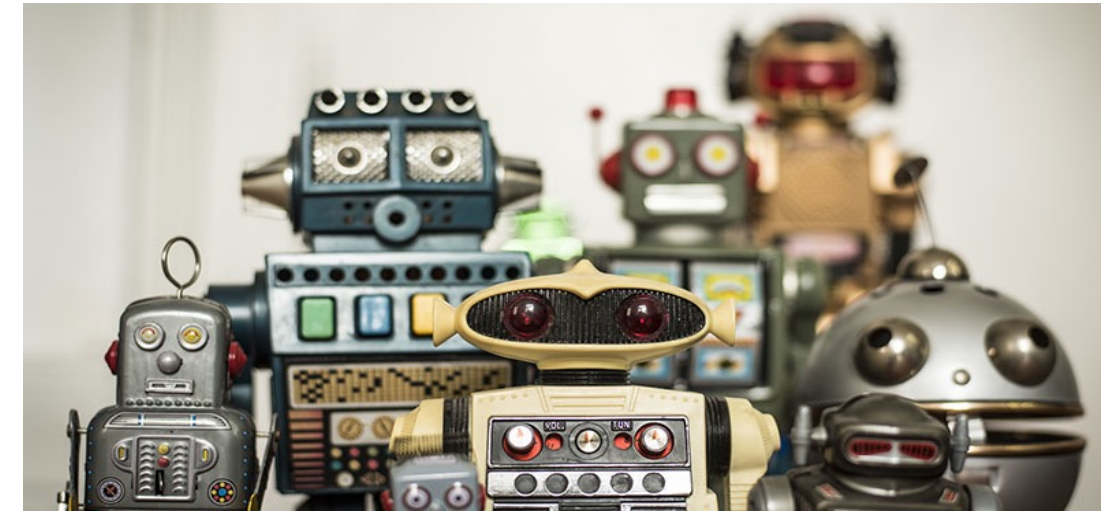
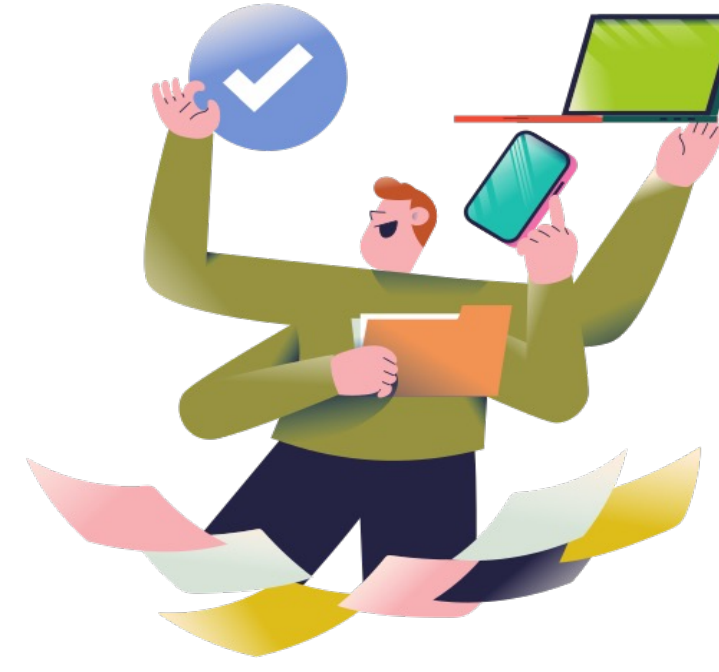


