# CORRELATION IN EXTENSIVE-FORM GAMES: SADDLE-POINT FORMULATION AND BENCHMARKS\*

#### ARXIV PREPRINT

#### Gabriele Farina

Computer Science Department Carnegie Mellon University gfarina@cs.cmu.edu

#### **Chun Kai Ling**

Computer Science Department Carnegie Mellon University chunkail@cs.cmu.edu

#### Fei Fang

Institute for Software Research Carnegie Mellon University feif@cs.cmu.edu

#### **Tuomas Sandholm**

Computer Science Department, CMU
Strategic Machine, Inc.
Strategy Robot, Inc.
Optimized Markets, Inc.
sandholm@cs.cmu.edu

October 27, 2019

#### **Abstract**

While Nash equilibrium in extensive-form games is well understood, very little is known about the properties of *extensive-form correlated equilibrium (EFCE)*, both from a behavioral and from a computational point of view. In this setting, the strategic behavior of players is complemented by an external device that privately recommends moves to agents as the game progresses; players are free to deviate at any time, but will then not receive future recommendations. Our contributions are threefold. First, we show that an EFCE can be formulated as the solution to a bilinear saddle-point problem. To showcase how this novel formulation can inspire new algorithms to compute EFCEs, we propose a simple subgradient descent method which exploits this formulation and structural properties of EFCEs. Our method has better scalability than the prior approach based on linear programming. Second, we propose two benchmark games, which we hope will serve as the basis for future evaluation of EFCE solvers. These games were chosen so as to cover two natural application domains for EFCE: conflict resolution via a mediator, and bargaining and negotiation. Third, we document the qualitative behavior of EFCE in our proposed games. We show that the social-welfare-maximizing equilibria in these games are highly nontrivial and exhibit surprisingly subtle *sequential* behavior that so far has not received attention in the literature.

## 1. Introduction

Nash equilibrium (NE) (Nash, 1950), the most seminal concept in non-cooperative game theory, captures a multi-agent setting where each agent is selfishly motivated to maximize their own payoff. The assumption underpinning NE is that the interaction is completely *decentralized*: the behavior of each agent is not regulated by any external orchestrator. Contrasted with the other—often utopian—extreme of a fully managed interaction, where an external dictator controls the behavior of each agent so that the whole system moves to a desired state, the social welfare that can be achieved by NE is generally lower, sometimes dramatically so (Koutsoupias & Papadimitriou, 1999; Roughgarden & Tardos, 2002). Yet, in many

<sup>\*</sup>This paper was accepted for publication at NeurIPS 2019.

realistic interactions, some intermediate form of centralized control can be achieved. In particular, in his landmark paper, Aumann (1974) proposed the concept of *correlated equilibrium* (CE), where a mediator (the *correlation device*) can *recommend* behavior, but not *enforce it*. In a CE, the correlation device is constructed so that the agents—which are still modeled as fully rational and selfish just like in an NE—have no incentive to deviate from the private recommendation. Allowing correlation of actions while ensuring selfishness makes CE a good candidate solution concept in multi-agent and semi-competitive settings such as traffic control, load balancing (Ashlagi et al., 2008), and carbon abatement (Ray & Gupta, 2009), and it can lead to win-win outcomes.

In this paper, we study the natural extension of correlated equilibrium in *extensive-form* (i.e., sequential) games, known as extensive-form correlated equilibrium (EFCE) (von Stengel & Forges, 2008). Like CE, EFCE assumes that the strategic interaction is complemented by an external mediator; however, in an EFCE the mediator only privately reveals the recommended next move to each *acting player*, instead of revealing the whole plan of action throughout the game (i.e., recommended move at *all* decision points) for each player at the beginning of the game. Furthermore, while each agent is free to defect from the recommendation at any time, this comes at the cost of future recommendations.

While the properties of correlation in *normal-form* games are well-studied, they do not automatically transfer to the richer world of sequential interactions. It is known in the study of NE that sequential interactions can pose different challenges, especially in settings where the agents retain private information. Conceptually, the players can strategically adjust to dynamic observations about the environment and their opponents as the game progresses. Despite tremendous interest and progress in recent years for computing NE in sequential interactions with private information, with significant milestones achieved in poker games (Bowling et al., 2015; Brown & Sandholm, 2017; Moravčík et al., 2017; Brown & Sandholm, 2019b) and other large, real-world domains, not much has been done to increase our understanding of (extensive-form) correlated equilibria in these settings.

**Contributions** Our primary objective with this paper is to spark more interest in the community towards a deeper understanding of the behavioral and computational aspects of EFCE.

- In Section 3 we show that an EFCE in a two-player general-sum game is the solution to a bilinear saddle-point problem (BSPP). This conceptual reformulation complements the EFCE construction by von Stengel & Forges (2008), and allows for the development of new and efficient algorithms. As a proof of concept, by using our reformulation we devise a variant of projected subgradient descent which outperforms linear-programming(LP)-based algorithms proposed by von Stengel & Forges (2008) in large game instances.
- In Section 5 we propose two benchmark games; each game is parametric, so that these games can scale in size as desired. The first game is a general-sum variant of the classic war game *Battleship*. The second game is a simplified version of the *Sheriff of Nottingham* board game. These games were chosen so as to cover two natural application domains for EFCE: conflict resolution via a mediator, and bargaining and negotiation.
- By analyzing EFCE in our proposed benchmark games, we show that even if the mediator cannot enforce behavior, it can induce significantly higher social welfare than NE and successfully deter players from deviating in at least two (often connected) ways: (1) using certain sequences of actions as 'passcodes' to verify that a player has not deviated: defecting leads to incomplete or wrong passcodes which indicate deviation, and (2) inducing opponents to play punitive actions against players that have deviated from the recommendation, if such a deviation is detected. Crucially, both deterrents are unique to sequential interactions and do not apply to non-sequential games. This corroborates the idea that the mediation of sequential interactions is a qualitatively different problem than that of non-sequential games and further justifies the study of EFCE as an interesting direction for the community. To our knowledge, these are the first experimental results and observations on EFCE in the literature.

The source code for our game generators and subgradient method is published online<sup>2</sup>.

<sup>&</sup>lt;sup>2</sup>https://github.com/Sandholm-Lab/game-generators https://github.com/Sandholm-Lab/efce-subgradient

# 2. Preliminaries

Extensive-form games (EFGs) are sequential games that are played over a rooted game tree. Each node in the tree belongs to a player and corresponds to a decision point for that player. Outgoing edges from a node v correspond to actions that can be taken by the player to which v belongs. Each terminal node in the game tree is associated with a tuple of payoffs that the players receive should the game end in that state. To capture imperfect information, the set of vertices of each player is partitioned into *information sets*. The vertices in a same information set are indistinguishable to the player that owns those vertices. For example, in a game of Poker, a player cannot distinguish between certain states that only differ in opponent's private hand. As a result, the strategy of the player (specifying which action to take) is defined on the information sets instead of the vertices. For the purpose of this paper, we only consider *perfect-recall* EFGs. This property means that each player does not forget any of their previous action, nor any private or public observation that the player has made. The perfect-recall property can be formalized by requiring that for any two vertices in a same information set, the paths from those vertices to the root of the game tree contain the exact same sequence of actions for the acting player at the information set.

A pure normal-form strategy for Player i defines a choice of action for every information set that belongs to i. A player can play a mixed strategy, i.e., sample from a distribution over their pure normal-form strategies. However, this representation contains redundancies: some information sets for Player i may become unreachable reachable after the player makes certain decisions higher up in the tree. Omitting these redundancies leads to the notion of reduced-normal-form strategies, which are known to be strategically equivalent to normal-form strategies (e.g., (Shoham & Leyton-Brown, 2009) for more details). Both the normal-form and the reduced-normal-form representation are exponentially large in the size of the game tree.

Here, we fix some notations. Let Z be the set of terminal states (or equivalently, outcomes) in the game and  $u_i(z)$  be the utility obtained by player i if the game terminates at  $z \in Z$ . Let  $\Pi_i$  be the set of pure reduced-normal-form strategies for Player i. We define  $\Pi_i(I)$ ,  $\Pi_i(I,a)$  and  $\Pi_i(z)$  to be the set of reduced-normal-form strategies that (a) can lead to information set I, (b) can lead to I and prescribes action a at information set I, and (c) can lead to the terminal state z, respectively. We denote by  $\Sigma_i$  the set of information set-action pairs (I,a) (also referred to as sequences), where I is an information set for Player i and a is an action at set I. For a given terminal state z let  $\sigma_i(z)$  be the last (I,a) pair belonging to Player i encountered in the path from the root of the tree to z.

Extensive-Form Correlated Equilibrium Extensive-form correlated equilibrium (EFCE) is a solution concept for extensive-form games introduced by von Stengel & Forges (2008). Like in the traditional correlated equilibrium (CE), introduced by Aumann (1974), a *correlation device* selects private signals for the players before the game starts. These signals are sampled from a correlated distribution  $\mu$ —a joint probability distribution over  $\Pi_1 \times \Pi_2$ —and represent recommended player strategies. However, while in a CE the recommended moves for the whole game tree are privately revealed to the players when the game starts, in an EFCE the recommendations are revealed incrementally as the players progress in the game tree. In particular, a recommended move is only revealed when the player reaches the decision point in the game for which the recommendation is relevant. Moreover, if a player ever deviates from the recommended move, they will stop receiving recommendations. To concretely implement an EFCE, one places recommendations into 'sealed envelopes' which may only be opened at its respective information set. Sealed envelopes may implemented using cryptographic techniques (see (Dodis et al., 2000) for one such example).

In an EFCE, the players know less about the set of recommendations that were sampled by the correlation device. The benefits are twofold. First, the players can be more easily induced to play strategies that hurt them (but benefit the overall social welfare), as long as "on average" the players are indifferent as to whether or not to follow the recommendations: the set of EFCEs is a *superset* of that of CEs. Second, since the players observe less, the set of probability distributions for the correlation device for which no

<sup>&</sup>lt;sup>3</sup>Other CE-related solution concepts in sequential games include the agent-form correlated equilibrium (AFCE), where agents continue to receive recommendations even upon defection, and normal-form coarse CE (NFCCE). NFCCE does not allow for defections during the game, in fact, before the game starts, players must decide to commit to following *all* recommendations upfront (before receiving them), or elect to receive none.

player has an incentive to deviate can be described succinctly in certain classes of games: von Stengel & Forges (2008, Theorem 1.1) show that in two-player, perfect-recall extensive-form games with no chance moves, the set of EFCEs can be described by a system of linear equations and inequalities of polynomial size in the game description. On the other hand, the same result cannot hold in more general settings: von Stengel & Forges (2008, Section 3.7) also show that in games with more than two players and/or chance moves, deciding the existence of an EFCE with social welfare greater than a given value is NP-hard. It is important to note that this last result only implies that the characterization of the set of *all* EFCEs cannot be of polynomial size in general (unless P = NP). However, the problem of finding *one* EFCE can be solved in polynomial time: Huang (2011) and Huang & von Stengel (2008) show how to adapt the *Ellipsoid Against Hope* algorithm (Papadimitriou & Roughgarden, 2008; Jiang & Leyton-Brown, 2015) to compute an EFCE in polynomial time in games with more than two players and/or with chance moves. Unfortunately, that algorithm is only theoretical, and known to not scale beyond extremely small instances (Leyton-Brown, 2019).

