



Safe Probabilistic Planning in Human-Robot Interaction Using Conformal Risk Control

Jake Gonzales¹, Kazuki Mizuta², Karen Leung^{2*}, Lillian Ratliff^{1*}

¹Department of Electrical and Computer Engineering, ²Department of Aeronautics & Astronautics
University of Washington, * Equal advising

Project Website:



Motivation

Problem: Safe robot planning around humans, where the robot must reason about and adapt to uncertain, complex human behavior from a learned stochastic policy trained on real-world pedestrian data.

Challenges: In multi-agent systems, agents must adapt to complex non-stationary environments that arise from each agent updating their policies without making overly conservative distributional assumptions.

Applications: Human-robot interaction, autonomous vehicles, warehouse robotics, and general multi-agent interactions under uncertainty.

Contributions

- A framework that combines control-theoretic safety constraints with distribution-free uncertainty quantification for safe human-robot interactions.
- Theoretical analysis establishing the connection conformal risk control and control barrier functions for high-probability safety guarantees.
- An assumption-light algorithm for dynamic uncertainty adaptation of *risk* margins through a safety filter.

Control Barrier Functions (CBF)

Control affine dynamical system: $\dot{x} = f(x) + g(x)u$

Safe set: $\mathcal{S} := \{x \in \mathcal{X} | h(x) \geq 0\}$

Definition (Robust Control Barrier Functions)[1]

$h : \mathcal{X} \rightarrow \mathbb{R}$ is a robust CBF if $\exists K \in \text{class-}\mathcal{K}_\infty$ such that

$$\sup_{u \in \mathcal{U}} \{L_f h(x) + L_g h(x)u + K(h(x))\} \geq \eta, \quad \forall x \in \mathcal{X}$$

where $L_f h(x) = \nabla h(x)^\top f(x)$ and $L_g h(x) = \nabla h(x)^\top g(x)$
 η ensures safety over the entire interval $[t_k, t_{k+1}]$

Key: we can use tools from control theory to give us safety guarantees

Conformal Risk Control

- We use conformal risk control [2] to quantify uncertainty in the stochastic human behavior model through a *safety margin* parameter.
- Loss function: $\mathcal{L}_i(\lambda) = \max\{0, |\mathcal{B}_i(x_k, u_k) - \hat{\mathcal{B}}_i(x_k, u_k)| - \lambda\}$
- Conformal risk control:

