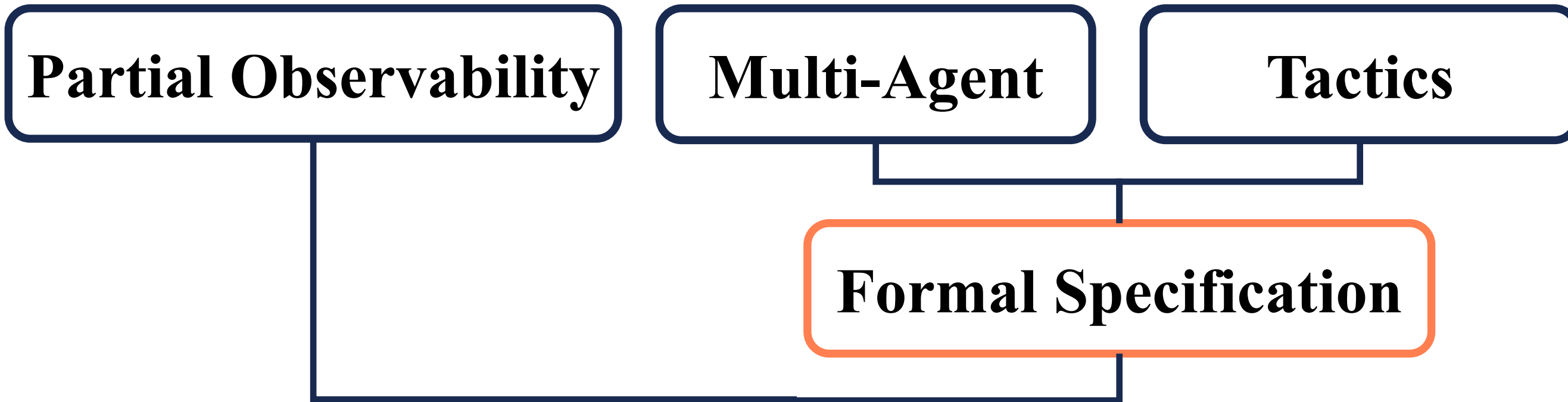


Research Question



How can **decentralized policies** be learned from formal specifications in partially observable multi-agent settings?

Formal Specification vs. Shaped Rewards

Mathematical rigor

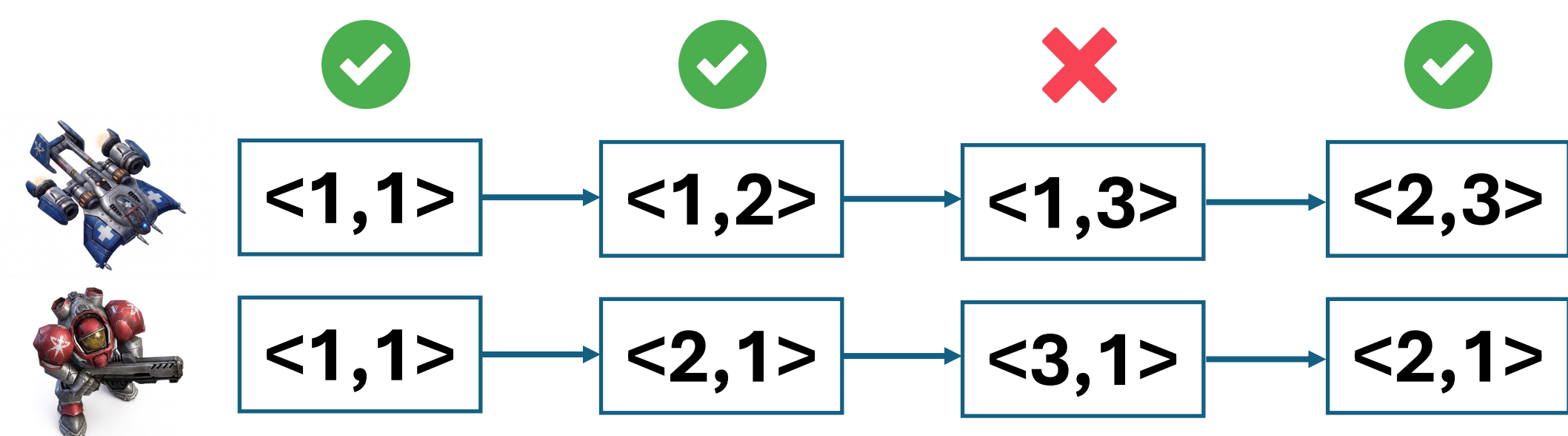
Expressiveness in specifying objectives and constraints

