

# Safe Interactive Autonomy for Multi-Agent Systems

Yiwei Lyu

Carnegie Mellon University  
Pittsburgh, PA 15213, USA  
yiweilyu@andrew.cmu.edu

## Abstract

It is envisioned that in the near future autonomous systems such as multi-agent systems, will co-exist with humans, e.g., autonomous vehicles will share roads with human drivers. These safety-critical scenarios require formally provable safety guarantees so that the robots will never collide with humans or with each other. It is challenging to provide such guarantees in the real world due to the stochastic environments and inaccurate models of heterogeneous agents including robots and humans. My PhD research investigates decision-making algorithm design for provably-correct safety guarantees in mixed multi-agent systems.

## Safe Control for Autonomous Driving with Non-Cooperative Vehicles under Uncertainty

The first component of my PhD research focuses on controller design for autonomous vehicles that can react to non-cooperative human-driven cars safely, especially under uncertainty. My work utilizes tools from probabilistic theory, control theory and optimization to explicitly reason about the impact of motion uncertainty among multi-vehicle systems on the control-based safe behavior design for an ego autonomous vehicle, and provide formal guarantees in terms of multi-vehicle collision avoidance and feasibility of our designed safe vehicle control framework. Building safety assurance for autonomous multi-agent systems in the real world is a challenging problem, as many existing works on multi-agent collision avoidance often (1) assume perfect information about robot position and motion for pre-computed safety guarantee, or (2) rely on overly conservative safe behavior to account for uncertainty due to imperfect robot information, which often makes the formulated safety-constrained optimization-based control problem infeasible even there exists one (Wang, Ames, and Egerstedt 2017).

To address the challenges above, in my work (Lyu, Luo, and Dolan 2021) I presented a novel bi-level control barrier function (CBF)(Ames et al. 2019)-based control framework with integrated chance-constraints to enable the ego autonomous vehicle to exhibit behaviors minimally deviated from the nominal control policy with varying conservativeness to accommodate different levels of uncertainties

between itself and other non-cooperative vehicles, while still enjoying a tighter probabilistic safety guarantee. This framework can be adapted to any existing multi-agent control policy to provide a formally provable safety guarantee with bounded probability. More importantly, the co-optimization between conservativeness and feasibility in my work provides a new way to characterize different degrees of safe behavior from aggressive agent to conservative agent, by explicitly specifying how fast an agent is willing to react when approaching to the boundary of the minimum safety distance. This finding generalizes our vehicle behavior design when interacting with drivers with different driving styles. For example, in a ramp merging scenario, the autonomous vehicle tends to be more conservative when yielding to an aggressive merging vehicle even at its normal speed. In my follow-up work (Van Koeveering et al. 2022), we further extend the proposed safe control method to high-relative degree systems that consider more realistic vehicle dynamics.

## Collision Avoidance for Heterogeneous Socially-Aware Multi-Agent Systems

Although in the work above we have derived safety-critical control for the autonomous agent that can interact with other non-cooperative agents in a collision-free manner under uncertainty, those agents are assumed to be moving obstacles that will not react to the ego agent. This may not be ideal as agents such as human drivers will also make decisions based on the ego agent's behavior. Hence how identifying human drivers' *reactive* behaviors becomes critical as this could affect the safe behavior design for the ego agent.

In the work (Lyu, Luo, and Dolan 2022a) as my first step towards this direction, I proposed to use machine learning techniques such as ridge linear regression to learn the behavior pattern of other agents from observations of their interactions with their surrounding agents. To ensure the method is generalizable, I focus on how to extend safe controllers' expressivity in describing different *cooperative* safe behaviors among heterogeneous agents. The heterogeneity captures the possibly various driving styles different drivers could have. With our extended safe controller, the ego vehicle is able to model other vehicles' underlying safe controllers with observed data. Given the learned information of other vehicles' underlying controllers, we presented an

interactive safe controller which provides greater flexibility for the ego vehicle to better coordinate with other heterogeneous robots with improved efficiency, while enjoying formally provable safety guarantees. We demonstrated how the ego vehicle adapts its own driving behavior in terms of the degree of conservative or aggressive based on the learned style of its opponent during a competitive task.