Research Question



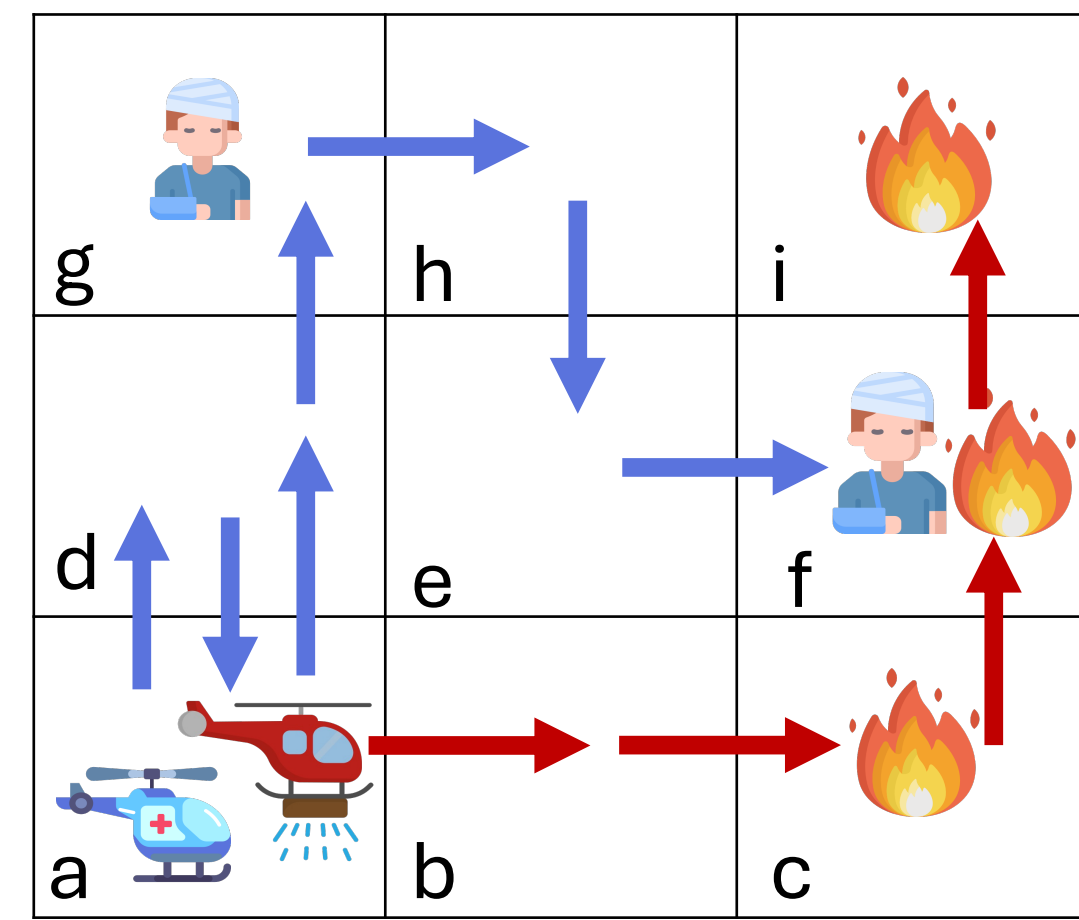
Multi-Objective

Multi-Agent

How to shape
Optimal Reward Functions?