# 3. Extensive-Form Correlated Equilibria as Bilinear Saddle-Point Problems

Our objective for this section is to cast the problem of finding an EFCE in a two-player game as a bilinear saddle-point problem, that is a problem of the form  $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} x^T Ay$ , where  $\mathcal{X}$  and  $\mathcal{Y}$  are compact convex sets. In the case of EFCE,  $\mathcal{X}$  and  $\mathcal{Y}$  are convex polytopes that belong to a space whose dimension is polynomial in the game tree size. This reformulation is meaningful:

- From a conceptual angle, it brings the problem of computing an EFCE closer to several other solution concepts in game theory that are known to be expressible as BSPP. In particular, the BSPP formulation shows that an EFCE can be viewed as a NE in a two-player zero-sum game between a *deviator*, who is trying to decide how to best defect from recommendations, and a *mediator*, who is trying to come up with an incentive-compatible set of recommendations.
- From a geometric point of view, the BSPP formulation better captures the combinatorial structure of the problem:  $\mathcal{X}$  and  $\mathcal{Y}$  have a well-defined meaning in terms of the input game tree. This has algorithmic implications: for example, because of the structure of  $\mathcal{Y}$  (which will be detailed later), the inner maximization problem can be solved via a single bottom-up game-tree traversal.
- From a computational standpoint, it opens the way to the plethora of optimization algorithms (both general-purpose and those specific to game theory) that have been developed to solve BSPPs. Examples include Nesterov's excessive gap technique (Nesterov, 2005), Nemirovski's mirror prox algorithm (Nemirovski, 2004) and regret-methods based methods such as mirror descent, follow-the-regularized-leader (e.g., (Hazan, 2016)), and CFR and its variants (Zinkevich et al., 2007; Farina et al., 2019; Brown & Sandholm, 2019a).

Furthermore, it is easy to show that by dualizing the inner maximization problem in the BSPP formulation, one recovers the linear program introduced by von Stengel & Forges (2008) (we show this in Appendix A). In this sense, our formulation subsumes the existing one.

**Triggers and Deviations** One effective way to reason about extensive-form correlated equilibria is via the notion of *trigger agents*, which was introduced (albeit used in a different context) in Gordon et al. (2008) and Dudik & Gordon (2009):

**Definition 1.** Let  $\hat{\sigma} := (\hat{I}, \hat{a}) \in \Sigma_i$  be a sequence for Player i, and let  $\hat{\mu}$  be a distribution over  $\Pi_i(\hat{I})$ . A  $(\hat{\sigma}, \hat{\mu})$ -trigger agent for Player i is a player that follows all recommendations given by the mediator unless they get recommended  $\hat{a}$  at  $\hat{I}$ ; in that case, the player 'gets triggered', stops following the recommendations and instead plays based on a pure strategy sampled from  $\hat{\mu}$  until the game ends.

A correlated distribution  $\mu$  is an EFCE if and only if any trigger agent for Player i can get utility at most equal to the utility that Player i earns by following the recommendations of the mediator at all decision points. In order to express the utility of the trigger agent, it is necessary to compute the probability of the game ending in each of the terminal states. As we show in Appendix B, this can be done concisely by partitioning the set of terminal nodes in the game tree into three different sets. In particular, let  $Z_{\hat{I},\hat{a}}$  be the

set of terminal nodes whose path from the root of the tree contains taking action  $\hat{a}$  at  $\hat{I}$  and let  $Z_{\hat{I}}$  be the set of terminal nodes whose path from the root passes through  $\hat{I}$  and are *not* in  $Z_{\hat{I},\hat{a}}$ . We have

**Lemma 1.** Consider  $a(\hat{\sigma}, \hat{\mu})$ -trigger agent for Player 1, where  $\hat{\sigma} = (\hat{I}, \hat{a})$ . The value of the trigger agent, defined as the expected difference between the utility of the trigger agent and the utility of an agent that always follows recommendations sampled from correlated distribution  $\mu$ , is computed as

$$v_{1,\hat{\sigma}}(\mu,\hat{\mu}) := \sum_{z \in Z_{\hat{I}}} u_1(z)\xi_1(\hat{\sigma};z)y_{1,\hat{\sigma}}(z) - \sum_{z \in Z_{\hat{I},\hat{\sigma}}} u_1(z)\xi_1(\sigma_1(z);z),$$

where 
$$\xi_1(\hat{\sigma};z) := \sum_{\pi_1 \in \Pi_1(\hat{\sigma})} \sum_{\pi_2 \in \Pi_2(z)} \mu(\pi_1, \pi_2)$$
 and  $y_{1,\hat{\sigma}}(z) := \sum_{\hat{\pi}_1 \in \Pi_1(z)} \hat{\mu}(\hat{\pi}_1)$ .

(A symmetric result holds for Player 2, with symbols  $\xi_2(\hat{\sigma};z)$  and  $y_{2,\hat{\sigma}}(z)$ .) It now seems natural to perform a change of variables, and pick distributions for the random variables  $y_{1,\hat{\sigma}}(\cdot),y_{2,\hat{\sigma}}(\cdot),\xi_1(\cdot;\cdot)$  and  $\xi_2(\cdot;\cdot)$  instead of  $\mu$  and  $\hat{\mu}$ . Since there are only a polynomial number (in the game tree size) of combinations of arguments for these new random variables, this approach allows one to remove the redundancy of realization-equivalent normal-form plans and focus on a significantly smaller search space. In fact, the definition of  $\xi=(\xi_1,\xi_2)$  also appears in (von Stengel & Forges, 2008), referred to as (sequence-form) correlation plan. In the case of the  $y_{1,\hat{\sigma}}$  and  $y_{2,\hat{\sigma}}$  random variables, it is clear that the change of variables is possible via the sequence form (von Stengel, 2002); we let  $Y_{i,\hat{\sigma}}$  be the sequence-form polytope of feasible values for the vector  $y_{i,\hat{\sigma}}$ . Hence, the only hurdle is characterizing the space spanned by  $\xi_1$  and  $\xi_2$  as  $\mu$  varies across the probability simplex. In two-player perfect-recall games with no chance moves, this is exactly one of the merits of the landmark work by von Stengel & Forges (2008). In particular, the authors prove that in those games the space of feasible  $\xi$  can be captured by a polynomial number of linear constraints. In more general cases the same does not hold (see second half of Section 2), but we prove the following (Appendix C):

**Lemma 2.** In a two-player game, as  $\mu$  varies over the probability simplex, the joint vector of  $\xi_1(\cdot;\cdot)$ ,  $\xi_2(\cdot;\cdot)$  variables spans a convex polytope  $\mathcal{X}$  in  $\mathbb{R}^n$ , where n is at most quadratic in the game size.

**Saddle-Point Reformulation** According to Lemma 1, for each Player i and  $(\hat{\sigma}, \hat{\mu})$ -trigger agent for them, the value of the trigger agent is a biaffine expression in the vectors  $y_{i,\hat{\sigma}}$  and  $\xi_i$ , and can be written as  $v_{i,\hat{\sigma}}(\xi_i,y_{i,\hat{\sigma}})=\xi_i^{\top}A_{i,\hat{\sigma}}y_{i,\hat{\sigma}}-b_{i,\hat{\sigma}}^{\top}\xi_i$  for a suitable matrix  $A_{i,\hat{\sigma}}$  and vector  $b_{i,\hat{\sigma}}$ , where the two terms in the difference correspond to the expected utility for deviating at  $\hat{\sigma}$  according to the (sequence-form) strategy  $y_{i,\hat{\sigma}}$  and the expected utility for not deviating at  $\hat{\sigma}$ . Given a correlation plan  $\xi=(\xi_1,\xi_2)\in\mathcal{X}$ , the maximum value of any deviation for any player can therefore be expressed as

$$v^*(\xi) := \max_{\{i, \hat{\sigma}, y_{i, \hat{\sigma}}\}} v_{i, \hat{\sigma}}(\xi_i, y_{i, \hat{\sigma}}) = \max_{i \in \{1, 2\}} \max_{\hat{\sigma} \in \Sigma_i} \max_{y_{\hat{\sigma}} \in Y_{\hat{\sigma}}} \{\xi_i^{\top} A_{i, \hat{\sigma}} y_{i, \hat{\sigma}} - b_{i, \hat{\sigma}}^{\top} \xi_i\}.$$

We can convert the maximization above into a continuous linear optimization problem by introducing the multipliers  $\lambda_{i,\hat{\sigma}} \in [0,1]$  (one per each Player  $i \in \{1,2\}$  and trigger  $\hat{\sigma} \in \Sigma_i$ ), and write

$$v^*(\xi) = \max_{\{\lambda_{i,\hat{\sigma}}, z_{i,\hat{\sigma}}\}} \sum_{i} \sum_{\hat{\sigma}} \xi_i^\top A_{i,\hat{\sigma}} z_{i,\hat{\sigma}} - \lambda_{i,\hat{\sigma}} b_{i,\hat{\sigma}}^\top \xi_i,$$

where the maximization is subject to the linear constraints  $[C_1] \sum_{i \in \{1,2\}} \sum_{\hat{\sigma} \in \Sigma_i} \lambda_{i,\hat{\sigma}} = 1$  and  $[C_2] z_{i,\hat{\sigma}} \in \lambda_{i,\hat{\sigma}} Y_{i,\hat{\sigma}}$  for all  $i \in \{1,2\}, \hat{\sigma} \in \Sigma_i$ . These linear constraints define a polytope  $\mathcal{Y}$ .

A correlation plan  $\xi$  is an EFCE if an only if  $v_{i,\hat{\sigma}}(\xi,y_{i,\hat{\sigma}}) \leq 0$  for every trigger agent, i.e.,  $v^*(\xi) \leq 0$ . Therefore, to find an EFCE, we can solve the optimization problem  $\min_{\xi \in \mathcal{X}} v^*(\xi)$ , which is a bilinear saddle point problem over the convex domains  $\mathcal{X}$  and  $\mathcal{Y}$ , both of which are convex polytopes that belong to  $\mathbb{R}^n$ , where n is at most quadratic in the input game size (Lemma 2). If an EFCE exists, the optimal value should be non-positive and the optimal solution is an EFCE (as it satisfies  $v^*(\xi) \leq 0$ ). In fact, since EFCE's always exist (as EFCEs are supersets of CEs (von Stengel & Forges, 2008)), and one can select triggers to be terminal sequences for Player 1, the optimal value of the BSPP is always 0. The BSPP can be interpreted as the NE of a zero-sum game between the *mediator*, who decides on a suitable correlation plan  $\xi$  and a deviator who selects the  $y_{i,\hat{\sigma}}$ 's to maximize each  $v_{i,\hat{\sigma}}(\xi_i,y_{i,\hat{\sigma}})$ . The value of this game is always 0.

