



But Why Does It Work? A Rational Protocol Design Treatment of Bitcoin

Christian Badertscher¹ , Juan Garay², Ueli Maurer¹, Daniel Tschudi³ ,
and Vassilis Zikas⁴

¹ ETH Zurich, Zürich, Switzerland

{christian.badertscher,maurer}@inf.ethz.ch

² Texas A&M University, College Station, USA
garay@tamu.edu

³ Aarhus University, Aarhus, Denmark
tschudi@cs.au.dk

⁴ University of Edinburgh and IOHK, Edinburgh, UK
vassilis.zikas@ed.ac.uk

Abstract. An exciting recent line of work has focused on formally investigating the core cryptographic assumptions underlying the security of Bitcoin. In a nutshell, these works conclude that Bitcoin is secure if and only if the majority of the mining power is honest. Despite their great impact, however, these works do not address an incisive question asked by positivists and Bitcoin critics, which is fuelled by the fact that Bitcoin indeed works in reality: Why should the real-world system adhere to these assumptions?

In this work we employ the machinery from the Rational Protocol Design (RPD) framework by Garay *et al.* [FOCS 2013] to analyze Bitcoin and address questions such as the above. We show that under the natural class of incentives for the miners' behavior—i.e., rewarding them for adding blocks to the blockchain but having them pay for mining—we can reserve the honest majority assumption as a fallback, or even, depending on the application, completely replace it by the assumption that the miners aim to maximize their revenue.

Our results underscore the appropriateness of RPD as a “rational cryptography” framework for analyzing Bitcoin. Along the way, we devise significant extensions to the original RPD machinery that broaden its applicability to cryptocurrencies, which may be of independent interest.

1 Introduction

Following a number of informal and/or *ad hoc* attempts to address the security of Bitcoin, an exciting recent line of work has focused on devising a rigorous cryptographic analysis of the system [2, 13, 14, 27]. At a high level, these works

D. Tschudi—Work done while author was at ETH Zurich.

V. Zikas—Work done in part while the author was at RPI.

start by describing an appropriate model of execution, and, within it, an abstraction of the original Bitcoin protocol [23] along with a specification of its security goals in terms of a set of intuitive desirable properties [13, 14, 27], or in terms of a functionality in a simulation-based composable framework [2]. They then prove that (their abstraction of) the Bitcoin protocol meets the proposed specification under the assumption that the majority of the computing power invested in mining bitcoins is by devices which mine according to the Bitcoin protocol, i.e., *honestly*. This assumption of *honest majority* of computing power—which had been a folklore within the Bitcoin community for years underlying the system’s security—is captured by considering the parties who are not mining honestly as controlled by a central adversary who coordinates them trying to disrupt the protocol’s outcome.

Meanwhile, motivated by the fact that Bitcoin is an “economic good” (i.e., BTCs are exchangeable for national currencies and goods) a number of works have focused on a rational analysis of the system [7–9, 15, 22, 24, 28–32]. In a nutshell, these works treat Bitcoin as a game between the (competing) rational miners, trying to maximize a set of utilities that are postulated as a natural incentive structure for the system. The goal of such an analysis is to investigate whether or not, or under which assumptions on the incentives and/or the level of collaboration of the parties, Bitcoin achieves a stable state, i.e., a game-theoretic equilibrium. However, despite several enlightening conclusions, more often than not the prediction of such analyses is rather pessimistic. Indeed, these results typically conclude that, unless assumptions on the amount of honest computing power—sometimes even stronger than just majority—are made, the induced incentives result in plausibility of an attack to the Bitcoin mining protocol, which yields undesired outcomes such as forks on the blockchain, or a considerable slowdown.

Yet, to our knowledge, no fork or substantial slowdown that is attributed to rational attacks has been observed to date, and the Bitcoin network keeps performing according to its specification, even though mining pools would, in principle, be able to launch collaborative attacks given the power they control.¹ In the game-theoretic setting, this mismatch between the predicted and observed behavior would be typically interpreted as an indication that the underlying assumptions about the utility of miners in existing analysis do not accurately capture the miners’ rationale. Thus, two main questions still remain and are often asked by Bitcoin skeptics:

Q1. How come Bitcoin is not broken using such an attack?

Or, stated differently, why does it work and why do majorities not collude to break it?

Q2. Why do honest miners keep mining given the plausibility of such attacks?

¹ We refer to forks of the Bitcoin chain itself, not to forks that spin-off a new currency.

In this work we use a rigorous cryptographic reasoning to address the above questions. In a nutshell, we devise a rational-cryptography framework for capturing the economic forces that underly the tension between honest miners and (possibly colluding) deviating miners, and explain how these forces affect the miners’ behavior. Using this model, we show how natural incentives (that depend on the expected revenue of the miners) in combination with a high monetary value of Bitcoin, can explain the fact that Bitcoin is not being attacked in reality *even though* majority coalitions are in fact possible. In simple terms, we show how natural assumptions about the miners’ incentives allow to substitute (either entirely or as a fallback assumption) the honest-majority assumption. To our knowledge, this is the first work that formally proves such rational statements that do not rely on assumptions about the adversary’s computing power. We stress that the incentives we consider depend solely on costs and rewards for mining—i.e., mining (coinbase) and transaction fees—and, in particular, we make no assumption that implicitly or explicitly deters forming adversarial majority coalitions.

What enables us to address the above questions is utilizing the Rational Protocol Design (RPD) methodology by Garay *et al.* [11] to derive stability notions that closely capture the idiosyncrasies of coordinated incentive-driven attacks on the Bitcoin protocol. To better understand how our model employs RPD to address the above questions, we recall the basic ideas behind the framework.

Instead of considering the protocol participants—in our case, the Bitcoin miners—as rational agents, RPD considers a meta-game, called the *attack game*. The attack game in its basic form is a two-agent zero-sum extensive game of perfect information with a horizon of length two, i.e., two sequential moves.² It involves two players, called the *protocol designer* D—who is trying to come up with the best possible protocol for a given (multi-party) task—and the *attacker* A—who is trying to come up with the (polynomial-time) strategy/adversary that optimally attacks the protocol. The game proceeds in two steps: First, (only) D plays by choosing a protocol for the (honest) players to execute; A is informed about D’s move and it is now his term to produce his move. The attacker’s strategy is, in fact, a cryptographic adversary that attacks the protocol proposed by the designer.

The incentives of both A and D are described by utility functions, and their respective moves are carried out with the goal of maximizing these utilities.³ In a nutshell, the attacker’s utility function rewards the adversary proportionally to how often he succeeds in provoking his intended breach, and depending on its severity. Since the game is zero-sum, the designer’s utility is the opposite of the attacker’s; this captures the standard goal of cryptographic protocols, namely, “taming” the adversary in the best possible manner.

Based on the above game, the RPD framework introduces the following natural security notion, termed *attack-payoff security*, that captures the quality of

² This is often referred to as a *Stackelberg game* in the game theory literature [26].

³ Notice, however, the asymmetry: The designer needs to come up with a protocol based on speculation of what the adversary’s move will be, whereas the attacker plays after being informed about the actual designer’s move, i.e., about the protocol.

a protocol Π for a given specification when facing incentive-driven attacks aiming to maximize the attacker’s utility. Informally, attack-payoff security ensures that the adversary is not willing to attack the protocol Π in any way that would make it deviate from its ideal specification. In other words, the protocol is secure against the class of strategies that maximize the attacker’s utility. In this incentive-driven setting, this is the natural analogue of security against malicious adversaries.⁴ For cases where attack payoff security is not feasible, RPD proposes the notion of *attack-payoff optimality*, which ensures that the protocol Π is a best response to the best attack.

A useful feature of RPD (see below) is that all definitions build on Canetti’s simulation-based framework (either the standalone framework [5] or the UC framework [6]), where they can be easily instantiated. In fact, there are several reasons, both at the intuitive and technical levels, that make RPD particularly appealing to analyze complex protocols that are already running, such as Bitcoin. First, RPD supports adaptive corruptions which captures the scenario of parties who are currently running their (mining) strategy changing their mind and deciding to attack. This is particularly useful when aiming to address the likelihood of insider attacks against a protocol which is already in operation. For the same reason, RPD is also suitable for capturing attacks induced by compromised hardware/software and/or bribing [4] (although we will not consider bribing here). Second, the use of a central adversary as the attacker’s move ensures that, even though we are restricting to incentive-driven strategies, we allow full collaboration of cheaters. This allows, for example, to capture mining pools deciding to deviate from the protocol’s specification.

At the technical level, using the attack-game to specify the incentives takes away many of the nasty complications of “rational cryptography” models. For example, it dispenses with the need to define cumbersome computational versions of equilibrium [10, 17–19, 21, 25], since the actual rational agents, i.e., \mathcal{D} and \mathcal{A} , are not computationally bounded. (Only their actions need to be PPT machines.) Furthermore, as it builds on simulation-based security, RPD comes with a composition theorem allowing for regular cryptographic subroutine replacement. The latter implies that we can analyze protocols in simpler hybrid-worlds, as we usually do in cryptography, without worrying about whether or not their quality or stability will be affected once we replace their hybrids by corresponding cryptographic implementations.

Our contributions. In this work, we apply the RPD methodology to analyze the quality of Bitcoin against incentive-driven attacks, and address the existential questions posted above. As RPD is UC-based, we use the Bitcoin abstraction as a UC protocol and the corresponding Bitcoin ledger functionality from [2] to capture the goal/specification of Bitcoin. As argued in [2], this functionality captures all the properties that have been proposed in [13, 27].

We define a natural class of incentives for the attacker by specifying utilities which, on one hand, reward him according to Bitcoin’s standard reward

⁴ In fact, if we require this for any arbitrary utility function, then the two notions—attack-payoff security and malicious security—coincide.

mechanisms (i.e., block rewards and transaction fees) for blocks permanently inserted in the blockchain by adversarial miners, and, on the other hand, penalize him for resources that he uses (e.g., use of mining equipment and electricity). In order to overcome the inconsistency of rewards being typically in Bitcoins and costs being in real money, we introduce the notion of a *conversion rate* CR converting reward units (such as BTC) into mining-cost units (such as US Dollar). This allows us to make statements about the quality of the protocol depending on its value measured in a national currency.

We then devise a similar incentive structure for the designer, where, again, the honest parties are (collectively) rewarded for blocks they permanently insert into the blockchain, but pay for the resources they use. What differentiates the incentives of the attacker from the designer’s is that the latter is utmost interested in preserving the “health” of the blockchain, which we also reflect in its utility definition. Implicit in our formulation is the assumption that the attacker does not gain reward from attacking the system, unless this attack has a financial gain.⁵

Interestingly, in order to apply the RPD methodology to Bitcoin we need to extend it in non-trivial ways, to capture for example non-zero-sum games—as the utility of the designer and the attacker are not necessarily opposites—and to provide stronger notions of security and stability. In more detail, we introduce the notion of *strong attack payoff security*, which mandates that the attacker will stick to playing a passive strategy, i.e., stick to Bitcoin (but might abuse the adversary’s power to delay messages in the network). We also introduce the natural notion of *incentive compatibility* (IC) which mandates that both the attacker and the designer will have their parties play the given protocol. Observe that incentive compatibility trivially implies strong attack payoff security, and the latter implies the standard attack payoff security from the original RPD framework assuming the protocol is at least correct when no party deviates. These extensions to RPD widen its applicability and might therefore be of independent interest. We note that although we focus on analysis of Bitcoin here, the developed methodology can be adapted to analyze other main-stream cryptocurrencies.

