# MEDIUM: Certified Robust Learning for Multi-Agent Planning and Control

PI: Yiannis Kantaros, Dept. of Electrical and Systems Engineering, Washington University in St. Louis
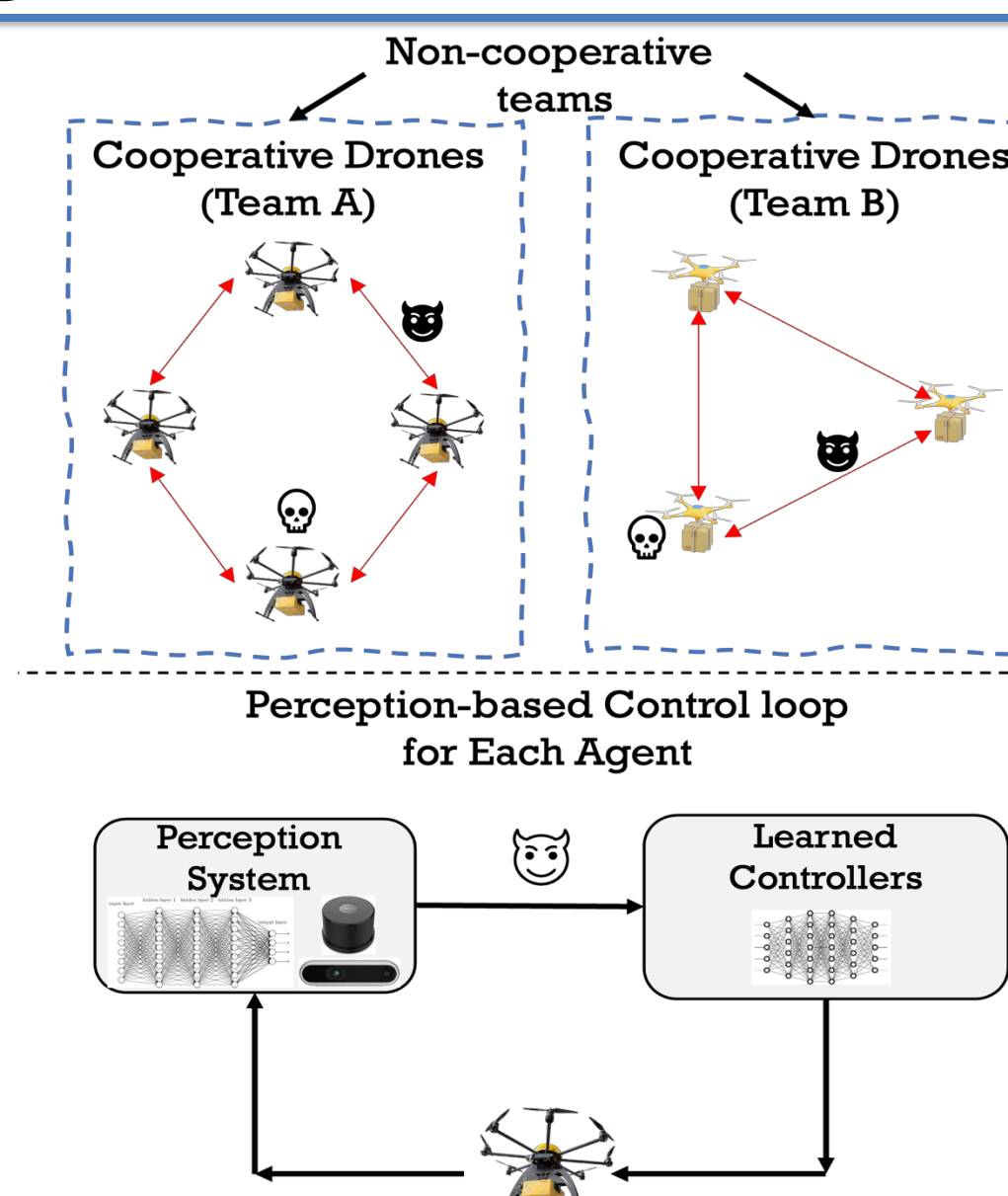
co-PI: Yevgeniy Vorobeychik, Dept. of Computer Science and Engineering , Washington University in St. Louis

co-PI: Hussein Sibai, Dept. of Computer Science and Engineering , Washington University in St. Louis