Finally, we can enforce a minimum lower bound  $\tau$  on the sum of players' utility by introducing an additional variable  $\lambda_{sw} \in [0, 1]$  and maximizing the new convex objective

$$v_{\text{sw}}^{*}(\xi) := \max_{\lambda_{\text{sw}} \in [0,1]} \left\{ (1 - \lambda_{\text{sw}}) \cdot v^{*}(\xi) + \lambda_{\text{sw}} \left[ \tau - \sum_{z \in Z} u_{1}(z)\xi_{1}(z;z) - \sum_{z \in Z} u_{2}(z)\xi_{2}(z;z) \right] \right\}. \tag{1}$$

# 4. Computing an EFCE using Subgradient Descent

(von Stengel & Forges, 2008) show that a SW-maximizing EFCE of a two-player game without chance may be expressed as the solution of an LP and solved using generic methods such as the simplex algorithm or interior-point methods. However, this does not scale to large games as these methods require storing and inverting large matrices. Another way of computing SW-maximizing EFCEs was provided by (Dudik & Gordon, 2009). However, their algorithm assumes that sampling from correlation plans is possible using the Monte Carlo Markov chain algorithm and does not factor in convergence of the Markov chain. Furthermore, even though their formulation generalizes beyond our setting of two-player games without chance, our gradient descent method admits more complex objectives. In particular, it allows the mediator to maximize over general concave objectives (in correlation plans) instead of only linear objectives with potentially some regularization. Here, we showcase the benefits of exploiting the combinatorial structure of the BSPP formulation of Section 3 by proposing a simple algorithm based on subgradient descent; in Section 6 we show that this method scales better than commercial state-of-the-art LP solver in large games.

For brevity, we only provide a sketch of our algorithm, which computes a feasible EFCE; the extension to the slightly more complicated objective  $v_{\rm sw}^*(\xi)$  (Equation 1) is straightforward—see Appendix D for more details. First, observe that the objective  $v^*(\xi)$  is convex since it is the maximum of linear functions of  $\xi$ . This suggests that we may perform subgradient descent on  $v^*$ , where the subgradients are given by  $\partial/\partial\xi\,v^*(\xi)=A_{i^*,\hat\sigma^*}y_{i^*,\hat\sigma^*}^*-b_{i,\hat\sigma^*}$ , where  $(i^*,\hat\sigma^*,y_{i^*,\hat\sigma^*}^*)$  is a triplet which maximizes the objective function  $v^*(\xi)$ . The computation of such a triplet can be done via a straightforward bottom-up traversal of the game tree. In order to maintain feasibility (that is,  $\xi\in\mathcal{X}$ ), it is necessary to project onto  $\mathcal{X}$ , which is challenging in practice because we are not aware of any distance-generating function that allows for efficient projection onto this polytope. This is so even in games without chance (where  $\xi$  can be expressed by a polynomial number of constraints (von Stengel & Forges, 2008)). Furthermore, iterative methods such as Dykstra's algorithm, add a dramatic overhead to the cost of each iterate.

To overcome this hurdle, we observe that in games with no chance moves, the set  $\mathcal{X}$  of correlation plans—as characterized by von Stengel & Forges (2008) via the notion of consistency constraints—can be expressed as the intersection of three sets: (i)  $\mathcal{X}_1$ , the sets of vectors  $\xi$  that only satisfy consistency constraints for Player 1; (ii)  $\mathcal{X}_2$ , the sets of vectors  $\xi$  that only satisfy consistency constraints for Player 2; and (iii)  $\mathbb{R}^n_+$ , the non-negative orthant.  $\mathcal{X}_1$  and  $\mathcal{X}_2$  are polytopes defined by equality constraints only. Therefore, an exact projection (in the Euclidean sense) onto  $\mathcal{X}_1$  and  $\mathcal{X}_2$  can be carried out efficiently by precomputing a suitable factorization the constraint matrices that define  $\mathcal{X}_1$  and  $\mathcal{X}_2$ . In particular, we are able to leverage the specific combinatorial structure of the constraints that form  $\mathcal{X}_1$  and  $\mathcal{X}_2$  to design an efficient and parallel sparse factorization algorithm (see Appendix D for the full details). Furthermore, projection onto the non-negative orthant can be done conveniently, as it just amounts to computing a component-wise maximum between  $\xi$ and the zero vector. Since  $\mathcal{X} = \mathcal{X}_1 \cap \mathcal{X}_2 \cap \mathbb{R}_+^n$ , and since projecting onto  $\mathcal{X}_1, \mathcal{X}_2$  and  $\mathbb{R}_+^n$  individually is easy, we can adopt the recent algorithm proposed by (Wang & Bertsekas, 2013) designed to handle exactly this situation. In that algorithm, gradient steps are interlaced with projections onto  $\mathcal{X}_1$ ,  $\mathcal{X}_2$  and  $\mathbb{R}^n_+$  in a cyclical manner. This is similar to projected gradient descent, but instead of projecting onto the intersection of  $\mathcal{X}_1$ ,  $\mathcal{X}_2$  and  $\mathbb{R}_+^n$  (which we believe to be difficult), we project onto just one of them in round-robin fashion. This simple method was shown to converge by (Wang & Bertsekas, 2013). However, no convergence bound is currently known.

## 5. Introducing the First Benchmarks for EFCE

In this section we introduce the first two benchmark games for EFCE. These games are naturally parametric so that they can scale in size as desired and hence used to evaluate different EFCE solvers. In addition, we show that the EFCE in these games are interesting behaviorally: the correlation plan in social-welfare-

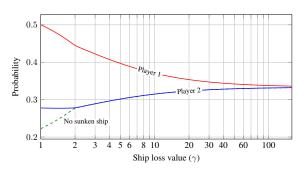
maximizing EFCE is highly nontrivial and even seemingly counter-intuitive. We believe some of these induced behaviors may prove practical in real-world scenarios and hope our analysis can spark an interest in EFCEs and other equilibria in sequential settings.

# 5.1. Battleship: Conflict Resolution via a Mediator

In this section we introduce our first proposed benchmark game to illustrate the power of correlation in extensive-form games. Our game is a general-sum variant of the classic game Battleship. Each player takes turns to secretly place a set of ships  $\mathcal S$  (of varying sizes and value) on separate grids of size  $H\times W$ . After placement, players take turns firing at their opponent—ships which have been hit at all the tiles they lie on are considered destroyed. The game continues until either one player has lost all of their ships, or each player has completed r shots. At the end of the game, the payoff of each player is computed as the sum of the values of the opponent's ships that were destroyed, minus  $\gamma$  times the value of ships which they lost, where  $\gamma \geq 1$  is called the  $loss\ multiplier$  of the game. The  $social\ welfare\ (SW)$  of the game is the sum of utilities to all players.

In order to illustrate a few interesting feature of social-welfare-maximizing EFCE in this game, we will focus on the instance of the game with a board of size  $3 \times 1$ , in which each player commands just 1 ship of value and length 1, there are 2 rounds of shooting per player, and the loss multiplier is  $\gamma=2$ . In this game, the social-welfare-maximizing *Nash* equilibrium is such that each player places their ship and shoots uniformly at random. This way, the probability that Player 1 and 2 will end the game by destroying the opponent's ship is 5/9 and 1/3 respectively (Player 1 has an advantage since they act first). The probability that both players will end the game with their ships unharmed is a meagre 1/9. Correspondingly, the maximum SW reached by any NE of the game is -8/9.

In the EFCE model, it is possible to induce the players to end the game with a peaceful outcome—that is, no damage to either ship—with probability  $^5/18$ , 2.5 times of the probability in NE, resulting in a much-higher SW of  $^{-13}/18$ . Before we continue with more details as to how the mediator (correlation device) is able to achieve this result in the case where  $\gamma=2$ , we remark that the benefit of EFCE is even higher when the loss multiplier  $\gamma$  increases: Figure 1 (left) shows, as a function of  $\gamma$ , the probability with which Player 1 and 2 terminate the game by sinking their opponent's ship, if they play according to the SW-maximizing EFCE. For all values of  $\gamma$ , the SW-maximizing NE remains the same while with a mediator, the probability of reaching a peaceful outcome increases as  $\gamma$  increases, and asymptotically gets closer to  $^1/3$  and the gap between the expected utility of the two players vanishes. This is remarkable, considering Player 1's advantage for acting first.



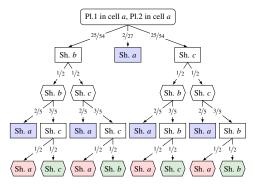


Figure 1: (Left) Probabilities of players sinking their opponent when the players play according to the SW-maximizing EFCE. For  $\gamma \geq 2$ , the probability of the game ending with no sunken ship and the probability of Player 2 sinking Player 1 coincide. (Right) Example of a playthrough of Battleship assuming both players are recommended to place their ship in the same position a. Edge labels represents the probability of an action being recommended. Squares and hexagons denote actions taken by Players 1 and 2 respectively. Blue and red nodes represent cases where Players 1 and 2 sink their opponent, respectively. The *Shoot* action is abbreviated 'Sh.'.

We now resume our analysis of the SW-maximizing EFCE in the instance where  $\gamma=2$ . In a nutshell, the correlation plan is constructed in a way that players are recommended to deliberately miss, and deviations

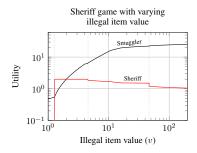
from this are punished by the mediator, who reveals to the opponent the ship location that was recommended to the deviating player. First, the mediator recommends the players a ship placement that is sampled uniformly at random and independently for each players. This results in 9 possible scenarios (one per possible ship placement) in the game, each occurring with probability 1/9. Due to the symmetric nature of ship placements, only two scenarios are relevant: whether the two players are recommended to place their ship in the same spot, or in different spots. Figure 1 (right) shows the probability of each recommendation from the mediator in the former case, assuming that the players do not deviate. The latter case is symmetric (see Appendix E for details). Now, we explain the first of the two methods in which the mediator compels non-violent behavior. We focus on the first shot made by Player 1 (i.e., the root in Figure 3). The mediator suggests that Player 1 shoot at the Player 2's ship with a low 2/27 probability, and deliberately miss with high probability. One may wonder how it is possible for this behavior to be incentive-compatible (that is, what are the incentives that compel Player 1 into not defecting), since the player may choose to randomly fire in any of the 2 locations that were *not* recommended, and get almost 1/2 chance of winning the game immediately. The key is that if Player 1 does so and does not hit the opponent's ship, then the mediator can punish him by recommending that Player 2 shoot in the position where Player 1's was recommended to place their ship. Since players value their ships more than destroying their opponents', the player is incentivized to avoid such a situation by accepting the recommendation to (most probably) miss. We see the first example of deterrent used by the mediator: inducing the opponent to play punitive actions against players that have deviated from the recommendation, if ever that deviation can be detected from the player. A similar situation arises in the first move of Player 2, where Player 2 is recommended to *deliberately* miss, hitting each of the 2 empty spots with probability 1/2. A more detailed analysis is available in Appendix E.