The ability to define **tactics** for achieving objectives

Specifications as Hyperproperties

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

"Is the distance between **Marine** and **Medivac** always less than 3?"



- We encode tactics as **HyperLTL** formulas, enabling specification-guided MARL under partial observability.

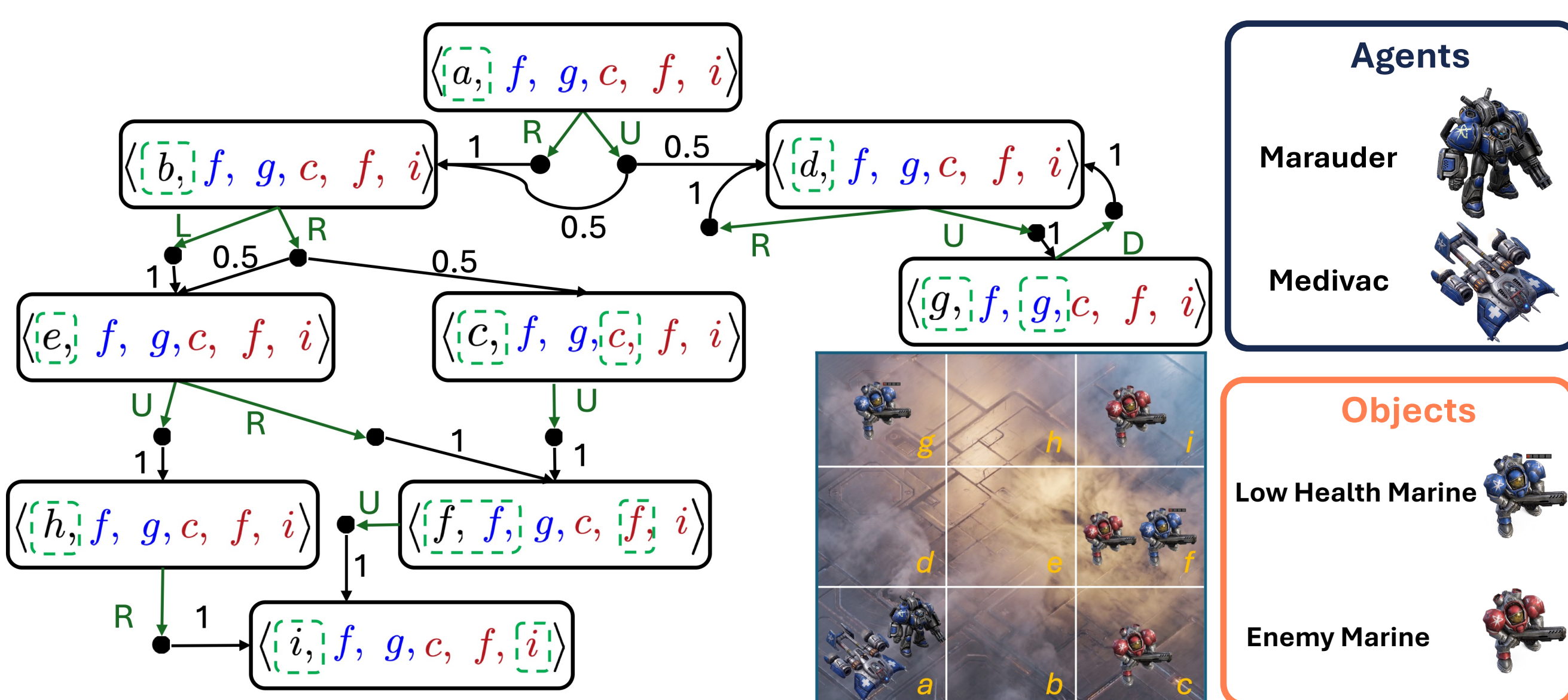
Problem Statement

Given a POMDP \mathcal{M} and a HyperLTL formula $\varphi = Q_1\tau_1 \dots Q_n\tau_n \cdot \psi$, our goal is to identify a tuple of n policies $\langle \pi_1^*, \dots, \pi_n^* \rangle$, such that:

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

where $\mathcal{D}_{\pi_1}, \dots, \mathcal{D}_{\pi_n}$ are the distributions over a set of histories drawn by policies π_1, \dots, π_n , and $\bigcup_{h \in (\mathcal{H}_i \sim \mathcal{D}_{\pi_i})} \{ \arg \max_{\zeta \in \mathcal{Z}^*} \Phi(h, \zeta) \}$ is a set of paths \mathcal{Z}_i associated with policy π_i .

Illustrative Example



HyperLTL Formula: $\forall \tau_1 \exists \tau_2. \text{Dist}(\text{pos}_{\tau_1}, \text{pos}_{\tau_2}) < 3 \ \mathcal{U} \ (\text{pos}_{\tau_1} = i)$