co-PI: Ning Zhang, Dept. of Computer Science and Engineering , Washington University in St. Louis
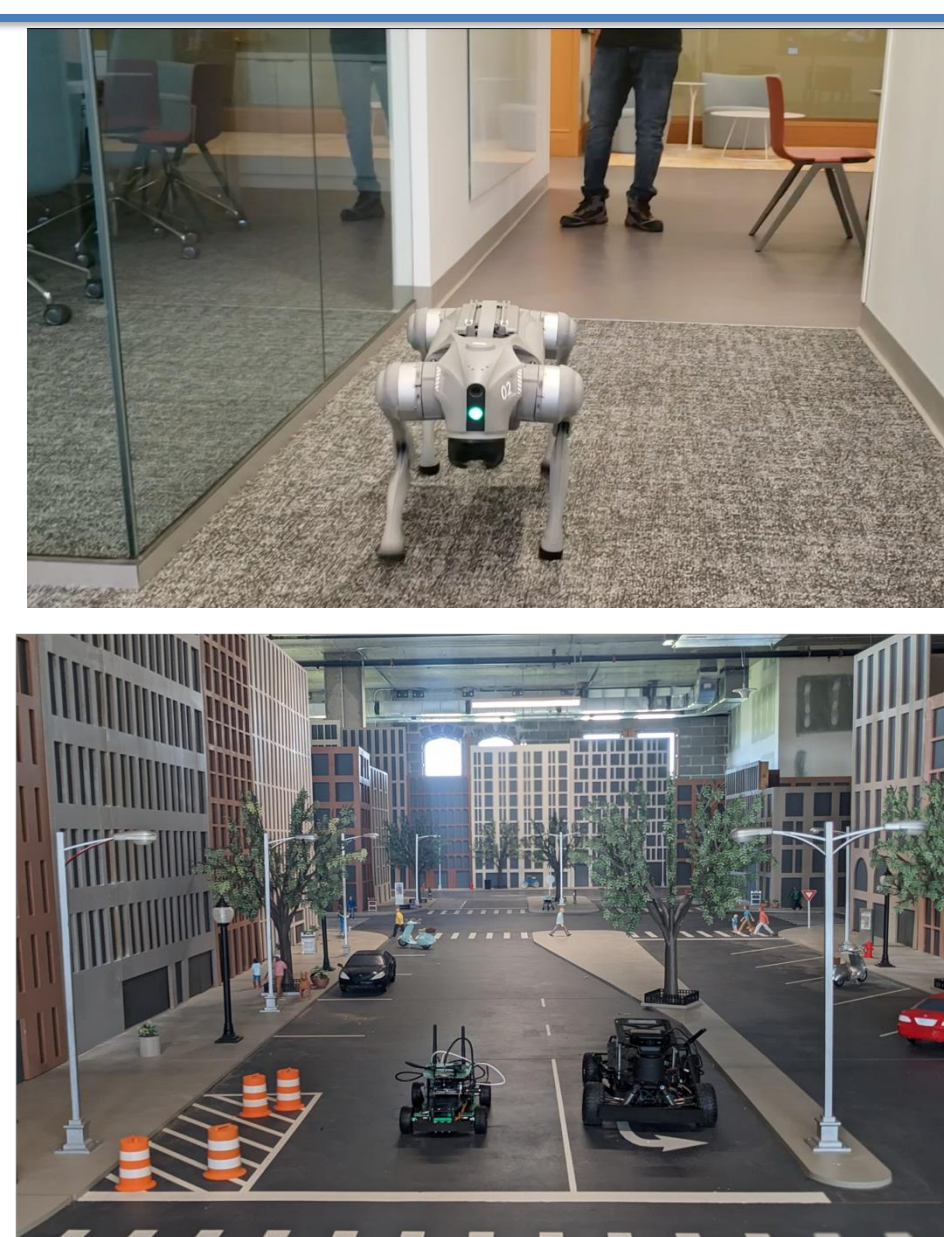
**Motivation:** Learning-based decision-making algorithms, such as Deep Reinforcement Learning or Neural Model Predictive Control, have been used to control multi-agent systems due to their generalization and real-time performance benefits. However, these algorithms lack robustness to imperceptible input perturbations. Despite impressive progress towards addressing this, existing methods lack rigorous safety and robustness guarantees. Also, lack of robustness becomes more pronounced in multi-agent settings, due to a larger surface of vulnerabilities, which has received significantly less research attention.

