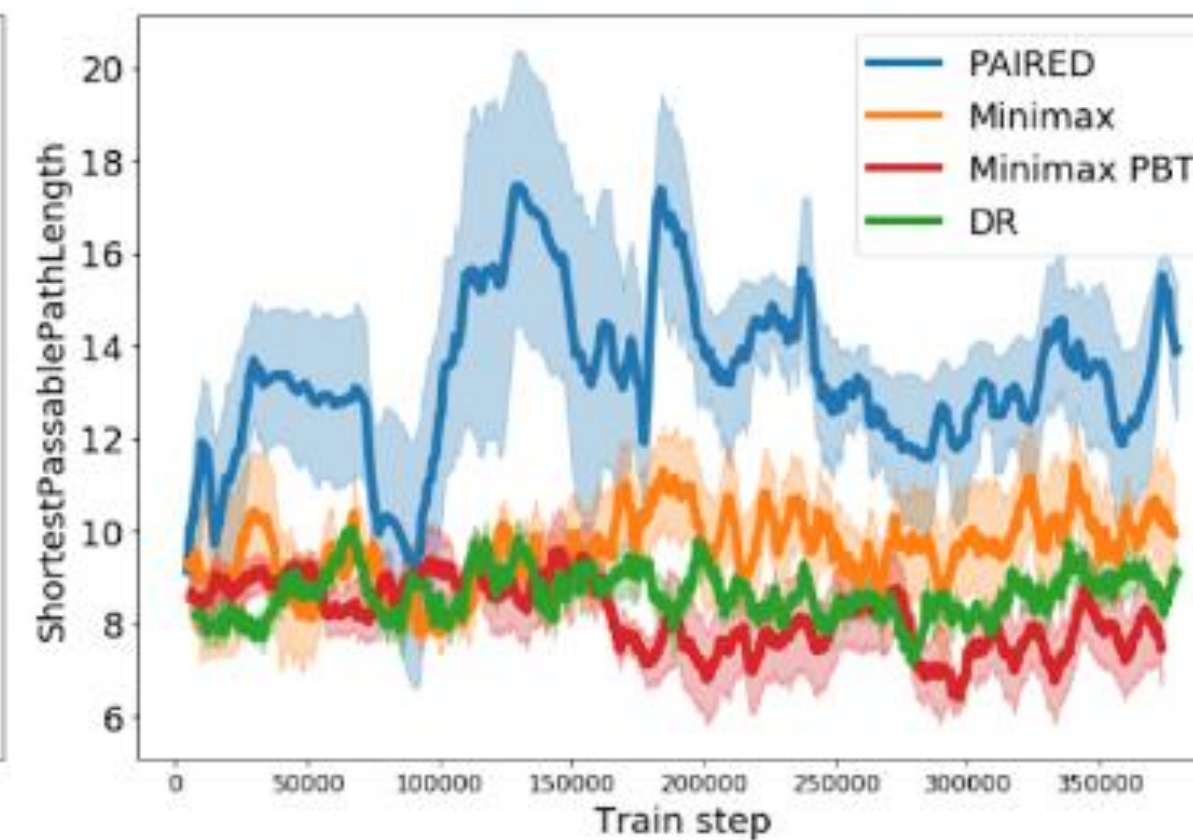
Generated environments



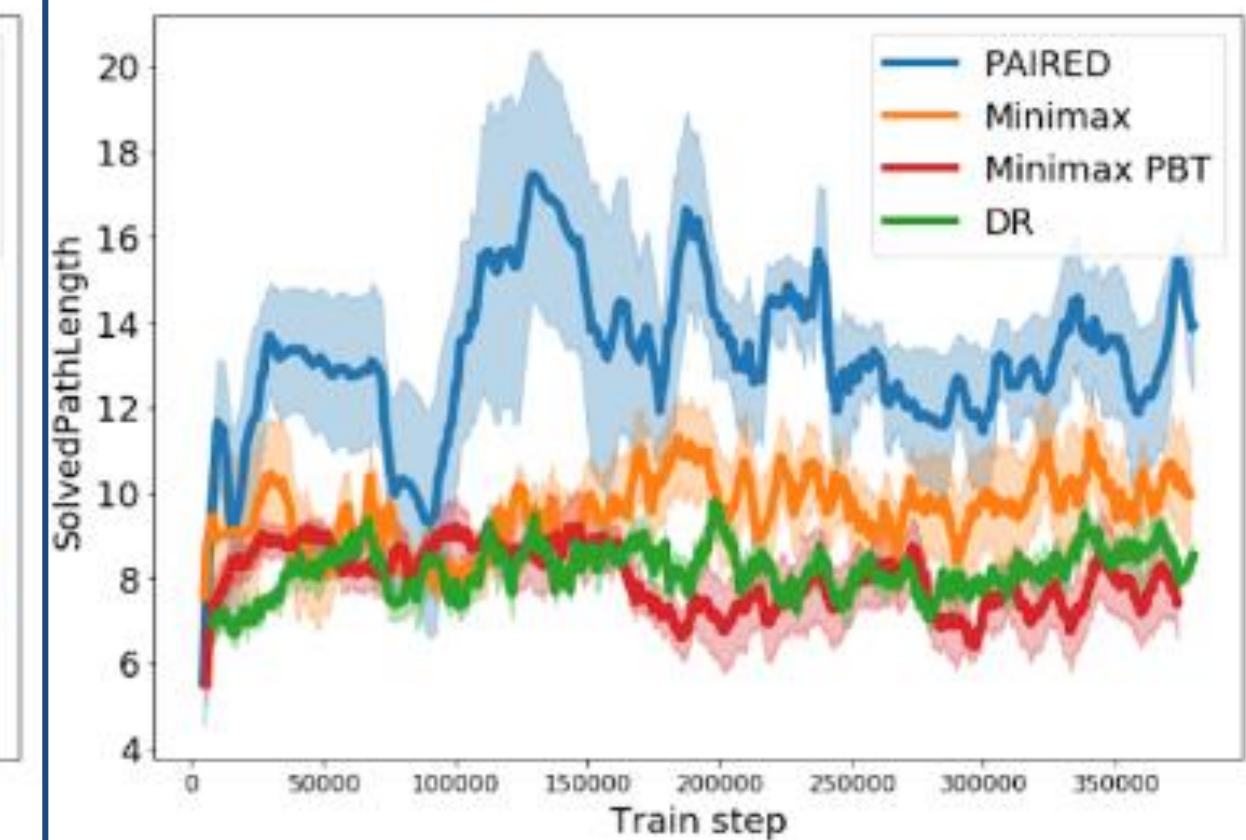
(a) Number of blocks



(b) Distance to goal



(c) Passable path length



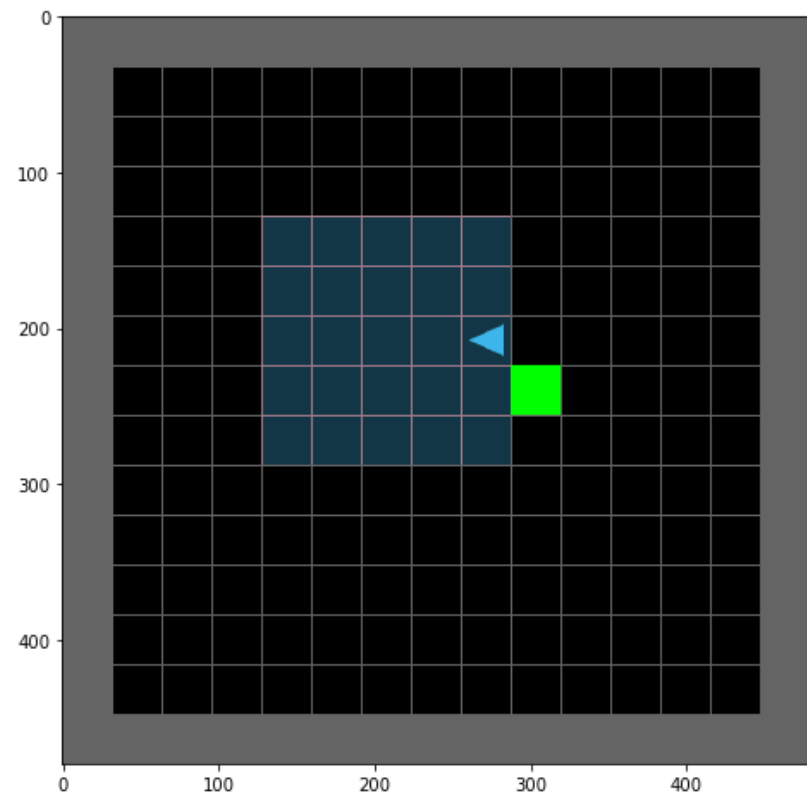
(d) Solved path length

Agent learning

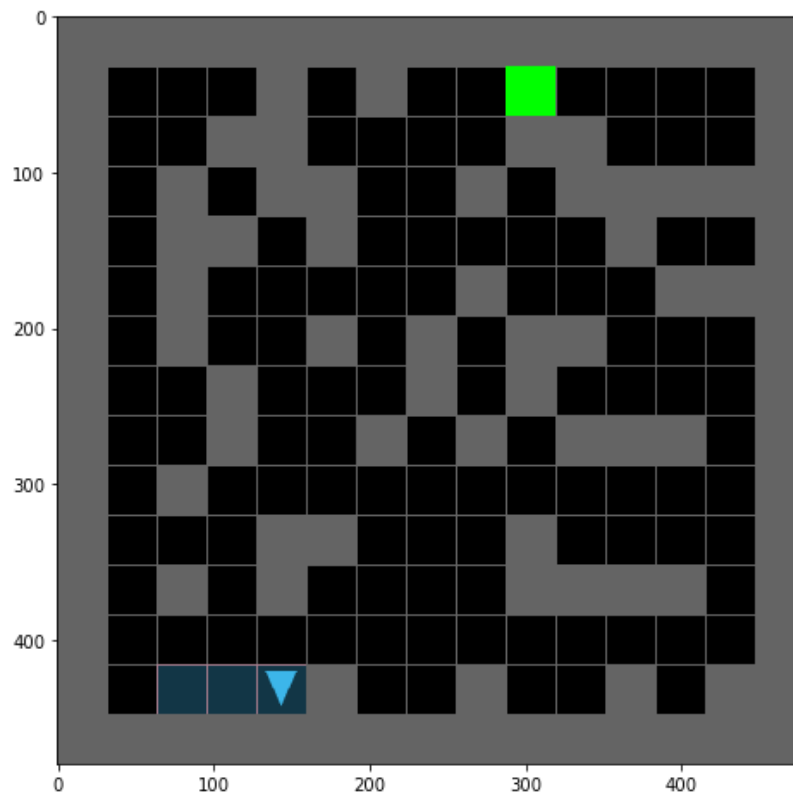
- Shortest path length of generated & solved mazes increases although agents were never trained on this objective