Samples \mathcal{Z}_{τ_1} : $h_1^1 = (a) \xrightarrow{U} (d) \xrightarrow{D} (g, g) \xrightarrow{D} (d) \xrightarrow{R} (d)$
 $h_2^2 = (a) \xrightarrow{U} (b) \xrightarrow{L} (e) \xrightarrow{R} (f, f, f) \xrightarrow{U} (i, i)$

Computing the satisfaction probability of HyperLTL specifications:
 $\text{Traces}(\{ \arg \max_{\zeta \in \mathcal{Z}^*} \Phi(h_1^1, \zeta) \}, \mathcal{Z}_{\tau_2}) \models \varphi_{\text{exp}} \wedge \text{Traces}(\{ \arg \max_{\zeta \in \mathcal{Z}^*} \Phi(h_2^2, \zeta) \}, \mathcal{Z}_{\tau_2}) \models \varphi_{\text{exp}} \Rightarrow h_2^2 \text{ serves as a witness for } \tau_2 \Rightarrow \mathbb{P} = 1$

Note that changing the formula to $\forall \psi$ reduces the satisfaction probability to 0.75, showing that the $\exists \psi$ formula expands the policy search space.

HYPOLE: Hyperproperty-Guided Multi-Agent Reinforcement Learning under Partial Observation

Arshia Rafieioskouei, Tzu-Han Hsu, Matthew Lucas, Borzoo Bonakdarpour



International Conference On Machine Learning



Under the **veil of uncertainty**, **HYPOLE** guides agents to learn **effective tactics**.

Our Solution

Step 1: We apply **Skolemization** to resolve **quantifier alternations** in a HyperLTL formula.

