# Learning responsibility allocations for multi-agent interactions: A differentiable optimization approach with control barrier functions

Isaac Remy[1], David Fridovich-Keil[2], and Karen Leung[1,3]

*Abstract*— From autonomous driving to package delivery, ensuring safe yet efficient multi-agent interaction is challenging as the interaction dynamics are influenced by hard-to-model factors such as social norms and contextual cues. Understanding these influences can aid in the design and evaluation of socially-aware autonomous agents whose behaviors are aligned with human values. In this work, we seek to codify factors governing safe multi-agent interactions via the lens of *responsibility*, i.e., an agent's willingness to deviate from their desired control to accommodate safe interaction with others. Specifically, we propose a data-driven modeling approach based on control barrier functions and differentiable optimization that efficiently learns agents' responsibility allocation from data. We demonstrate on synthetic and real-world datasets that we can obtain an interpretable and quantitative understanding of how much agents adjust their behavior to ensure the safety of others given their current environment.

## I. INTRODUCTION

Safe navigation and collision avoidance with others come naturally to humans, but so much of our decision-making is governed by *social norms*, which depend on hard-to-model factors such as context and interaction history. For example, how do we precisely quantify factors that influence a driver's decision whether to slow down for an overtaking vehicle? On the one hand, end-to-end approaches are powerful in capturing complex and nuanced interaction dynamics but lack interpretability in explaining the interactions [1]–[3]. On the other hand, handcrafted model-based approaches, while fully interpretable, may miss capturing nuanced corner-cases and subtle interactions [4], [5]. Naturally, data-driven model-based methods present an attractive middle-ground for providing interpretability while capturing nuanced interactions.

In this work, we present a data-driven model-based approach to modeling social norms underlying multi-agent interactions. Specifically, we model social norms through the lens of *responsibility*, which we define as an agent's willingness to deviate from its desired control to accommodate safe interactions with others (visualized in Fig. 1), and present a computationally efficient technique to infer agents' responsibility allocations given information about their surroundings. Core to our approach is the combination of Control Barrier Functions (CBF) and differentiable optimization which provides interpretable inductive bias and the capability to extract interpretable quantities governing multi-
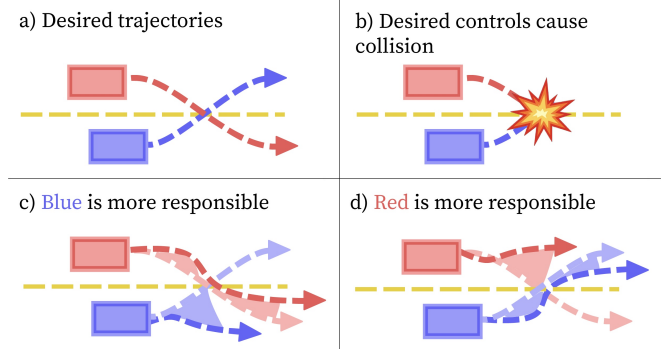


Fig. 1: In a) and b), two cars are swapping lanes on a highway, but their desired controls lead to collision. In c) and d), we see how the agents may deviate from their ideal trajectories, according to two different responsibility allocations. In this work we wish to solve the *inverse* problem: how do we infer these allocations from data?

agent interactions from data. Our responsibility learning framework is both applicable to guiding socially-aware robot policy construction and offline evaluation and data analysis.

**Contributions.** The contributions of this paper are four-fold: **(i)** We propose a novel mathematical formalization of describing responsibility allocations for multi-agent interactions, which is based upon Control Barrier Functions. **(ii)** We present a computationally efficient technique to learn agents' responsibility allocations from data. Computational efficiency is afforded by combining differentiable optimization techniques with modern deep learning and automatic differentiation tools. **(iii)** We introduce the concept of *symmetric responsibility* and provide a tractable approach to learning symmetric responsibility allocation models from data. We demonstrate the benefit of symmetric responsibility in improving data efficiency. **(iv)** We demonstrate the efficacy of our responsibility learning method on synthetic and human interaction data and highlight its capability to provide interpretable insights into the interaction.

## II. RELATED WORK

Recent works consider learning and instilling social awareness into autonomous planning algorithms to improve coordination between agents and, therefore, boost the overall safety and efficiency of the multi-agent system [6]–[9]. The idea is to include another agent's objective function into the autonomous agent's planning objective. The weighting between the two objective terms, referred to as the social value orientation (SVO), determines the output behavior, ranging from prosocial, egotistic, altruistic, competitive, and masochistic. The SVO weightings can be learned from data [6] to understand how much an agent considers the welfare

of others. Although these works do not frame their approach in terms of responsibility explicitly, the SVO weightings can potentially encode responsibility as they describe how much one agent considers the welfare of others. However, a few issues arise. First, the objectives can include other considerations beyond safety, making interpreting the SVO weightings challenging. Second, the SVO paradigm applies to only two agents, and it is unclear how it can extend to multiple agents. Third, it relies on knowledge of other agents' reward functions, which is often hard to obtain, especially for multi-agent settings.

There are a few recent works that explicitly consider the concept of responsibility in terms of multiple agents satisfying a shared collision avoidance constraint. Optimal reciprocal collision avoidance (ORCA) is a popular collision avoidance strategy, and the original concept assumes agents share equal responsibility in avoiding the velocity obstacle [10]. Various works have changed the responsibility allocation, either via heuristics or based on agents' states [11]–[13]. More recently, [14] proposes a variable responsibility optimal reciprocal collision avoidance (VR-ORCA) which computes a responsibility allocation that minimizes each agent's deviation from their preferred velocity. This choice of responsibility allocation results in what is more convenient for agents rather than based on what social norms would dictate. Although the concept of responsibility is explicitly modeled, the proposed formulation is not compatible with learning from data to mimic how humans naturally interact with one another. Lastly, similar to our proposed approach, [15] also frames responsibility allocations using CBFs but instead views responsibility as a parameter denoting how much each agent contributes to a joint CBF control constraint. In contrast, our responsibility parameter is a value describing how much agents deviate from their desired control—interesting future work could lie in more rigorous comparison of different responsibility formulations.

