

Learning Mixed Strategies in Trajectory Games

Peters, L.; Fridovich-Keil, David; Ferranti, L.; Stachniss, Cyril; Alonso-Mora, Javier; Laine, Forrest

DOI

[10.15607/RSS.2022.XVIII.051](https://doi.org/10.15607/RSS.2022.XVIII.051)

Publication date

2022

Document Version

Final published version

Published in

Proceedings Robotics: Science and System XVIII

Citation (APA)

Peters, L., Fridovich-Keil, D., Ferranti, L., Stachniss, C., Alonso-Mora, J., & Laine, F. (2022). Learning Mixed Strategies in Trajectory Games. In K. Hauser, D. Shell, & S. Huang (Eds.), *Proceedings Robotics: Science and System XVIII* (Robotics: Science and Systems). Robotics Science and Systems (RSS). <https://doi.org/10.15607/RSS.2022.XVIII.051>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Learning Mixed Strategies in Trajectory Games

Lasse Peters* David Fridovich-Keil† Laura Ferranti* Cyrill Stachniss‡ Javier Alonso-Mora* Forrest Laine§

* Delft University of Technology, Netherlands † University of Texas at Austin, USA

‡ University of Bonn, Germany § Vanderbilt University, USA

{ l.peters, l.ferranti, j.alonsomora }@tudelft.nl,

dfk@utexas.edu, cyrill.stachniss@igg.uni-bonn.de, forrest.laine@vanderbilt.edu

Abstract—In multi-agent settings, game theory is a natural framework for describing the strategic interactions of agents whose objectives depend upon one another’s behavior. Trajectory games capture these complex effects by design. In competitive settings, this makes them a more faithful interaction model than traditional “predict then plan” approaches. However, current game-theoretic planning methods have important limitations. In this work, we propose two main contributions. First, we introduce an offline training phase which reduces the online computational burden of solving trajectory games. Second, we formulate a *lifted* game which allows players to optimize multiple candidate trajectories in unison and thereby construct more competitive “mixed” strategies. We validate our approach on a number of experiments using the pursuit-evasion game “tag.”

I. INTRODUCTION

Trajectory optimization techniques have become increasingly common in motion planning. So long as vehicle dynamics, design objectives, and safety constraints satisfy mild regularity conditions, a motion planning problem may be encoded as a nonlinear program and solved efficiently to a locally-optimal solution. The widespread successes of trajectory optimization have sparked growing interest in similar techniques for multi-agent, noncooperative decision-making and motion planning. In this context, game theory offers an elegant mathematical framework for modeling the strategic interactions of rational agents with distinct interests. By reasoning about interactions with others as a *trajectory game*, an autonomous agent can plan future decisions while accounting for the strategic reactions of others.

Since they involve multiple players with distinct, potentially competing objectives, trajectory games can be far more complex to solve than single-agent trajectory optimization problems. Recent algorithmic advances make solving trajectory games tractable in some instances [10, 15]. Still, they remain fundamentally more challenging to solve than single-agent problems, and consequently, trajectory games have not been widely adopted in the robotics community.

Perhaps more importantly, however, equilibrium solutions to trajectory games do not always exist. Nonexistence arises even in extremely simple static games such as rock-paper-scissors, in which neither player wishes to commit to a fixed, deterministic action which could be exploited by its opponent. Unsurprisingly, the same phenomenon can arise in more complex trajectory games. For example, consider

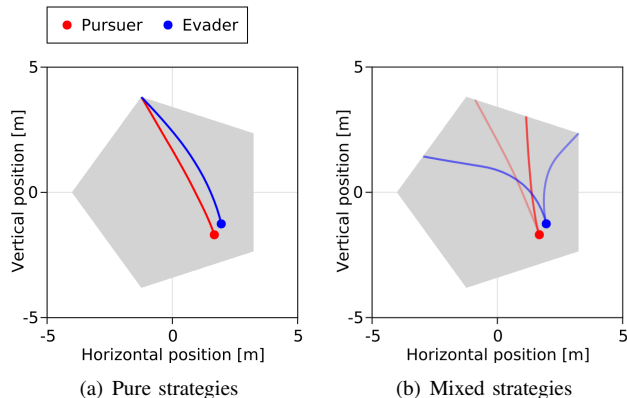


Fig. 1: A zero-sum game of tag played between two agents with planar point-mass dynamics in a pentagonal environment. (a) In pure strategies, players are bound to deterministic behavior, and the evader is quickly captured. (b) Our approach lifts the strategy space to learn more competitive, mixed strategies, i.e., distributions over multiple trajectory candidates per player. The opacity of each trajectory in (b) encodes probability of selecting that learned candidate.

the game of tag shown in Figure 1, where the red pursuer wishes to catch the blue evader. Here, if the evader chooses a single, deterministic trajectory, it will certainly be caught by a rational pursuer. In the context of small, discrete games such as rock-paper-scissors, these non-existence issues are commonly avoided by allowing players to “mix” their actions, i.e., to choose an action at random from a distribution of their choice. This distribution is called a “mixed strategy,” in contrast to the choice of a single deterministic action or “pure strategy.” However, for continuous trajectory games it can be difficult to represent mixed strategies. Hence, it is common to regularize players’ objectives in order to encourage the existence of pure solutions.¹ For example, in Figure 1(a) each player is penalized for large accelerations, leading to an equilibrium in which the evader is cornered by the pursuer.

With these issues in mind, this paper introduces the following key contributions:

¹Regularizing players’ control inputs to ensure the existence of equilibria is well-established in the literature on dynamic games and robust control [5], [4].

- 1) a principled method for reducing the online computation needed to solve trajectory games via the introduction of an offline training phase, and
- 2) a formulation of *lifted* games over multiple trajectory candidates, which admit a natural class of high-performance mixed strategies.

Together, these contributions enable efficient and reliable online trajectory planning for autonomous agents in noncooperative settings, such as the tag example of Figure 1. We validate our methods in a suite of Monte Carlo studies, in which we demonstrate that lifting gives rise to mixed strategies as shown in Figure 1(b), providing a significant competitive advantage in both open-loop and receding-horizon play. Our method’s reliable convergence and its ability to explicitly account for constraints enables training from scratch within only a few minutes of simulated self-play. Once fully trained, learning can be disabled and our method generates mixed strategies within 2 ms for the tag example in Figure 1.

II. RELATED WORK

Our contributions build upon recent work in trajectory optimization and game-theoretic planning, and bear a close relationship with work in learning motion primitives and implicit differentiation. We discuss these relationships in further detail below.

A. Trajectory Optimization

Trajectory optimization refers to a finite-horizon optimal control problem in which a robot seeks a sequence of control inputs which minimize a performance criterion [30]. It is common to use trajectory optimization for model predictive control (MPC), whereby a robot quickly re-optimizes a new sequence of control inputs as new sensor data becomes available [8]. While a host of trajectory optimization techniques have been proposed in recent years, most common algorithms build upon the iterative linear-quadratic regulator [28, 29, 42] and differential dynamic programming [21, 41, 46]. In turn, these may be understood as specific approximations to standard algorithms in nonlinear programming (NLP), such as sequential quadratic programming [6, 33]. As discussed below, this fundamental NLP representation underlies the proposed approach for multi-agent trajectory games.

B. Trajectory Games

Recent work has sought to generalize the aforementioned trajectory optimization techniques to address multi-agent, competitive planning. Here, each player seeks to minimize an individual performance criterion subject to constraints arising from, e.g., dynamics and actuator limits. The objectives and constraints for different players may, in general, depend upon the trajectories of others. Solutions to these problems are characterized by equilibrium points at which all players’ strategies are unilaterally optimal.

The theoretical underpinnings of *dynamic* games were established in the context of state feedback [5, 20, 36, 37]. However, computational methods were historically limited

to highly-structured problems such as those found in linear robust control [4, 18]. Recent work on iterative linear-quadratic methods [15, 24, 25] extends these ideas to more general games such as those found in noncooperative robotic planning.

Closely related problems have also been studied in the context of *static* games. Here, equilibrium points are found by treating the trajectory of each player as a single action, and assuming the players choose these actions simultaneously [5]. This results in a Generalized Nash Equilibrium Problem (GNEP), for which general-purpose solution methods exist [11, 12, 13]. Several domain-specific solvers have been developed to exploit the structure of trajectory games, ranging from augmented Lagrangian [26] to iterated best response methods [43, 44]. Still, these methods can have a high computational burden in challenging settings.