Zero-shot transfer

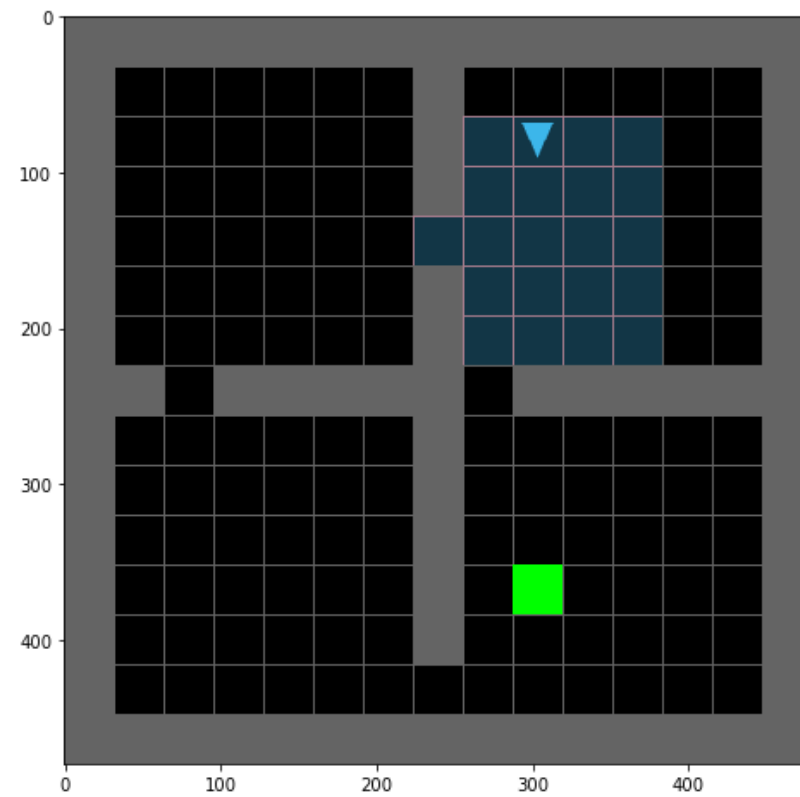
Empty



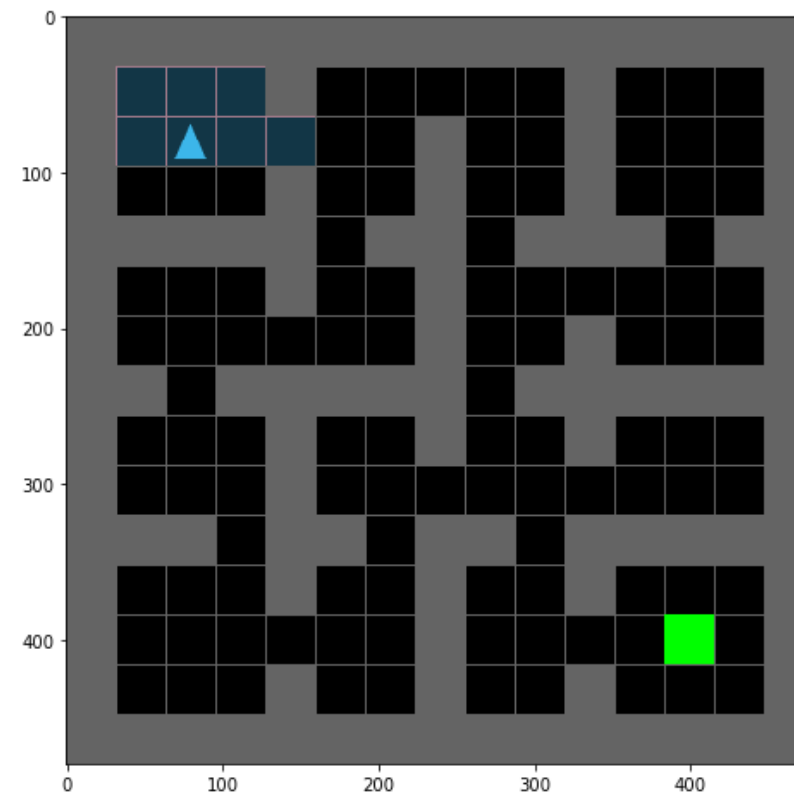
50 Blocks



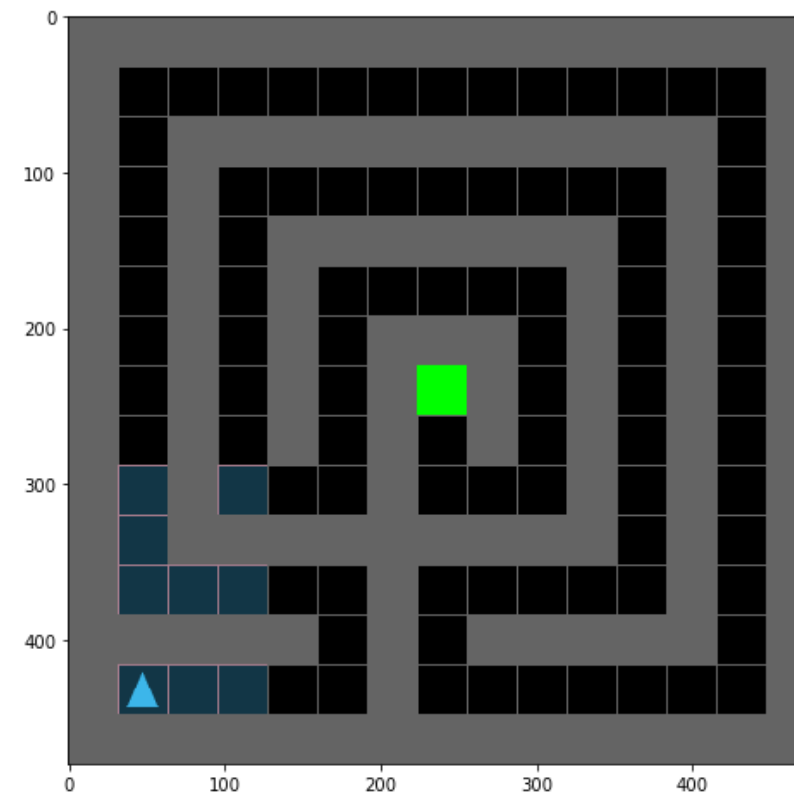
4 Rooms



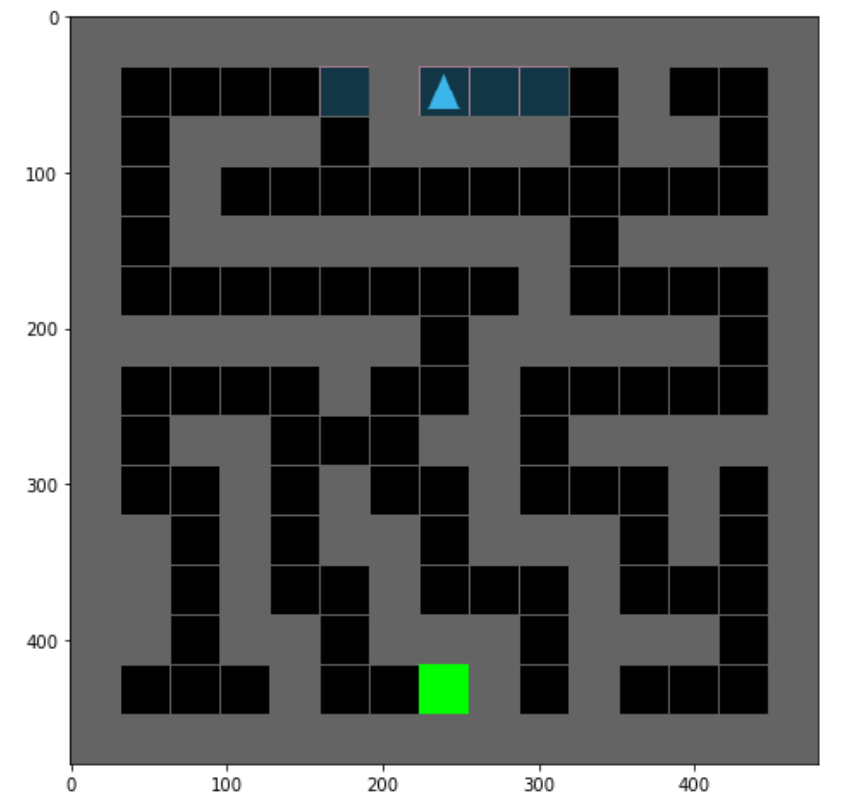
16 Rooms



Labyrinth

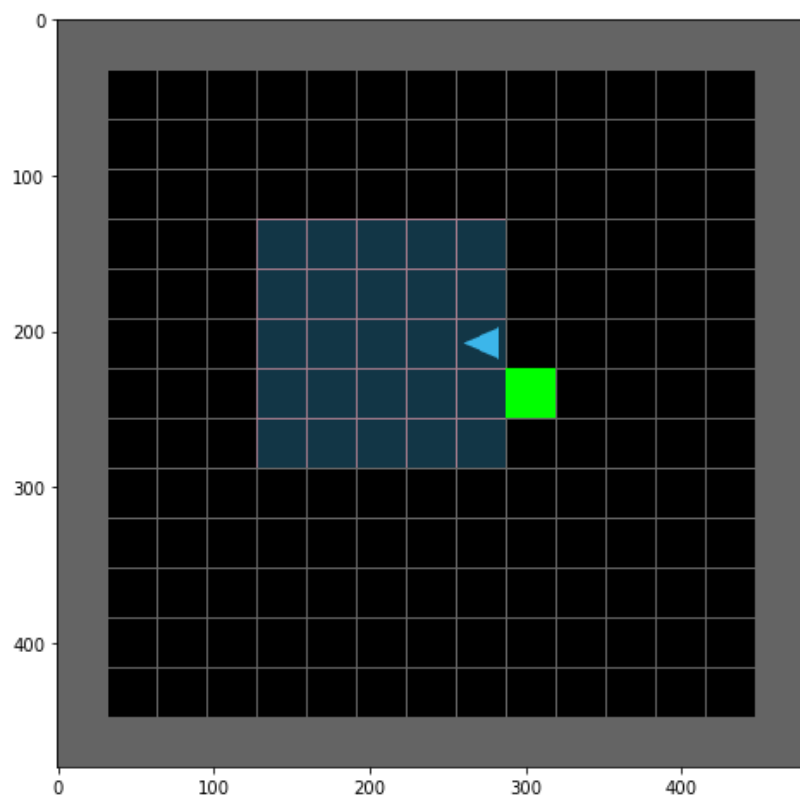


Maze

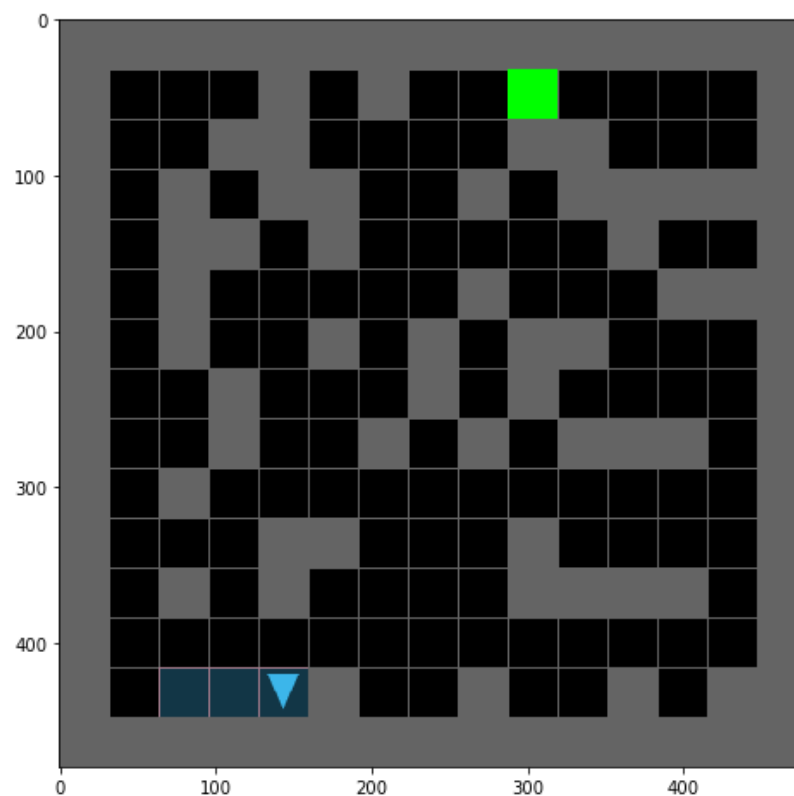


Zero-shot transfer

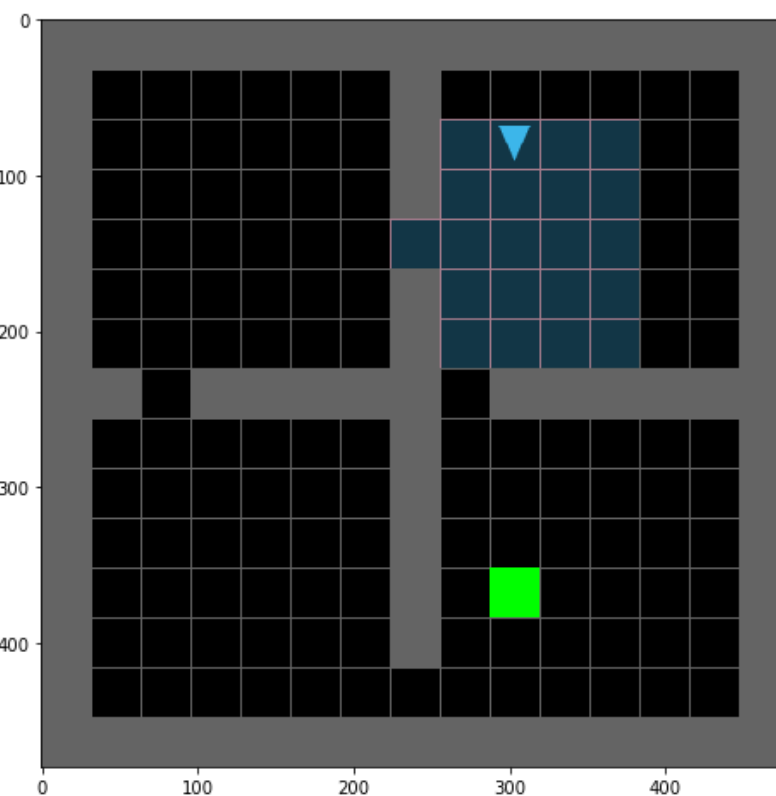
Empty



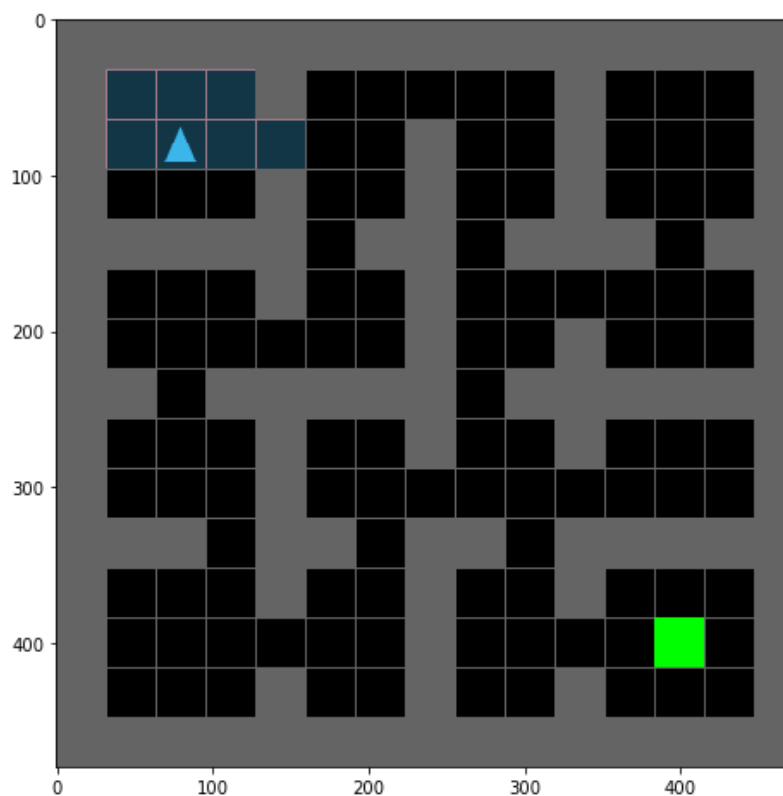
50 Blocks



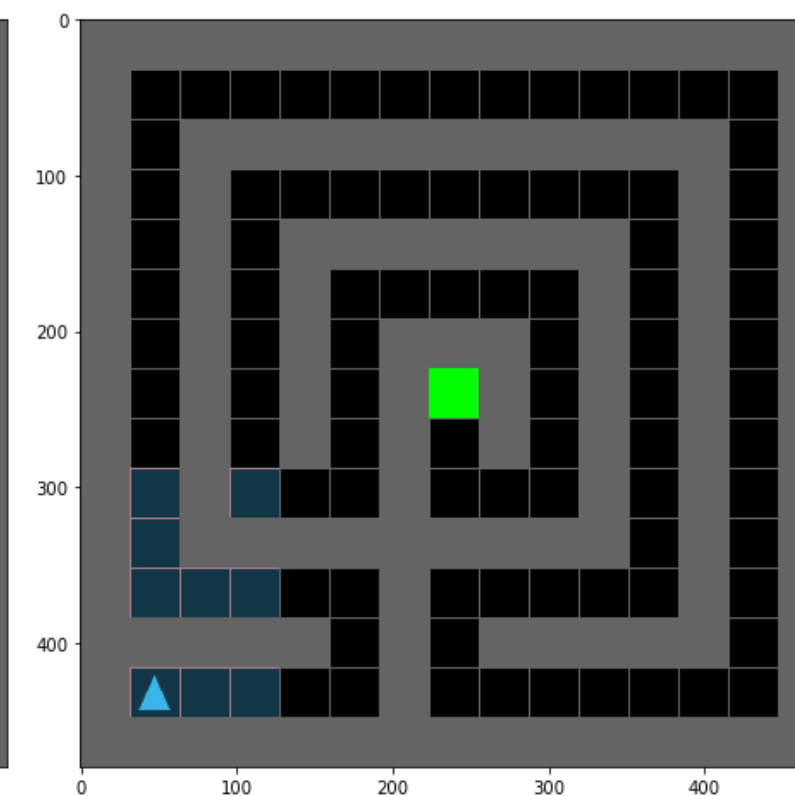
4 Rooms



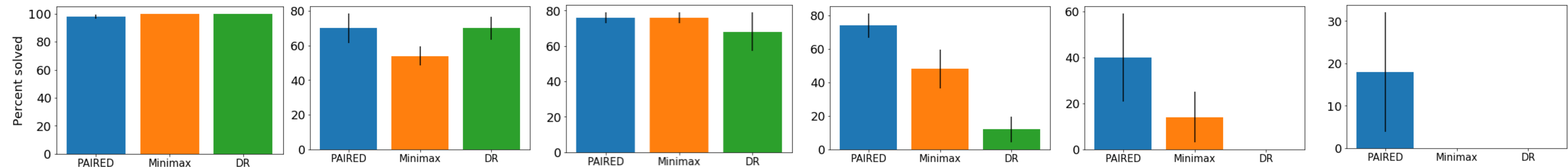
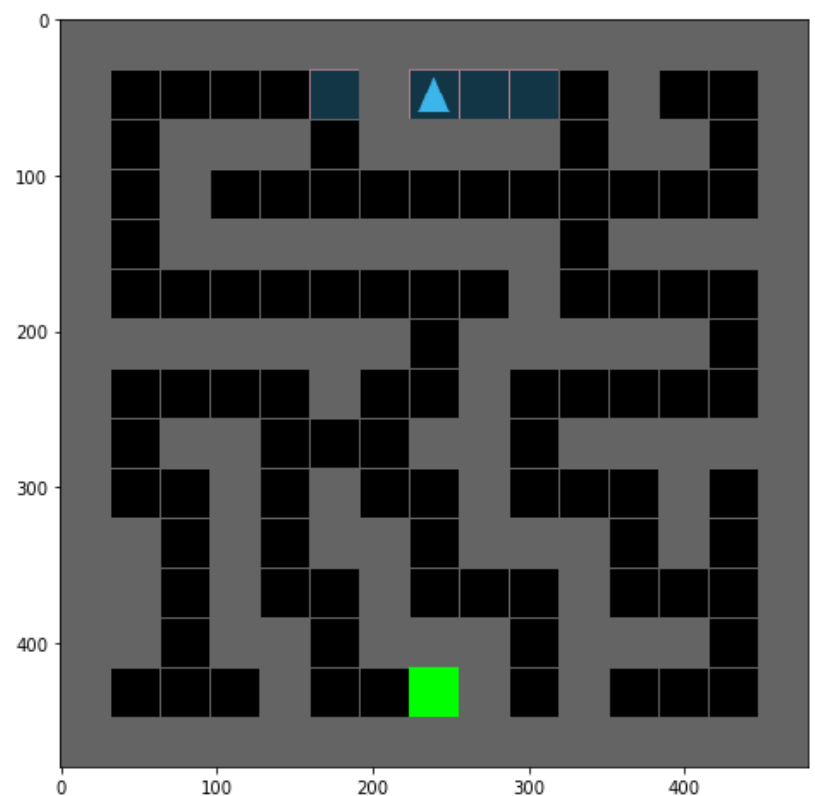
16 Rooms



Labyrinth



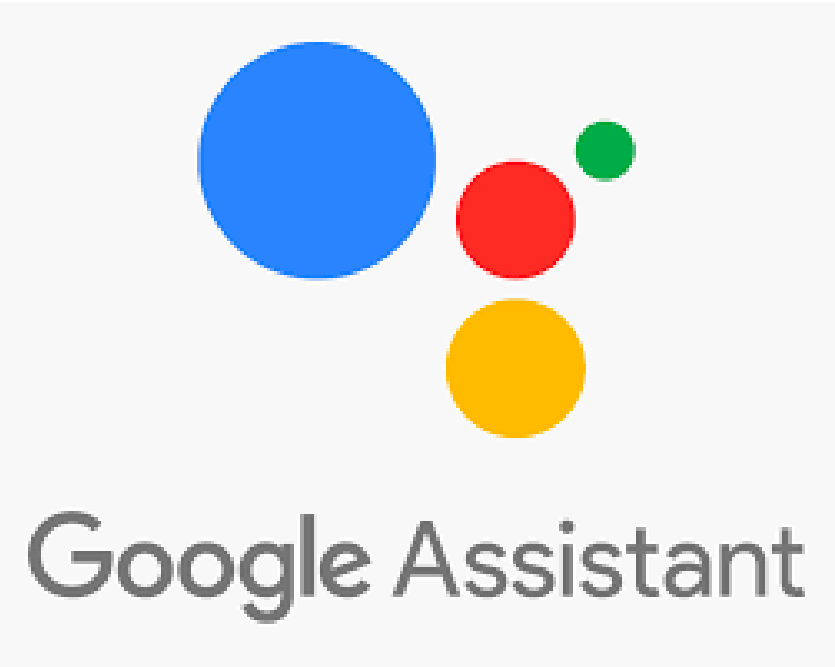
Maze



PAIRED
Minimax
Domain Randomization


See videos at <http://bit.ly/pairedvids!>

Environment Generation for Web Navigation



Goal: agent that can automatically **navigate real web pages** to complete tasks for users

Your Cart



Logitech - MX Keys
Advanced Wireless
Illuminated Keyboard
- Black

Pickup at **Gilroy**

Ready for pickup in **1hr**

Best Buy Curbside Pickup is available in checkout

FREE Shipping to **95112**

Get it by **Thu, Aug 6**

Faster shipping options are also available in checkout

1

Remove

Save

\$99.99

2-Year Accidental Geek Squad Replacement

★★★★★ 5.0 (4)

\$14.99

Add to Cart

Checkout

✓ Delivery

Arrives by Wed, Jul 22

Edit

Subtotal (1 item)

\$529.00

Delivery

Free

Est. taxes & fees
(Based on 48104)

\$31.74

Est. total

\$560.74

See item details

+

2

Enter delivery address

Required field *

First name*

Street address*

Last name*

Apt, suite, etc (optional)

Phone number* (in case of delivery questions)

City*

Ann Arbor

Email address for order notification*

State*

Michigan

ZIP code*

48104

☒ Email me about hot items, great savings, and more.

Continue

3

Enter payment method

Leave feedback

© 2020 Walmart Inc.

Do not sell my personal information

Request my personal information

DICK'S

CONTACT INFO

First Name

Last Name

Email

Phone (111) 111-1111

Apply my ScoreCard

BILLING & SHIPPING ADDRESS

☒ My Billing and Shipping are the same

Street Address

ZIP Code Only (EX: 12345)

DELIVERY METHODS

Choose a delivery method for your product(s):

Nike Men's Air Max Oketo Shoes

Qty: 1 | White/Red, 10.0, Medium

Est. Delivery: Thu 3/12 - Mon 3/16

\$ 54.99

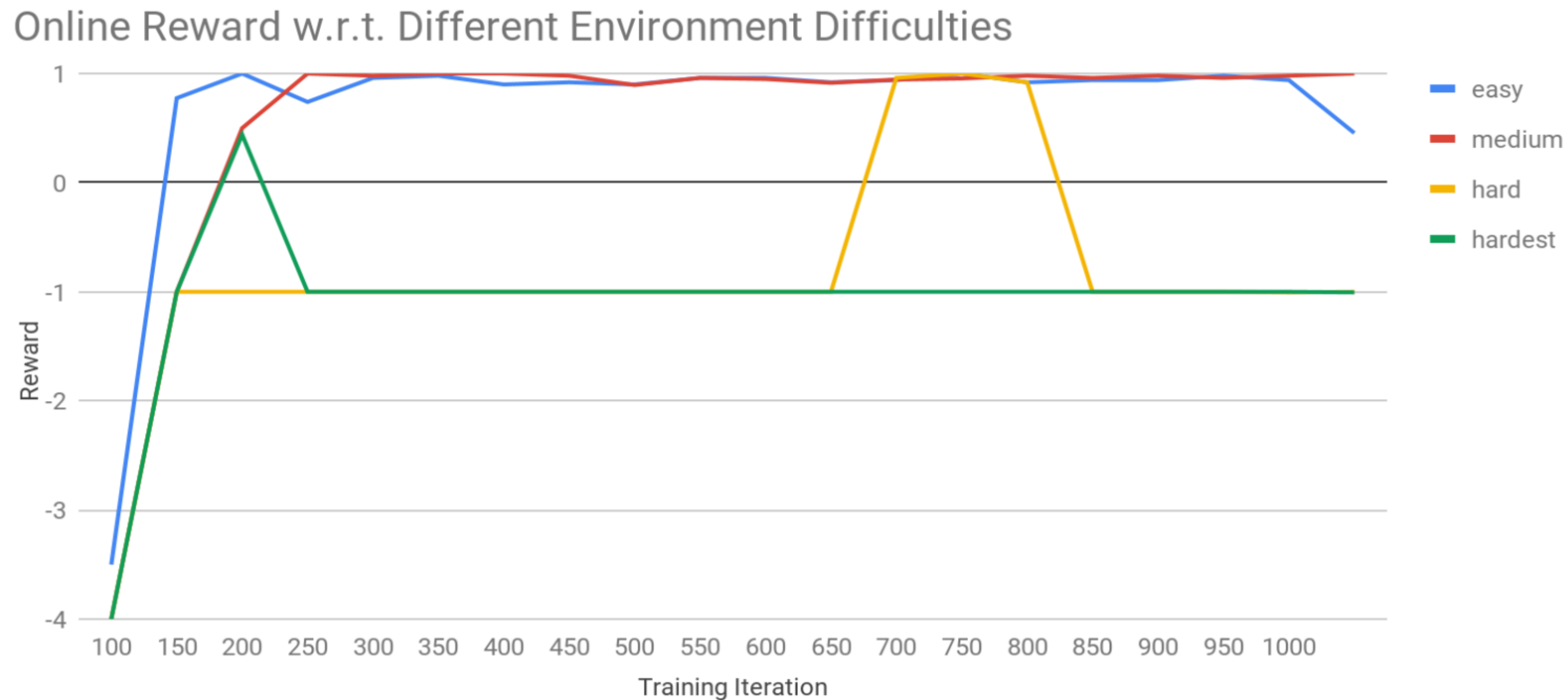
Standard

\$ 5.99

Environment Generation for Zero-Shot Compositional Reinforcement Learning. Gur, Jaques, Miao, Choi, Malta, Tiwari, Lee, Faust (2021). Neural Information Processing Systems (NeurIPS).

27

Web navigation



- **Issue:** initial approach involved **hand-programming** curriculum of easy, medium, hard sites
- **Idea:** PAired adversary generates a curriculum of websites

Web navigation - generated environments

Number of passengers

From

Continue

Deal of the Day
Gaming workstation
[Get it today!](#)

Login and Checkout

Payment
☐ Credit Card
☐ Debit Card

Continue

(a) Early training

Address

Continue

Last Name

To

First Name

Continue

(b) Mid training

To

Last Name

First Name

Address

Full name

Payment
☐ Credit Card
☐ Debit Card

From

Continue

(c) Late training

HOME

Username

Password

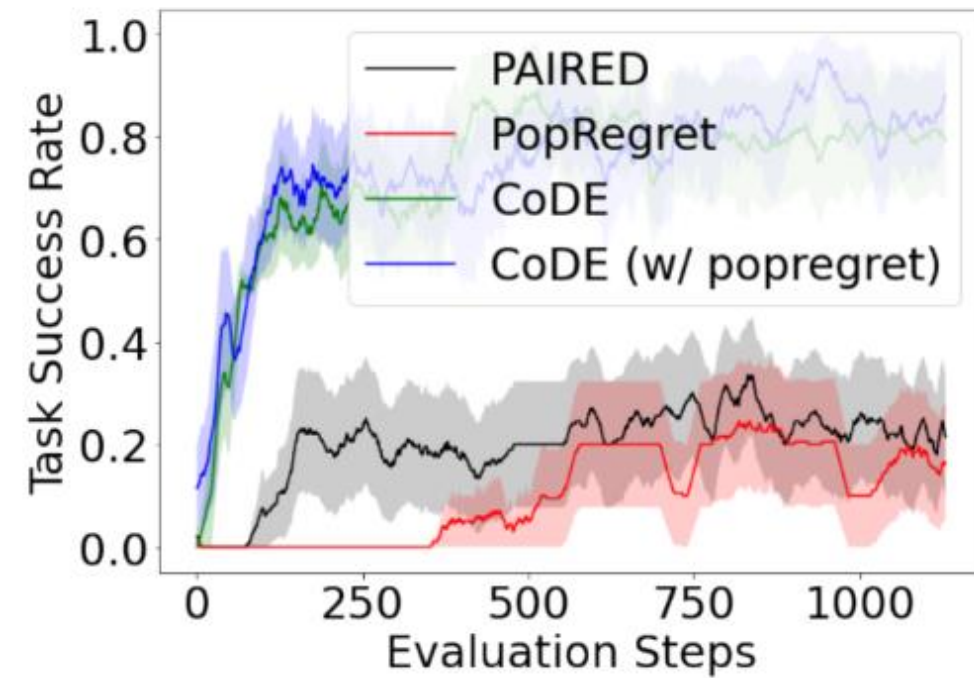
☐ Remember me
☐ Stay logged in
Enter Captcha

[Forgot user name.](#)
[Forgot password.](#)

