

Announcements

- Homework 5 Deadline **extended to Friday (March 28, 11:59pm)**.
- No TA hours tomorrow, normal TA hours resume next week.
- Your Full Proposal + Lit Review are due April 4th!



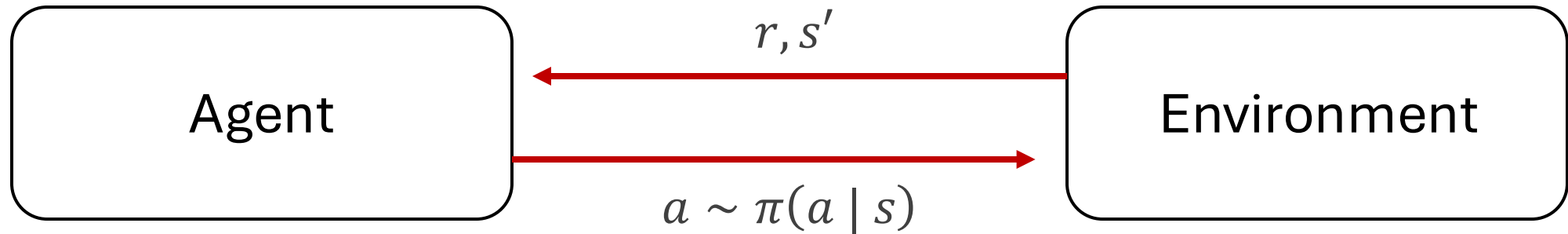
Multi-Agent Reinforcement Learning

Guest Lecturer: Connor Mattson

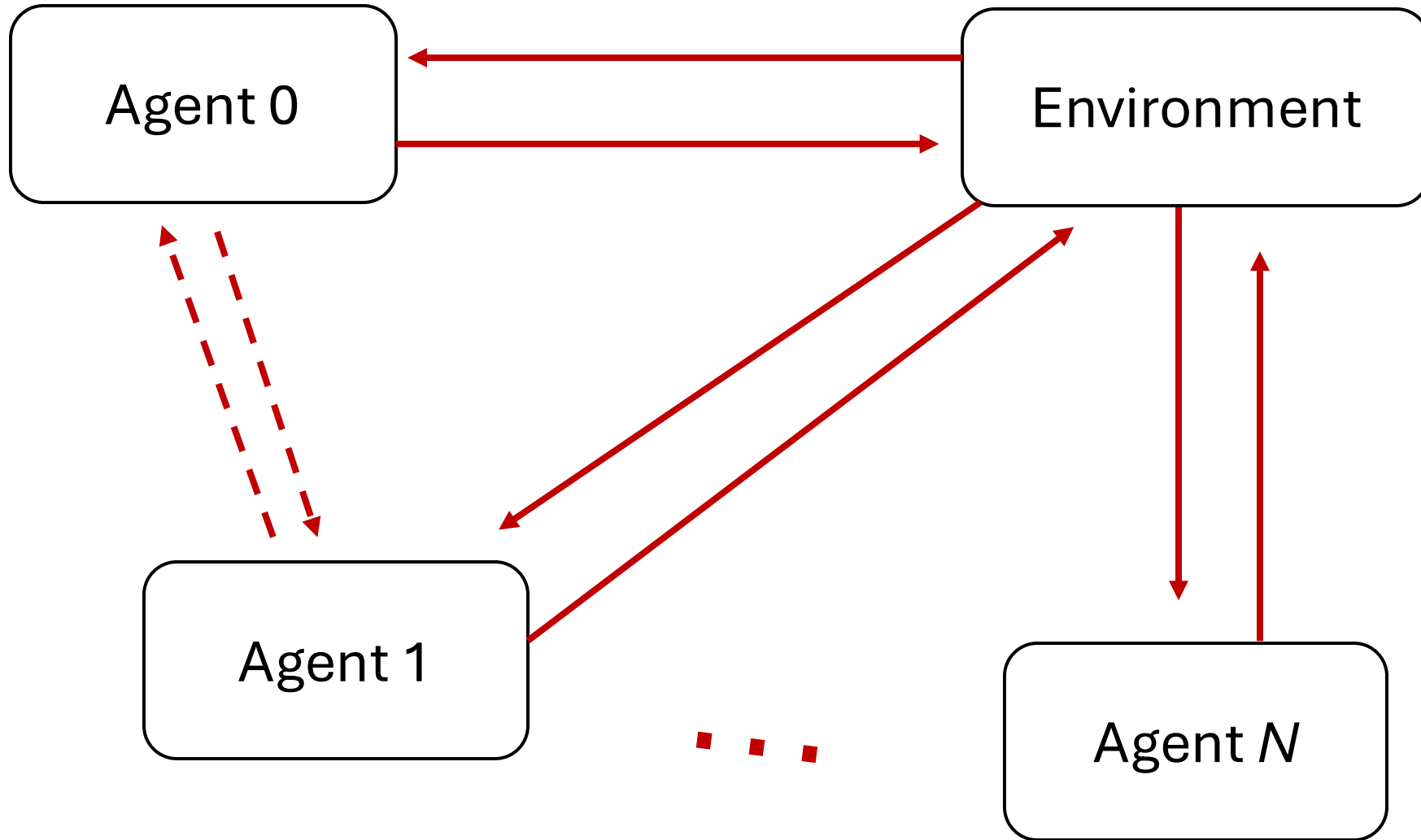
CS5955/CS6955: Advanced Artificial Intelligence, Spring 2025

What are Multi-Agent Systems?

Recall our simple single-agent interaction loop!

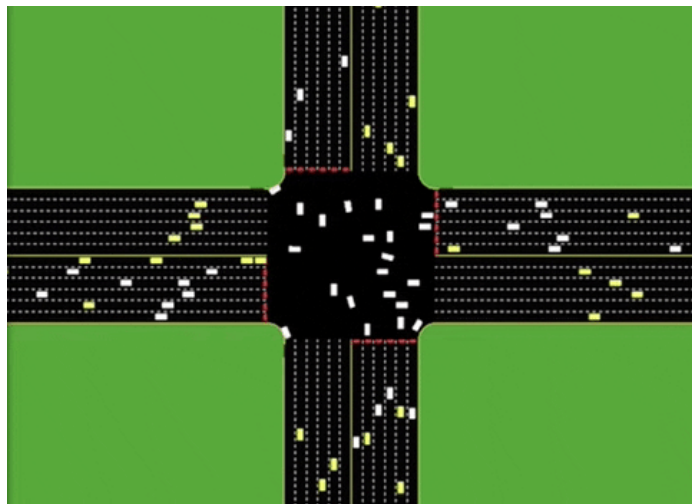


What are Multi-Agent Systems?

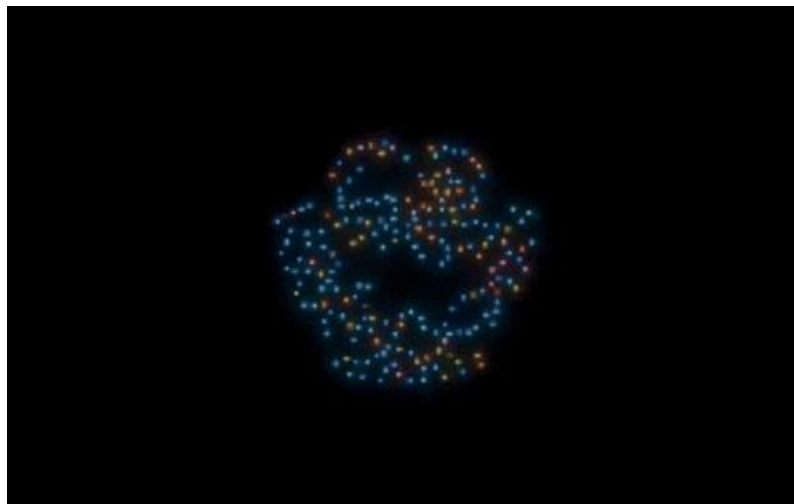


What information is being passed along the red arrows?

Examples of Multi-Agent Applications



Autonomous Driving
+ Traffic Control



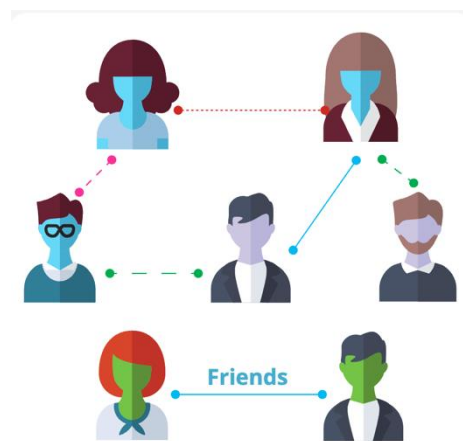
UAVs + Drones



Coordinated Robot Motion



n-Player Games



Psychology + Sociology



Human Multi-Robot Interaction

Amazon Robotics



Why Study Multi-Agent Systems?

Some of humanities greatest feats are multi-agent problems!

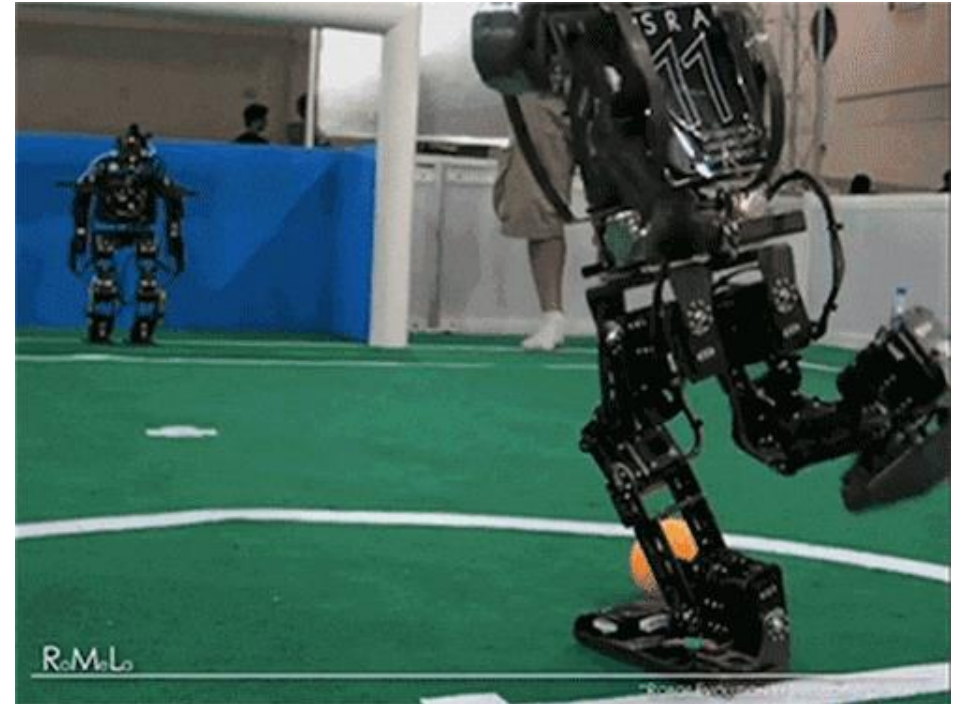


Examples:

Team Sports, Government, International Relations, Academic Research, Game Theory, Resource Allocation, etc.