Regardless of the equilibrium definition (dynamic or static), solving trajectory games is fundamentally harder than solving single-agent trajectory optimization problems, if for nothing else but the increased problem dimension. The number of decision variables scales linearly with the number of players involved, and even with proper handling of sparsity, computation generally scales cubically with the number of players [15]. In Section IV-A, we introduce an offline training phase for trajectory games which effectively reduces the online computational burden to that of solving a trajectory optimization problem for each player in parallel.

C. Motion Primitives

In this work, we introduce a *trajectory lifting* technique, which may be understood in the context of motion primitives [22]. As we discuss in Section IV, this reformulation endows each player with a distribution over finitely many trajectory candidates, which may be learned. However, learning trajectories, or motion primitives, is also meaningful in the context of a single agent, and recent work has proposed this concept in the contexts of quadcopter navigation [9] and robot manipulation [40]. In this light, the present paper may be viewed as a multi-agent generalization of these techniques. Additionally, our work constitutes an adaptive, learning-enabled generalization of the multi-agent motion primitive games formulated for autonomous racing in [31]. In Section IV-B, we show that these trajectory primitives may be learned efficiently with first-order optimization.

D. Differentiable Optimization

To improve learning efficiency, we employ implicit differentiation to propagate derivative information through all steps of our proposed trajectory lifting approach. Recent works in end-to-end neural architectures for autonomous driving have developed specialized network layers that embed optimization problems [1, 2, 3]. Like these methods, we obtain derivatives of players’ game values with respect to learnable parameters by implicitly differentiating through the first order optimality conditions for all players in a lifted trajectory game.

III. FORMULATION

We develop our approach in the context of games played between two² agents over time, in which each agent's motion is characterized by a smooth discrete-time dynamical system. That is, we model agent i 's motion as the temporal evolution of its state $\bar{x}_i(t) \in \mathbb{R}^n$ and control input $\bar{u}_i(t) \in \mathbb{R}^m$ over discrete time-steps $t \in \{1, \dots, T\}$ with $\bar{x}_i(t+1) = F(\bar{x}_i(t), \bar{u}_i(t))$ for differentiable vector field $F(\cdot, \cdot)$.

Taking an egocentric approach, we investigate using *model-predictive game play* (MPGP) [15, 26, 44] as a method by which each player can plan strategically while accounting for the predicted reactions of its opponent. MPGP constitutes a natural analogue to MPC [8] for *noncooperative*, multi-agent settings. That is, at regular intervals the 'ego' agent formulates a finite-horizon trajectory game between itself and its opponent. The equilibrium of this game specifies optimal trajectories for both players; the ego agent begins to execute its equilibrium trajectory, and the procedure repeats after a short time once the players have moved.

The finite-horizon trajectory games formulated at each planning interval can be modeled as a pair of coupled optimization problems, as is common in the literature [5, 24]:

$$\begin{aligned} \text{OPT}_1(\tau_2, x_1) &:= \arg \min_{\tau_1} f_1(\tau_1, \tau_2) \\ \text{s.t.} \quad &\tau_1 \in \mathcal{K}_1(x_1) \end{aligned} \quad (1a)$$

$$\begin{aligned} \text{OPT}_2(\tau_1, x_2) &:= \arg \min_{\tau_2} f_2(\tau_1, \tau_2) \\ \text{s.t.} \quad &\tau_2 \in \mathcal{K}_2(x_2) \end{aligned} \quad (1b)$$

The decision variables τ_i for each player $i \in \{1, 2\}$ represent discrete-time state-control trajectories starting from initial configuration x_i . Therefore, the constraint sets $\mathcal{K}_i(x_i)$ represent the set of all trajectories satisfying dynamic constraints, control limits, etc. Note that these sets need not be compact or convex, and that the players' constraint sets are independent of one another's trajectory. In contrast, the differentiable cost functions $f_i(\tau_1, \tau_2)$ in each problem can depend upon *both* players' trajectories. Thus, the f_i can encode preferences such as goal-reaching and collision-avoidance. In particular, since constraints are decoupled, we assume that any aspect of interaction in the game is modeled via the cost functions and not through constraints.

As discussed in Section II-B, existing methods to find local equilibrium solutions of (1) include iterative best response [44, 45] and iterative linear-quadratic methods [10, 15, 24, 26]. A Nash equilibrium for Game (1) starting from initial configuration (x_1, x_2) is defined to be a pair of trajectories, (τ_1^*, τ_2^*) , satisfying

$$\tau_1^* \in \text{OPT}_1(\tau_2^*, x_1) \quad \text{and} \quad \tau_2^* \in \text{OPT}_2(\tau_1^*, x_2). \quad (2)$$

Nash equilibrium points encode rational strategic play for both players, and hence serve as a natural solution concept in trajectory games (1). For this reason, most recent MPGP methods

[10, 15, 24, 26, 45] aim to compute a Nash equilibrium of the trajectory game (1). As a practical matter, however, Nash equilibria can be intractable to compute and modern methods often settle for *local* equilibria, in which players' trajectories are only locally optimal.

Several important issues arise when employing an MPGP approach. The first is that solving for a Nash equilibrium—even a local Nash—is harder than solving for a locally optimal trajectory (as would be done in the single-agent setting of MPC). Not only is the search space larger due to the inclusion of both players' trajectory variables, but potential complications are also introduced by agents' different and potentially conflicting objectives. As in MPC, real-world applications depend upon our ability to compute solutions to (1) quickly; unfortunately though, this increased complexity can make MPGP unsuitable for real-time applications.

The second issue is that a Nash equilibrium point may not even exist for Game (1), particularly when one or both of the subproblems (1a) and (1b) are non-convex. Relatedly, even if a Nash equilibrium does exist, it may not be unique. Consequently, in MPGP an agent may spend significant computational effort searching for an equilibrium point that does not exist. Worse, non-uniqueness implies that even if an agent finds an equilibrium, the opponent's predicted Nash trajectory may not be representative of its true strategy.

To make these issues more concrete, consider the following "toy" variant of the tag game in Figure 1. Let τ_1 and τ_2 be scalars, $f_1(\tau_1, \tau_2) = \|\tau_1 - \tau_2\|_2^2 = -f_2(\tau_1, \tau_2)$, and $\mathcal{K}_1 = \mathcal{K}_2 = [-1, 1]$. Here, the pursuer (Player 1) and evader (Player 2) choose positions in the interval $[-1, 1]$. By inspection, we may verify that no Nash equilibrium exists. With additional regularization, however, this example can be modified to admit *local* equilibria. With f_1 defined as above, if we redefine the function $f_2(\tau_1, \tau_2) = -\|\tau_1 - \tau_2\|_2^2 - \|\tau_2\|_2^2$, two local equilibrium points result: $(-1, -1)$ and $(1, 1)$. Unfortunately, the locality of these equilibria causes a significant problem: if Player 1 computed one of these equilibria, and Player 2 computed the other, the resulting pairing of actions, e.g. $(-1, 1)$, would have a significantly different outcome for the players than what occurs at either local equilibrium.

IV. APPROACH

We propose a novel *lifted* trajectory game formulation which ameliorates the complexity and existence/uniqueness issues discussed in Section III.

A. Reducing Run-time Computation

To begin, we propose a technique for offloading the complexity introduced by multi-agent interactions to an offline training phase. The result of this pre-training is that at run-time, only a single-agent trajectory optimization problem remains for each player, and these problems can be solved in parallel. To do so, we introduce auxiliary trajectory references ξ_i for each Player i , and with a slight abuse of notation,

²Although we limit our discussion to two players, our formulation may be extended to the general case. For further discussion, refer to Section IV-C.

reformulate Game (1) as:

$$\text{OPT}_1(\xi_2, x_1, x_2) := \arg \min_{\xi_1} f_1(\tau_1, \tau_2) \quad (3a)$$

$$\text{OPT}_2(\xi_1, x_1, x_2) := \arg \min_{\xi_2} f_2(\tau_1, \tau_2) \quad (3b)$$

Here, the decision variables ξ_i and the initial states x_i determine trajectory variables $\tau_i = \text{TRAJ}_i(\xi_i, x_i)$, which we presume to have the form:

$$\begin{aligned} \text{TRAJ}_i(\xi_i, x_i) &:= \arg \min_{\tau} \frac{1}{2} \|G_i \tau - \xi_i\|_2^2 + \frac{1}{2} \|H_i \tau\|_2^2 \\ \text{s.t.} \quad &\tau \in \mathcal{K}_i(x_i). \end{aligned} \quad (4)$$

The first term of the cost functions in problem (4) enables ξ_i to serve as a reference for τ_i . For example, if $\tau_i = [X_i^T U_i^T]^T$, with X_i and U_i representing the state and control variables of the trajectory, then G_i could be $[0 \ I]$, giving ξ_i the interpretation of a control reference signal. Alternatively, ξ_i could represent a reference for the terminal state of the trajectory. The second term allows regularization of the trajectory, which may be needed if the reference and constraint sets are otherwise insufficient to isolate solutions.