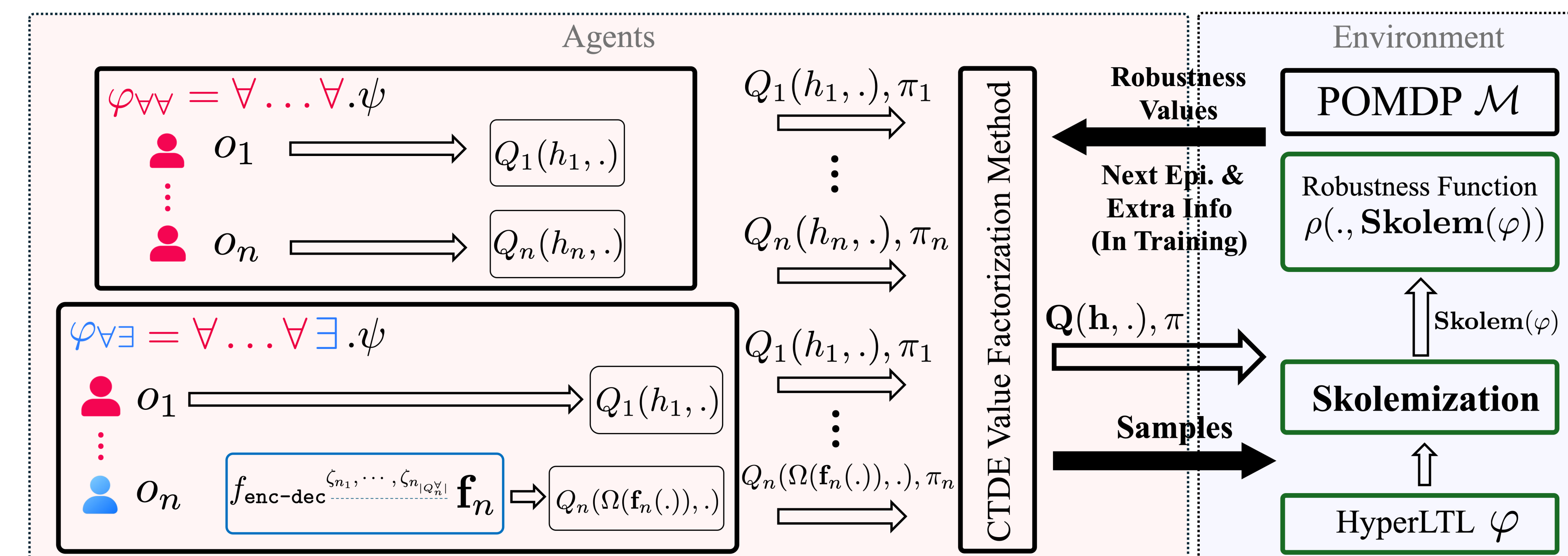
$$\text{Skolem}(\varphi) = \underbrace{\exists f_i(\tau_{i_1}, \dots, \tau_{i_{|Q_i^\exists|}})}_{\text{for each } i \in Q^\exists} \cdot \underbrace{\forall \tau_j}_{\text{for each } j \in Q^\forall} \cdot \text{Skolem}(\psi)$$

Step 2: We use **quantitative semantics** to construct robustness functions that transform skolemized HyperLTL formula into real-valued rewards (real-valued robustness values).

