# Sending Spies as Insurance Against Bitcoin Pool Mining Block Withholding Attacks

Isamu Okada[1,2], Hannelore De Silva[1], and Krzysztof Paruch[1(✉)]

[1] Research Institute for Cryptoeconomics, Vienna University of Economics,
Welthandelsplatz 1, 1020 Vienna, Austria
`krzysztof.paruch@wu.ac.at`
[2] Faculty of Business Administration, Soka University,
1-236 Tangi, Hachioji, Tokyo 192-8577, Japan
`https://www.wu.ac.at/en/cryptoeconomics`, `https://www.soka.ac.jp/`

**Abstract.** Theoretical studies show that a block withholding attack is a considerable weakness of pool mining in Proof-of-Work consensus networks. Several defense mechanisms against the attack have been proposed in the past with a novel approach of sending sensors suggested by Lee and Kim in 2019. In this work we extend their approach by including mutual attacks of multiple pools as well as a deposit system for miners forming a pool. In our analysis we show that block withholding attacks can be made economically irrational when miners joining a pool are required to provide deposits to participate which can be confiscated in case of malicious behavior. We investigate minimal thresholds and optimal deposit requirements for various scenarios and conclude that this defense mechanism is only successful, when collected deposits are not redistributed to the miners.

**Keywords:** Blockchain · Bitcoin · Proof-of-Work · Pool mining · Block withholding attack · Game theory · Agent-based simulation

## 1 Introduction

### 1.1 Bitcoin, Pool Mining and Block Withholding Attacks

The Bitcoin network relies on the Proof-of-Work (PoW) consensus mechanism for its security where computational power provided by miners which are economically incentivized to participate to achieve a collectively agreed state of the system is utilized. Since reward mechanisms are stochastic agents aim to reduce their payout volatility by forming mining pools orchestrated by pool managers. All pool earnings are equally distributed among miners depending on their contributions allowing them to create a stable source of income with identical expected payoff but much lower variance compared to the solo mining case.

In a pool mining system multiple miners collaboratively find a *Nonce* number satisfying the following condition: $h(Merkle+h(PreBlock)+Nonce) < D$ where $h(\cdot)$ is a hash function, *Merkle* refers to the merkle root, *PreBlock* indicates the Block ID, and $D$ is a threshold which indicates the computation difficulty. An eligible *Nonce* number is a full solution of Proof-of-Work (fPoW). Miners in a pool agree to share their fPoW reward with all other participating miners in the pool as the merkle root includes a coinbase transaction which funnels rewards to the pool manager who distributes them among all miners after taking a management fee.

An attack vector in form of the block withholding attack (BWA) (Rosenfeld (2011)), (Zhu et al. (2018)), (Eyal and Sirer (2018)) is a serious threat to the pool mining approach and with this the whole legitimacy of the PoW consensus. In this attack the miner who finds a fPoW does not report it to the pool manager - either wanting to monopolize the reward for themselves or operates within another pool and aims to harm their competitors. For provability of provided efforts a pool mining system adopts an indirect index to ensure miners do not withhold correct solutions. This is accomplished via reporting of partial solutions of Proof-of-Work (pPoW) which solve the hash equation above for a $D'$ with $D' >> D$. Since pPoW are much easier to find they provide a statistically corresponding signature for the contributed efforts over time. A pool manager can reward his miners in relation to their pPoW solutions and keep the fPoW rewards for himself.

## 1.2    Literature Review and Previous Approaches

(Rosenfeld (2011)) is one of the first to analyse bitcoin pooled mining reward systems. He introduces terminology in the comparison of solo mining versus pool mining techniques and describes the main benefit of pool mining in reducing earnings variance. He treats different reward systems and associated score-based methods for reward distributions and outlines potential attack vectors *pool hopping* and *block withholding* for one malicious actor. Rosenfeld proposes two solutions which of one is a *pop quiz* that identifies not truthfully reporting agents and the other one is *oblivious shares* which necessitates a change in the protocol logic.