In Appendix A, we prove that for any stationary point (τ_1, τ_2) of Game (1), there exists a stationary point (ξ_1, ξ_2) of Game (3) such that for both players i , $\tau_i = \text{TRAJ}_i(\xi_i, x_i)$. This implies that no stationary points are “lost” in the reformulation from (1) to (3). Furthermore, we discuss practical methods to guarantee that all computed stationary points of (3) result in stationary points of (1). This implies that no spurious stationary points are “introduced” in the reformulation.

With this reformulation, it is now possible to offload a significant amount of computation to an offline training phase. To do so, we propose training a *reference generator* for each player, denoted by the function $\pi_{\theta_i}(x_1, x_2)$, which maps both player’s initial states (x_1, x_2) to reference ξ_i . Generator π_{θ_i} is parameterized by θ_i and, e.g., may be a multi-layer perceptron as described in Section IV-D. Given a data set³ of initial MPPG configurations $D := \{x_1^k, x_2^k\}_{k=1}^d$, we train the reference generators (π_{θ_1} and π_{θ_2}) by solving the following game offline:

$$\text{GEN}_1(\theta_2, D) := \arg \min_{\theta_1} \frac{1}{d} \sum_{k=1}^d f_1(\tau_1^k, \tau_2^k), \quad (5a)$$

$$\text{GEN}_2(\theta_1, D) := \arg \min_{\theta_2} \frac{1}{d} \sum_{k=1}^d f_2(\tau_1^k, \tau_2^k). \quad (5b)$$

Similar to Game (3), each trajectory τ_i^k appearing in (5) is a function of θ_i, x_1^k and x_2^k , via the relationships

$$\begin{aligned} \tau_1^k &= \text{TRAJ}_1(\pi_{\theta_1}(x_1^k, x_2^k), x_1^k), \\ \tau_2^k &= \text{TRAJ}_2(\pi_{\theta_2}(x_1^k, x_2^k), x_2^k). \end{aligned} \quad (6)$$

³ D need not be constructed laboriously; in Section V-E we show that it can even be accumulated during online operation.

A Nash equilibrium for Game (5) can be found by simultaneous gradient descent over each player’s reference generator parameters, θ_1 and θ_2 . Simultaneous gradient play is widely used in adversarial machine learning, and is particularly important in both generative adversarial networks [17] and multi-agent reinforcement learning [14]. Here, each player’s parameter θ_i is iteratively updated as $\theta_i \leftarrow \theta_i - \delta\theta_i$, where

$$\begin{aligned} \delta\theta_1 &= \frac{\alpha_1}{d} \nabla_{\theta_1} \sum_{k=1}^d f_1(\text{TRAJ}_1(\theta_1, x^k), \text{TRAJ}_2(\theta_2, x^k)) \\ \delta\theta_2 &= \frac{\alpha_2}{d} \nabla_{\theta_2} \sum_{k=1}^d f_2(\text{TRAJ}_1(\theta_1, x^k), \text{TRAJ}_2(\theta_2, x^k)) \end{aligned} \quad (7)$$

Note that in (7), we use the shorthand $x^k \equiv (x_1^k, x_2^k)$, and although we abbreviate the arguments to the TRAJ functions, they should be interpreted exactly as in (6). The values α_1 and α_2 are learning rates used for the respective reference generators. To compute these gradients, we must differentiate through each player’s objective f_i and through each TRAJ_i . We have assumed *a priori* that the f_i were differentiable. To differentiate through the trajectory optimization step of (4), we follow a procedure similar to what is outlined in [1, 2, 3].

Assuming that offline gradient play converges to a Nash equilibrium over the training set D , and that the resulting trajectory generators generalize to instances of (3) defined by configurations (x_1, x_2) not included in D , then an approximate equilibrium solution to Game (1), denoted by $(\hat{\tau}_1^*, \hat{\tau}_2^*)$ can be found via the following evaluations:

$$\begin{aligned} \xi_1 &= \pi_{\theta_1}(x_1, x_2), & \xi_2 &= \pi_{\theta_2}(x_1, x_2) \\ \hat{\tau}_1^* &= \text{TRAJ}_1(\xi_1, x_1), & \hat{\tau}_2^* &= \text{TRAJ}_2(\xi_2, x_2) \end{aligned} \quad (8)$$

Hence, at run-time, solving this reformulated game only requires evaluating the reference generators and solving the optimization problems TRAJ_i to compute the corresponding trajectories. These problems can be solved in parallel. Furthermore, since trajectories are generated according to (4), each player’s constraints defined by $\mathcal{K}_i(x_i)$ are guaranteed to be satisfied. Thus, if the reference generator does not generalize well, the only negative consequence is suboptimality (but not infeasibility).⁴

In summary, by pre-training a reference generator for each player offline, the run-time concerns of MPPG can be alleviated. Unfortunately, however, potential issues persist due to the possible non-existence or non-uniqueness of Nash equilibrium solutions. To address this concern, we introduce a concept we refer to as *strategy lifting*.

B. Lifted Trajectory Games

Rather than endowing each player with a single reference and its resulting trajectory, we allow each player to choose

⁴Recall that each player’s constraints do not depend upon the trajectory of the other player. Extension to this more complex case is possible, but beyond the scope of this paper.

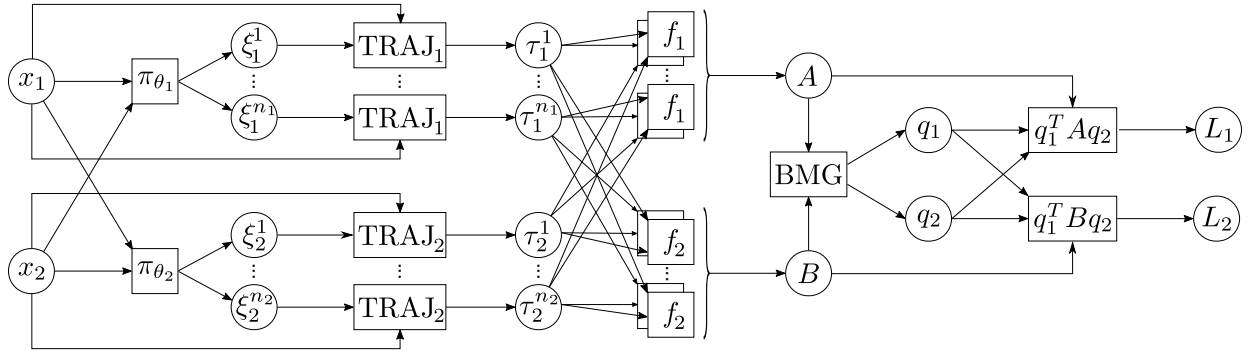


Fig. 2: Overview of our proposed lifted game solver using reference generators. Generators π_{θ_i} for both players are trained jointly to minimize their respective average losses L_i (10a). At run-time deployment of this pipeline, the pre-trained generators produce references ξ_i which approximate a Nash equilibrium of (9). These references evaluate to equilibrium motion plan candidates τ_i and mixing strategies q_i for each player. When used in an ego-centric MPGP fashion, e.g. for Player 1, (τ_1, q_1) serves as a distribution over ego motion plans, and (τ_2, q_2) constitutes a probabilistic opponent prediction.

among multiple independent references according to the equilibrium solution of the bimatrix game formulated below:

$$\text{OPT}_1^{\text{lifted}}(\xi_2, x_1, x_2) := \arg \min_{\xi_1} L_1(\xi_1, \xi_2) \quad (9a)$$

$$\text{OPT}_2^{\text{lifted}}(\xi_1, x_1, x_2) := \arg \min_{\xi_2} L_2(\xi_1, \xi_2) \quad (9b)$$

where the dependence of L_1 and L_2 on ξ_1 and ξ_2 is made explicit through the following relationships:

$$L_1 = q_1^T A q_2, \quad L_2 = q_1^T B q_2 \quad (10a)$$

$$A_{i,j} = f_1(\tau_1^i, \tau_2^j), \quad B_{i,j} = f_2(\tau_1^i, \tau_2^j) \quad (10b)$$

$$\tau_1^i = \text{TRAJ}_1(\xi_1^i, x_1), \quad i \in N_1 \quad (10c)$$

$$\tau_2^j = \text{TRAJ}_2(\xi_2^j, x_2), \quad j \in N_2 \quad (10d)$$

$$(q_1, q_2) = \text{BMG}(A, B). \quad (10e)$$

Here, $N_1 := \{1, \dots, n_1\}$, and $N_2 := \{1, \dots, n_2\}$, where n_1 and n_2 are the number of trajectories for Player 1 and Player 2, respectively. Specifically, each variable τ_1^i represents one of n_1 trajectories that Player 1 optimizes over (and similar for Player 2). The reference variables $\xi_i := (\xi_i^1, \dots, \xi_i^{n_i})$ are now collections of trajectory references, with each ξ_i^j associated to τ_i^j . The function $\text{BMG}(A, B)$ maps cost matrices A and B to mixed equilibrium strategies for the resultant bimatrix game. Specifically, $\text{BMG}(A, B)$ returns a point $(q_1^*, q_2^*) \in \mathbb{R}^{n_1} \times \mathbb{R}^{n_2}$ such that

$$\begin{aligned} (q_1^*)^T A q_2^* &\leq q_1^T A q_2^*, \quad \forall q_1 \in \Delta^{n_1-1}, \\ (q_1^*)^T B q_2^* &\leq (q_1^*)^T B q_2, \quad \forall q_2 \in \Delta^{n_2-1}. \end{aligned} \quad (11)$$