## III. PROBLEM FORMULATION

Consider a system with $N$ agents where $x \in \mathcal{X}$ denotes the state describing all agents, $u = (u_1, u_2, ..., u_N) \in \mathcal{U}_1 \times \mathcal{U}_2 \times ... \times \mathcal{U}_N$ denote the control inputs of each agent, and $\dot{x} = f(x, u)$ denotes the dynamics of the multi-agent system. Let $g(x, u) \geq 0$ denote a collision avoidance constraint shared amongst all agents. In real-world settings, humans take control actions to accomplish their desired task and avoid collision with one another. But *how much does each agent contribute to ensuring the collision avoidance constraint is satisfied?* In other words, *how responsible* is each agent in taking control actions that ensures that $g(x, u) \geq 0$ is satisfied. Suppose $\pi(\cdot)$ is a policy describing how humans behave around each other. *How does $\pi$ encode each agent's responsibility allocation in ensuring the collision avoidance constraint is satisfied?*

While a human's true behavior policy is extremely complex, we seek to develop an interpretable lens on $\pi$ to describe each agent's responsibility allocation. Let $\gamma$ denote the responsibility allocation for each agent. Then we seek to answer the following two questions: **(i)** *What is an interpretable model describing the dependency of $\pi$ on $\gamma$?* **(ii)** *Given such a model, how can $\gamma$ be tractably inferred from human interaction data?*

## IV. DEFINING RESPONSIBILITY ALLOCATIONS

In this section, we define responsibility as an agent's inclination to deviate from their desired control to satisfy collision avoidance constraints. At a high level, suppose agents have a desired control action they would like to take, e.g., maintain speed in their lane, or change lanes. Then, a *more (less) responsible* agent is willing to deviate more (less) from their desired control to respect collision avoidance constraints. We model agents' willingness to deviate from their desired control action as an optimization problem based on a safety filter derived from a control barrier function (CBF). Our use of CBFs specifically as a model for human safety is motivated by their simple but descriptive nature— CBFs let us describe human collision avoidance constraints concisely, intuitively, and with relatively minor assumptions. We first introduce CBFs and describe their use as a safety filter, and then formally define responsibility allocations.

### A. Control barrier function (CBFs)

Control Barrier Functions (CBFs) are mathematical tools used in control theory to ensure the safety of dynamic systems while achieving desired control objectives. Consider continuous-time dynamics, $\dot{x} = f(x, u)$ where $x \in \mathcal{X} \subset \mathbb{R}^n$, $u \in \mathcal{U} \subset \mathbb{R}^m$, and $f : \mathcal{X} \times \mathcal{U} \to \mathbb{R}^n$ is locally Lipschitz continuous. A CBF is defined as follows.

*Definition 1 (Control Barrier Function (CBF) [16]):* Given the aforementioned system dynamics, consider a continuously differentiable function $b : \mathcal{X} \to \mathbb{R}$ and a set $\mathcal{C} \subset \mathcal{X}$ where $\mathcal{C} = \{x \in \mathcal{X} \mid b(x) \geq 0\}$. Then $b$ is a valid CBF if there exists a class $\mathcal{K}_\infty$ function $\alpha$ such that,

$$\sup_{u \in \mathcal{U}} \nabla b(x)^T f(x, u) + \alpha(b(x)) \geq 0, \quad \forall x \in \mathcal{X}. \quad (1)$$

Intuitively, $b(\cdot)$ is a measure of safety, and a valid CBF ensures that there exists a feasible control preventing $b(x(t))$ from decreasing faster than a rate of $-\alpha(b(x))$ along system trajectories and that $\dot{b}(x(t)) = 0$ along the boundary of $\mathcal{C}$. It follows that if $b(\cdot)$ is a valid CBF for a set $\mathcal{C}$, then any locally Lipschitz controller $k : \mathcal{X} \to \mathcal{U}$ satisfying $\nabla b(x)^T f(x, u) + \alpha(b(x)) \geq 0$ renders the set $\mathcal{C}$ forward invariant [17].

### B. CBF safety filters

Given the forward invariance property of CBFs, we can compute a set of control inputs that prevent the system from exiting $\mathcal{C}$. The control set is simply the set of controls that satisfy the CBF inequality in (1), i.e.,

$$\mathcal{U}^{\text{safe}}(x; b, \alpha) = \{u \in \mathcal{U} \mid \nabla b(x)^T f(x, u) + \alpha(b(x)) \geq 0\}. \quad (2)$$

*Problem 1 (CBF safety filter):* A safety filter will take the desired control, $u^{\text{des}}$, computed from an upstream planning/control module, and project it into the safe control set $\mathcal{U}^{\text{safe}}(x; b, \alpha)$ defined in (2), i.e.

$$u^\star = \underset{u \in \mathcal{U}}{\arg\min} \|u - u^{\text{des}}\|_2^2 \text{ s.t. } \nabla b(x)^T f(x, u) + \alpha(b(x)) \geq 0.$$

*Remark 1:* If the dynamics are control affine, i.e., $\dot{x} = f(x) + g(x)u$ and $u \in \mathcal{U}$ are box constraints, i.e. $u_{\min} \leq$

$u \leq u_{\max}$, then Prob. 1 is a convex quadratic program (QP) and can be solved efficiently [18], [19].

*Remark 2:* For a valid CBF, Prob. 1 is always feasible. However, finding a valid CBF is often nontrivial. Instead, a desired CBF can be chosen and slack variables can be introduced (demonstrated in Prob. 2) to encourage agents to satisfy the CBF constraint as best as possible within the control limits. In the context of this work, the loss of guarantees is not an issue as we are analyzing how much deviation each agent makes from their desired control rather than directly synthesizing guaranteed safe control strategies.

*Remark 3:* Definition 1 assumes we have a relative degree 1 system. To account for higher relative degree systems (e.g., double integrator systems, with $b(x)$ a function only of position variables), we can employ high order CBFs [20] which results in a slightly more complicated CBF inequality that remains linear in control for control affine systems.

### C. Responsibility allocations via a CBF safety filter lens

In this work, we consider using CBF safety filters as models to describe *multi-agent* interactions and use them for system analysis and inference to gain insight into how agents cooperate to maintain safety. The insights gained into multi-agent interactions can be used for offline (crash) analysis, measuring social acceptability of a policy, or to adapt an agent's behavior online.

Consider a *multi-agent* control affine dynamical system where $\mathbf{x} = (x_1, ..., x_N)$ is the joint state describing $N$ agents, and each agent has control input $u_i$.

$$\dot{\mathbf{x}} = \tilde{f}(\mathbf{x}) + \sum_{i=1}^{N} g_i(\mathbf{x})u_i \qquad (3)$$

Given a CBF $b(\cdot)$ for the multiagent system, Prob. 2 describes the direct application of the CBF safety filter problem (Prob. 1) onto the joint dynamics.

