

# Sequential Rationality in Cryptographic Protocols

Ronen Gradwohl  
Kellogg School of Management  
Northwestern University  
Evanston, IL, USA  
Email: r-gradwohl@kellogg.northwestern.edu

Noam Livne  
Department of CS and Applied Mathematics  
Weizmann Institute of Science  
Rehovot, Israel  
Email: noam.livne@weizmann.ac.il

Alon Rosen  
School of Computer Science  
IDC Herzliya  
Herzliya, Israel  
Email: alon.rosen@idc.ac.il

**Abstract**—Much of the literature on rational cryptography focuses on analyzing the strategic properties of cryptographic protocols. However, due to the presence of computationally-bounded players and the asymptotic nature of cryptographic security, a definition of sequential rationality for this setting has thus far eluded researchers.

We propose a new framework for overcoming these obstacles, and provide the first definitions of computational solution concepts that guarantee sequential rationality. We argue that natural computational variants of subgame perfection are too strong for cryptographic protocols. As an alternative, we introduce a weakening called threat-free Nash equilibrium that is more permissive but still eliminates the undesirable “empty threats” of non-sequential solution concepts.

To demonstrate the applicability of our framework, we revisit the problem of implementing a mediator for correlated equilibria (Dodis-Halevi-Rabin, Crypto’00), and propose a variant of their protocol that is sequentially rational for a non-trivial class of correlated equilibria. Our treatment provides a better understanding of the conditions under which mediators in a correlated equilibrium can be replaced by a stable protocol.

**Index Terms**—game theory; cryptography;

## I. INTRODUCTION

A recent line of research has considered replacing the traditional cryptographic modeling of adversaries with a game-theoretic one. Rather than assuming arbitrary *malicious* behavior, participants are viewed as being self-interested, *rational* entities that wish to maximize their own profit, and that would deviate from a protocol’s prescribed instructions if and only if it is in their best interest to do so.

Such game theoretic modeling is expected to facilitate the task of protocol design, since rational behavior may be easier to handle than malicious behavior. It also has the advantage of being more realistic in that it does not assume that some of the parties honestly follow the protocol’s instructions, as is frequently done in cryptography.

The interplay between cryptography and game theory can also be beneficial to the latter. For instance, using tools from secure computation, it has been shown how to transform games in the mediated model into games in the unmediated model.

But regardless of whether one analyzes cryptographic protocols from a game theoretic perspective or whether one uses protocols to enhance game theory, it is clear that the results are meaningful only if one provides an adequate framework for such analyses.

## A. Computational Nash Equilibrium

Applying game-theoretic reasoning in a cryptographic context consists of modeling interaction as a *game*, and designing a protocol that is in *equilibrium*. The game specifies the model of interaction, as well as the utilities of the various players as a function of the game’s outcome. The protocol lays out a specific plan of action for each player, with the goal of realizing some pre-specified task. Once a protocol has been shown to be in equilibrium, rational players are expected to follow it, thus reaching the desired outcome.

A key difficulty in applying game-theoretic reasoning to the analysis of cryptographic protocols stems from the latter’s use of computational infeasibility. Whereas game theory places no bounds on the computational ability of players, in cryptography it is typically assumed that players are computationally bounded. Thus, in order to retain the meaningfulness of cryptographic protocols, it is imperative to restrict the set of strategies that are available to protocol participants. This gives rise to a natural analog of Nash equilibrium (NE), referred to as *computational Nash equilibrium* (CNE): any polynomial-time computable deviation of a player from the specified protocol can improve her utility by only a negligible amount (assuming other players stick to the prescribed strategy).

Consider, for example, the following (two-stage, zero-sum) game (related to a game studied by Ben-Sasson et al. [3] and Fortnow and Santhanam [6]), which postulates the existence of a one-way permutation  $f : \{0, 1\}^n \mapsto \{0, 1\}^n$ .

**Example 1.1:** (One-way permutation game):

- 1)  $P_1$  chooses some  $x \in \{0, 1\}^n$ , and sends  $f(x)$ .
- 2)  $P_2$  sends a message  $z \in \{0, 1\}^n$ .
- 3)  $P_2$  wins (gets payoff 1) if  $z = x$  (and gets -1 otherwise).

In classical game theory, in all NE of this game  $P_2$  wins, since there always exists some  $z$  such that  $z = x$ . However, in the computational setting, the following is a CNE: both players choose their messages uniformly at random (resulting in an expected loss for  $P_2$ ). This is true because if  $P_2$  chooses  $z$  at random, then  $P_1$  can never improve his payoff by not choosing at random. If  $P_1$  chooses  $x$  at random, then by the definition of a one-way permutation, any computationally-bounded strategy  $\sigma_2$  of  $P_2$  will be able to guess the value of  $x$  with at most negligible (in  $n$ ) probability. Thus, the expected utility of  $P_2$  using  $\sigma_2$  is negligible, and so he loses at most that much by sticking to his CNE strategy (i.e. picking some  $z$  at random).

## B. Computational Subgame Perfection

The notion of CNE serves as a first stepping stone towards a game-theoretic treatment of cryptographic protocols. However, protocols are typically *interactive*, and CNE does not take their sequential nature into consideration.

In traditional game theory interaction is modeled via extensive games. The most basic equilibrium notion in this setting is *subgame perfect equilibrium* (SPE), which requires players' strategies to be in NE at any point of the interaction, regardless of the history of prior actions taken by other players. Basically, this ensures that players will not reconsider their actions as a result of reaching certain histories (a.k.a. "empty threats").

As already noted in previous works (cf. [16], [19], [23]), it is not at all clear how to adapt SPE to the computational setting. A natural approach would be to require the strategies to be CNE at every possible history. However, if we condition on the history, then this means that *different* machines can and will do much better than the prescribed equilibrium strategy. For example, in the one-way permutation game of Example 1.1, given any message history, a machine  $M$  can simply have the correct inverse hardwired.

Although this requirement can be relaxed to ask that the prescribed strategy should be better than any other fixed machine on all inputs, this again may be too strong, since a fixed machine can always do better on some histories. Therefore, it seems that we must accept the following: for any machine  $M$ , with *high probability* over possible message histories, the prescribed strategy does at least as well as  $M$ . However, it turns out that this approach also fails to capture our intuitive understanding of a computational SPE (CSPE). Consider the following (two-stage) variant of the one-way permutation game from Example 1.1:

**Example 1.2:** (Modified one-way permutation game):

- 1)  $P_1$  chooses some  $x \in \{0, 1\}^n$ , and sends  $f(x)$ .
- 2)  $P_2$  sends a message  $z \in \{0, 1\}^n$ .
- 3) If exactly one of  $P_1$  and  $P_2$  send message 0, both players get payoff  $-2$ . If both players send message 0, both players get payoff  $+2$ . Otherwise,  $P_2$  wins (with payoff  $+1$ ) if and only if  $z = x$ , and the non-winning player loses (with payoff  $-1$ ).

Using a similar argument to the one applied in Section I-A, it can be shown that the strategies in which both players choose a message uniformly at random from  $\{0, 1\}^n \setminus \{0\}$  satisfy the above "probabilistic" variant of CSPE. However, this equilibrium does not match our intuitive understanding of SPE:  $P_1$  will prefer to send message 0 regardless of  $P_2$ 's strategy, knowing that  $P_2$  will then respond with 0 as well. The threat of playing uniformly from all other messages is empty, and hence should not be admitted by the definition.<sup>1</sup>

The examples above are rather simple, so it is reasonable to expect that issues arising in their analyses are inherent in many other cryptographic protocols. This raises the question

<sup>1</sup>We note that a simple change to the payoffs yields a game whose empty threat is more "typical": For the case in which both players send message 0, let  $P_2$ 's payoff be  $-3/2$ .

of whether a computational variant of SPE is at all attainable in a cryptographic setting.

At the heart of this question is the fact that essentially any cryptographic protocol carries some small (but positive) probability of being broken. This means that, while there may be a polynomial-time TM that can "perform well" on the *average* message history, there is no single TM that will do better than *all* other TMs on every history (as for any history there exists some TM that has the corresponding "secret information" hardwired).