In (11), Δ^k is the k -simplex, representing the space of valid parameters for a categorical distribution over $k+1$ elements. Note that when $n_1 = n_2 = 1$, Game (9) reduces exactly to Game (3), since $\text{BMG}(A, B) \equiv (1, 1)$.

Continuous games, such as (1), may suffer from non-existence of equilibrium points, but when those games are *separable*, a mixed strategy equilibrium is known to exist with finite support [16, 38]. This theoretical result motivates the lifting (9) of our reference-based formulation of Game (1).

For comparative purposes, Game (9) is presented analogously to (3), i.e. without any explicit dependence on reference generators. Nevertheless, generators can be trained analogously to (5), using a similar simultaneous gradient procedure. As before, solving (9) or an analogous version of (5) via gradient play requires that each of the function evaluations in (10) are differentiable in their arguments. It has already been discussed how each of these functions are differentiable, with the exception of the bimatrix game in (10e). We discuss in Appendix B how this function is also differentiable.

A summary of the lifted game solver that utilizes reference generators for reduced online computation is provided in Figure 2. With this computation graph, the cost of approximating solutions to Game (9) is that of evaluating the two generator calls, solving the resultant $n_1 + n_2$ trajectory optimization problems (in parallel, if warranted), and solving a bimatrix game formed by considering all combinations of player trajectories.

C. Extension to Many-Player Games

We reiterate that, although we present this formulation in the two player setting, generalizations to larger games are straightforward. In this case, each player would consider multiple trajectory candidates, and a cost *tensor* would be created for each player, representing the costs for all possible combinations of players' trajectories. A finite Nash equilibrium could be identified over these cost tensors to compute the equilibrium mixing weights q_i [34], and computed by solving a nonlinear, mixed complementarity program [11, 23]. Note that the majority of computation required to construct these cost tensors can be trivially parallelized, making our framework particularly promising for many-player settings. We defer further study of such games to future work.

D. Implementation

We implement the lifted game solver depicted in Figure 2 in the Julia programming language [7]. For the experiments conducted in this work, reference generators π_{θ_i} are realized

as multi-layer perceptrons, trajectory optimization problems TRAJ_i are solved via OSQP [39], and bimatrix games are solved using a custom implementation of the Lemke-Howson algorithm [27].

In order to facilitate back-propagation of gradients through this computation graph, we utilize the auto-differentiation tool Zygote [19]. For those components that cannot be efficiently differentiated automatically, namely TRAJ_i and BMG in Figure 2, we provide custom gradient rules via the implicit function theorem, c.f. [1, 2, 3] and Appendix B. Our implementation can be found at <https://lasse-peters.net/pub/lifted-games>.

V. RESULTS

We have presented a novel formulation of lifted trajectory games in which learned reference generators facilitate the efficient online computation of mixed strategies. In this section, we evaluate the performance of our proposed lifted game solver on variants of the “tag” game shown in Figure 1 and described below in Section V-A. Concretely, we aim to quantify the utility of learning trajectory references rather than choosing them *a priori* (Section V-B), characterize the equilibria identified by trajectory lifting (Section V-C), evaluate the performance of trajectory lifting in head-to-head decentralized competition (Section V-D), and demonstrate our method’s capacity for online training in receding horizon MPPG (Section V-E). Our supplementary material includes a video summarizing these results.

A. Environment: The Tag Game

We validate our methods in a two-player tag game, illustrated in Figure 1. Here, each player’s trajectory τ_i follows time-discretized planar double-integrator dynamics $\ddot{p}_i = u_i$, where $p_i \in \mathbb{R}^2$ is understood to represent horizontal and vertical position in the plane. The set $\mathcal{K}_i(x_i)$ then encompasses all dynamically-feasible trajectories that also satisfy input saturation limits and state constraints. In particular, we require that positions remain within a closed set, such as the pentagon illustrated in Figure 1, and that speeds remain below a fixed magnitude. These choices yield *linear* constraints, so that (4) becomes a quadratic program. We note, however, that our approach does not rely upon this convenient structure and is compatible with more general embedded nonlinear programs.

For the purposes of this example, we shall designate Player 1 to be the “pursuer” and Player 2 to be the “evader.” Hence, the pursuer’s objective $f_1(\tau_1, \tau_2)$ measures the average distance between players’ trajectories over time and is regularized by the difference in control effort between the two players to ensure the existence of at least local pure Nash equilibria for the original game (1). The evader’s objective is $f_2(\tau_1, \tau_2) = -f_1(\tau_1, \tau_2)$. Since the tag game has zero-sum cost structure, throughout the following evaluations we only report the cost for the pursuer and refer to this quantity as the *game value*. Furthermore, unless otherwise stated, we use an input reference signal ξ_i for all players in (3) and (9).

B. The Importance of Learning Trajectory Candidates

Without lifting, it is still possible to approximate mixed strategies for the trajectory game by discretizing the trajectory space (e.g., via sampling [31]). We compare to a sampling-based mixed-strategy baseline to study the isolated effects of learning in a lifted space.

Setup. We instantiate an evader with $n_2 = 20$ pre-sampled trajectory references. To strengthen the evader, we ensure that these samples cover a large region of the trajectory space. To that end, in this experiment (only) we use ξ_i as a reference for Player i ’s goal state rather than their input sequence. We compare the pursuer’s performance for two different schemes of generating trajectory candidates. The non-learning baseline samples $n_1 \in \{1, \dots, 20\}$ pursuer trajectory references from the same distribution as the evader. Our method computes the pursuer strategy by performing gradient play on (9a) to *learn* 2 trajectory candidates via the goal reference parameterization. The mixed Nash equilibrium (q_1, q_2) over the players’ trajectory candidates is computed according to (10). We evaluate both methods for 50 random initial conditions, and record the game value for each trial.

Discussion. Figure 3 summarizes the results of this experiment. As shown, the baseline steadily improves its performance with increasing numbers of sampled trajectory references to mix over. However, even with 20 trajectory samples, it cannot match the performance of our approach with only two learned candidates. Moreover, learning only a few trajectory references drastically reduces the number of trajectory optimizations and, consequently, the size of the bimatrix game in (10).

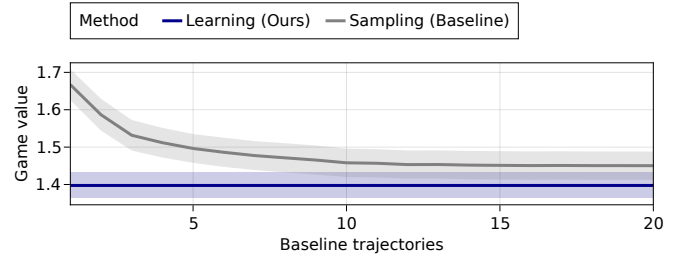


Fig. 3: Comparison of game value for both sampled and learned pursuer trajectories. Lines trace the sample mean over 50 randomized trials, and the surrounding ribbons denote the SEM. On the horizontal axis, we vary the number of sampled *baseline* trajectories n_1 , while fixing *our* approach to learn only 2 trajectories.

C. Convergence and Characteristics of Lifted Equilibria

In this experiment, we analyze *mixed* strategies found by our lifted solver and compare them to *pure* strategies computed by a non-lifting baseline. We shall demonstrate that both approaches reliably converge to different equilibria, and characterize these differences.