Having laid out the model, we then use it to analyze Bitcoin. We start our analysis with the simpler case where the utilities do not depend on the messages—i.e., transactions—that are included into the blocks of the blockchain: when permanently inserting a block into the blockchain, a miner is just rewarded with a fixed block-reward value. This can be seen as corresponding to the Bitcoin backbone abstraction proposed in [13], but enriched with incentives to mine blocks. An interpretation of our results for this setting, listed below, is that they address blockchains that are not necessarily intended to be used as cryptocurrency ledgers. Although arguably this is not the case for Bitcoin, our analysis already reveals several surprising aspects, namely, that in this setting one does not need to rely on honest majority of computing power to ensure the quality

⁵ In particular, a fork might be provoked by the attacker only if it is expected to increase his revenue.

of the system. Furthermore, these results offer intuition on what is needed to achieve stability in the more complete case, which also incorporates transaction fees. Summarizing, we prove the following statements for this backbone-like setting, where the contents of the blocks do not influence the player’s strategies (but the rewards and costs do):

- Bitcoin is strongly attack-payoff secure, i.e., no coordinated coalition has an incentive to deviate from the protocol, provided that the rest of the parties play it. Further, this statement holds no matter how large the coalition (i.e., no matter how large the fraction of corrupt computing power) and no matter how high the conversion rate is. This means that in this backbone-like setting we can fully replace the assumption of honest majority of computing power by the above intuitive rational assumption.⁶
- If the reward for mining a block is high enough so that mining is on average profitable, then the Bitcoin protocol is even incentive-compatible with respect to local deviations. In other words, not only colluding parties (e.g., mining pools) do not have an incentive to deviate, but also the honest miners have a clear incentive to keep mining. Again, this makes no honest-majority assumption. Furthermore, as a sanity check, we also prove that this is not true if the conversion rate drops so that miners expect to be losing revenue by mining. The above confirms the intuition that after the initial bootstrapping phase where value is poured into the system (i.e., CR becomes large enough), such a ledger will keep working according to its specification for as long as the combination of conversion rate and block-reward is high enough.

With the intuition gained from the analysis in the above idealized setting, we next turn to the more realistic setting which closer captures Bitcoin, where block contents are messages that have an associated fee. We refer to these messages as *transactions*, and use the standard restrictions of Bitcoin on the transaction fees: every transaction has a maximum fee and the fee is a multiple of the minimum division.⁷ We remark that in all formal analyses [2, 13, 27] the transactions are considered as provided as inputs by an explicit environment that is supposed to capture the application layer that sits on top of the blockchain and uses it. As such, the environment will also be responsible for the choice of transaction fees and the distribution of transactions to the miners. For most generality, we do not assume as in [13, 27] that all transactions are communicated by the environment to all parties via a broadcast-like mechanism, but rather that they are distributed (i.e., input) by the environment to the miners, individually, who might then forward them using the network (if they are honest) or not. This more realistic transaction-submission mechanism is already explicit in [2].

We call this model that incorporates both mining rewards and transaction fees into the reward of the miner for a block as the *full-reward* model. Interestingly, this model allows us to also make predictions about the Bitcoin era when

⁶ It should be noted though that our analysis considers, similarly to [2, 13, 27], a fixed difficulty parameter. The extension to variable difficulty is left as future research.

⁷ For Bitcoin the minimum division is 1 satoshi = 10^{-8} BTC, and there is typically a cap on fees [3].

the rewards for mining a block will be much smaller than the transaction fees (or even zero).

We stress that transactions in our work are dealt with as messages that have an explicit fee associated with them, rather than actions which result in transferring BTCs from one miner to another. This means that other than its associated fee, the contents of a transaction does not affect the strategies of the players in the attack game. This corresponds to the assumption that the miners, who are responsible for maintaining the ledger, are different than the users, which, for example, translate the contents of the blocks as exchanges of cryptocurrency value, and which are part of the application/environment. We refer to this assumption as *the miners/users separation principle*. This assumption is explicit in all existing works, and offers a good abstraction to study the incentives for maintaining the ledger—which is the scope of our work—separately from the incentives of users to actually use it. Note that this neither excludes nor trivially deters “forking” by a sufficiently powerful (e.g., 2/3 majority) attacker; indeed, if some transaction fees are much higher than all others, then such an attacker might fork the network by extending both the highest and the second highest chain with the same block containing these high-fee transactions, and keep it forked for sufficiently long until he cashes out his rewards from both forks.

In this full-reward model, we prove the following statements:

- First, we look at the worst-case environment, i.e., the one that helps the adversary maximize its expected revenue. We prove that in this model Bitcoin is still incentive compatible, hence also strongly attack payoff secure. In fact, the same is true if the environment makes sure that there is a sufficient supply of transactions to the honest miners and to the adversary, such that the fees are high enough to build blocks that reach exactly the maximal rewarding value (note that not necessarily the same set of transactions have to be known to the participants). For example, as long as many users submit transactions with the heaviest possible fee (so-called *full-fee transactions*), then the system is guaranteed to work without relying on an honest majority of miners. In a sense, the users can control the stability of the system through transaction fees.
- Next, we investigate the question of whether or not the above is true for arbitrary transaction-fee distributions. Not surprisingly, the answer here is negative, and the protocol is not even attack-payoff secure (i.e. does not even achieve its specification). The proof of this statement makes use of the above sketched forking argument. On the positive side, our proof suggests that in the honest-majority setting where forking is not possible (except with negligible probability), the only way the adversary is incentivized to deviate from the standard protocol is to withhold the transactions he is mining on to avoid risking to lose the fees to honest parties.

Interpreting the above statements, we can relax the assumption for security of Bitcoin from requiring an honest majority to requiring long-enough presence of sufficiently many full-fee transactions, with a fallback to honest majority.

- Finally, observing that the typically large pool of transactions awaiting validation justifies the realistic assumption that there is enough supply to the network (and given the high adoption, this pool will not become small too fast), we can directly use our analysis, to propose a possible modification which would help Bitcoin, or other cryptocurrencies, to ensure incentive compatibility (hence also strong attack-payoff security) in the full-reward model in the long run: The idea is to define an exact cumulative amount on fees (or overall reward) to be allowed for each block. If there are enough high-fee transactions, then the blocks are filled up with transactions until this amount is reached. As suggested by our first analysis with a simple incentive structure, ensuring that this cap is non-decreasing would be sufficient to argue about stability; however, it is well conceivable that such a bound could be formally based on supply-and-demand in a more complex and economy-driven incentive structure and an interesting future research direction is to precisely define such a proposal together with the (economical) assumptions on which the security statements are based. We note that the introduction of such a rule would typically only induce a “soft fork,” and would, for a high-enough combination of conversion rate and reward bound, ensure incentive compatibility even when the flat reward per block tends to zero and the main source of rewards would be transaction fees, as it is the plan for the future of Bitcoin.

2 Preliminaries

In this section we introduce some notation and review the basic concepts and definitions from the literature, in particular from [11] and [2] that form the basis of our treatment. For completeness, an expanded version of this review can be found in the full version [1]. Our definitions use and build on the simulation-based security definition by Canetti [6]; we assume some familiarity with its basic principles.

Throughout this work we will assume an (at times implicit) security parameter κ . We use ITM to denote the set of *probabilistic polynomial time (PPT)* interactive Turing machines (ITMs). We also use the standard notions of *negligible*, *noticeable*, and *overwhelming* (e.g., see [16]) where we denote negligible (in κ) functions as $\text{negl}(\kappa)$. Finally, using standard UC notation we denote by $\text{EXEC}_{\Pi, \mathcal{A}, \mathcal{Z}}$ (resp. $\text{EXEC}_{\mathcal{F}, \mathcal{S}, \mathcal{Z}}$) the random variable (ensemble if indexed by κ) corresponding to the output of the environment \mathcal{Z} witnessing an execution of protocol Π against adversary \mathcal{A} (resp. an ideal evaluation of functionality \mathcal{F} with simulator \mathcal{S}).

2.1 The RPD Framework

The RPD framework [11] captures incentive-driven adversaries by casting attacks as a *meta-game* between two rational players, the protocol designer D and the attacker A , which we now describe. The game is parameterized by a (multi-party) functionality \mathcal{F} known to both agents D and A which corresponds to the ideal

goal the designer is trying to achieve (and the attacker to break). Looking ahead, when we analyze Bitcoin, \mathcal{F} will be a ledger functionality (cf. [2]). The designer D chooses a PPT protocol Π for realizing the functionality \mathcal{F} from the set of all probabilistic and polynomial-time (PPT) computable protocols.⁸ D sends Π to A who then chooses a PPT adversary \mathcal{A} to attack protocol Π . The set of possible terminal histories is then the set of sequences of pairs (Π, \mathcal{A}) as above.

Consistently with [11], we denote the corresponding attack game by $\mathcal{G}_{\mathcal{M}}$, where \mathcal{M} is referred to as the *attack model*, which specifies all the public parameters of the game, namely: (1) the functionality, (2) the description of the relevant action sets, and (3) the utilities assigned to certain actions (see below).

Stability in RPD corresponds to a refinement of a *subgame-perfect equilibrium* (cf. [26, Definition 97.2]), called ϵ -*subgame perfect equilibrium*, which considers as solutions profiles in which the parties' utilities are ϵ -close to their best-response utilities (see [11] for a formal definition). Throughout this paper, we will only consider $\epsilon = \text{negl}(\kappa)$; in slight abuse of notation, we will refer to $\text{negl}(\kappa)$ -*subgame perfect equilibrium* simply as *subgame perfect*.

The utilities. The core novelty of RPD is in how utilities are defined. Since the underlying game is zero-sum, it suffices to define the attacker's utility. This utility depends on the goals of the attacker, more precisely, the security breaches which he succeeds to provoke, and is defined, using the simulation paradigm, via the following three-step process:

First, we modify the ideal functionality \mathcal{F} to obtain a (possibly weaker) ideal functionality $\langle \mathcal{F} \rangle$, which explicitly allows the attacks we wish to model. For example, $\langle \mathcal{F} \rangle$ could give its simulator access to the parties' inputs. This allows to score attacks that aim at input-privacy breaches.

Second we describe a scoring mechanism for the different breaches that are of interest to the adversary. Specifically, we define a function v_A mapping the joint view of the relaxed functionality $\langle \mathcal{F} \rangle$ and the environment \mathcal{Z} to a real-valued *payoff*. This mapping defines the random variable (ensemble) $v_A^{\langle \mathcal{F} \rangle, \mathcal{S}, \mathcal{Z}}$ as the result of applying v_A to the views of $\langle \mathcal{F} \rangle$ and \mathcal{Z} in a random experiment describing an ideal evaluation with ideal-world adversary \mathcal{S} ; in turn, $v_A^{\langle \mathcal{F} \rangle, \mathcal{S}, \mathcal{Z}}$ defines the *attacker's (ideal) expected payoff* for simulator \mathcal{S} and environment \mathcal{Z} , denoted by $U_{\mathcal{I}_A}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z})$, so the expected value of $v_A^{\langle \mathcal{F} \rangle, \mathcal{S}, \mathcal{Z}}$. The triple $\mathcal{M} = (\mathcal{F}, \langle \mathcal{F} \rangle, v_A)$ constitutes the *attack model*.

The third and final step is to use $U_{\mathcal{I}_A}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z})$ to define the attacker's utility, $u_A(\Pi, \mathcal{A})$, for playing an adversary \mathcal{A} against protocol Π , as the expected payoff of the “best” simulator that successfully simulates \mathcal{A} in its (\mathcal{A} 's) favorite environment. This best simulator is the one that translates the adversary's breaches against Π into breaches against the relaxed functionality $\langle \mathcal{F} \rangle$ in a faithful manner, i.e., so that the ideal breaches occur only if the adversary really makes them necessary for the simulator in order to simulate. As argued in [11], this corresponds to the simulator that minimizes the attacker's utility. Formally, for a functionality $\langle \mathcal{F} \rangle$ and a protocol Π , denote by $\mathcal{C}_{\mathcal{A}}$ the class of simulators that

⁸ Following standard UC convention, the protocol description includes its hybrids.

are “good” for \mathcal{A} , i.e., $\mathcal{C}_{\mathcal{A}} = \{\mathcal{S} \in \text{ITM} \mid \forall \mathcal{Z} : \text{EXEC}_{\Pi, \mathcal{A}, \mathcal{Z}} \approx \text{EXEC}_{\langle \mathcal{F} \rangle, \mathcal{S}, \mathcal{Z}}\}$.⁹ Then the attacker’s (expected) utility is defined as:

$$u_{\mathcal{A}}(\Pi, \mathcal{A}) = \sup_{\mathcal{Z} \in \text{ITM}} \left\{ \inf_{\mathcal{S} \in \mathcal{C}_{\mathcal{A}}} \left\{ U_{I^{\mathcal{A}}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}) \right\} \right\}.$$

For \mathcal{A} and Π with $\mathcal{C}_{\mathcal{A}} = \emptyset$, the utility is ∞ by definition, capturing the fact that we only want to consider protocols which at the very least implement the relaxed (i.e., explicitly breachable) functionality $\langle \mathcal{F} \rangle$. Note that as the views in the above experiments are in fact random variable ensembles indexed by the security parameter κ , the probabilities of all the relative events are in fact functions of κ , hence the utility is also a function of κ . Note also that as long as $\mathcal{C}_{\mathcal{A}} = \emptyset$ is non-empty, for each value of κ , both the supremum and the infimum above exist and are finite and reachable by at least one pair $(\mathcal{S}, \mathcal{Z})$, provided the scoring function assigns finite payoffs to all possible transcripts (for $\mathcal{S} \in \mathcal{C}_{\mathcal{A}}$) (cf. [11]).