This state of affairs calls for an alternative approach. While such an approach should be meaningful enough to express strategic considerations in an interactive setting, it should also be sufficiently weak to be realizable. As demonstrated above, any approach for tackling this challenge should explicitly address the associated probability of error and take asymptotics into consideration.

## II. OUR RESULTS

We propose a new framework for guaranteeing sequential rationality in a computational setting. Our starting point is a weakening of subgame perfection, called *threat-free Nash equilibrium*, that is more permissive, but still eliminates the undesirable empty threats of non-sequential solution concepts.

To cast our new solution concept into the computational setting, we develop a methodology that enables us to "translate" arguments that involve computational infeasibility into a purely game theoretic language. This translation enables us to argue about game theoretic concepts directly, abstracting away complications that are related to computation.

In order to demonstrate the applicability of our framework, we revisit the problem of implementing a mediator for correlated equilibria [5], and propose a protocol that is sequentially rational for a non-trivial class of correlated equilibria (see Section II-C for details).

We emphasize that this version of the paper is a summary of [11], and we strongly recommend that the interested reader turn to the latter for proofs and more elaborate discussions.

### A. Threat-Free Nash Equilibria

We introduce *threat-free Nash equilibria* (TFNE), a weakening of subgame perfection whose objective is to capture strategic considerations in an interactive setting. Loosely speaking, a pair of strategies in an extensive game is a TFNE if it is a NE, and if in addition no player is facing an empty threat at any history.

The problem of empty threats is the following: in a NE of an extensive game, it is possible that a player plays sub-optimally at a history that is reached with probability 0. The other player may strategically choose to deviate from his prescribed strategy and arrive at that history, knowing that this will cause the first player to play an optimal response rather than the prescribed one. In an SPE this problem is eliminated by requiring that no player can play sub-optimally at any history, and so no other player will strategically deviate and take advantage of this.

The main observation leading to the definition of TFNE is that the above requirement may be too strong a condition to eliminate such instability: if an optimal response of a player *decreases* the utility of the other, then this other player would not want to strategically deviate. By explicitly ruling out this possibility, the instability caused by empty threats is eliminated, despite the equilibrium notion being more permissive than subgame perfection.

To make this precise, we give the first formal definition of an empty threat in extensive games. The definition is regressive: Roughly speaking, a player  $i$  is facing a threat at a history if there is some deviation at that history, along with a threat-free continuation from that history onwards, so that  $i$  increases his overall expected payoff when the players play this new deviation and continuation.

We note that the notion of TFNE is strong enough to eliminate the undesirable strategy of playing randomly in the modified OWP game from Example 1.2 – in the full version [11] we show that in any computational TFNE of this game the second player outputs 0 after history 0.

### B. Strategy-Filters and Tractable Strategies

To cast the definition of TFNE into a computational setting, we map the given protocol into a sequence of extensive games using *strategy-filters* that map computable strategies into their “strategic representation” (the strategic representation corresponds to the strategy effectively played by a given interactive Turing machine). We can then apply pure game theoretic solution concepts, and in particular our newly introduced concept of TFNE, to understand the strategic behavior of players.

Similarly to the definition of CNE, the computational treatment departs from the traditional game theoretic treatment in two crucial ways. First of all, our definition is framed *asymptotically* (in order to capture computational infeasibility), whereas traditional game-theory is framed for finitely sized games. Second, it allows for a certain *error probability*. This is an artifact of the (typically negligible) probability with which the security of essentially any cryptographic scheme can be broken.

Given a cryptographic protocol, we consider a corresponding sequence of extensive games. The sequence is indexed by a security parameter  $k$  and an error parameter  $\varepsilon$ . For each game, we “constrain” the strategies available to players to be a subset of those that can be generated by PPT players in the protocol. Intuitively, the game indexed by  $(k, \varepsilon)$  contains those strategies that run in time polynomial in  $k$  and “break crypto” with probability at most  $\varepsilon$ . We also require that strategy-filters be *PPT-covering*: that for any polynomially-small  $\varepsilon$ , every PPT is eventually a legal strategy, far enough into the sequence of extensive games.

Using this framework we formalize computational threat-free Nash equilibrium (CTFNE). To the best of our knowledge this is the first attempt at analyzing sequential strategic reasoning in the presence of computational infeasibility.

### C. Applications

Our treatment provides a powerful tool for arguing about the strategic behavior of players in a cryptographic protocol. It also enables us to isolate sequential strategic considerations that are suitable for use in cryptographic protocols (so that the solution concept is not too weak and not too strong).

We revisit the general problem of implementing a mediator for correlated equilibria [5], and propose a protocol that is sequentially rational for a non-trivial class of correlated equilibria. In particular, our protocol is in a CTFNE for correlated equilibria that are convex combinations of Nash equilibria and that are “undominated”: There does not exist any convex combination of Nash equilibria for which both players get a strictly higher expected payoff.

Our treatment explores the conditions under which mediators in a correlated equilibrium can be replaced by a stable protocol, and sheds light on some structural properties of such equilibria.

Finally, we prove a general theorem that identifies sufficient conditions for a TFNE in extensive games. Namely, we show that if an undominated NE has the additional property that no player can harm the other by a unilateral deviation, then that NE must also be threat-free.

### D. Related Work

This paper contributes to the growing literature on rational cryptography. Many of the papers in this line of research, such as [5], [14], [15], [1], [9], [20], [21], [16], [18], [19], [17], [23], [22], [2], [10], explore various solution concepts for cryptographic protocols viewed as games (often in the context of rational secret-sharing). Aside from the works of Lepinski et al. [15], [20], Ong et al. [23], and Gradwohl [10], who work in a different model<sup>2</sup>, all prior literature has considered solution concepts that are non-sequential. More specifically, they all use variants of NE such as strict NE, NE with stability to trembles, and everlasting equilibrium.

An additional related work is that of Halpern and Pass [13], in which the authors present a general framework for game theory in a setting with computational cost. While their approach to computational limitations is more general than ours, they only address NE. Finally, Fortnow and Santhanam [6] study a different framework for games with computational limits, but also only in the context of NE.

### E. Future Work

One potential application of our new definition is an analysis of rational secret-sharing protocols. For some ideas about why known gradual release protocols satisfy a solution concept that is related to but slightly weaker than CTFNE, see the full version of this paper [11].

There are numerous other compelling problems left for future work. The first problem is to extend our definition to games with simultaneous moves. While we do offer a

<sup>2</sup>More specifically, [15], [20] make strong physical assumptions, [23] assume the existence of a fraction of honest (non-rational) players, and [23], [10] work in an information-theoretic setting.

partial extension tailored to the problem of implementing a mediator, the problem of defining CTFNE for general games with simultaneous moves is open. Such a definition would be particularly useful for a sequential analysis of protocols with a simultaneous channel. Another natural extension of the definition is to multiple players, as opposed to 2. Such an extension comes with its own challenges, particularly with regard to the possibility of collusion. A third extension is to incorporate the threat-freeness property with stronger variants of NE, such as stability with respect to trembles, strict NE, or survival of iterated elimination of dominated strategies. Finally, we would like to find more applications for our definition. One particularly interesting problem is to extend our results on the implementation of mediators to a larger class of correlated equilibria.

### III. GAME THEORY DEFINITIONS

#### A. Extensive Games

Informally, a game in extensive form can be described as a game tree in which each node is owned by some player and edges are labeled by legal actions. The game begins at the root, and at each step follows the edge labeled by the action chosen by the current node's owner. Utilities of players are given at the leaves of the tree. More formally, we have the following standard definition of extensive games (see, for example, Osborne and Rubinstein [24]):

**Definition 3.1 (Extensive game):** A 2-person extensive game is a tuple  $\Gamma = (H, P, A, u)$  where

- $H$  is a set of (finite) history sequences such that the empty word  $\epsilon \in H$ . A history  $h \in H$  is terminal if  $\{a : (h, a) \in H\} = \emptyset$ . The set of terminal histories is denoted  $Z$ .
- $P : (H \setminus Z) \rightarrow \{1, 2\}$  is a function that assigns a “next” player to every non-terminal history.
- $A$  is a function that, for every non-terminal history  $h \in H \setminus Z$ , assigns a finite set  $A(h) = \{a : (h, a) \in H\}$  of available actions to player  $P(h)$ .
- $u = (u_1, u_2)$  is a pair of payoff functions  $u_i : Z \mapsto \mathbb{R}$ .