*Problem 2 (Unweighted multiagent CBF safety filter):*

$$u_{1:N}^{\text{exec}} = \operatorname*{argmin}_{u_1,...,u_N} \quad \sum_{i=1}^{N} \|u_i - u_i^{\text{des}}\|_2^2$$

$$\text{s.t.} \quad \nabla b(\mathbf{x})^T \left[ \tilde{f}(\mathbf{x}) + \sum_{i=1}^{N} g_i(\mathbf{x})u_i \right] + \alpha(b(\mathbf{x})) \geq 0$$

$$u_1 \in \mathcal{U}_1, ..., u_N \in \mathcal{U}_N.$$

Notice that Prob. 2 has *equal* weighting in penalizing each agent for deviating from their desired control input, implying no agent is more or less inclined to ensure the safety constraint is satisfied. We modify Prob. 2 by adding a coefficient $\gamma_i$ for each agent to influence the relative penalty in deviating away from the desired control. As such, we propose the following definition of responsibility allocation.

*Definition 2 (Responsibility allocation):* Let there be $N$ agents interacting, and their joint dynamics are governed by (3). Let $b(\cdot)$ be a CBF for the joint system describing the collision set. Each agent $i \in \{1, 2, \ldots, N\}$ has a desired control input $u_i^{\text{des}}$ that is projected into the safe control set via a projection mapping proj. The projection mapping proj depends on the corresponding state $\mathbf{x}$, CBF $b$, class $\mathcal{K}_\infty$ function $\alpha$, and the responsibility allocation vector

$\boldsymbol{\gamma} = [\gamma_1, \gamma_2, ..., \gamma_N]$ with $\mathbf{1}^T\boldsymbol{\gamma} = 1$ and $0 \leq \gamma_i \leq 1$, which determines how much deviation from the desired control each agent is willing to make to satisfy (2). The projection mapping proj can be computed by solving the optimization problem described in Prob. 3 (differences from Prob. 2 are highlighted in blue).

*Problem 3 (Responsible multiagent CBF safety filter):*

$$\text{proj}(u_{1:N}^{\text{des}}; \mathbf{x}, b, \alpha, \boldsymbol{\gamma}) \coloneqq$$

$$\operatorname*{argmin}_{u_1,...,u_N,\epsilon} \quad \sum_{i=1}^{N} \left( \gamma_i \|u_i - u_i^{\text{des}}\|_2^2 + \beta_1\|u_i\|_2^2 \right) + \beta_2\epsilon^2$$

$$\text{s.t.} \quad \nabla b(\mathbf{x})^T \left[ \tilde{f}(\mathbf{x}) + \sum_{i=1}^{N} g_i(\mathbf{x})u_i \right] + \alpha(b(\mathbf{x})) \geq -\epsilon$$

$$u_1 \in \mathcal{U}_1, ..., u_N \in \mathcal{U}_N$$

$$\epsilon > 0.$$

*Remark 4:* The slack variable $\epsilon > 0$ reflects the discussion in Remark 2. The regularization term $\|u_i\|_2^2$ ensures a unique solution for cases in which Agent $i$ bears total responsibility for constraint satisfaction, i.e. $\gamma_i = 0$.

*Example 1 (Two-agent 1D single integrator):* Consider a two-agent 1D single integrator system with CBF,

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}, \; b(x_1, x_2) = (x_1 - x_2)^2 - 1, \; \alpha(x) = x.$$

Let $\gamma$ denote the responsibility allocation of Agent 1. Fig. 2 shows the solution to Prob. 3 for various values of $\gamma$. We see that when $\gamma = 0$, Agent 1 takes full responsibility by deviating from its desired control value to satisfy the CBF constraint (green region) while Agent 2 has no deivation. As $\gamma$ increases, Agent 2 begins to deviate more and Agent 1 less, until at $\gamma = 1$, Agent 2 is fully responsible.
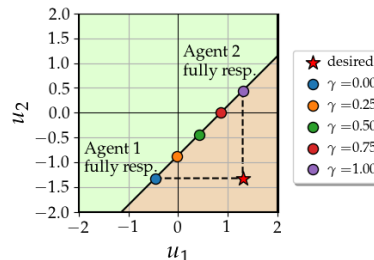


Fig. 2: Solutions to the two-agent CBF filter problem with different responsibility allocation values for Agent 1. The green and orange regions are the sets of feasible and infeasible controls, respectively.

### V. INFERRING RESPONSIBILITY FROM DATA

The previous section described how a choice of responsibility allocation $\boldsymbol{\gamma}$ influences each agent's deviation from their desired control input. But *how do we choose responsibility allocation values?* Human expertise can qualitatively describe trends (e.g., the agent behind should yield to the agent in front), but data quantitatively informs what the responsibility allocation values should be. To this end, we propose a data-driven inverse optimization approach to learn responsibility allocations from interaction data.

### A. Bi-level optimization formulation

The goal is to find a value for $\boldsymbol{\gamma}$ such that the solution to Prob. 3 matches an interaction dataset as closely as possible. Suppose we have a dataset describing $N$-agent interactions

**Algorithm 1** Single gradient step to solve responsibility allocation inference

**Require:** $\widetilde{\gamma}$ (unnormalized responsibility allocation)
**Require:** $\delta$ (step size)
**Require:** $\mathcal{D} = \{(\mathbf{x}^{(k)}, u_{1:N}^{(k)})\}_{k=1,...,K}$ (dataset)
1: $\gamma = \text{softmax}(\widetilde{\gamma})$ (ensuring $\mathbf{1}^T\gamma = 1$ and $\gamma \in [0,1]$)
2: $\widetilde{U}_{1:N} = \text{batch}(\widetilde{u}_{1:N}(\gamma, \mathbf{x}^{\text{data}}); \mathcal{D})$
3: $\ell(\gamma) = \sum_{(u_{1:N}^{\text{data}}, \widetilde{u}_{1:N}) \in \mathcal{D} \times \widetilde{U}_{1:N}} \Delta(u_{1:N}^{\text{data}}, \widetilde{u}_{1:N})$
4: $\widetilde{\gamma} \leftarrow \widetilde{\gamma} - \delta \frac{\partial \ell}{\partial \gamma} \frac{\partial \gamma}{\partial \widetilde{\gamma}}$

---

$\mathcal{D} = \{(\mathbf{x}^{(k)}, u_{1:N}^{(k)})\}_{k=1,...,K}$, and a *desired* policy $u_{1:N}^{\text{des}} = \pi_{1:N}^{\text{des}}(\mathbf{x})$ producing the agents' desired control input. For example, each agent may wish to ignore the others and follow a nominal trajectory. Then, for a choice of CBF parameters $b$ and $\alpha$, we seek to find a responsibility allocation vector $\gamma$ such that $u_{1:N}^{\text{data}} \approx \text{proj}(\pi_{1:N}^{\text{des}}(\mathbf{x}^{\text{data}}); \mathbf{x}^{\text{data}}, b, \alpha, \gamma)$ for some choice of similarity loss metric $\Delta$ (e.g., L1, L2, Huber loss). Since the projection mapping proj contains an optimization problem (see Prob. 3), the resulting regression problem takes on a bi-level structure.