Motivating Example: A wildfire scenario

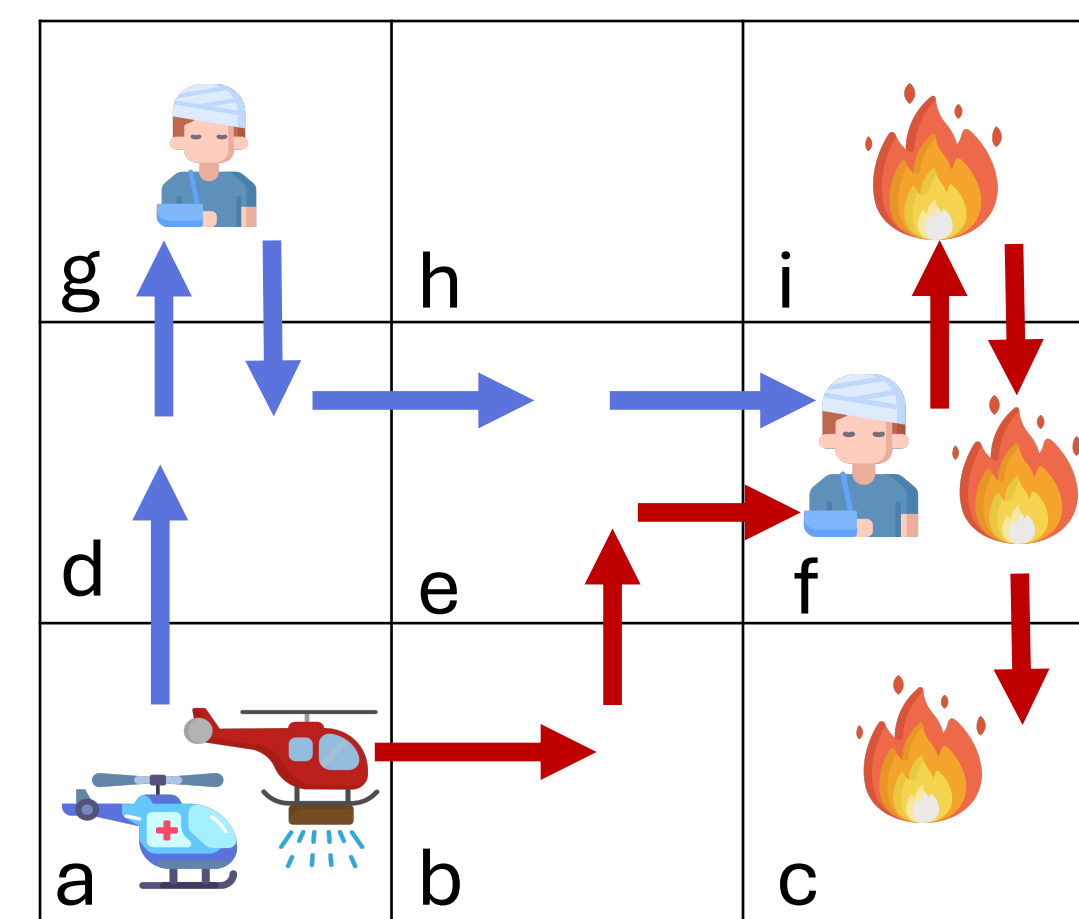
A **possible** control policy with reward function R1



R1(Extinguish fire) = +50
R1(Rescue victim) = +10
R1(Our of range) = -100
R1(MED in fire) = -100

a-b-c-f-i
a-d-a-d-g-h-e-f

An **optimal** control policy with reward function R2



R2(Extinguish fire) = +10
R2(Rescue victim) = +50
R2(Our of range) = -100
R2(MED in fire) = -100