**Definition 3.2 (Behavioral strategy):** Behavioral strategies of players in an extensive game are collections  $\sigma_i = (\sigma_i(h))_{h:P(h)=i}$  of independent probability measures, where  $\sigma_i(h)$  is a probability measure over  $A(h)$ .

For any extensive game  $\Gamma = (H, P, A, u)$ , any player  $i$ , and any history  $h$  satisfying  $P(h) = i$ , we denote by  $\Sigma_i(h)$  the set of all probability measures over  $A(h)$ . We denote by  $\Sigma_i$  the set of all strategies  $\sigma_i$  of player  $i$  in  $\Gamma$ . For each profile  $\sigma = (\sigma_1, \sigma_2)$  of strategies, define the outcome  $O(\sigma)$  to be the probability distribution over terminal histories that results when each player  $i$  follows strategy  $\sigma_i$ .

#### B. Nash Equilibrium

Each profile of strategies yields a distribution over outcomes, and we are interested in profiles that guarantee the players some sort of optimal outcomes. There are many solution concepts that capture various meanings of “optimal”, and one of the most basic is the Nash equilibrium (NE).

**Definition 3.3 (Nash equilibrium (NE)):** An  $\varepsilon$ -Nash equilibrium of an extensive game  $\Gamma = (H, P, A, u)$  is a profile  $\sigma^*$  of strategies such that for each player  $i$ ,

$$E[u_i(O(\sigma^*))] \geq E[u_i(O(\sigma_{-i}^*, \sigma_i))] - \varepsilon$$

for every strategy  $\sigma_i$  of player  $i$ . It is a NE if the above holds for  $\varepsilon \leq 0$  and a strict NE if it holds for some  $\varepsilon < 0$ .

One of the premises behind the stability of profiles that are in an  $\varepsilon$ -NE is that players will not bother to deviate for a mere gain of  $\varepsilon$ . For applications in cryptography we will generally have  $\varepsilon$  be some negligible function, and this corresponds to our understanding that we do not care about negligible gains.

#### C. Constrained Games

In the standard game theory literature, where there are no computational constraints on the players, the available strategies  $\sigma_i$  of player  $i$  are all possible collections  $(\sigma_i(h))_{h:P(h)=i}$ , where  $\sigma_i(h)$  is an arbitrary distribution over  $A(h)$ . In our setting, however, we will only consider strategies that can be implemented by computationally bounded ITMs. This requires being able to constrain players' strategies to a strict subset of the possible strategies. Given a pair  $T = (T_1, T_2)$  of such sets we can then define a constrained version of a game, in which only strategies that belong to these sets are considered.

**Definition 3.4 (Constrained game):** Let  $\Gamma = (H, P, A, u)$  be an extensive game, and let  $T = (T_1, T_2)$  where  $T_i \subseteq \bigotimes_{h:P(h)=i} \Sigma_i(h)$  for each  $i \in \{1, 2\}$ . The  $T$ -constrained version of  $\Gamma$  is the game in which the only allowed strategies for player  $i$  belong to  $T_i$ .

This definition enables us to capture restrictions that might result from requiring strategies to be implementable by polynomial time ITMs. NE of constrained games are defined similarly to regular NE, except that players' strategies and deviations must be from the constraint sets.

### IV. THREAT-FREE NASH EQUILIBRIUM

Our starting point is the inadequacy of subgame perfection in capturing sequential rationality in a computational context. As argued in Section I-B, it is unreasonable to require computationally-bounded players to play optimally at every node of a game. In particular, in cryptographic settings this requires breaking the security of the protocol, which is assumed impossible under the computational constraints.

A possible idea might be to require that players “play optimally at every node of the game, under their computational constraints.” However, this idea cannot be interpreted in a sensible way. Computational constraints must be defined “globally,” and thus the notion of playing optimally under some computational constraint on a particular history is senseless. In particular, for any history of some cryptographic protocol, there is a small machine that plays optimally on this specific history *unconditionally* (and breaks “cryptographic challenges” appearing in this history, by having the solutions hardwired). This machine is efficient, and so meets essentially any computational constraint. So, while under computational constraints every machine fails on cryptographic challenges

in most histories, for every history there is a machine that succeeds. We thus assume that a player chooses his machine before the game starts, and cannot change his machine later.

### A. A New Solution Concept

In light of the above discussion, it seems like the solution concept we are looking for has to reconcile between the following seemingly conflicting properties:

- 1) It implies an optimal strategy for the players *under their computational constraints*, which implies *non-optimal* play on certain histories.
- 2) It does not allow empty threats, thus implying “sequential rationality.”

The crucial observation behind our definition is that in order to rule out empty threats, one does not necessarily need to require that players play optimally at *every* node, because not every non-optimal play carries a threat to other players. In fact, in a typical cryptographic protocol, the security of each player is *building* on other players not playing optimally (because playing optimally would mean breaking the security of the protocol). Thus, a player’s “declaration” to play non-optimally does not necessarily carry a threat: the other players may even gain from it. More generally, even in non-cryptographic protocols, at least in 2-player perfect information games, we can use the following observation: in any computational challenge, either a player gains from the other not playing optimally, or, if he does not gain, he can avoid introducing that computational challenge to the other player.

Following the above observation, we introduce a new solution concept for extensive games. The new solution concept requires that players be in NE, and moreover, that no player impose an empty threat on the other. At the same time, it does not require players to play optimally at every node. In other words, players may (declare to) play non-optimally on non-equilibrium support, yet this declaration of non-optimal play does not carry an empty threat. We call our new solution concept TFNE, for threat-free Nash equilibrium.

To make the above precise, we introduce a formal definition of an empty threat. An empty threat occurs when a player threatens to play “non-rationally” on some history in order to coerce the other player to avoid this history. Crucially, empty threats are such that, had the threatened *not* believed the threat, had he deviated accordingly, and had the threatening player played “rationally”, the threatened player would have benefitted. To rephrase our intuition: a player faces an empty threat with respect to some strategy profile if by deviating from his prescribed strategy, and having the other player react “rationally”, he improves his payoff (in comparison with sticking to the prescribed strategy and having the other player react “rationally” from then on.)

But what does it mean for the other player to react “rationally”? The other player may assume, recursively, that the first player will play a best response, and will not carry out empty threats against him, and so on, leading to a regressive definition.

### B. Vanilla Version

Before giving the general definition of TFNE that we will use, we present a simpler version that has no slackness parameter and that works for games without constrained strategies.

For a player  $i$  and a history  $h$ , two strategies  $\sigma_i$  and  $\pi_i$  are *equivalent for player  $i$  on  $h$*  if  $P(h) = i$  and  $\sigma_i(h) = \pi_i(h)$ , or  $P(h) \neq i$ . Two strategies *differ only on the subgame  $h$*  if they are equivalent on every non-terminal history that does not have  $h$  as a prefix. Formally, they are equivalent on every history in  $H \setminus \{h' \in H : h' = h \circ h'' \text{ for some } h''\}$ . For a history  $h \in H$  a strategy  $\sigma$  and a distribution  $\tau = \tau(h)$  on  $A(h)$ , define the set  $\text{Cont}(h, \sigma, \tau)$  as

$$\left\{ \pi : (\pi \text{ differs from } \sigma \text{ only on } h) \ \& \ (\pi(h) = \tau(h)) \right\}.$$

We now proceed to define a threat. For simplicity, we will do so for generic games, in which each player’s possible payoffs are distinct. For such games, the set  $\text{Cont}(h, \sigma, \tau)$  always contains exactly one “threat-free” element (defined below).