Remark 1 (Event-based utility [11]). In many applications, including those in our work, meaningful payoff functions have the following, simple representation: Let (E_1, \dots, E_ℓ) denote a vector of (typically disjoint) events defined on the views (of \mathcal{S} and \mathcal{Z}) in the ideal experiment corresponding to the security breaches that contribute to the attacker’s utility. Each event E_i is assigned a real number γ_i , and the payoff function $v_{\mathcal{A}}^{\vec{\gamma}}$ assigns, to each ideal execution, the sum of γ_i ’s for which E_i occurred. The ideal expected payoff of a simulator is computed according to our definition as

$$U_{I^{\mathcal{A}}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}) = \sum_{E_i \in \vec{E}, \gamma_i \in \vec{\gamma}} \gamma_i \Pr[E_i],$$

where the probabilities are taken over the random coins of \mathcal{S} , \mathcal{Z} , and $\langle \mathcal{F} \rangle$.

Building on the above definition of utility, [11] introduces a natural notion of security against incentive-driven attackers. Intuitively, a protocol Π is *attack-payoff secure* in a given attack model $\mathcal{M} = (\mathcal{F}, \cdot, v_{\mathcal{A}})$, if the utility of the best adversary against this protocol is the same as the utility of the best adversary in attacking the \mathcal{F} -hybrid “dummy” protocol, which only relays messages between \mathcal{F} and the environment.

Definition 1 (Attack-payoff security [11]). *Let $\mathcal{M} = (\mathcal{F}, \langle \mathcal{F} \rangle, v_{\mathcal{A}}, v_{\mathcal{D}})$ be an attack model inducing utilities $u_{\mathcal{A}}$ and $u_{\mathcal{D}}$ on the attacker and the designer, respectively,¹⁰ and let $\phi^{\mathcal{F}}$ be the dummy \mathcal{F} -hybrid protocol. A protocol Π is attack-payoff secure for \mathcal{M} if for all adversaries $\mathcal{A} \in \text{ITM}$,*

$$u_{\mathcal{A}}(\Pi, \mathcal{A}) \leq u_{\mathcal{A}}(\phi^{\mathcal{F}}, \mathcal{A}) + \text{negl}(\kappa).$$

⁹ This class is finite for every given value of the security parameter, Π , and \mathcal{A} .

¹⁰ In [11], by default $u_{\mathcal{D}} = -u_{\mathcal{A}}$ as the game is zero-sum.

Intuitively, this security definition accurately captures security against an incentive-driven attacker, as in simulating an attack against the dummy \mathcal{F} -hybrid protocol, the simulator never needs to provoke any of the “breaching” events. Hence, the utility of the best adversary against Π equals the utility of an adversary that does not provoke any “bad event.”

2.2 A Composable Model for Blockchain Protocols

In [2], Badertscher *et al.* present a universally composable treatment of the Bitcoin protocol, $\Pi^{\mathbb{B}}$, in the UC framework. Here we highlight the basic notions and results and refer to the full version [1] for details.

The Bitcoin ledger. The ledger functionality $\mathcal{G}_{\text{LEDGER}}^{\mathbb{B}}$ maintains a ledger state **state**, which is a sequence of state blocks. A state block contains (application-specific) content values—the “transactions.” For each honest party p_i , the ledger stores a pointer to a state block—the head of the state from p_i ’s point of view—and ensures that pointers increase monotonically and are not too far away from the head of the state (and that it only moves forward). Parties or the adversary might submit transactions, which are first validated by means of a predicate $\text{ValidTx}_{\mathbb{B}}$, and, if considered valid, are added to the functionality’s buffer. At any time, the $\mathcal{G}_{\text{LEDGER}}^{\mathbb{B}}$ allows the adversary to propose a candidate next-block for the state. However, the ledger enforces a specific *extend policy* specified by an algorithm ExtendPolicy that checks whether the proposal is compliant with the policy. If the adversary’s proposal does not comply with the ledger policy, ExtendPolicy rejects the proposal. The policy enforced by the Bitcoin ledger can be succinctly summarized as follows:

- *Ledger’s growth.* Within a certain number of rounds the number of added blocks must not be too small or too large.
- *Chain quality.* A certain fraction of the proposed blocks must be mined honestly and those blocks satisfy special properties (such as including all recent transactions).
- *Transaction liveness.* Old enough (and valid) transactions are included in the next block added to the ledger state.

The Bitcoin protocol. In [2] it was proved that (a [13]-inspired abstraction of) Bitcoin as a synchronous-UC protocol [20], called the *ledger protocol* and denoted by $\Pi^{\mathbb{B}}$, realizes the above ledger. $\Pi^{\mathbb{B}}$ uses blockchains to store a sequence of transactions. A *blockchain* \mathcal{C} is a (finite) sequence of blocks $\mathbf{B}_1, \dots, \mathbf{B}_\ell$. Each *block* \mathbf{B}_i consist of a *pointer* \mathbf{s}_i , a *state block* \mathbf{st}_i , and a *nonce* \mathbf{n}_i . string. The chain $\mathcal{C}^{\uparrow k}$ is \mathcal{C} with the last k blocks removed. The *state* $\vec{\mathbf{st}}$ of the blockchain $\mathcal{C} = \mathbf{B}_1, \dots, \mathbf{B}_\ell$ is defined as a sequence of its state blocks, i.e., $\vec{\mathbf{st}} := \mathbf{st}_1 || \dots || \mathbf{st}_\ell$.

The validity of a blockchain $\mathcal{C} = \mathbf{B}_1, \dots, \mathbf{B}_\ell$ where $\mathbf{B}_i = \langle \mathbf{s}_i, \mathbf{st}_i, \mathbf{n}_i \rangle$ is decided by a predicate $\text{isvalidchain}_D(\mathcal{C})$. It combines two types of validity: *chain-level*, aka syntactic, validity—which, intuitively requires that valid blocks need to be solving a proof-of-work-type puzzle for a hash function $\mathbf{H} : \{0, 1\}^* \rightarrow \{0, 1\}^\kappa$

and difficulty \mathfrak{d} —and *state-level*, aka semantic, validity, which specifies whether the block’s contents, i.e., transactions, are valid, with respect to a blockchain-specific predicate.

The Bitcoin protocol $\Pi^{\mathfrak{B}}$ is executed in a hybrid world where parties have access to a random oracle functionality \mathcal{F}_{RO} (modeling the hash function \mathbf{H}), a multicast asynchronous network using channels with bounded delay $\mathcal{F}_{\text{N-MC}}$, and a global clock $\mathcal{G}_{\text{CLOCK}}$. Each party maintains a (local) current blockchain. It receives the transactions from the environment (and circulates them), and adds newly received valid transactions to a block that is then mined-on using the algorithm extendchain_D . The idea of the algorithm is to find a proof of work—by querying the random oracle \mathcal{F}_{RO} —which allows to extend the local chain with a valid block. After each mining attempt the party uses the network to multicast their current blockchain. Parties always adopt the longest chain that they see starting from a pre-agreed genesis block. The protocol (implicitly) defines the ledger state to be a certain prefix of the contents of the longest chain held by each party. More specifically, if a party holds a valid chain \mathcal{C} that encodes the sequence of state blocks $\vec{\mathfrak{st}}$, then the ledger state is defined to be $\vec{\mathfrak{st}}^{\lceil T}$, i.e., the party outputs a prefix of the encoded state blocks of its local longest chain. T is chosen such that honest parties output a consistent ledger state.

The flat model of computation. In this paper, we state the results in the synchronous flat model (with fixed difficulty) by Garay *et al.* [13]. This means we assume a number of parties, denoted by n , that execute the Bitcoin protocol $\Pi^{\mathfrak{B}}$, out of which t parties can get corrupted. For simplicity, the network $\mathcal{F}_{\text{N-MC}}$ guarantees delivery of messages sent by honest parties in round r to be available to any other party at the onset of round $r + 1$. Moreover, every party will be invoked in every round and can make at most one “calculation” query to the random oracle \mathcal{F}_{RO} in every round (and an unrestricted number of “verification” queries to check the validity of received chains)¹¹, and use the above diffusion network $\mathcal{F}_{\text{N-MC}}$ once in a round to send and receive messages. To capture these restrictions in a composable treatment, the real-world assumptions are enforced by means of a “wrapper” functionality, $\mathcal{W}_{\text{flat}}$, which adequately restricts access to $\mathcal{G}_{\text{CLOCK}}$, \mathcal{F}_{RO} and $\mathcal{F}_{\text{N-MC}}$ as explained in [2].

Denote by ρ the fraction of dishonest parties (i.e., $t = \rho \cdot n$) and define $p := \frac{\mathfrak{d}}{2^n}$ which is the probability of finding a valid proof of work via a fresh query to \mathcal{F}_{RO} (where \mathfrak{d} is fixed but sufficiently small, depending on n). Let $\alpha^{\text{flat}} = 1 - (1 - p)^{(1-\rho) \cdot n}$ be the mining power of the honest parties, and $\beta^{\text{flat}} = p \cdot (\rho \cdot n)$ be the mining power of the adversary.

Theorem 1. *Consider $\Pi^{\mathfrak{B}}$ in the $\mathcal{W}_{\text{flat}}(\mathcal{G}_{\text{CLOCK}}, \mathcal{F}_{\text{RO}}, \mathcal{F}_{\text{N-MC}})$ -hybrid world. If, for some $\lambda > 1$, the honest-majority assumption*

$$\alpha^{\text{flat}} \cdot (1 - 4\alpha^{\text{flat}}) \geq \lambda \cdot \beta^{\text{flat}}$$

¹¹ This fine-grained round model with one hash query was already used by Pass *et al.* [27]. The extension to a larger, constant upper bound of calculation queries per round as in [13] is straightforward for the results in this work.

holds in any real-world execution, then protocol Π^B UC-realizes $\mathcal{G}_{\text{LEDGER}}^B$ for some specific range of parameters (given in [2]).

3 Rational Protocol Design of Ledgers

In this section we present our framework for rational analysis of the Bitcoin protocol. It uses as basis the framework for *rational protocol design* (and analysis—RPD framework for short) by Garay *et al.* [11], extending it in various ways to better capture Bitcoin’s features. (We refer to Sect. 2 and to the full version for RPD’s main components and security definitions.) We note that although our analysis mainly focuses on Bitcoin, several of the extensions have broader applicability, and can be used for the rational analysis of other cryptocurrencies as well.

RPD’s machinery offers the foundations for capturing incentive-driven attacks against multi-party protocols for a given specification. In this section we show how to tailor this methodology to the specific task of protocols aimed to securely implement a public ledger. The extensions and generalizations of the original RPD framework we provide add generic features to the RPD framework, including the ability to capture non-zero-sum attack games—which, as we argue, are more suitable for the implementation of a cryptocurrency ledger—and the extension of the class of events which yield payoff to the attacker and the designer.

The core hypothesis of our rational analysis is that the incentives of an attacker against Bitcoin—which affect his actions and attacks—depend only on the possible earnings or losses of the parties that launch the attack. We do not consider, for example, attackers that might create forks just for the fun of it. An attacker might create a “fork” in the blockchain if he expects to gain something by doing so. In more detail, we consider the following events that yield payoff (or inflict a cost) for running the Bitcoin protocol:

- *Inserting a block into the blockchain.* It is typical of cryptocurrencies that when a party manages to insert a block into the ledger’s state, then it is rewarded for the effort it invested in doing so. In addition, it is typical in such protocols that the contents of the blocks (usually transactions) have some *transaction fee* associated with them. (For simplicity, in our initial formalization (Sects. 3 and 4) we will ignore transaction fees in our formal statements, describing how they are extended to also incorporate also such fees in Sect. 5.)
- *Spending resources to mine a block.* These resources might be the electricity consumed for performing the mining, the investment on mining hardware and its deterioration with time, etc.