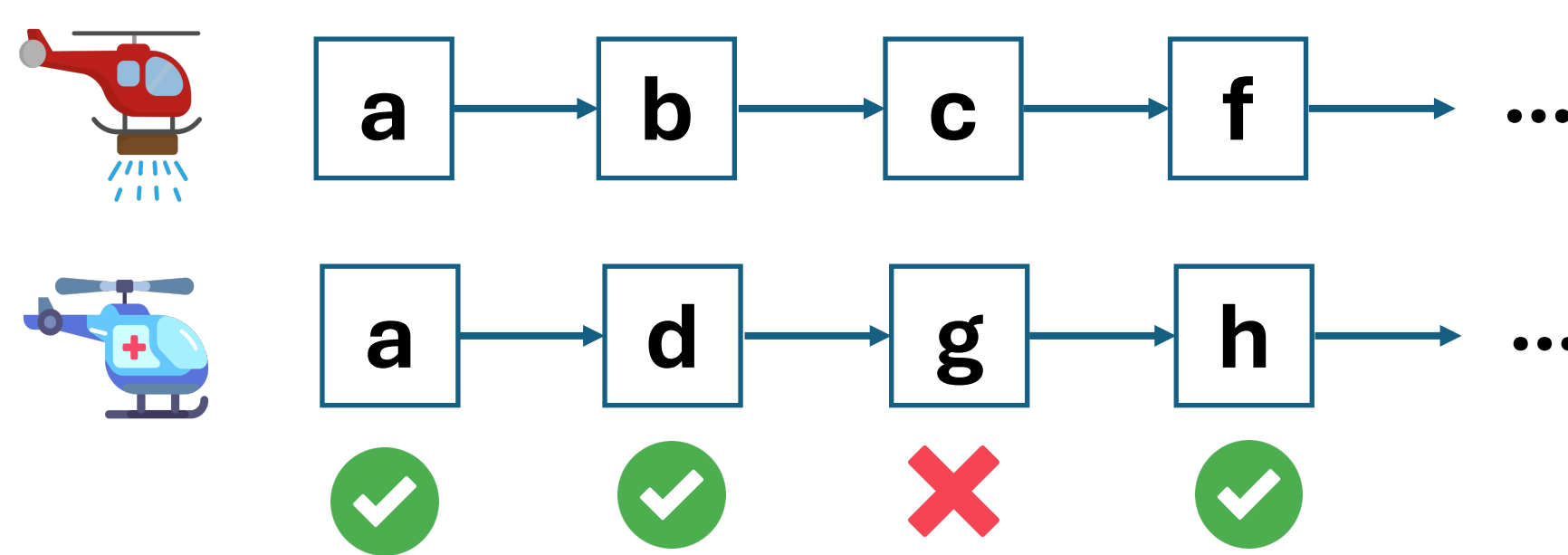
a-d-e-f-i-f-c
a-d-g-d-e-f

Specifications as Hyperproperties

- Hyperproperties** characterize requirements over sets of execution traces, allowing the specification of behaviors of multi-agent.

Example:

“Is the distance between **FF** and **Med** always less than 3 cells?”



- We use hyperproperties, expressed in the temporal logic **HyperLTL**, to achieve **specification-guided RL** for multi-agent w.r.t multi-objective and relational constraints.

The HyperLTL specification for the wildfire scenario:

$$\varphi_{\text{Rescue}} \triangleq \forall \tau_1. \exists \tau_2. (\psi_{\text{fire}} \wedge \psi_{\text{save}} \wedge \psi_{\text{dist}} \wedge \psi_{\text{safe}})$$

(Extinguish fire) $O_1 : \psi_{\text{fire}} \triangleq \Diamond(i_{\tau_1}) \wedge \Diamond(f_{\tau_1}) \wedge \Diamond(c_{\tau_1})$

(Rescue victim) $O_2 : \psi_{\text{save}} \triangleq \Diamond(j_{\tau_2}) \wedge \Diamond(f_{\tau_2})$

(Our of range) $C_1 : \psi_{\text{dist}} \triangleq \Box(|\text{Location}_{\tau_1} - \text{Location}_{\tau_2}| < 3)$

(MED in fire) $C_2 : \psi_{\text{safe}} \triangleq (\neg i_{\tau_2} \mathcal{U} i_{\tau_1}) \wedge (\neg f_{\tau_2} \mathcal{U} f_{\tau_1}) \wedge (\neg c_{\tau_2} \mathcal{U} c_{\tau_1})$

HYPRL: Reinforcement Learning of Control Policies for Hyperproperties

Tzu-Han Hsu*, Arshia Rafieioskouei*, Borzoo Bonakdarpour
Michigan State University



Problem Statement

Given an MDP \mathcal{M} with unknown transitions and a HyperLTL formula φ of the form $\mathbb{Q}_1 \tau_1 \dots \mathbb{Q}_n \tau_n. \psi$, our goal is to identify a tuple of n policies $\langle \pi_1^*, \dots, \pi_n^* \rangle$, such that:

$$\langle \pi_i^* \rangle_{i \in \{1, \dots, n\}} \in \left[\arg \max_{\langle \pi_i \rangle} \mathbb{P} \left[\langle \text{Traces}(\mathcal{Z}_{\tau_i} \sim \mathcal{D}_{\pi_i}) \rangle \models \varphi \right] \right]_{i \in \{1, \dots, n\}}$$