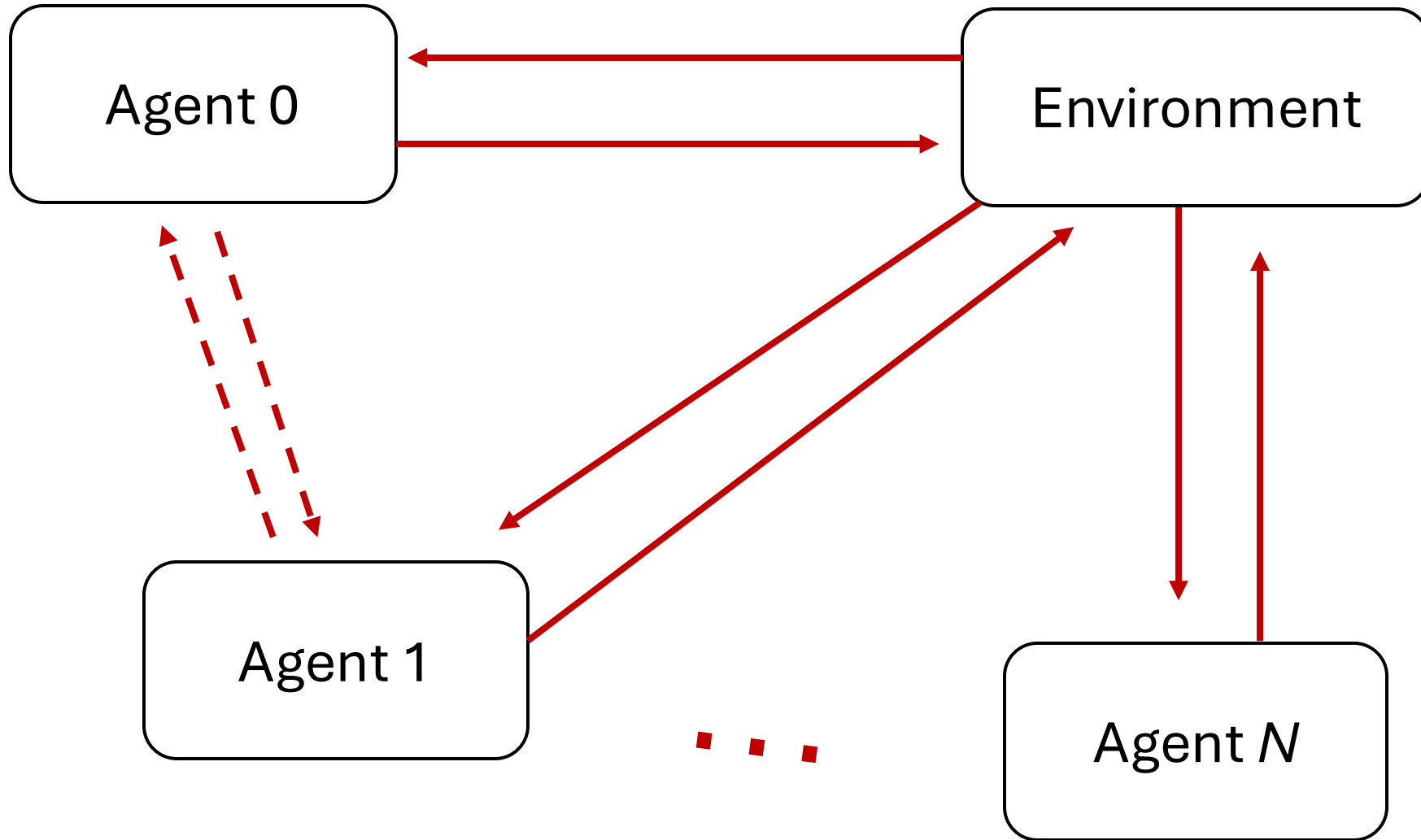
Videos: Fox News, RoboCup

Robots/Agents should also possess these capabilities!



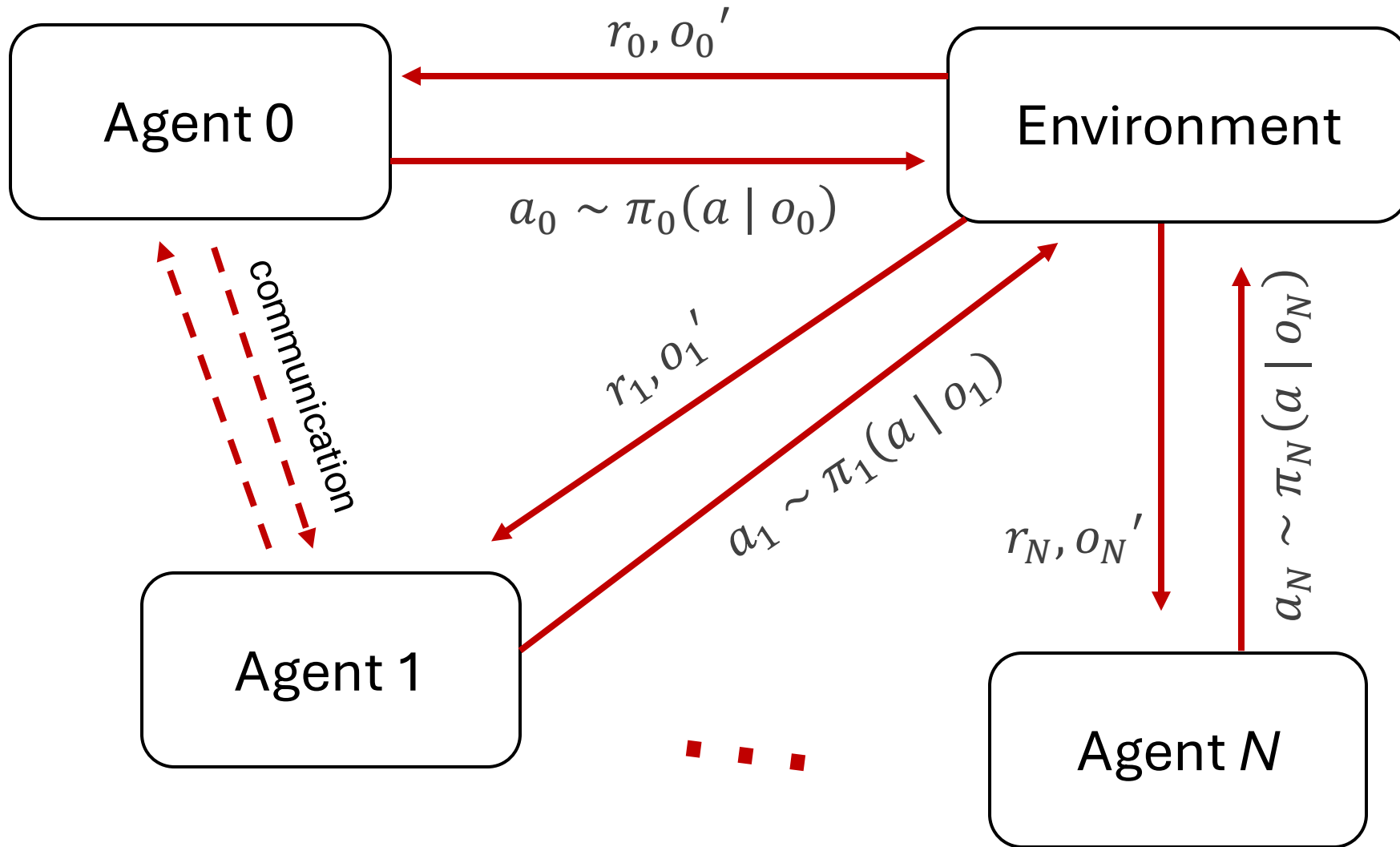
(we have a long way to go!)

What are Multi-Agent Systems?



What information is being passed along the red arrows?

What are Multi-Agent Systems?



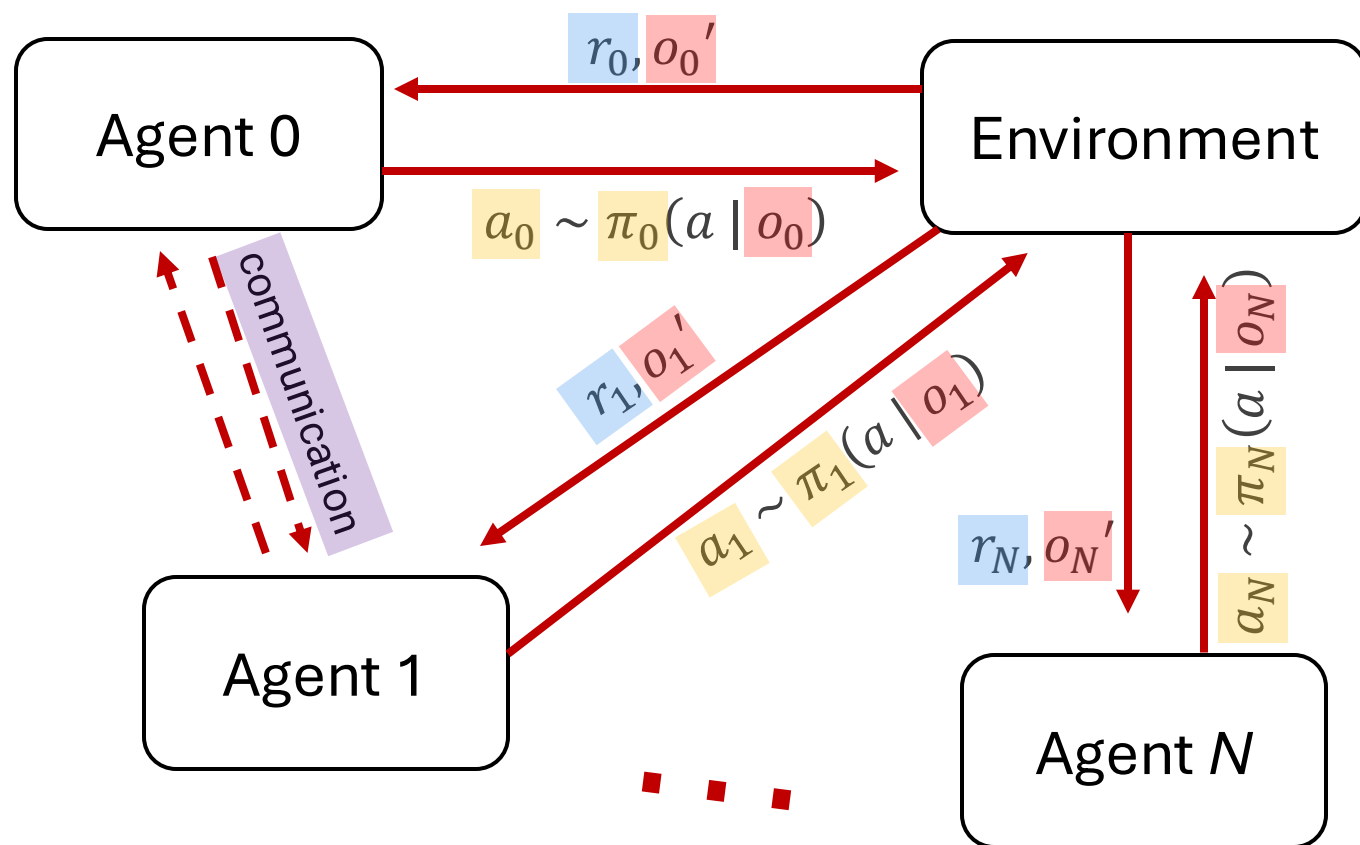
What do you notice?

Rewards?

Actions/Policies?

State/Observations?

How does MARL differ from single-agent RL?



Separate State Signals: Agent 0 and Agent 1 receive o_0' and o_1' , respectively.

Separate Reward Signals: Agent 0 and Agent 1 receive r_0 and r_1 , respectively.

Separate Policies: π_i can be different for each agent (so a_i are different!)

Communication: Agents can communicate with each other (if allowed)

Parallels to Single-Agent RL

Same Desired Outcome: Policies that maximize an expected return.

Assumptions: We don't know transition dynamics but have reward signals.

Same Definitions of Value: We can model/approximate $V(s)$ and $Q(s, a)$.

MARL is often the expansion of Single-Agent RL to trickier problems!