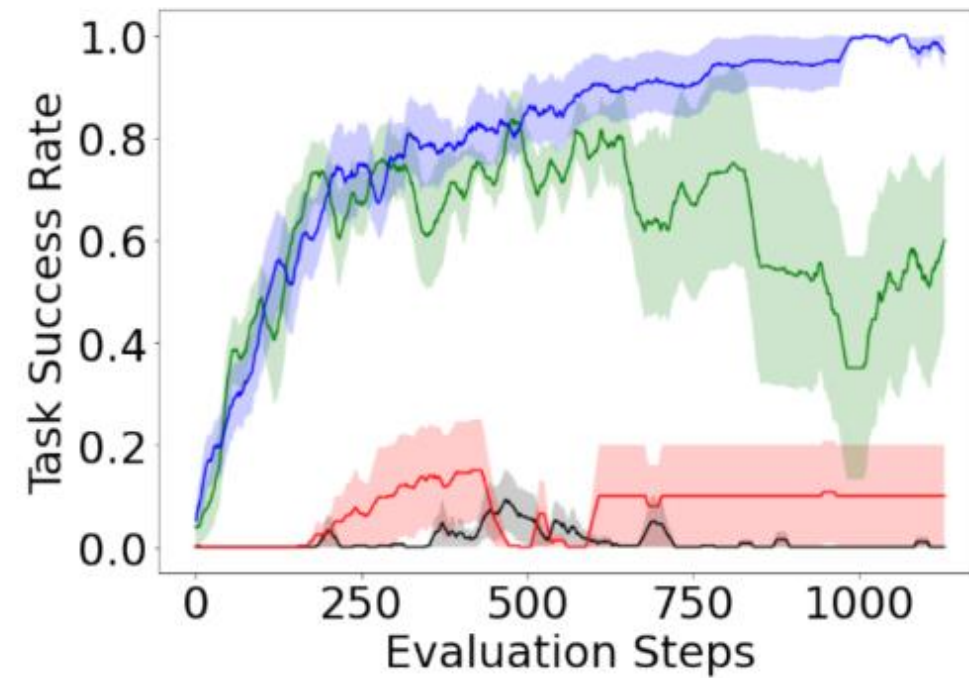
Continue

(d) Test

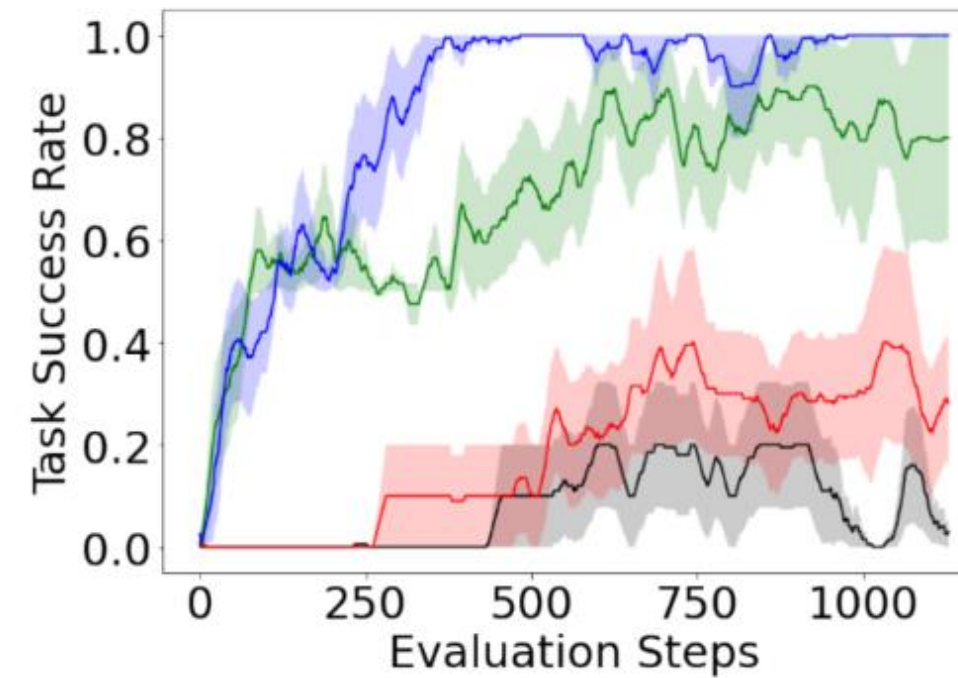
Web navigation - task success



(a) Payment

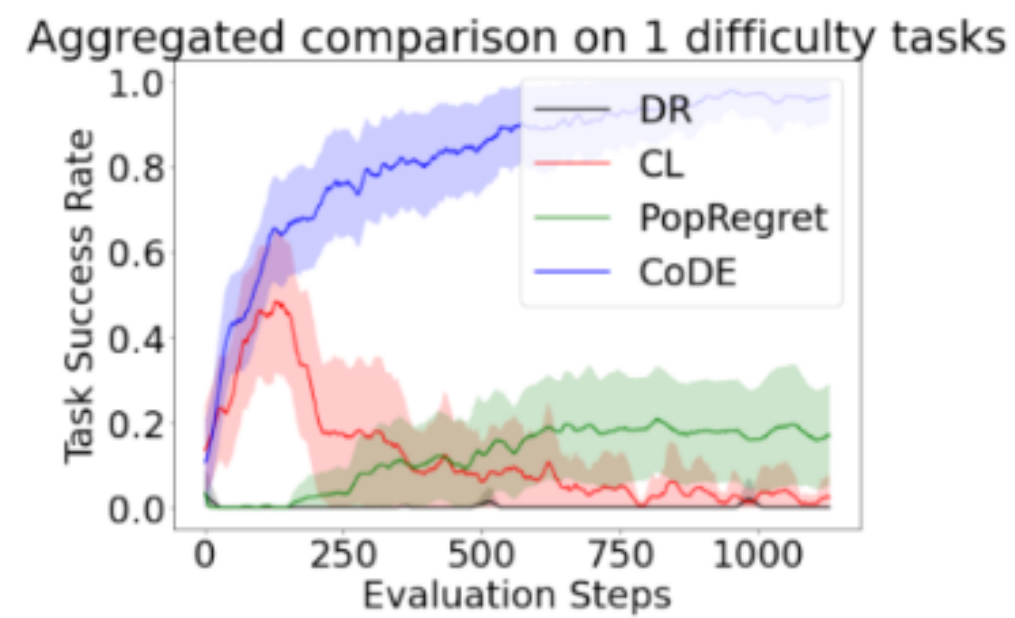


(b) Shopping

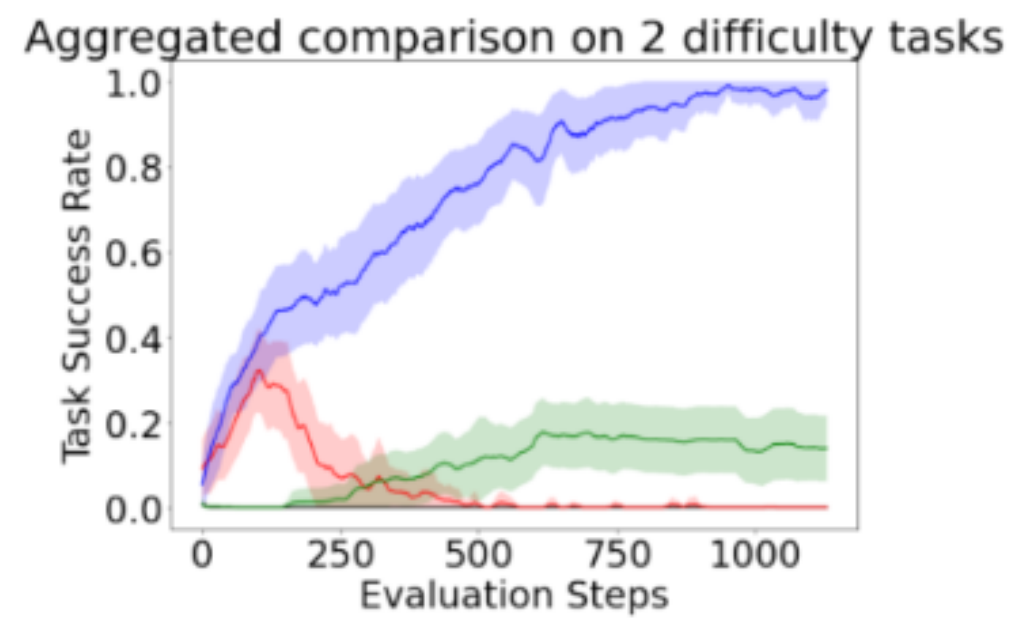


(c) Flight booking

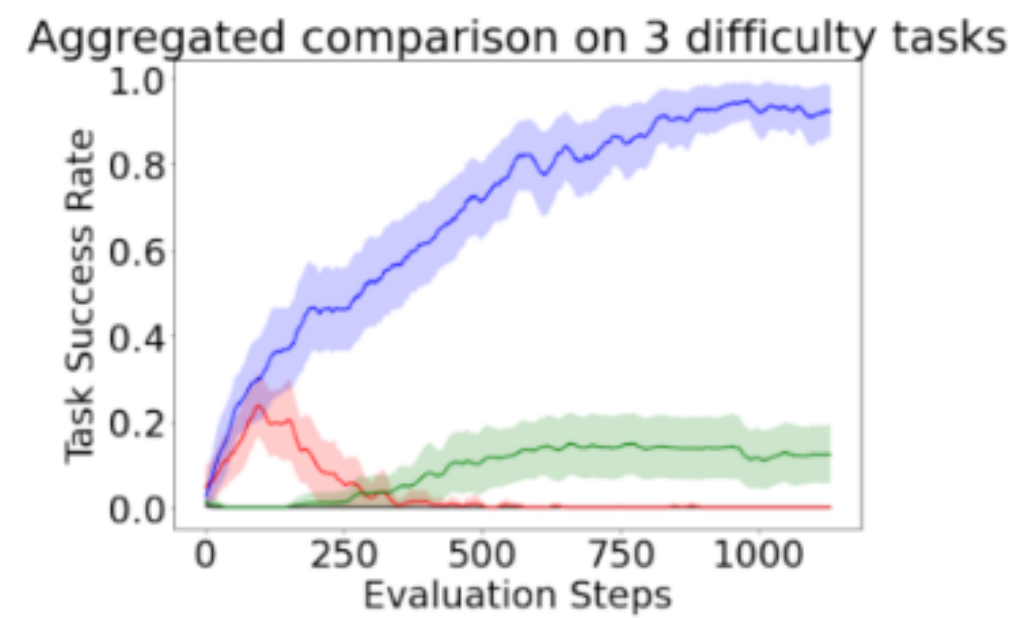
Deployed to real websites
• **93%+ accuracy on initial tasks**



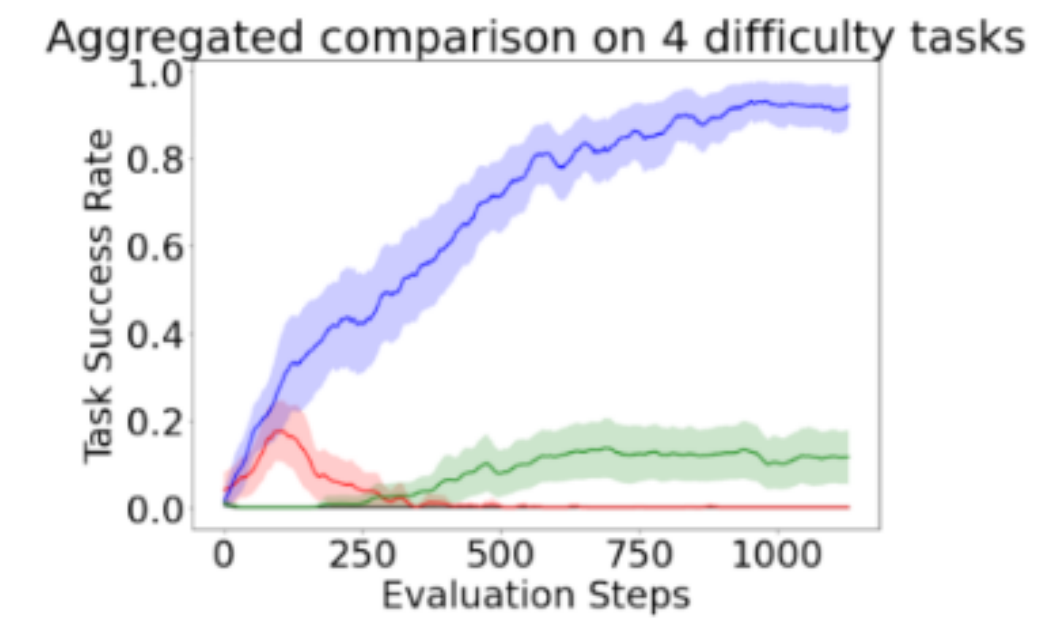
(a) Difficulty level 1



(b) Difficulty level 2



(c) Difficulty level 3



(d) Difficulty level 4

- **4x more successful than SOTA prior work**
- Reaches more than **95% task success** across all difficulty levels

Conclusions

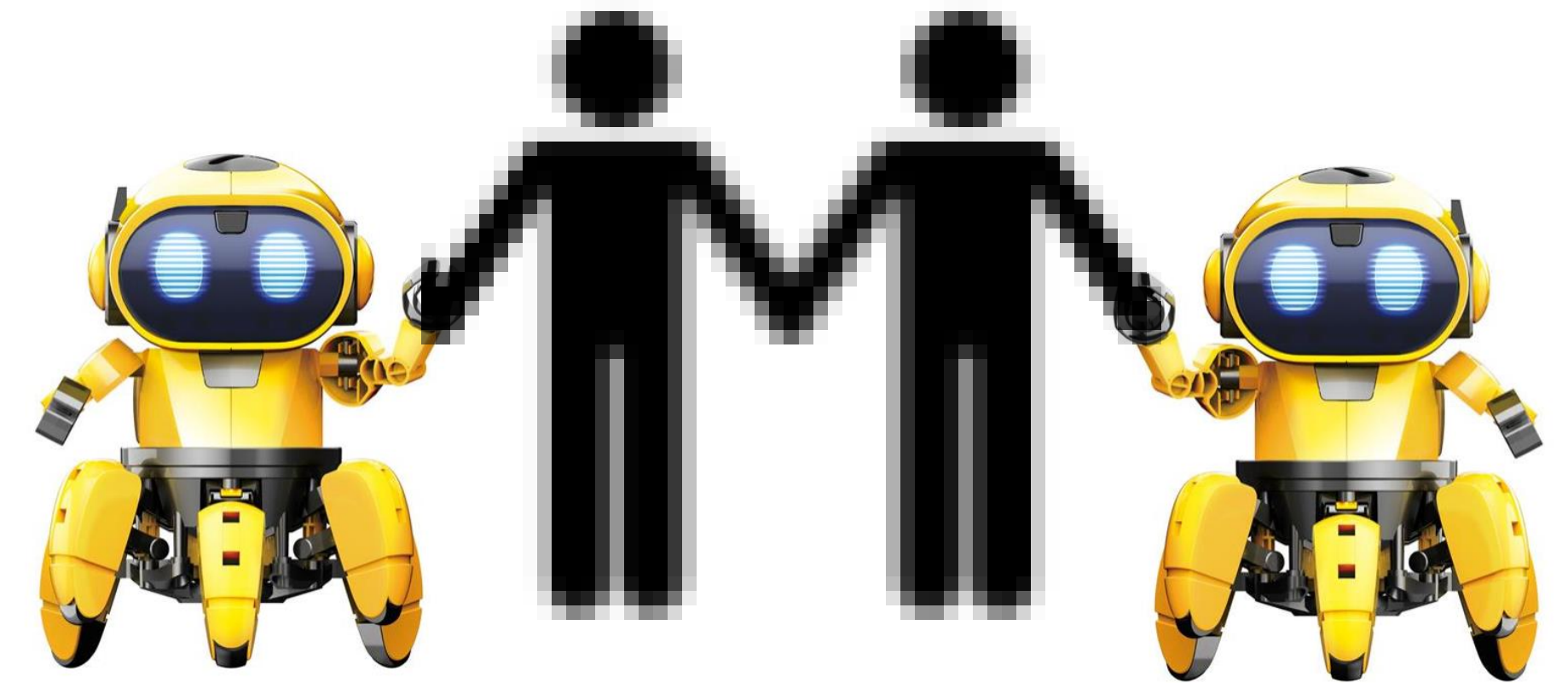
- **Multi-agent training** can be an effective tool for improving **learning** and **generalization** *even for single-agent settings*
- Advantages of **PAIRED**:
 - The adversary builds **feasible yet challenging** environments that make the protagonist more robust
 - **Regret** minimization creates an **automatic curriculum**

Outline

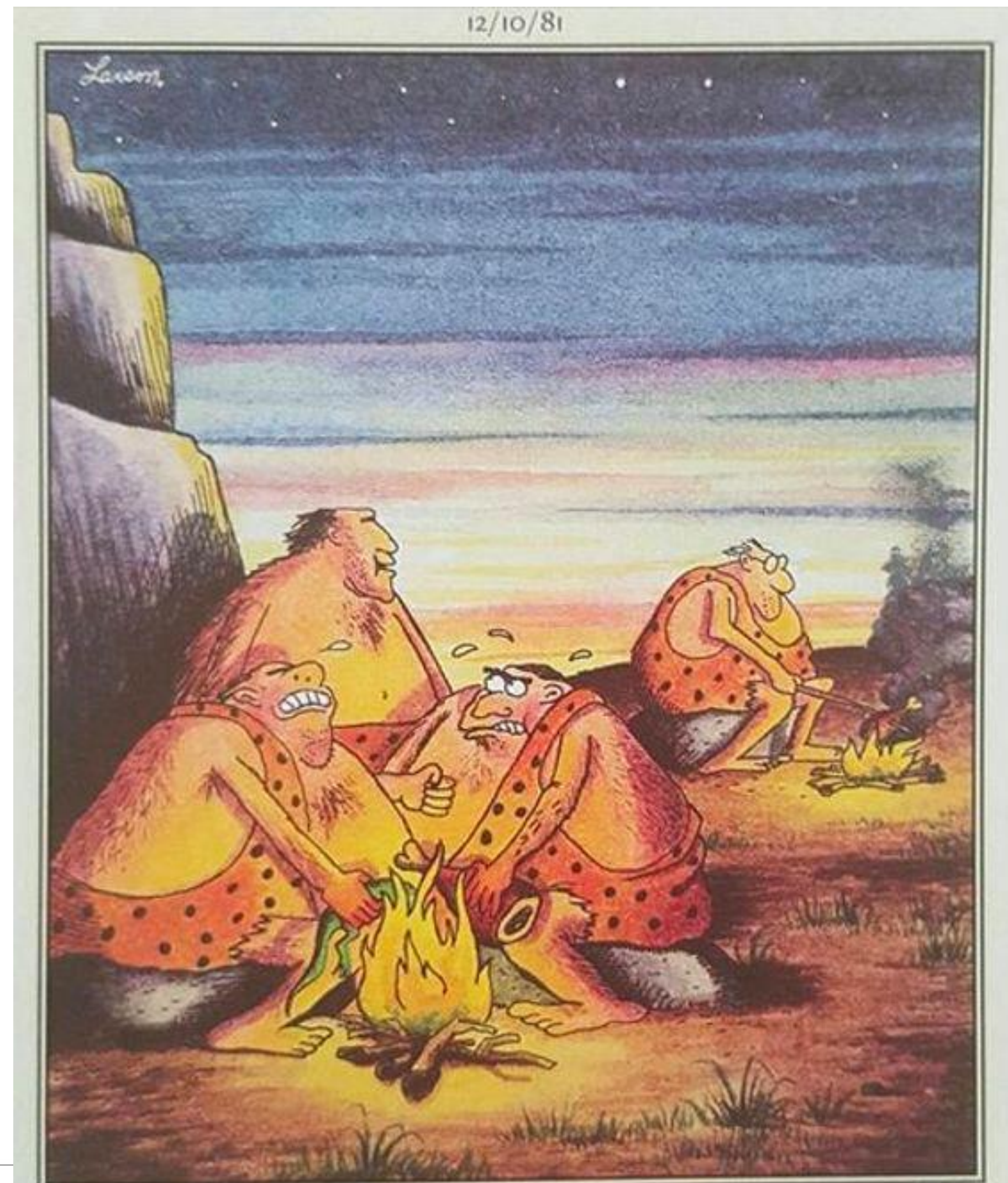
Multi-agent...

1. Emergent complexity

2. Social Learning



How social learning can accelerate learning



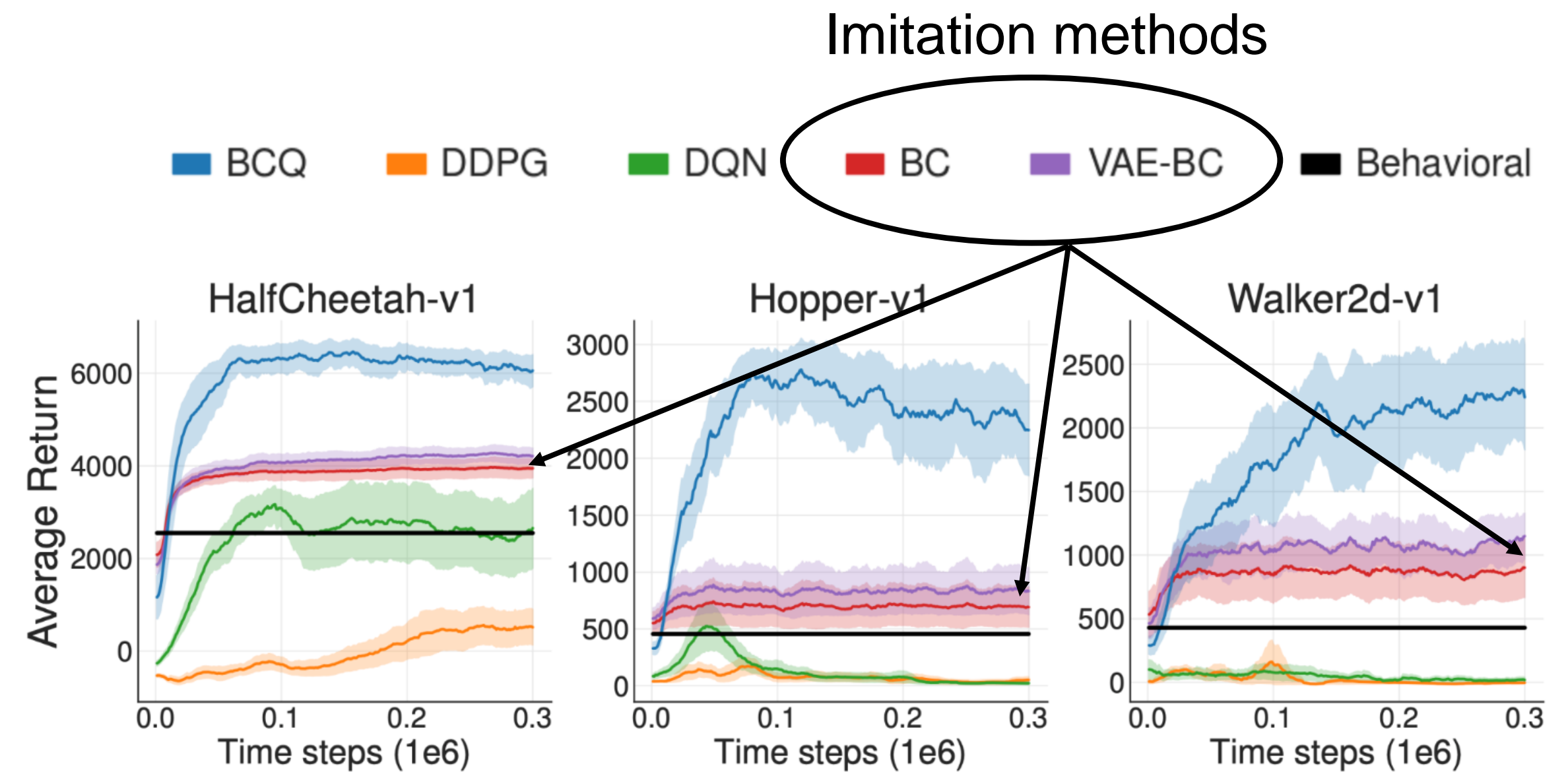
“Hey! Look what Zog do!”

Social learning enables you to
“stand on the shoulders of giants”

Why not just do imitation learning?

= supervised learning on $\langle s, a \rangle$ expert dataset

- Requires access to **specially curated** expert trajectories
- **Copy data exactly**; perform poorly if data comes from imperfect experts
- Resulting policies are brittle, **do not generalize**



(d) Imperfect demonstrations performance

From Fujimoto, Scott, David Meger, and Doina Precup. "Off-policy deep reinforcement learning without exploration." International Conference on Machine Learning. PMLR, 2019.