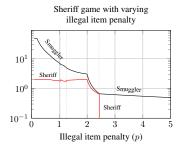
#### 5.2. Sheriff: Bargaining and Negotiation

Our second proposed benchmark is a simplified version of the Sheriff of Nottingham board game. The game models the interaction of two players: the *Smuggler*—who is trying to smuggle illegal items in their cargo—and the Sheriff—who is trying to stop the Smuggler. At the beginning of the game, the Smuggler secretly loads his cargo with  $n \in \{0, \dots, n_{\text{max}}\}$  illegal items. At the end of the game, the Sheriff decides whether to inspect the cargo. If the Sheriff chooses to inspect the cargo and finds illegal goods, the Smuggler must pay a fine worth  $p \cdot n$  to the Sheriff. On the other hand, the Sheriff has to compensate the Smuggler with a utility s if no illegal goods are found. Finally, if the Sheriff decides not to inspect the cargo, the Smuggler's utility is  $v \cdot n$  whereas the Sheriff's utility is 0. The game is made interesting by two additional elements (which are also present in the board game): bribery and bargaining. After the Smuggler has loaded the cargo and before the Sheriff chooses whether or not to inspect, they engage in r rounds of bargaining. At each round  $i = 1, \dots, r$ , the Smuggler tries to tempt the Sheriff into not inspecting the cargo by proposing a bribe  $b_i \in \{0, \dots b_{\max}\}$ , and the Sheriff responds whether or not they would accept the proposed bribe. Only the proposal and response from round r will be executed and have an impact on the final payoffs—that is, all but the r-th round of bargaining are non-consequential and their purpose is for the two players to settle on a suitable bribe amount. If the Sheriff accepts bribe  $b_r$ , then the Smuggler gets  $p \cdot n - b_r$ , while the Sheriff gets  $b_r$ . See Appendix F for a formal description of the game.

We now point out some interesting behavior of EFCE in this game. We refer to the game instance where  $v=5, p=1, s=1, n_{\text{max}}=10, b_{\text{max}}=2, r=2$  as the baseline instance.

Effect of v,p and s. First, we show what happens in the baseline instance when the item value v, item penalty p, and Sheriff compensation (penalty) s are varied in isolation over a continuous range of values. The results are shown in Figure 2. In terms of general trends, the effect of the parameter to the Smuggler is fairly consistent with intuition: the Smuggler benefits from a higher item value as well as from higher sheriff penalties, and suffers when the penalty for smuggling is increased. However, the finer details are much more nuanced. For one, the effect of changing the parameters not only is non-monotonic, but also discontinuous. This behavior has never been documented and we find it rather counterintuitive. More counterintuitive observations can be found in Appendix F. Effect of  $n_{\text{max}}$ ,  $b_{\text{max}}$ , and r. Here, we try to empirically understand the impact of n and p on the SW maximizing equilibrium. As before we set p = 1, p = 1 and vary p and p simultaneously while keeping p = 1 and vary p and p simultaneously while keeping p = 2 and vary p and p = 3 and vary p = 3 and vary p = 4 and vary p = 5 and vary p = 5 and vary p = 6 and vary p = 6 and vary p = 7 and vary p = 7 and vary p = 8 and vary p = 9 and vary p = 9 and vary p = 1 and vary p = 1 and vary p = 1 and vary p = 2 and vary p = 3 and vary p = 6 and vary p = 6 and vary p = 6 and vary p = 7 and vary p = 7 and vary p = 8 and vary p = 9 and va





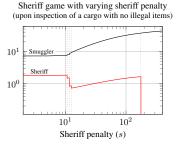


Figure 2: Utility of players with varying v, p and s for the SW-maximizing EFCE. We verified that these plots are not the result of equilibrium selection issues.

example, consider the case when  $b_{\max}=2, n_{\max}=2, r=1$  (shown in blue in Table 1, right) where the payoffs are (8.0, 2.0). This achieves the maximum attainable social welfare by smuggling  $n_{\max}=2$  items and having the Sheriff accept a bribe of 2. When  $n_{\max}$  is increased to 5 (red entry in the table), the payoffs to both players drop significantly, and even more so when  $n_{\max}$  increases further. While counter-intuitive, this behavior is consistent in that the Smuggler may not benefit from loading 3 items every time he was recommended to load 2; the Sheriff reacts by inspecting more, leading to lower payoffs for both players.

That behavior is avoided by increasing the number of rounds r: by increasing to r=2 (entry shown in purple), the behavior disappears and we revert to achieving a social welfare of 10 just like in the instance with  $n_{\rm max}=2, r=1$ . With sufficient bargaining steps, the Smuggler, with the aid of the mediator, is able to convince the Sheriff that they have complied with the recommendation by the mediator. This is because the mediator spends the first r-1 bribes to give a 'passcode' to the Smuggler so that the Sheriff can verify

| $n_{ m max}$ | r = 1        | r=2                   | r = 3        |
|--------------|--------------|-----------------------|--------------|
| 1            | (3.00, 2.00) | (3.00, 2.00)          | (3.00, 2.00) |
| 2            | (8.00, 2.00) | (8.00, 2.00)          | (8.00, 2.00) |
| 5            | (2.28, 1.26) | ( <b>8.00, 2.00</b> ) | (8.00, 2.00) |
| 10           | (1.76, 0.93) | (7.26, 1.82)          | (8.00, 2.00) |

Table 1: Payoffs for (Smuggler, Sheriff) in the SW-maximizing EFCE.

compliance—if an 'unexpected' bribe is suggested, then the Smuggler must have deviated, and the Sheriff will inspect the cargo as punishment. With more rounds, it is less likely that the Smuggler will guess the correct passcode. See also Appendix F for additional insights.

#### 6. Experimental Evaluation

Even our proof-of-concept algorithm based on the BSSP formulation and subgradient descent, introduced in Section 3, is able to beat LP-based approaches using the commercial solver Gurobi (Gurobi Optimization, 2018) in large games. This confirms known results about the scalability of methods for computing NE, where in the recent years first-order methods have affirmed themselves as the only algorithms that are able to handle large games.

We experimented on Battleship over a range of parameters while fixing  $\gamma=2$ . All experiments were run on a machine with 64 cores and 500GB of memory. For our method, we tuned step sizes based on multiples of 10. In Table 2, we report execution times when all constraints (feasibility and deviation) are violated by no greater than  $10^{-1}$ ,  $10^{-2}$  and  $10^{-3}$ . Our method outperforms the LP-based approach for larger games. However, while we outperform the LP-based approach for accuracies up to  $10^{-3}$ , Gurobi spends most of its time reordering variables and preprocessing and its solution converges faster for higher levels of precision; this is expected of a gradient-based method like ours. On very large games with more than 100 million variables, both our method and Gurobi fail—in Gurobi's case, it was due to a lack of memory while in our case, each iteration required nearly an hour which was prohibitive. The main bottleneck in our method was the projection onto  $\mathcal{X}_1$  and  $\mathcal{X}_2$ . We also experimented on the Sheriff game and obtained similar findings (Appendix H).

| (H, W) | r | Ship   | #Actions |       | #Relevant  | Time (LP) $10^{-1}$ $10^{-2}$ $10^{-3}$ |        |        | Time (ours) |                       |      |  |  |
|--------|---|--------|----------|-------|------------|---|--------|--------|-------------|-----------------------|------|--|--|
|        |   | length | PII      | F1 Z  | seq. pairs | 10 -                                    | 10 -   | 10 °   | 10 -        | 10 -                  | 10 " |  |  |
| (2, 2) | 3 | 1      | 741      | 917   | 35241      | 2s                                      | 2s     | 2s     | 1s          | 2s                    | 3s   |  |  |
| (3, 2) | 3 | 1      | 15k      | 47k   | 3.89M      | 3m 6s                                   | 3m 17s | 3m 24s | 8s          | 34s                   | 52s  |  |  |
|        |   |        |          |       | 26.4M      |   |        |        |             |                       |      |  |  |
| (3, 2) | 4 | 2      | 970k     | 2.27M | 111M       | — out of memory <sup>†</sup> —          |        |        | — did       | — did not achieve ‡ — |      |  |  |

Table 2: #Relevant seq. pairs is the dimension of  $\xi$  under the compact representation of (von Stengel & Forges, 2008). For LPs, we report the fastest of Barrier, Primal and Dual Simplex, and 3 different formulations (Appendix G). † Gurobi went out of memory and was killed by the system after  $\sim 3000$  seconds † Our method requires 1 hour per iteration and did not achieve the required accuracy after 6 hours.

#### 7. Conclusions

In this paper, we proposed two parameterized benchmark games in which EFCE exhibits interesting behaviors. We analyzed those behaviors both qualitatively and quantitatively, and isolated two ways through which a mediator is able to compel the agents to follow the recommendations. We also provided an alternative saddle-point formulation of EFCE and demonstrated its merit with a simple subgradient method which outperforms standard LP based methods.

We hope that our analysis will bring attention to some of the computational and practical uses of EFCE, and that our benchmark games will be useful for evaluating future algorithms for computing EFCE in large games.

# Acknowledgments

This material is based on work supported by the National Science Foundation under grants IIS-1718457, IIS-1617590, and CCF-1733556, and the ARO under award W911NF-17-1-0082. Gabriele Farina is supported by a Facebook fellowship. Co-authors Ling and Fang are supported in part by a research grant from Lockheed Martin.

## References

Ashlagi, I., Monderer, D., and Tennenholtz, M. On the value of correlation. *Journal of Artificial Intelligence Research*, 33:575–613, 2008.

Aumann, R. Subjectivity and correlation in randomized strategies. *Journal of Mathematical Economics*, 1: 67–96, 1974.

Bowling, M., Burch, N., Johanson, M., and Tammelin, O. Heads-up limit hold'em poker is solved. *Science*, 2015.

Brown, N. and Sandholm, T. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, Dec. 2017.

Brown, N. and Sandholm, T. Solving imperfect-information games via discounted regret minimization. In *AAAI*, 2019a.

Brown, N. and Sandholm, T. Superhuman AI for multiplayer poker. *Science*, 365(6456):885-890, 2019b. ISSN 0036-8075. doi: 10.1126/science.aay2400. URL https://science.sciencemag.org/content/365/6456/885.

Crawford, V. P. and Sobel, J. Strategic information transmission. *Econometrica: Journal of the Econometric Society*, pp. 1431–1451, 1982.

Dodis, Y., Halevi, S., and Rabin, T. A cryptographic solution to a game theoretic problem. In *Annual International Cryptology Conference*, pp. 112–130. Springer, 2000.