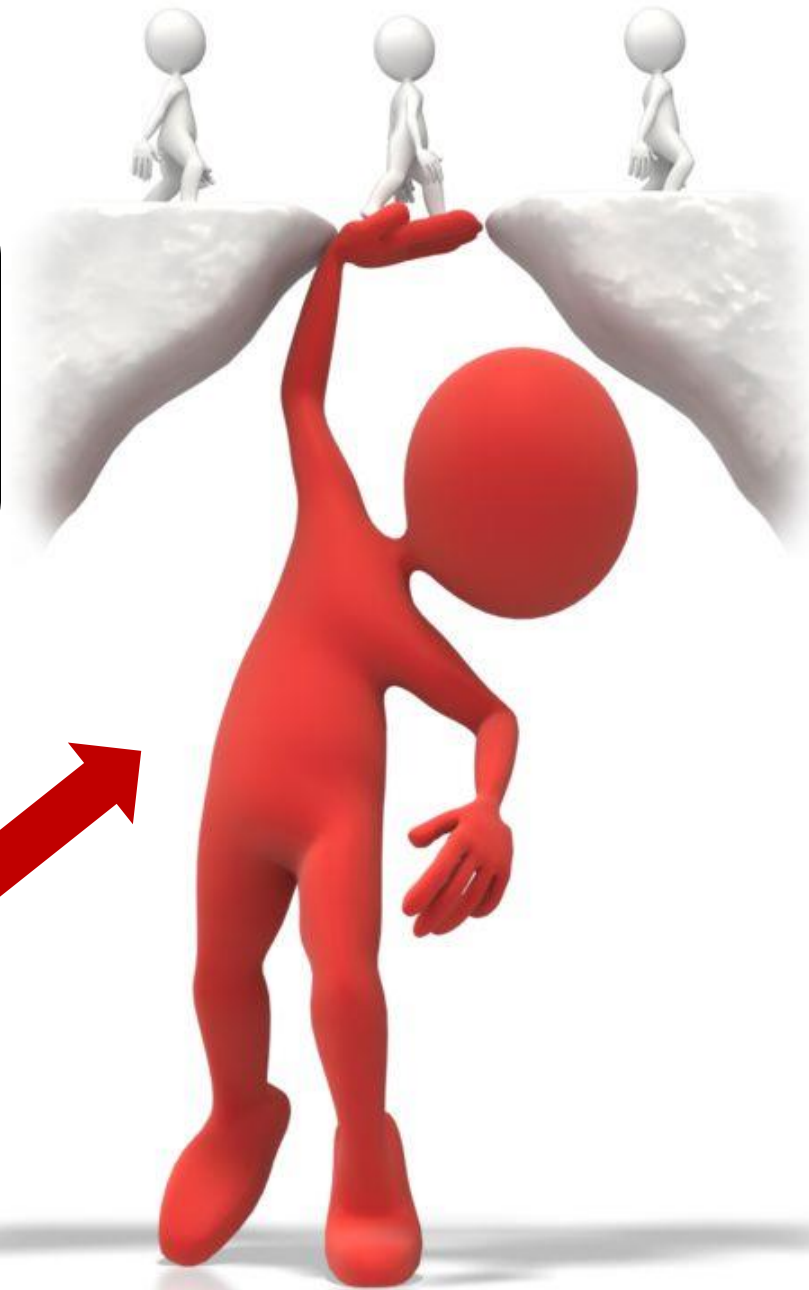
Single-Agent Learning	Multi-Agent Learning	Single-Agent Learning	Multi-Agent Learning	Single-Agent Learning	Multi-Agent Learning
Behavioral Cloning	Imitation-MA (2007) Coord-MAIL (2017)	Q-Learning	IQL (1993)	TD3	MA-TD3 (2020)
Advanced Behavioral Cloning	GMA-BC (2018)	DQN	MA-DQL (2015) VDN (2017) QMIX (2018)	Model-Based RL	MB-MARL (2022)
Multi-Armed Bandits	Game Theory (Normal-Form Games)	Policy-Gradient (REINFORCE)	Multi-Agent PG (2020)	Offline RL	MARL-IGL (2023)
Markov Decision Processes	Game Theory (Markov Games)	A2C/A3C	SEAC (2020)	RLHF	MARLHF (2024)
Exact Solutions for MDPs	MARL-Textbook	PPO	MAPPO (2022)	Inverse RL	MAIRL (2021)
Value-Based RL	MARL-Textbook	SAC	Dec-SAC (2021)	LLM Agents	Generative Agents (2023)
TD Learning	MARL-Textbook	DDPG	MADDPG (2020)

We could spend an **entire semester** talking about MARL! See NYU's [TRGY 8103](#).

Single Agent
Reinforcement
Learning

Multi-Agent
Reinforcement
Learning

Game Theory









Normal-Form Games (Bandits++)

Definition 2 (Normal-form game) *A normal-form game consists of:*

- *Finite set of agents $I = \{1, \dots, n\}$*
- *For each agent $i \in I$:*
 - *Finite set of actions A_i*
 - *Reward function $\mathcal{R}_i : A \rightarrow \mathbb{R}$, where $A = A_1 \times \dots \times A_n$*






Most Popular Normal-Form Game: Rock, Paper, Scissors!

- A single-shot, simultaneous, full-information game
- **Single shot:** Played once
- **Simultaneous:** Both actions are revealed at once
- **Full-information:** Everything about the game is known (in this case, the costs / rewards)

		Player 2		
				
Player 1				
				
				

Most Popular Normal-Form Game: Rock, Paper, Scissors!

- A single-shot, simultaneous, full-information game
- **Single shot:** Played once
- **Simultaneous:** Both actions are revealed at once
- **Full-information:** Everything about the game is known (in this case, the costs / rewards)

		Player 2		
				
Player 1		$(0,0)$	$(-1,1)$	$(1,-1)$
		$(1,-1)$	$(0,0)$	$(-1,1)$
		$(-1,1)$	$(1,-1)$	$(0,0)$

Normal-Form Games (Bandits++)

	R	P	S
R	0,0	-1,1	1,-1
P	1,-1	0,0	-1,1
S	-1,1	1,-1	0,0

(a) Rock-Paper-Scissors

Zero-Sum Game

	A	B
A	10	0
B	0	10

(b) Coordination Game

Common Reward

	C	D
C	-1,-1	-5,0
D	0,-5	-3,-3

(c) Prisoner's Dilemma

Mixed Reward

We won't go over the various ways to solve NFGs from evaluative feedback, but this is a well studied area of game theory!

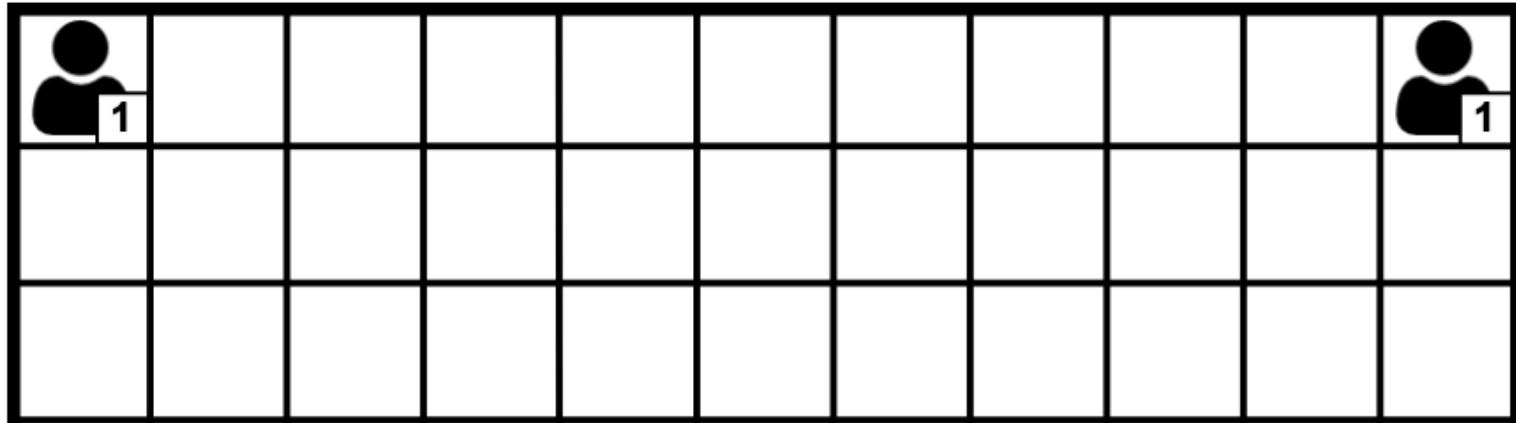
Markov Games (MDPs++)

Definition 3 (Stochastic game) *A stochastic game consists of:*

- *Finite set of agents $I = \{1, \dots, n\}$*
- *Finite set of states S . with subset of terminal states $\bar{S} \subset S$*
- *For each agent $i \in I$:*
 - *Finite set of actions A_i*
 - *Reward function $\mathcal{R}_i : S \times A \times S \rightarrow \mathbb{R}$, where $A = A_1 \times \dots \times A_n$*
- *State transition probability function $\mathcal{T} : S \times A \times S \rightarrow [0, 1]$ such that*

$$\forall s \in S, a \in A: \sum_{s' \in S} \mathcal{T}(s, a, s') = 1 \quad (3.1)$$

Markov Games (MDPs++)



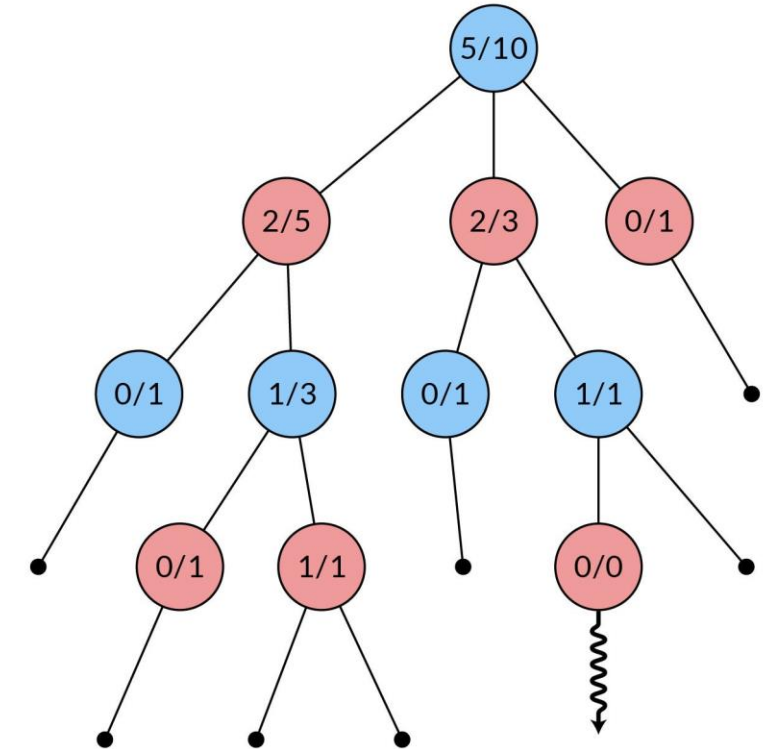
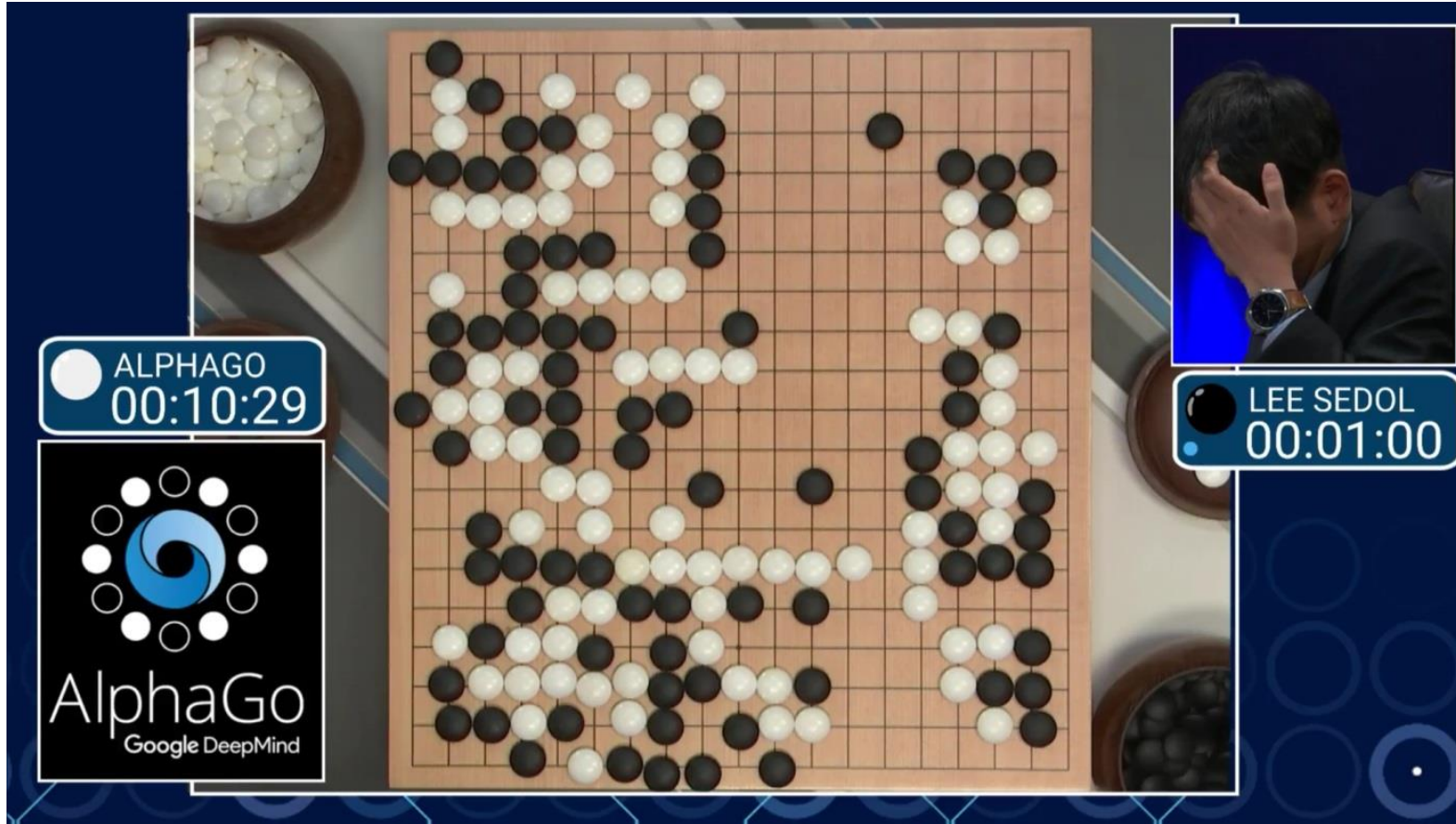
Each agent has four actions:
{UP, LEFT, DOWN, RIGHT}

A step forward in the game is
a **joint action**

$$a \in A = A_0 \times A_1 \times \dots \times A_N$$

- **Sequential Decision Making:** state changes over time (could be finite or infinite horizon).
- **Simultaneous:** Both agents take actions at the same time.
- **System Dynamics:** The actions of both agents (joint action) determine the subsequent state.