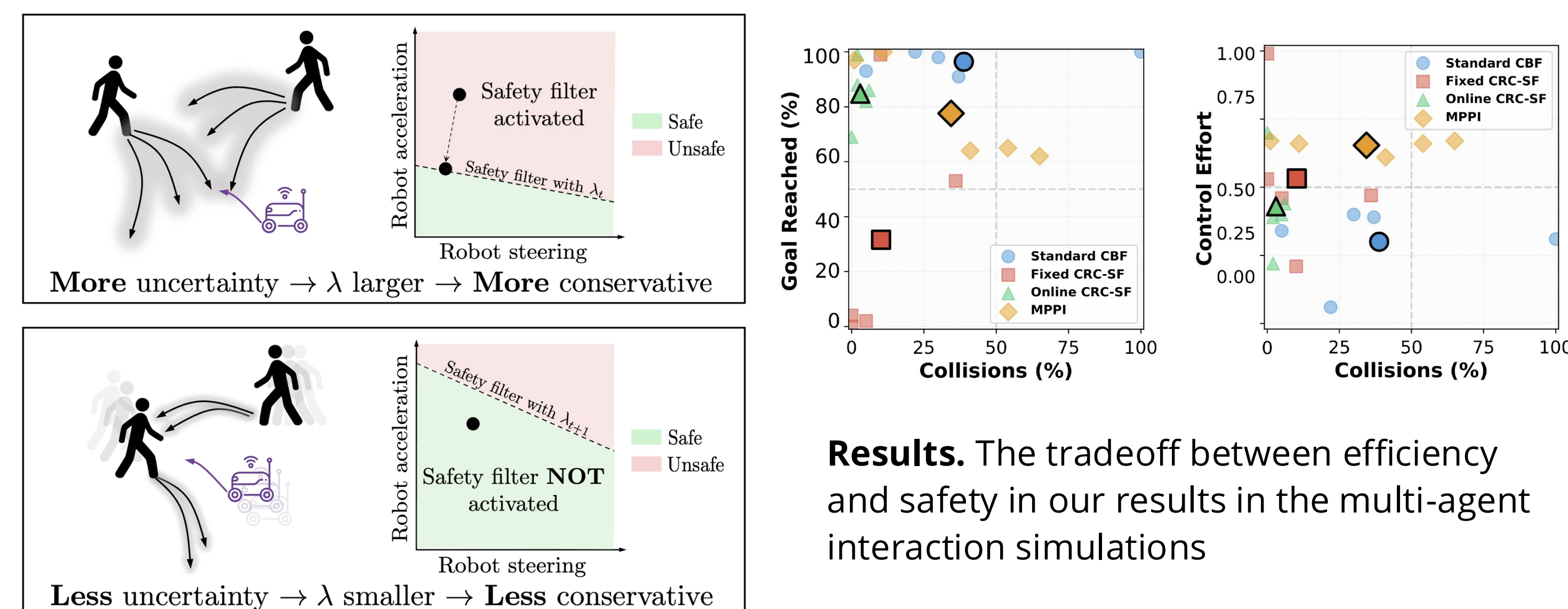
$$\hat{\lambda}_k = \inf \left\{ \lambda : \frac{1}{n_w + 1} \sum_{i=1}^{n_k} w_i \mathcal{L}_i(\lambda) + \frac{B}{n_w + 1} \leq \alpha \right\}, \quad n_w = \sum_{i=1}^{n_k} w_i$$

worst-case bound
empirical risk risk level geometrically decaying weights

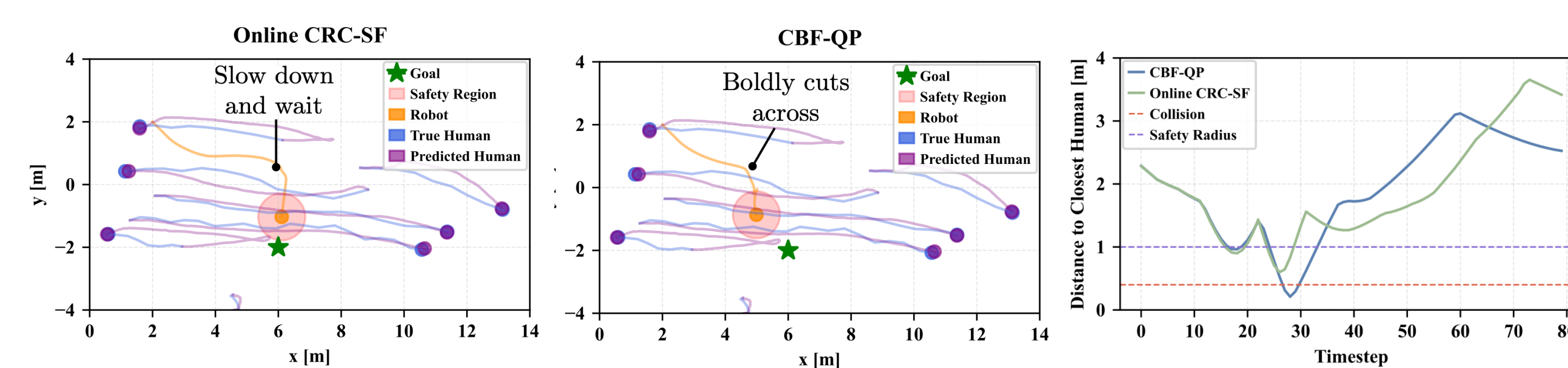
Gives the guarantee: $\mathbb{E}[\mathcal{L}(\lambda_k)] \leq \alpha + \beta$
total variation distance

Overview

Framework Overview. We develop a risk-aware adaptive safety filter that dynamically adjusts robot conservativeness based on the uncertainty of human behavior.