- Dudik, M. and Gordon, G. J. A sampling-based approach to computing equilibria in succinct extensive-form games. In *UAI*, pp. 151–160. AUAI Press, 2009.
- Farina, G., Kroer, C., and Sandholm, T. Online convex optimization for sequential decision processes and extensive-form games. In *AAAI Conference on Artificial Intelligence*, 2019.
- Gordon, G. J., Greenwald, A., and Marks, C. No-regret learning in convex games. In *Proceedings of the* 25<sup>th</sup> international conference on Machine learning, pp. 360–367. ACM, 2008.
- Gurobi Optimization, L. Gurobi optimizer reference manual, 2018. URL http://www.gurobi.com.
- Hazan, E. Introduction to online convex optimization. Foundations and Trends in Optimization, 2016.
- Huang, W. *Equilibrium computation for extensive games*. PhD thesis, London School of Economics and Political Science, January 2011.
- Huang, W. and von Stengel, B. Computing an extensive-form correlated equilibrium in polynomial time. In *International Workshop On Internet And Network Economics (WINE)*, pp. 506–513. Springer, 2008.
- Jiang, A. X. and Leyton-Brown, K. Polynomial-time computation of exact correlated equilibrium in compact games. *Games and Economic Behavior*, 91:347–359, 2015.
- Koutsoupias, E. and Papadimitriou, C. Worst-case equilibria. In *Symposium on Theoretical Aspects in Computer Science*, 1999.
- Leyton-Brown, K. Personal communication, 2019.
- Moravčík, M., Schmid, M., Burch, N., Lisý, V., Morrill, D., Bard, N., Davis, T., Waugh, K., Johanson, M., and Bowling, M. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 2017.
- Nash, J. Equilibrium points in n-person games. *Proceedings of the National Academy of Sciences*, 36: 48–49, 1950.
- Nemirovski, A. Prox-method with rate of convergence O(1/t) for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 2004.
- Nesterov, Y. Excessive gap technique in nonsmooth convex minimization. *SIAM Journal of Optimization*, 2005.
- Papadimitriou, C. H. and Roughgarden, T. Computing correlated equilibria in multi-player games. *Journal of the ACM*, 55(3):14, 2008.
- Ray, I. and Gupta, S. S. Technical Report, 2009.
- Roughgarden, T. and Tardos, É. How bad is selfish routing? *Journal of the ACM (JACM)*, 49(2):236–259, 2002.
- Shoham, Y. and Leyton-Brown, K. *Multiagent systems: Algorithmic, game-theoretic, and logical foundations.* Cambridge University Press, 2009.
- von Stengel, B. Efficient computation of behavior strategies. Games and Economic Behavior, 1996.
- von Stengel, B. Computing equilibria for two-person games. In Aumann, R. and Hart, S. (eds.), *Handbook of game theory*, volume 3. North Holland, Amsterdam, The Netherlands, 2002.
- von Stengel, B. and Forges, F. Extensive-form correlated equilibrium: Definition and computational complexity. *Mathematics of Operations Research*, 33(4):1002–1022, 2008.
- Wang, M. and Bertsekas, D. P. Incremental constraint projection-proximal methods for nonsmooth convex optimization. *SIAM J. Optim.(to appear)*, 2013.
- Zinkevich, M., Bowling, M., Johanson, M., and Piccione, C. Regret minimization in games with incomplete information. In *NIPS*, 2007.

# A. Recovering the Linear Program of (von Stengel & Forges, 2008)

Recall the continuous version of the primal version of the inner maximization problem which was obtained by adding the multipliers  $\lambda_{i,\hat{\sigma}} \in [0,1]$ .

$$\begin{split} \max_{\lambda, z_{i, \hat{\sigma}}} \quad & \sum_{i \in \{1, 2\}} \sum_{\hat{\sigma} \in \Sigma_{i}} \xi_{i}^{\top} A_{i, \hat{\sigma}} z_{i, \hat{\sigma}} - \lambda_{i, \hat{\sigma}} b_{i, \hat{\sigma}}^{\top} \xi_{i} \\ \text{such that} \quad & \sum_{i \in \{1, 2\}} \sum_{\hat{\sigma} \in \Sigma_{i}} \lambda_{i, \hat{\sigma}} = 1 \\ & \lambda_{i, \hat{\sigma}} \geq 0 \\ & z_{i, \hat{\sigma}} \in \lambda_{i, \hat{\sigma}} Y_{i, \hat{\sigma}}, \qquad \forall i \in \{1, 2\}, \hat{\sigma} \in \Sigma_{i} \end{split}$$

where  $z_{i,\hat{\sigma}}$  may be seen as the sequence form representation of a game rooted at a particular information set of player i, and scaled by the factor  $\lambda_{i,\hat{\sigma}}$ . By expanding the sequence form constraints which define  $Y_{i,\hat{\sigma}}$ , we get

$$\begin{split} \max_{\lambda,z} \quad & \sum_{i \in \{1,2\}} \sum_{\hat{\sigma} \in \Sigma_i} \xi_i^\top A_{i,\hat{\sigma}} z_{i,\hat{\sigma}} - \lambda_{i,\hat{\sigma}} b_{i,\hat{\sigma}}^\top \xi_i \\ \text{such that} \quad & \sum_{i \in \{1,2\}} \sum_{\hat{\sigma} \in \Sigma_i} \lambda_{i,\hat{\sigma}} = 1 \\ & \lambda_{i,\hat{\sigma}} \geq 0 \\ & z_{i,\hat{\sigma}} \geq 0 \\ & F_{i,\hat{\sigma}} z_{i,\hat{\sigma}} - \lambda_{\hat{\sigma}} f_{i,\hat{\sigma}} = 0, \qquad \forall i \in \{1,2\}, \hat{\sigma} \in \Sigma_i \end{split}$$

where  $F_{i,\hat{\sigma}}$  and  $f_{i,\hat{\sigma}}$  are sequence form constraint matrices rooted at the information set  $\hat{I}$  containing  $\hat{\sigma}$ , with the only difference that instead of having the 'empty sequence' be equal to 1, we require that all actions belonging to  $\hat{I}$  sum to  $\lambda_{i,\hat{\sigma}}$ . We are now in a position to take duals; the only non-zero elements on the right hand side of the constraints are from the sum-to-one constraints over  $\lambda_{i,\hat{\sigma}}$ . This give s the following dual

$$\begin{split} \min_{u,\nu_i(\hat{\sigma},\cdot)} & u \\ \text{such that} & F_{i,\hat{\sigma}}^T \nu_i(\hat{\sigma},\cdot) \geq A_{i,\hat{\sigma}}^T \xi_i & \forall i \in \{1,2\}, \hat{\sigma} \in \Sigma_i \\ & u - \nu_i(\hat{\sigma},\hat{I}) \geq -b_{i,\hat{\sigma}}^T \xi_i & \forall i \in \{1,2\}, \hat{\sigma} = (\hat{I},\hat{a}) \in \Sigma_i, \end{split}$$

where u and  $\nu(\hat{\sigma}, \cdot)$  are free in sign. Combining this with the outer minimization over  $\xi_i$  gives us the linear program by (von Stengel & Forges, 2008), up to a change in variable names and conventions.

## **B.** Derivation of Probabilities over Terminal States

In order to express the utility of a trigger agent, it is necessary to compute the probability of the game ending in each of the terminal states. Before that, we will review the notation introduced in earlier sections in more detail.

- Z be the set of terminal states (or equivalently, outcomes) in the game, and  $z \in Z$  is some terminal state.
- $u_i(z)$  be the utility obtained by player i if the game terminates at some terminal state  $z \in Z$ .
- $\Pi_i$  be the set of pure reduced-normal-form strategies for Player i. We also require notation for subsets of  $\Pi_i$ , namely,
  - $\Pi_i(I)$ , is the set of reduced-normal-form strategies that can lead to information set I (which belongs to player i) assuming that the other player acts to do so as well. This is equivalent (assuming no zero-chance nodes or disconnected game trees) to saying that all reduced-normal-from strategies in  $\Pi_i(I)$  have *some* action which belongs to information set I.
  - $\Pi_i(I, a)$  is the set of reduced normal form strategies which will lead to information set I and recommend the action a in I. This is equivalent to the set of reduced normal form strategies which contain a as part of their recommendation (this set is typically a subset of  $\Pi_i(I, a)$ ).

- $\Pi_i(z)$  is the set of reduced-normal-form strategies which can lead to the terminal state z (assuming the other player players to do so). This is equivalent to the set of reduced-normal-form strategies which contain the  $\sigma=(I,a)$  pair where  $\sigma=(I,a)$  is the unique last information set-action pair which has to be encountered by player i before the terminal state z.
- $\Sigma_i$  the set of information set-action pairs (I, a) (also known as *sequences*), where I is an information set for Player i and a is an action at set I.
- $\sigma_i(z)$  is the last (I, a) pair belonging to Player i encountered before some terminal state  $z \in Z$ .

We are interested in characterizing the random variable  $t_{\hat{\sigma}}:\Pi_1\times\Pi_2\times\Pi_1(\hat{I})\to Z$  that maps a triple of reduced-normal-form strategies  $(\pi_1,\pi_2,\hat{\pi}_1)$  to the terminal state of the game that is reached when Player 1 is a  $\hat{\sigma}$ -trigger agent and Player 2 follows all recommendations. That is, we want to find the probabillity of terminating at each  $z\in Z$  for a  $\hat{\sigma}$ -trigger agent, given the mediator's joint distribution  $\mu$  over reduced normal form strategies and the trigger strategy  $\hat{\mu}$  for the deviating player, which we will assume to be Player 1 without loss of generality. For each trigger  $\hat{\sigma}$ , the terminal leaves may be partitioned into the following 3 sets.

•  $Z_{\hat{\sigma}}$  (or equivalently  $Z_{\hat{I},\hat{a}}$ ) is the set of terminal nodes that are descendants of the trigger  $\hat{\sigma}=(\hat{I},\hat{a})$ . In order for the game to end in one of these terminal nodes, it is necessary that the recommendation device recommended to Player 1 the trigger sequence  $\hat{\sigma}$ , and therefore the agent must have deviated. Furthermore, Player 2 must have been recommended the terminal sequence  $\sigma_2(z)$  corresponding to the terminal state, and finally  $\hat{\pi}_1$  must be compatible with  $\sigma_1(z)$ . We can capture all these constraints concisely by saying that the sampled  $(\pi_1,\pi_2,\hat{\pi}_1)$  must be such that  $\pi_1\in\Pi_1(\hat{\sigma}),\,\pi_2\in\Pi_2(z)$  and  $\hat{\pi}_1\in\Pi_1(z)$ . Therefore the probability that a  $\hat{\sigma}$  trigger agent terminates at some  $z\in Z_{\hat{\sigma}}$  is given by,

$$\mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}} = z \in Z_{\hat{\sigma}}] = \left(\sum_{\substack{\pi_1 \in \Pi_1(\hat{\sigma}) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2)\right) \left(\sum_{\hat{\pi}_1 \in \Pi_1(z)} \hat{\mu}_1(\hat{\pi}_1)\right),$$

where the first term in the product is the probability that Player 2 plays to z and Player 1 gets triggered, and the second term is the probability that the deviation strategy from Player 1 upon getting triggered is one that reaches z.