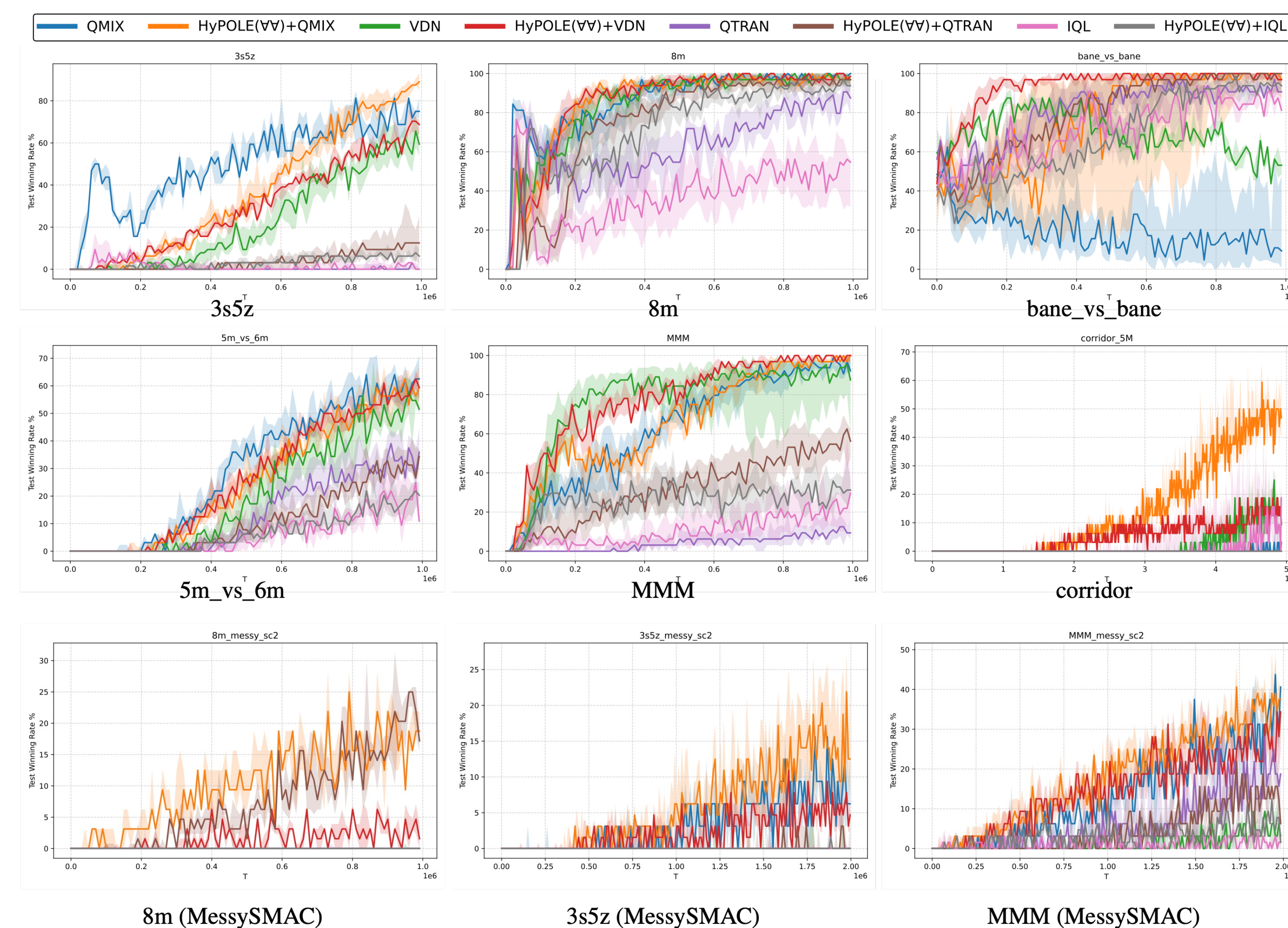
Step 3: We use **CTDE algorithms** (VDN/QMIX/QTRAN) to learn decentralized policies:

- We lift a shared POMDP environment to a Dec-POMDP setting compatible with CTDE methods.
- We learn optimal policies for agents with universal quantifiers. For each $j \in Q^\forall$, $\pi_j^*(h_{j[0:k]})$.
- Since existentially quantified agents require the traces of preceding universally quantified agents as input to their Skolem functions, we train an encoder-decoder model on replay-buffer samples during training and use it to construct policies for each $i \in Q^\exists$, $\pi_i^*(\Omega(\mathcal{F}_i(\text{Tr}(\zeta_{i_1[0:k]}), \dots, \text{Tr}(\zeta_{i_{|Q_i^\forall|}[0:k]}))))$.
- CTDE algorithms aim to approximate the joint value function $Q^*(\mathbf{h}, \mathbf{a}) = \max_{\pi} Q^\pi(\mathbf{h}, \mathbf{a})$. If the CTDE algorithm satisfies the IGM property and learns the optimal joint value function, then the induced decentralized policies $\langle \pi_i \rangle_{i \in Q^\exists}$ and $\langle \pi_j \rangle_{j \in Q^\forall}$ are optimal as well.