Setup. We perform a Monte Carlo study in which we randomly sample 20 initial states of the tag game. On each sample, we invoke two solvers which perform gradient play on different

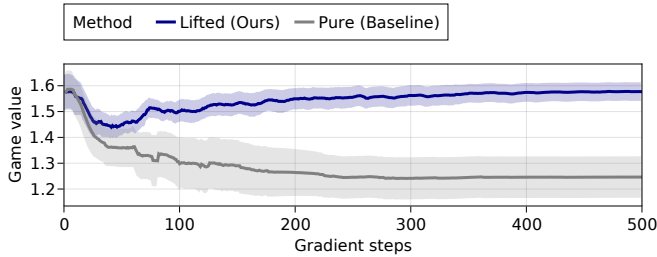


Fig. 4: Equilibrium value convergence averaged over 20 initial states. Ribbons indicate the SEM. The baseline approximates Nash equilibria in *pure* strategies via gradient play on (3). Our method approximates Nash equilibria in *lifted* strategies via gradient play on (9).

strategy spaces. The baseline solver is restricted to pure strategies as in Game (3).⁵ Our method utilizes lifting to find mixed strategies which solve Game (9). In each iteration of gradient play, we record the game value.

Discussion. Figure 4 shows the reliable convergence of both methods in this Monte Carlo study. Since both players learn competitively via simultaneous gradient play, the game value ought not to evolve monotonically; an equilibrium is reached when neither player can improve its strategy unilaterally. At convergence, we observe that the mixed strategies found by our lifting procedure result in a higher game value. This higher value implies that, by operating in a lifted strategy space, the evader can secure a greater average distance between itself and the pursuer.

This gap in value may be understood intuitively by examining the strategy profiles for each method shown in Figure 1. In Figure 1(a), players are restricted to pure strategies, and a rational pursuer can exploit the evader’s deterministic choice of trajectory. In contrast, our proposed lifting formulation allows the evader to mix over multiple trajectory candidates, c.f. Figure 1(b), making its motion less predictable and hence increasing the chance of escaping the pursuer. In response, the pursuer also mixes between two trajectory candidates. However, each of the pursuer’s candidates must account for the full distribution of evader trajectories; hence, the pursuer plans to turn less aggressively than the evader.

In this experiment, we have studied a *centralized* setting in which each method computes strategies for both players from a single game. Therefore, the results presented above are only suitable to characterize the solution points of (3) and (9), but do not justify conclusions about the competitive performance of these solutions in decentralized settings, such as MPGP. In the next section, we extend our analysis to settings in which the opponent’s decision-making process is unknown.

D. Competitive Evaluation Against Non-Lifted Strategies

This experiment is designed to examine the performance of both pure (Baseline) and lifted (Ours) strategies in *de-*

⁵Such pure Nash solutions could also be found using iterated best response [44], iterative linear-quadratic methods [10], or mixed complementarity methods [11].

TABLE I: Open-loop competition.

Pursuer	Evader	
	Lifted	Pure
Lifted	1.577 ± 0.021	1.502 ± 0.022
Pure	1.672 ± 0.022	1.370 ± 0.027

centralized head-to-head competition. For this purpose, we perform two additional Monte Carlo studies which simulate tournaments among players in each strategy class.

Note that, in contrast to previous experiments, here, player strategies are not computed as the solution to a single, centralized game. Rather, each player is oblivious to their opponent’s decision making process and solves its own version of the game from a known initial state over a finite time interval.

1) *Open-Loop Competition:* To begin, we evaluate both methods in open-loop on a fixed, 20-step time horizon.

Setup. For this Monte Carlo study, we randomly sample 100 initial states. For every sampled state, we invoke pure and lifted game solvers twice with randomly sampled initial strategies; once to obtain pursuer strategies, and once to obtain evader strategies.⁶ For all possible solver pairings on these 100 state samples we record the resultant value of the competing strategies; i.e., if Player i chooses trajectory τ_i , we record $f_1(\tau_1, \tau_2)$.

Discussion. Table I summarizes the mean and the standard error of the mean (SEM) of the resultant game value for this open-loop tournament. The evader has a clear incentive to utilize lifted strategies, since they secure the highest game value irrespective of the solution technique used by the pursuer. The best response of the pursuer is then also to play a lifted strategy to minimize value within this column. Hence, (Ours, Ours) is the unique Nash equilibrium in this meta game between solvers.

Additionally, observe that the baseline pursuer performs very well against the baseline evader, as deterministic evasion strategies can always be exploited by a rational pursuer. However, the tournament value reported in the bottom right of Table I is inconsistent with the equilibrium value for the baseline found earlier in Figure 4. This discrepancy suggests that players in this decentralized setup find different local solutions depending on the initialization of the baseline solver. Hence, random initialization effectively makes even a pure strategy evader slightly unpredictable, thereby allowing it to attain a higher average value. By contrast, the value of the lifted strategy computed by our method (top left, Table I) closely agrees with the equilibrium value computed in Figure 4, which indicates that non-uniqueness of solutions is not an issue for our approach.⁷

⁶This initialization procedure avoids leaking information about players’ decision making processes to one another.

⁷This close agreement in value suggests that our method identifies global (rather than local) Nash equilibria which satisfy the so-called *ordered interchangeability property* [5]. Unfortunately, as in continuous optimization, it is generally intractable to properly verify that these solutions are global.

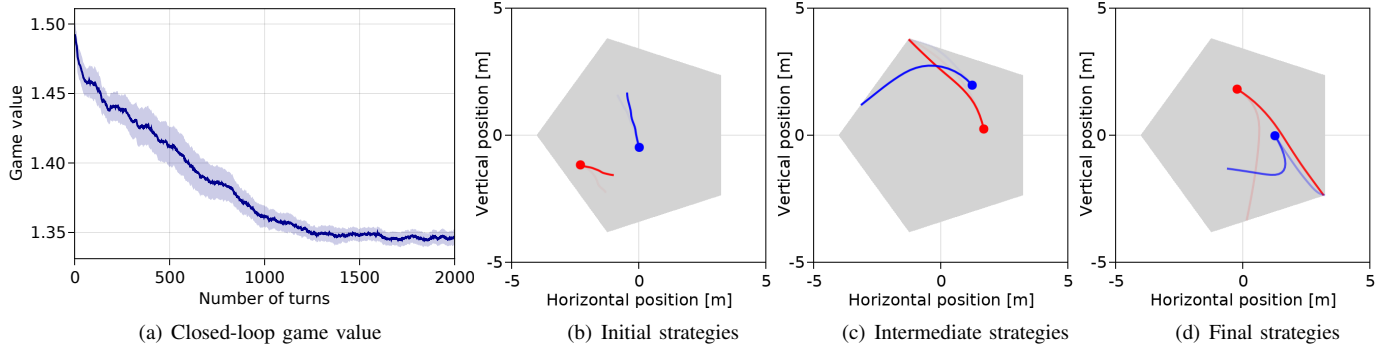


Fig. 5: Training a reference generator for the lifted game solver in Figure 2 in simulated self-play. Transparency denotes the mixing probabilities associated with each trajectory which arise from BMG. (a) Game value mean and SEM over a moving window of 500 turns. (b-d) Strategies at different phases of learning.

TABLE II: Receding-horizon competition.

Pursuer	Evader	
	Lifted	Pure
Lifted	1.360 ± 0.003	1.289 ± 0.005
Pure	1.463 ± 0.004	0.903 ± 0.009

2) *Receding-Horizon Competition*: As MPPG is naturally applied in a receding horizon fashion, we replicate the previous Monte Carlo study in that setting.

Setup. For each of 5 state samples, we simulate receding-horizon competitions for all possible solver pairings. As before, we use a planning horizon of 20 time steps for all players, and in order to simulate latency, we only allow players to update their plans every 9 time steps. Each simulation terminates once players have updated their strategy for the 500th time. From one such trial, we compute the game value by evaluating the pursuer’s objective f_1 on the entire *closed-loop* trajectories of both players. Note that, in contrast to previous experiments, here we use pre-trained reference generators for all solvers as described in Section IV-A to accelerate the computation in this large simulation.

Discussion. Table II summarizes both the mean and SEM for the resultant game value in this receding-horizon Monte Carlo tournament. Overall, we observe the same patterns as in the open-loop setting: lifting is the dominant strategy for the evader and the corresponding best response for the pursuer. However, the game values found for this receding horizon setting are generally lower than in open-loop. By replanning in receding horizon, the pursuer can react to the evader’s decision before the distance between them grows very large.

E. Learning in Receding-Horizon Self-Play

Finally, we demonstrate that a lifted game solver with trajectory generators, as shown in Figure 2, can be rapidly trained from scratch in simulated self-play.