**Key Challenges:** Consider a multi-agent system and user-specified task and safety requirements $\phi$ (using e.g., formal languages or reward-based functions). How to design (and/or verify) learning-based controllers that are robust/safe (w.r.t. to satisfaction of $\phi$) in the presence of **(i)** perceptual noise; **(ii)** mis-calibrated confidence in predictions; **(iii)** adversarial communications; **(iv)** agent failures; **(v)** non-cooperative agents sharing the same workspace?



Non-cooperative teams — Cooperative Drones (Team A) / Cooperative Drones (Team B)

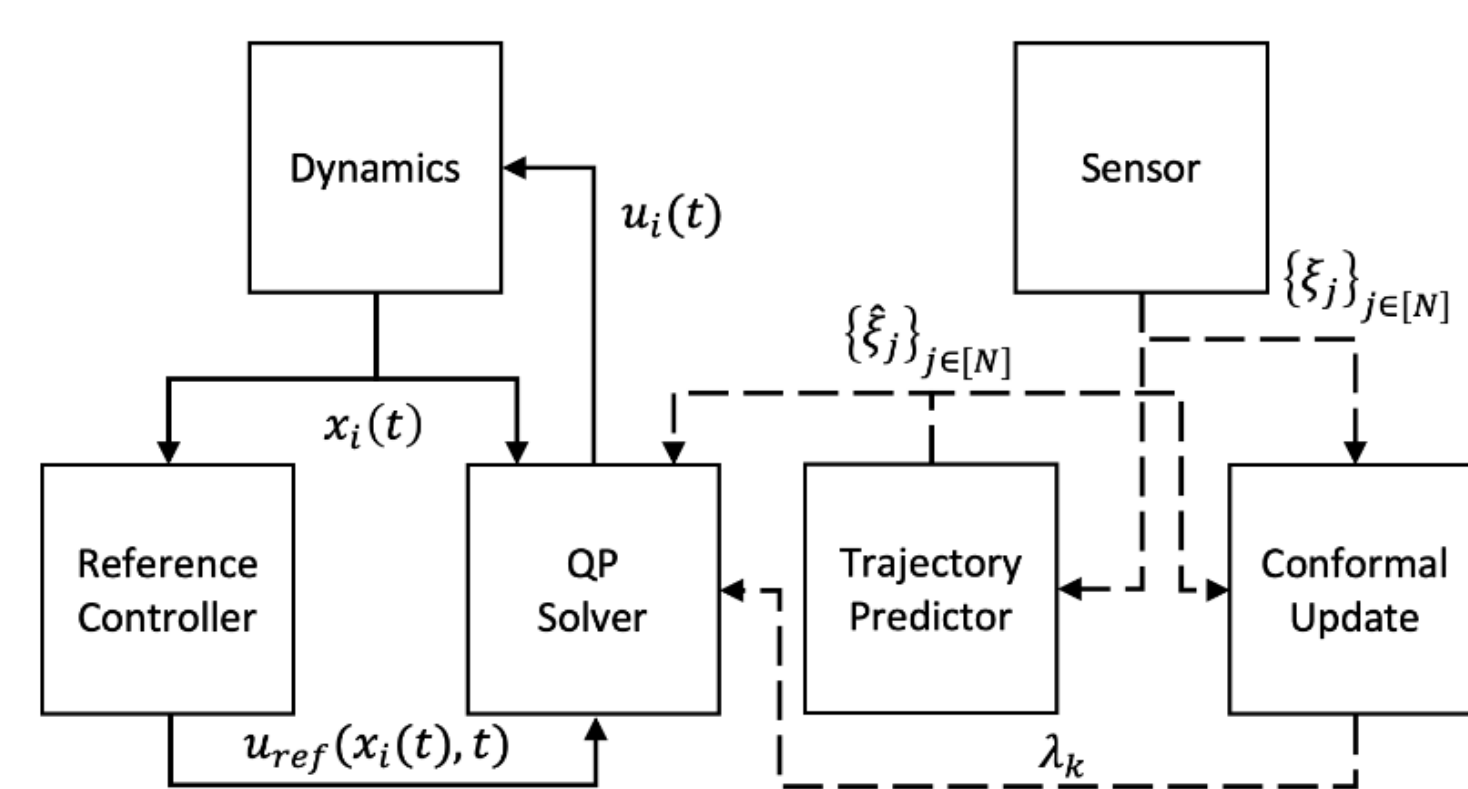Perception-based Control loop for Each Agent — Perception System / Learned Controllers