Competition (2-player games)



We already talked about strategies for learning turn-based 2-player games when we covered AlphaGo and AlphaZero!

Cooperation!

For today, let's focus on the *collaborative* case.

All agents still have independent reward signals

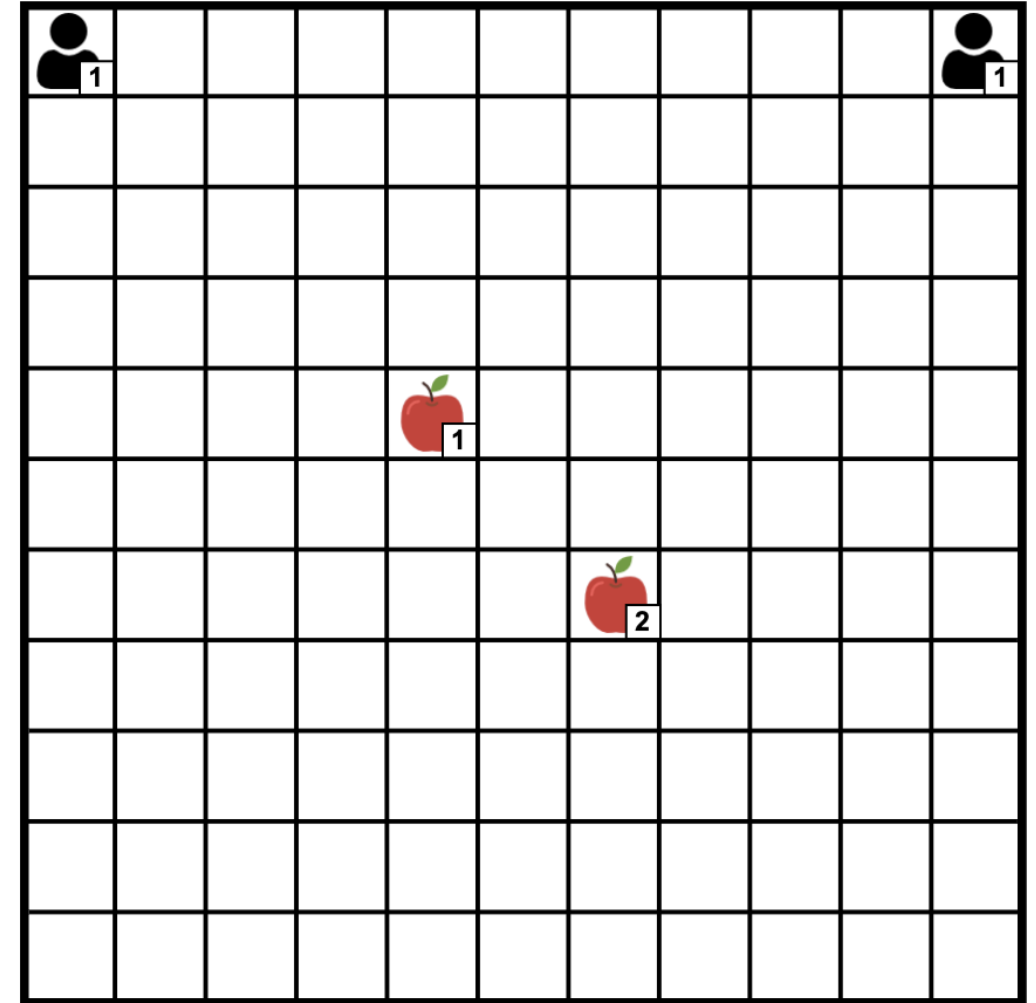
$$r_0, r_1, \dots, r_N$$

Our objective will be to **maximize the collective welfare** of the entire system.

$$\max_{(\pi_0, \pi_1, \dots, \pi_N)} E_{(\pi_0, \pi_1, \dots, \pi_N)} \left[\sum_{t=0}^{\infty} \sum_{i=0}^N \gamma^t R_i(s_t, \pi_i(s_t), s_{t+1}) \right]$$

Apple Picking Task

- N agents must collect all apples, but some apples cannot be picked alone and require teamwork!
- Each agent has a “level” that indicates how skilled they are at picking apples.
- Each apple has a “level” that indicates how difficult it is to pick that apple.
- An apple can only be collected if the sum of agent levels simultaneously picking it is greater than or equal to the level of the apple.
- Discrete action space with 6 actions.
 - $A = \{\text{UP, DOWN, LEFT, RIGHT, COLLECT, NO_OP}\}$
- Reward is +1 for all agents every time an apple is picked.



Simplest Solution: Convert MARL to RL

The simplest way to apply RL algorithms to multi-agent settings is to **reduce the problem to a single-agent problem** and apply traditional RL techniques.

Central Learning

- Apply one single-agent RL algorithm to **control all agents centrally**.
- A central policy is learned over the joint action space $A = (A_1 \times A_2 \times \dots \times A_N)$.

Independent Learning

- Apply single-agent RL algorithms to each agent **independently**.
- Agents do not explicitly consider or represent each other!

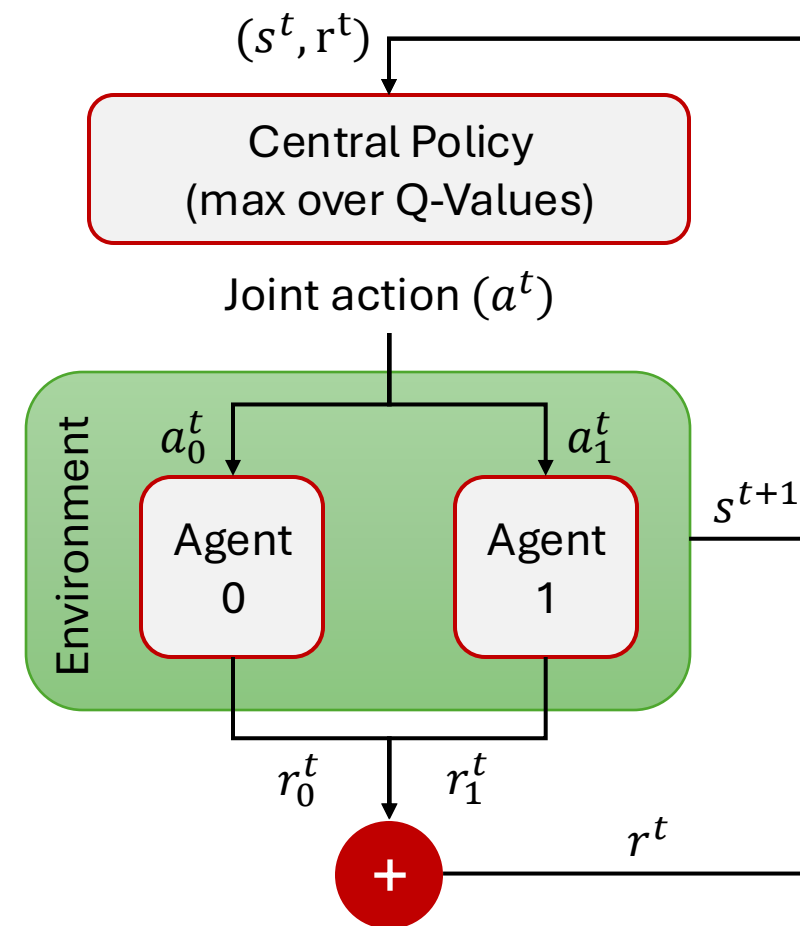
Centralized Learning: Algorithm

Algorithm 4 Central Q-learning (CQL) for stochastic games

- 1: Initialize: $Q(s, a) = 0$ for all $s \in S$ and $a \in A = A_1 \times \dots \times A_n$
- 2: Repeat for every episode:
- 3: **for** $t = 0, 1, 2, \dots$ **do**
- 4: Observe current state s^t
- 5: With probability ϵ : choose random joint action $a^t \in A$
- 6: Otherwise: choose joint action $a^t \in \arg \max_a Q(s^t, a)$
- 7: Apply joint action a^t , observe rewards r_1^t, \dots, r_n^t and next state s^{t+1}
- 8: Transform r_1^t, \dots, r_n^t into scalar reward r^t
- 9: $Q(s^t, a^t) \leftarrow Q(s^t, a^t) + \alpha [r^t + \gamma \max_{a'} Q(s^{t+1}, a') - Q(s^t, a^t)]$

We learn q-values over the space of joint-actions, so the algorithm doesn't know that we are training multiple agents.

Single-agent RL doesn't know how to process multiple reward signals, so we must transform the joint reward into a scalar reward.



Independent Learning: Algorithm

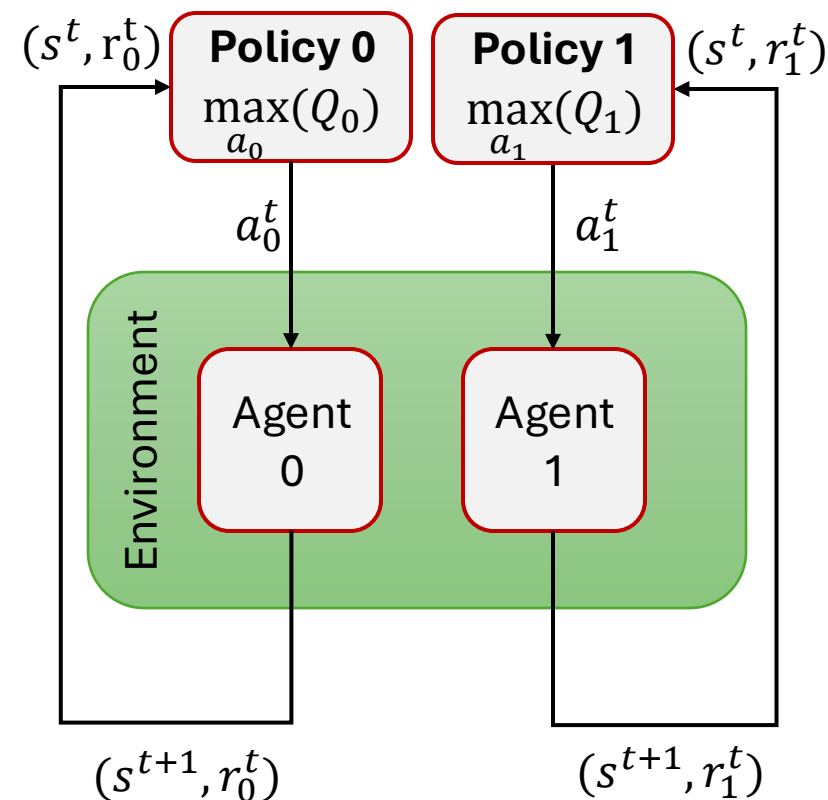
Algorithm 5 Independent Q-learning (IQL) for stochastic games