Setup. We repeat the following experiment 10 times. For each player, we randomly sample an initial state x_i and initialize their reference generator π_{θ_i} with parameters θ_i sampled from

a uniform distribution. Subsequently, we simulate receding-horizon learning over 2500 turns with the lifted game solver in the loop. In contrast to the setup used in Section V-D2, here, we do *not* use pre-trained reference generators. Instead, the network parameters are updated on the fly using gradient descent. That is, at every turn, we first perform a forward pass through the computation graph of Figure 2 to compute a lifted strategy profile, followed by a backwards pass to compute a gradient step on each player’s reference generator parameters θ_i . For each experiment, we record players’ strategies as well as the game value over a moving window of 500 turns.

Discussion. Figure 5 summarizes the results of lifted learning in self-play. Initially, the untrained reference generators cause both players to move haphazardly, c.f. Figure 5(b). As learning progresses, players become more competitive, resulting in purposeful, dynamic maneuvers, c.f. Figure 5(c). Within approximately 1500 turns, learning converges, the game value stabilizes, and the solver has learned to generate highly competitive mixed strategies as shown in Figure 5(d).

Note that, throughout the learning procedure, state and input constraints are explicitly enforced in the TRAJ step of the pipeline in Figure 2. Moreover, since our proposed pipeline is end-to-end differentiable, it provides a strong learning signal. Therefore, training in simulated self-play over 2500 turns can be performed in less than three minutes on a standard laptop. Then, once the reference generators π_{θ_i} are fully trained, learning can be disabled, and a forward pass on the pipeline in Figure 2 can be computed with an average run-time of 2 ms. In summary, these results indicate that our method learns quickly and reliably, making it well-suited for online learning in real systems with embedded computational hardware.

VI. CONCLUSION

In this paper, we have proposed two key contributions to the field of noncooperative, multi-agent motion planning. First, we have introduced a principled technique to reduce the online computational complexity of solving these trajectory games. Second, we extended this approach to optimize over a richer, probabilistic class of *lifted* strategies for each

player. Taken together, these innovations facilitate efficiently-computable, high-performance online trajectory planning for multiple autonomous agents in competitive settings. Moreover, our method directly accounts for problem constraints and hence guarantees that learned trajectories satisfy these constraints whenever they are feasible.

While our formulations readily extend to games with many players and arbitrary cost structure, we demonstrate our results in a two-player, zero-sum game of tag. We validate our approach in extensive Monte Carlo studies, in which we observe rapid and reliable convergence to solutions which outperform those which emerge in the original, non-lifted strategy space.

Finally, we showcase our approach in online learning, where each player solves lifted trajectory games in a receding time horizon. Despite the additional complexity present in this setting—e.g., non-stationary training data and potential limit cycles—our method converges reliably to competitive mixed strategies. These initial results are extremely encouraging, and future work should investigate online learning and adaptation in noncooperative settings more extensively. In particular, we note that our method is limited to so-called *open-loop* information structures, in which each agent in a trajectory game must choose future control inputs as a function only of the current state. We believe that the incorporation of feedback structures at the trajectory-level will be an exciting direction for future research.

ACKNOWLEDGMENTS

This work was supported in part by the National Police of the Netherlands. All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors. L. Ferranti received support from the Dutch Science Foundation NWO-TTW within the Veni project HARMONIA (nr. 18165).

REFERENCES

- [1] Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and J Zico Kolter. Differentiable convex optimization layers. *Proc. of the Advances in Neural Information Processing Systems (NIPS)*, 32, 2019.
- [2] Brandon Amos. *Differentiable optimization-based modeling for machine learning*. PhD thesis, Carnegie Mellon University, 2019.
- [3] Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks. In *Proc. of the Int. Conf. on Machine Learning (ICML)*, pages 136–145. PMLR, 2017.
- [4] Tamer Başar and Pierre Bernhard. *H-infinity optimal control and related minimax design problems: a dynamic game approach*. Springer Science & Business Media, 2008.
- [5] Tamer Başar and Geert Jan Olsder. *Dynamic noncooperative game theory*, volume 23. Society for Industrial and Applied Mathematics (SIAM), 1999.
- [6] Dimitri P Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):334–334, 1997.
- [7] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B. Shah. Julia: A fresh approach to numerical computing. *SIAM Review (SIREV)*, 59(1):65–98, 2017.
- [8] Francesco Borrelli, Alberto Bemporad, and Manfred Morari. *Predictive control for linear and hybrid systems*. Cambridge University Press, 2017.
- [9] Efe Camci and Erdal Kayacan. Learning motion primitives for planning swift maneuvers of quadrotor. *Autonomous Robots*, 43(7):1733–1745, 2019.
- [10] Bolei Di and Andrew Lamperski. Newton’s method and differential dynamic programming for unconstrained nonlinear dynamic games. In *Proceedings of the Conference on Decision Making and Control (CDC)*. IEEE, 2019.
- [11] Steven P Dirkse and Michael C Ferris. The path solver: a nonmonotone stabilization scheme for mixed complementarity problems. *Optimization methods and software*, 5(2):123–156, 1995.
- [12] Francisco Facchinei and Christian Kanzow. Generalized nash equilibrium problems. *Annals of Operations Research*, 175(1):177–211, 2010.
- [13] Francisco Facchinei, Andreas Fischer, and Veronica Piccialli. Generalized nash equilibrium problems and newton methods. *Mathematical Programming*, 117(1-2):163–194, 2009.
- [14] Jakob N Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and Igor Mordatch. Learning with opponent-learning awareness. *arXiv preprint arXiv:1709.04326*, 2017.
- [15] David Fridovich-Keil, Ellis Ratner, Lasse Peters, Anca D. Dragan, and Claire J. Tomlin. Efficient iterative linear-quadratic approximations for nonlinear multi-player general-sum differential games. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*. IEEE, 2020.
- [16] I Glicksberg and Oliver Gross. 9. notes on games over the square. In *Contributions to the Theory of Games (AM-28), Volume II*, pages 173–182. Princeton University Press, 2016.
- [17] I.J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative Adversarial Networks. *arXiv preprint*, 2014. URL <https://arxiv.org/pdf/1406.2661>.
- [18] Michael Green and David JN Limebeer. *Linear robust control*. Courier Corporation, 2012.
- [19] Michael Innes. Don’t unroll adjoint: Differentiating ssa-form programs. *arXiv preprint arXiv:1810.07951*, 2018.
- [20] Rufus Isaacs. Differential games i-iv. Technical report, RAND CORP SANTA MONICA CA SANTA MONICA, 1954-1955.
- [21] D.H. Jacobson and D.Q. Mayne. *Differential Dynamic Programming*. Modern analytic and computational methods in science and mathematics. American Elsevier Publishing Company, 1970.
- [22] Oussama Khatib, Sean Quinlan, and David Williams. Robot planning and control. *Journal on Robotics and*