Problem Formulation



We formulate the human-robot interaction as follows: $x := (x_R, x_H) \quad u := (u_R, u_H)$

$$\begin{aligned} \text{(CT)} \quad \dot{x}(t) &= f(x(t)) + B_R u_R(t) + B_H u_H(t) \\ \text{(ZOH)} \quad u(t) &:= u_k, \forall t \in [t_k, t_{k+1}] \end{aligned} \quad \rightarrow \quad \begin{aligned} \text{(DT)} \quad x_{k+1} &= f(x_k) + B_R u_{R,k} + B_H u_{H,k} \\ &:= F_d(x_k, u_k) \end{aligned}$$

Offline Calibration:

- Collect interaction data using learned human behavior model and nominal robot policy
- Compute barrier certificates (ground truth and predicted)
- Calibrate safety margins using conformal risk control (CRC)
- Train uncertainty adaptation model

Online Deployment:

- Initialize start and end positions
- Sample human predictions $\hat{u}_H \sim P(u_H | x_{0:k})$
- Update safety margin based on context
- Solve the CRC safety filter
- Apply control to robot system
- Repeat at each time step

Theoretical Guarantees

Theorem 1 (CRC-CBF Safety Guarantee). *Our approach ensures the robot stays safe with high probability by dynamically adapting the safety margin λ based on prediction uncertainty. Specifically, the safe control set*

$$C_\lambda = \{u_R \mid \hat{\mathcal{B}}(\hat{x}, \hat{u}) \geq \hat{\lambda} + \epsilon\}$$

guarantees $\Pr(h(x_{k+1}) \geq 0) \geq 1 - \gamma$ at each timestep.