Social Learning with Multi-Agent RL

- Works in more **naturalistic** settings. When other agents:
 - Are **not motivated to teach you**
 - May or may not have **relevant expertise**
- **Do not have to copy exactly**
- Social learners can **learn how to acquire** info from other agents
 - **Generalize/adapt** to new environments



<https://www.sfgate.com/traffic/>

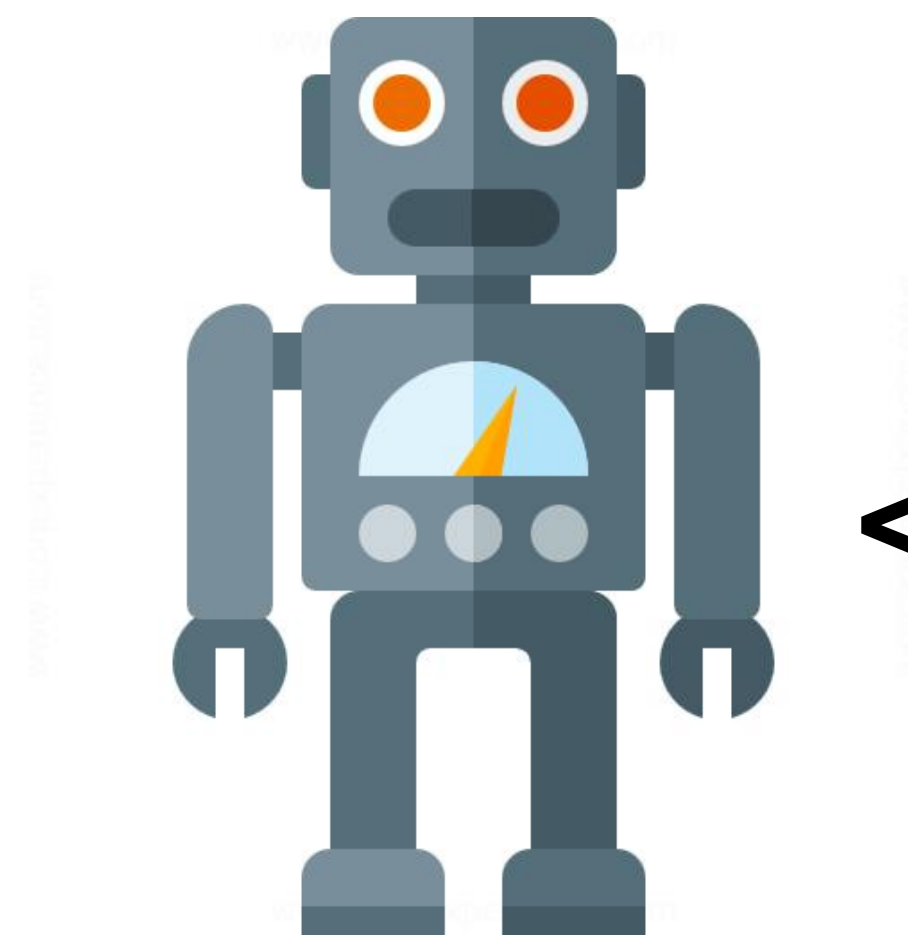
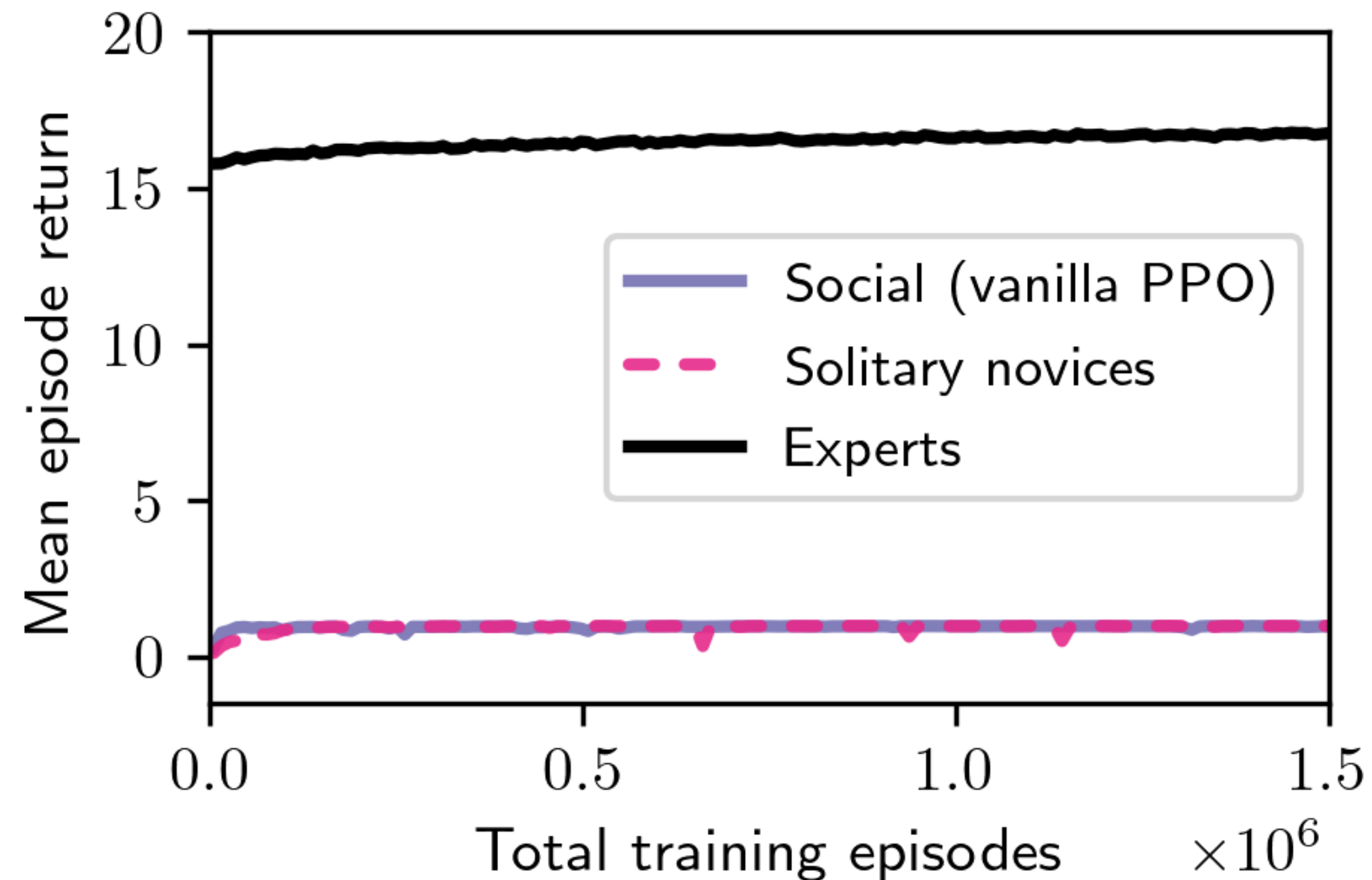
Landolfi, N. C., & Dragan, A. D. (2018, October). Social cohesion in autonomous driving. In 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (pp. 8118-8125). IEEE.

Learning Social Learning

Multi-agent env with partial observability and no privileged access to other agents' states or actions (unlike imitation learning)



**Kamal
Ndousse**

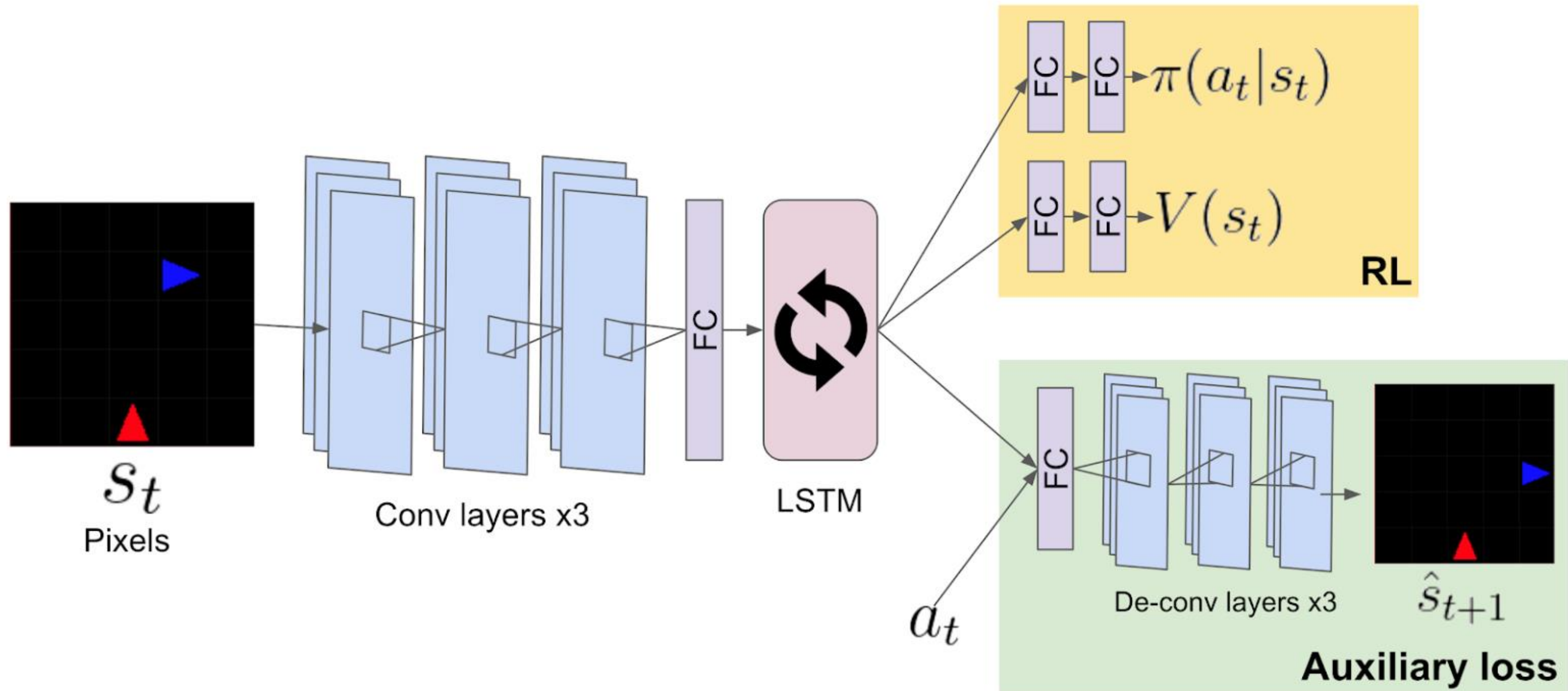


<

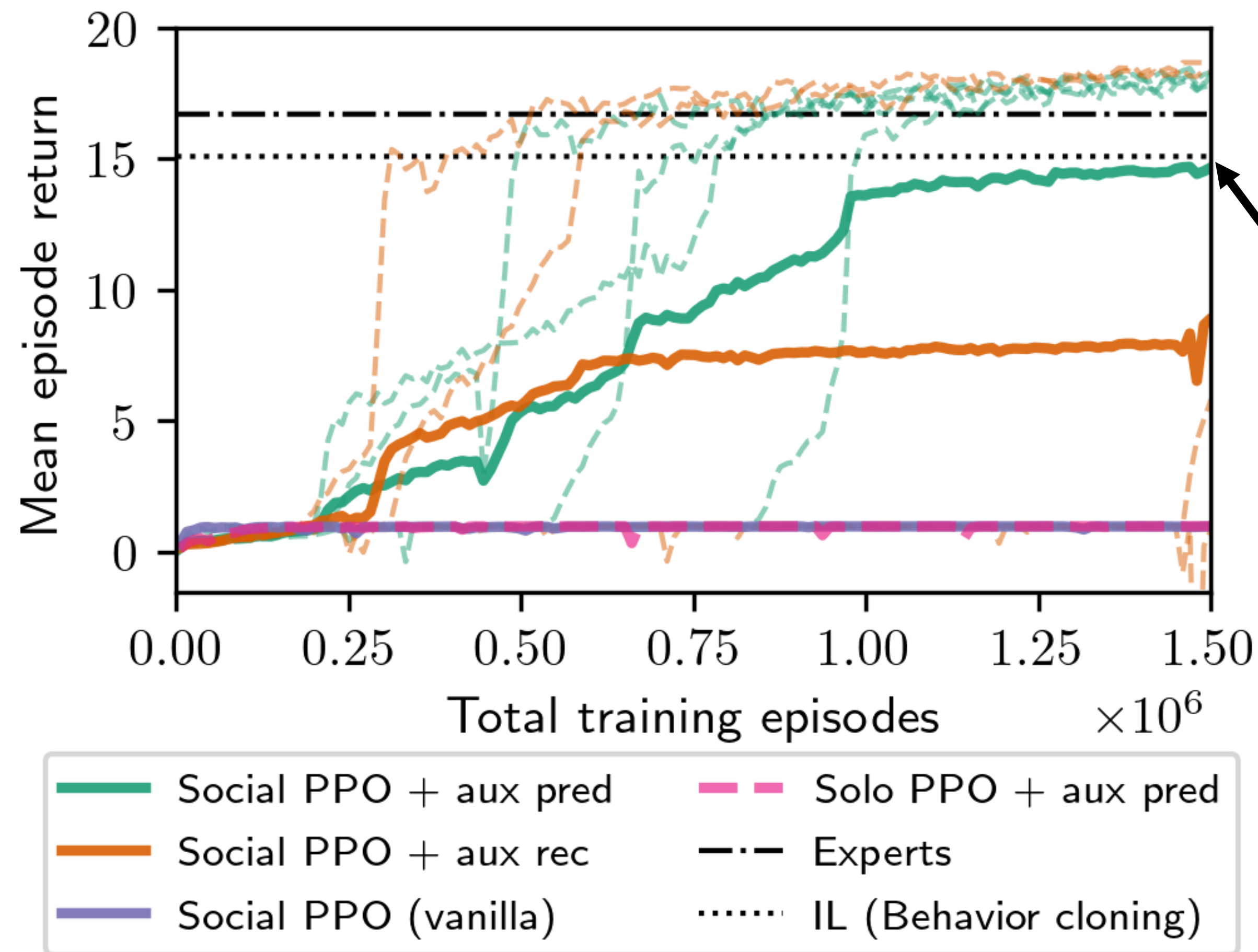


Vanilla RL agents fail to learn from experts in their environment

Augment agents with auxiliary prediction loss



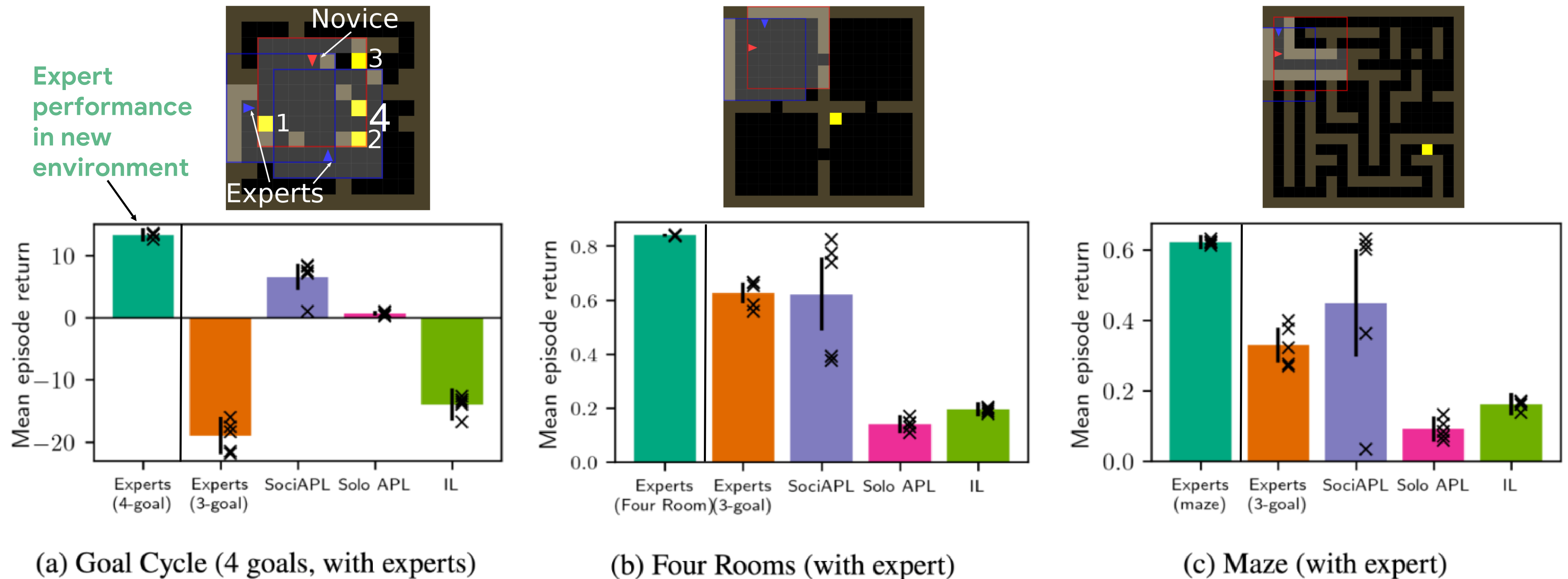
Agents are able to use social learning



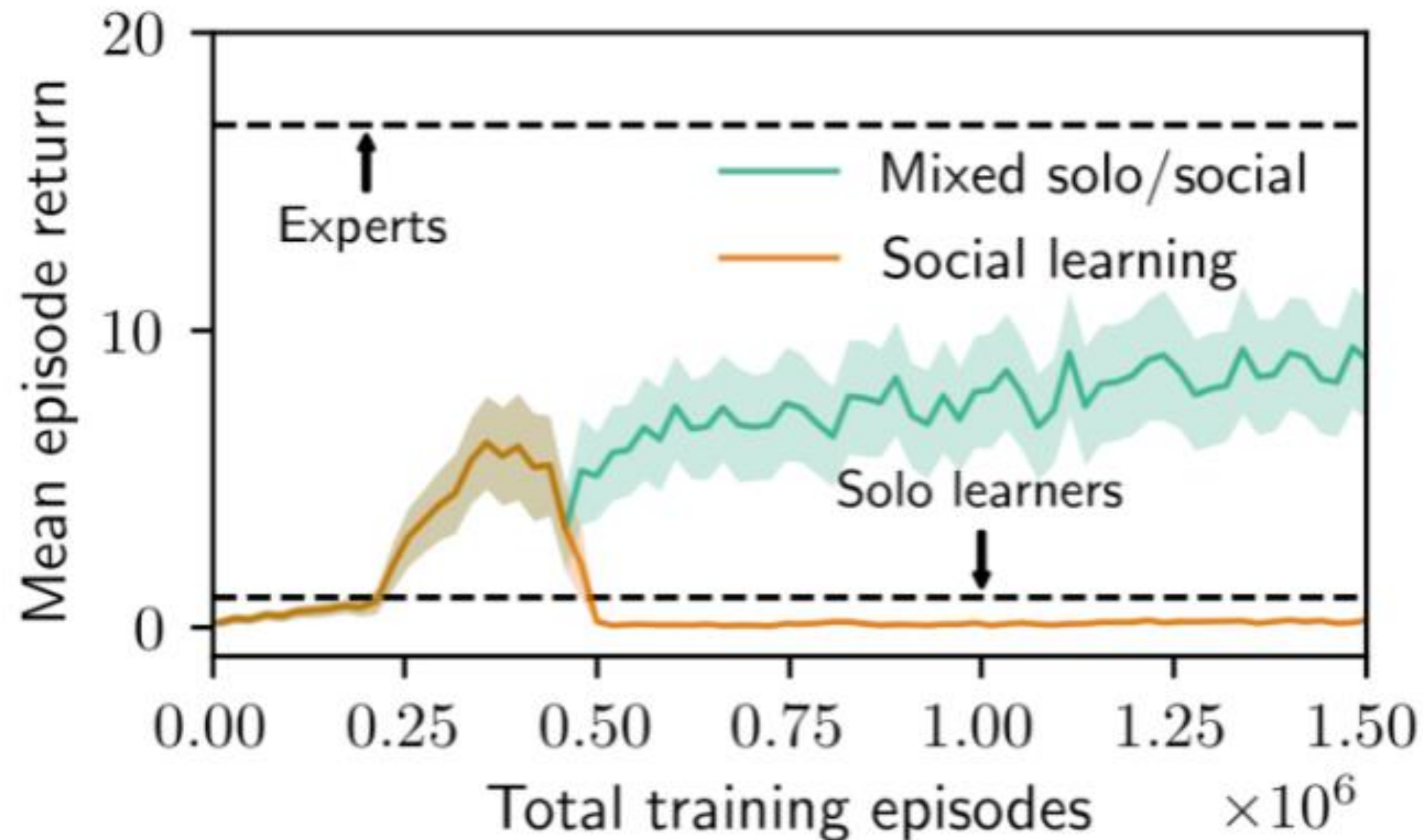
Same as imitation learning (IL) but no access to $\langle s, a, s' \rangle$ of other agents

Social Learners generalize to new envs

Social learners generalize better than **RL experts** or **imitation learners**



Social Learning can benefit performance when alone



(b) Transfer to solo 3-Goal

- Social learners discover skills that enable them to perform better when **alone**
 - If trained with a mix of solo & social episodes
- Outperform agents that were always trained alone

Learning social learning helps with



**Rapidly adapting
to new
environments**



**Learning complex
behavior**

Improving Social Learning with PsiPhi-Learning



Angelos
Filos

No-reward
demos
<s,a,s'>





$$\mathbb{E}^{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t, a_t \right] = Q^{\pi}(s_t, a_t) = \Psi^{\pi}(s_t, a_t)^T \mathbf{w}$$

$$\Psi^{\pi}(s_t, a_t) = \mathbb{E}^{\pi} \left[\sum_{k=0}^{\infty} \gamma^k \Phi(s_{t+k}, a_{t+k}) \right]$$

Future accumulated state features

Preference vector


$$\begin{bmatrix} 15 \\ 2 \end{bmatrix}$$


$$\begin{bmatrix} 1 & -1 \end{bmatrix} = \begin{bmatrix} 15(1) + 2(-1) \end{bmatrix} = 1.3$$

Successor Features^{1,2}

¹Peter Dayan. Improving generalization for temporal difference learning: The successor representation. Neural Computation, 5(4):613–624, 1993.

²Barreto, A., Dabney, W., Munos, R., Hunt, J. J., Schaul, T., Van Hasselt, H., & Silver, D. (2016). Successor features for transfer in reinforcement learning. arXiv preprint arXiv:1606.05312.

Combine other agents' policies with ego-agent's preferences

- Use Generalized Policy Improvement (GPI)²

Ego-agent

$$Q^{ego}(s, a) = \Psi^{ego}(s, a)^T \mathbf{w}^{ego}$$

Other agents

$$Q^1(s, a) = \Psi^1(s, a)^T \mathbf{w}^1$$

$$Q^2(s, a) = \Psi^2(s, a)^T \mathbf{w}^2$$

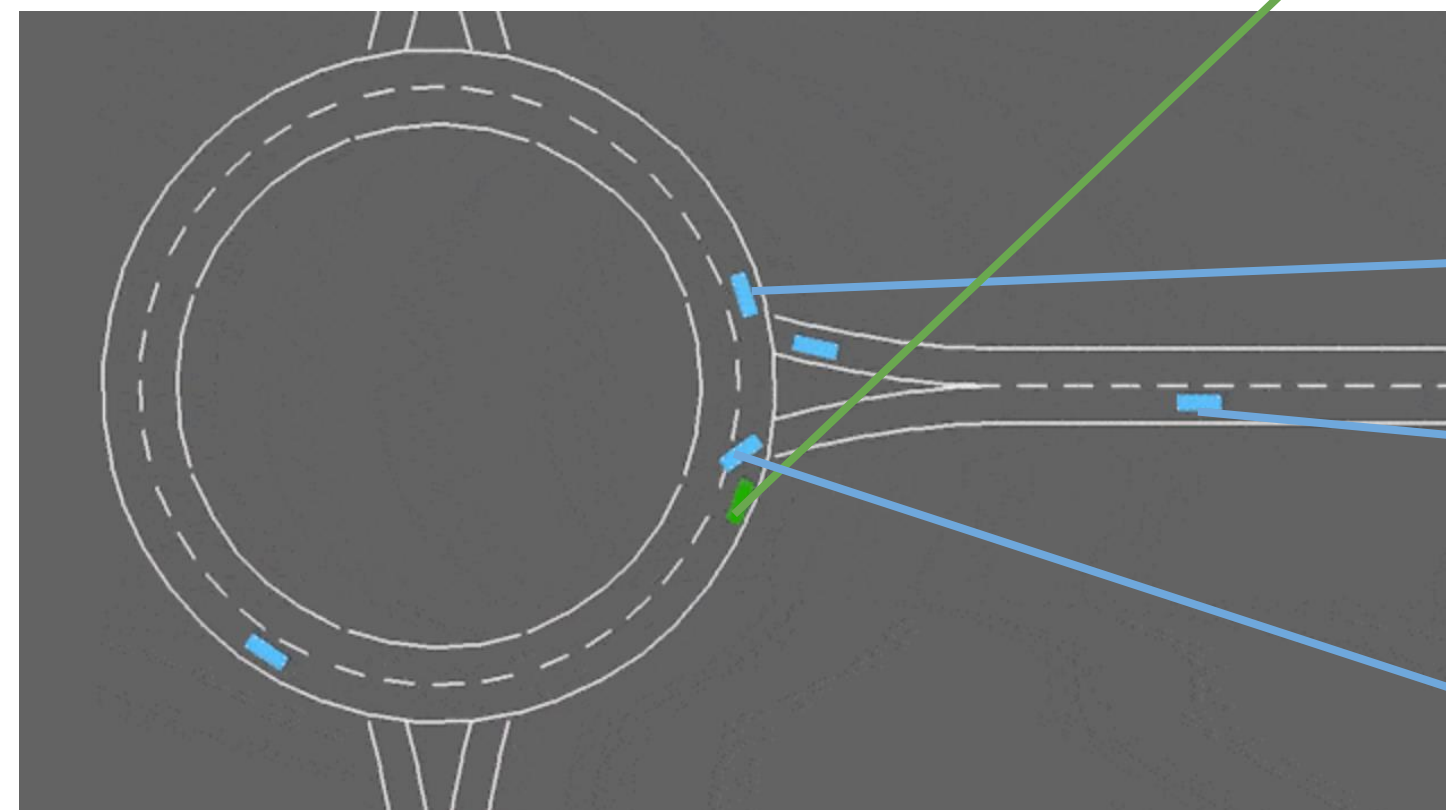
...

$$Q^k(s, a) = \Psi^k(s, a)^T \mathbf{w}^k$$

\mathbf{w}^{ego}

$$\pi^{ego}(s) =$$

$$\arg \max_a \max_k \Psi^k(s, a)^T \mathbf{w}^{ego}$$



Multi-agent environment
(autonomous driving)

Problem: why should other agents' Ψ be compatible with my \mathbf{w} ?