HYPOLE Overview

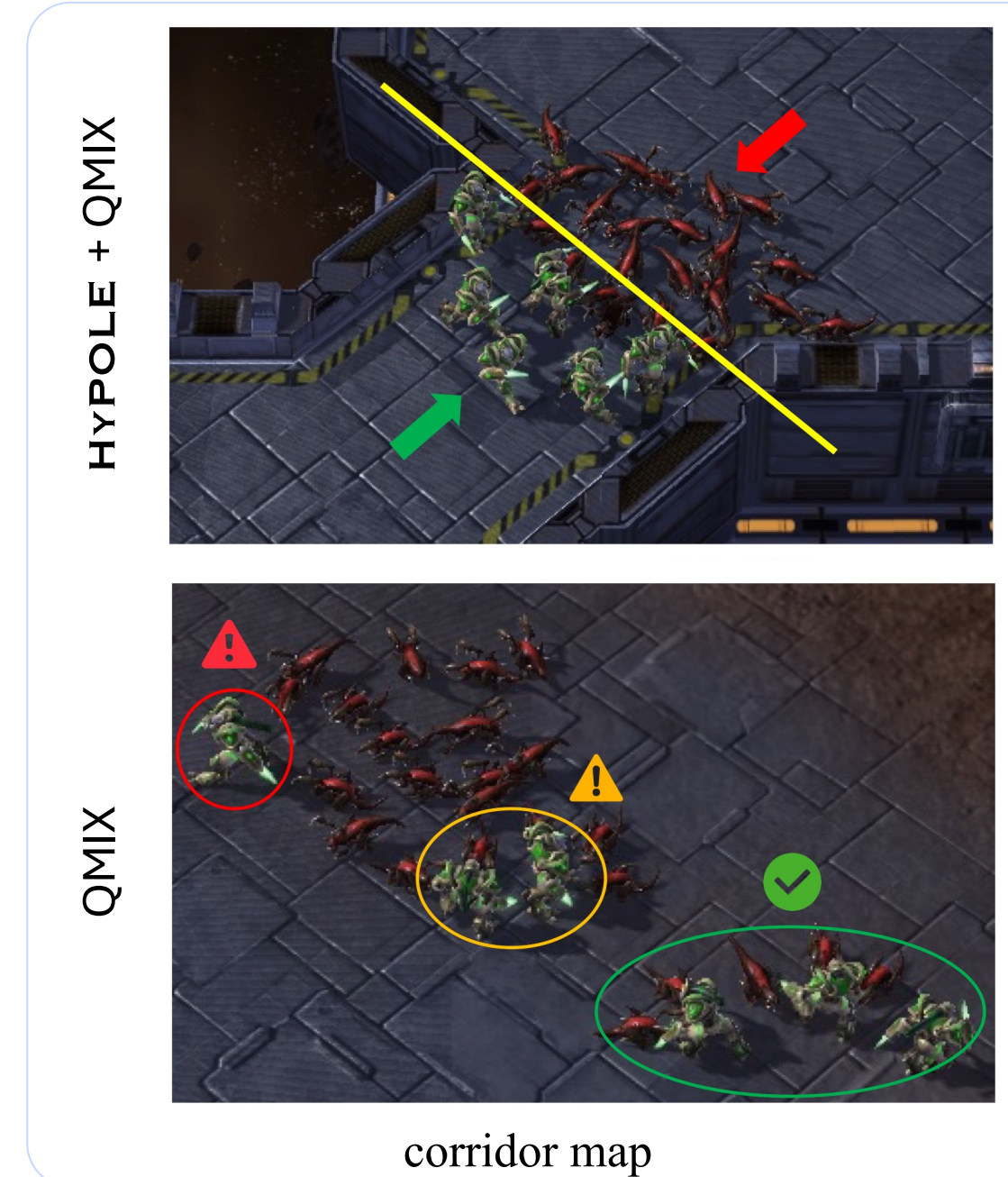


HYPOLE vs. Shaped Rewards

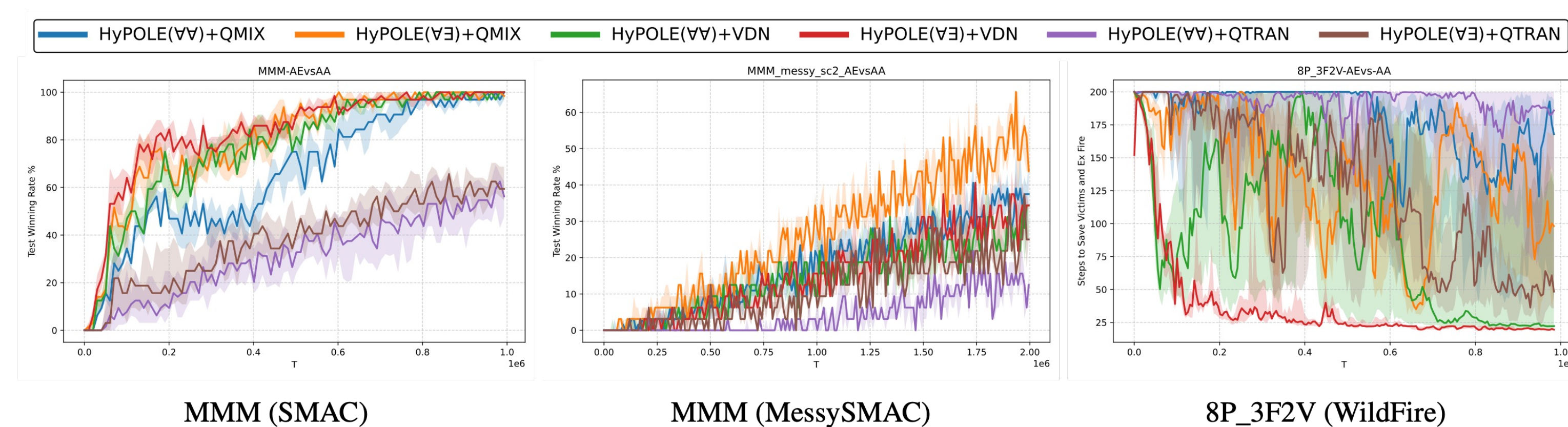


Take Away:

Our experiments show that encoding effective tactics as HyperLTL formulas improves MARL baselines over shaped rewards. **HYPOLE+QIQL** outperforms QTRAN on 8m, while **HYPOLE+QMIX** and **HYPOLE+VDN** achieve up to **90%** and **55%** gains in win rate on the challenging **bane_vs_bane** map. On **corridor**, **HYPOLE+QMIX** achieves up to a **60%** gain by encoding the defensive tactic of holding the choke point in HyperLTL. Even in the noisier and more stochastic **MessySMAC** environment, **HYPOLE** improves MARL baselines by up to **25%** on 8m, where the original baselines struggle to achieve nonzero win rates. Notably, **HYPOLE+VDN** achieves up to a **30%** gain.