**Definition 4.1 (Threat):** *Let  $\Gamma = (H, P, A, u)$  be an extensive game with distinct payoffs. Let  $\sigma$  be a strategy profile, and let  $h \in H$ . Player  $i = P(h)$  is facing a threat at history  $h$  with respect to  $\sigma$  if there exists a distribution  $\tau = \tau(h)$  over  $A(h)$  such that the unique  $\pi \in \text{Cont}(h, \sigma, \tau)$  and  $\pi' \in \text{Cont}(h, \sigma, \sigma)$  that are threat-free on  $h$  satisfy*

$$E[u_i(O(\pi))] > E[u_i(O(\pi'))],$$

where strategy  $\pi$  is said to be threat free on  $h$  if for all  $h' \neq h$  satisfying  $h \circ h' \in H$  player  $P(h \circ h')$  is not facing a threat at  $h \circ h'$  with respect to  $\pi$ .

Note that if  $h$  is such that for all  $a \in A(h)$  it holds that  $h \circ a \in Z$ , then any profile  $\pi$  is threat free on  $h$ .

**Definition 4.2 (Threat-free Nash equilibrium):** *Let  $\Gamma = (H, P, A, u)$  be an extensive game. A strategy profile  $\sigma^*$  is said to be in threat-free Nash equilibrium (TFNE) if:*

- 1)  $\sigma^*$  is a NE of  $\Gamma$ , and
- 2) for any  $h \in H$ , player  $P(h)$  is not facing a threat at history  $h$  with respect to  $\sigma^*$ .

Note that in every profile that is in a TFNE, the effective play matches some SPE profile (more precisely, there is an SPE profile that yields the exact same distribution on outcomes). This and other properties of threats and TFNE are formalized in the companion paper to this work [12].

In the definition of a threat we used the fact that  $\text{Cont}(h, \sigma, \tau)$  and  $\text{Cont}(h, \sigma, \sigma)$  each contain exactly one profile that is threat-free on  $h$ . To show that this must be the case, we have the following proposition (see full version [11] for a proof), which is not unlike the fact that generic games have unique subgame perfect equilibria.

**Proposition 4.3:** *For any extensive game  $\Gamma = (H, P, A, u)$ , strategy profile  $\sigma$ , player  $i$ , history  $h \in H \setminus Z$  with  $P(h) = i$ , and distribution  $\tau$  over  $A(h)$ , the set  $\text{Cont}(h, \sigma, \tau)$  contains exactly one profile that is threat-free on  $h$ .*

### C. Round-Parameterized Version

We will use our general definition of TFNE in games that are induced by cryptographic protocols. We assume that in these games players alternate moves, and thus there is a natural notion of the “rounds” in the game: Player  $i$  makes a move in round 1, then player  $-i$  makes a move in round 2, and so on until the end of the game.

For the general definition, we introduce a few modifications to vanilla version:

- We add a slackness parameter  $\varepsilon$ . This is necessary for our applications in order to handle the probability of error inherent in almost all cryptographic protocols.
- We allow players to be threatened at rounds, rather than just specific histories. This is needed because when we add the slackness parameter, a player might be threatened at a set of histories, where the weight of each individual threat does not exceed the slackness parameter, but the overall weight does.
- Finally, for a player to be threatened, we require that he improve on *all* threat-free continuations  $\pi$ . The reason we need this is that in the general case, there may be more than one  $\pi$  that is threat-free. If a player deviates from his prescribed behavior, he cannot choose *which* (threat-free) continuation will be played.

The definitions below make use of the notion of a round  $R$  strategy of player  $i$ : This is simply a function mapping every history  $h$  that reaches round  $R$  to a distribution over  $A(h)$ .

For a round  $R \in \mathbb{N}$  we let  $\sigma_i(R)$  represent player  $i$ 's round  $R$  strategy implied by  $\sigma$ . Let  $\sigma(R) = (\sigma_1(R), \sigma_2(R))$ , and let  $\text{Cont}(\sigma(1), \dots, \sigma(R)) \stackrel{\text{def}}{=} \left\{ \pi \in T : \pi(S) = \sigma(S) \ \forall S \leq R \right\}$ .

**Definition 4.4 ( $\varepsilon$ -threat):** Let  $\Gamma = (H, P, A, u)$  be an extensive game with constraints  $T = (T_1, T_2)$ . Let  $\varepsilon \geq 0$ , let  $\sigma \in T$  be a strategy profile, and let  $R \in \mathbb{N}$ . Player  $i = P(R)$  is facing an  $\varepsilon$ -threat at round  $R$  with respect to  $\sigma$  if there exists a round  $R$  strategy  $\tau = \tau(R)$  for player  $i$  such that

- (i) the set  $\text{Cont}(\sigma(1), \dots, \sigma(R-1), \tau(R))$  is non-empty, and
- (ii) for all  $\pi \in \text{Cont}(\sigma(1), \dots, \sigma(R-1), \tau(R))$  and  $\pi' \in \text{Cont}(\sigma(1), \dots, \sigma(R))$  that are  $\varepsilon$ -threat-free on  $R$

$$\mathbb{E}[u_i(O(\pi))] > \mathbb{E}[u_i(O(\pi'))] + \varepsilon,$$

where strategy  $\pi$  said to be  $\varepsilon$ -threat-free on  $R$  if for all rounds  $S > R$  it holds that player  $P(S)$  is not facing an  $\varepsilon$ -threat at round  $S$  with respect to  $\pi$ .

Note that if  $R$  is the last round of the game, then any profile  $\pi \in T$  is  $\varepsilon$ -threat-free on  $R$ . Using Definition 4.4, we can now define an  $\varepsilon$ -TFNE.

**Definition 4.5 ( $\varepsilon$ -threat-free Nash equilibrium):** Let  $\Gamma = (H, P, A, u)$  be an extensive game with constraints  $T = (T_1, T_2)$ . A strategy profile  $\sigma^* \in T$  is said to be in  $\varepsilon$ -threat-free Nash equilibrium ( $\varepsilon$ -TFNE) if:

- 1)  $\sigma^*$  is an  $\varepsilon$ -NE of  $\Gamma$ , and
- 2) for any round  $R$  of  $\Gamma$ , player  $P(R)$  is not facing an  $\varepsilon$ -threat at round  $R$  with respect to  $\sigma^*$ .

As is the case for Definition 4.1, Definition 4.4 (and hence Definition 4.5) would not be (semantically) well-defined if either one of the sets  $\text{Cont}(\sigma(1), \dots, \sigma(R-1), \tau(R))$  or  $\text{Cont}(\sigma(1), \dots, \sigma(R))$  would not contain at least one profile  $\pi$  that is  $\varepsilon$ -threat-free on  $R$ . The following proposition (whose proof is deferred to [11]) shows that this can never be the case.

**Proposition 4.6:** Let  $\Gamma = (H, P, A, u)$  be an extensive game with constraints  $T = (T_1, T_2)$ . Let  $\varepsilon \geq 0$ , let  $\sigma \in T$  be a strategy profile, and let  $R$  be a round of  $\Gamma$ . For any round  $R$  strategy  $\tau = \tau(R)$  for player  $i = P(R)$ , if the set  $\text{Cont}(\sigma(1), \dots, \sigma(R-1), \tau(R))$  is nonempty then it contains at least one profile  $\pi$  that is  $\varepsilon$ -threat-free on  $R$ .

## V. THE COMPUTATIONAL SETTING

In the following we explain how to use the notion of TFNE for cryptographic protocols. In Section V-A we describe how to view a cryptographic protocol as a sequence of extensive games. In Section V-B we show how to translate the behavior of an interactive TM to a sequence of strategies. In Section V-C we show how to express computational hardness in a game-theoretic setting. Finally, in Section V-D we give our definition of computational TFNE.

### A. Protocols as Sequences of Games

When placing cryptographic protocols in the framework of extensive games, the possible messages of players in a protocol correspond to the available actions in the game tree, and the prescribed instructions correspond to a strategy in the game.