•  $Z_{\hat{I}}$  is the set of terminal states that are descendant of any sequence in  $\hat{I}$ , except  $\hat{\sigma}$ . In order for the game to reach this terminal state, recommendations issued to Player 1 by the correlation device must have been such that Player 1 reached  $\hat{I}$ . There are two cases: either the correlation device recommended  $\hat{\sigma}$  at  $\hat{I}$ , or it did not. In the former case, Player 1 started deviating (using the sampled reduced-normal-form plan  $\hat{\pi}_1$ ); hence, in this case it must be  $\hat{\pi}_1 \in \Pi_1(z)$ . In the latter case, Player 1 does not deviate from the recommendation, and therefore it must be  $\pi_1 \in \Pi_1(z)$ . Either way, Player 2 must have been recommended the terminal sequence z corresponding to the terminal state z; that is,  $\pi_2 \in \Pi_2(z)$ . Collecting all these constraints, it must be

$$(\pi_1, \pi_2, \hat{\pi}_1) \in \Pi_1(\hat{\sigma}) \times \Pi_2(z) \times \Pi_1(z) \cup \Pi_1(z) \times \Pi_2(z) \times \Pi_1(\hat{I}).$$

Using the fact that the two cases as to whether or not Player 1 was recommended  $\hat{\sigma}$  or not at  $\hat{I}$  are disjoint, we can write

$$\mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}} = z \in Z_{\hat{I}}] = \left(\sum_{\substack{\pi_1 \in \Pi_1(\hat{\sigma}) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2)\right) \left(\sum_{\hat{\pi}_1 \in \Pi_1(z)} \hat{\mu}_1(\hat{\pi}_1)\right) + \left(\sum_{\substack{\pi_1 \in \Pi_1(z) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2)\right).$$

The first term in the summation may be understood as the probability that the agent was triggered and its deviation was to play something other than  $\hat{\sigma}$ . The second term is that probability that the agent was not triggered and the game simply terminates at z based on  $\mu$ .

• Finally,  $Z_{-\hat{I}}$  is the set of terminal nodes that are neither in  $Z_{\hat{\sigma}}$  nor in  $Z_{\hat{I}}$ . If the game has ended in any terminal state that belongs to  $Z_{-\hat{I}}$ , Player 1 has not deviated from the recommended strategy,

since they have never even reached the trigger information set,  $\hat{I}$ . Hence, in this case it must be  $(\pi_1, \pi_2) \in \Pi_1(z) \times \Pi_2(z)$ . Hence,

$$\mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}} = z \in Z_{-\hat{I}}] = \sum_{\substack{\pi_1 \in \Pi_1(z) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2).$$

With the above, we can finally express the constraint that no deviation strategy  $\hat{\mu}$  can lead to a higher utility for Player 1 than simply following each recommendation. Indeed, for all  $\hat{\mu}$ , the utility of the trigger agent is expressed as

$$\sum_{z \in Z} u_1(z) \mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}} = z],$$

where the correct expression for  $\mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}}=z]$  must be selected depending on whether  $z\in Z_{\hat{\sigma}},\,z\in Z_{\hat{I}}$  or  $z\in Z_{-\hat{I}}$ . On the other hand, the utility of an agent that follows all recommendations is

$$\sum_{z \in Z} u_1(z) \mathbb{P}_{\mu,\hat{\mu}}[\pi_1 \in \Pi_1(z), \pi_2 \in \Pi_2(z)] = \sum_{z \in Z} \left( u_1(z) \sum_{\substack{\pi_1 \in \Pi_1(z) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2) \right).$$

Therefore, following all recommendations is a best response for the  $\hat{\sigma}$ -trigger agent if and only if  $\mu$  is chosen so that

$$\sum_{z \in Z} u_1(z) \left( \mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}} = z] - \sum_{\substack{\pi_1 \in \Pi_1(z) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2) \right) \le 0 \qquad \forall \hat{\mu} \in \Delta^{|\Pi_1(\hat{I})|}. \tag{2}$$

The crucial observation is that all the probabilities  $\mathbb{P}_{\mu,\hat{\mu}}[t=z]$  defined above can be expressed via the following quantities:

$$y_{1,\hat{\sigma}}(z) := \sum_{\hat{\pi}_1 \in \Pi_1(z)} \hat{\mu}_1(\hat{\pi}_1) \quad \forall z \in Z; \qquad \xi_1(\sigma_1; z) := \sum_{\substack{\pi_1 \in \Pi_1(\sigma_1) \\ \pi_2 \in \Pi_2(z)}} \mu(\pi_1, \pi_2) \quad \forall \sigma_1 \in \Sigma_1, z \in Z.$$

For example, for all  $z \in Z_{\hat{I}}$  we can write

$$\mathbb{P}_{\mu,\hat{\mu}}[t_{\hat{\sigma}} = z] = \xi_1(\hat{\sigma}; z) y_{i,\hat{\sigma}}(z) + \xi_1(\sigma_1(z); z).$$

When deviations relative to Player 2 are brought into the picture, the following two sets of symmetric quantities also become relevant:

$$y_{2,\hat{\sigma}}(z) := \sum_{\hat{\pi}_2 \in \Pi_2(z)} \hat{\mu}_1(\hat{\pi}_2) \quad \forall z \in Z; \qquad \xi_2(\sigma_2; z) := \sum_{\substack{\pi_1 \in \Pi_1(z)) \\ \pi_2 \in \Pi_2(\sigma_2)}} \mu(\pi_1, \pi_2) \quad \forall \sigma_2 \in \Sigma_2, z \in Z.$$

It would now seem natural to perform a change of variables, and pick (correlated) distributions for the random variables  $y_{1,\hat{\sigma}}(\cdot), y_{2,\hat{\sigma}}(\cdot), \xi_1(\cdot; \cdot)$  and  $\xi_2(\cdot; \cdot)$  instead of  $\mu, \hat{\mu}_1$  and  $\hat{\mu}_2$ . Since there are only a polynomial number (in the game tree size) of combinations of arguments for these new random variables, this approach would allow one to remove the redundancy of realization-equivalent normal-form plans and focus on a polynomially-small search space. In the case of the random variables  $y_{1,\hat{\sigma}}$  and  $y_{2,\hat{\sigma}}$ , it is clear that the change of variables is possible via the sequence form (von Stengel, 2002). Therefore, the only difficulty is in characterizing the space spanned by  $\xi_1$  and  $\xi_2$  as  $\mu$  varies across the probability simplex. In two-player perfect-recall games with no chance moves, this is exactly the merit of the landmark work by von Stengel & Forges (2008). In particular, the authors prove that in those games the space of feasible  $\xi_1, \xi_2$  can be captured by a polynomial number of linear constraints.

## C. Proof of Lemma 2

The vectors of entries  $\xi_1(\cdot;\cdot)$ ,  $\xi_2(\cdot;\cdot)$  are obtained from  $\mu$  via a linear mapping. Hence, the set of values that can be assumed by  $\xi$  is the image of the probability simplex via a linear mapping. Since images of polytopes via linear functions are polytopes, the lemma holds.

# D. Details of Our Subgradient Method

First, observe that the objective  $v^*(\xi)$  is convex since it is the maximum of linear functions of  $\xi$ . This suggests that we may perform subgradient descent on  $v^*$ , where the subgradients are given by

$$\partial/\partial \xi \, v^*(\xi) = A_{i^*,\hat{\sigma}^*} y_{i^*,\hat{\sigma}^*}^* - b_{i,\hat{\sigma}^*}, \tag{3}$$

where  $(i^*, \hat{\sigma}^*, y^*_{i^*, \hat{\sigma}^*})$  is a triplet which maximizes the objective function  $v^*(\xi)$ . The computation of such a triplet is a straightforward bottom-up traversal of the game tree. In order to maintain feasibility (that is,  $\xi \in \mathcal{X}$ ), it is necessary to project onto  $\mathcal{X}$ , which is challenging in practice, because we are not aware of any distance-generating function which allows for efficient projection onto this polytope. This is so even in games without chance (where  $\xi$  can be expressed by a polynomial number of constraints (von Stengel & Forges, 2008)). Furthermore, iterative methods such as Dykstra's algorithm, add an dramatic overhead to the cost of each iterate.

To overcome this hurdle, we observe that in games with no chance moves, the set  $\mathcal{X}$  of correlation plans—as characterized by von Stengel & Forges (2008) via the notion of consistency constraints—can be expressed as the intersection of three sets: (i)  $\mathcal{X}_1$ , the sets of vectors  $\xi$  that only satisfy consistency constraints for Player 1; (ii)  $\mathcal{X}_2$ , the sets of vectors  $\xi$  that only satisfy consistency constraints for Player 2, respectively; and (iii)  $\mathbb{R}_{+}^{n}$ , the non-negative orthant.  $\mathcal{X}_{1}$  and  $\mathcal{X}_{2}$  are polytopes defined by equality constraints only. Therefore, an exact projection (in the Euclidean sense) onto  $\mathcal{X}_1$  and  $\mathcal{X}_2$  can be carried out efficiently by precomputing a suitable factorization the constraint matrices that define  $\mathcal{X}_1$  and  $\mathcal{X}_2$ . In particular, we are able to leverage the specific combinatorial structure of the constraints that form  $\mathcal{X}_1$  and  $\mathcal{X}_2$  to design an efficient and parallel sparse factorization algorithm (see Appendix D for the full details). Furthermore, projection onto the non-negative orthant can be done conveniently, as it just amounts to computing a component-wise maximum between  $\xi$  and the zero vector. Since  $\mathcal{X} = \mathcal{X}_1 \cap \mathcal{X}_2 \cap \mathbb{R}^n$ , and since projecting onto  $\mathcal{X}_1, \mathcal{X}_2$  and  $\mathbb{R}^n$ individually is easy, we can adopt the recent algorithm proposed by (Wang & Bertsekas, 2013) designed to handle exactly this situation. In that algorithm, gradient steps are interlaced with projections onto  $\mathcal{X}_1, \mathcal{X}_2$ and  $\mathbb{R}^n_+$  in a cyclical manner. This is similar to projected gradient descent, but instead of projecting onto the intersection of  $\mathcal{X}_1$ ,  $\mathcal{X}_2$  and  $\mathbb{R}_+^n$  (which we believe to be difficult), we project onto just one of them in round-robin fashion. This simple method was shown to converge by (Wang & Bertsekas, 2013), however, no convergence bound is currently known.

#### **D.1.** Factorization of constraints over $\mathcal{X}$

von Stengel & Forges (2008) showed that a  $\xi$  may be represented compactly as a 2-dimensional matrix, with dimensions equal to the sequence form representation (von Stengel, 1996) of each player, where one is only interested in entries corresponding to relevant sequence pairs (von Stengel & Forges (2008) for details). Then, the aforementioned constraints (i) and (ii) defining  $\mathcal{X}_1$  and  $\mathcal{X}_2$  are equivalent to the sequence form constraints for each row and column respectively. Constraint (iii) ensures that the entries of  $\xi$  are non-negative and that the entry for the empty sequence pair is 1.