*Problem 4 (Responsibility allocation inference):*

$$\min_{\gamma} \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}^i, u_{1:N}^i) \in \mathcal{D}} \Delta(u_{1:N}^i, \widetilde{u}_{1:N}^i(\gamma))$$

$$\text{s.t. } \widetilde{u}_{1:N}^i(\gamma) = \underbrace{\text{proj}(\pi_{1:N}^{\text{des}}(\mathbf{x}^i); \mathbf{x}^i, b, \alpha, \gamma)}_{\text{See Prob. 3}},$$

$$\mathbf{1}^T\gamma = 1, \quad 0 \leq \gamma_i \leq 1.$$

Prob. 4 consists of *multiple* inner problems, one for each data point. This problem structure promotes the use of parallelization as an efficient solution method. Next, we describe how we can efficiently solve Prob. 4.

### B. Leveraging differentiable optimization

We exploit recent results in differentiable optimization [19], [21]–[24] and automatic differentiation tools [25]–[27] to tractably solve Prob. 4 using gradient descent. We first present the gradient descent algorithm and then discuss how this can be solved efficiently.

The procedure for taking a single gradient descent step is described in Alg. 1. Line 1 ensures that the constraints on $\gamma$ are met rather than employing a projected gradient descent. Line 2 describes solving many instances of Prob. 3 (i.e., batching the computation) and collating the solutions as $\widetilde{U}_{1:N}$ which is used in computing the loss in line 3. Line 4 performs the update step with step size $\delta$. The gradient $\frac{\partial \ell}{\partial \gamma} = \frac{\partial \ell}{\partial \gamma} \frac{\partial \gamma}{\partial \widetilde{\gamma}}$ requires differentiating through *all* the instances of Prob. 3 with respect to $\gamma$.

Referring back to Remark 1, if the dynamics are control affine, then Prob. 3 is a quadratic program. This means that (i) the problem is convex and can be efficiently solved, and (ii) we can differentiate through convex programs [19], [21]–[24], and (iii) with the recent advancements in automatic differentiation tooling and batching, such as JAX [25], we can efficiently solve many instances of a QP and differentiate through them in parallel [19].

### C. Context-dependent responsibility

Sec. V-B described the algorithm for finding the responsibility allocation parameter $\gamma$ using gradient descent. That formulation assumed $\gamma$ is fixed for each agent throughout the entire dataset. However, in general, the responsibility allocation is not agent-specific, but rather state- or context-dependent. Given that we are solving Prob. 4 using gradient descent, Alg. 1 can easily extend to the case where $\gamma = \text{softmax}(h_\theta(\mathbf{x}, e))$ where $h_\theta : \mathcal{X}^N \times \mathcal{E} \rightarrow [0,1]^N$ is a differentiable function parameterized by $\theta$, such as a neural network, and depends on the state $\mathbf{x}$ and environment information $e$, such as map/road information. The softmax function ensures the $\gamma$ constraints are satisfied. As such, we can perform gradient descent on the parameters $\theta$ instead.

### D. Enforcing symmetry in responsibility allocation

In learning state and/or context-dependent responsibility allocations, there may be some symmetries we would like to exploit. In particular, the assigned number ordering of agents (i.e., who is labeled Agent 1, Agent 2,..., Agent $N$) should not matter. We introduce a *symmetry responsibility*.

*Definition 3 (Symmetric responsibility):* For an $N$ agent system with joint state space $\mathcal{X}^N$, suppose some assignment ordering $(1, 2, ..., N)$. Let $\gamma : \mathcal{X}^N \rightarrow [0,1]^N$ be the responsibility allocation function, and $\gamma_i(\mathbf{x})$ (with $\mathbf{x} = [x_1, ..., x_N]$) denote the responsibility allocation for agent $i$. Then, we say the system has symmetric responsibility allocation if the following properties hold:

$$\gamma_i(\mathbf{x}) = \gamma_i(\sigma(\mathbf{x})) \, \forall \sigma \in \mathcal{P}_i, \, \forall i = 1, ..., N \quad (4)$$

$$\gamma_i(\mathbf{x}) = \gamma_j(\mathbf{x}_{i \leftrightarrow j}) \, \forall i, j = 1, ..., N \quad (5)$$

$$\sum_{i=1}^{N} \gamma_i(\mathbf{x}) = 1 \quad (6)$$

where $\mathcal{P}_i$ is the set of all permutations (for $N$ agents) that keep index $i$ fixed, and $\mathbf{x}_{i \leftrightarrow j}$ denotes that the entries at indices $i$ and $j$ are swapped.

The first equation (4) describes that the responsibility allocation for agent $i$ is invariant when the assignment ordering of other agents were permuted. The second equation (5) describes that for a given agent, their responsibility allocation should remain the same regardless of their assigned numbering. The third property (6) ensures that the sum of all agents' responsibilities is equal to one. We can construct a function satisfying the symmetric responsibility property.

First, let us define $\sigma : \mathcal{X}^N \rightarrow \mathcal{X}^N$ as a permutation function that permutes the input ordering. We further write $\sigma_{i \rightarrow j}$ to denote any permutation that moves index $i$ to $j$.

*Proposition 1:* Let $\widetilde{\phi} : \mathcal{X}^N \rightarrow \mathbb{R}$ be any (differentiable) function, e.g., a neural network. Then we define $\gamma_i$,

$$\gamma_i(\mathbf{x}) = \frac{\exp(\phi(\sigma_{i \rightarrow 1}(\mathbf{x})))}{\sum_{j=1}^{N} \exp(\phi(\sigma_{j \rightarrow 1}(\mathbf{x})))}, \quad \phi(\mathbf{x}) = \sum_{\sigma \in \mathcal{P}_1} \widetilde{\phi}(\sigma(\mathbf{x})) \quad (7)$$

Then $\gamma(\mathbf{x}) = [\gamma_1(\mathbf{x}), ..., \gamma_N(\mathbf{x})]$ satisfies the symmetric responsibility properties.
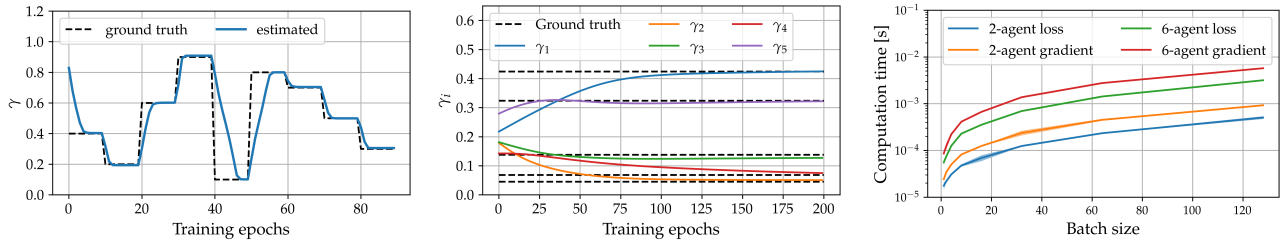
(a) Time-varying $\gamma$ for a 2-agent system.     (b) $\gamma$ for each agent in a 6-agent system.     (c) Computation time as a function of batch size.