The protocol is parameterized by a security parameter  $k \in \mathbb{N}$ . The set of possible messages in the protocol, as well as its prescribed instructions, typically depend on this  $k$ . Assigning for each  $k$  and each party a payoff for every outcome, a protocol naturally induces a sequence  $\Gamma^{(k)} = (H^{(k)}, P^{(k)}, A^{(k)}, u^{(k)})$  of extensive games, where:

- $H^{(k)}$  is the set of possible *transcripts* of the protocol (sequences of messages exchanged between the parties). A history  $h \in H^{(k)}$  is *terminal* if the prescribed instructions of the protocol instruct the player whose turn it is to play next to halt on input  $h$ .
- $P^{(k)} : (H^{(k)} \setminus Z^{(k)}) \rightarrow \{1, 2\}$  is a function that assigns a “next” player to every non-terminal history.
- $A^{(k)}$  is a function that assigns to every non-terminal history  $h \in H^{(k)} \setminus Z^{(k)}$  a set  $A^{(k)}(h) = \{m : (h \circ m) \in H^{(k)}\}$  of possible protocol messages to player  $P^{(k)}(h)$ .<sup>3</sup>
- $u^{(k)} = (u_1^{(k)}, u_2^{(k)})$  is a vector of payoff functions  $u_i^{(k)} : Z^{(k)} \rightarrow \mathbb{R}$ .

A sequence  $\Gamma = \{\Gamma^{(k)}\}_{k \in \mathbb{N}}$  of games defined as above is referred to as a *computational game*.

### B. Strategic Representation of Interactive Machines

Protocols are defined in terms of *interactive Turing machines* (ITMs). (See [8].) Thus, the prescribed behavior for

<sup>3</sup>We can interpret “disallowed” messages in the protocol as abort, and define “abort” as a possible protocol message. This will imply that every execution of the protocol corresponds to some history in the game.

each player is defined via an ITM, and any possible deviation of this player corresponds to choosing a different ITM. In order to argue about the protocol in a game-theoretic manner we formalize, using game-theoretic notions, the strategic behavior implied by ITMs. We believe this formalization is necessary for our treatment or any game-theoretic analysis of ITMs, in particular because, to the best of our knowledge, it has never been done before. The full formalization is deferred to [11], and has the following (informally stated) conclusion: The strategic behavior of an ITM for player  $i$  in a protocol may be seen as a collection of independent distributions on actions, one for each of player  $i$ 's histories that are reached with positive probability given the ITM of player  $i$  and some strategy profile of the other players. We refer to this collection as the behavioral reduced strategy induced by the ITM.

### C. Computational Hardness in the Game-Theoretic Setting

The security of cryptographic protocols stems from the assumption on the limitation of the computational power of the players. In our strategic analysis of games, we also expect to deduce the (sequential) equilibrium from this limitation. However, because protocols are parameterized by a security parameter, a strategic analysis of protocols requires dealing with a *sequence* of games rather than a single game. While relating to the sequence of games is crucial in order to express computational hardness (as this hardness is defined in an asymptotic manner), this raises a new difficulty: How do we extend the definition of TFNE to sequences of games?

Our approach insists on analyzing empty threats for *individual* games. Thus, our solution concept reflects a hybrid approach that relates to a protocol both as a family of *individual, extensive games* and as a *sequence of normal-form games*. To eliminate empty threats one must relate to the *interactive* aspect of each *individual* game (as this is the setting where threats are defined). In order to claim players are playing optimally under their computational constraints, one must think of the protocol as a *sequence of one-shot* games (because computational hardness is meaningful only when players are required to choose their machines in advance, and as the traditional notion of hardness is stated asymptotically).

1) *Strategy-filters*: When considering computational games  $\Gamma = \{\Gamma^{(k)}\}_{k \in \mathbb{N}}$ , the computational bounds on the players will be expressed by restricting the space of available strategies for the players. The available sequences of reduced strategies for the players will be exactly those that can be played by the ITMs that meet the computational bound on the players. In our case we will consider PPT ITMs.

While on the one hand every PPT ITM fails on cryptographic challenges for large enough values of the security parameter  $k$  (under appropriate assumptions), on the other hand, PPT ITMs can have arbitrarily large size and thus arbitrarily much information hardwired, and so for every  $k$  there is a PPT ITM that breaks the cryptographic challenges with security parameter  $k$ . In our analysis, we would like to “filter” machines according to their ability to break cryptographic challenges for specific  $k$ 's, and allow using them only in games

that correspond to large enough  $k$ 's, where these machines fail (and in particular, cannot use hard-wiring to solve the cryptographic challenges).

To this end, we define the notion of *strategy-filter*. For each value  $k$  of the security parameter and value  $\varepsilon$ , a strategy-filter maps the ITM  $M$  to either  $\perp$  or to its strategic representation, according to whether  $M(1^k)$  violates level of security  $\varepsilon$  or does not (respectively).

**Definition 5.1 (Strategy-filter):** Let  $\Gamma = \{\Gamma^{(k)}\}_{k \in \mathbb{N}}$  be a computational game and let  $i$  be a player. A strategy-filter is a sequence  $F_i = \{F_i^{(k)} : \mathcal{M} \times [0, 1] \rightarrow \Sigma_i^{(k)} \cup \{\perp\}\}_{k \in \mathbb{N}}$  such that for every ITM  $M$ , every  $k \in \mathbb{N}$  and every  $\varepsilon \in [0, 1]$ , it holds that either  $F_i^{(k)}(M, \varepsilon) = \perp$ , or  $F_i^{(k)}(M, \varepsilon) = \sigma_i^{(k)}$ , where  $\sigma_i^{(k)}$  is the strategic representation of the machine  $M(1^k, \cdot)$ .

A strategy-filter is meaningful if it allows us to reason about all reduced strategies that are considered to be feasible, in our case PPT implementable reduced strategies, and in particular does not filter them out.

**Definition 5.2 (PPT-covering filter):** A strategy-filter  $F_i$  is said to be PPT-covering if for every PPT ITM  $M$  and any positive polynomial  $p(\cdot)$  there exists  $k_0$  such that for all  $k \geq k_0$ , it holds that  $F_i^{(k)}(M, 1/p(k)) \neq \perp$ .

Typically, protocols have the following security guarantee (under computational assumptions): for every  $i$ , every PPT ITM  $M$  of  $P_i$  and every polynomial  $p(\cdot)$ , there exists  $k_0$  such that for any  $k \geq k_0$ , the ITM  $M$  does not break level of security  $1/p(k)$  in the protocol with security parameter  $k$ . Such a protocol will naturally have a PPT-covering filter, where if  $F_i^{(k)}(M, \varepsilon) \neq \perp$  then the reduced strategy  $F_i^{(k)}(M, \varepsilon)$  “does not break level of security  $\varepsilon$  in the game  $\Gamma^{(k)}$ .”

2) *Tractable Reduced Strategies*: As reflected above, the asymptotic nature of defining security does not determine any level of security for any  $k$ . Rather, it dictates that any PPT ITM “eventually fails in violating  $1/p(k)$  security” for any  $p(\cdot)$  (where “eventually” means for large enough  $k$ ). Thus, we follow the same approach in our game theoretic analysis: roughly speaking, our solution concept requires that  $\varepsilon$ -security will imply  $\varepsilon$ -stability for any  $k$  (rather than requiring a particular level of stability for each  $k$ ). More formally, we require that for any  $k$  and any  $\varepsilon$ , the game induced by the protocol with security parameter  $k$  be in  $\varepsilon$ -TFNE, given that the available strategies for the players are those that do not break level of security  $\varepsilon$ . Thus, for any pair  $(k, \varepsilon)$  we will consider the game  $\Gamma^{(k)}$  with available reduced strategies restricted to those that guarantee  $\varepsilon$ -security. The following definition derives from a PPT-covering filter, for each such game, the set of available reduced strategies for each player.

**Definition 5.3 (Tractable reduced strategies):** Let  $F_i$  be a PPT-covering filter. For every  $k \in \mathbb{N}$  and  $\varepsilon \in [0, 1]$  we define the set  $T_{i, \varepsilon}^{(k)}(F_i)$  of  $(k, \varepsilon)$ -tractable reduced strategies for player  $i \in \{1, 2\}$  as

$$\{F_i^{(k)}(M, \varepsilon) \mid M \text{ is a PPT ITM and } F_i^{(k)}(M, \varepsilon) \neq \perp\}.$$

Whenever  $F_i$  will be understood from the context, we will write  $T_{i, \varepsilon}^{(k)}$  to mean  $T_{i, \varepsilon}^{(k)}(F_i)$ .