Our Solutions to the Main Challenges

- We apply **Skolemization** to resolve **quantifier alternations** in a HyperLTL formula.

$$\text{Skolem}(\varphi) = \underbrace{\exists \mathbf{f}_i(\tau_{i_1}, \dots, \tau_{i_{|\mathbb{Q}_i^\forall|}})}_{\text{for each } i \in \mathbb{Q}^\exists} \cdot \underbrace{\forall \tau_j}_{\text{for each } j \in \mathbb{Q}^\forall} \cdot \text{Skolem}(\psi)$$

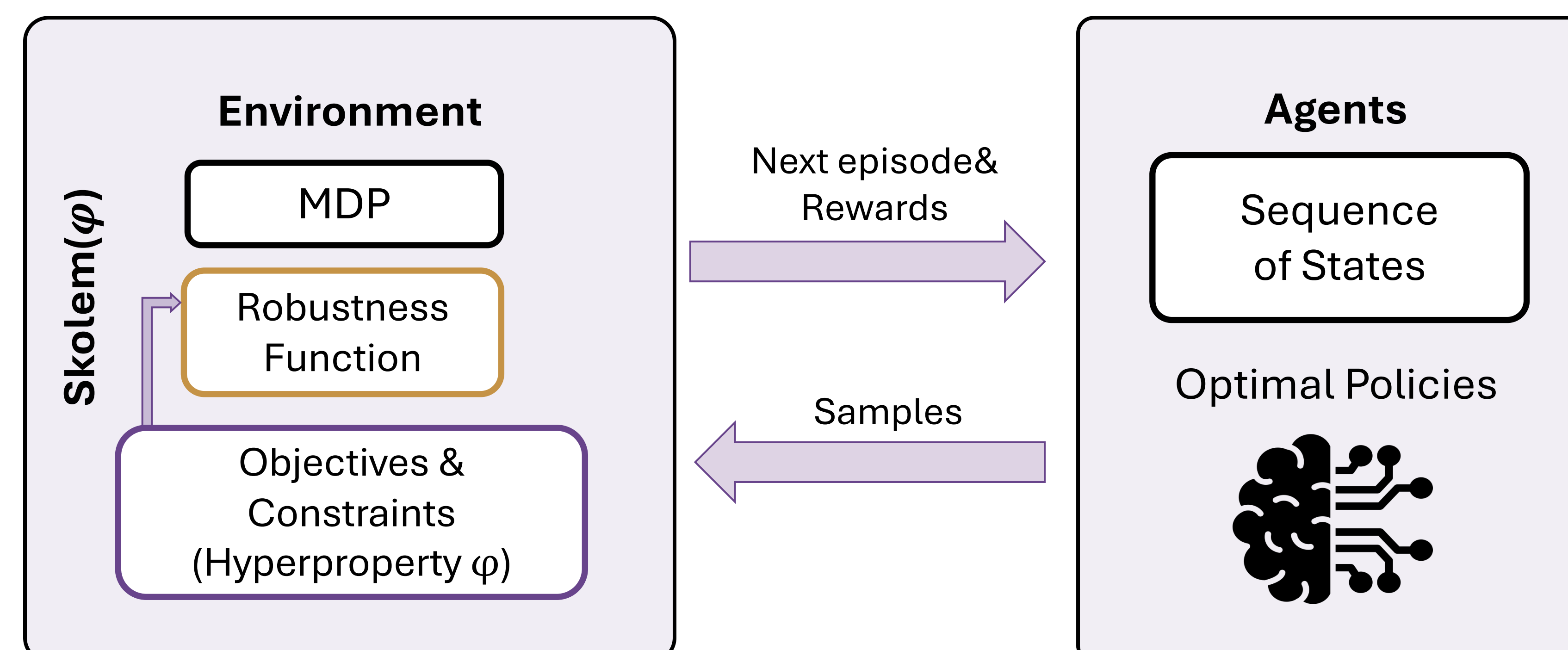
- We define **quantitative semantics** by min-max to **interpret temporal satisfaction**.