Observe that projection (based on L2 distance) on  $\mathcal{X}_1$  and  $\mathcal{X}_2$  individually can be decomposed a series of disjoint projections (either on rows or columns) and thus computed in parallel. We now show that L2-projection of each individual row/column over the sequence form constraints (von Stengel & Forges, 2008) may be done efficiently. Let F and f be matrices and vectors corresponding to the sequence form constraints Fx - f = 0. Here, F is a (sparse) matrix of size #information sets  $\times$  #sequences which contains entries in  $\{-1,0,1\}$  and f is a vector containing 1's or 0's. Each information set in F corresponds to the 'flow' constraints for an information set, with a coefficient of -1 for the unique parent sequence leading to that

information set, and a coefficient of -1 for all sequences immediately following that information set. <sup>4</sup> Given a vector w, the projection onto the affine space given by Fx - f is given by the optimization problem

$$\min_{x} \qquad \frac{1}{2} ||x - w||_{2}^{2}$$
s.t. 
$$Fx - f = 0$$

The closed form solution may be found using Lagrange multipliers, and is given by

$$x^* = F^T (FF^T)^{-1} (f - Fw) + w,$$

Since F is sparse, the main difficulty in computing  $x^*$  is overcome if we can efficiently compute  $(FF^T)^{-1}q$  for any vector q.

**Lemma 3.** Let F be the sequence form constraint matrix. Computing  $(FF^T)^{-1}q$  may be done efficiently.

*Proof.* The key here is to exploit the structure of  $FF^T$ . Observe that  $FF^T$  is symmetric, positive-definite and has dimension equal to the number of information sets. Furthermore,  $FF^T$  may be expressed in closed form:

$$(FF^T)_{ij} = \begin{cases} -1 & i \text{ is the direct parent/child of } j \\ 1, & i \text{ is the sibling of } j \\ 1+\# \text{ actions at } i, & i = j \\ 0, & \text{otherwise} \end{cases},$$

where i,j above are information sets, and i being the parent of j means that there is some action in i which can lead to information set j (without any other information set from the same player in between), and i being the sibling of j means that the (unique) sequence leading to i and j are the same. Observe that  $FF^T$  is almost, but not quite tree-structured. However, it is sparse and more importantly, has 0 fill-in if we order variables in a bottom-up fashion in the player's game tree. That is, we treat  $FF^T$  as a graph with information sets as vertices, then repeatedly remove vertices (information sets) in a bottom-up fashion, while forming cliques with all neighbors of the removed vertex. Note that due to the structure of  $FF^T$ , we will not introduce any new edges. In other words, performing Gaussian elimination on  $(FF^T)$  may be done without introducing additional non-zero entries. If the number of maximum number of actions that an information set may have is upper bounded by a constant  $a_{\max}$ , then eliminating a single variable will only require time quadratic in  $a_{\max}$ . This means that computing  $(FF^T)^{-1}q$  may be done efficiently when  $x_{\max}$  is small.

**Remark.** Lemma 3 and the fact that L2 projections onto sequence form constraints can be done efficiently may be of separate interest to researchers beyond the scope of EFCEs.

In practice, we precompute a sparse Cholesky factor of  $FF^T$ . From the previous discussion, the Cholesky factors are guaranteed to be sparse and easily stored. Withe the Cholesky decomposition of  $FF^T$ , finding  $(FF^T)^{-1}q$  becomes straightforward. This precomputation is done once per trigger-sequence  $\hat{\sigma}$ , since the set of relevant sequence pairs for each trigger sequence (i.e., the location of non-zero entries in the matrix representing  $\xi$ ) differs. This precomputation step is trivially parallel. In our experiments, computing the Cholesky factors was rarely the bottleneck (although we do include this timing when evaluating our method)

#### **D.2. Social Welfare Maximization**

Observe that Equation (1) may be rewritten in the form of  $(1 - \lambda_{sw})v_{sw}^*(\xi) - \lambda_{sw}b^T\xi$  for a suitable vector b. Hence, the gradient for the modified objective is given by

$$\partial/\partial\xi\;v^*(\xi) = \begin{cases} A_{i^*,\hat{\sigma}^*}y_{i^*,\hat{\sigma}^*}^* - b_{i,\hat{\sigma}^*} & v^*(\xi) \geq \kappa(\xi) \\ -b & \text{otherwise} \end{cases},$$

 $<sup>^{4}</sup>$ In our implementation, f need not have this restriction, but it is included her to be more in line with the classic work of (von Stengel, 1996).

where  $\kappa(\xi) = \tau - \sum_{z \in Z} u_1(z) \xi_1(z;z) - \sum_{z \in Z} u_2(z) \xi_2(z;z)$ , the difference between  $\tau$  and the social welfare obtained from  $\xi$ .

# E. Battleship Game

#### E.1. Extended Description of the Game

A game of Battleship is parameterized by a tuple  $(H, W, S, r, \gamma)$ , where

- the integers  $H, W \ge 1$  define the height and width of the playing field for each player;
- S is an ordered list containing ship descriptions  $s_i$  for each player. Each description is a pair  $s_i = (\ell_i, v_i)$ , where  $\ell_i$  is the length of the *i*-th ship and  $v_i$  is its value;
- $r \ge 1$  is the number of rounds in the game;
- $\gamma \geq 1$  is a *loss multiplier* that controls the relative value of a losing versus destroying ships.

The game proceeds in 2 phases: *ship placement* and *shooting*. During the ship placement phase, the players (starting with Player 1) take turns placing their ships on their playing field. The players must place all their ships, in the same order in which they appear in S, on the playing field. The ship placement phase ends when all ships have been placed. We remark that the players' playing fields are separate: in other words, there are two playing fields of dimensions  $H \times W$ , one per player. The ships may be placed either horizontally or vertically on each player's grid (playing field); all ships must lie entirely within the playing field and may not overlap with other ships the player has already placed. Finally, the locations of a player's ships is private information for each player.

In the shooting phase, players take turns firing at each other; Player 1 starts first. This is done by selecting a pair of integer coordinates (x,y) that identify a cell within the playing field. After taking a shot, the player is told if the shot was a hit, that is, the selected cell (x,y) is occupied by a ship of the opponent, or if it is a miss, that is, (x,y) does not contain an opponent's ship. If all cells covered by a ship have been shot at, the ship is destroyed and this fact is announced. Note that the identity of the ship which was hit or sunk is not revealed; players only know that some ships was hit or sunk. The game ends when r shots have been made by each player, or if one player has lost all their ships, whichever comes first. At the end of the game, each player's payoff is computed as follows: for each opponent's ship that the player has destroyed, the player receives a payoff equal to the value v of that ships; for each ship that the player has lost to the opponent, the player incurs a negative payoff equal to v0 that is the value of the ship times the loss multiplier v1. Note that when v1 the game is general sum.

Since  $\gamma \geq 1$ , this asymmetric model describes situations where players are encouraged to destroy other ships, but are ultimately more protective of their own assets. The loss multiplier  $\gamma$  governs this gap; a higher value of  $\gamma$  makes so that each player values their ships more than destroying others. Note that when  $\gamma = 1$ , we obtain a zero-sum version of battleships (with varying scores for each ship).

For the remainder of the discussion, we define the *social welfare* (SW) of any outcome to be the sum of payoffs of each player. We will demonstrate that with the aid of a mediator (the correlation device), the social welfare of the optimal correlated equilibria are dramatically higher than the social welfare of even the best Nash equilibrium. In other words, the mediator leads to significantly less destructive outcomes, and leads to more frequent ties where the players sometimes agree to deliberately miss their opponents, while still retaining incentive-compatibility and rationality in the standard game-theoretic sense.

## E.2. Analysis of Social-Welfare-Maximizing EFCE

We analyze one social-welfare-maximizing EFCE in the same small instance of Battleship as the previous section. The mediator in this EFCE recommends the players a ship placement that is sampled uniformly at random and independently for each players. This results in 9 possible scenarios (one per possible ship placement) in the game, each occurring with probability 1/9. Due to the symmetric nature of ship

placements, only two scenarios are relevant: whether the two players are recommended to place their ship in the same spot, or in different spots. Figure 3 details the strategy of the mediator in each of these two scenarios, assuming that the players do not deviate. Note that the game trees in Figure 3 are parametric on the recommended ship placements a and b; all 9 possible ship placements can be recovered from Figure 3 by setting a and b to appropriate values in  $\{1, 2, 3\}$ .

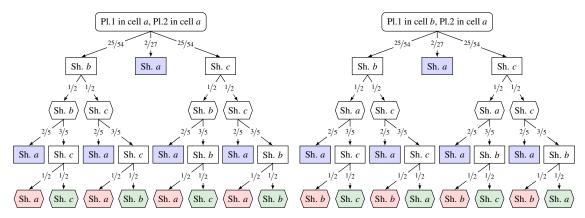


Figure 3: Example of a playthrough of Battleship assuming both players were recommended to place their ship in a (left), or that Player 1 and 2 were recommended to place their ships in a and b respectively (right). For both pictures, the numbers along each edge denote probabilities of each action being recommended; no edge is shown for actions recommended with zero probability. Squares and hexagons denote actions taken by Players 1 and 2 respectively. Similarly, blue and red nodes represent cases where Players 1 and 2 sink their opponent's ship, respectively. Green leaf nodes are where the game results in no ship loss. The *Shoot* action is abbreviated to 'Sh.'

For both game trees, note that the correlation device suggests that Player 1 shoot at the Player 2's ship with a low  $^2/^2$ 7 probability, and deliberately miss with high probability. As hinted in earlier sections, this type of recommendation is key to understanding why the EFCE succeeds in promoting less destructive outcomes. One may wonder why this behavior is incentive-compatible (that is, what are the incentives that compel Player 1 into not defecting), since the player may choose to randomly fire in any of the 2 locations that were *not* recommended, and get almost  $^1/^2$ 2 chance of winning the game immediately. The key is that if Player 1 does so and does not hit the opponent's ship, then the mediator can punish him by recommending that Player 2 shoot at the location of Player 1's ship. Since players value their ships more than destroying their opponents, the player is incentivized to avoid such a situation by accepting the recommendation to (most probably) miss.

A similar situation arises in the first move of Player 2. Here, Player 2 is recommended to *deliberately* miss, hitting each of the 2 empty spots with probability  $^{1}/_{2}$ . If he deviates and attempts to destroy Player 1's ship, then he risks the mediator revealing his location to his opponent if his shot misses; this risk is enough to keep Player 2 'in line'. The second move of Player 1 (third shot of the full game) bears a similar ideas. Here, Player 1 is recommended to hit Player 2's ship with probability  $^{2}/_{5}$ . Similar to his first shot, Player 1 may deviate and fire at the remaining location and enjoy  $^{3}/_{5}$  chance of winning the game out right. Yet, this behavior is discouraged, since in the  $^{2}/_{5}$  chance that he misses the shot (i.e., the recommendation was in fact, the correct location of Player 2's ship), then his location would be revealed by the mediator and he loses the next round. Again, this threat from the mediator encourages peaceful behavior, even though the recommendation to Player 1 reveals a more accurate 'posterior' of Player 2's ship location, as compared to the uniform distribution of  $^{1}/_{2}$ . While making these recommendations, the mediator ensures that Player 2 has a uniform distribution of Player 1's ship location, meaning that even though Player 2 has the final move, he may not do better than guessing at uniform at this stage.