#### D. Computational TFNE

We can now define our computational variant of TFNE. Roughly, the definition requires that there exist a family of PPT compatible constraints such that for any  $k$  and any  $\varepsilon$ , the strategies played by the machines on input security parameter  $k$  are in  $\varepsilon$ -TFNE in the game indexed by  $(k, \varepsilon)$ .

**Definition 5.4 (Computational TFNE):** Let  $\Gamma$  be a computational game. A pair of PPT machines  $(M_1, M_2)$  is said to be in a computational threat-free Nash equilibrium (CTFNE) of  $\Gamma$  if there exists a pair of PPT-covering filters  $(F_1, F_2)$  such that for every  $k, \varepsilon$  for which  $F_1^{(k)}(M_1, \varepsilon)$  and  $F_2^{(k)}(M_2, \varepsilon)$  are tractable the profile  $(F_1^{(k)}(M_1, \varepsilon), F_2^{(k)}(M_2, \varepsilon))$  constitutes an  $\varepsilon$ -TFNE in the  $(T_{1, \varepsilon}^{(k)}, T_{2, \varepsilon}^{(k)})$ -constrained version of  $\Gamma^{(k)}$ .

### VI. CORRELATED EQUILIBRIA WITHOUT A MEDIATOR

In one of the first papers to consider the intersection between game theory and cryptography, Dodis, Halevi and Rabin proposed an appealing methodology for implementing a correlated equilibrium in a 2-player normal-form game without making use of a mediator [5]. Under standard hardness assumptions, they showed that for any 2-players normal-form game  $\Gamma$  and any correlated equilibrium  $\sigma$  for  $\Gamma$ , there exists a new 2-player extensive-form “extended game”  $\Gamma'$  and a CNE  $\sigma'$  for  $\Gamma'$ , such that  $\sigma$  and  $\sigma'$  achieve the same payoffs for the players. However, as already pointed out by Dodis et al., their protocol lacks a satisfactory analysis of its sequential nature - the resulted “extended game” is an extensive-form game, while the solution concept they use, CNE, is not strong enough for these games.

In the following, we extend the definition of CTFNE to allow handling this setting (that is, we define CTFNE for extensive games with simultaneous moves at the leaves), give some justification for our new definition, and then provide a new protocol for removing the mediator that achieves CTFNE in a wide class of correlated equilibria that are in the convex hull of Nash equilibria (see definition below).

#### A. TFNE for Games with Simultaneous Moves at the Leaves

For a formal definition of extensive game with simultaneous moves see Osborne and Rubinstein [24]. In order to adjust our definition for extensive games with simultaneous moves, we notice that when a player deviates on a history with a simultaneous move, he cannot expect the other to react to this deviation (because they both play at the same time). However, in order to argue that a profile is rational, we still need to require that for every simultaneous move in the equilibrium support, each player is playing a “best response” given the other player’s prescribed behavior. This means the prescribed behavior for the players should form some kind of equilibrium for normal-form games. In our case, the prescribed behavior will form a NE. The question of what should a CTFNE profile prescribe in off-equilibrium-support histories is more delicate: Clearly, in order to claim that the profile is “rational”, again we need some kind of equilibrium for normal-form games. But in this case one can argue that after one player deviated,

the other player cannot assume the deviating player will play his prescribed behavior in the simultaneous move (as he is already not following his prescribed behavior). However, we argue that it is in fact still rational to assume the deviating player will play his prescribed behavior. The justification for this claim is essentially the same as the justification for the rationality of NE. Once there is a prescribed behavior that is a NE, each player knows the other has no incentive to deviate, and so he also has no incentive to deviate.

Thus, our new definition of TFNE for extensive games with simultaneous moves at the leaves (abbreviated GSML), is essentially the same as the original definition, except that (i) we require a profile in TFNE to prescribe a NE in any terminal leaf, and (ii) in the definition of a threat we do not allow a player to assume the other will deviate from his strategy in any NE. In order to formally modify our definition of TFNE to achieve (ii), essentially we would need to define the only threat-free continuation on a leaf to be the one that assigns to the players the actions in the prescribed NE (which expresses the idea that a player is not allowed to assume the other will deviate from his strategy in any NE).

However, we adopt an equivalent, simpler convention. Given a GSML  $\Gamma$  and a profile  $\sigma$  that assigns a NE at every simultaneous move, we look at a slightly modified game  $\Gamma'$ : All simultaneous moves are removed, and instead at each leaf where a simultaneous move was removed each player is assigned his expected payoff in the corresponding NE for that leaf. Note that the modified game is now a regular extensive game with *no* simultaneous moves. We then “prune” the profile to remove all the distributions on actions on all simultaneous leaves and denote the resulting profile  $\sigma'$ . We say that  $\sigma$  is a TFNE in  $\Gamma$  if  $\sigma'$  is a TFNE in  $\Gamma'$ . We call  $\Gamma'$  and  $\sigma'$  the *pruned representation* of  $\Gamma$  and  $\sigma$ .

The definition of CTFNE for GSML is derived from the above definition of TFNE for GSML, similarly to the derivation of CTFNE from TFNE in the non-simultaneous case.

It seems that for general GSML’s our definition is too strong, because in certain cases it is computationally intractable to compute the assigned NE in every leaf. While we do not yet know how to relax our definition to apply to these cases, we believe our definition, when met, is sufficient.

#### B. Our Protocol

For a non-trivial class of correlated equilibria, we show how to modify the DHR protocol to achieve CTFNE. Our basic idea is to use Nash equilibria as “punishments” for aborting players. That is, if there is a NE that assigns to a player a payoff at most his expected payoff when not aborting, then assigning this NE in case he aborts serves as a punishment and yields that the player has no incentive to abort. In the following we characterize a family of correlated equilibria for which we can use the aforementioned punishing technique, and prove that for this family we can remove the mediator while achieving CTFNE.

We say that a correlated equilibrium  $\pi$  is a *convex combination of Nash equilibria* if  $\pi$  is induced by a distribution on

(possibly mixed) Nash equilibria. (The set of such distributions is sometimes referred to as the *convex hull of Nash equilibria*.) Note that any such distribution is a correlated equilibrium (CE), but the converse is not true.

Let  $\pi$  be a correlated equilibrium for a two-player game  $\Gamma$  that is a convex combination of a set  $N$  of NEs. We say that  $\pi$  is *weakly Pareto optimal* if there does not exist a different CE  $\rho$  in the convex hull of  $N$  for which both  $\mathbb{E}[u_1(O(\rho))] > \mathbb{E}[u_1(O(\pi))]$  and  $\mathbb{E}[u_2(O(\rho))] > \mathbb{E}[u_2(O(\pi))]$ .

We say that a distribution is *samplable* if there exists a probabilistic TM that halts on every infinite randomness vector, and can sample it. This is equivalent to requiring that all probabilities can be expressed in binary (assuming we work over  $\{0, 1\}$ ). Note that every distribution can be approximated arbitrarily accurately by a samplable distribution.

**Theorem 6.1:** *Assume there exists a non-interactive computationally binding commitment scheme. Let  $\pi$  be a weakly Pareto optimal correlated equilibrium for a two-player game  $\Gamma$  that is a samplable convex combination  $\Pi$  of some set of samplable Nash equilibria. Then there exists an extended extensive game and a profile that achieves the same expected payoffs as  $\pi$  and is a CTFNE.*

*Proof:* Let  $\Pi$  be as above. Since  $\Pi$  is samplable, the common denominator of all probabilities in  $\Pi$  is a power of two. Thus, we can assume  $\Pi$  is a *uniform* distribution on a sequence of Nash equilibria that may contain repetitions, where the length of the sequence is a power of two. Let  $2^\ell$  be the length of that sequence, and let  $(\pi_{0^\ell}, \dots, \pi_{1^\ell})$  be that sequence. Note that the distribution  $\pi$  can now be generated by first choosing uniformly at random a string  $r$  in  $\{0, 1\}^\ell$ , and then choosing a pair of actions according to  $\pi_r$ .