**Scientific Impact**

Provide the first theoretically-grounded algorithms to *certify*, *verify*, and *train* safe and robust multi-agent learning-enabled decision-making algorithms against (i)-(iv).
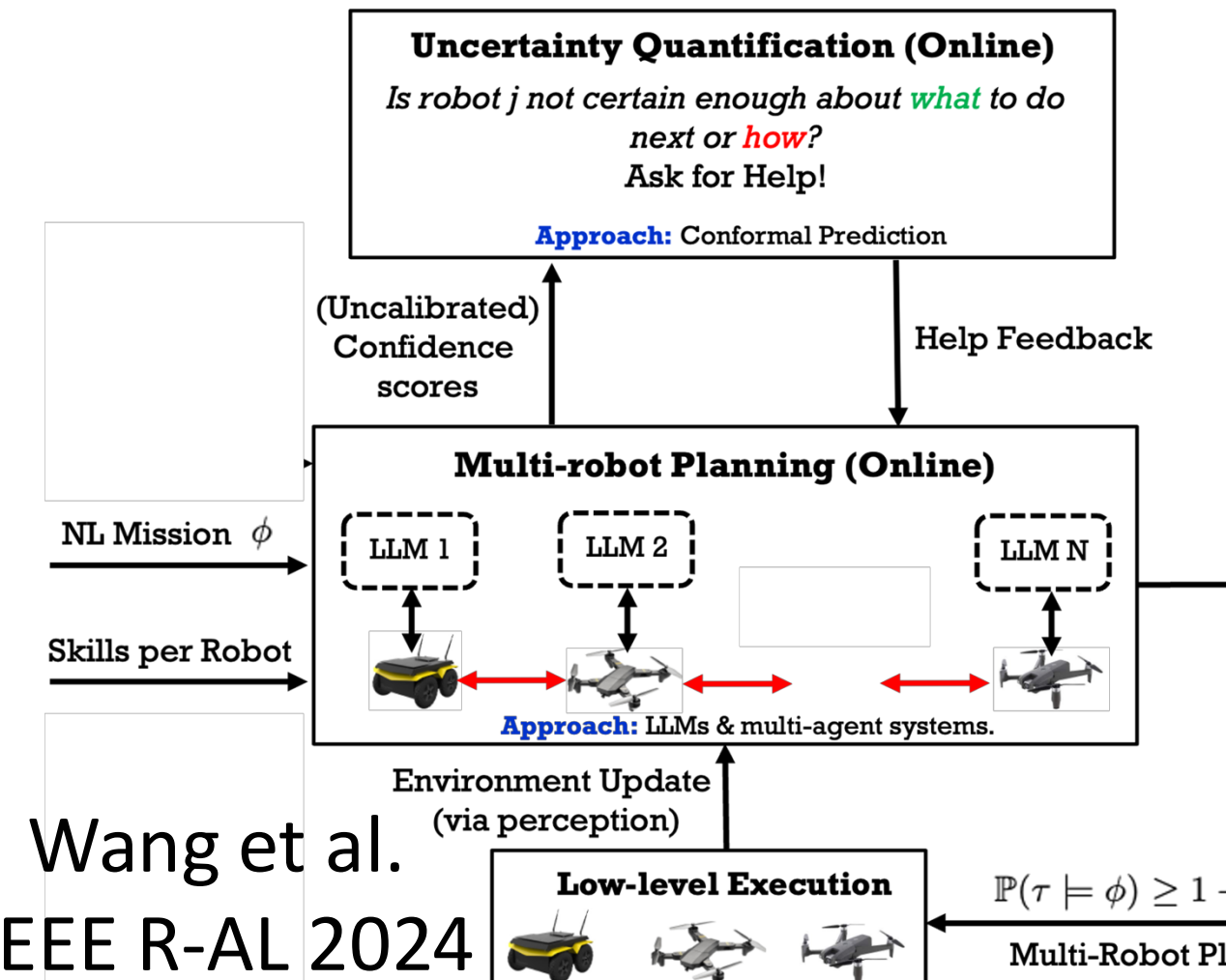


## New Contributions:

### Safe Decentralized ML-based Multi-agent Control



Huriot et al ICRA 2025

### Safe LLM-based Planning for Multi-Robot Systems

Uncertainty Quantification (Online)
*Is robot j not certain enough about what to do next or how?*
Ask for Help!
**Approach:** Conformal Prediction

(Uncalibrated) Confidence scores — Help Feedback

**Multi-robot Planning (Online)**
NL Mission $\phi$ / Skills per Robot
**Approach:** LLMs & multi-agent systems.
Environment Update (via perception)
Low-level Execution — $\mathbb{P}(\tau \models \phi) \geq 1 - \alpha$ — Multi-Robot Plan