Online CRC Safety Filter

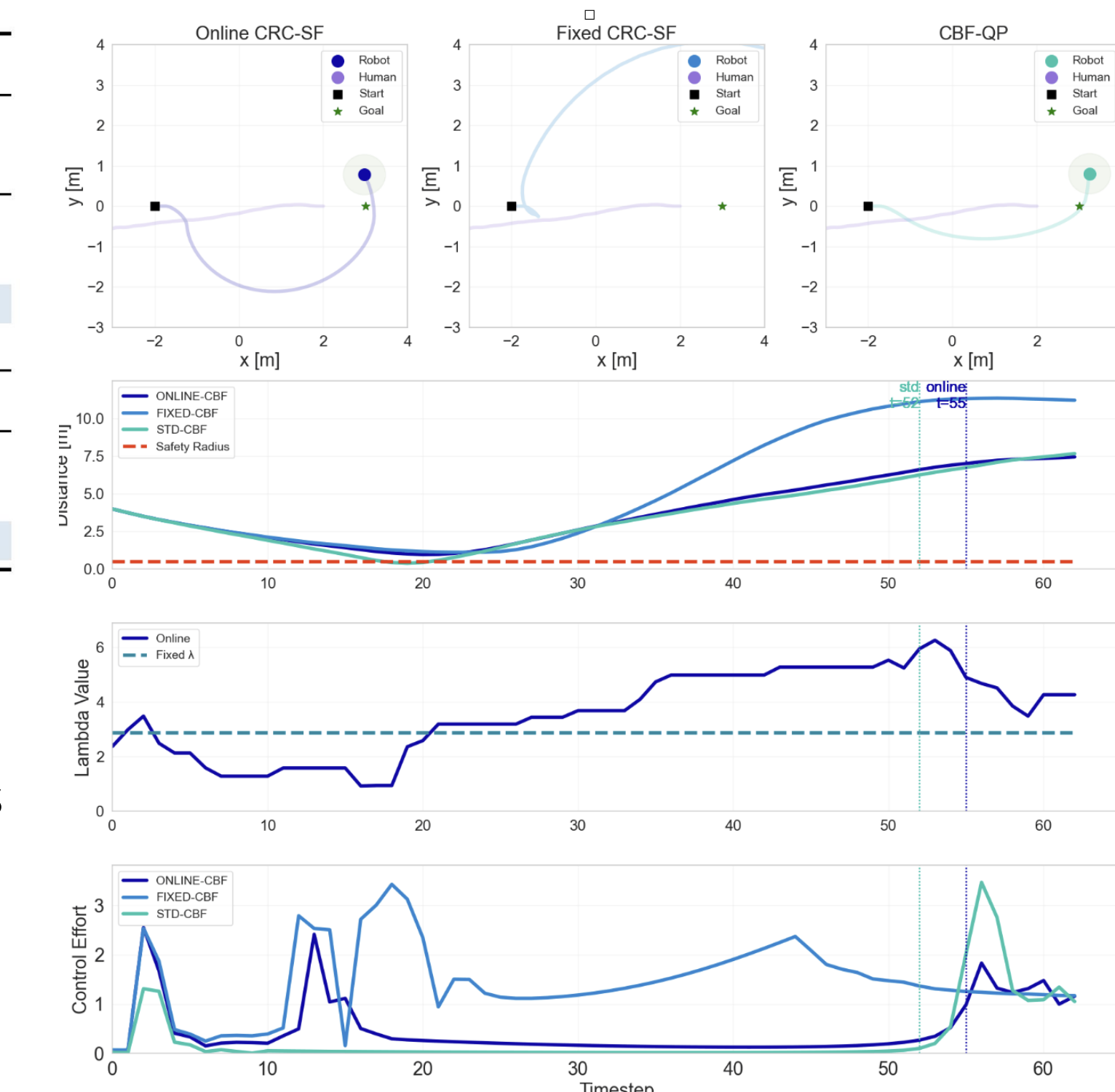
$$\begin{aligned} \min_{u_R} \quad & \|u_R - u_{\text{nom}}\|_2^2 \\ \text{such that} \quad & \Pr(u_R \in C_\lambda(\hat{x}, \hat{u}_H)) \geq 1 - \gamma \end{aligned}$$

Experimental Results

Multi-Agent Scenario					
Method	Coll. (%)↓	Safety Viol. (%)↓	Goal (%)↑	Control Effort↓	Control Smooth.↓
CBF-QP	38.8	81.4	96.4	0.30	0.12
Fixed CRC-SF	10.2	49.0	31.6	0.53	0.23
Online CRC-SF	3.0	53.2	84.8	0.43	0.19
MPPI	34.4	81.2	77.6	0.65	0.46

Single-Agent Scenario				
CBF-QP	16.0	57.0	100.0	0.31
Fixed CRC-SF	0.0	1.0	14.0	1.35
Online CRC-SF	2.0	15.0	78.0	0.59

Table: Performance comparison between multi-agent and single-agent human-robot navigation scenarios. The multi-agent results represent averages across five different test configurations, with each configuration run 100 times. Single-agent is 100 head-on simulations.



Limitations and Future Work

- **Limitations:** Our approach suffers from distribution shift in online deployment due to the dependence on offline data and since we learn a neural network to estimate the risk/safety margin we lose the full theoretical guarantee of safety.
- **Future Work:** Further testing on diverse scenarios and applications, quantifying error in the online estimation model for the safety margin, and more principled methods to both leverage offline data and update the parameter online.

[1] Dean et al. Guaranteeing Safety of Learned Perception Modules via Measurement Robust Control Barrier Functions, *Conference on Robot Learning (CoRL)*, 2020.

[2] Farinhas et al. Non-exchangeable Conformal Risk Control, *International Conference on Learning Representations (ICLR)*, 2024.

Acknowledgement: UW + Amazon Science Hub

Contact: jakegonz@uw.edu