// Algorithm controls agent i

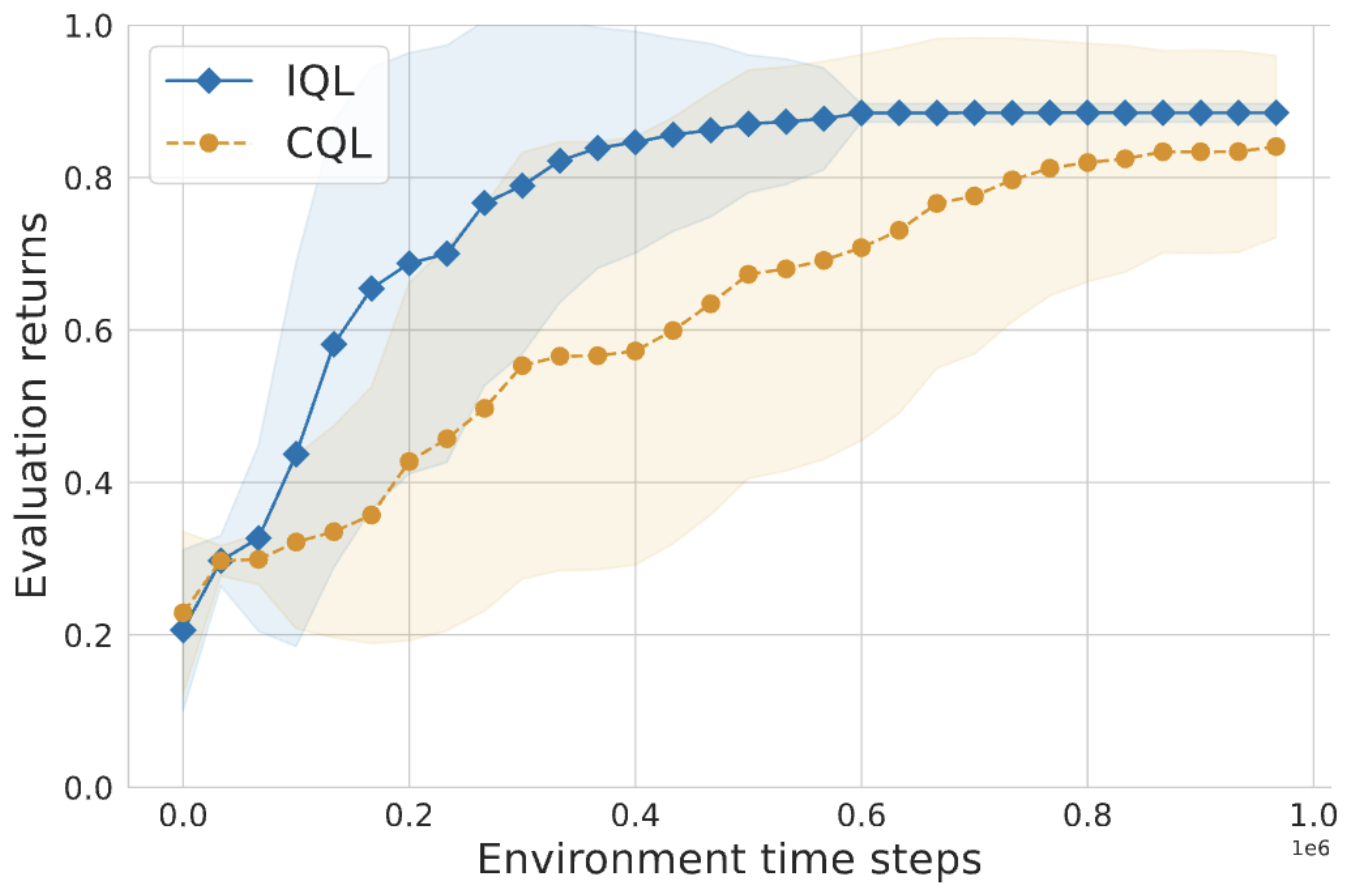
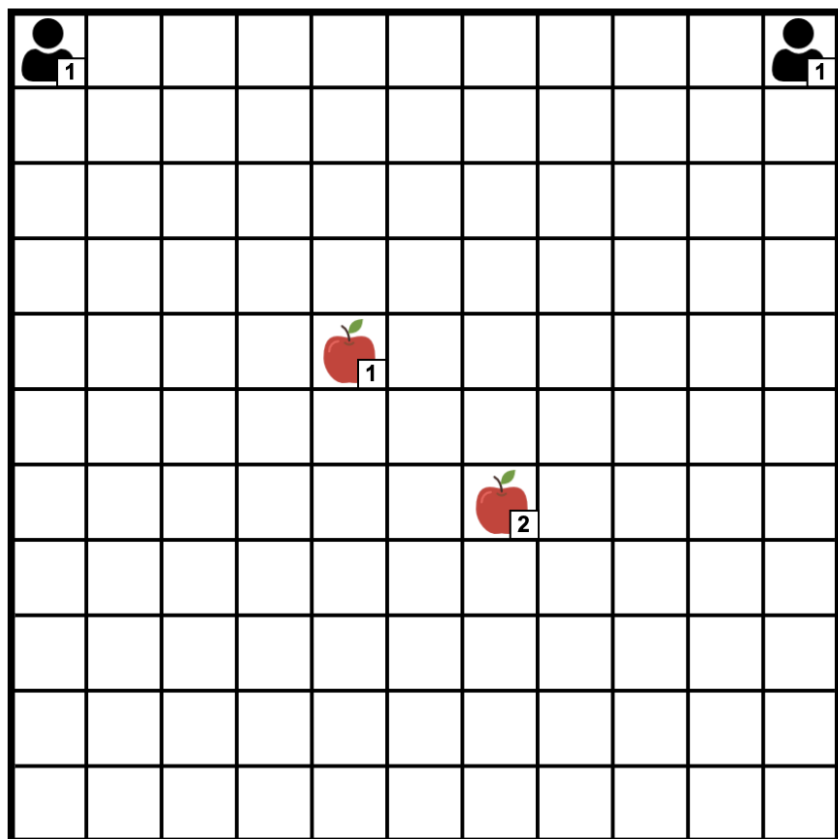
- 1: Initialize: $Q_i(s, a_i) = 0$ for all $s \in \mathcal{S}$, $a_i \in A_i$
- 2: Repeat for every episode:
- 3: **for** $t = 0, 1, 2, \dots$ **do**
- 4: Observe current state s^t
- 5: With probability ϵ : choose random action $a_i^t \in A_i$
- 6: Otherwise: choose action $a_i^t \in \arg \max_{a_i} Q_i(s^t, a_i)$
- 7: (meanwhile, other agents $j \neq i$ choose their actions a_j^t)
- 8: Observe own reward r_i^t and next state s^{t+1}
- 9: $Q_i(s^t, a_i^t) \leftarrow Q_i(s^t, a_i^t) + \alpha [r_i^t + \gamma \max_{a_i'} Q_i(s^{t+1}, a_i') - Q_i(s^t, a_i^t)]$

IQL doesn't have any information about the other agents.

It solely tries to model its own reward interactions in its Q table.



CQL + IQL: Results



Quiz: Why does IQL learn faster than CQL?

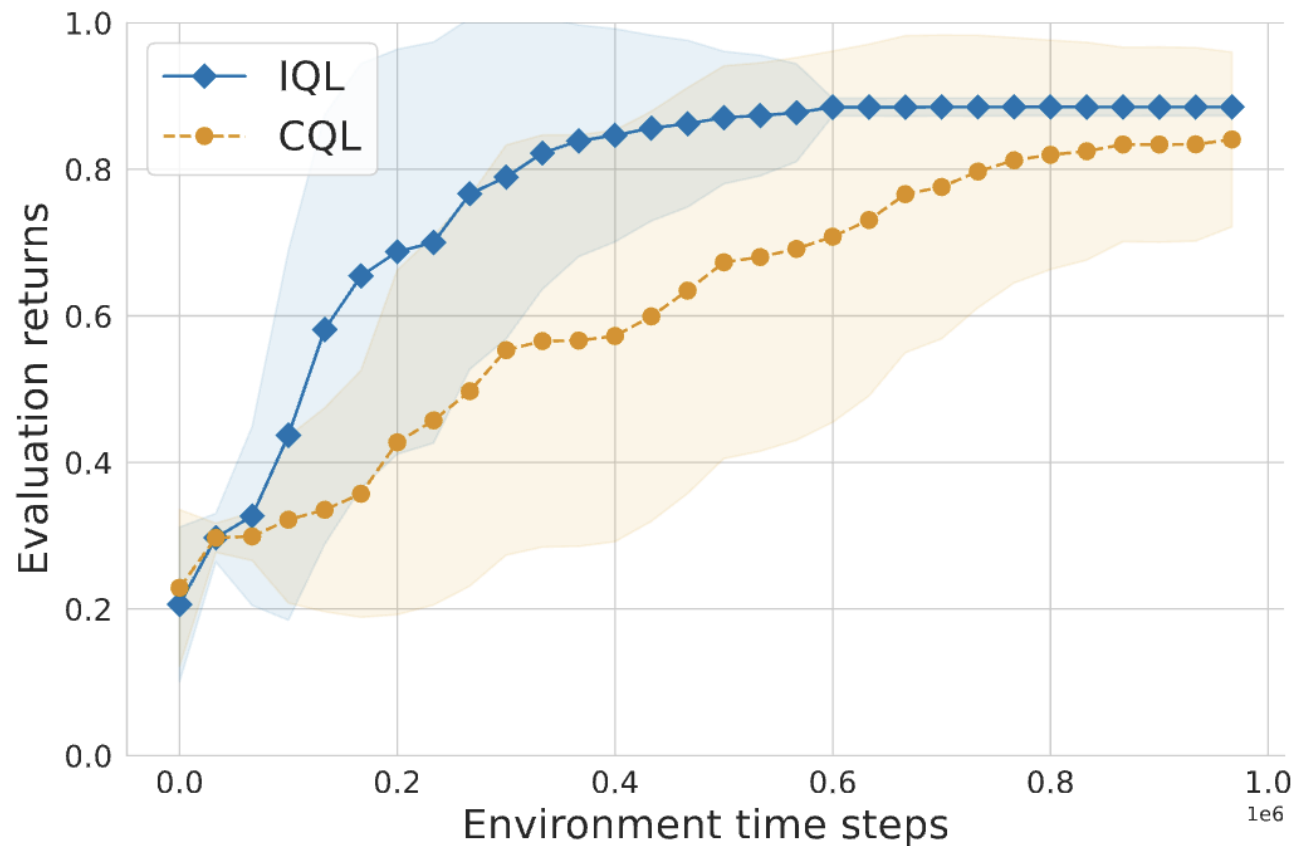
Algorithm 4 Central Q-learning (CQL) for stochastic games

- 1: Initialize: $Q(s, a) = 0$ for all $s \in S$ and $a \in A = A_1 \times \dots \times A_n$
- 2: Repeat for every episode:
- 3: **for** $t = 0, 1, 2, \dots$ **do**
- 4: Observe current state s^t
- 5: With probability ϵ : choose random joint action $a^t \in A$
- 6: Otherwise: choose joint action $a^t \in \arg \max_a Q(s^t, a)$
- 7: Apply joint action a^t , observe rewards r_1^t, \dots, r_n^t and next state s^{t+1}
- 8: Transform r_1^t, \dots, r_n^t into scalar reward r^t
- 9: $Q(s^t, a^t) \leftarrow Q(s^t, a^t) + \alpha [r^t + \gamma \max_{a'} Q(s^{t+1}, a') - Q(s^t, a^t)]$

Algorithm 5 Independent Q-learning (IQL) for stochastic games

// Algorithm controls agent i

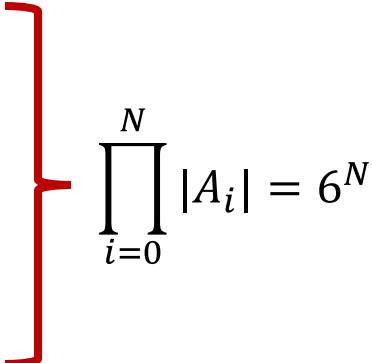
- 1: Initialize: $Q_i(s, a_i) = 0$ for all $s \in S, a_i \in A_i$
- 2: Repeat for every episode:
- 3: **for** $t = 0, 1, 2, \dots$ **do**
- 4: Observe current state s^t
- 5: With probability ϵ : choose random action $a_i^t \in A_i$
- 6: Otherwise: choose action $a_i^t \in \arg \max_{a_i} Q_i(s^t, a_i)$
- 7: (meanwhile, other agents $j \neq i$ choose their actions a_j^t)
- 8: Observe own reward r_i^t and next state s^{t+1}
- 9: $Q_i(s^t, a_i^t) \leftarrow Q_i(s^t, a_i^t) + \alpha [r_i^t + \gamma \max_{a_i'} Q_i(s^{t+1}, a_i') - Q_i(s^t, a_i^t)]$



Centralized Learning: Challenges

- Works great for low N! But exploration suffers exponentially as N grows.