- Autonomous Systems (RAS)*, 21(3):249–261, 1997.
- [23] Forrest Laine. TensorGames, 2022. URL <https://github.com/4estlaine/TensorGames.jl>.
 - [24] Forrest Laine, David Fridovich-Keil, Chih-Yuan Chiu, and Claire Tomlin. The computation of approximate generalized feedback nash equilibria. *arXiv preprint arXiv:2101.02900*, 2021.
 - [25] Forrest Laine, David Fridovich-Keil, Chih-Yuan Chiu, and Claire Tomlin. Multi-hypothesis interactions in game-theoretic motion planning. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, pages 8016–8023. IEEE, 2021.
 - [26] Simon Le Cleac’h, Mac Schwager, and Zachary Manchester. ALGAMES: A fast solver for constrained dynamic games. In *Proc. of Robotics: Science and Systems (RSS)*, 2020.
 - [27] Carlton E Lemke and Joseph T Howson, Jr. Equilibrium points of bimatrix games. *Society for Industrial and Applied Mathematics (SIAM)*, 12(2):413–423, 1964.
 - [28] Weiwei Li and Emanuel Todorov. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *ICINCO*, pages 222–229. Citeseer, 2004.
 - [29] Weiwei Li and Emanuel Todorov. Iterative linearization methods for approximately optimal control and estimation of non-linear stochastic system. *International Journal of Control*, 80(9):1439–1453, 2007.
 - [30] Daniel Liberzon. *Calculus of variations and optimal control theory*. Princeton University Press, 2011.
 - [31] Alexander Liniger and John Lygeros. A noncooperative game approach to autonomous racing. *IEEE Transactions on Control Systems Technology*, 28(3):884–897, 2019.
 - [32] Katta G Murty and Feng-Tien Yu. *Linear complementarity, linear and nonlinear programming*, volume 3. Citeseer, 1988.
 - [33] Jorge Nocedal and Stephen Wright. *Numerical optimization*. Springer Verlag, 2006.
 - [34] Christos H Papadimitriou and Tim Roughgarden. Computing equilibria in multi-player games. In *SODA*, volume 5, pages 82–91. Citeseer, 2005.
 - [35] Daniel Ralph and Stephan Dempe. Directional derivatives of the solution of a parametric nonlinear program. *Mathematical programming*, 70(1):159–172, 1995.
 - [36] Alan Wilbor Starr and Yu-Chi Ho. Further properties of nonzero-sum differential games. *Journal of Optimization Theory and Applications*, 3(4):207–219, 1969.
 - [37] Alan Wilbor Starr and Yu-Chi Ho. Nonzero-sum differential games. *Journal of optimization theory and applications*, 3(3):184–206, 1969.
 - [38] Noah D Stein, Asuman Ozdaglar, and Pablo A Parrilo. Separable and low-rank continuous games. *International Journal of Game Theory*, 37(4):475–504, 2008.
 - [39] Bartolomeo Stellato, Goran Banjac, Paul Goulart, Alberto Bemporad, and Stephen Boyd. OSQP: An operator splitting solver for quadratic programs. *Mathematical Programming Computation*, 12(4):637–672, 2020.
 - [40] Freek Stulp, Evangelos Theodorou, Mrinal Kalakrishnan, Peter Pastor, Ludovic Righetti, and Stefan Schaal. Learning motion primitive goals for robust manipulation. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 325–331. IEEE, 2011.
 - [41] Yuval Tassa, Nicolas Mansard, and Emo Todorov. Control-limited differential dynamic programming. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, pages 1168–1175, 2014.
 - [42] Emanuel Todorov and Weiwei Li. A generalized iterative LQG method for locally-optimal feedback control of constrained nonlinear stochastic systems. In *Proc. of the IEEE American Control Conference (ACC)*, pages 300–306. IEEE, 2005.
 - [43] Mingyu Wang, Zijian Wang, John Talbot, J Christian Gerdes, and Mac Schwager. Game-theoretic planning for self-driving cars in multivehicle competitive scenarios. *IEEE Trans. on Robotics (TRO)*, 2021.
 - [44] Zijian Wang, Riccardo Spica, and Mac Schwager. Game theoretic motion planning for multi-robot racing. *Distributed Autonomous Robotic Systems*, pages 225–238, 2019.
 - [45] Zijian Wang, Riccardo Spica, and Mac Schwager. Game theoretic motion planning for multi-robot racing. In *Distributed Autonomous Robotic Systems*, pages 225–238. Springer, 2019.
 - [46] Zhaoming Xie, C Karen Liu, and Kris Hauser. Differential dynamic programming with nonlinear constraints. In *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, pages 695–702, 2017.

APPENDIX A

EQUIVALENCE OF GAME (1) AND GAME (3)

In this section we establish an equivalence result between (1) and (3). We prove this for a particular interpretation of the reference variables ξ_i , and forms of $G_i, H_i, \mathcal{K}_i(x_i)$, noting that similar results can be established for other settings.

For each Player i , consider the instance of (4) in which $G_i := I$ and $H_i := 0$, representing the identity and zero matrices of appropriate dimension. Furthermore, assume that the set $\mathcal{K}_i(x_i) := \{\tau : lb_i \leq g_i(\tau) \leq ub_i\}$, for some vector-valued and twice-differentiable function g_i , and lower and upper bounds lb_i and ub_i . It is assumed that a suitable constraint qualification applies to this constraint set, such as the LICQ [33]. This implies that ξ_i has dimension equal to that of the decision variable τ , and the objective of (4) is to find a trajectory $\tau \in \mathcal{K}_i(x_i)$ which is as close as possible to ξ_i as measured by the ℓ_2 -norm.

Theorem 1. *In the setting as stated above,*

- 1) *For any stationary point (τ_1, τ_2) of Game (1), there exists a stationary point (ξ_1, ξ_2) of Game (3) such that $\tau_i = \text{TRAJ}_i(\xi_i, x_i)$ for all players i .*
- 2) *For any stationary point (ξ_1, ξ_2) of Game (3) satisfying $\xi_i \in \mathcal{K}_i(x_i)$, the trajectories $\tau_i = \text{TRAJ}_i(\xi_i, x_i)$ constitute a stationary point for (1).*

Proof: To prove this result, we first make explicit the definition of a stationary point for (1) and (3). A stationary point for (1) is a point (τ_1, τ_2) , such that for both players i ,

$$d^\top \nabla_{\tau_i} f_i(\tau_1, \tau_2) \geq 0, \forall d \in T_{\mathcal{K}_i}(\tau_i). \quad (12)$$

Here, $T_{\mathcal{K}_i}(\tau_i)$ is the set of linearized feasible directions with respect to constraint set $\mathcal{K}_i(x_i)$ at τ_i , which because we have assumed a suitable constraint qualification, is equivalent to the tangent cone at this point [33]. Specifically, at a feasible point τ , let $\mathcal{I}_l(\tau) := \{j : g_{i,j}(\tau) = lb_j\}$, and $\mathcal{I}_u(\tau) := \{j : g_{i,j}(\tau) = ub_j\}$. Then

$$T_{\mathcal{K}_i}(\tau) := \{d : d^\top \nabla g_{i,j}(\tau) \geq 0, j \in \mathcal{I}_l(\tau), \\ d^\top \nabla g_{i,j}(\tau) \leq 0, j \in \mathcal{I}_u(\tau)\} \quad (13)$$

A stationary point for (3) is a point (ξ_1, ξ_2) such that

$$(\nabla_{\xi_i} \tau_i \cdot d)^\top \nabla_{\tau_i} f_i(\tau_1, \tau_2) \geq 0, \forall d, \quad (14)$$

where $\tau_i = \text{TRAJ}_i(\xi_i, x_i)$. Note that $(\nabla_{\xi_i} \tau_i \cdot d) := (\nabla_{\xi_i} \text{TRAJ}_i(\xi_i, x_i) d)$ as appearing above is the directional derivative of $\text{TRAJ}_i(\xi_i, x_i)$ with respect to changes of ξ in the direction d . This directional derivative is defined to be e , where e solves the following quadratic program [35]:

$$\begin{aligned} \min_e \quad & \frac{1}{2} e^\top Q_1 e + d^\top Q_2 e \\ \text{s.t.} \quad & e \in \mathcal{C}_{\lambda_i}(\tau_i), \end{aligned} \quad (15)$$

where $Q_1 := I - \nabla_{\tau, \tau}^2 (g(\tau_i)^\top \lambda_i)$, $Q_2 := -I$, λ_i are the dual variables associated with the primal solution τ_i to $\text{TRAJ}_i(\xi_i, x_i)$, and $\mathcal{C}_{\lambda_i}(\tau)$ is the critical cone to the constraint set $g_i(\tau)$ with respect to λ_i at τ :

$$\mathcal{C}_{\lambda_i}(\tau) := \{d : \bar{lb}_{i,j} \leq d^\top \nabla g_{i,j}(\tau) \leq \bar{ub}_{i,j}\}. \quad (16)$$

The bounds $\bar{lb}_{i,j}$ and $\bar{ub}_{i,j}$ are defined as:

$$(\bar{lb}_{i,j}, \bar{ub}_{i,j}) := \begin{cases} (0, 0) & j \in \mathcal{I}_l(\tau) \ \& \ \lambda_{i,j} > 0 \\ (0, \infty) & j \in \mathcal{I}_l(\tau) \ \& \ \lambda_{i,j} = 0 \\ (-\infty, \infty) & i \notin (\mathcal{I}_l(\tau) \cup \mathcal{I}_u(\tau)) \\ (-\infty, 0) & j \in \mathcal{I}_u(\tau) \ \& \ \lambda_{i,j} > 0 \\ (0, 0) & j \in \mathcal{I}_u(\tau) \ \& \ \lambda_{i,j} < 0 \end{cases} \quad (17)$$

Now, to prove 1), we show that $\xi_i = \tau_i$ satisfies the claim. It follows directly that $\tau_i = \text{TRAJ}_i(\xi_i, x_i)$. It can be verified that, because $\tau_i \in \mathcal{K}_i(x_i)$ by definition, then all constraints appearing in TRAJ_i are only weakly active, implying $\lambda_i = 0$. This implies that the constraint set appearing in (15) is precisely the tangent cone (13). Therefore, for all directions d , $\nabla_{\xi_i} \text{TRAJ}_i(\tau_i, x_i) \cdot d \in T_{\mathcal{K}_i}(\tau_i)$, which by (12), implies that (14) holds, establishing the result.