$$\begin{aligned} \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi) &= \rho_{\min} \text{ if } \text{Tr}(\zeta_{[\ell:k]}) = \epsilon \text{ and } \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi) \text{ otherwise.} \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \text{true}) &= \rho_{\max} \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), f(L(s_\ell) < c)) &= c - f(L(s_\ell)) \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \neg \psi) &= -\rho(\text{Tr}(\zeta_{[\ell:k]}), \psi) \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \bigcirc \psi) &= \rho(\text{Tr}(\zeta_{[\ell+1:k]}), \psi) \text{ if } (k > \ell). \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \bigvee \psi) &= \min_{i \in [\ell, k]} \rho(\text{Tr}(\zeta_{[i:k]}), \psi) \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \bigwedge \psi) &= \max_{i \in [\ell, k]} \rho(\text{Tr}(\zeta_{[i:k]}), \psi) \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_1 \wedge \psi_2) &= \min(\rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_1), \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_2)) \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_1 \vee \psi_2) &= \max(\rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_1), \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_2)) \\ \rho(\text{Tr}(\zeta_{[\ell:k]}), \psi_1 \mathcal{U} \psi_2) &= \max_{i \in [\ell, k]} \left(\min \left(\rho(\text{Tr}(\zeta_{[i:k]}), \psi_2), \min_{j \in [\ell, i]} \rho(\text{Tr}(\zeta_{[j:i]}), \psi_1) \right) \right) \end{aligned}$$

- We **inductively construct policies** to handle **dependencies** among multi-agent.

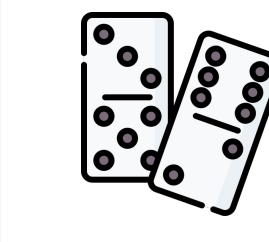
- For each $j \in \mathbb{Q}^\forall$, $\pi_j^*(\zeta_{j[0:k]}) \triangleq \mathcal{NN}_j^*(s_k)$
- For each $i \in \mathbb{Q}^\exists$, $\pi_i^*(\zeta_{i[0:k]}) \triangleq \mathcal{NN}_i^*(\mathbf{f}_i(\text{Tr}(\zeta_{i_1[0:k]}), \dots, \text{Tr}(\zeta_{i_{|\mathbb{Q}_i^\forall|}[0:k]})))$