Number of Agents in Apple Picking (N)	Size of Joint Action Space (A)
2	36
3	216
4	1296
5	7776


$$\prod_{i=0}^N |A_i| = 6^N$$

- Physical Limitations. In robotics/real world tasks we cannot assume that we have centralized control over all agents (or even if we do there are major communication/latency issues!)
- Requires the scalarization of reward signals, which is sometimes not possible.

Independent Learning

- It performs fairly compared to a lot of SOTA MARL algorithms!
- All single-agent algorithms have been independently applied to MARL.

Algorithm 17 Independent deep Q-networks

- 1: Initialize n value networks with random parameters $\theta_1, \dots, \theta_n$
 - 2: Initialize n target networks with parameters $\bar{\theta}_1 = \theta_1, \dots, \bar{\theta}_n = \theta_n$
 - 3: Initialize a replay buffer for each agent D_1, D_2, \dots, D_n
 - 4: **for** time step $t=0, 1, 2, \dots$ **do**
 - 5: Collect current observations o_1^t, \dots, o_n^t
 - 6: **for** agent $i=1, \dots, n$ **do**
 - 7: With probability ϵ : choose random action a_i^t
 - 8: Otherwise: choose $a_i^t \in \arg \max_{a_i} Q(h_i^t, a_i; \theta_i)$
 - 9: Apply actions (a_1^t, \dots, a_n^t) ; collect rewards r_1^t, \dots, r_n^t and next observations $o_1^{t+1}, \dots, o_n^{t+1}$
 - 10: **for** agent $i=1, \dots, n$ **do**
 - 11: Store transition $(h_i^t, a_i^t, r_i^t, h_i^{t+1})$ in replay buffers D_i
 - 12: Sample random mini-batch of B transitions $(h_i^k, a_i^k, r_i^k, h_i^{k+1})$ from D_i
 - 13: **if** s^{k+1} is terminal² **then**
 - 14: Targets $y_i^k \leftarrow r_i^k$
 - 15: **else**
 - 16: Targets $y_i^k \leftarrow r_i^k + \gamma \max_{a_i' \in A_i} Q(h_i^{k+1}, a_i'; \bar{\theta}_i)$
 - 17: Loss $\mathcal{L}(\theta_i) \leftarrow \frac{1}{B} \sum_{k=1}^B \left(y_i^k - Q(h_i^k, a_i^k; \theta_i) \right)^2$
 - 18: Update parameters θ_i by minimizing the loss $\mathcal{L}(\theta_i)$
 - 19: In a set interval, update target network parameters $\bar{\theta}_i$
-

Algorithm 18 Independent REINFORCE

- 1: Initialize n policy networks with random parameters ϕ_1, \dots, ϕ_n
 - 2: Repeat for every episode:
 - 3: **for** time step $t=0, 1, 2, \dots, T-1$ **do**
 - 4: Collect current observations o_1^t, \dots, o_n^t
 - 5: **for** agent $i=1, \dots, n$ **do**
 - 6: Sample actions a_i^t from $\pi(\cdot | h_i^t; \phi_i)$
 - 7: Apply actions (a_1^t, \dots, a_n^t) ; collect rewards r_1^t, \dots, r_n^t and next observations $o_1^{t+1}, \dots, o_n^{t+1}$
 - 8: **for** agent $i=1, \dots, n$ **do**
 - 9: Loss $\mathcal{L}(\phi_i) \leftarrow -\frac{1}{T} \sum_{t=0}^{T-1} \left(\sum_{\tau=t}^{T-1} \gamma^{\tau-t} r_i^\tau \right) \log \pi(a_i^t | h_i^t; \phi_i)$
 - 10: Update parameters ϕ_i by minimizing the loss $\mathcal{L}(\phi_i)$
-

Algorithm 19 Independent A2C with synchronous environments

- 1: Initialize n actor networks with random parameters ϕ_1, \dots, ϕ_n
 - 2: Initialize n critic networks with random parameters $\theta_1, \dots, \theta_n$
 - 3: Initialize K parallel environments
 - 4: **for** time step $t=0 \dots$ **do**
 - 5: Batch of observations for each agent and environment:
$$\begin{bmatrix} o_1^{t,1} \dots o_1^{t,K} \\ \vdots \\ o_n^{t,1} \dots o_n^{t,K} \end{bmatrix}$$
 - 6: Sample actions $\begin{bmatrix} a_1^{t,1} \dots a_1^{t,K} \\ \vdots \\ a_n^{t,1} \dots a_n^{t,K} \end{bmatrix} \sim \pi(\cdot | h_1^t; \phi_1), \dots, \pi(\cdot | h_n^t; \phi_n)$
 - 7: Apply actions; collect rewards $\begin{bmatrix} r_1^{t,1} \dots r_1^{t,K} \\ \vdots \\ r_n^{t,1} \dots r_n^{t,K} \end{bmatrix}$ and observations $\begin{bmatrix} o_1^{t+1,1} \dots o_1^{t+1,K} \\ \vdots \\ o_n^{t+1,1} \dots o_n^{t+1,K} \end{bmatrix}$
 - 8: **for** agent $i=1, \dots, n$ **do**
 - 9: **if** $s^{t+1,k}$ is terminal **then**
 - 10: Advantage $Adv(h_i^{t,k}, a_i^{t,k}) \leftarrow r_i^{t,k} - V(h_i^{t,k}; \theta_i)$
 - 11: Critic target $y_i^{t,k} \leftarrow r_i^{t,k}$
 - 12: **else**
 - 13: Advantage $Adv(h_i^{t,k}, a_i^{t,k}) \leftarrow r_i^{t,k} + \gamma V(h_i^{t+1,k}; \theta_i) - V(h_i^{t,k}; \theta_i)$
 - 14: Critic target $y_i^{t,k} \leftarrow r_i^{t,k} + \gamma V(h_i^{t+1,k}; \theta_i)$
 - 15: Actor loss $\mathcal{L}(\phi_i) \leftarrow \frac{1}{K} \sum_{k=1}^K Adv(h_i^{t,k}, a_i^{t,k}) \log \pi(a_i^{t,k} | h_i^{t,k}; \phi_i)$
 - 16: Critic loss $\mathcal{L}(\theta_i) \leftarrow \frac{1}{K} \sum_{k=1}^K \left(y_i^{t,k} - V(h_i^{t,k}; \theta_i) \right)^2$
 - 17: Update parameters ϕ_i by minimizing the actor loss $\mathcal{L}(\phi_i)$
 - 18: Update parameters θ_i by minimizing the critic loss $\mathcal{L}(\theta_i)$
-

Independent Learning: Challenges

- Much better scaling! Can be deployed easily on physical systems!
- Suffers from *non-stationary* environment dynamics.

Single Agent Dynamics

$$P(s' | s, a) = T(s' | s, a)$$

Probability of subsequent states stays constant over time.

Multi Agent Dynamics

$$P_i(s' | s, a_i) = \sum_{\mathbf{a}_{-i}} T(s' | s, \langle a_i, \mathbf{a}_{-i} \rangle) \left[\prod_{j \neq i} \pi_j(a_j | s^t) \right]$$

Changes to environment state are not guaranteed to be the result of my actions

- Depending on reward formulation, it suffers from poor *credit assignment*.

The Credit Assignment Problem

Temporal Credit Assignment

Which of my **past actions** were responsible for my earned reward?

Multi-Agent Credit Assignment

Which **agent** was responsible for my earned reward?

Recall that each agent its own reward signal

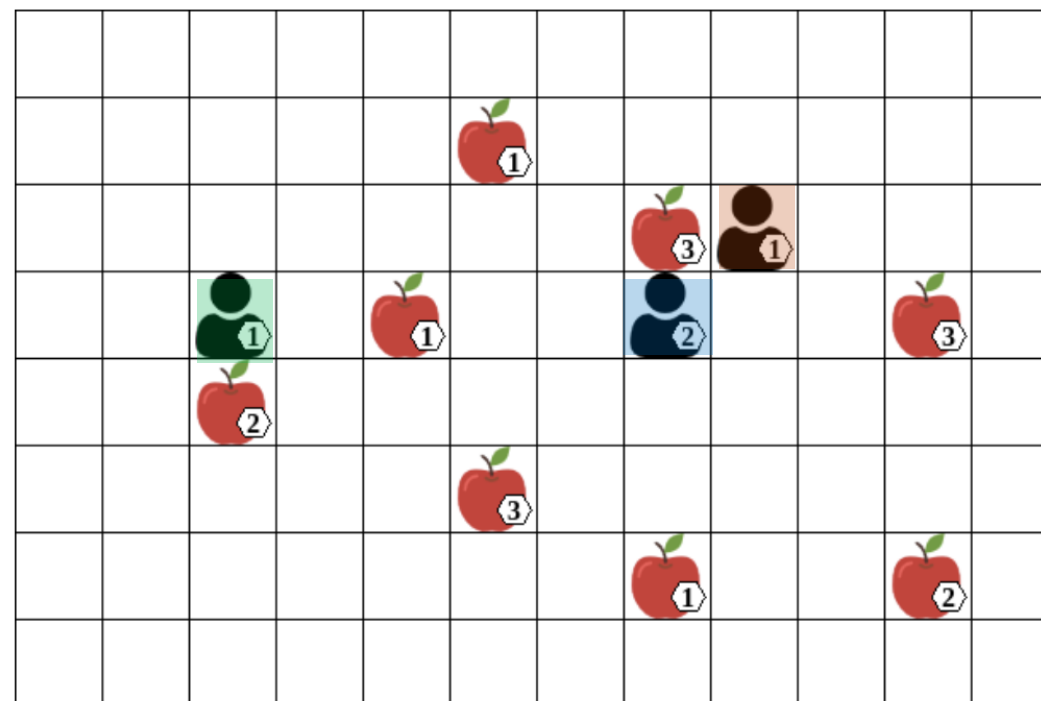
$$r_0, r_1, \dots, r_N$$

What happens when we shift to a common reward?
(i.e. all agents receive the same reward signal at time t)

$$r_0 = r_1 = \dots = r_N$$

Now the **green agent** gets some reward for the apple
that the **orange** and **blue** agents did!

Tasks are often formulated as common reward!



Middle Ground? CTDE Methods!

Can we create MARL algorithms that can still be deployed on a decentralized system while overcoming obstacles in scaling, stationarity, and credit assignment?

Independent Learning: Decentralized Training, Decentralized Execution

Central Learning: Centralized Training, Centralized Execution

Centralized Training, Decentralized Execution.

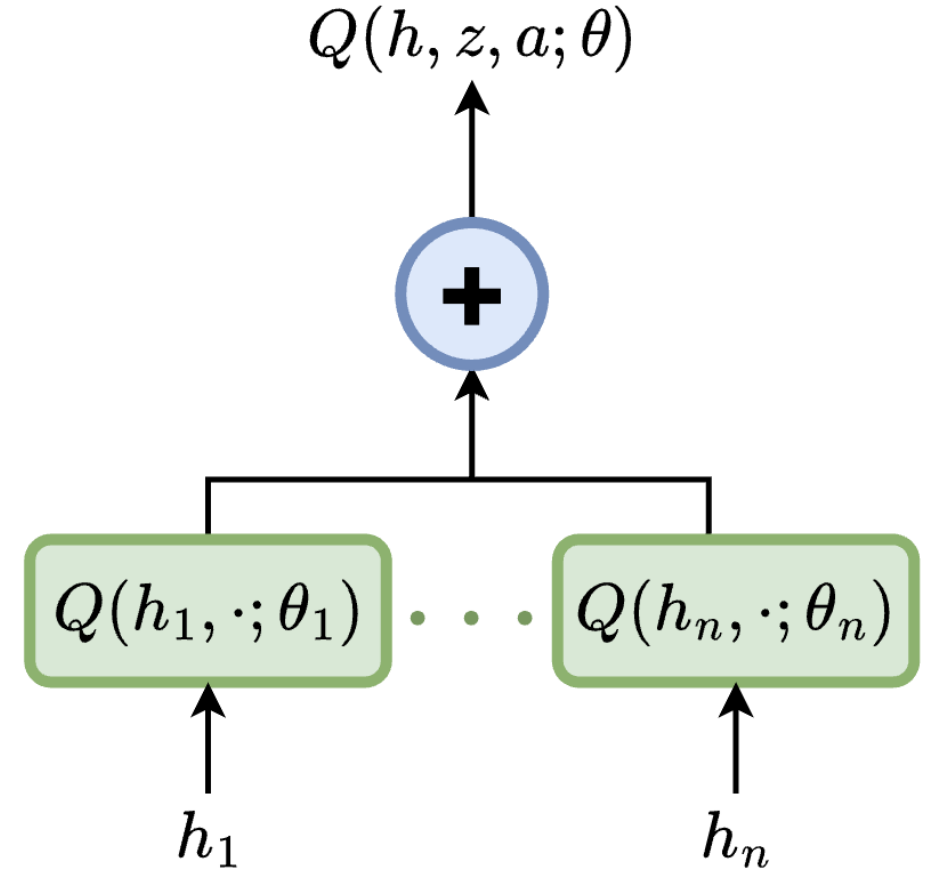
Utilize *privileged centralized information* at training time to determine how agent's actions affect the rewards of the system and share experiences between agents.