To prove 2), we simply note that if $\xi_i \in \mathcal{K}_i(x_i)$, then $\tau_i = \xi_i$. Furthermore, in this setting $\lambda_i = 0$ as before, and therefore the critical cone appearing in (15) is again equivalent to the tangent cone (13). This implies that the directional derivative $(\nabla_{\xi_i} \tau_i \cdot d)$ is defined to simply be the projection of the direction d into the tangent cone at τ_i . The set of all directions d mapped

through this projection results precisely in $T_{\mathcal{K}_i}(\tau_i)$. Therefore, the conditions (14) imply (12) for this setting, implying our result. ■

The result as stated in Theorem 1 does not imply that an *arbitrary* stationary point found for (3) corresponds to a stationary point for (1), since it may be that either of the references $\xi_i \notin \mathcal{K}_i(x_i)$. For such reference points, it is possible that for some direction d the expression in (14) holds with equality, yet the expression in (12) is violated. This situation results in “sticky constraints,” in which a descent direction exists for $f_i(\tau_1, \tau_2)$, yet that direction is not in the range of $\nabla_{\xi_i} \text{TRAJ}_i(\xi_i, x_i)$, i.e. small changes to the reference are not enough to release τ_i away from the active constraint boundaries.

To address this issue, we propose a modest regularization scheme to eliminate the possibility of reference stationary points of (3) which do not correspond to trajectory stationary points of (1). One such approach could be to enforce constraints in (3) such that $\xi_i \in \mathcal{K}_i(x_i)$. This, however, would render the reformulation from (1) to (3) pointless. Instead, we impose a simple regularization in the objectives of each player in (3). Namely, instead of minimizing over $f_i(\tau_1, \tau_2)$ w.r.t. ξ_i , we minimize over

$$f_i(\tau_1, \tau_2) + \|(g(\xi_i) - ub)_+ + (lb - g(\xi_i))_+\|_2^2, \quad (18)$$

where $(\cdot)_+ := \max(\cdot, 0)$.

Note that this introduced regularization is exact, and has precisely the effect of eliminating any stationary points for (3) in which $\xi_i \notin \mathcal{K}_i(x_i)$. If the regularization term is non-zero, then necessarily from the definition of the directional derivative (15), the gradient of the regularization component is in the null-space of $\nabla_{\xi_i} \text{TRAJ}_i(\xi_i, x_i)$. This implies the regularization can be driven to zero without changing the resultant solution τ_i . This is true irrespective of the scale factor multiplying the regularization term. Furthermore, if $\xi_i \in \mathcal{K}_i(x_i)$, then the regularization term is zero, and has no effect on stationary points of the un-regularized game (3).

We note that the particular choice of regularization (18) is only applicable for the interpretation of the references ξ_i made throughout this section. For more general parameterizations of the reference, as discussed in the main text, a suitable regularization is the norm of inequality constraint multipliers associated with the solution of $\text{TRAJ}_i(\xi_i, x_i)$. The use of this dual-variable regularization is effective at eliminating the spurious stationary points for (3), so long as the parameterization of the reference is rich enough such that for any ξ_i and associated τ_i, λ_i , there exists directions d in which the ξ_i can be perturbed and the directional derivative of τ_i is 0, and the directional derivative of $\lambda_{i,j}$ is negative for all j . This is true, for example, of the control signal reference used throughout this work.

Therefore, with use of the introduced regularization (18), the stationary points of Games (1) and (3) have a one-to-one correspondence, warranting the use of Game (3) in place of Game (1).

APPENDIX B
DIFFERENTIATING THROUGH BMG

The problem of finding q_1, q_2 which satisfy (11) (as is the task of the function BMG), can be equivalently expressed as the linear complementarity problem [32]

$$\begin{aligned} & \text{find } p_1, p_2 \\ & \text{s.t. } p_1 \geq 0 \perp \bar{A}p_2 \geq 1 \\ & \quad p_2 \geq 0 \perp \bar{B}^\top p_1 \geq 1. \end{aligned} \quad (19)$$

The solution (q_1, q_2) to the BMG are related to the solution to (19) by the relations

$$(q_1)_i = \frac{(p_1)_i}{\sum_k (p_1)_k}, \quad (q_2)_i = \frac{(p_2)_i}{\sum_k (p_2)_k}. \quad (20)$$

It is assumed that \bar{A} and \bar{B} are derived from the original matrices A, B , as the following. $\bar{A}_{i,j} := A_{i,j} + \alpha$, $\bar{B}_{i,j} := B_{i,j} + \beta$, for some positive constants α, β such that every element of \bar{A} and \bar{B} are strictly positive. Furthermore, the 1s appearing in the right-hand side of the constraints in (19) are assumed to represent vectors of appropriate dimension with each value equal to 1.

Consider some solution p_1, p_2 to (19) in which strict complementarity holds for each condition, e.g. either $p_{1,j} = 0$ or $(\bar{A}p_2)_j = 1$, but not both. For each $j \in \{1, 2\}$, Denote the index sets $\mathcal{I}_j^+ := \{i : (p_j)_i > 0\}$. Then let

$$\begin{aligned} p_1^+ &:= [p_1]_{\mathcal{I}_1^+}, & p_2^+ &:= [p_2]_{\mathcal{I}_2^+}, \\ \bar{A}^+ &:= [\bar{A}]_{\mathcal{I}_1^+, \mathcal{I}_2^+}, & \bar{B}^+ &:= [\bar{B}]_{\mathcal{I}_1^+, \mathcal{I}_2^+}. \end{aligned}$$

In words, p_1^+ is the vector formed by only considering the non-zero elements of p_1 , and \bar{B}^+ is the matrix formed by considering the columns specified by \mathcal{I}_1^+ and rows specified by \mathcal{I}_2^+ . By the strict complementarity, at equilibrium, it is that

$$\begin{bmatrix} 0 & \bar{A}^+ \\ (\bar{B}^+)^\top & 0 \end{bmatrix} \begin{bmatrix} p_1^+ \\ p_2^+ \end{bmatrix} = 1, \quad (21)$$

where, as before, the right-hand side 1 is a vector consisting of all 1s.

The values $p_j^- := [p_j]_i, i \notin \mathcal{I}_j^+$ are defined to be identically 0, and as such have 0 derivative with respect to the values \bar{A}, \bar{B} . The derivatives of remaining portion of the solution, p_1^+, p_2^+ , can be evaluated from (21). If the matrix on the left-hand-side of (21) is singular, then the resulting solution is in fact non-isolated (there exist a continuum of solutions satisfying (19)), and the derivatives of the solution are not defined. If the matrix is non-singular, then necessarily so are \bar{A}^+ and \bar{B}^+ , and the isolated solutions of p_1^+, p_2^+ are locally related to the matrices \bar{A}, \bar{B} as

$$\begin{aligned} p_1^+ &= (\bar{B}^+)^{-\top} 1, \\ p_2^+ &= (\bar{A}^+)^{-1} 1. \end{aligned} \quad (22)$$

In this form, the derivatives of each element of p_j^+ can be found by differentiating through the expressions (22). Combining the above results, the derivatives of the solution

vector p_1, p_2 with respect to the problem data A, B , can be established as the following:

$$\begin{aligned} \frac{\partial (p_1)_i}{\partial A_{j,k}} &:= 0, \\ \frac{\partial (p_1)_i}{\partial B_{j,k}} &:= \begin{cases} 0 & : i \notin \mathcal{I}_1^+ \\ -((\bar{B}^+)^{-\top} I_{k,j} p_1^+)_i & : \text{else} \end{cases}, \\ \frac{\partial (p_2)_i}{\partial A_{j,k}} &:= \begin{cases} 0 & : i \notin \mathcal{I}_2^+ \\ -((\bar{A}^+)^{-1} I_{j,k} p_2^+)_i & : \text{else} \end{cases}, \\ \frac{\partial (p_2)_i}{\partial B_{j,k}} &:= 0. \end{aligned} \quad (23)$$

Above, the term $I_{j,k}$ is used to refer to the matrix consisting of zero everywhere except at the (j, k) -th position, which has value 1.

When strict complementarity does not hold at the solution to (19), then only directional derivatives of the solution vectors exist w.r.t. the problem data. The various directional derivatives are found by, for each condition which does not hold with strict complementarity, making a selection on whether that index should be included the sets \mathcal{I}_j^+ or not. Then proceeding with the remainder of calculations, the result forms one of the directional derivative for the system. The directions for which this derivative is valid are defined to be those which make the directional derivative consistent with the selected index sets.

The derivatives of the elements of p_1 and p_2 with respect to the cost matrices are formed by propagating the derivatives (23) through the relationships (20).