Remark 2 (The miners/users separation principle). We remark that the scope of our work is to analyze the security of cryptocurrencies against incentive-driven attacks by the miners, i.e., the parties that are responsible for maintaining the blockchain. In particular, consistently with [2, 13, 27] we shall consider the inputs to the protocol as provided by a (not-necessarily rational) environment, which in particular captures the users of the system. As a result, other than the

transaction fees, we will assume that the contents of the ledger do not affect the miners’ strategies, which we will refer to as the *miners/users separation principle*. This principle captures the case where the users do not collude with the miners—an assumption implicit in the above works. We leave the full rational analysis of the protocol, including application layers for future research.

There are several challenges that one needs to overcome in deriving a formal treatment of incentive-driven attacks against Bitcoin. First, the above reward and cost mechanisms are measured in different “units.” Specifically, the block reward is a cryptocurrency convention and would therefore be measured in the specific cryptocurrency’s units, e.g., BTCs in the case of the Bitcoin network. On the other hand, the cost for mining (e.g., the cost of electricity, equipment usage, etc.) would be typically measured in an actual currency. To resolve this mismatch—and refrain from adopting a specific currency—we introduce a variable **CR** which corresponds to the *conversion rate* of the specific cryptocurrency unit (e.g., BTCs) to the cost unit (e.g., euros or US dollars). As we shall see in the next section, using such an explicit exchange rate allows us to make statements about the quality of the Bitcoin network that depend on its price—as they intuitively should. For example, we can formally confirm high-level statements of the type: “Bitcoin is stable—i.e., miners have incentive to keep mining honestly—as long as its price is high enough” (cf. Sect. 4).

Furthermore, this way we can express all payoffs in terms of cost units: Assume that it takes r rounds for a miner (or a collection of miners) to insert a block into the state. Denote by **mcost** the cost for a single mining attempt (in our case a single RO query), and by **breward** the fraction of cryptocurrency units (e.g., BTCs) that is given as a reward for each mined block.¹² Then, the payoff for the insertion of a single block is $\mathbf{breward} \cdot \mathbf{CR} - q_r \cdot \mathbf{mcost}$, where q_r is the number of queries to the RO that were required to mine this block during r rounds.

The second challenge is with respect to *when* should a miner receive the reward for mining. There are several reasons why solving a mining puzzle—thereby creating a new block—does not necessary guarantee a miner that he will manage to insert this block into the blockchain, and therefore be rewarded for it, including the possibility of collisions—more than one miner solving the puzzle—or, even worse, adversarial interference—e.g., network delays or “selfish mining.” And even if the miner is the only one to solve the puzzle in a given round, he should only be rewarded for it if his block becomes part of the (permanent) state of the blockchain—the so-called blockchain’s “common prefix.”

To overcome this second challenge we rely on the RPD methodology. In particular, we will use the ideal experiment where parties have access to the global ledger functionality, where we can clearly identify the event of inserting a block into the state, and decide, by looking into the state, which miner added which block.¹³

¹² Currently, for the Bitcoin network, this is 1/4 of the original reward (12.5 BTCs).

¹³ In [2], each block of the state includes the identifier of the miner who this block is attributed to.

In order to formalize the above intuitions and apply the RPD methodology to define the utilities in the attack game corresponding to implementing a ledger against an incentive-driven adversary, we need to make some significant adaptations and extensions to the original framework, which is what we do next. We then (Sect. 3.2) use the extended framework to define the attack-model for the Bitcoin protocol, and conclude the section by giving appropriate definitions of security and stability in this model.

3.1 Extending the RPD Framework

We describe how to extend the model from [11] to be able to use it in our context.

Black-box simulators. The first modification is adding more flexibility to how utilities are defined. The original definition of ideal payoff $U_{IA}^{(\mathcal{F})}(\mathcal{S}, \mathcal{Z})$ computes the payoff of the simulator using the joint view of the environment and the functionality. This might become problematic when attempting to assign cost to resources used by the adversary—the RO queries in our scenario, for example. Indeed, these queries are not necessarily in this joint view, as depending on the simulator, one might not be able to extract them.¹⁴ To resolve this we modify the definition to restrict it to black-box simulators, resulting in \mathcal{C}_A being the class of simulators that use the adversary as a black box. This will ensure that the queries to the RO are part of the interaction of the simulator with its adversary, and therefore present in the view of the simulator. Further, we include this part of the simulator’s view in the definition of the scoring function v_A , which is defined now as a mapping from the joint view of the relaxed functionality $\langle \mathcal{F} \rangle$, the environment \mathcal{Z} , and the simulator \mathcal{S} to a real-valued *payoff*.

Non-zero-sum attack games. The second modification is removing the assumption that the attack game is zero-sum. Indeed, the natural incentive of the protocol designer in designing a Ledger protocol is not to optimally “tame” its attacker—as in [11]—but rather to maximize the revenue of the non-adversarially controlled parties while keeping the blockchain healthy, i.e., free of forks. This is an important modification as it captures attacks in which the adversary preserves his rate of blocks inserted into the state, but slows down the growth of the state to make sure that honest miners accrue less revenue in any time interval. For example, the so called “selfish-mining” strategy [9] provokes a slowdown since honest mining power is invested into mining on a chain which is not the longest one (as the longest chain is kept private as long as possible by the party that does the selfish-mining).

To formally specify the utility of the designer in such a non-zero-sum attack game, we employ a similar reasoning as used in the original RPD framework for defining the attacker’s utility. The first step, relaxing the functionality, can be omitted provided that we relaxed it sufficiently in the definition of the attacker’s utility. In the second step, we define the scoring mechanism for the incentives

¹⁴ Indeed, in the ideal simulation of the Bitcoin protocol presented in [2], there is no RO in the ideal world.

of the designer as a function $v_{\mathcal{D}}$ mapping the joint view of the relaxed functionality $\langle \mathcal{F} \rangle$, the environment \mathcal{Z} , and the simulator \mathcal{S} to a real-valued *payoff*, and define the designer’s (*ideal*) *expected payoff* for simulator \mathcal{S} with respect to the environment \mathcal{Z} as

$$U_{\mathcal{D}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}) = E(v_{\mathcal{D}}^{\langle \mathcal{F} \rangle, \mathcal{S}, \mathcal{Z}}),$$

where $v_{\mathcal{D}}^{\langle \mathcal{F} \rangle, \mathcal{S}, \mathcal{Z}}$ describes (as a random variable) the payoff of \mathcal{D} allocated by \mathcal{S} in an execution using directly the functionality $\langle \mathcal{F} \rangle$.

The third and final step is the trickiest. Here we want to use the above ideal expected payoff to define the expected payoff of a designer using protocol Π when the attacker is playing adversary \mathcal{A} . In order to ensure that our definition is consistent with the original definition in [11]—which applied to (only) zero-sum games—we need to make sure that the utility of the designer increases as the utility of the attacker decreases and vice versa. Thus, to assign utility for the designer to a strategy profile (Π, \mathcal{A}) , we will use the same simulators and environments that were used to assign the utility for the attacker. Specifically, let $\mathbb{S}_{\mathcal{A}}$ denote the class of simulators that are used to formulate the utility of the adversary, and let $\mathbb{Z}_{\mathcal{A}}$ denote the class of environments that maximize this utility for simulators in $\mathbb{S}_{\mathcal{A}}$ ¹⁵, then

$$\mathbb{S}_{\mathcal{A}} = \left\{ \mathcal{S} \in \mathcal{C}_{\mathcal{A}} \text{ s.t. } \sup_{\mathcal{Z} \in \text{ITM}} \{U_{\mathcal{I}_{\mathcal{A}}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z})\} = u_{\mathcal{A}}(\Pi, \mathcal{A}) \right\} \quad (1)$$

and

$$\mathbb{Z}_{\mathcal{A}} = \left\{ \mathcal{Z} \in \text{ITM} \text{ s.t. for some } \mathcal{S} \in \mathbb{S}_{\mathcal{A}} : U_{\mathcal{I}_{\mathcal{A}}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z})\} = u_{\mathcal{A}}(\Pi, \mathcal{A}) \right\}. \quad (2)$$

It is easy to verify that this choice of simulator respects the utilities being opposite in a zero-sum game as defined in [11], thereby preserving the results following the original RPD paradigm.

Lemma 1. *Let $v_{\mathcal{D}} = -v_{\mathcal{A}}$ and let $U_{\mathcal{D}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z})$ defined as above. For some $\mathcal{S} \in \mathbb{S}_{\mathcal{A}}$ and some $\mathcal{Z} \in \mathbb{Z}_{\mathcal{A}}$, define $u_{\mathcal{D}}(\Pi, \mathcal{A}) := U_{\mathcal{D}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z})$. Then $u_{\mathcal{D}}(\Pi, \mathcal{A}) = -u_{\mathcal{A}}(\Pi, \mathcal{A})$.*

Proof. Since $v_{\mathcal{D}} = -v_{\mathcal{A}}$, we have that for all $\mathcal{Z}, \mathcal{S} \in \text{ITM}$,

$$U_{\mathcal{D}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}) = -U_{\mathcal{I}_{\mathcal{A}}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}). \quad (3)$$

However, by definition, since $\mathcal{S} \in \mathbb{S}_{\mathcal{A}}$, we have

$$u_{\mathcal{A}}(\Pi, \mathcal{A}) = U_{\mathcal{I}_{\mathcal{A}}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}) \stackrel{3}{=} -U_{\mathcal{D}}^{\langle \mathcal{F} \rangle}(\mathcal{S}, \mathcal{Z}) = -u_{\mathcal{D}}(\Pi, \mathcal{A}).$$

□

The above lemma confirms that for a zero-sum attack game we can take any pair $(\mathcal{S}, \mathcal{Z}) \in \mathbb{S}_{\mathcal{A}} \times \mathbb{Z}_{\mathcal{D}}$ in the definition of $u_{\mathcal{D}}(\Pi, \mathcal{A})$ and it will preserve the zero-sum property (and hence all the original RPD results). This is so because all these

¹⁵ Recall that as argued in Sect. 2.1, these sets are non-empty provided $\mathcal{C}_{\mathcal{A}} \neq \emptyset$.

simulators induce the same utility $-u_A(\Pi, \mathcal{A})$ for the designer. However, for our case of non-zero-sum games, each of those simulator/environment combinations might induce a different utility for the designer. To choose the one which most faithfully translates the designer’s utility from the real to the ideal world we use the same line of argument as used in RPD for defining the attacker’s utility: The best (i.e., the most faithful) simulator is the one which always rewards the designer whenever his protocol provokes some profitable event; in other words, the one that maximizes the designer’s expected utility. Similarly, the natural environment is the one that puts the protocol in its worst possible situation, i.e., the one that minimizes its expected gain; indeed, such an environment will ensure that the designer is guaranteed to get his allocated utility. The above leads to the following definition for the designer’s utility in non-zero-sum games:

$$u_D(\Pi, \mathcal{A}) := \inf_{\mathcal{Z} \in \mathbb{Z}_A} \left\{ \sup_{\mathcal{S} \in \mathbb{S}_A} \left\{ U_{I^p}^{(\mathcal{F})}(\mathcal{S}, \mathcal{Z}) \right\} \right\}.$$

For completeness, we set $u_D(\Pi, \mathcal{A}) = -\infty$ if $\mathcal{C}_A = \emptyset$, i.e., if the protocol does not even achieve the relaxed functionality. This is not only intuitive—as $\mathcal{C}_A = \emptyset$ means that the designer chose a protocol which does not even reach the relaxed goal—but also analogous to how RPD defines the attacker’s utility for protocols that do not achieve their relaxed specification.¹⁶

Finally, the attack model for non-zero-sum games is defined as the quadruple $\mathcal{M} = (\mathcal{F}, \langle \mathcal{F} \rangle, v_A, v_D)$.