²Barreto, A., Dabney, W., Munos, R., Hunt, J. J., Schaul, T., Van Hasselt, H., & Silver, D. (2016). Successor features for transfer in reinforcement learning. arXiv preprint arXiv:1606.05312.

$\Psi\Phi$ -Learning

$$\begin{array}{ccc} \begin{array}{c} \text{Banana} \\ \text{Apple} \end{array} \begin{bmatrix} 15 \\ 2 \end{bmatrix} & \Psi^\pi(s_t, a_t) = \mathbb{E}^\pi \left[\sum_{k=0}^{\infty} \gamma^k \Phi(s_{t+k}, a_{t+k}) \right] & \begin{array}{c} \text{Banana} \\ \text{Apple} \end{array} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ \text{Policy-specific} & & \text{Agent-agnostic} \end{array}$$

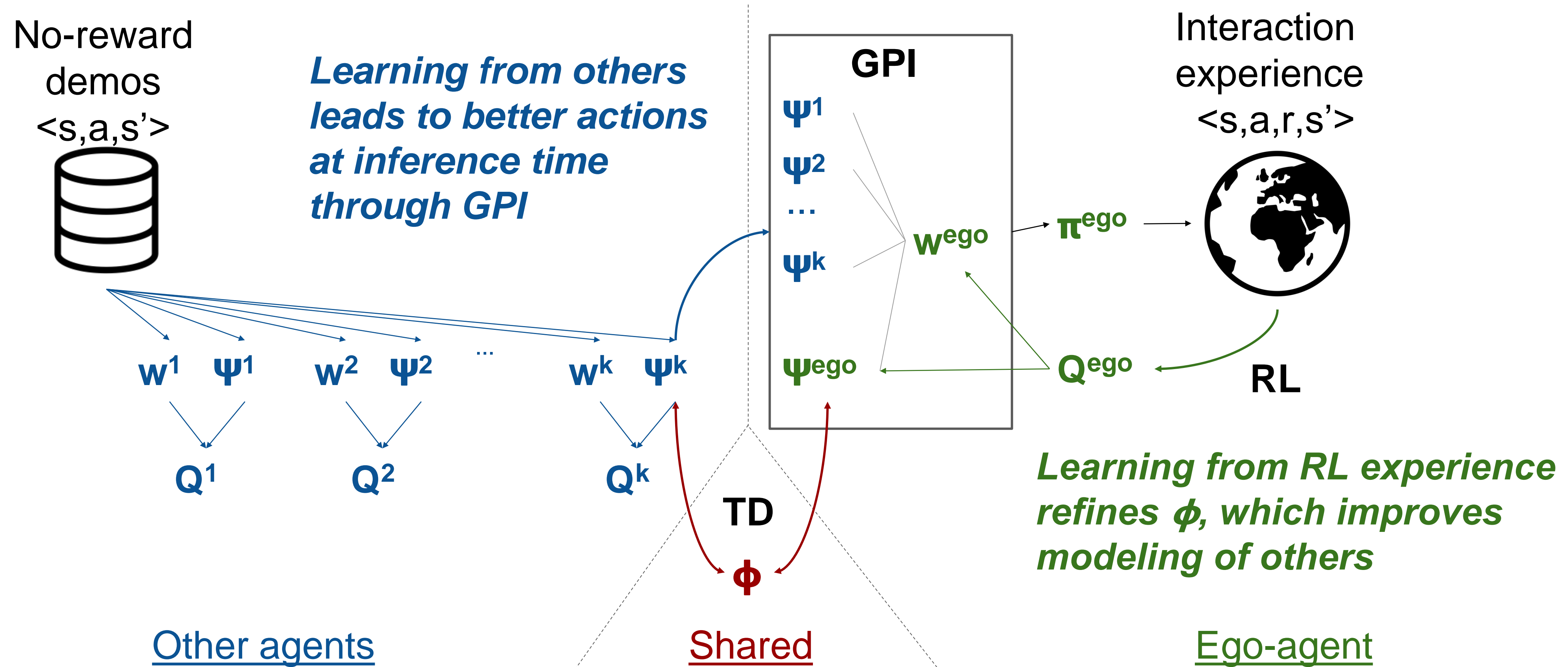
Solution: learn a shared feature representation Φ that explains all agents' behavior

- Each agent's Ψ^k should be built from the same Φ
- Learn Φ with **Inverse Temporal Difference** loss:

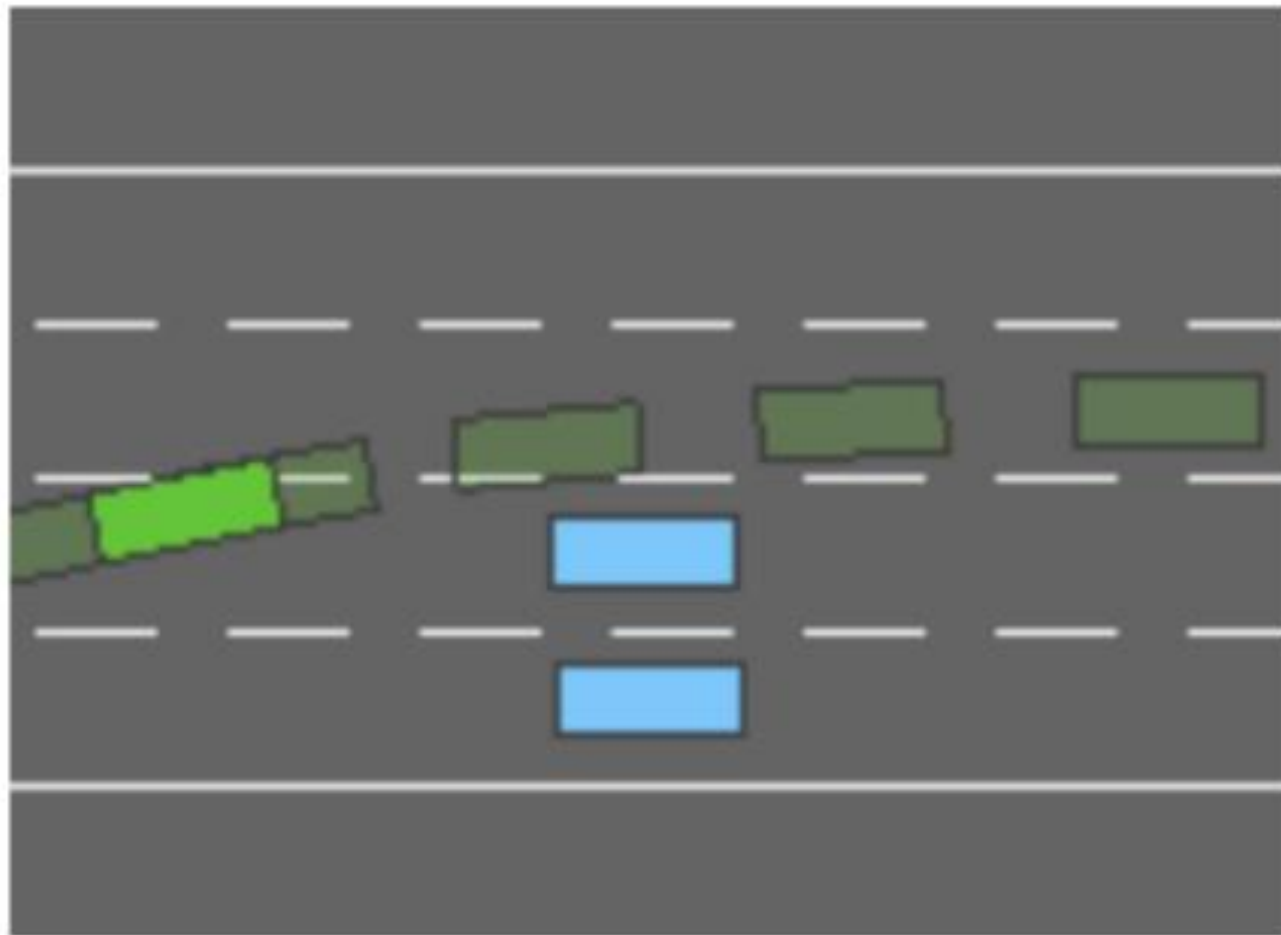
$$L_{ITD} = \|\Phi(s_t, a_t) + \gamma \Psi(s_{t+1}, a_{t+1}) - \Psi(s_t, a_t)\|$$

- Given Φ , learn each agents' Ψ^k with supervised behavior cloning loss
- Update Φ with rewards observed by ego-agent: $r^{ego}(s, a) = \Phi(s, a)^T \mathbf{w}^{ego}$

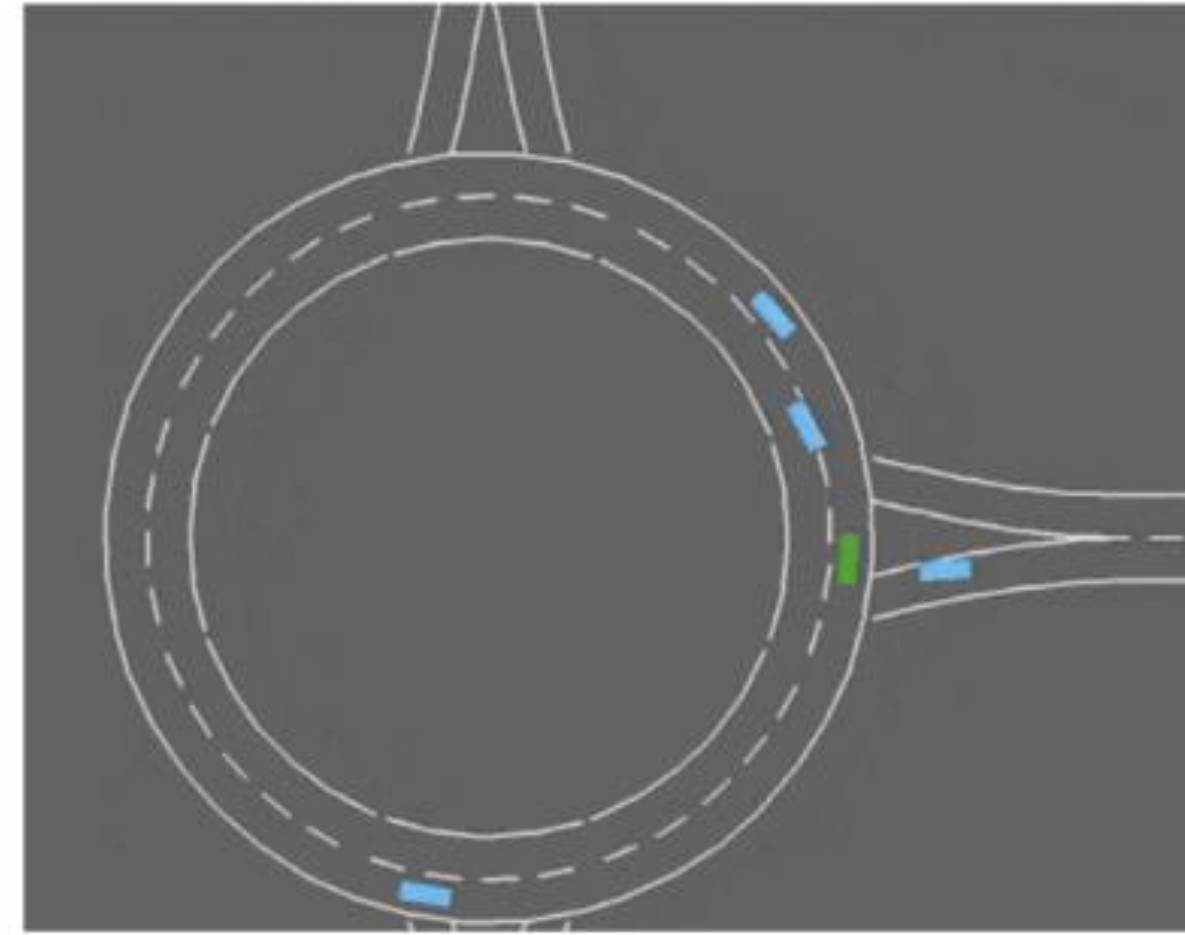
$\Psi\Phi$ -Learning



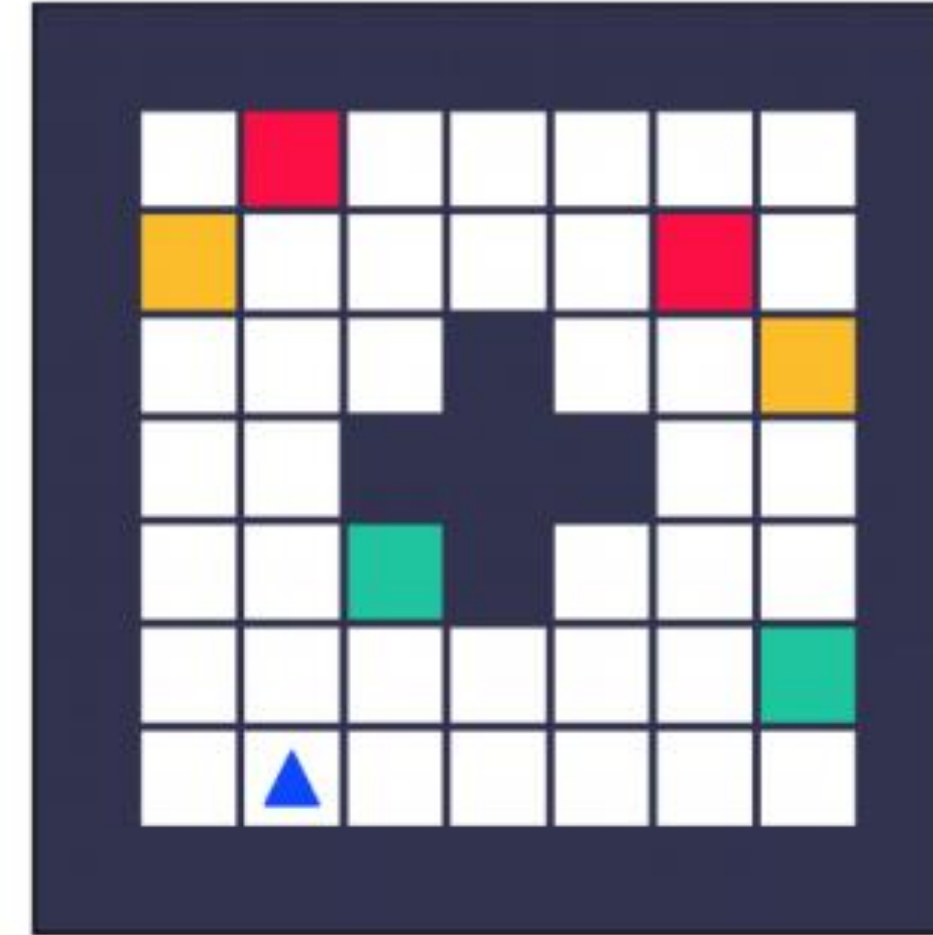
$\Psi\Phi$ -Learning experiments



(a) Highway



(b) Roundabout



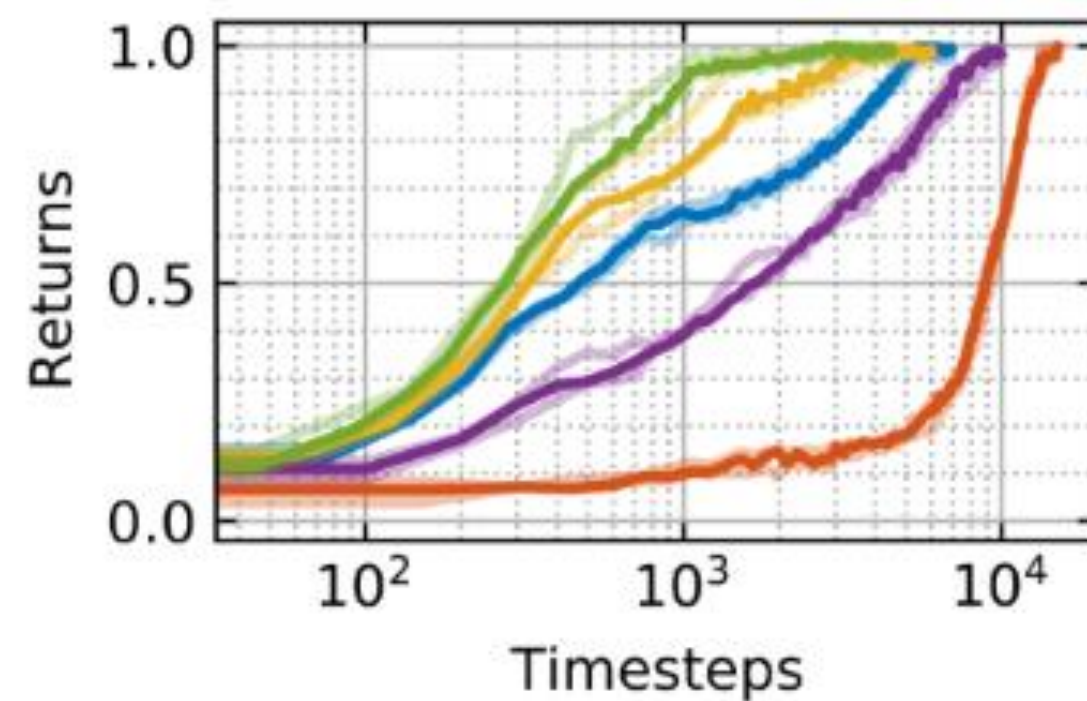
(c) CoinGrid



(d) Fruit Bot

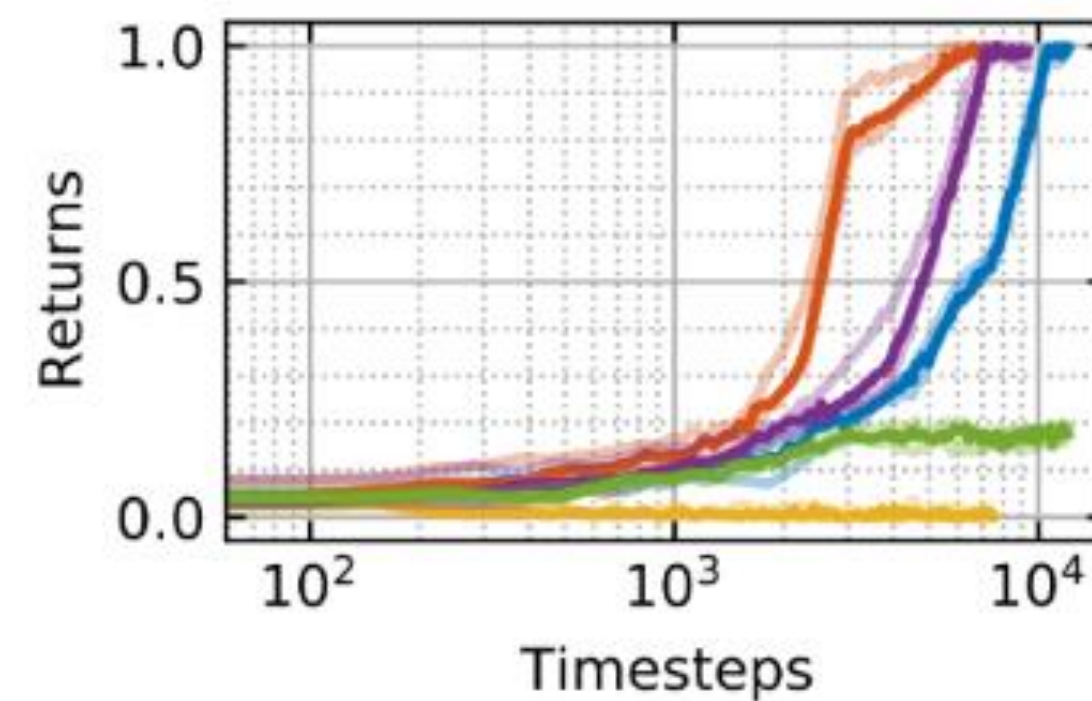
$\Psi\Phi$ -Learning improves on RL tasks