HYPRL Architecture

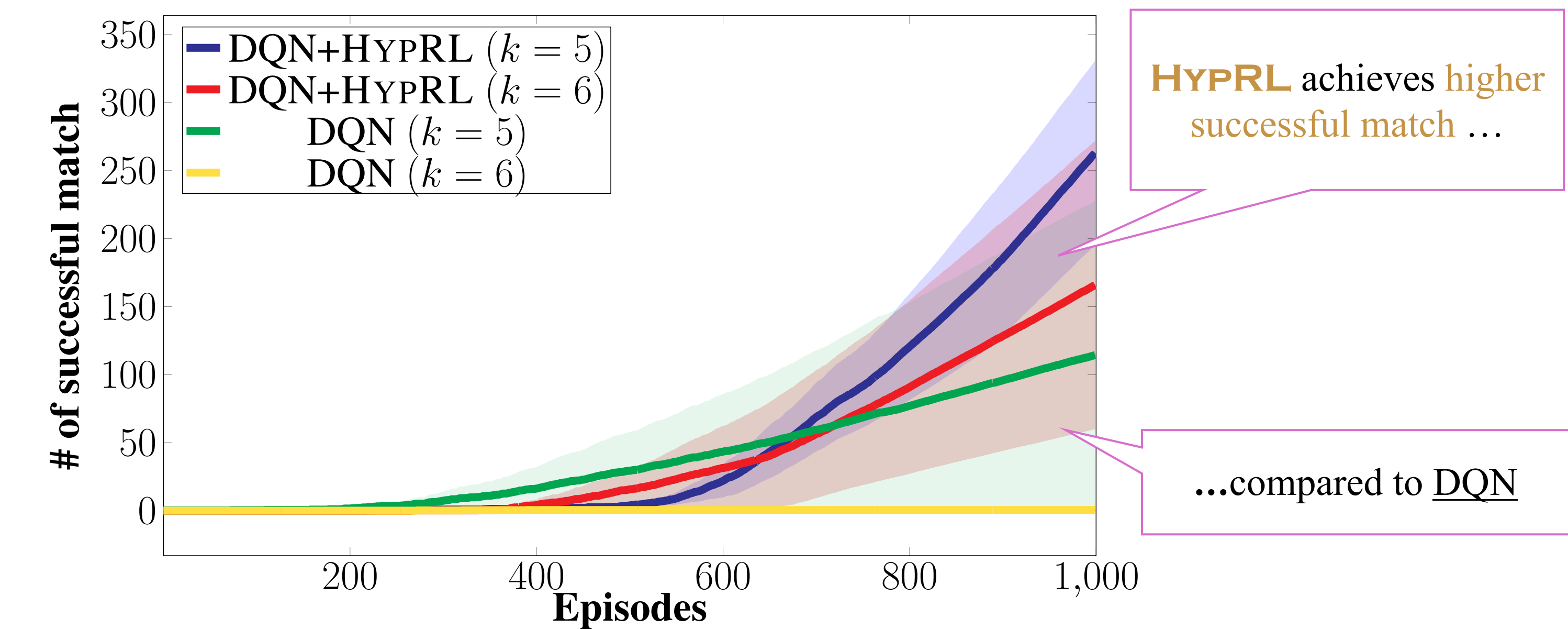


Case Studies & Results

We investigate **4 popular RL benchmarks** and compared with DQN, PPO, CQ-learning, and Shielding.



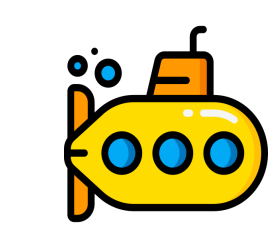
Post Correspondence Problem (PCP)



Wildfire Scenario

HYPRL achieves all objectives while obeying relational constraints in **fewer steps** compared to PPO.

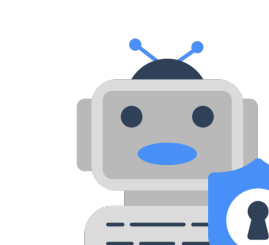
Size	Method	Dist	Steps O_1	Steps O_2
3^2	PPO	2.5 ± 0.01	33.43 ± 4.1	787.03 ± 31.8
	PPO + HYPRL	2.30 ± 0.03	18.940 ± 1.1	143.550 ± 1.1
5^2	PPO	4.2 ± 0.01	62.7 ± 9.7	S/B
	PPO + HYPRL	2.1 ± 0.05	59.50 ± 14.3	8057.5 ± 121.4
8^2	PPO	11.2 ± 0.03	16801.8 ± 2144.0	S/B
	PPO + HYPRL	6.94 ± 0.07	4149.6 ± 1743.1	386.2 ± 80.5
10^2	PPO	10.9 ± 0.01	S/B	29023.6 ± 976.4
	PPO + HYPRL	5.3 ± 0.10	21272.8 ± 3579.0	570.3 ± 52.2