**Remark.** It is important to note that Figure 3 does not convey the full information of the correlated plans. Crucially, it does not show the consequences suffered if a player deviates from his recommended strategy—in this case, the deviating player stops receiving recommendations and risks having his ship's

location revealed to the opponent. These 'counterfactual' scenarios may be counter-intuitive but are key to understanding how SW-maximizing EFCEs achieve their purpose.

## F. Sheriff Game

## F.1. Extended Version of the Game

The Sheriff game is described by the the parameters  $v, p, s \in \mathbb{R}^+, n_{\max}, b_{\max}, r \in \mathbb{N}$ . The parameters  $v, p, s \geq 0$  describe the value of *each* illegal item, the penalty that the Smuggler has to pay for *each* discovered illegal item, and the compensation that the Sheriff pays to the Smuggler in the case of a false alarm. At the beginning of the game, the Smuggler loads  $n \in \{0, \dots, n_{\max}\}$  items into his cargo. The amount of goods loaded is unknown to the Sheriff. The game then proceeds for  $r \geq 1$  rounds of bargaining. Each round comprises two steps. First, the Smuggler offers a bribe  $b_t \in \{0, \dots, b_{\max}\}$  to the Sheriff, where  $t \leq r$  is the round of bargaining. After that, the Sheriff responds with 'Yes' or 'No'.

All actions are public knowledge, except for the selection of cargo contents, which only the Smuggler knows. In the final step, we compute the payoffs to players. The outcome of the game is decided by the *last* step of bargaining. In particular, the first r-1 rounds of bargaining have no explicit bearing on the outcome of the game, except for purposes of coordination. The payoffs for each outcome are:

- 1. Sheriff accepts the bribe. The Smuggler's gets  $n \cdot v b_r$ , and the Sheriff's gets the bribe offered  $b_r$ .
- 2. Sheriff inspects and discovers illegal items. The Smuggler is fined and gets a payoff of  $-n \cdot p$  while the Sheriff gets a payoff of  $n \cdot p$ .
- 3. Sheriff chooses to inspect and does not find illegal items. The Smuggler receives a compensation of s, while the Sheriff gets -s.

The objective of the mediator is to maximize social welfare in the space EFCEs. Ideally, this will involve the Smuggler bringing in goods and the Sheriff accepting bribes – any other outcome would simply be zero-sum, since it no goods will be successfully smuggled and money only changes hands between players. A qualitative description of the welfare maximizing equilibrium is not obvious, since the game contains elements of both lying and bargaining.

**Remark.** The communication in the bargaining steps is at a superficial level similar to that in *cheap talk* (Crawford & Sobel, 1982), where costless and non-binding signals are transmitted between players. However, in our setting, the signals are transmitted in the middle of the game as opposed to just at the beginning. More importantly, the presence of the mediator during the phase of bargaining bestows more uses for the signals—in particular, the mediator may be able to take punitive measures against players who deviate from recommendations, since future recommendations will be withheld from players who deviate. We will illustrate the importance of this at the end of Appendix F.2.

## F.2. Effect of Additional Rounds of Bargaining (r)

| $n_{max}$ | r=1          | r = 2        | r = 3        | r = 4        |
|-----------|--------------|--------------|--------------|--------------|
| 1         | (4.00, 1.00) | (4.00, 1.00) | (4.00, 1.00) | (4.00, 1.00) |
| 2         | (1.24, 0.19) | (4.00, 1.00) | (4.00, 1.00) | (4.00, 1.00) |
| 5         | (0.89, 0.11) | (1.11, 1.00) | (4.00, 1.00) | (4.00, 1.00) |
| 10        | (0.82, 0.00) | (0.84, 1.00) | (3.62, 1.00) | (4.00, 1.00) |

Table 3: Payoffs for (Smuggler, Sheriff) when players play according to the SW-maximizing EFCE in the Sheriff game with  $b_{max} = 2$  (right).

We illustrate the effect of the non-consequential bribes with two small settings, where  $v=5, p=1, s=1, n_{\text{max}}=3, b_{\text{max}}=2, r\in\{1,2\}$ . Examples of SW-maximizing equilibria are shown in Figure 4 and Figure 5. <sup>5</sup>

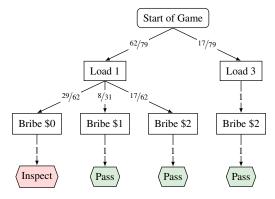


Figure 4: Example of a playthrough of the Sheriff game with r=1. Edge labels correspond to action probabilities, edges with 0 probability are omitted. Squares and hexagons denote actions taken by Players 1 and 2 respectively, while green and red nodes denote the Sheriff choosing to pass or inspect.

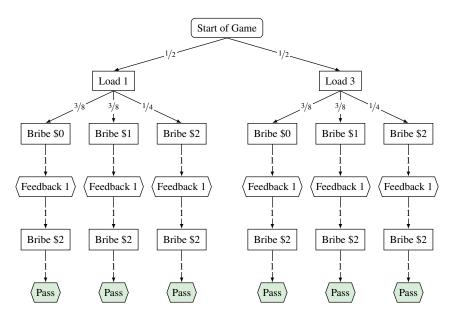


Figure 5: Example of a playthrough of the Sheriff game with r=2. Edge labels correspond to action probabilities, edges with 0 probability are omitted. Squares and hexagons denote actions taken by Players 1 and 2 respectively, while green and red nodes denote the Sheriff choosing to pass or inspect.

The SW maximizing EFCE yields payoffs of (3.89, 1.43) and (8.0, 2.0) for r=1 and r=2 respectively. We will first consider the case where r=2 (Figure 5. Here, what occurs happens along the equilibrium path is straightforward. The Smuggler loads in 1 or 3 items with equal probability. Next, he offers a (non-consequential) bribe of either 0, 1, or 1. Then, he receives some feedback of 1, and proceeds to offer a bribe of 1, which the Sheriff gladly accepts. The payoffs to players is 1, 1, and 1, a

The underlying mechanism is in fact fairly straightforward and mirrors the idea in the modified signalling game of (von Stengel & Forges, 2008). Assume that a random number is chosen uniformly from  $\{0, \ldots, b_{\text{max}}\}$ . This acts as a 'passcode' which the Sheriff expects from the Smuggler in the first round.

<sup>&</sup>lt;sup>5</sup>As with the analysis of Battleship, note that this only shows interactions of players on the equilibrium path, that is, the graph omits what would happen if some player deviated.

This passcode forms part of the correlated plan, and will eventually be revealed to the Smuggler assuming he did not deviate when selecting the number of illegal items (recall that the sequential nature of the EFCE means that the recommended amount to bribe is not revealed until the Smuggler loads the cargo with the recommended number of items.) In other words, the first (non-consequential) bribe may be used as a signal which hints to the Sheriff if the Smuggler has deviated—if it is not equal to the passcode, the Smuggler must have deviated somewhere. On the other hand, a deviating Smuggler may successfully guess the passcode with probability no greater than  $1/(b_{\text{max}}+1)$ ; if the number of signals  $b_{\text{max}}$  is sufficiently large, then it is near impossible to guess the code. Using these tools, the mediator is able to engineer a 'deviation detector' which checks if the Smuggler ever deviated. Note, however, that unlike the Signaling game, the Sheriff is not able to glean exactly how what was recommended (in this case, the number of items in the cargo); he is only able to deduce if the player deviated from the recommendation (in this case, this would be load either 1 or 3 items).

Issuing threats to the Smuggler becomes straightforward with this deviation detector. If the Sheriff knows the Smuggler is lying, he employs a 'grim trigger' for the rest of the game—in this case, the Sheriff opts to inspect all of the player's cargo, regardless of the bribe offered in the second round. The Smuggler could also be pretending to bring in illegal goods, i.e., by loading 0 items and hoping that he would guess the *incorrect* passcode, resulting in the Sheriff making a false accusation. However, because the Smuggler's payoff for deceiving the Sheriff in this manner is just 1, he remains incentivized to stick to the recommendations, which guarantees him a payoff of either 3 of 13.

We now make the following hypotheses. First, the effect of additional bargaining rounds r is that the chance of randomly guessing the passcode is reduced. If there are r rounds, then there are  $(b+1)^{r-1}$  different possible signals that the Smuggler could have sent to the Sheriff through the first r-1 rounds. When r=1, this class of correlation plans fails since the bribe by the Smuggler serves both as the answer to the 'secret question' and as the actual bribe to be offered. This aliasing of roles is what leads to a lower payoff; the risk of sending an incorrect passcode is not sufficiently high to dissuade the Smuggler from deviating.

# G. 3 LP formulations for computing EFCEs

Refer to the dualized problem in Appendix A. Observe that u is the value of the maximum deviation over all  $\hat{\sigma}$ —when all incentive constraints are met, u should be non-positive. We propose 3 different formulations.

- Min-Deviation: what was presented in Appendix A.
- Feas-Deviation: instead of minimizing u in the objective, replace that by a hard constraint that  $u \leq 0$ .
- Maximum-SW: formulate the LP similar to Feas-Deviation, but with the SW-maximizing objective.

## H. Additional Experiments on Sheriff Game

The results for the Sheriff game were run using the parameters  $p=1, v=5, b_{\max}=3, r=5, s=4$  while varying the maximum number of items that can be smuggled  $n_{\max}$ . The time required for the error to drop below a certain threshold is reported for both Gurobi and our subgradient method. The results are reported in Table 4. As before, we observe that our method can outperforms Gurobi if lower levels of accuracy are

|                    | #Actions |      | #Relevant    | Time (LP)<br>2 1 0.75 0.5 |      |      |      | Time (ours) |     |      |     |
|--------------------|----------|------|--------------|---------------------------|------|------|------|-------------|-----|------|-----|
| $n_{\mathrm{max}}$ | Pl 1     | Pl 2 | seq. pairs   | 2                         | 1    | 0.75 | 0.5  | 2           | 1   | 0.75 | 0.5 |
| 6                  | 131k     | 37k  | 6.5M<br>8.4M | 723                       | 723  | 723  | 743  | 42          | 42  | 44   | 88  |
| 8                  | 168k     | 37k  | 8.4M         | 1187                      | 1223 | 1223 | 1223 | 63          | 69  | 102  | 333 |
| 10                 | 206k     | 37k  | 10 <b>M</b>  | 1662s                     | 1774 | 1774 | 1829 | 83          | 240 | 298  | N/A |

Table 4: #Seq. pairs is the dimension of  $\xi$  under the compact representation of (von Stengel & Forges, 2008). For LPs, we report the fastest of Barrier, Primal and Dual Simplex, and 3 different formulations (Appendix G). Our subgradient method did not manage to achieve an accuracy of 0.5 after 1 hour of running.

desired. However, it was observed that for higher levels of accuracy, Gurobi requires significantly less

time, if our method converges at all. This is because Gurobi spends the majority of its time performing preprocessing steps.