— $\Psi\Phi\mathbf{L}$ (ours) — RL — BC — RL + BC-Aux — SQILv2



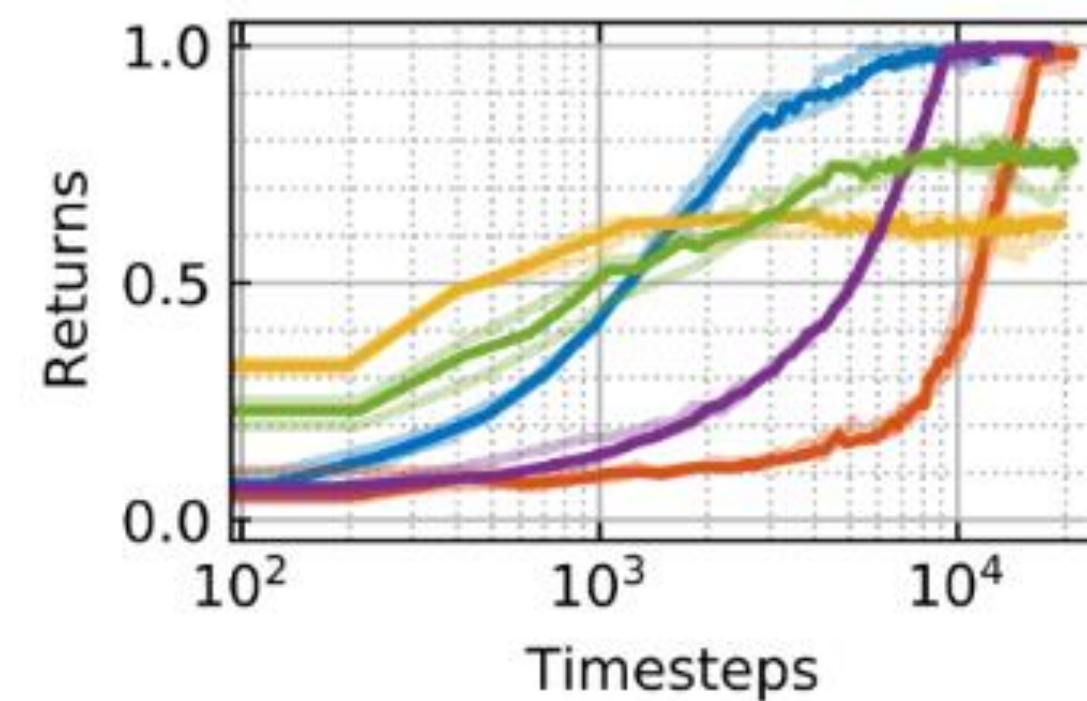
(a) Highway: Single-task

As good as IL
when all other
agents relevant



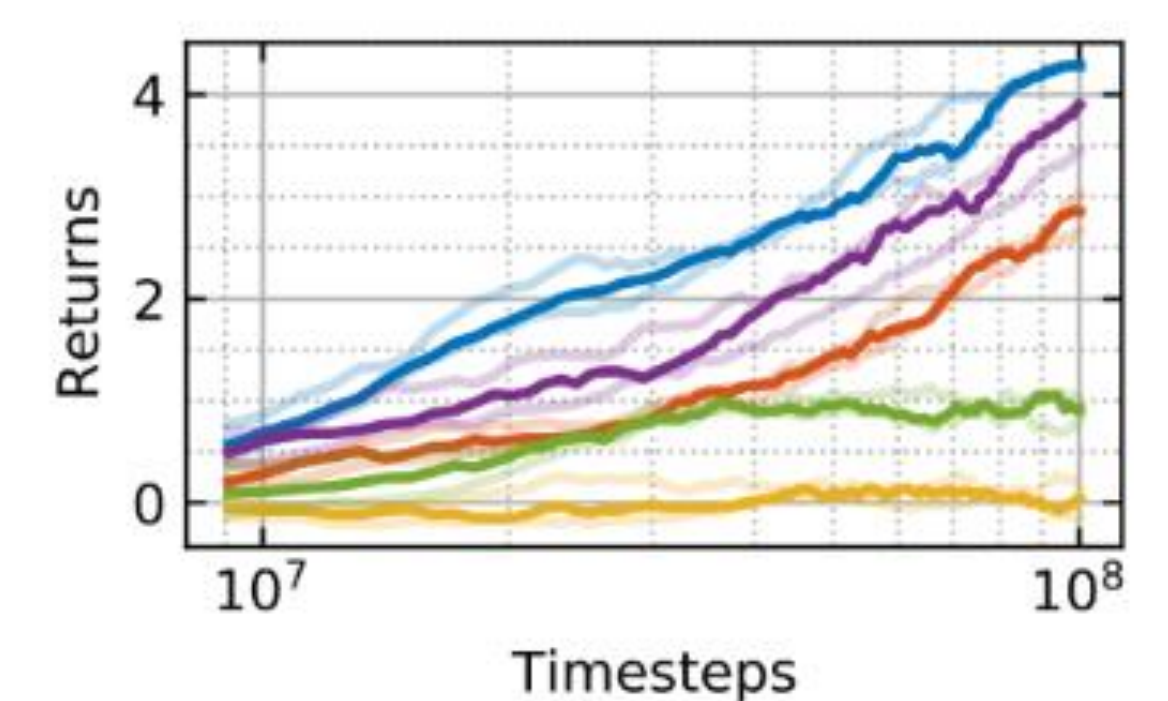
(b) Highway: Adversarial

As good as RL
when all other
agents irrelevant



(c) Highway: Multi-task

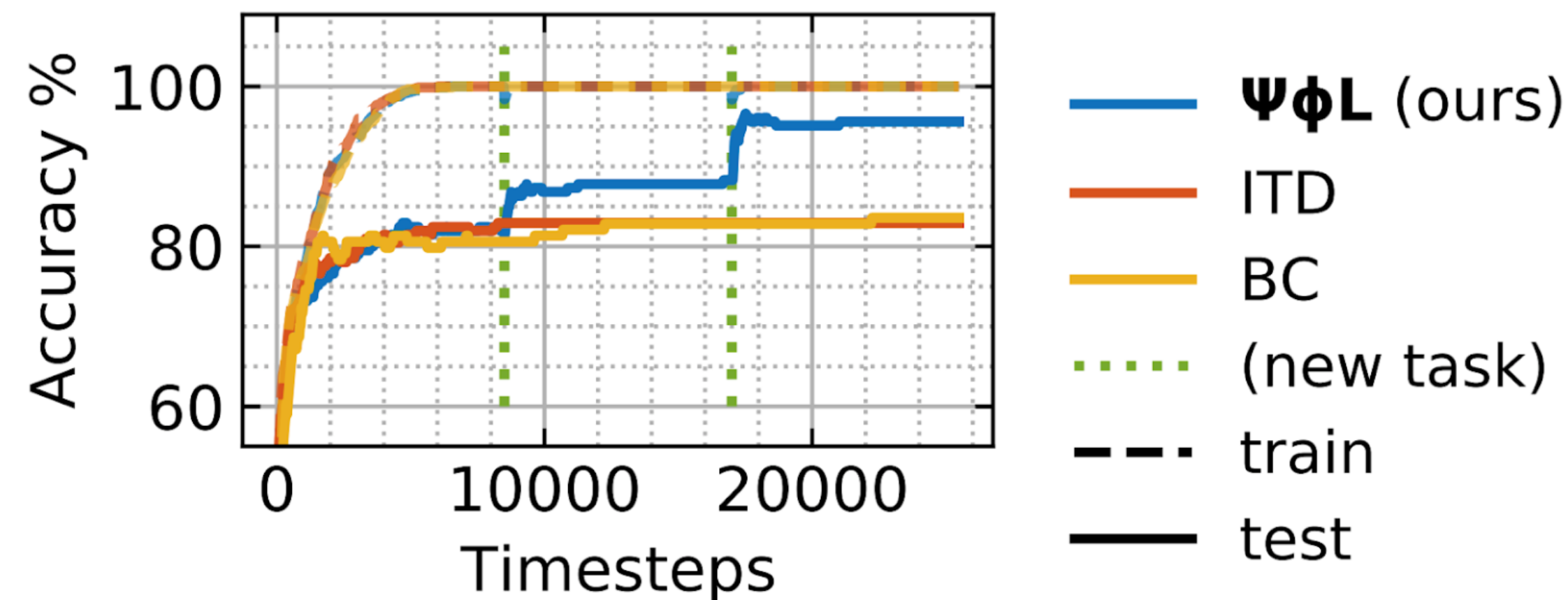
Better than both in multi-task
setting



(d) FruitBot

$\Psi\Phi$ -Learning improves modeling of other agents

Predicting actions:



Inverse RL:

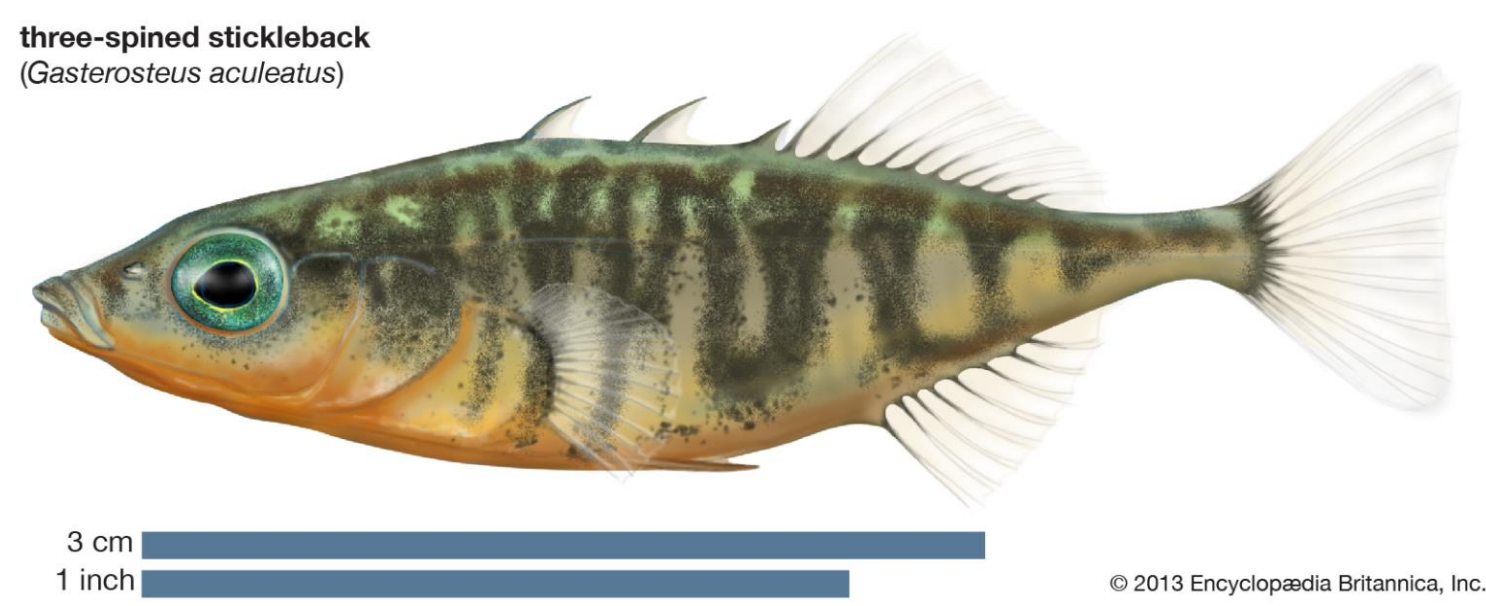
(predict rewards,
train on them)

Methods	Roundabout ^{DQN}	CoinGrid ^{DQN}	FruitBot ^{PPO}
BC ^{†♣} (Pomerleau, 1989)	0.81±0.02	0.69±0.06	0.37±0.02
SQIL ^{†♣} (Reddy et al., 2019)	0.85±0.02	0.64±0.05	0.35±0.03
GAIL ^{†◇} (Ho & Ermon, 2016)	0.77±0.07	0.73±0.02	0.31±0.02
ITD [◇] (ours, cf. Section 3.1)	0.92±0.01	0.77±0.03	0.35±0.04

$\Psi\Phi$ -Learning improves generalization

Methods	0-shot				1-shot				100-shot			
	R+G	R-G	-R+G	-R-G	R+G	R-G	-R+G	-R-G	R+G	R-G	-R+G	-R-G
SQILv2♣ (Reddy et al., 2019)	1.0±0.0	0.0±0.0	0.0±0.0	-1.0±0.0	1.0±0.0	0.0±0.0	0.0±0.0	-1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0
$\Psi\Phi$ -learning ◇ (ours, cf. Section 3.2)	1.0±0.0	0.2±0.1	0.2±0.1	-0.4±0.2	1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0	1.0±0.0

Few-shot transfer in CoinGrid



Conclusion

Social learning is a powerful mechanism that can help RL agents:



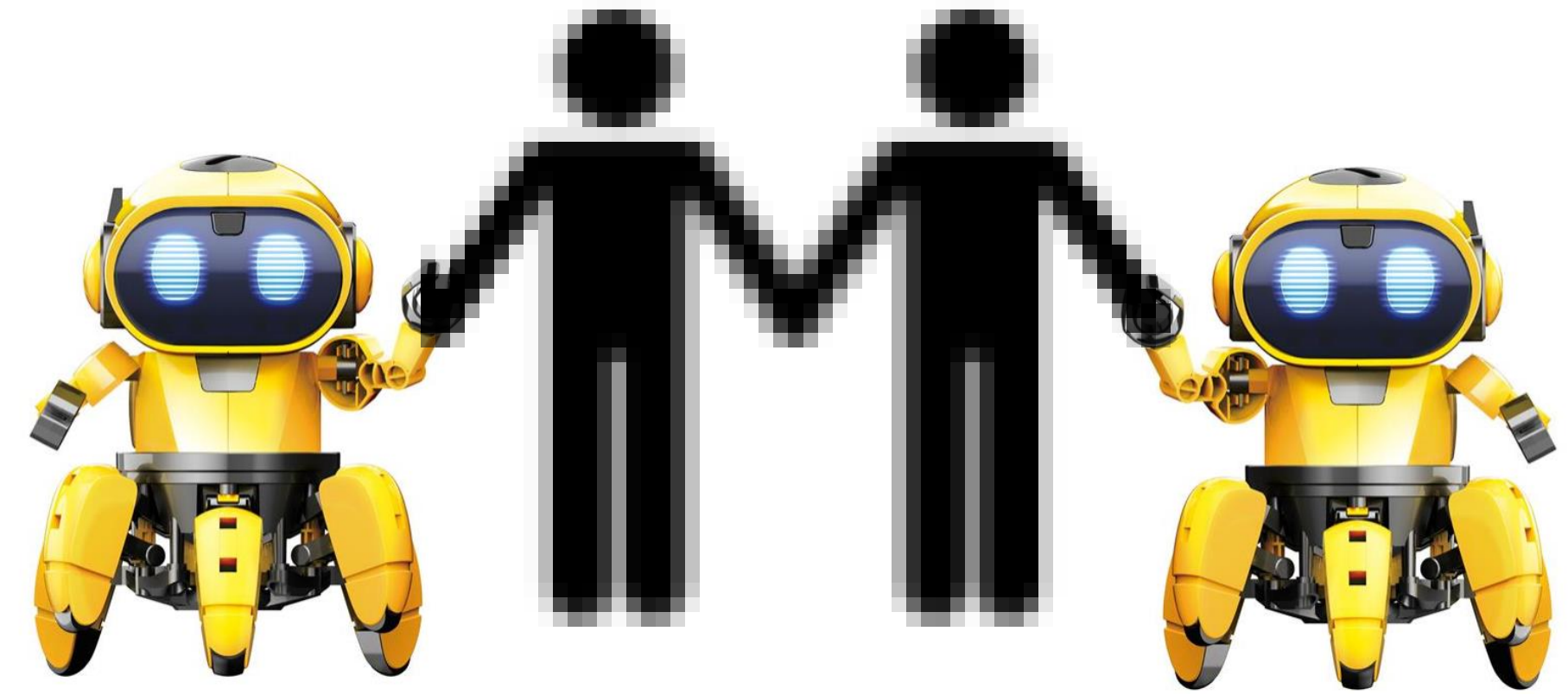
**Rapidly adapt to
new environments**



**Learn complex
behavior**

Social Reinforcement Learning improves...

- **Coordination** with other agents
- **Human-AI interaction**
- **Learning** complex behavior
- **Generalization** to new environments



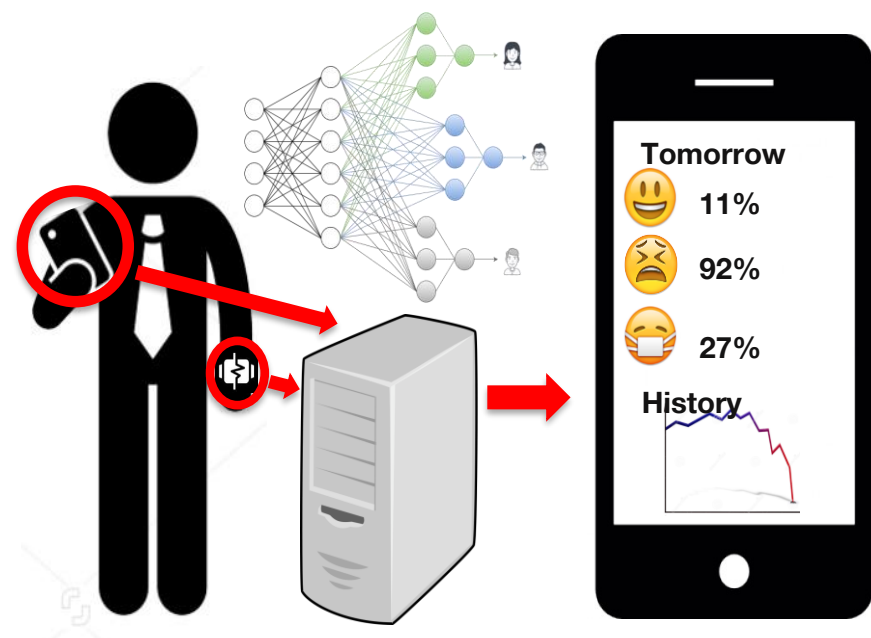
Extra slides

Questions?

Human-AI

Multi-agent

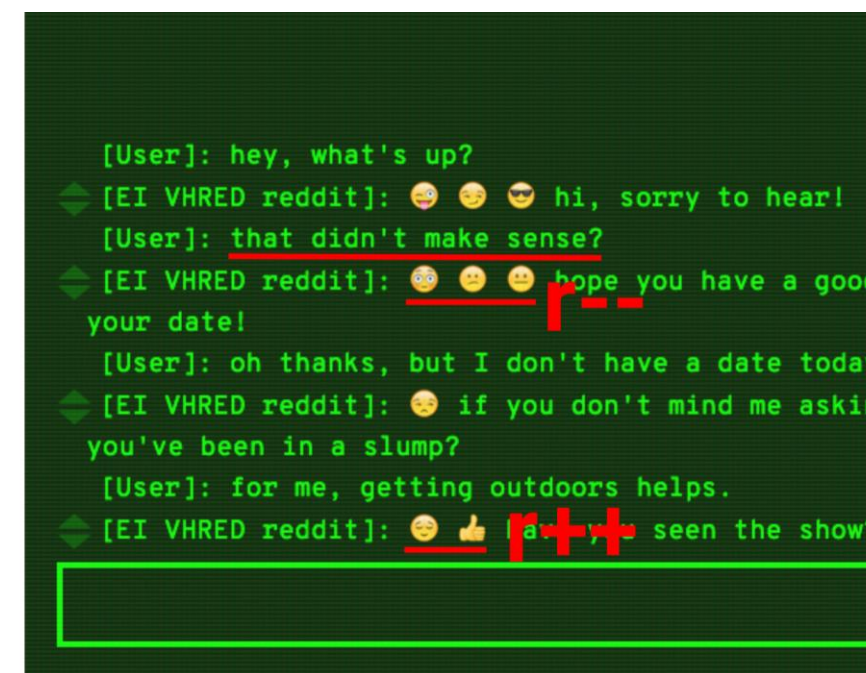
Detecting social
& affective cues



[Personalized Multitask Learning](#). Taylor*, Jaques*, et al. **IEEE Transactions on Affective Computing (TAC) 2020 Best Paper**.

[Predicting bonding from facial expressions and body language](#). Jaques et al. **IVA 2016**

Learning from
human social cues



[Human-centric Dialog Training via Offline RL](#). Jaques*, Shen*, et al. **EMNLP 2021**.

[Human Evaluation of Dialog Systems](#). Ghandeharioun*, Shen*, Jaques*, et al. **NeurIPS 2020**.

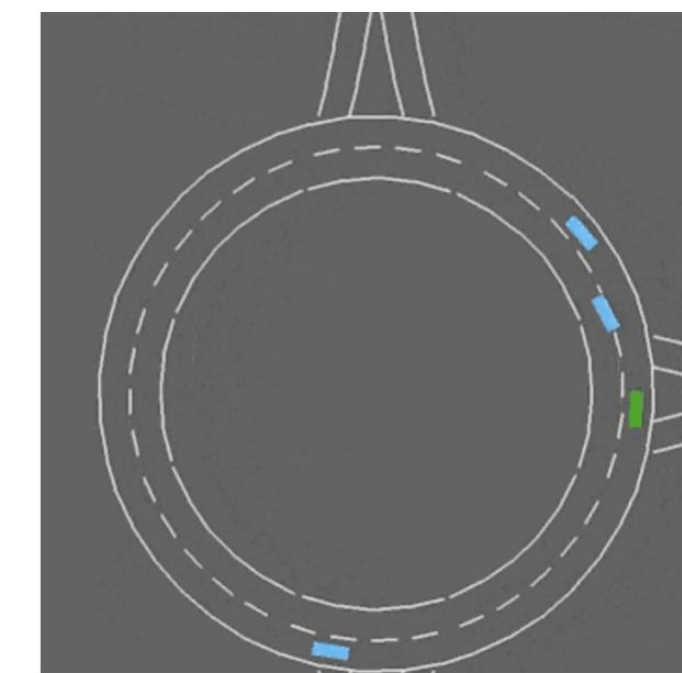
[Hierarchical RL for Dialog](#). Saleh*, Jaques*, et al. **AAAI Oral (top 8%)**

Coordination



[Social influence](#). Jaques et al. **ICML 2019 Best Paper Honourable mention (top 0.26%)**.

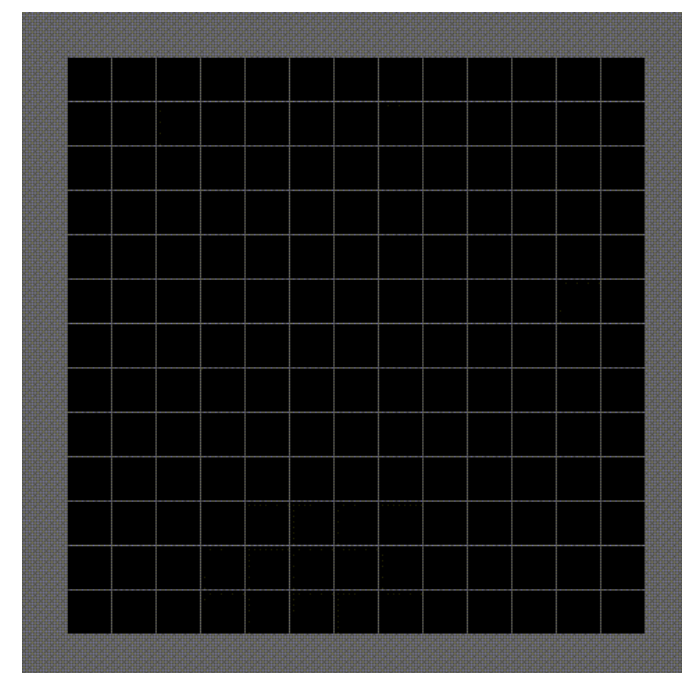
Multi-agent
social learning



[Emergent Social Learning from MARL](#). Ndousse, ..., Jaques. **ICML 2021 & NeurIPS CoopAI workshop Best Paper 2020**

[PsiPhi-Learning](#). Filos, ... Jaques*, Farquhar*. **ICML 2021 oral (top 3%)**.

Emergent
Complexity



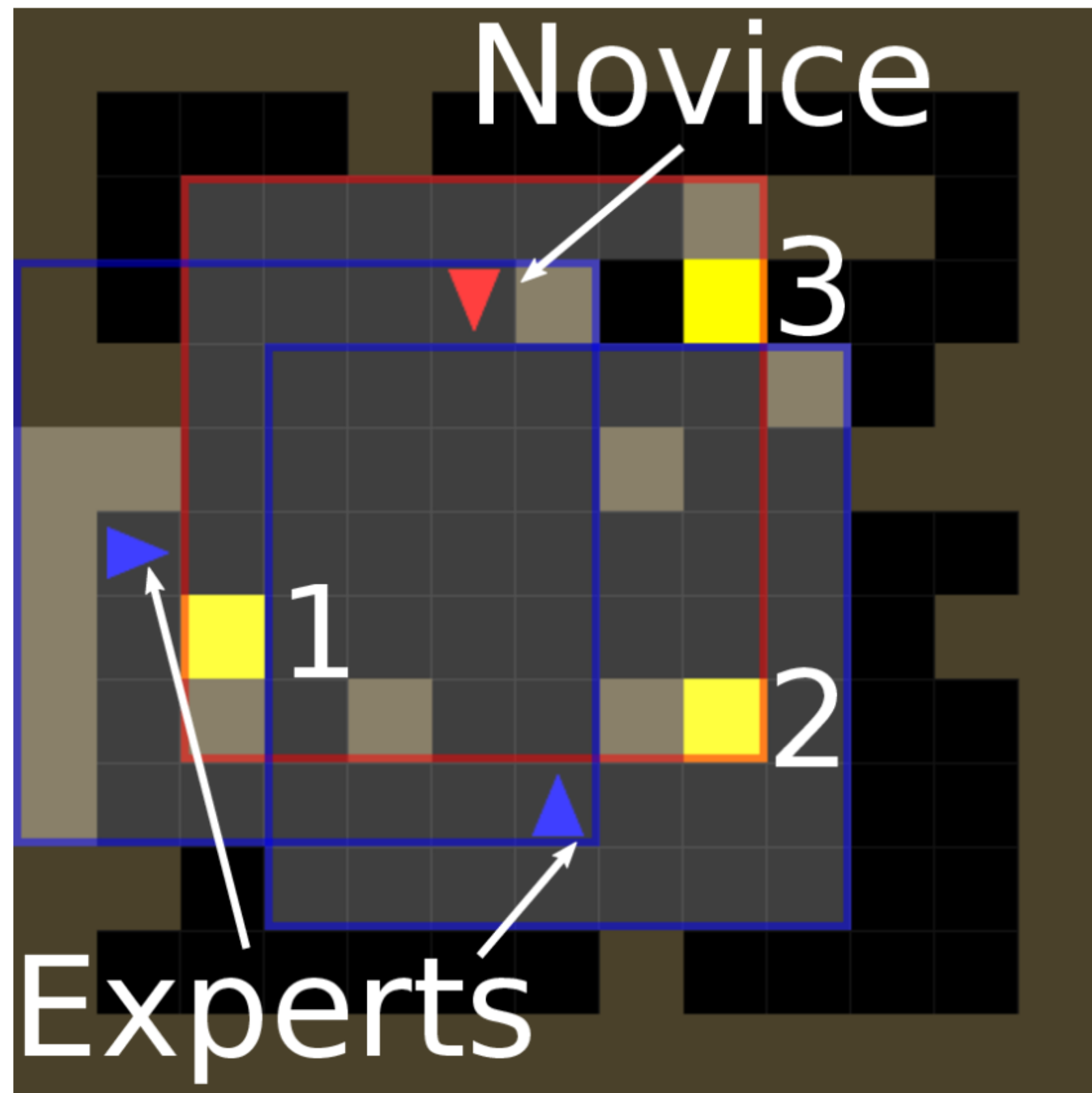
[PAIRED](#). Dennis*, Jaques*, et al. **NeurIPS 2020 oral (top 1%)**.