Deep Sea Treasure

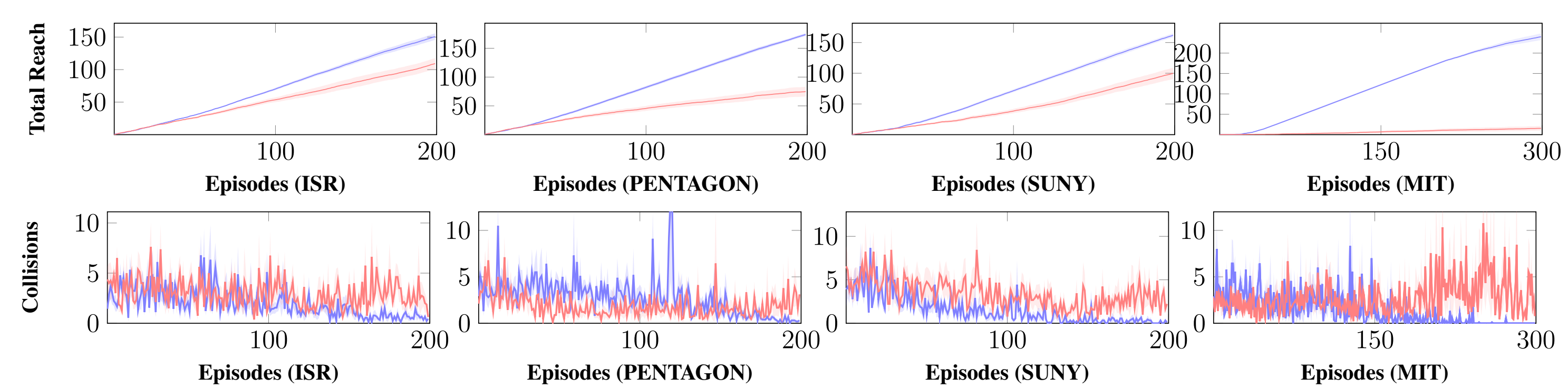
HYPRL finds **more treasures** in **less steps**, compared to PPO.

Method	Epi.	$\sum \text{Treasures}$	$\sum \text{Steps}$	Method	Epi.	$\sum \text{Treasures}$	$\sum \text{Steps}$
PPO	500	4.31 ± 1.2	0.17 ± 0.04	DQN	500	1.08 ± 0.2	0.05 ± 0.00
PPO + HYPRL	500	22.93 ± 2.2	0.91 ± 0.08	DQN + HYPRL	500	4.12 ± 1.4	0.14 ± 0.05
PPO	1000	3.22 ± 0.8	0.12 ± 0.03	DQN	1000	1.45 ± 0.2	0.02 ± 0.00
PPO + HYPRL	1000	23.97 ± 2.2	1.12 ± 0.08	DQN + HYPRL	1000	4.43 ± 0.8	0.21 ± 0.03



Safe RL

HYPRL achieves goal in less steps while maintaining low collision rate compared to CQ-learning and Shielding.



Maps	No. Agents	CQ		CQ + Shield <small>ElSayed-Aly et al. (2021)</small>		CQ + HYPRL	
		Steps	Collisions	Steps	Collisions	Steps	Collisions
ISR	2	27.95 ± 7.4	0.19 ± 0.1	17.40 ± 2.2	0.00 ± 0.0	7.58 ± 0.3	0.25 ± 0.2
Pentagon		36.46 ± 7.7	0.28 ± 0.1	75.20 ± 12.6	0.00 ± 0.0	11.90 ± 4.6	0.53 ± 0.5
SUNY		11.99 ± 0.5	0.01 ± 0.0	11.50 ± 0.3	0.00 ± 0.0	12.48 ± 0.6	0.00 ± 0.0
MIT		41.28 ± 8.5	0.20 ± 0.1	33.46 ± 3.4	0.00 ± 0.0	23.20 ± 0.5	0.00 ± 0.0
ISR	3	98.79 ± 0.8	12.68 ± 3.8	S/B	0.00 ± 0.0	74.18 ± 5.1	7.78 ± 1.0
Pentagon		97.15 ± 2.4	16.46 ± 7.2	S/B	0.00 ± 0.0	78.82 ± 1.7	10.92 ± 1.4
SUNY		84.89 ± 7.9	0.63 ± 0.2	82.35 ± 4.1	0.00 ± 0.0	44.95 ± 8.3	0.71 ± 0.4
MIT		96.96 ± 1.8	2.83 ± 1.3	S/B	0.00 ± 0.0	71.53 ± 7.7	1.58 ± 0.7

Take aways!

- Hyperproperties** naturally express both objectives and relational constraints in multi-agent RL.
- HYPRL** automatically converts any HyperLTL specification into a reward function.
- HYPRL** consistently outperforms baseline reward designs.