3.2 Bitcoin in the RPD Framework

Having formulated the above extensions to the RPD framework, we are ready to apply the methodology to analyze Bitcoin.

Basic foundations. We explain in more depth on how to implement the core steps of RPD. First, we define the Ledger functionality from [2] as Bitcoin’s ideal goal (see Sect. 2.2). Following the three steps of the methodology, we start by defining the relaxed version of the Ledger, denoted as $\mathcal{G}_{\text{WEAK-LEDGER}}^{\text{B}}$. Informally, the relaxed Ledger functionality operates as the original ledger with the following modifications:

The state is a tree: Instead of storing a single ledger state `state` as a straight-line blockchain-like structure, $\mathcal{G}_{\text{WEAK-LEDGER}}^{\text{B}}$ stores a tree `state-tree` of state blocks where for each node the direct path from the root defines a possible ledger state that might be presented to any of the honest miners. The functionality maintains for each registered party $p_i \in \mathcal{P}$ a pointer `pti` to a node in the tree which defines p_i ’s current-state view. Furthermore, instead of restricting the adversary to only be able to set the state “slackness” to be not larger than a specific parameter, $\mathcal{G}_{\text{WEAK-LEDGER}}^{\text{B}}$ offers the command `SET-POINTER` which allows the adversary to set the pointers of honest parties

¹⁶ Recall that RPD sets $u_A(\Pi, \mathcal{A}) = \infty$ if \mathcal{A} cannot be simulated, i.e., if $\mathcal{C}_A = \emptyset$.

within `state-tree` with the following restriction: The pointer of an honest party can only be set to a node whose distance to the root is at least the current-pointer node’s.

Relaxed validity check of transactions: All submitted transactions are accepted into the buffer `buffer` without validating against `state-tree`. Moreover, transactions in `buffer` which are added to `state-tree` are not removed as they could be reused at another branch of `state-tree`.

Ability to create forks: This relaxation gives the simulator the explicit power to create a fork on the ledger’s state. This is done as follows: The command `NEXT-BLOCK`—which, recall, allows the simulator to propose the next block—is modified to allow the simulator to extend an arbitrary leaf of a sufficiently long rooted path of `state-tree`. Thus, when `state-tree` is just a single path, this command operates as in the original ledger from [2]. Additionally, in the relaxed ledger, the simulator is also allowed to add the next block to an intermediate, i.e., non-leaf node of `state-tree`. This is done by using an extra command `FORK` which, other than extending the chain from the indicated block provides the same functionality as `NEXT-BLOCK`.

Relaxed state-extension policy: As explained in Sect. 2.2, the extend policy is a compliance check that the ledger functionality performs on blocks that the simulator proposes to be added to the ledger’s state. This is to ensure that they satisfy certain conditions. This is the mechanism which the ledger functionality uses to enforce, among others, common generic-ledger properties from the literature, such as the chain quality or the chain growth properties, and for Bitcoin ledgers the transaction-persistence/stability properties [13, 27]. of the ledger state, or on transaction persistence/stability [13]. The relaxed ledger uses a much more permissive extend policy, denoted as `weakExtendPolicy`, derived from `ExtendPolicy` with the following modifications: Intuitively, in contrast to `ExtendPolicy`, the weaker version does not check if the adversary inserts too many or too few blocks, and it does not check if all old-enough transactions have been included. There is also no check of whether enough blocks are mined by honest parties, i.e., that there are enough blocks with coin-base transactions from honest parties. In other words, `weakExtendPolicy` does not enforce any concrete bounds on the chain quality or the chain growth properties of the ledger state, or on transaction persistence/stability. It rather ensures basic validity criteria of the resulting ledger state.

More formally, instead of `state`, it takes `state-tree` and a pointer `pt` as input. It first computes a valid default block \vec{N}_{df} which can be appended at the longest branch of `state-tree`. It then checks if the proposed blocks \vec{N} can be safely appended to the node `pt` (to yield a valid state). If this is the case it returns (\vec{N}, pt) ; otherwise it returns \vec{N}_{df} and a pointer to the leaf of the longest branch in `state-tree`.

The formal description of the relaxed ledger functionality is found in the full version [1]. This completes the first step of the RPD methodology.

The second step is defining the scoring function. This is where our application of RPD considerably deviates from past works [11, 12]. In particular, those works consider attacks against generic secure multi-party computation protocols, where the ideal goal is the standard secure function evaluation (SFE) functionality (cf. [6]). The security breaches are breaking correctness and privacy [11] or breaking fairness [12]. These can be captured by relaxing the SFE functionality to allow the simulator to request extra information (breaking privacy), reset the outputs of honest parties to a wrong value (breaking correctness), or cause an abort (breaking fairness.) The payoff function is then defined by looking at events corresponding to whether or not the simulator provokes these events, and the adversary is given payoff whenever the best simulator is forced to provoke them in order to simulate the attack.

However, attacks against the ledger that have as an incentive increasing the revenue of a coalition are not necessarily specific events corresponding to the simulator sending special “break” commands. Rather, they are events that are extracted from the joint views (e.g., which blocks make it to the state and when). Hence, attacks to the ledger correspond to the simulator implicitly “tweaking” its parameters. Therefore, in this work we take the following approach to define the payoffs of the attacker and designer. In contrast to the RPD examples in [11, 12], which use explicit events that “downgrade” the ideal functionality for defining utility, we directly use more intuitive events defined on the joint view of the environment, the functionality, and the simulator. The reason is that as we have assumed that the only rationale is to increase one’s profit, the incentives in case of cryptocurrencies are as follows: whenever a block is mined, the adversary gets rewarded. A “security breach” is relevant if (and only if) the adversary can get a better reward by doing so.

Defining concrete utility functions. Defining the utility functions lies at the core of a rational analysis of a blockchain protocol like Bitcoin. The number of aspects that one would like to consider steers the complexity of a concrete analysis, the ultimate goal being to reflect exactly the incentive structure of the actual blockchain ecosystem. Our extended RPD framework for blockchain protocols provides a guideline to defining utility functions of various complexity and to conduct the associated formal analysis. Recall that the utility functions are the means to formalize the assumed underlying incentive structure. As such, our approach is extensible: if certain relevant properties or dynamics are identified or believed (such as reflecting a doomsday risk of an attacker or a altruistic motivation of honest miners), one can enrich the incentive structure by reflecting the associated events and rewards in the utility definition, or by making the costs and rewards time-dependent variables. The general goal of this line of research on rational aspects of cryptocurrencies is to eventually arrive at a more detailed model and, if the assumptions are reasonable, to have more predictive models for reality.

Below we define a first, relatively simple incentive model to concretely showcase our methodology. We conduct the associated rational analysis in the next

section and observe that, although being a simplified model, we can already draw interesting conclusions from such a treatment.

Utility of the attacker. Informally, this particular utility is meant to capture the average revenue of the attacker. Consider the following sequence of events defined on the views of the environment, the relaxed ledger functionality, and the black-box simulator of the entire experiment (i.e., until the environment halts) for a given adversary \mathcal{A} :

1. For each pair $(q, r) \in \mathbb{N}^2$ define event $W_{q,r}^{\mathcal{A}}$ as follows: The simulator makes q mining queries in round r , i.e., it receives q responses on different messages to the RO in round r .¹⁷
2. For each pair $(b, r) \in \mathbb{N}^2$ define event $I_{b,r}^{\mathcal{A}}$ as follows: The simulator inserts b blocks into the state of the ledger in round r , such that all these blocks were previously queries to the (simulated) random oracle by the adversary. More formally, $I_{b,r}^{\mathcal{A}}$ occurs if the function extend policy (of the weak ledger) is successfully invoked and outputs a sequence of b non-empty blocks (to be added to the state), where for each of these blocks the following properties hold: (1) The block has appeared in the past in the transcript between the adversary and the simulator, and (2) the contents of the block have appeared on this transcript prior to the block's first appearance, as a query from the adversary to its (simulated) RO. We note in passing that this event definition ensures that the simulator (and therefore also the adversary) does not earn reward by adaptively corrupting parties after they have done the work/query to mine a block but before their block is added into the state. In other words, the adversary only gets rewarded for state blocks which corrupted parties mined while they were already under the adversary's control.

Now, using the simplified event-based utility definition (Remark 1) we define the attacker's utility for a strategy profile (Π, \mathcal{A}) in the attack game as:¹⁸

$$u_{\mathcal{A}}^{\mathfrak{B}}(\Pi, \mathcal{A}) = \sup_{Z \in \text{ITM}} \left\{ \inf_{\mathcal{S}^{\mathcal{A}} \in \mathcal{C}_{\mathcal{A}}} \left\{ \sum_{(b,r) \in \mathbb{N}^2} b \cdot \text{breward} \cdot \text{CR} \cdot \Pr[I_{b,r}^{\mathcal{A}}] - \sum_{(q,r) \in \mathbb{N}^2} q \cdot \text{mcost} \cdot \Pr[W_{q,r}^{\mathcal{A}}] \right\} \right\}.$$

We remark that although the above sums are in principle infinite, in any specific execution these sums will have only as many (non-zero) terms as the number of rounds in the protocol. Indeed, if the experiment finishes in r' rounds then for any $r > r'$, $\Pr[I_{b,r}^{\mathcal{A}}] = \Pr[W_{q,r}^{\mathcal{A}}] = 0$ for all $b \in \mathbb{N}$. Furthermore, we assume that breward , CR and mcost are $O(1)$, i.e., independent of the security parameter.

¹⁷ Observe that since our ideal world is the $\mathcal{G}_{\text{clock}}$ -hybrid synchronous world, the round structure is trivially extracted from the simulated ideal experiment by the protocol definition and the clock value. Furthermore, the adversary's mining queries can be trivially extracted by its interaction with the black-box simulator.

¹⁸ Recall that we assume synchronous execution as in [2] where the environment gets to decide how many rounds it wishes to witness.

The above expression can be simplified to the following more useful expression. Let B^A denote the random variables corresponding to the number of blocks contributed to the ledger’s state by adversarial miners and Q^A denote the number of queries to the RO performed by adversarial miners (throughout the execution of the random experiment). Then the adversary’s utility can be described using the expectations of these random variables as follows:

$$u_A^{\mathbb{B}}(\Pi, \mathcal{A}) = \sup_{\mathcal{Z} \in \text{ITM}} \left\{ \inf_{\mathcal{S}^A \in \mathcal{C}_A} \left\{ \text{breward} \cdot \text{CR} \cdot E(B^A) - \text{mcost} \cdot E(Q^A) \right\} \right\}.$$

Utility of the designer. Since the game is not zero-sum we also need to formally specify the utility of the protocol designer. Recall that we have assumed that, analogously to the attacker, the designer accrues utility when honest miners insert a block into the state, and spends utility when mining—i.e., querying the RO. In addition, what differentiates the incentives of the designer from that of an attacker is that his most important goal is to ensure the “health” of the blockchain, i.e., to avoid forks. To capture this, we will assign a cost for the designer to the event the simulator is forced to request the relaxed ledger functionality to fork, which is larger than his largest possible gain. This yields the following events that are relevant for the designer’s utility.

1. For each pair $(q, r) \in \mathbb{N}^2$ define $W_{q,r}^{\Pi}$ as follows: The honest parties, as a set, make q mining queries in round r .¹⁹
2. For each pair $(b, r) \in \mathbb{N}^2$ define $I_{b,r}^{\Pi}$ as follows: The honest parties jointly insert b blocks into the state of the ledger in round r ; that is, the simulator inserts b blocks into the state of the ledger in round r , such that for each of these blocks, at least one of the two properties specified in the above definition of $I_{b,r}^A$ does not hold.²⁰
3. For each $r \in \mathbb{N}$ define K_r as follows: The simulator uses the FORK command in round r .