[Environment Generation for Zero-Shot Compositional RL](#). Gur, Jaques, et al. **NeurIPS 2021**

Social learning environment



Kamal
Ndousse



(a) Goal Cycle

- **Partial observability** and **no privileged access** to other agents states or actions
- Information (correct goal ordering) is **hidden**¹. Incorrect **exploration expensive**.
- Agents change colour with the recent average of their rewards (**prestige cue**)

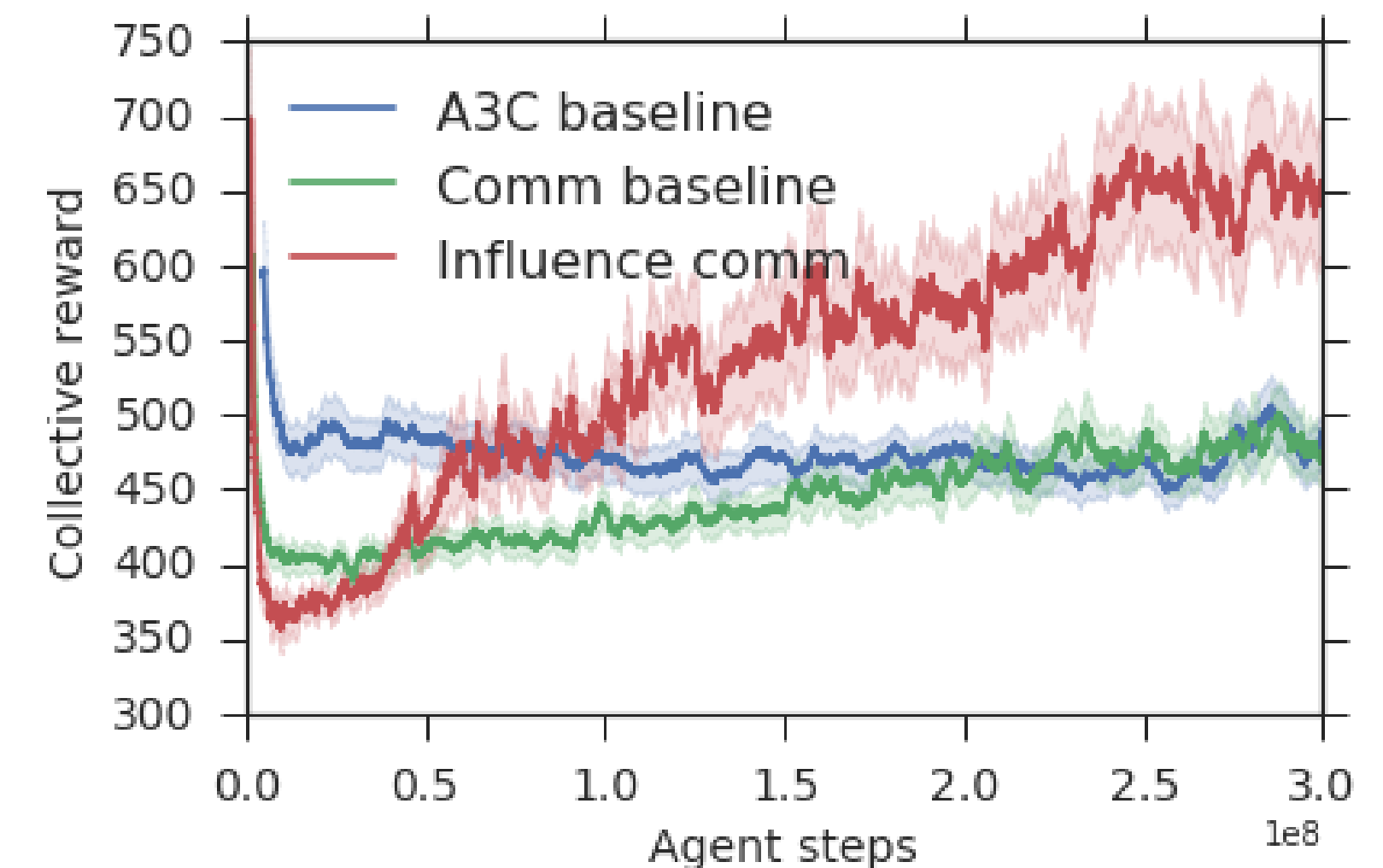
¹Hiding information to evoke social learning proposed in Borsa, D., Piot, B., Munos, R., & Pietquin, O. (2017). Observational learning by reinforcement learning.

Coordinate with other agents... via social influence

Give agents an intrinsic social reward for having a **causal influence** on other agents' actions. Like **social empowerment**.

$$D_{KL}[p(a_{t+1}^B | s_t^A, a_t^A) || p(a_{t+1}^B | s_t^A)]$$

- Enhances **cooperation** in social dilemmas
- Leads to **emergent communication**
- Works without centralized control or privileged access to other agent's states/rewards



Improve human-AI interaction.. in dialog

✕ Close Chat and Rate

[User]: hey, what's up?

◆ [EI VHRED reddit]: 🤔 😊 😎 hi, sorry to hear!

[User]: that didn't make sense?

◆ [EI VHRED reddit]: 🤔 😞 😐 hope you have a good day, and good luck on your date! **r--**

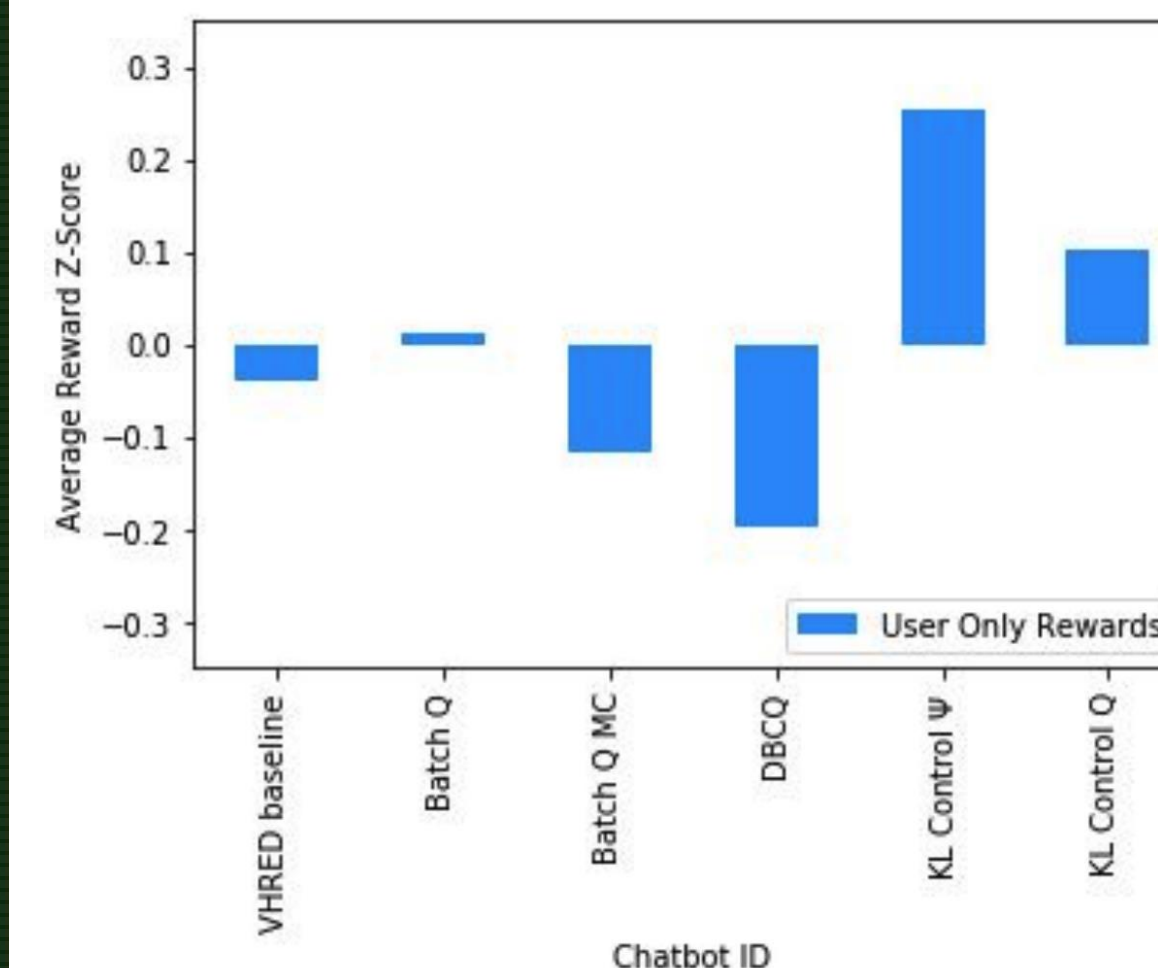
[User]: oh thanks, but I don't have a date today.

◆ [EI VHRED reddit]: 😞 if you don't mind me asking, what do you do when you've been in a slump?

[User]: for me, getting outdoors helps.

◆ [EI VHRED reddit]: 😊 👍 have you seen the show? **r++**

SEND



Learning from human-AI interaction... with faces

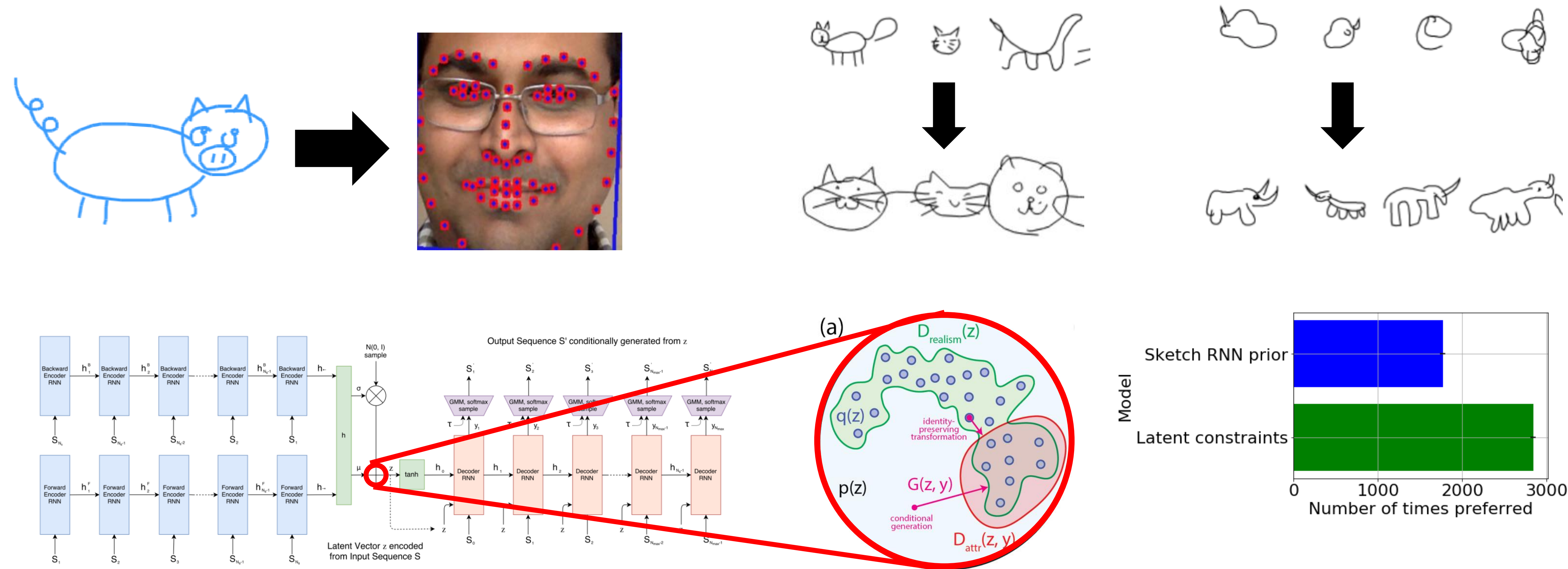
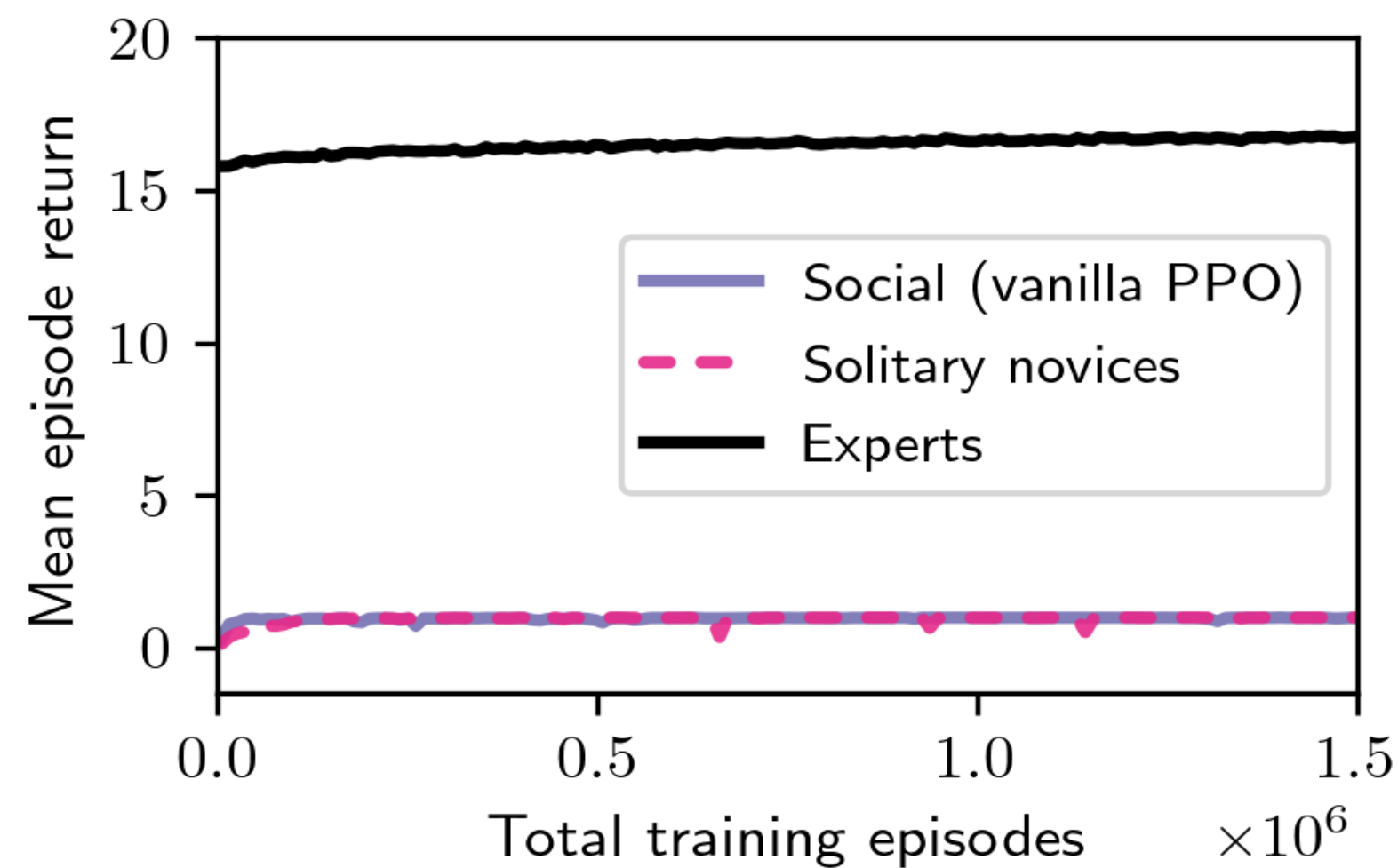


Figure 2: Schematic diagram of sketch-rnn.

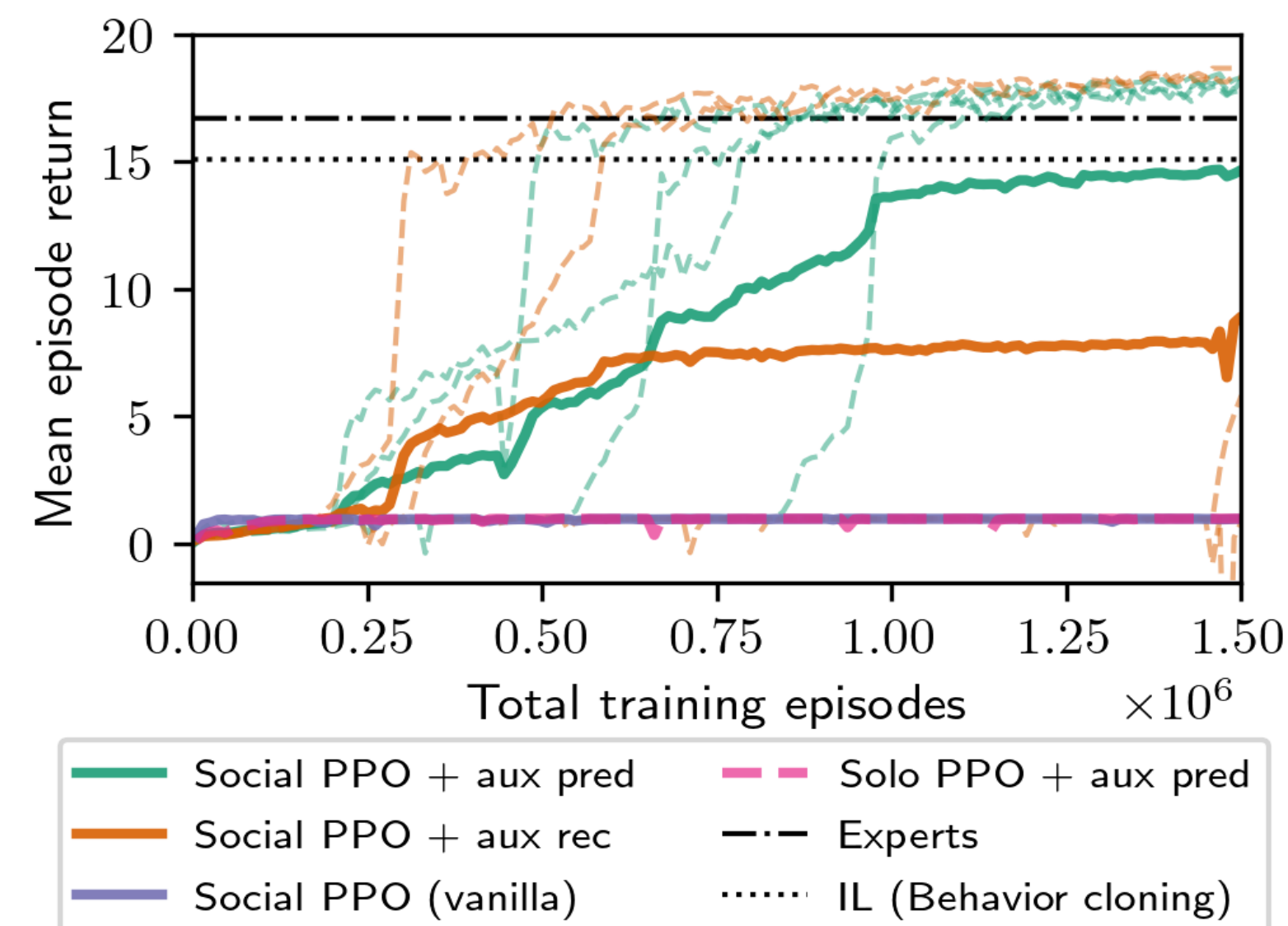
Learning via social awareness: Improving a deep generative sketching model with facial feedback. Jaques, McCleary, Engel, Ha, Bertsch, Picard, Eck (2018). International Conference on Representation Learning (ICLR) workshop.