Fig. 3: Learning $\gamma$ with synthetically generated data. After a few epochs, the estimated responsibility allocation value $\gamma$ converges to the ground truth value. The computation time scales linearly as the batch size increases.
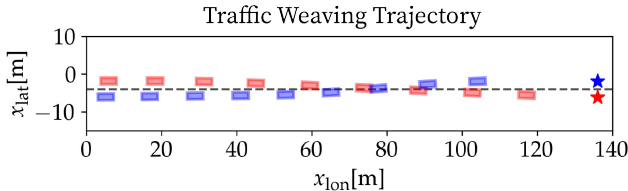


Fig. 4: An example trajectory from the traffic weaving dataset. Two cars (red and blue rectangles) start in adjacent lanes, and need to change lanes safely to reach their lateral goals (stars).

*Proof:* First, we prove that $\gamma_i$ satisfies (4) for all $i$. By construction, $\exp(\phi(\sigma_{i\rightarrow1}(\mathbf{x})))$ is constant under permutations from $\mathcal{P}_i$. We need to show that the denominator of $\gamma_i$ is invariant under any permutation. For any permutation $\sigma$, denote $\zeta(j)$ as the index that maps to $j$. We write $\mathbf{x}_{k\rightarrow l}$ to highlight that index $k$ has been moved to index $l$, with other indices potentially also shuffled. For any $\sigma \in \mathcal{P}$, we have,

$$\sum_{j=1}^{N} \exp(\phi(\sigma_{j\rightarrow1}(\sigma(\mathbf{x})))) = \sum_{j=1}^{N} \exp(\phi(\sigma_{j\rightarrow1}(\mathbf{x}_{\zeta(j)\rightarrow j})))$$
$$= \sum_{j=1}^{N} \exp(\phi(\sigma_{\zeta(j)\rightarrow1}(\mathbf{x})))$$
$$= \underbrace{\sum_{k=1}^{N} \exp(\phi(\sigma_{k\rightarrow1}(\mathbf{x})))}_{\text{Denominator of } \gamma_i}$$

Leveraging the fact that the denominator is invariant under permutations, property (5) follows,

$$\gamma_j(\mathbf{x}_{i\leftrightarrow j}) = \frac{\exp(\phi(\sigma_{j\rightarrow1}(\mathbf{x}_{i\leftrightarrow j})))}{\sum_{k=1}^{N} \exp(\phi(\sigma_{k\rightarrow1}(\mathbf{x}_{i\leftrightarrow j})))}$$
$$= \frac{\exp(\phi(\sigma_{i\rightarrow1}(\mathbf{x})))}{\sum_{k=1}^{N} \exp(\phi(\sigma_{k\rightarrow1}(\mathbf{x})))}$$
$$= \gamma_i(\mathbf{x}).$$

Property (6) is met via use of the softmax function. ∎

Constructing a symmetric responsibility allocation function amounts to defining a single function $\widetilde{\phi} : \mathcal{X}^N \rightarrow \mathbb{R}$ and applying all possible $N-1$ permutations on it. As such, this construction can become prohibitively expensive as $N$ increases. In this work, we study cases with small $N$ values for which the computation is tractable, and defer investigating scalable approaches for large $N$ in future work.

### E. Symmetry in relative coordinates for two-agent systems

In the case of a two-agent setting, it may be more convenient to express the system in terms of *relative dynamics* rather than concatenating the states. Let $\mathbf{r} = (r_1, ..., r_n) \in \mathcal{R}$ denote the relative coordinate. The state space formulation from Sec. V-D must be transferred to the relative coordinate frame. Swapping the assignment of agents is equivalent to negating the relative coordinates. Thus we have $\gamma_1(\mathbf{r}) + \gamma_1(-\mathbf{r}) = 1$. Essentially, rather than requiring invariance under permutations, we require invariance under negation.

We can similarly construct a function satisfying this complementary symmetry for relative coordinates. Let $\phi : \mathcal{R} \rightarrow \mathbb{R}$ be any function, e.g., neural network. Then,

$$\gamma_1(\mathbf{r}) = \frac{1}{2}(1 + \tanh(\phi(\mathbf{r}) - \phi(-\mathbf{r})) \tag{8}$$

satisfies the analogous symmetry conditions.

### VI. EXPERIMENTS ON SYNTHETIC DATASET

We evaluate our approach on a synthetically generated dataset to verify that we can recover the ground truth values. We consider a 2-agent system with 1D single integrator dynamics (see Example 1) and a 6-agent system with 2D double integrator dynamics. We use a CBF that keeps the distance between the (closest) agents at least 1 unit apart. Note that the double integrator is a relative degree two system, so we employ high order CBFs. As mentioned in Remark 3, the high order CBF inequality requires computing additional terms for the CBF inequality constraint, but otherwise, the overall approach remains the same.

For each system, we randomly sample 128 state and desired control pairs and then, for some choice of $\gamma$, we apply the responsibility-aware CBF filter described in Prob. 3, and then add zero-mean, Gaussian noise (variance of 0.1) to the solutions. Then we execute Alg. 1 to learn $\gamma$. As shown in Figs. 3a and 3b, for a random initial guess, our estimated responsibility allocation $\gamma$ quickly converges to the ground truth value, even if the ground truth is time-varying (see Fig. 3a). We use a mini-batch size of 8 and a step size of $\delta = 0.005$ and $\delta = 0.05$ for the 2-agent and 6-agent problems, respectively. Additionally, Fig. 3c presents computation times for calculating the loss and gradient for Prob. 4. We see that our approach has great potential for real-time applications, such as estimating responsibility allocation online, which would only require small batches of data that are streaming in real-time. The experiments were performed on a MacBook Pro with a M2 processor.