$\forall \forall$ vs. $\forall \exists$ Specifications



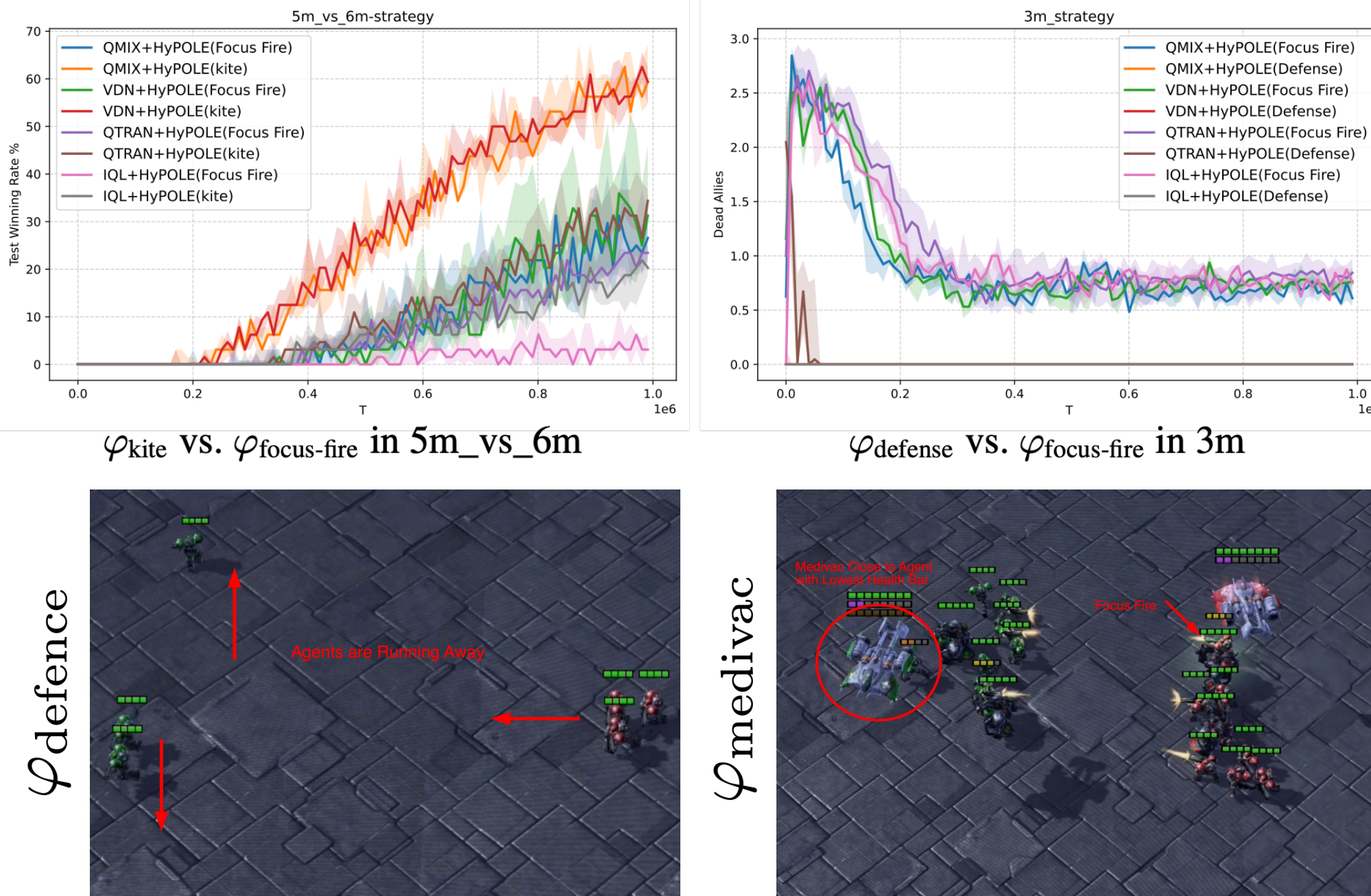
Take Away:

Changing the HyperLTL quantifier structure from $\forall \forall$ to $\forall \exists$ improves MARL when agents have strong interdependencies. In **MMM**, assigning an existential quantifier to the medivac agent yields up to **40%** win-rate gains for **HYPOLE($\forall \exists$)+QMIX** in SMAC and up to **20%** in MessySMAC compared to **HYPOLE($\forall \forall$)** with similar gains observed for QTRAN.

Tactics as HyperLTL Formulas

Take Away:

HYPOLE can encode different tactics through HyperLTL specifications. In **5m_vs_6m**, the kiting specification consistently outperforms the focus-fire specification, indicating that kiting is the more effective tactic for this scenario. In **3m**, a defensive specification encourages agents to retreat from enemies and reduce casualties.



Caution:

HYPOLE is sensitive to the design of the HyperLTL formula. When we remove the subformula that encourages agents to shoot enemies from the focus-fire specification, the win rate drops to zero, highlighting the importance of encoding effective tactics.