Wang et al. IEEE R-AL 2024


D: Fire Hydrant 1 / B: Fire Hydrant 2, Traffic Cone 3 / C: Traffic Cone 1, 2 / A: EMPTY
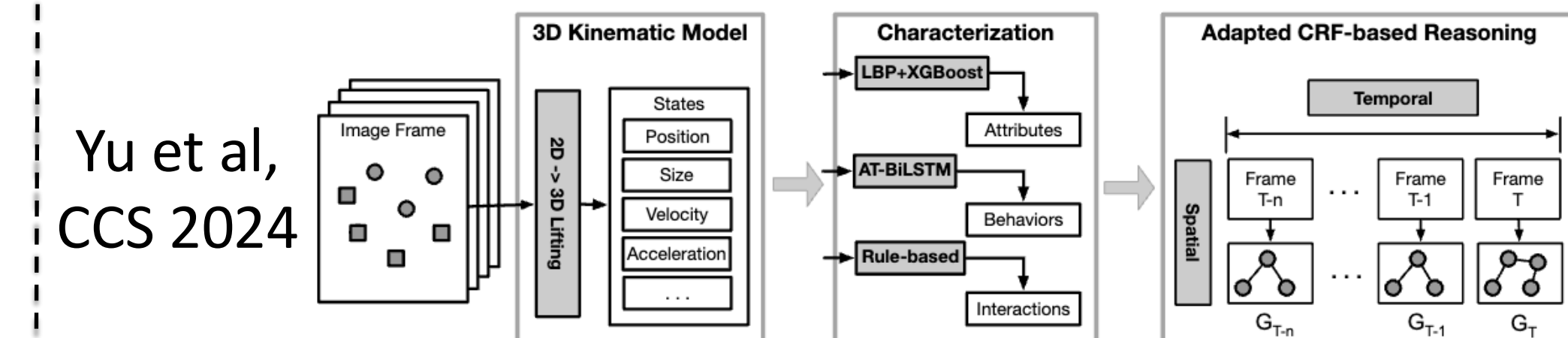
### Verified Safe RL-based Control

*Train controllers that can be efficiently verified for safety.*

$$\max_{\theta} \mathcal{J}(\pi_\theta)$$
$$\text{s.t. } \mathcal{J}_{C_i}(\pi_\theta) \leq d_i, \quad \forall i \in [m] \quad \text{(empirically satisfied)}$$
$$s_t \in \mathcal{S}_{\text{safe}}, \quad \forall t \in [K] \quad \text{(mathematically verified)}$$
$$s_{t+1} = F(s_t, a_t), a_t = \pi_\theta(s_t), s_0 \in \mathcal{S}_0 \subseteq \mathcal{S}_{\text{safe}}$$

| | Verified-50(↑) | | | Verified-50(↑) |
|---|---|---|---|---|
| PPO-Lag | 0.0 | | PPO-Lag | 72.8 |
| PPO-PID | 0.0 | | PPO-PID | 72.0 |
| CAP | 57.1 | Obstacle | CAP | 73.3 | Vehicle |
| MBPPO | 0.0 | Avoidance | MBPPO | 82.6 | Avoidance |
| CBF-RL | 0.0 | (quadrotor) | CBF-RL | 73.0 | |
| RESPO | 0.0 | | RESPO | 74.5 | |
| VSRL | 100.0 | | VSRL | 100.0 | |

Wu et al, NeurIPS 2024

### Defenses against Physical Adversarial Examples in Autonomous Systems

Yu et al, CCS 2024


3D Kinematic Model — States (Position, Size, Velocity, Acceleration) — Characterization (LBP+XGBoost, Attributes, Behaviors, Rule-based, Interactions) — Adapted CRF-based Reasoning (Temporal, Frame T-n, Frame T-1, Frame T)

### Reactive Multi-Robot Planning to Robot Skill Failures

Kalluraya et al (under review)



**Broader Impacts:**

- Enable more reliable and robust learning-enabled multi-agent systems that can safely perform complex tasks in uncertain, adversarial, and dynamic environments.

- Applications: delivery, transportation, manufacturing, search-and-rescue.
- Research opportunities to K12, UG, MS, PhD students and the WashU Robotics Club.

- Design new graduate courses (e.g., *Learning and Planning in Robotics, Trustworthy Autonomy*)
- Release open-source software and demonstrations