The utility of the designer is then defined similarly to the attacker’s, where we denote by \mathbb{S}_A the class of simulators that assign to the adversary his actual utility (cf. Eq. 1):

$$u_D^{\mathbb{B}}(\Pi, \mathcal{A}) = \inf_{\mathcal{Z} \in \mathbb{Z}} \left\{ \sup_{\mathcal{S}^A \in \mathbb{S}_A} \left\{ \sum_{(b,r) \in \mathbb{N}^2} b \cdot \text{CR} \cdot (\text{breward} \cdot \Pr[I_{b,r}^{\Pi}] - 2^{\text{polylog}(\kappa)} \cdot \Pr[K_r]) - \sum_{(q,r) \in \mathbb{N}^2} q \cdot \text{mcost} \cdot \Pr[W_{q,r}^{\Pi}] \right\} \right\}.$$

¹⁹ Note that although there is no RO in the ideal model of [2], whenever a miner would make such a query in the Bitcoin protocol, the corresponding dummy party sends a special MAINTAIN-LEDGER command to the Ledger functionality, making it possible for us to count the mining queries also in the ideal world.

²⁰ By definition, these two properties combined specify when the adversary should be considered the recipient of the reward.

At first glance, the choice of $2^{\text{polylog}(\kappa)}$ might seem somewhat arbitrary. However, it is there to guarantee that if the ledger state forks (recall that this reflects a violation of the common-prefix property) with noticeable probability, then the designer is punished with this super-polynomially high penalty to make his expected payoff negative as κ grows. On the other hand, if the probability of such a fork is sufficiently small (e.g. in the order of $2^{-\Omega(\kappa)}$), then the loss in utility is made negligible. This, combined with the fact that our stability notions will render negligible losses in the utility irrelevant, will allow the designer the freedom to provide slightly imperfect protocols, i.e., protocols where violations of the common-prefix property occur with sufficiently small probability.

We will denote by $\mathcal{M}^{\mathbb{B}}$ the Bitcoin attack model which has $\mathcal{G}_{\text{LEDGER}}^{\mathbb{B}}$ as the goal, $\langle \mathcal{G}_{\text{LEDGER}}^{\mathbb{B}} \rangle$ as the relaxed functionality, and scoring functions for the attacker and designer inducing utilities $u_{\mathbb{A}}^{\mathbb{B}}$ and $u_{\mathbb{D}}^{\mathbb{B}}$, respectively.

3.3 Attack-Payoff Security and Incentive Compatibility

The definition of the respective utilities for designer and attacker completes the specification of an attack game. Next, we define the appropriate notions of security and stability as they relate to Bitcoin and discuss their meaning.

We start with *attack-payoff security* [11], which, as already mentioned, captures that the adversary would have no incentive to make the protocol deviate from a protocol that implements the ideal specification (i.e., from a protocol that implements the ideal [non-relaxed] ledger functionality), and which is useful in arguing about the resistance of the protocol against incentive-driven attacks. However, in the context of Bitcoin analysis, one might be interested in achieving an even stronger notion of incentive-driven security, which instead of restricting the adversary to strategies that yield payoff as much as the ideal ledger $\mathcal{G}_{\text{LEDGER}}^{\mathbb{B}}$ from [2] would, restricts him to play in a coordinated fashion but *passively*, i.e., follow the mining procedure mandated by the Bitcoin protocol, including announcing each block as soon as it is found, but ensure that no two corrupt parties try to solve the same puzzle (i.e., use the same nonce).

One can think of the above strategy as corresponding to cooperating mining-pools which run the standard Bitcoin protocol. Nonetheless, as the adversary has control over message delays, he is able to make sure that whenever he finds a new block in the same round as some other party, his own block will be the one propagated first²¹, and therefore the one that will be added to the blockchain. Note that a similar guarantee is not there for honest miners as in the event of collisions—two miners solve a puzzle in the same round—the colliding miners have no guarantee about whose block will make it. We will refer to such an adversary that sticks to the Bitcoin mining procedure but makes sure his blocks are propagated first as *front running*.

Definition 2 (Front-running, passive mining adversary). *The front-running adversarial strategy \mathcal{A}_{fr} is specified as follows: Upon activation in round*

²¹ This can be thought of as a “rushing” strategy with respect to network delays.

$r > 0$, \mathcal{A}_{fr} activates in a round-robin fashion all its (passively) corrupted parties, say p_1, \dots, p_t . When party p_i generated some new message to be sent through the network, \mathcal{A}_{fr} immediately delivers m to all its recipients.²² In addition, upon any activation, any message submitted to the network $\mathcal{F}_{\text{N-MC}}$ by an honest party is maximally delayed.

Note that there might be several front-running, passive mining strategies, depending on which parties are corrupted and (in case of adaptive adversaries) when. We shall denote the class of all such adversary strategies by \mathbb{A}_{fr} . We are now ready to provide the definition of (strong) attack-payoff security for Bitcoin. The definition uses the standard notion of *negl-best-response* strategy from game theory: Consider a two-player game with utilities u_1 and u_2 , respectively. A strategy for m_1 of p_1 is *best response* to a strategy m_2 of p_2 if for all possible strategies m'_1 , $u_1(m'_1, m_2) \leq u_1(m_1, m_2) + \text{negl}(\kappa)$. For conciseness, in the sequel we will refer to *negl-best-response* simply as *best-response* strategies.

Definition 3. A protocol Π is strongly attack-payoff secure for attack model \mathcal{M}^{B} if for some $\mathcal{A} \in \mathbb{A}_{\text{fr}}$ the attacker playing \mathcal{A} is a (negl-)best-response to the designer playing Π .

Remark 3. It is instructive to see that for such a weak class of adversaries the usual blockchain properties hold with very nice parameters²³: first, the common-prefix property is satisfied except with negligible probability (as no intentional forks are provoked by anyone). Second, the fraction of honest blocks (in an interval of say k blocks) is roughly $\frac{\alpha}{\alpha+\beta} \stackrel{p \ll 1}{\approx} \frac{(1-\rho)np}{(1-\rho)np + \rho np} = (1-\rho)$ and thus, in expectation, the chain quality corresponds to the relative mining power of honest parties. Finally, since the adversary does contribute his mining power to the main chain, the number of rounds it takes for the chain to grow by k blocks is in expectation $\frac{k}{\alpha+\beta} \stackrel{p \ll 1}{\approx} \frac{k}{np}$.

Security thus means that if the honest parties stick to their protocol then the adversary has no incentive to deviate. However, unlike in [11], where the game is zero-sum, in a non-zero-sum setting it does not imply that the *designer* has an incentive to stick to the protocol. This means that the definition is useful to answer the question whether, assuming the network keeps mining, some of the miners have an incentive to deviate from the protocol, but it does not address the question of why the *honest miners* would keep mining. To address this question, we adopt the notion of *incentive compatibility* (IC).

Informally, a protocol being incentive-compatible means that both the attacker and the designer are willing to stick to it. In other words, it is strongly attack-payoff secure—i.e., the adversary will run it if the honest parties do—and

²² I.e., \mathcal{A}_{fr} sets the delay of the corresponding transmissions to 0.

²³ Recall the notation introduced in Sect. 2.2: n denotes the number of parties, ρ the fraction of corrupted parties, α and β denote honest and dishonest mining power, respectively, and p is the probability of a fresh RO-query to return a correct PoW solution.

if the adversary plays it passively (and front-running), then the honest miners will have an incentive to follow the protocol—i.e., the protocol is the designer’s best response to a passive front-running adversary. We note that requiring IC for Bitcoin for the class of all possible protocols would imply a proof that Bitcoin is not only a protocol that the miners wish to follow, but also that there is no other protocol that they would rather participate in instead. This is clearly too strong a requirement, even more so in the presence of results [13, 28] that argue that there are alternative “fairer” blockchain protocols which improve on the miners’ expected revenue. Thus, we can only hope to make such statements for a subclass of possible protocols, and therefore devise a version of IC which is parameterized by the set of all acceptable deviations (i.e., alternative protocols) \mathbb{D} . For full generality, we also parameterize it with respect to the class of acceptable adversaries \mathbb{A} , but stress that all statements in this work are for the class of all (PPT) adversaries.

Towards providing the formal definition of IC, we first give the straightforward restriction of equilibrium (in our case, subgame-perfect equilibrium) to a subset of strategies.

Definition 4. *Let \mathbb{D} and \mathbb{A} be sets of possible strategies for the designer and the attacker, respectively. We say that a pair $(\Pi, \mathcal{A}) \in (\mathbb{D}, \mathbb{A})$ is a (\mathbb{D}, \mathbb{A}) -subgame perfect equilibrium in the attack game defined by model \mathcal{M} , if it is a $(\text{negl}(\kappa))$ -subgame-perfect equilibrium on the restricted attack game where the set of all possible deviations of the designer (resp., the attacker) is \mathbb{D} (resp., \mathbb{A}).*

The formal definition of (parameterized) IC is then as follows:

Definition 5. *Let Π be a protocol and \mathbb{D} be a set of polynomial-time protocols that have access to the same hybrids as Π . We say that Π is \mathbb{D} -incentive compatible (\mathbb{D} -IC for short) in the attack model \mathcal{M} iff for some $\mathcal{A} \in \mathbb{A}_{fr}$, (Π, \mathcal{A}) is a (\mathbb{D}, ITM) -subgame-perfect equilibrium in the attack game defined by \mathcal{M} .*

4 Analysis of Bitcoin Without Transaction Fees

In this section, we present our RPD analysis of Bitcoin for the concrete incentive structure defined in the previous section. We note that this incentive structure does not, in particular, reflect rewards that stem from transaction fees and hence the reward per block is constant. First, in Sect. 4.1, we prove that Bitcoin is strongly attack-payoff secure—i.e., if the designer plays it, the attacker is better off sticking to it as well (but in a front-running fashion). The result is independent of the distribution of computing power to honest vs adversarial miners and independent of the conversion rate or the values of `breward` and `mcost`.

Subsequently, in Sect. 4.2, we investigate the role of mining costs vs conversion rate vs block rewards for the stability (i.e., IC) of Bitcoin in the presence of such incentive-driven coordinated coalitions (e.g., utility-maximizing mining pools.) We devise conditions on these values that either make the utility of honest parties negative—hence make playing the Bitcoin protocol a sub-optimal choice

of the protocol designer, or yield high enough utility for mining that makes Bitcoin *optimal* among all possible deviations from the standard protocol that are still compatible with the Bitcoin network (i.e., produce valid blockchains); combining this with the results from Sect. 4.1, we deduce that for this latter range of parameters Bitcoin is incentive-compatible.

4.1 Attack-Payoff Security of Bitcoin (Without Fees)

The attack-payoff security of Bitcoin without fees is stated in the following theorem.

Theorem 2. *The Bitcoin protocol is strongly attack-payoff secure in the attack model \mathcal{M}^B .*

Proof. The theorem follows as a direct corollary of the following general lemma.

Lemma 2. *Given any adversarial strategy, there is a front-running, semi-honest mining adversary \mathcal{A} that achieves better utility. In particular, the adversarial strategy \mathcal{A} makes as many RO-queries per round as allowed by the real-world restrictions, and one environment that maximizes its utility is the environment \mathcal{Z} that activates \mathcal{A} as the first ITM in every round until \mathcal{A} halts.*

Proof intuition. The proof of the lemma consists of three steps. First, we analyze Bitcoin in the real world. By invoking the subroutine-replacement theorem from [11, Theorem 6], we are able to work in a hybrid world where we can easily compute the relevant values, such as the number of blocks an adversary can mine in a given interval of rounds (the hybrid world is the so-called state-exchange hybrid world of [2]). Second, we show by a generic argument that this real-world analysis is sufficient to compute the payoffs for the attacker (which is defined on the transcript in the ideal world). Last but not least, we make a case distinction whether the adversary has expected utility smaller than zero (in which he does not corrupt any party and does not participate in the network), or whether mining Bitcoin is profitable for the attacker. In both cases, we prove that for any attacker \mathcal{A} , we can devise a front-running and semi-honest mining adversary which gets higher utility. The formal proof of the lemma is found in the full version [1]. \square

4.2 Incentive Compatibility of Bitcoin (Without Fees)

We proceed by investigating how the IC of Bitcoin depends on the relation between rewards and the conversion rate. Concretely, we describe a sufficient condition for IC (Theorem 4) and a condition that makes it non-IC (Theorem 3). We start with the negative result, which, informally, says that if the expected costs are too high with respect to the expected rewards, then Bitcoin is not IC (although it is strongly attack-payoff secure as proved above). As above, we denote by p the probability of solving a proof of work (and hence being a candidate to extend the ledger state) using one query to the random oracle (or equivalently, that a query to the state-exchange functionality successfully extends a state).