### VII. EXPERIMENTS ON TRAFFIC-WEAVING INTERACTION

We demonstrate the efficacy of our responsibility allocation learning approach on a traffic-weaving interaction
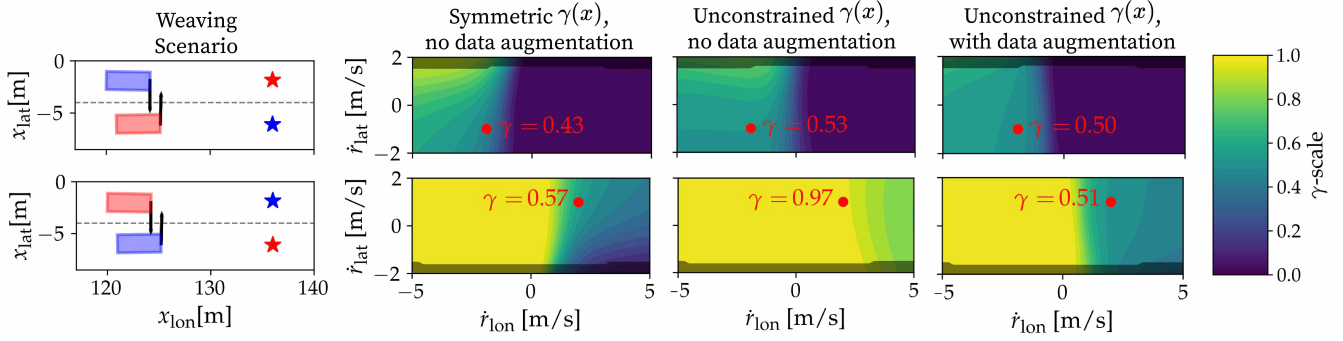
Fig. 5: Comparison of $\gamma$ landscapes for the red car using symmetric and unconstrained models. In the scenario, both cars want to change lanes (stars denote each agent's desired goal). The car in the lower lane is moving faster. The bottom row is the same scenario but with the agents swapped. The black arrows in the left plot represent each agent's desired control vector, and the red point in the contour plots corresponds to the depicted scenarios. The shaded regions denote relative states for which the CBF filter is not active.
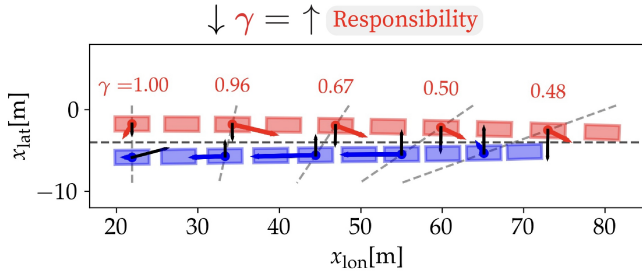


Fig. 6: Learned responsibility allocations from a single traffic weaving trajectory. The shown $\gamma$ values correspond to the red agent, and the black and colored arrows represent each agent's desired and actual control, respectively (magnitudes enlarged for visual clarity).

dataset, containing trajectories of two human drivers (in a driving simulator) quickly swapping lanes within a short distance, reminiscent of merging at highway on/off ramps [28]. An example is shown in Fig. 4. This traffic-weaving dataset presents an interesting case study as it contains diverse driving styles from different drivers with no explicit driving rules describing how drivers yield to each other, in contrast to, say, entering a roundabout where cars in the roundabout have the right of way.

### A. Experimental set up

#### 1) Traffic-weaving dynamics and CBF construction

We model each agent with double integrator dynamics, and consider their relative state, $\mathbf{r} = [r_{\text{lon}}, r_{\text{lat}}, \dot{r}_{\text{lon}}, \dot{r}_{\text{lat}}] = \mathbf{x}_2 - \mathbf{x}_1$ where $\mathbf{x}_i = [x_{\text{lon}}, x_{\text{lat}}, \dot{x}_{\text{lon}}, \dot{x}_{\text{lat}}]^T$ is the state of agent $i$, and lon and lat denotes the longitudinal and lateral directions respectively. With this relative double integrator dynamics, we select a CBF $b(\mathbf{r}) = \frac{r_{\text{lon}}^2}{a_1^2} + \frac{r_{\text{lat}}^2}{a_2^2} - 1$. We set $a_1 = 9.22$ and $a_2 = 1.76$ based on the car's sizes. Due to the double integrator being relative degree two, we employ a high-order CBF (see Remark 3).

#### 2) Desired control.

Our proposed approach requires knowledge of each agent's *desired control*. Unfortunately, this information is not typically known. In this paper, we use a desired control scheme handcrafted from simple heuristics and empirical observation, leaving exploration of learning agents' desired control policies for future work.

For an agent at state $\mathbf{x}$, its desired lateral control is,

$$\pi_{\text{lat}}^{\text{des}}(\mathbf{x};\ x_{\text{lat}}^{\text{des}}) = -(x_{\text{lon}} + \zeta)\beta \tanh(\alpha(x_{\text{lat}} - x_{\text{lat}}^{\text{des}})), \quad (9)$$

where $\zeta$, $\beta$ and $\alpha$ are tunable parameters, set as: $\zeta = 4.7$, $\beta = 0.022$, $\alpha = 0.8$, and $x_{\text{lat}}^{\text{des}}$ is the agent's desired lateral position, i.e., the center of the other lane. Intuitively, the agent's lateral control/steering will get more drastic the longer they take to change lanes.

An agent's desired longitudinal control is,

$$\pi_{\text{lon}}^{\text{des}}(\mathbf{r} = \mathbf{x}' - \mathbf{x}) = \begin{cases} 0 & r_{\text{lon}} > 0 \\ \frac{-\kappa}{2}(\tanh(r_{\text{lon}}\dot{r}_{\text{lon}}) - 1) & r_{\text{lon}} \leq 0 \end{cases} \quad (10)$$

where an agent, if ahead of the other agent, is encouraged to speed up to achieve a larger relative longitudinal distance up to some threshold. Otherwise, if the agent is behind the other agent, it will maintain speed. The parameter $\kappa$ influences the limits of the desired control, which we set to $\kappa = 2$.

#### 3) Dataset preparation.

As mentioned previously, the traffic-weaving dataset consists of diverse behaviors owing to the various driving styles of the driving participants. As such, we consider four different subsets of the dataset to perform our analysis. **(D1)** Single traffic-weaving trajectory. **(D2)** Trajectories where both cars start side-by-side at equal initial longitudinal speed. **(D3)** Trajectories where one car starting behind but traveling faster overtakes the other car. **(D4)** Full dataset.

Before training, we consider two types of data augmentation: **(A1)** mirroring the agents' lateral states and controls across the lane divider ($x_{\text{lat}}, \dot{x}_{\text{lat}}, u_{\text{lat}}, u_{\text{lat}}^{\text{des}} \rightarrow -x_{\text{lat}}, -\dot{x}_{\text{lat}}, -u_{\text{lat}}, -u_{\text{lat}}^{\text{des}}$), and **(A2)** swapping the states of the agents ($\mathbf{r} \rightarrow -\mathbf{r}$). We apply augmentation A1 to every data subset to remove bias in which lane an agent starts. We apply augmentation A2 on dataset D1 to verify the benefits of a symmetric $\gamma$-function; we expand on this in Sec. VII-B. Note: A2 is the "data augmentation" referred to in Fig. 5.