Multi-agent social learning

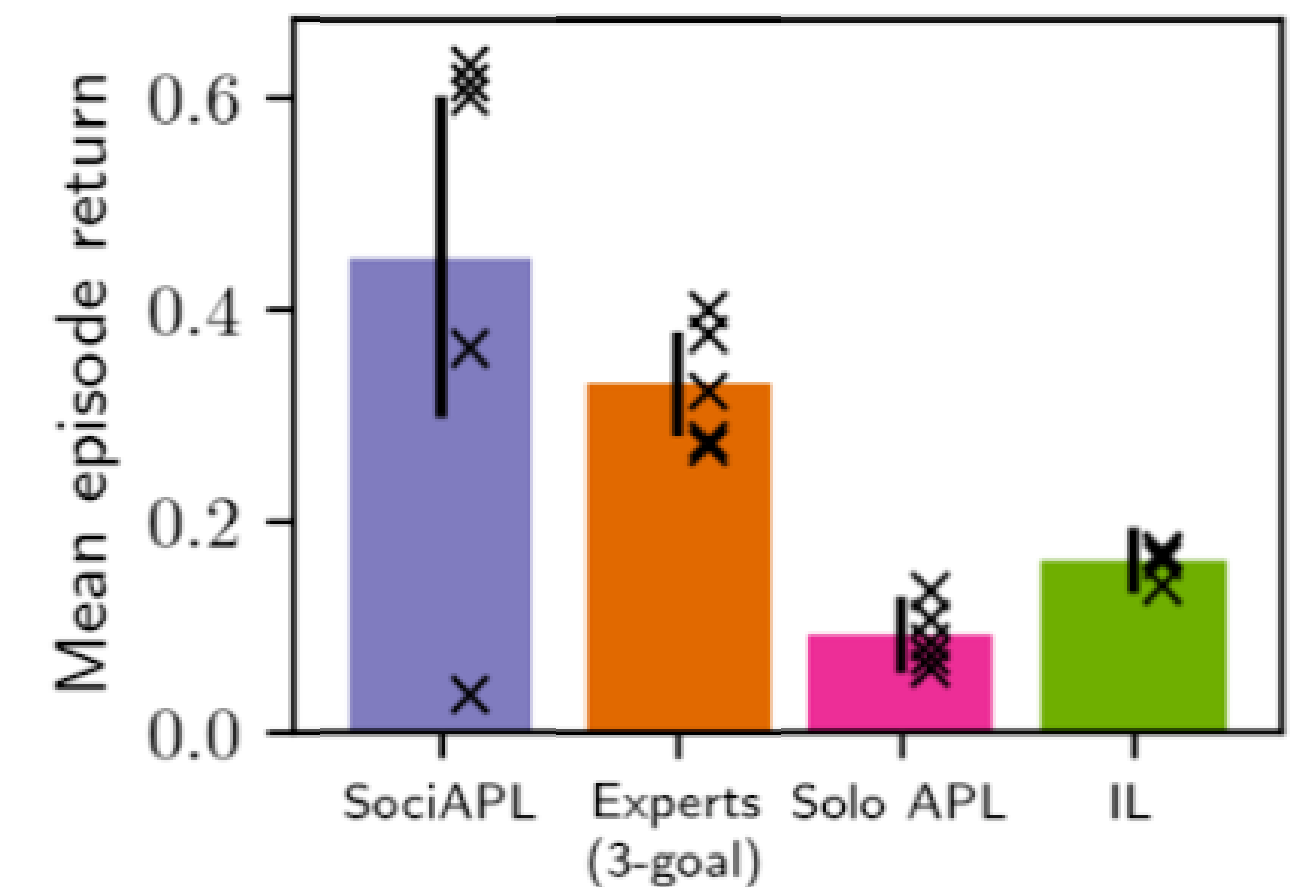
Model-free RL agents do not learn from experts in their environment



Fix with model-based auxiliary loss



Social learners generalize better than **RL experts** or **imitation learners**



Improving Social Learning with PsiPhi-Learning

- Model both other agents and RL agent with **successor features**

$$Q^k(s, a) = \Psi^k(s, a)^T \mathbf{w}^k$$

- Learned a **shared basis for behavior**

$$\Psi^k(a_t, s_t) = \mathbb{E}^\pi \left[\sum_{t'=t}^T \phi(a_{t'}, s_{t'}) \right]$$

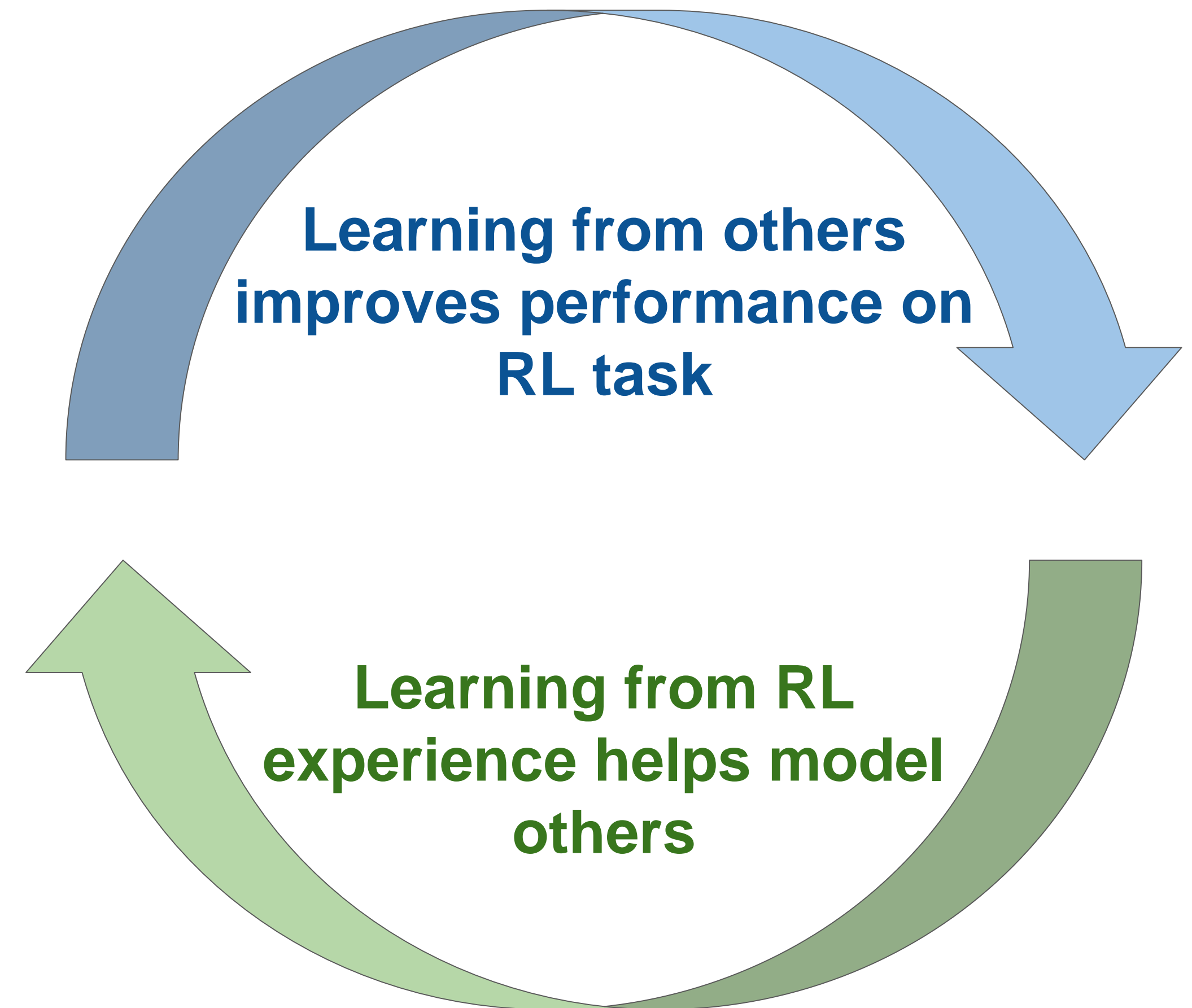


Figure 1 displays the performance of different algorithms across four tasks, showing Returns (Y-axis) versus Timesteps (X-axis, logarithmic scale).

Legend:

- $\Psi\Phi\mathcal{L}$ (ours) (Blue line)
- RL (Red line)
- BC (Yellow line)
- RL + BC-Aux (Purple line)
- SQLv2 (Green line)

(a) Highway: Single-task

(b) Highway: Adversarial

(c) Highway: Multi-task

(d) FruitBot

...multi-task demos

[illegible]

Imitation

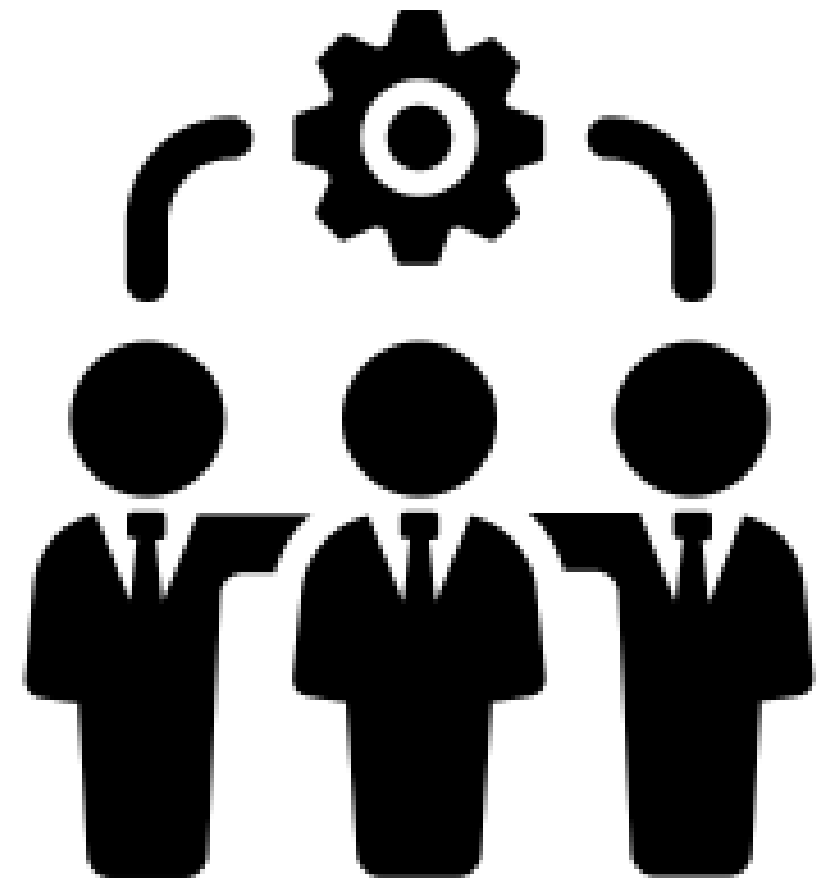
[illegible]

Social learning

Learning from other intelligent agents in your environment



Interaction,
feedback



Cooperation



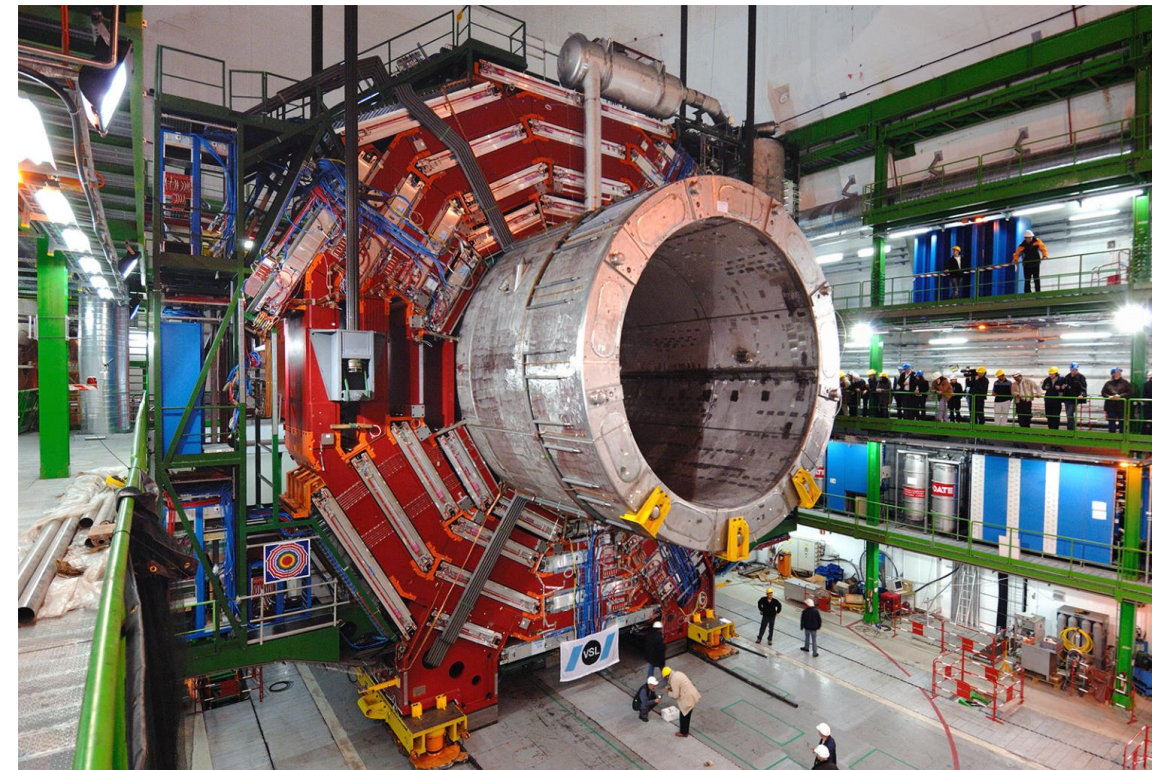
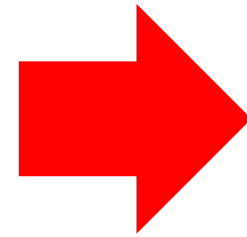
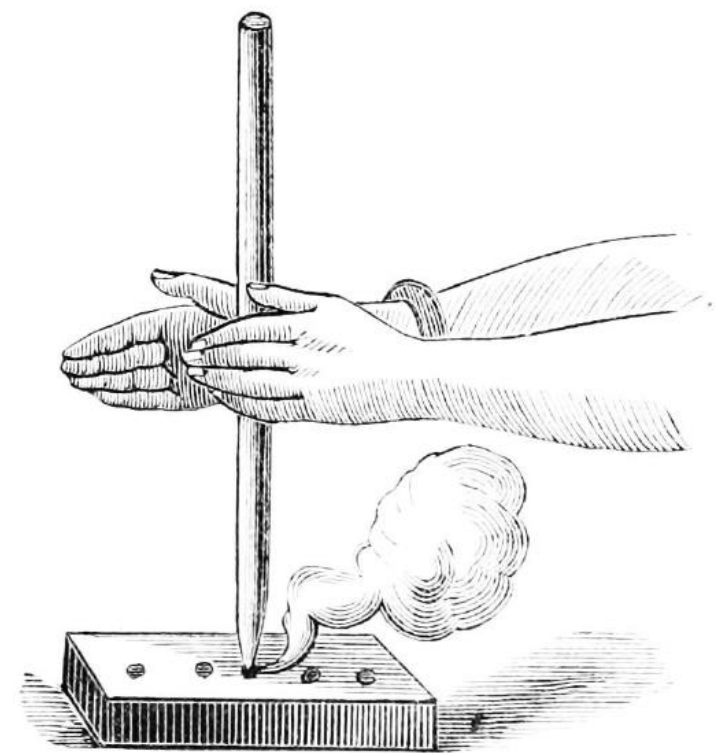
Competition



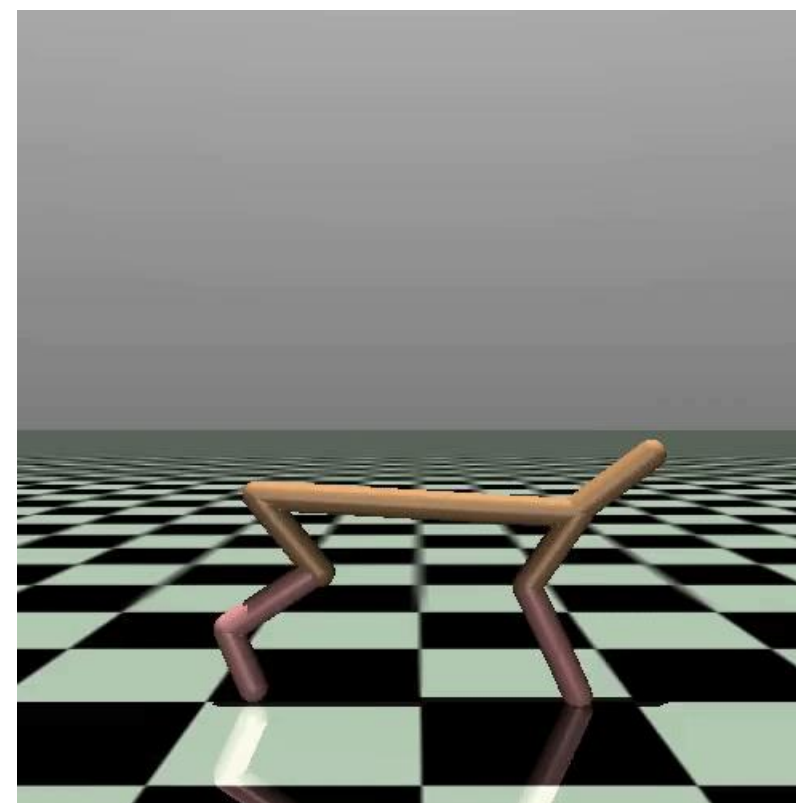
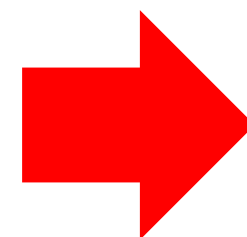
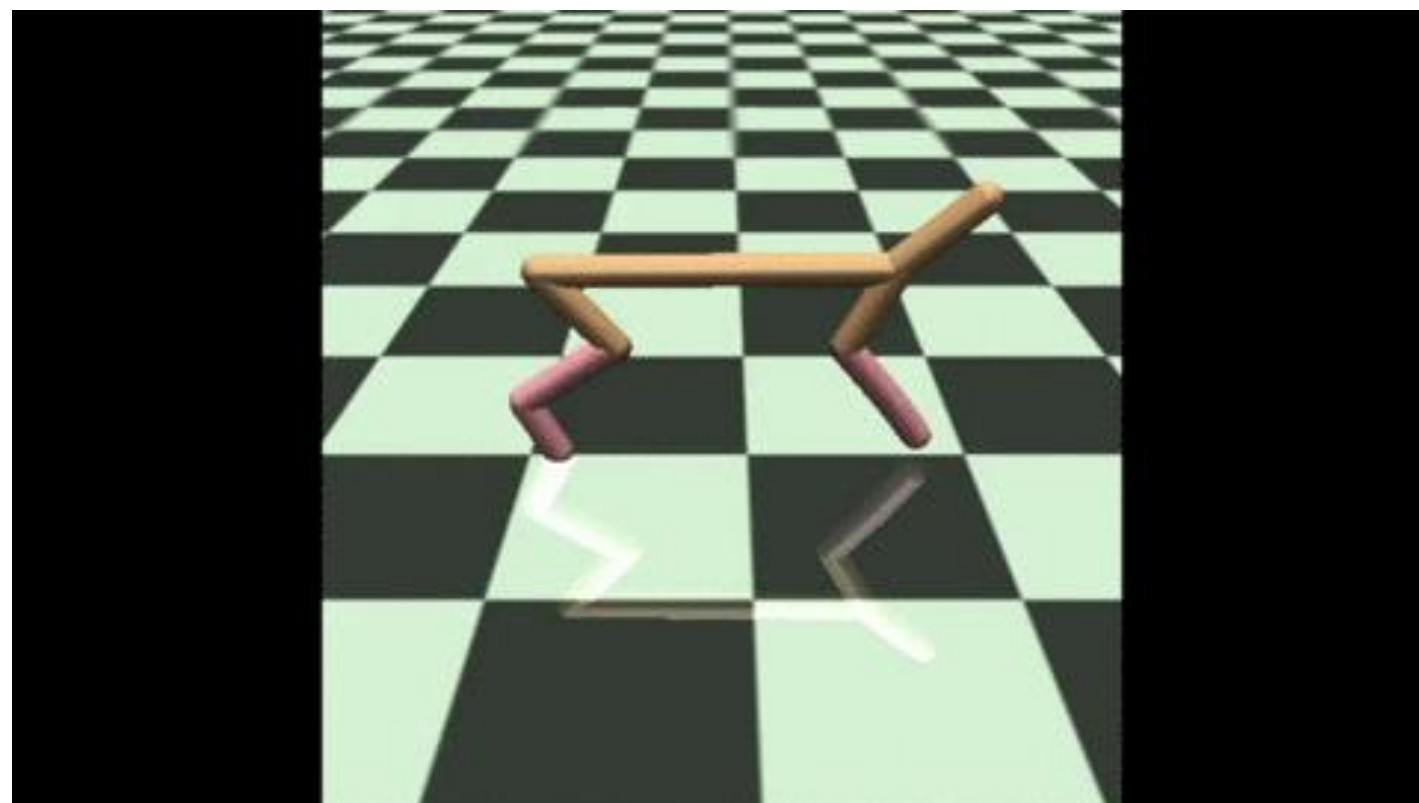
Co-existence

Benefits of social learning for RL

- Individual exploration is expensive, time consuming



- Social learning lets you “stand on the shoulders of giants”



Benefits of social learning for RL

- Individual exploration is dangerous



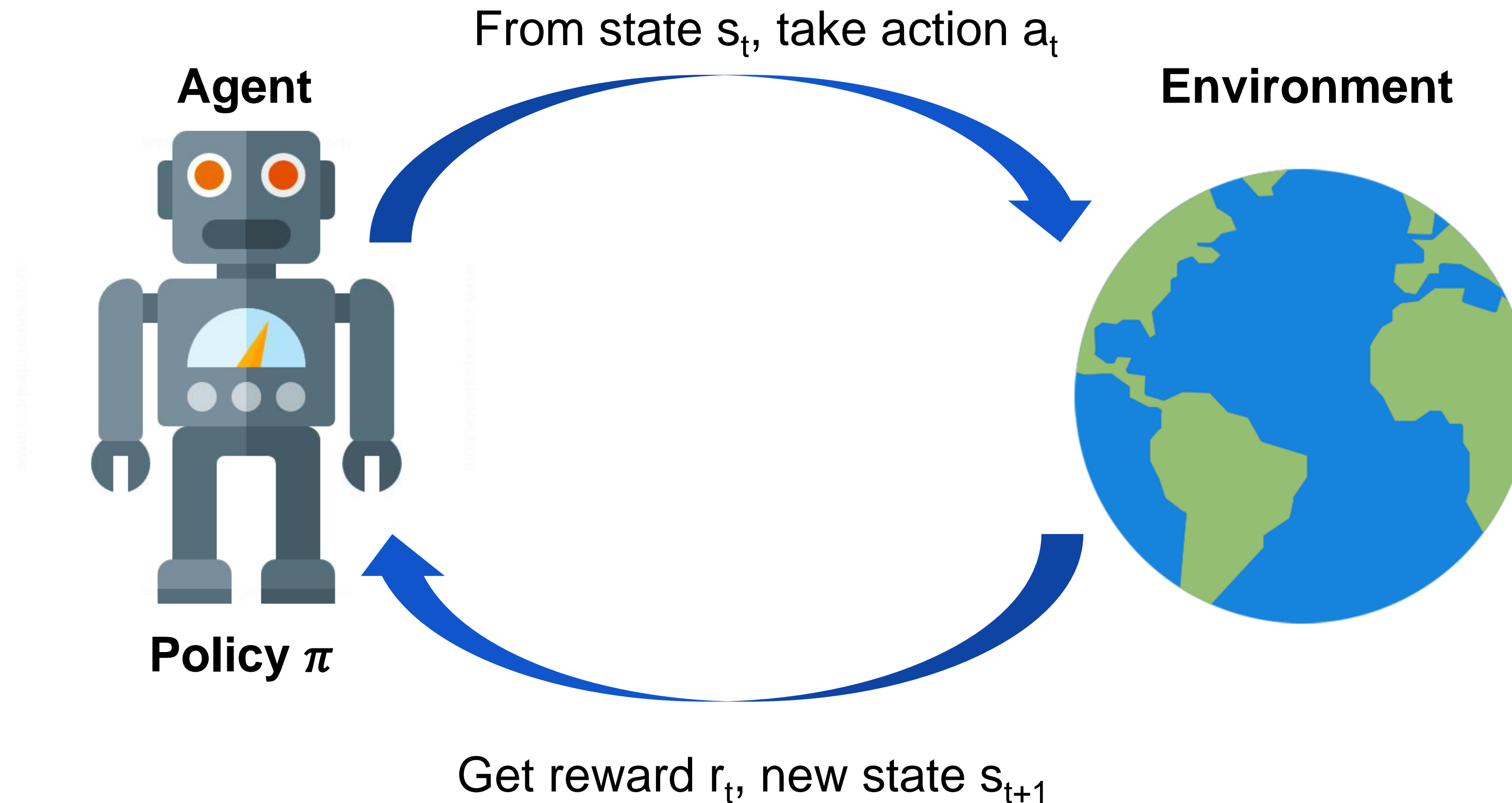
¹<https://www.gpsworld.com/autonomous-car-hits-autonomous-robot-in-bizarre-collision/>

²<https://www.sciencealert.com/video-captures-self-driving-tesla-hitting-and-killing-a-robot-in-las-vegas>



From Koide, K., & Miura, J. (2016). Identification of a specific person using color, height, and gait features for a person following robot. *Robotics and Autonomous Systems*, 84, 76-87.

Reinforcement Learning



Goal: maximize discounted future reward

$$R(\tau) = \sum_{t=0}^T \gamma^t r_t$$

γ = discount factor

- **Sequential decision making:** interact with environment at each timestep t
- **Trial and error learning.** When to **explore** vs. **exploit**