Theorem 3. For $n > 0$ and $\mathbf{breward} \cdot \mathbf{CR} < \frac{\mathbf{mcost}}{p}$ the Bitcoin protocol is not incentive compatible.

The proof is a straightforward calculation of the utility for the designer per round. Under the above condition, this expectation is less than 0, since they spend (on average) more on queries than what the reward compensates. Hence, the best response would be a protocol that does nothing.

While the above condition implies that the Bitcoin protocol is not a stable solution for all choices of the rewards, costs, and \mathbf{CR} , we next provide conditions under which the standard Bitcoin protocol is in fact a stable solution in the attack game. For this, we need to compare it to arbitrary alternative strategies that produce valid blocks for the Bitcoin network. Informally, our condition for IC requires that \mathbf{CR} and $\mathbf{breward}$ are sufficiently higher than the costs.

Theorem 4. Consider the real world consisting of the random oracle functionality \mathcal{F}_{RO} , the diffusion network $\mathcal{F}_{\text{N-MC}}$, and the clock $\mathcal{G}_{\text{CLOCK}}$, and let $\mathcal{W}_{\text{flat}}(\cdot)$ be the wrapper that formalizes the restrictions of the flat model.²⁴ Consider the class $\mathbb{I}_{\text{invalidchain}_{H,d}(\cdot)}$ of protocols Π that are defined for the $\mathcal{W}_{\text{flat}}(\mathcal{G}_{\text{CLOCK}}, \mathcal{F}_{\text{RO}}, \mathcal{F}_{\text{N-MC}})$ -hybrid world and which are compatible with the Bitcoin network, i.e., which obey the following two restrictions:

1. With probability 1, the real-world transcript (i.e., the real-world UC-execution of Π , any environment and adversary) does not contain a chain \mathcal{C} with $\text{invalidchain}_{H,d}(\mathcal{C}) = 0$ and this chain was an output to the network from an uncorrupted protocol instance running Π .
2. Upon input $(\text{READ}, \text{sid})$ to a protocol instance, the return value is $(\text{READ}, \text{sid}, \vec{\text{st}}^{\uparrow T})$ (for some integer T), where $\vec{\text{st}}^{\uparrow T}$ denotes the prefix of the state $\vec{\text{st}}$ encoded in the longest valid chain \mathcal{C} received by this protocol instance.

With respect to the class $\mathbb{I}_{\text{invalidchain}_{H,d}(\cdot)}$, the Bitcoin protocol is an incentive-compatible choice for the protocol designer if Bitcoin is profitable as in Lemma 3, i.e., if we are in the region $\mathbf{breward} \cdot \mathbf{CR} > \frac{n \cdot \mathbf{mcost}}{p}$, and if

$$\mathbf{breward} \cdot \mathbf{CR} > \frac{\mathbf{mcost}}{p \cdot (1-p)^{n-1}}. \quad (4)$$

Remark 4. Formula 4 constitutes a stronger requirement than the mere condition that mining should be profitable (which we treat separately in Lemma 3 for completeness). The theorem says that the probability that a fixed miner is uniquely successful stands in a reasonable relation to the mining cost and block rewards to achieve a stable solution. While Bitcoin would already yield positive utility to the protocol designer in the case of $\mathbf{breward} \cdot \mathbf{CR} > \frac{n \cdot \mathbf{mcost}}{p}$, we have for large n , $\frac{\mathbf{mcost}}{p} \cdot n \leq \frac{\mathbf{mcost}}{p} \cdot (\frac{1}{1-p})^{n-1}$ (for $p \in (0, 1)$).

²⁴ Recall from [2] that we model restrictions by using functionality wrappers. The above implemented restrictions correspond to the so-called flat model of Bitcoin, where each party gets one query to the random oracle per round and can send and receive one vector of messages in each round.

Proof intuition. The proof follows by demonstrating, in a sequence of claims, that the actual choices of the Bitcoin protocol (i.e., our abstraction of it) are optimal under the conditions of the theorem. This includes proving that the assumed resources cannot be employed in a way that would yield better payoff to the protocol designer. Intuitively, if the protocol has to be compatible with the Bitcoin network (i.e., it has to produce valid chains with probability 1), and invest its resources to achieve the optimum reward vs. query ratio in a setting where it knows it is running against front-running adversary running Bitcoin (such as mining pools). Optimality under the theorem’s condition follows by deducing a couple of useful properties from the fact that the protocol has to work potentially independently (per round) and by computing (and maximizing) the distribution of the possible query-vs.-reward ratios. The formal proof is found in the full version [1]. \square

We note that the above conditions are not necessarily tight. Thus one might wonder whether we can prove or disprove their tightness, and in the latter case investigate tight conditions for the statements to hold. We conclude this section with the following lemma which serves as first partial attempt to investigate this gap. The lemma implies that there might be potential to prove (partial) IC even for values of the parameters that fall in the gap between the above theorems. We leave the thorough investigation of this gap in terms of stability as a future research direction.

Lemma 3. *If $\text{breward} \cdot CR > \frac{n \cdot \text{mcost}}{p}$ then the Bitcoin protocol yields, with overwhelming probability, a positive utility for the protocol designer in the presence of front-running adversaries, i.e., the Bitcoin protocol is profitable in such a setting.*

5 Analysis of Bitcoin with Transaction Fees

Recall that in our formal treatment a chain \mathcal{C} encodes a ledger state $\vec{\text{st}}$. A ledger state is a sequence of individual state-blocks, i.e., $\vec{\text{st}} = \text{st}_1 || \dots || \text{st}_\ell$. In addition, each state-block $\text{st} \in \vec{\text{st}}$ (except the genesis state) of the state encoded in the blockchain has the form $\text{st} = \text{Blockify}(\vec{N})$ where \vec{N} is a vector of transactions, i.e., $\vec{N} = \text{tx}_1, \dots, \text{tx}_k$. A transaction tx_i can be seen as the abstract content of a block. Our above analysis assumes that the contents of the blocks do not affect the incentives of the attacker and the designer. In the real-world execution of Bitcoin, however, this is not the case as the contents of the blocks are money-like transactions and have transaction fees associated with them. We model these using positive-valued function $\text{tx} \mapsto f(\text{tx})$ mapping individual transactions to a positive real value that are integer multiples of 1 *Satoshi* (equals 10^{-8} Bitcoin).²⁵ For sake of brevity, we will also denote by $\hat{f}(\text{st}) := \sum_{\text{tx} \in \text{st}} f(\text{tx})$ the sum of all fees contained in the state block st . The fees have to be considered when defining the utilities in a rational analysis since they are added to the (flat) block reward and the total sum is given as a reward to the miner who inserts the block into

²⁵ Note that this modeling aspect is not sensitive to the basic unit of measurement.

the ledger state. Hence, this section treats the case where overall block rewards can be a dynamic quantity. In fact, the plan for Bitcoin is to eventually drop the block rewards at which point mining will be incentivized exclusively by the associated transaction fees. In this section we study the security and stability of the Bitcoin network incorporating also such fees.

5.1 Utility Functions with Fees

We first have to change the definition of the utility functions to incorporate that the attacker and the designer receive a different reward when inserting a block into the ledger state. The difference are the transactions fees. To this end, we first introduce a set $\mathcal{T}_{\mathcal{Z}}$ which contains all transactions that are submitted by the environment (and in particular not by the adversary), and then define the relevant events to capture fees in our model.²⁶

- In an execution, let $\mathcal{T}_{\mathcal{Z}}$ be the set of transactions such that $\mathbf{tx} \in \mathcal{T}_{\mathcal{Z}}$ if and only if \mathbf{tx} first appeared as an input from the environment (i.e., the first occurrence of \mathbf{tx} is in a command (SUBMIT, \mathbf{tx}) in this execution).
- For each $(\mu, r) \in \mathbb{N}^2$ the event $F_{r,\mu}^A$ is defined as follows: $F_{r,\mu}^A$ denotes the event that the total sum of the transaction fees $f(\mathbf{tx})$ of all $\mathbf{tx} \in \mathcal{T}_{\mathcal{Z}}$ contained in the blocks that the adversary adds to the state in round r is equal to $\mu \cdot 10^{-8} \cdot \text{CR}$ cost units.²⁷
- For each $(\mu, r) \in \mathbb{N}^2$ let the event $F_{r,\mu}^D$ be defined as follows: $F_{r,\mu}^D$ is the event that the total sum of the transaction fees $f(\mathbf{tx})$ of all $\mathbf{tx} \in \mathcal{T}_{\mathcal{Z}}$ contained in the blocks that the honest miners (jointly) add to the state in round r is equivalent to $\mu \cdot 10^{-8} \cdot \text{CR}$ cost units.

Since it is the environment that decides on the block-content, the sum of the fees in each block is effectively a random variable whose distribution is induced by the environment. The utilities of the attacker and designer that incorporate fees are defined as follows (we use $\hat{u}_A^{\mathbb{B}}$ and $\hat{u}_D^{\mathbb{B}}$ to denote the utilities when fees are added to the incentives):

$$\hat{u}_A^{\mathbb{B}}(\Pi, \mathcal{A}) = \sup_{\mathcal{Z} \in \text{ITM}} \left\{ \inf_{\mathcal{S}^A \in \mathcal{C}_A} \left\{ \text{breward} \cdot \text{CR} \cdot E(B^A) - q \cdot \text{mcost} \cdot E(Q^A) + \sum_{(\mu, r) \in \mathbb{N}^2} \mu \cdot 10^{-8} \cdot \text{CR} \cdot \Pr[F_{r,\mu}^A] \right\} \right\}$$

²⁶ Note that we assume that only transactions submitted by the environment can yield fees, since the environment models “the application layer”. In particular, if the adversary creates a transaction on his own and includes it in his next mined block, then this should not assign him any payoff.

²⁷ Recall that CR is the conversion of one cryptocurrency unit (e.g., one bitcoin) to one cost unit (e.g., one US dollar).

and

$$\hat{u}_D^{\mathbb{B}}(\Pi, \mathcal{A}) = \inf_{\mathcal{Z} \in \mathcal{Z}} \left\{ \sup_{\mathcal{S}^{\mathcal{A}} \in \mathcal{S}_{\mathcal{A}}} \left\{ \sum_{(b,r) \in \mathbb{N}^2} b \cdot \text{CR} \cdot (\text{breward} \cdot \Pr[F_{b,r}^{\mathbb{D}}] - 2^{\text{polylog}(\kappa)} \cdot \Pr[K_r]) \right. \right. \\ \left. \left. - \sum_{(q,r) \in \mathbb{N}^2} q \cdot \text{mcost} \cdot \Pr[W_{q,r}^{\mathbb{D}}] + \sum_{(\mu,r) \in \mathbb{N}^2} \mu \cdot 10^{-8} \cdot \text{CR} \cdot \Pr[F_{r,\mu}^{\mathbb{D}}] \right\} \right\}.$$

Note that the multiplicative factor 10^{-8} is there to allow us to set μ to the integer multiple of one Satoshi that the fee yields. We will denote by $\hat{\mathcal{M}}^{\mathbb{B}}$ the Bitcoin attack model which has $\mathcal{G}_{\text{LEDGER}}^{\mathbb{B}}$ as the goal, $\langle \mathcal{G}_{\text{LEDGER}}^{\mathbb{B}} \rangle$ as the relaxed functionality, and scoring functions for the attacker and designer inducing utilities $\hat{u}_A^{\mathbb{B}}$ and $\hat{u}_D^{\mathbb{B}}$.

Upper bounds on fees and total reward for blocks. In reality, transaction fees and the overall reward of a block are naturally bounded (either by size limits or by restricting the total value of the system).²⁸ In the following, we assume that for all \mathbf{tx} , $f(\mathbf{tx}) \leq \text{max}_{\text{fee}}$, and that the sum of fees per block is bounded, yielding an upper bound on the total profit per block: For all state blocks \mathbf{st} we require that $\text{breward} + \hat{f}(\mathbf{st}) \leq \text{max}_{\text{block}}$, where max_{fee} and $\text{max}_{\text{block}}$ are (strictly) positive multiples of one Satoshi.