#### 4) Model parameterization and training.

We parameterize the responsibility allocation function as a neural network $\phi$ that takes in the relative state $\mathbf{r}$ and outputs the responsibility allocation. Our $\phi(\mathbf{r})$ function is a small multi-layer perception with three linear hidden layers, hidden
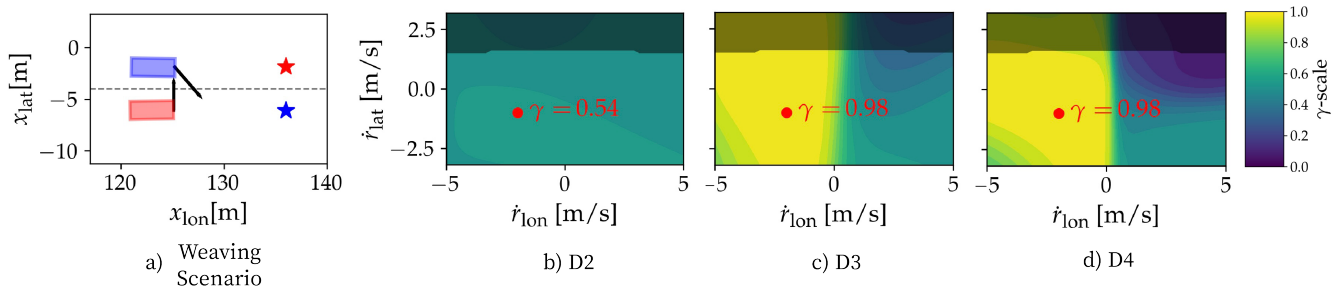
Fig. 7: Comparison of $\gamma$ landscapes when learned using different subsets of the traffic weaving dataset. In the displayed scenario, the red car is moving faster than the blue car, and the black arrows represent their desired controls as they move to their goals.

layer size of sixteen, and $\tanh$ activation functions, whose final output is passed through the function in (8) to enforce symmetry. Our setup used JAX [25] and Equinox [29]. We trained on each dataset for 3000 epochs with batch sizes of 16 for the single weaving trajectory and 256 for the other datasets, and used the ADAM optimizer [30] with a learning rate of $10^{-3}$. Our loss function $\Delta$ from Prob. 4 is the Huber loss, and the regularization/slack variables $\beta_1$ and $\beta_2$ from Prob. 1 are tuned to 0.1 and 600 respectively.

### B. Training on a single trajectory (D1)

We first apply our learning method to the D1 dataset for a sanity check and validate our intuition. In this single traffic-weaving trajectory (similar to Fig. 4), the car in the upper lane starts 2m/s faster and 3m behind the lower car, before overtaking the lower car and switching lanes first.

#### 1) Benefits of using a symmetric model

Fig. 5 shows the learned responsibility allocation for the red car $\gamma_{\text{red}}$ (and the blue car's responsibility allocation is $\gamma_{\text{blue}} = 1 - \gamma_{\text{red}}$). The top and bottom rows are in the same scenario except for the cars swap places. This means we expect $\gamma_{\text{red}}^{\text{top scene}} = \gamma_{\text{blue}}^{\text{bottom scene}}$. This is precisely the symmetric responsibility property, and we see that, as expected, a symmetric model satisfies this property exactly whereas the unconstrained model does not.

Additionally, note that the results with the symmetric model require no augmenting of swapped players' states, and we see that it achieves very similar results to the unconstrained model using this data augmentation. The unconstrained model without data augmentation has poor results when the cars are in a configuration unseen in the original trajectory. These results highlight the data efficiency benefits of using a symmetric model—data augmentation may be unnecessary, a benefit when working with more complex scenarios where data augmentation is less trivial. Moving forward, we use a symmetric responsibility allocation model.

#### 2) Interpreting $\gamma$ over a trajectory

Fig. 6 shows how the value of $\gamma$ changes during the trajectory as the two cars swap lanes. The evolution of $\gamma$ is consistent with the alignment between the agents' actual (colored) control and desired (black) control vectors; initially, the red car has low responsibility as it does not yet need to change lanes. However, after overtaking the blue car, its responsibility increases because it becomes more eager to swap lanes but cannot because of the blue car. The blue

car, meanwhile, has to brake early in the trajectory when it wants to move laterally (hence the greater responsibility early on), as it must accommodate the red car overtaking it. This example illustrates how we can learn intuitive responsibility allocations for a basic traffic weaving trajectory.

### C. Training on data with same initial conditions (D2)

We apply our learning method to the D2 dataset to see what happens when there are diverse and contradictory behaviors within the dataset. As the initial conditions place the cars on equal footing at the start, there is ambiguity as to who passes whom, resulting in multimodal behavior within the dataset. Learning responsibility allocations in this setting is substantially more difficult; as shown by Fig. 7b. The learned responsibility allocation is constant—regardless of their relative state, both agents are predicted to be equally responsible. This performance drop under multimodal interactions highlights a limitation in our approach, necessitating future work on a *probabilistic* extension.

### D. Training on data when rear, faster car overtakes (D3)

We apply our learning method to the D3 dataset to investigate if it can successfully capture the trend that the slower car yields to the faster car. In Fig. 7c, we see that with the same longitudinal position, the faster car has a much larger $\gamma$ value, meaning they are less responsible for maintaining safety, consistent with behavior in the dataset where the slower car yields to the faster car.

### E. Training on the whole dataset (D4)

Finally, we apply our learning methods to the entire traffic-weaving dataset (D4). Fig. 7d shows that it learned a responsibility allocation model similar to the model trained on D3. Despite D4 containing more diverse trajectories, there is a bias towards cars starting behind but traveling faster to overtake rather than yield to the slower but front car [28]. We see that our model can capture the bias in the dataset.

### F. Discussion

From the results of our different training scenarios on the traffic weaving dataset, we draw three key takeaways:

1) Our responsibility model and inference technique was able to efficiently perform on real human trajectory data, assigning interpretable quantities to observed behaviors. This affirms that we can transfer our definition of responsibility allocation from theory to practice.

2) Our approach learned nontrivial responsibility allocations on data with consistent bias (e.g., a significant portion of data having the faster car overtake the slower car), but struggled with data containing multimodal behaviors. This suggests a need to investigate probabilistic extensions to capture multimodality.

3) Identifying the "best" desired control policy is nontrivial. While our handcrafted nonlinear control policy for the traffic weaving data lets us learn intuitive responsibility allocations, how we can do better is unclear; investigating principled approaches to learning desired policies from data is a logical next step.

## VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we introduced a novel responsibility allocation framework that distills ill-defined social norms underpinning human interactions into an interpretable quantity. We presented a mathematical formulation of responsibility rooted in Control Barrier Functions and a computationally efficient approach to learn responsibility allocations from data via differentiable optimization. We also proposed an approach to learning symmetric responsibility allocations which can help with data efficiency. We demonstrated the efficacy of our algorithm on synthetic and human interaction data, and results indicate that we can gain interpretable insights into the social dynamics of multi-agent interactions. For future work, we seek to **(i)** investigate principled ways to construct the desired control policy, e.g. via learned generative models, **(ii)** develop a probabilistic extension to this framework to account for multimodal interactions, and **(iii)** examine the use of responsibility allocations for guiding robot policy construction.

## REFERENCES

[1] N. Nayakanti, R. Al-Rfou, A. Zhou, K. Goel, K. S. Refaat, and B. Sapp, "Wayformer: Motion forecasting via simple and efficient attention networks," in *Proc. IEEE Conf. on Robotics and Automation*, 2023.

[2] B. Ivanovic, K. Leung, E. Schmerling, and M. Pavone, "Multimodal deep generative models for trajectory prediction: A conditional variational autoencoder approach," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 295–302, 2021.

[3] K. Mizuta and K. Leung, "CoBL-Diffusion: Diffusion-based conditional robot planning in dynamic environments using control barrier and lyapunov functions," in *IEEE/RSJ Int. Conf. on Intelligent Robots & Systems*, 2024.

[4] D. Sadigh, S. Sastry, S. A. Seshia, and A. D. Dragan, "Planning for autonomous cars that leverage effects on human actions," in *Robotics: Science and Systems*, 2016.

[5] C. I. Mavrogiannis, P. Alves-Olivera, W. Thomason, and R. A. Knepper, "Social momentum: Design and evaluation of a framework for socially competent robot navigation," *ACM Transactions on Human-Robot Interaction*, vol. 37, no. 4, 2021.

[6] W. Schwarting, A. Pierson, J. Alonso-Mora, S. Karaman, and D. Rus, "Social behavior for autonomous vehicles," *Proceedings of the National Academy of Sciences*, vol. 116, no. 50, pp. 24 972–24 978, 2019.

[7] B. Toghi, R. Valiente, D. Sadigh, R. Pedarsani, and Y. P. Fallah, "Cooperative autonomous vehicles that sympathize with human drivers," in *IEEE/RSJ Int. Conf. on Intelligent Robots & Systems*, 2021.

[8] B. Toghi, R. Valiente, D. Sadigh, R. Pedarsani, and Y. Fallah, "Social coordination and altruism in autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 23, no. 12, pp. 24 791–24 804, 2022.

[9] L. Sun, W. Zhan, M. Tomizuka, and A. Dragan, "Courteous autonomous cars," in *IEEE/RSJ Int. Conf. on Intelligent Robots & Systems*, 2018.

[10] J. van den Berg, M. Lin, and D. Manocha, "Reciprocal velocity obstacles for real-time multi-agent navigation," in *Proc. IEEE Conf. on Robotics and Automation*, 2008.

[11] Z. Yin, J. Liu, and L. Wang, "Less-effort collision avoidance in virtual pedestrian simulation," in *ACM Int. Conf. on Artificial Intelligence and Computer Science*, 2019.

[12] Y. Luo, P. Cai, A. Bera, D. Hsu, W. S. Lee, and D. Manocha, "PORCA: Modeling and planning for autonomous driving among many pedestrians," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3418–3425, 2018.

[13] Y. Luo, Y. Cai, P. Lee, and D. Hsu, "GAMMA: A general agent motionprediction model for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3499–3506, 2022.

[14] K. Guo, D. Wang, T. Fan, and J. Pan, "VR-ORCA: Variable responsibility optimal reciprocal collision avoidance," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4520–4527, 2021.

[15] R. Cosner, Y. Chen, K. Leung, and M. Pavone, "Learning responsibility allocations for safe human-robot interaction with applications to autonomous driving," in *Proc. IEEE Conf. on Robotics and Automation*, 2023.

[16] A. D. Ames, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs with application to adaptive cruise control," in *Proc. IEEE Conf. on Decision and Control*, 2014.

[17] A. D. Ames, X. Xu, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs for safety critical systems," *IEEE Transactions on Automatic Control*, vol. 62, no. 8, pp. 3861–3876, 2017.

[18] B. Stellato, G. Banjac, P. Goulart, A. Bemporad, and S. Boyd, "OSQP: An operator splitting solver for quadratic programs," *Available at https://arxiv.org/abs/1711.08013*, 2017.

[19] K. Tracy and Z. Manchester, "On the differentiability of the primal-dual interior-point method," *Available at https://arxiv.org/abs/2406.11749*, 2024.

[20] W. Xiao and C. Belta, "High order control barrier functions," *IEEE Transactions on Automatic Control*, vol. 67, no. 7, pp. 3655–3662, 2021.

[21] B. Amos and J. Z. Kolter, "Optnet: Differentiable optimization as a layer in neural networks," in *Int. Conf. on Machine Learning*, 2017.

[22] A. Agrawal, B. Amos, S. Barratt, S. Boyd, S. Diamond, and Z. Kolter, "Differentiable convex optimization layers," in *Conf. on Neural Information Processing Systems*, 2019.

[23] A. Agrawal, S. Barratt, S. Boyd, and B. Stellato, "Learning convex optimization control policies," in *Learning for Dynamics & Control Conference*, 2020.

[24] A. Agrawal, S. Barratt, and S. Boyd, "Learning convex optimization models," *IEEE/CAA Journal of Automatica Sinica*, vol. 8, no. 8, pp. 1355–1364, 2021.

[25] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, "JAX: composable transformations of Python+NumPy programs," *Available at http://github.com/google/jax*, 2018.

[26] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in PyTorch," in *Conf. on Neural Information Processing Systems - Autodiff Workshop*, 2017.

[27] J. Revels, M. Lubin, and T. Papamarkou, "Forward-mode automatic differentiation in julia," *Available at https://arxiv.org/abs/1607.07892*, 2016.

[28] E. Schmerling, K. Leung, W. Vollprecht, and M. Pavone, "Multimodal probabilistic model-based planning for human-robot interaction," in *Proc. IEEE Conf. on Robotics and Automation*, 2018.

[29] P. Kidger and C. Garcia, "Equinox: neural networks in JAX via callable PyTrees and filtered transformations," in *Conf. on Neural Information Processing Systems*, 2021.

[30] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Int. Conf. on Learning Representations*, 2015.