Yet sometimes, the other agents may not be fully cooperative and their behavior cannot be represented by a single controller. This motivates another line of my work (Lyu, Luo, and Dolan 2022b) that uses social preferences as an abstraction to describe agents' varying degrees of cooperation levels in safe multi-agent interaction. In this work, drawn upon a sociology concept of Social Value Orientation, a metric to quantify individual selfishness, I proposed a novel concept of Responsibility-associated Social Value Orientation to express the intended relative social implications between pairwise agents. This information is used to redefine each agent's preferences in terms of how much responsibility it should take in the collision avoidance scenarios against its neighboring robots. With this concept, we presented a novel decentralized control framework for multi-agent systems, where agents with varying cooperation levels can safely coordinate with formal safety guarantees. The proposed method is applicable to large-scaled systems, as each agent only needs to abide by its own constraint without predicting how the other agents are going to move.

### **Ongoing: Risk-Aware Safe Interactive Behavior Planning**

Different from (Lyu, Luo, and Dolan 2022b) that focuses on how to design safe controllers given the information of agent social preferences, I have been investigating a principle approach to *assign* such social identities or preferences to agents based on *risk* so that they can make improved *collective* decisions to *minimize* potential safety violation for the *entire* multi-agent systems. There have been some works in risk-aware safe control, however, all of them are focusing on using risk as an indicator to revise a single robot's controller facing static obstacles. In this work, I propose to combine a concept from financial portfolio investment, Conditional Value at Risk and model-based safe control to quantify the cumulative risk the agents face in a crowded dynamic environment. Instead of only knowing whether the system satisfies the safety verification or not, this unique combination would enhance our knowledge from the safe control perspective, about to which extent, the system is safe or unsafe. I demonstrated how risk within a multi-agent system can be measured by explicitly reasoning over the propagation of each agent's safe behavior and its impact on the rest of the group. With that, I also showed a novel way to allocate responsibilities among agents based on the estimated level of risk, so that the decentralized multi-agent team could safely and efficiently execute their tasks.

### **Future Plan: Towards Human-Like Safe Decision Making and Control**

My ultimate goal is to build safe controllers that enable autonomous agents to interact with humans in a socially-aware manner. To that end, my future work would put more focus on adding human factors into controller design, so that the generated behavior would not only be more expressive, and risk-aware, but more human-like. The main difference between humans and robots is in how they make decisions. From the perspective of the logic of decision-making, the current interactive design allows the robots under our control to execute actions after taking human's possible future actions into account. Yet real humans think more than that, about consequences when making decisions. Can we mimic that so that robots can also make consequence-aware decisions, by taking actions to influence human decision-making toward a mutual goal? One of my ongoing attempts is to build such consequence-aware controllers for autonomous vehicles so that they would actively probe human reactions while ensuring safety to maximize information gain about the human model's belief. Then autonomous vehicles would be capable to influence human decisions in challenging scenarios like unsigned interactions while keeping everyone safe.

Another reason that humans are special is, humans are not completely rational as machines. Research in cognitive science has investigated human modeling for a long time, and different human models with bounded rationality, like Bozeman Model are studied. How could we integrate these human models into our safe-critical systems is still a challenge that I would like to tackle. In the future, I also plan to do some human studies to verify the system design.

### **References**

- Ames, A. D.; Coogan, S.; Egerstedt, M.; Notomista, G.; Sreenath, K.; and Tabuada, P. 2019. Control barrier functions: Theory and applications. In *2019 18th European control conference (ECC)*, 3420–3431. IEEE.
- Lyu, Y.; Luo, W.; and Dolan, J. M. 2021. Probabilistic safety-assured adaptive merging control for autonomous vehicles. In *International Conference on Robotics and Automation (ICRA)*, 10764–10770. IEEE.
- Lyu, Y.; Luo, W.; and Dolan, J. M. 2022a. Adaptive Safe Merging Control for Heterogeneous Autonomous Vehicles using Parametric Control Barrier Functions. In *The 33rd Intelligent Vehicles Symposium (IV)*. IEEE.
- Lyu, Y.; Luo, W.; and Dolan, J. M. 2022b. Responsibility-associated Multi-agent Collision Avoidance with Social Preferences. In *The 25th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE.
- Van Koeveering, S.; Lyu, Y.; Luo, W.; and Dolan, J. 2022. Provable Probabilistic Safety and Feasibility-Assured Control for Autonomous Vehicles using Exponential Control Barrier Functions. In *Intelligent Vehicles Symposium (IV)*.
- Wang, L.; Ames, A. D.; and Egerstedt, M. 2017. Safety barrier certificates for collisions-free multirobot systems. *IEEE Transactions on Robotics*, 33(3): 661–674.