Restrictions on the availability of transactions. So far in our treatment, the environment induces a distribution on the available transactions and is in principle unrestricted in doing so. For example, the set $\mathcal{T}_{\mathcal{Z}}$ is not bounded in size except by the running time of \mathcal{Z} . As will become apparent below in Theorem 5, putting no restrictions on the set $\mathcal{T}_{\mathcal{Z}}$ can still lead to meaningful statements that apply, for example, to applications that are believed to generate an (a priori) unbounded number of transactions. However, to model different kinds of scenarios that appear in the real world, we have to develop a language that allows us to speak about limited availability of transactions. To this end, we introduce parameterized environments $\mathcal{Z}^{\mathcal{D}}$. More precisely, let \mathcal{D} be an oracle which takes inputs $(\text{NEXTTXS}, r)$ and returns a vector $\vec{T}_r = (\mathbf{tx}_1, p_{i_1}), \dots, (\mathbf{tx}_k, p_{i_k})$. We say that an environment is \mathcal{D} -respecting, if, in every round r , the environment queries the oracle \mathcal{D} and only transactions $\mathbf{tx} \in \vec{T}_r$ are added to $\mathcal{T}_{\mathcal{Z}}$. We further require that \mathcal{Z} submits $(\text{SUBMIT}, \mathbf{tx}_i)$ to party p_k in round r if and only if $(\mathbf{tx}_i, p_k) \in \vec{T}_r$. For simplicity, we call \mathcal{D} simply a distribution. The utility for the attacker in such environments is taken to be the supremum as above, but only over all \mathcal{D} -respecting environments.

5.2 Analysis of Bitcoin (with Fees)

The following theorem says that if we look at unrestricted environments, then Bitcoin is still incentive compatible. This is a consequence of Theorems 2 and 4 and proven formally in the full version [1].

²⁸ For example, the number of total Bitcoins is limited and the block-size is bounded.

Theorem 5. *Consider arbitrary environments and let the sum of the transaction fees per block be bounded by $\max_{\text{block}} > 0$. Then the Bitcoin protocol is strongly attack-payoff secure in the attack model $\hat{\mathcal{M}}^{\mathcal{B}}$. It is further incentive-compatible with respect to the class of protocols that are compatible with the Bitcoin network under the same conditions as in Theorem 4), i.e., if*

$$\text{breward} \cdot \text{CR} > \frac{\text{mcost}}{p \cdot (1-p)^{n-1}}.$$

The previous statement is void in case the flat block reward is 0. However, for certain types of distributions \mathcal{D} , namely, the ones that provide sufficient high-fee transactions to the participants, it will remain in an equilibrium state. The statement is proven in the full version [1].

Theorem 6. *Consider distributions \mathcal{D} with the following property: In every round, \mathcal{D} outputs a vector of transactions such that any party gets as input a list of transactions to build a valid next state block \mathbf{st} to extend the longest chain and such that $\hat{f}(\mathbf{st}) = \max_{\text{block}}$ holds (where $\max_{\text{block}} > 0$). Then, with respect to \mathcal{D} -respecting environments, the Bitcoin protocol is strongly attack-payoff secure in the attack model $\hat{\mathcal{M}}^{\mathcal{B}}$. It is further incentive compatible with respect to the class of protocols that are compatible with the Bitcoin network (as defined in Theorem 4) if $\max_{\text{block}} \cdot \text{CR} > \frac{\text{mcost}}{p \cdot (1-p)^{n-1}}$.*

However, if an application cannot provide enough transactions, it becomes problematic, as the following counterexample shows.

Theorem 7. *There exist distributions \mathcal{D} such that the Bitcoin protocol is neither attack-payoff secure nor strongly attack-payoff secure with respect to \mathcal{D} -respecting environments.*

Proof. The proof is straightforward and follows from a general observation: assume there is just a single transaction in the network which has been received only by a corrupted party p_i . Then, the adversary does not publish this transaction to the network. If he does not, then he will be the one claiming the reward with probability one, which is his best choice. Hence, he does not follow the protocol (as the semi-honest front-running adversary would do) and hence it cannot be strongly attack-payoff secure.

Furthermore, the protocol is also not attack-payoff secure. If the honest-majority assumption does not hold, and thus an adversary can fork the ledger state, he would exercise his power to create a ledger state where it is a corrupted party who mines the block containing the only transaction in the system as this will yield better reward than simply mining on empty blocks. \square

Fallback security. Note that because cryptographic security trivially implies attack-payoff security for all possible environments and utilities, we can easily derive a fallback security notion: If the majority of miners mines honestly, then we get attack-payoff security; and even if this fails, we still get attack-payoff security under the assumption that the distribution of the fees and the relation between rewards vs costs vs conversion rate are as in Theorem 5 or 6.

References

1. Badertscher, C., Garay, J., Maurer, U., Tschudi, D., Zikas, V.: But why does it work? A rational protocol design treatment of bitcoin. Cryptology ePrint Archive, Report 2018/138 (2018). <https://eprint.iacr.org/2018/138>
2. Badertscher, C., Maurer, U., Tschudi, D., Zikas, V.: Bitcoin as a transaction ledger: a composable treatment. In: Katz, J., Shacham, H. (eds.) CRYPTO 2017. LNCS, vol. 10401, pp. 324–356. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-63688-7_11
3. Github: Bitcoin Core Version 0.12.0. Wallet: Transaction Fees. <https://github.com/bitcoin/bitcoin/blob/v0.12.0/doc/release-notes.md#wallet-transaction-fees>
4. Bonneau, J.: Why buy when you can rent? In: Clark, J., Meiklejohn, S., Ryan, P.Y.A., Wallach, D., Brenner, M., Rohloff, K. (eds.) FC 2016. LNCS, vol. 9604, pp. 19–26. Springer, Heidelberg (2016). https://doi.org/10.1007/978-3-662-53357-4_2
5. Canetti, R.: Security and composition of multiparty cryptographic protocols. J. Cryptol. **13**(1), 143–202 (2000)
6. Canetti, R.: Universally composable security: a new paradigm for cryptographic protocols. In: 42nd FOCS, pp. 136–145. IEEE Computer Society Press, October 2001
7. Carlsten, M., Kalodner, H.A., Weinberg, S.M., Narayanan, A.: On the instability of bitcoin without the block reward. In: Weippl, E.R., Katzenbeisser, S., Kruegel, C., Myers, A.C., Halevi, S. (eds.) ACM CCS 2016, pp. 154–167. ACM Press, October 2016
8. Eyal, I.: The miner’s dilemma. In: 2015 IEEE Symposium on Security and Privacy, pp. 89–103. IEEE Computer Society Press, May 2015
9. Eyal, I., Sirer, E.G.: Majority is not enough: bitcoin mining is vulnerable. In: Christin, N., Safavi-Naini, R. (eds.) FC 2014. LNCS, vol. 8437, pp. 436–454. Springer, Heidelberg (2014). https://doi.org/10.1007/978-3-662-45472-5_28
10. Fuchsbauer, G., Katz, J., Naccache, D.: Efficient rational secret sharing in standard communication networks. In: Micciancio, D. (ed.) TCC 2010. LNCS, vol. 5978, pp. 419–436. Springer, Heidelberg (2010). https://doi.org/10.1007/978-3-642-11799-2_25
11. Garay, J.A., Katz, J., Maurer, U., Tackmann, B., Zikas, V.: Rational protocol design: cryptography against incentive-driven adversaries. In: 54th FOCS, pp. 648–657. IEEE Computer Society Press, October 2013
12. Garay, J.A., Katz, J., Tackmann, B., Zikas, V.: How fair is your protocol? A utility-based approach to protocol optimality. In: Georgiou, C., Spirakis, P.G. (eds.) 34th ACM PODC, pp. 281–290. ACM, July 2015
13. Garay, J., Kiayias, A., Leonardos, N.: The bitcoin backbone protocol: analysis and applications. In: Oswald, E., Fischlin, M. (eds.) EUROCRYPT 2015. LNCS, vol. 9057, pp. 281–310. Springer, Heidelberg (2015). https://doi.org/10.1007/978-3-662-46803-6_10
14. Garay, J., Kiayias, A., Leonardos, N.: The bitcoin backbone protocol with chains of variable difficulty. In: Katz, J., Shacham, H. (eds.) CRYPTO 2017. LNCS, vol. 10401, pp. 291–323. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-63688-7_10
15. Gervais, A., Karame, G.O., Wüst, K., Glykantzis, V., Ritzdorf, H., Capkun, S.: On the security and performance of proof of work blockchains. In: Weippl, E.R., Katzenbeisser, S., Kruegel, C., Myers, A.C., Halevi, S. (eds.) ACM CCS 2016, pp. 3–16. ACM Press, October 2016

16. Goldreich, O.: Foundations of Cryptography: Volume 1, Basic Tools. Cambridge University Press, Cambridge (2003)
17. Gradwohl, R., Livne, N., Rosen, A.: Sequential rationality in cryptographic protocols. In: 51st FOCS, pp. 623–632. IEEE Computer Society Press, October 2010
18. Halpern, J.Y., Pass, R., Seeman, L.: Computational extensive-form games. In: EC (2016)
19. Katz, J.: Bridging game theory and cryptography: recent results and future directions. In: Canetti, R. (ed.) TCC 2008. LNCS, vol. 4948, pp. 251–272. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-78524-8_15
20. Katz, J., Maurer, U., Tackmann, B., Zikas, V.: Universally composable synchronous computation. In: Sahai, A. (ed.) TCC 2013. LNCS, vol. 7785, pp. 477–498. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-36594-2_27
21. Kol, G., Naor, M.: Games for exchanging information. In: Ladner, R.E., Dwork, C. (eds.) 40th ACM STOC, pp. 423–432. ACM Press, May 2008
22. Luu, L., Teutsch, J., Kulkarni, R., Saxena, P.: Demystifying incentives in the consensus computer. In: Ray, I., Li, N., Kruegel, C. (eds.) ACM CCS 2015, pp. 706–719. ACM Press, October 2015
23. Nakamoto, S.: Bitcoin: A Peer-to-Peer Electronic Cash System (2008). <http://bitcoin.org/bitcoin.pdf>
24. Nayak, K., Kumar, S., Miller, A., Shi, E.: Stubborn mining: generalizing selfish mining and combining with an eclipse attack. In: S&P (2016)
25. Ong, S.J., Parkes, D.C., Rosen, A., Vadhan, S.: Fairness with an honest minority and a rational majority. In: Reingold, O. (ed.) TCC 2009. LNCS, vol. 5444, pp. 36–53. Springer, Heidelberg (2009). https://doi.org/10.1007/978-3-642-00457-5_3
26. Osborne, M.J., Rubinstein, A.: A Course in Game Theory. MIT Press, Cambridge (1994)
27. Pass, R., Seeman, L., Shelat, A.: Analysis of the blockchain protocol in asynchronous networks. In: Coron, J.-S., Nielsen, J.B. (eds.) EUROCRYPT 2017. LNCS, vol. 10211, pp. 643–673. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-56614-6_22
28. Pass, R., Shi, E.: FruitChains: a fair blockchain. In: Schiller, E.M., Schwarzmann, A.A. (eds.) 36th ACM PODC, pp. 315–324. ACM, July 2017
29. Rosenfeld, M.: Analysis of bitcoin pooled mining reward systems. CoRR (2011)
30. Sapirshtein, A., Sompolinsky, Y., Zohar, A.: Optimal selfish mining strategies in bitcoin. In: Grossklags, J., Preneel, B. (eds.) FC 2016. LNCS, vol. 9603, pp. 515–532. Springer, Heidelberg (2017). https://doi.org/10.1007/978-3-662-54970-4_30
31. Schrijvers, O., Bonneau, J., Boneh, D., Roughgarden, T.: Incentive compatibility of bitcoin mining pool reward functions. In: Grossklags, J., Preneel, B. (eds.) FC 2016. LNCS, vol. 9603, pp. 477–498. Springer, Heidelberg (2017). https://doi.org/10.1007/978-3-662-54970-4_28
32. Teutsch, J., Jain, S., Saxena, P.: When cryptocurrencies mine their own business. In: Grossklags, J., Preneel, B. (eds.) FC 2016. LNCS, vol. 9603, pp. 499–514. Springer, Heidelberg (2017). https://doi.org/10.1007/978-3-662-54970-4_29