At execution time, ensure that all each policy can independently execute, with zero dependencies on centralized information.

Value Decomposition Networks

Algorithm 21 Value decomposition networks (VDN)

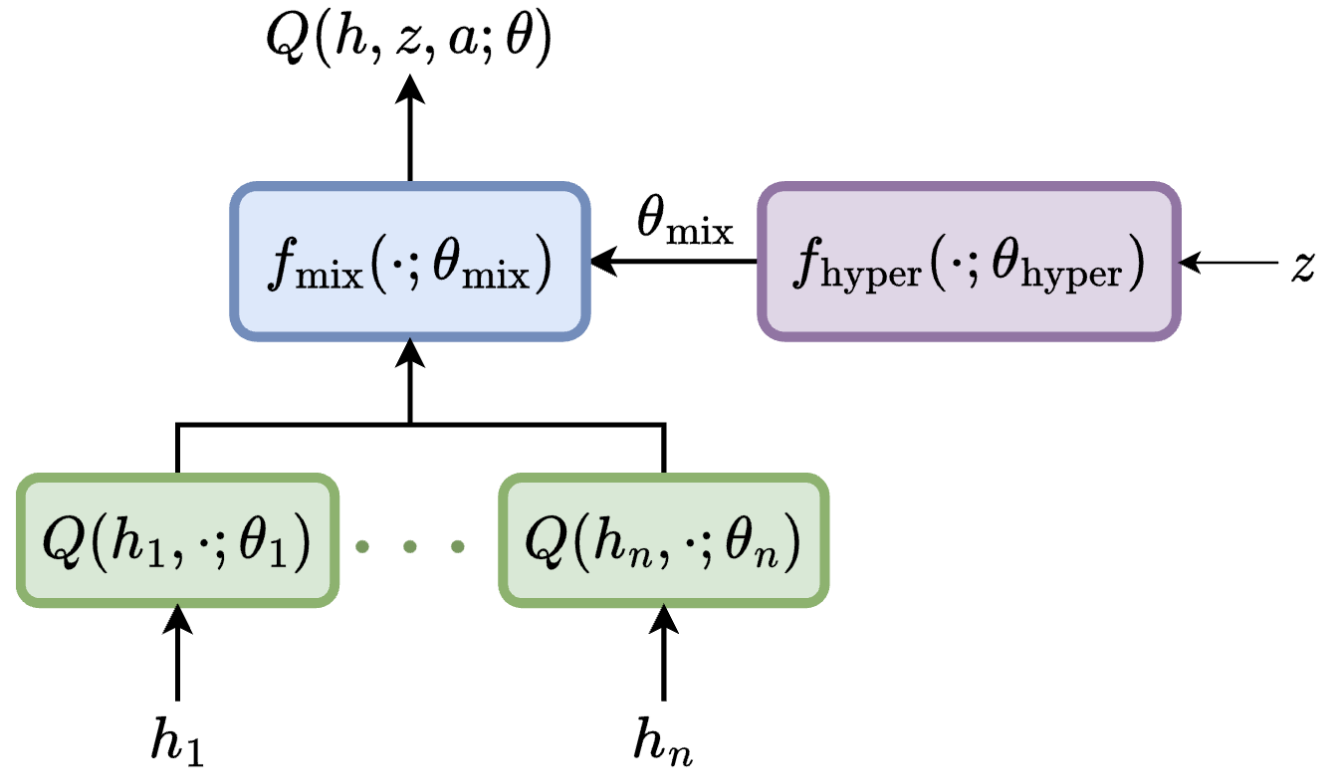
- 1: Initialize n utility networks with random parameters $\theta_1, \dots, \theta_n$
 - 2: Initialize n target networks with parameters $\bar{\theta}_1 = \theta_1, \dots, \bar{\theta}_n = \theta_n$
 - 3: Initialize a shared replay buffer D
 - 4: **for** time step $t = 0, 1, 2, \dots$ **do**
 - 5: Collect current observations o_1^t, \dots, o_n^t
 - 6: **for** agent $i = 1, \dots, n$ **do**
 - 7: With probability ϵ : choose random action a_i^t
 - 8: Otherwise: choose $a_i^t \in \arg \max_{a_i} Q(h_i^t, a_i; \theta_i)$
 - 9: Apply actions; collect shared reward r^t and next observations $o_1^{t+1}, \dots, o_n^{t+1}$
 - 10: Store transition (h^t, a^t, r^t, h^{t+1}) in shared replay buffer D
 - 11: Sample mini-batch of B transitions (h^k, a^k, r^k, h^{k+1}) from D
 - 12: **if** s^{k+1} is terminal **then**
 - 13: Targets $y^k \leftarrow r^k$
 - 14: **else**
 - 15: Targets $y^k \leftarrow r^k + \gamma \sum_{i \in I} \max_{a'_i \in A_i} Q(h_i^{k+1}, a'_i; \bar{\theta}_i)$
 - 16: Loss $\mathcal{L}(\theta) \leftarrow \frac{1}{B} \sum_{k=1}^B \left(y^k - \sum_{i \in I} Q(h_i^k, a_i^k; \theta_i) \right)^2$
 - 17: Update parameters θ by minimizing the loss $\mathcal{L}(\theta)$
 - 18: In a set interval, update target network parameters $\bar{\theta}_i$ for each agent i
-



QMIX: Non-Linear VDN

Algorithm 22 QMIX

- 1: Initialize n utility networks with random parameters $\theta_1, \dots, \theta_n$
- 2: Initialize n target networks with parameters $\bar{\theta}_1 = \theta_1, \dots, \bar{\theta}_n = \theta_n$
- 3: Initialize hypernetwork with random parameters θ_{hyper}
- 4: Initialize a shared replay buffer D
- 5: **for** time step $t=0, 1, 2, \dots$ **do**
- 6: Collect current centralized information z^t and observations o_1^t, \dots, o_n^t
- 7: **for** agent $i=1, \dots, n$ **do**
- 8: With probability ϵ : choose random action a_i^t
- 9: Otherwise: choose $a_i^t \in \arg \max_{a_i} Q(h_i^t, a_i; \theta_i)$
- 10: Apply actions; collect shared reward r^t , next centralized information z^{t+1} and observations $o_1^{t+1}, \dots, o_n^{t+1}$
- 11: Store transition $(h^t, z^t, a^t, r^t, h^{t+1}, z^{t+1})$ in shared replay buffer D
- 12: Sample mini-batch of B transitions $(h^k, z^k, a^k, r^k, h^{k+1}, z^{k+1})$ from D
- 13: **if** s^{k+1} is terminal **then**
- 14: Targets $y^k \leftarrow r^k$
- 15: **else**
- 16: Mixing parameters $\theta_{\text{mix}}^{k+1} \leftarrow f_{\text{hyper}}(z^{k+1}; \theta_{\text{hyper}})$
- 17: Targets $y^k \leftarrow r^k + \gamma f_{\text{mix}} \begin{pmatrix} \max_{a_1'} Q(h_1^{k+1}, a_1'; \bar{\theta}_1), \\ \vdots \\ \max_{a_n'} Q(h_n^{k+1}, a_n'; \bar{\theta}_n) \end{pmatrix}; \theta_{\text{mix}}^{k+1}$
- 18: Mixing parameters $\theta_{\text{mix}}^k \leftarrow f_{\text{hyper}}(z^k; \theta_{\text{hyper}})$
- 19: Value estimates $Q(h^k, z^k, a^k; \theta) \leftarrow f_{\text{mix}}(Q(h_1^k, a_1^k; \theta_1), \dots, Q(h_n^k, a_n^k; \theta_n); \theta_{\text{mix}}^k)$
- 20: Loss $\mathcal{L}(\theta) \leftarrow \frac{1}{B} \sum_{k=1}^B \left(y^k - Q(h^k, z^k, a^k; \theta) \right)^2$
- 21: Update parameters θ by minimizing the loss $\mathcal{L}(\theta)$
- 22: In a set interval, update target network parameters $\bar{\theta}_i$ for each agent i



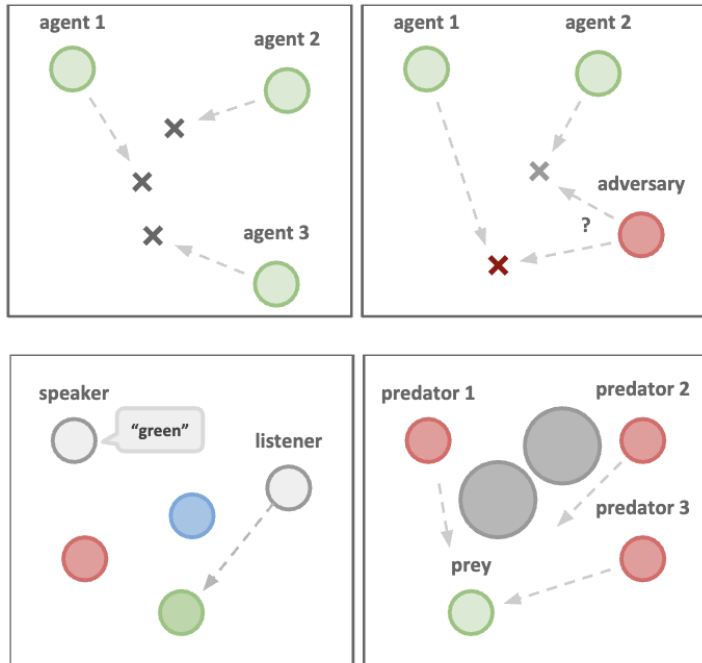
Other CTDE Extensions

Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments

Ryan Lowe*
McGill University
OpenAI

Yi Wu*
UC Berkeley

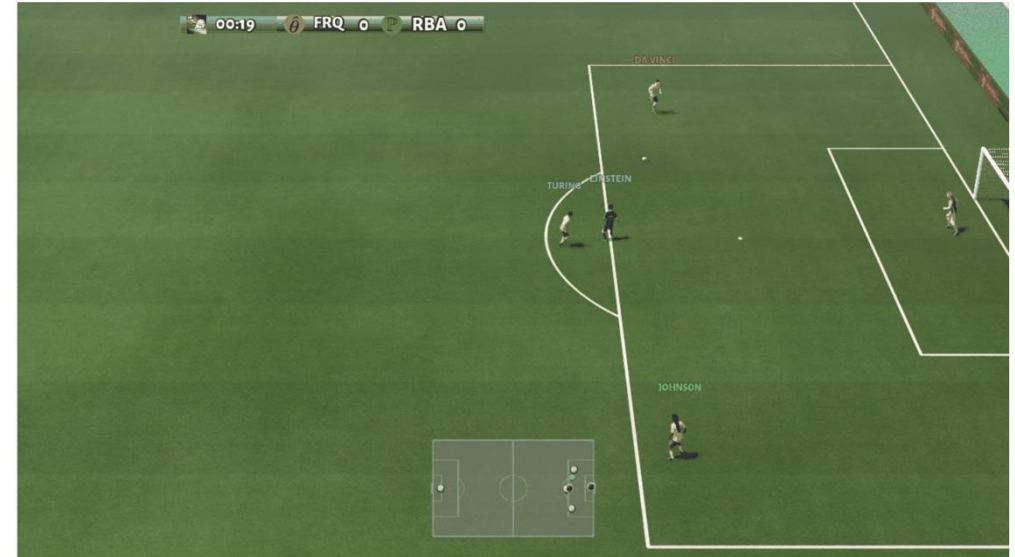
Aviv Tamar
UC Berkeley



The Surprising Effectiveness of PPO in Cooperative Multi-Agent Games

Chao Yu^{1‡}, Akash Velu^{2‡*}, Eugene Vinitsky^{2b}, Jiaxuan Gao¹,
Yu Wang^{1b}, Alexandre Bayen², Yi Wu^{13b}

¹ Tsinghua University ² University of California, Berkeley ³ Shanghai Qi Zhi Institute
[‡]zoe-yuchao@gmail.com, [‡]akashvelu@berkeley.edu



Current Benchmarks and Challenges

Where to Deploy?