Let  $\hat{\sigma}^i$  be the NE that assigns the worst payoff for  $P_i$  (this value represents the “severest punishment” for player  $i$ ).

Our protocol embeds a 2-party string sampling protocol, which is a simple generalization of the Blum coin flipping protocol [4] whose security is defined in the real-ideal model. (For more details on this model see [7].) The protocol consists of simply running the Blum protocol concurrently for a fixed number of times. It can be shown that this protocol too is secure in the real-ideal model.

We denote the ITMs playing the strategies of  $P_1, P_2$  by  $M_1, M_2$ , respectively.

- **Round 1:** Player 1 chooses uniformly at random a string  $r = (r_1, \dots, r_\ell)$  from  $\{0, 1\}^\ell$ , and sends  $c = (c_1 = \text{com}^{(k)}(r_1), \dots, c_\ell = \text{com}^{(k)}(r_\ell))$  to player 2 (player 1 also obtains  $(\text{decom}_1, \dots, \text{decom}_\ell)$ , where  $\text{decom}_i$  is a legal decommitment with respect to  $c_i$  and  $r_i$ ).
- **Round 2:** If Player 1 aborted, the assigned NE is  $\hat{\sigma}^1$ . Else, Player 2 chooses a uniformly random string  $r' = (r'_1, \dots, r'_\ell)$  from  $\{0, 1\}^\ell$ , and sends  $r'$  to player 1.
- **Round 3:** If Player 2 aborted, the assigned NE is  $\hat{\sigma}^2$ . Else, Player 1 sends  $((r_1, \text{decom}_1), \dots, (r_\ell, \text{decom}_\ell))$ .
- If Player 1 aborted, the assigned NE is  $\hat{\sigma}^1$ . Else, Player 2 verifies that  $\text{decom}_i$  is a legal decommitment with respect to  $c_i$  and  $r_i$  for  $1 \leq i \leq \ell$ . If the verification fails (which

is equivalent to an abort of Player 1, as it means Player 1 sent an illegal message), the assigned NE is  $\hat{\sigma}^1$ . Else, the assigned NE is  $\pi_{r \oplus r'}$  (where  $\oplus$  is bitwise exclusive-or).

We now show that the pair  $(M_1, M_2)$  forms a CTFNE for the protocol above. Let the sequence of games induced by the protocol be  $\{\tilde{\Gamma}^{(k)}\}_{k \in \mathbb{N}}$ . Denote the pruned representation of  $\tilde{\Gamma}^{(k)}$  by  $\Gamma^{(k)}$ . Let  $\tilde{\sigma}_1^{(k)}, \tilde{\sigma}_2^{(k)}$  be the strategies of  $P_1, P_2$  in the protocol with security parameter  $k$ , and let  $\sigma_1^{(k)}, \sigma_2^{(k)}$  be their pruned representations. Let  $\sigma^{(k)} = (\sigma_1^{(k)}, \sigma_2^{(k)})$ . We prove that  $\{\sigma^{(k)}\}$  is CTFNE in  $\{\Gamma^{(k)}\}$ , which, by the discussion above, implies that  $\{\tilde{\sigma}^{(k)}\}$  is CTFNE in  $\{\tilde{\Gamma}^{(k)}\}$ .

First we define the functions  $F_1^{(k)}$  and  $F_2^{(k)}$ . For any  $k$ , the function  $F_1^{(k)}$  never maps to  $\perp$  (this, roughly speaking, reflects the fact that the protocol is secure against an all-powerful player 1, which follows from the perfect binding of the commitment scheme). For  $F_2$  we use the following rule:  $F_2^{(k)}(M, \varepsilon) = \perp$  if and only if

$$\mathbb{E}[u_2^{(k)}(O(\sigma_1^{(k)}, \sigma_M^{(k)}))] > \mathbb{E}[u_2^{(k)}(O(\sigma^{(k)}))] + \varepsilon,$$

where  $\sigma_M^{(k)}$  is the strategic representation of machine  $M$  and  $\sigma_1^{(k)}$  is the strategic representation of machine  $M_1$ , both with security parameter  $k$ . In other words,  $P_2$  cannot unilaterally  $\varepsilon$ -improve in the  $(T_{1,\varepsilon}^{(k)}, T_{2,\varepsilon}^{(k)})$ -constrained version of  $\Gamma^{(k)}$ .

The fact that  $F_1$  is PPT-covering is straightforward. The fact that  $F_2$  is PPT covering follows from the security of the commitment scheme – see the full version for details [11].

Next, we show that for all  $k, \varepsilon$  for which  $F_1^{(k)}(M_1, \varepsilon) \neq \perp$  and  $F_2^{(k)}(M_2, \varepsilon) \neq \perp$  the profile  $(F_1^{(k)}(M_1, \varepsilon), F_2^{(k)}(M_2, \varepsilon))$  constitutes an  $\varepsilon$ -TFNE in the  $T = (T_{1,\varepsilon}^{(k)}, T_{2,\varepsilon}^{(k)})$ -constrained version of  $\Gamma^{(k)}$ . Let  $k, \varepsilon$  be as above, and let  $\sigma = (\sigma_1, \sigma_2) = (F_1^{(k)}(M_1, \varepsilon), F_2^{(k)}(M_2, \varepsilon))$ . Suppose  $P_1$  unilaterally  $\varepsilon$ -improves in the  $T$ -constrained version of  $\Gamma^{(k)}$ . From similar arguments as above we can assume  $P_1$  never aborts. But when  $P_1$  never aborts the outcome is exactly  $\pi$ , as the players are playing  $\pi_{r \oplus r'}$ , and  $r'$  is chosen uniformly at random. Suppose now that  $P_2$  unilaterally  $\varepsilon$ -improves in the  $T$ -constrained version of  $\Gamma^{(k)}$ . However, this is a contradiction to the constraints, that state that for any  $k$   $P_2$  cannot unilaterally  $\varepsilon$ -improve in the  $(T_{1,\varepsilon}^{(k)}, T_{2,\varepsilon}^{(k)})$ -constrained version of  $\Gamma^{(k)}$ .

Next, we show that no player is  $\varepsilon$ -threatened with respect to  $\sigma$  at any round of the  $T$ -constrained version of  $\Gamma^{(k)}$ . To this end, suppose towards a contradiction that some player is  $\varepsilon$ -threatened with respect to  $\sigma$ . We divide the proof into cases.

**Case 1 –  $P_1$  is facing an  $\varepsilon$ -threat in round 3:** In step 3 player 1 has exactly two options: He can (i) play honestly, send  $((r_1, \text{decom}_1), \dots, (r_\ell, \text{decom}_\ell))$  which he generated in round 1, and receive  $\mathbb{E}[u_1(O(\sigma))]$ , or he can (ii) abort and receive  $\mathbb{E}[u_1(O(\hat{\sigma}^1))]$ . The value  $\mathbb{E}[u_1(O(\hat{\sigma}^1))]$  is at most  $\mathbb{E}[u_1(O(\sigma))]$ , and so  $P_1$  cannot improve over  $\mathbb{E}[u_1(O(\sigma))]$ . Hence player 1 is not facing an  $\varepsilon$ -threat at round 3.

**Case 2 –  $P_2$  is facing an  $\varepsilon$ -threat in round 2:** We first note that for any round 1 strategy for  $P_1$  and round 2 strategy for  $P_2$ , the round strategy of playing honestly in round 3 for  $P_1$  is threat-free, since he cannot improve over that strategy

(again, since his only deviation is aborting, which gives him the worst possible NE). Thus, if  $P_2$  is  $\varepsilon$ -threatened at round 2, he has some round strategy that  $\varepsilon$ -improves over  $E[u_2(O(\sigma))]$  when  $P_1$  plays in round 3 (and 1) according to the protocol. This means that  $P_2$  unilaterally  $\varepsilon$ -improves, which contradicts the constraints (as well as the  $\varepsilon$ -NE).