What Challenges require Multi-Agent Learning right now?

How does deployment difficulty scale?



Current Benchmarks and Challenges

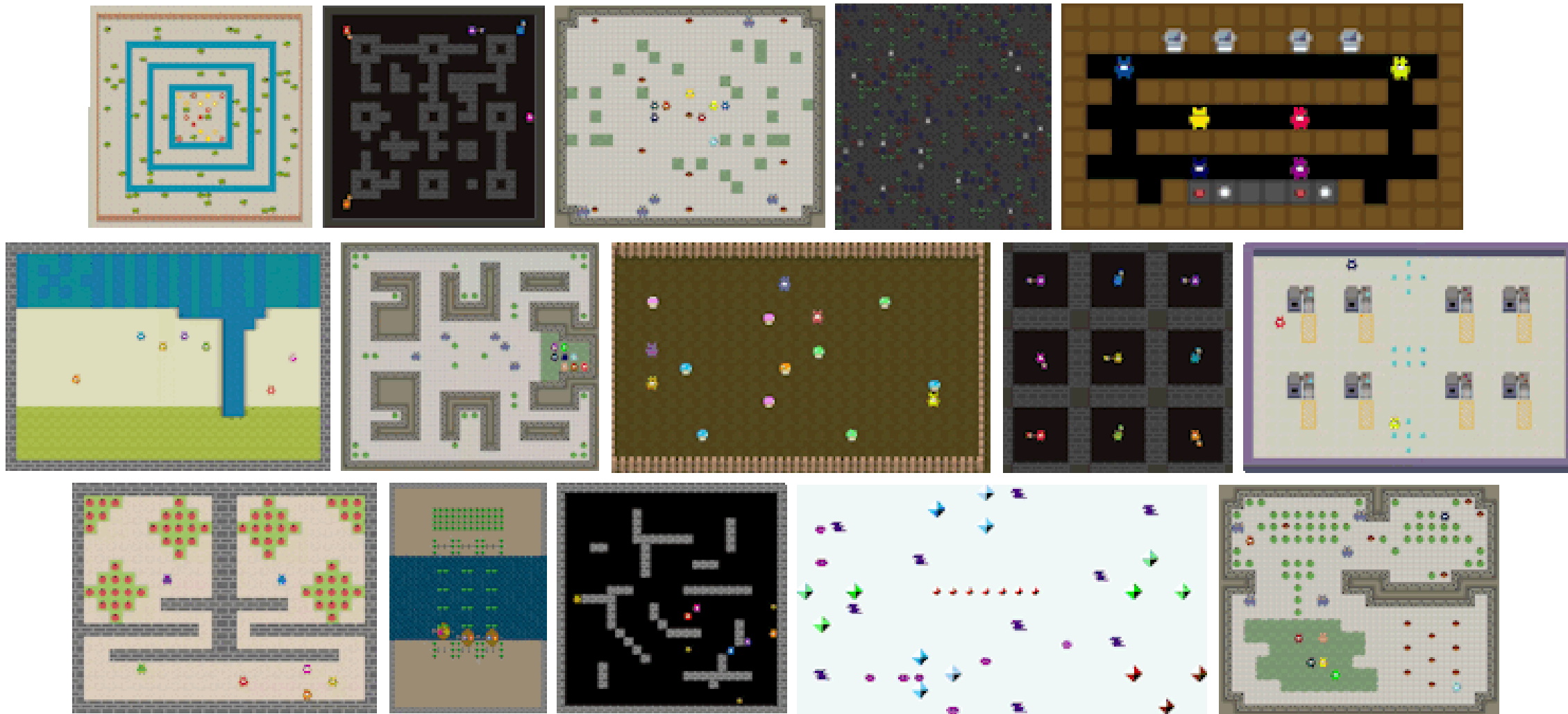
Article | Published: 30 October 2019

Grandmaster level in StarCraft II using multi-agent reinforcement learning

[Oriol Vinyals](#) ✉, [Igor Babuschkin](#), [Wojciech M. Czarnecki](#), [Michaël Mathieu](#), [Andrew Dudzik](#), [Junyoung Chung](#), [David H. Choi](#), [Richard Powell](#), [Timo Ewalds](#), [Petko Georgiev](#), [Junhyuk Oh](#), [Dan Horgan](#), [Manuel Kroiss](#), [Ivo Danihelka](#), [Aja Huang](#), [Laurent Sifre](#), [Trevor Cai](#), [John P. Agapiou](#), [Max Jaderberg](#), [Alexander S. Vechnevets](#), [Rémi Leblond](#), [Tobias Pohlen](#), [Valentin Dalibard](#), [David Budden](#), ... [David Silver](#) ✉



Current Benchmarks and Challenges

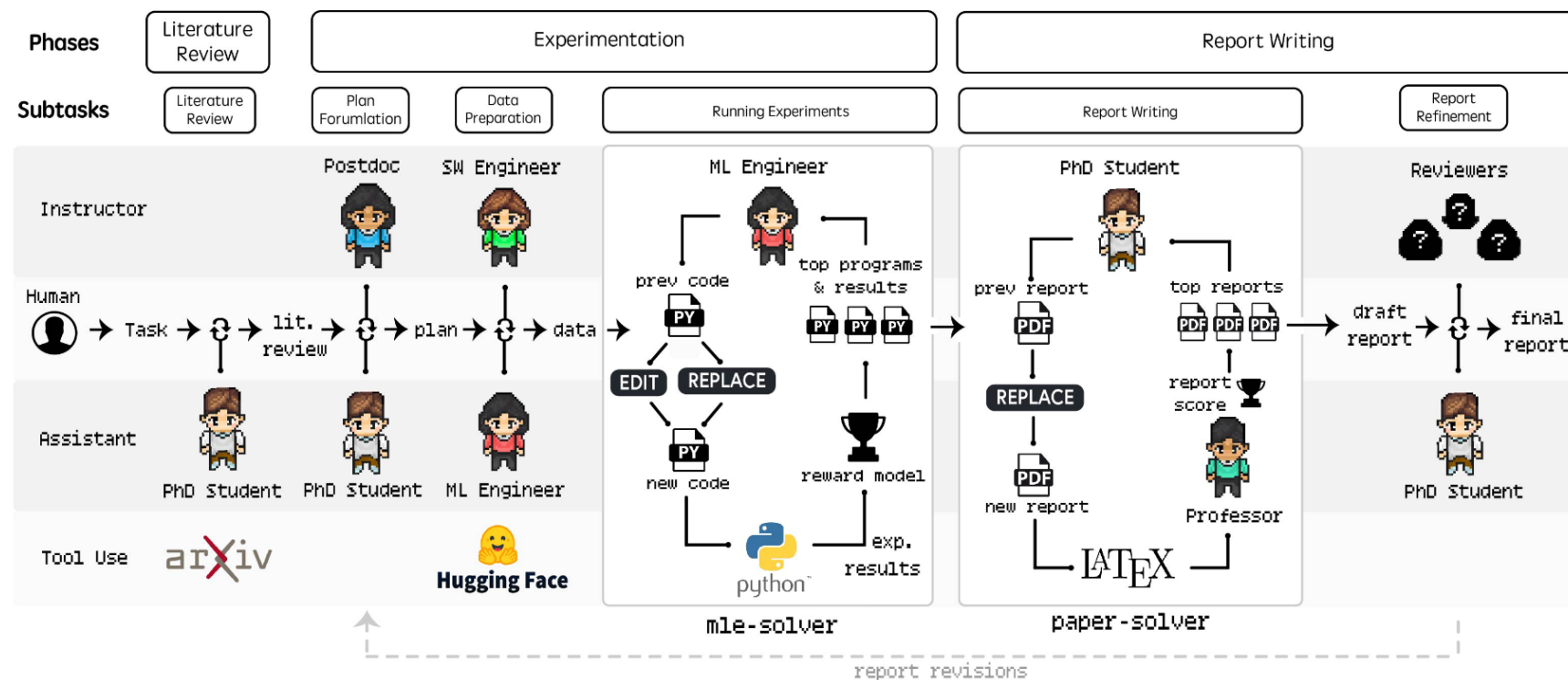


Current Explorations: LLMs

Agent Laboratory: Using LLM Agents as Research Assistants



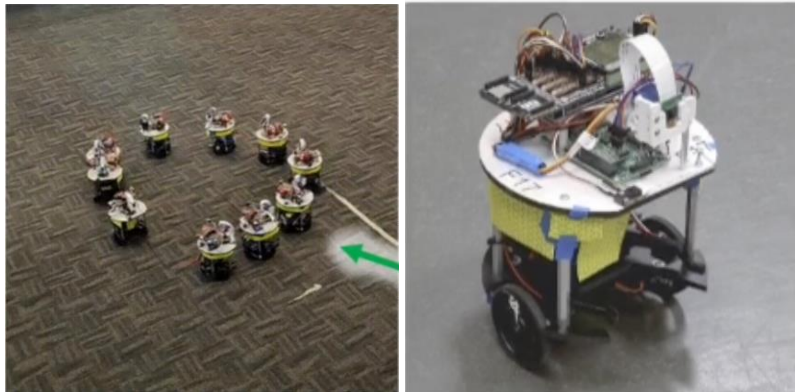
Agent Laboratory: Using LLM Agents as Research Assistants



Current Explorations: LLMs

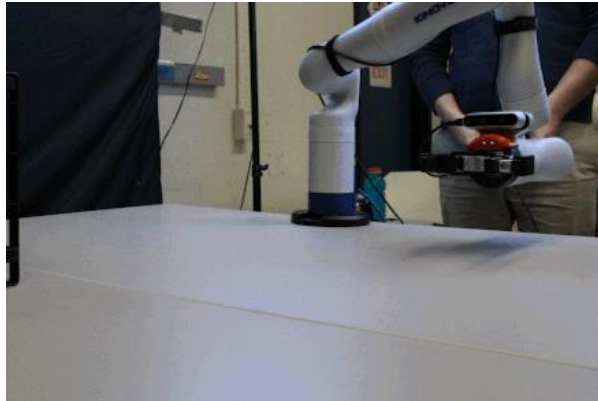


What does my PhD Research Look Like?



Robot Swarm Learning + Sim2Real

[ICONS'24], [CoRR'24], [AAMAS'25]



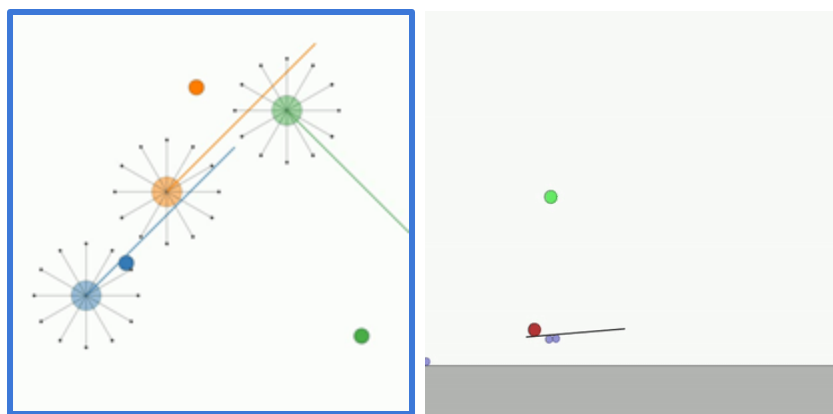
Assistive Robotics

[HRI'25], [SciData'24]



X-Embodiment Learning

[RLJ'24]



Multi-Agent Reinforcement Learning

[in preparation]



Repr. Learning + Behavior Discovery

[GECCO'23], [MRS'23] [AAMAS'25]

Multi-Agent Learning
is still a young field
and there's lots of
exciting work to do!



Let's chat!
c.mattson@utah.edu