**Case 3 –  $P_1$  is facing an  $\varepsilon$ -threat in round 1:** If  $P_1$  is  $\varepsilon$ -threatened in round 1, he has some round 1 strategy  $\tau(1)$  for which every  $\varepsilon$ -threat-free continuation  $\varepsilon$ -improves over every  $\varepsilon$ -threat-free continuation of  $\sigma_1(1)$ . We will describe an  $\varepsilon$ -threat-free continuation of  $\tau(1)$  and an  $\varepsilon$ -threat-free continuation of  $\sigma_1(1)$  that contradict this.

The  $\varepsilon$ -threat-free continuation of  $\sigma_1(1)$ : We established in case 2 that when  $P_1$  plays honestly in round 1, if  $P_2$  plays honestly in round 2 he is not  $\varepsilon$ -threatened. We also established there that  $P_1$  playing honestly in round 3 is always  $\varepsilon$ -threat-free. It follows that the continuation of both players playing honestly in rounds 2 and 3 is an  $\varepsilon$ -threat-free continuation of  $\sigma_1(1)$ . On this profile  $P_1$  receives  $E[u_1(O(\sigma))]$ .

The  $\varepsilon$ -threat-free continuation of  $\tau_1(1)$ : As we established in case 2, playing honestly in round 3 is always  $\varepsilon$ -threat-free for  $P_1$ . Now, note that there is no profile in which both players improve simultaneously – because all leaves are Nash equilibria, such a profile would be a distribution on Nash equilibria that contradicts the Pareto-optimality of  $\pi$ . Note also that because  $P_1$  receives the worst possible payoff when he aborts, it follows that he improves also conditioned on not aborting (as this can only help him). Thus, in any threat-free continuation of  $\tau(1)$ , conditioned on  $P_1$  not aborting in round 1,  $P_2$  again cannot improve over  $E[u_2(O(\sigma))]$ , as this again contradicts the Pareto-optimality of  $\pi$ . However, if  $P_2$  plays honestly in round 2 and then  $P_1$  plays honestly in round 3, then  $P_2$  receives exactly  $E[u_2(O(\sigma))]$  conditioned on  $P_1$  not aborting in round 1. It follows that this continuation is the best possible for  $P_2$ , and thus  $P_2$  is not  $\varepsilon$ -threatened in round 2 of this continuation. It follows that this continuation is  $\varepsilon$ -threat-free. However, in this continuation  $P_1$  receives  $E[u_1(O(\sigma))]$  conditioned on not aborting, and thus receives at most  $E[u_1(O(\sigma))]$  without the conditioning. ■

## VII. A GENERAL THEOREM

**Definition 7.1 (Weakly Pareto optimal):** A strategy profile  $\sigma \in T$  of an extensive game  $\Gamma = (H, P, A, u)$  with constraints  $T$  is weakly Pareto optimal if there does not exist a strategy profile  $\pi \in T$  for which both  $E[u_1(O(\pi))] > E[u_1(O(\sigma))]$  and  $E[u_2(O(\pi))] > E[u_2(O(\sigma))]$ .

**Definition 7.2 ( $\varepsilon$ -safe):** A strategy profile  $\sigma = (\sigma_1, \sigma_2) \in T$  of an extensive game  $\Gamma = (H, P, A, u)$  with constraints  $T = (T_1, T_2)$  is  $\varepsilon$ -safe if for each player  $i$ ,

$$E[u_{-i}(O(\sigma))] \geq E[u_{-i}(O(\sigma'_i, \sigma_{-i}))] - \varepsilon$$

for every strategy  $\sigma'_i \in T_i$  of player  $i$ .

We now have the following theorem and corollary, whose proofs appears in the full version [11].

**Theorem 7.3:** Let  $\Gamma = (H, P, A, u)$  be an extensive game with constraints  $T = (T_1, T_2)$ , and let  $\sigma = (\sigma_1, \sigma_2)$  be a weakly Pareto optimal  $\varepsilon$ -NE of  $\Gamma$  that is  $\varepsilon$ -safe. Then  $\sigma$  is an  $\varepsilon$ -TFNE of  $\Gamma$ .

**Corollary 7.4:** Let  $\Gamma = (H, P, A, u)$  be a zero-sum extensive game with constraints  $T = (T_1, T_2)$ , and let  $\sigma$  be an  $\varepsilon$ -NE of  $\Gamma$ . Then  $\sigma$  is an  $\varepsilon$ -TFNE of  $\Gamma$ .

The corollary follows from the observation that any  $\varepsilon$ -NE of a zero-sum game is both weakly Pareto optimal and  $\varepsilon$ -safe.

## ACKNOWLEDGMENTS

We thank Eddie Dekel, Oded Goldreich, Ehud Kalai, Eran Omri, and Gil Segev for helpful conversations, and the anonymous referees for careful reading and insightful comments.

## REFERENCES

- [1] I. Abraham, D. Dolev, R. Gonen, and J. Halpern. Distributed computing meets game theory: robust mechanisms for rational secret sharing and multiparty computation. In *25th PODC*, pages 53–62. ACM Press, 2006.
- [2] G. Asharov and Y. Lindell. Utility dependence in correct and fair rational secret sharing. In *Crypto 2009*, pages 559–576, 2009.
- [3] E. Ben-Sasson, A. Tauman-Kalai, E. Kalai. An approach to bounded rationality. *Advances in Neural Information Processing Systems*, 2007.
- [4] M. Blum. Coin flipping by telephone. In *CRYPTO81*, p. 1115, 1981.
- [5] Y. Dodis, S. Halevi, and T. Rabin. A cryptographic solution to a game theoretic problem. In *Crypto'00*, 2000.
- [6] L. Fortnow and R. Santhanam. Bounding rationality by discounting time. In *1st ICS*, 2010.
- [7] Oded Goldreich. *Computational Complexity: A Conceptual Perspective*. Cambridge University Press, 2008.
- [8] Oded Goldreich. *Foundations of Cryptography: Basic Tools*. Cambridge University Press, 2001.
- [9] S. D. Gordon and J. Katz. Rational secret sharing, revisited. In *5th SCN*, p. 229–241. Springer, 2006.
- [10] R. Gradwohl. Rationality in the full-information model. In *TCC 2010*.
- [11] R. Gradwohl, N. Livne, and A. Rosen. Sequential rationality in cryptographic protocols. Full version.
- [12] R. Gradwohl, N. Livne, and A. Rosen. Incredible threats. In preparation.
- [13] J. Y. Halpern and R. Pass: Game theory with costly computation. In *1st ICS*, 2010.
- [14] J. Halpern and V. Teague. Rational secret sharing and multiparty computation: Extended abstract. In *36th STOC*, pages 623–632, 2004.
- [15] S. Izmalkov, S. Micali, and M. Lepinski. Rational secure computation and ideal mechanism design. In *FOCS 2005*, 2005.
- [16] J. Katz. Bridging game theory and cryptography: Recent results and future directions. In *5th TCC*, pages 251–272, 2008.
- [17] J. Katz, G. Fuchsbaauer, and D. Naccache. Efficient rational secret sharing in the standard communication model. In *TCC 2010*.
- [18] G. Kol and M. Naor. Cryptography and game theory: Designing protocols for exchanging information. In *5th TCC*, pages 320–339, 2008.
- [19] G. Kol and M. Naor. Games for exchanging information. In *40th STOC*, pages 423–432, 2008.
- [20] M. Lepinski, S. Micali, and a. shelat. Collusion-free protocols. In *STOC 2005*, 2005.
- [21] A. Lysyanskaya and N. Triandopoulos. Rationality and adversarial behavior in multi-party computation. In *Crypto 2006*, p. 180197, 2006.
- [22] S. Micali and A. Shelat. Truly rational secret sharing. In *6th TCC*, pages 54–71, 2009.
- [23] S. J. Ong, D. Parkes, A. Rosen, and S. Vadhan. Fairness with an honest minority and a rational majority. In *6th TCC*, pages 36–53, 2009.
- [24] M. J. Osborne and A. Rubinstein. *A Course in Game Theory*. MIT Press, 1994.