(Eyal (2013)) formally shows that the Bitcoin mining protocol is not incentive-compatible as it allows for colluding miners to obtain a bigger reward than their share of contribution to the network. He captures mathematically the logic behind the BWA and further shows this design will lead to the formation of centralized colluding entities: the selfish mining pools. Eyal suggested solution requires an alternation of the protocol to resolve these threats.

(Luu et al. (2015)) use game-theoretic methodology to assess BWA in the context of Bitcoin mining pools. Their model allows to treat mining decisions as computational power splitting games where agents distribute their resources among competing pools conditioned upon the rewards they offer. The authors are able to show that selfish mining is always favorable in the long run, that optimal behavior is a stochastic mixed strategy and thus the network in total

always wastes resources for competition. In addition the authors introduce seven desired properties of countermeasures and evaluate a change in reward payoff under these considerations.

(Eyal (2015)) extends the previous research from one-miner attacks to collaborative attacks performed by malicious pool managers grouping their participants to infiltrate other pools. Eyal introduces a model of the pool game where a number of pools can attack each other for varying pool sizes and attack strategies. This results in the Miner's Dilemma where similar to the iterative prisoner's dilemma pool managers decide continuously whether to attack the competitor. The strategy not to attack is dominated but leads to long-term highest gains if cooperation is established. Eyal suggest two countermeasures which of one goes in line with Rosenfeld's proposal to change the protocol and establish a honeypot technique to expose attackers and the other one claims that in closed pools no block withholding is possible.

(Lee and Kim (2019)) discuss available countermeasures to BWA and categorize them according to properties introduced by (Luu et al. (2015)). Moreover they propose an additional countermeasure that does not require a change in protocol and performs better with respect to the other characteristics. Their suggested method includes sending sensors to competing pools to be able to detect infiltrations and punish attackers. The authors consider a model with one attacking pool and one defending pool and asymmetric behavior.

### 1.3   Contribution of This Paper

We theoretically explore effective protocol defensive methods against BWA when assuming rational pool managers by developing an agent based framework. We build our analysis on top of the model introduced by (Eyal (2015)) and extend the defending method suggested by (Lee and Kim (2019)) by incorporating both pools capabilities to attack and send sensors to their competitors. After our model confirms existing results additional improvements are suggested: the effects of punishment and deposit on the agent behaviors.

The model is developed in four versions. The first includes two mutual attackers (Model 1) and is able to replicate Eyal's miner's dilemma situation. This is extended to a model applying mutual punishment (Model 2) using Lee and Kim's (2019) idea. The next model introduces a deposit system for participation in a pool (Model 3). The final model shows the effects of collected deposits redistribution (Model 4). All models are supported with numerical experiments.

In all models two pool managers adopt the Pay Per Latest N Share method and maximize their expected payoffs by adjusting their strategies on how to infiltrate the other pool with spies. The analysis of more complex management compensation schemes which could also impact agent's decisions is set aside for future research. With this paper we contribute to the existing discussion on the topic of pool mining and BWA as outlined above.

At first we solve each model version analytically for Nash equilibria of the two managers by calculating their best responses to finally derive best response strategies from numerical simulations.

## 2   Results

### 2.1   Model Definition

We consider a model of two pools $A$ and $B$ and assume respective relative mining powers within the whole network $\alpha$ and $\beta$, where $0 < \alpha + \beta \leq 1$. With $m$ the total number of mining tasks $m\alpha$ tasks are calculated in Pool $A$ while $m\beta$ tasks in $B$. Each miner joining in either Pool $A$ or $B$ is assigned one task of same size resulting in *one task per miner*. We conceptualize the possibility of more powerful miners which occur in reality as if virtually dividing them into an accumulation of miners with the capability of calculating one task only. The consequences of this simplification is left for further research.

We further allow both pools to send spies to each other. We define $x$ and $y$ as the fraction of miners which are spies meaning that $A$ sends $s_A \equiv m\alpha x$ spies to $B$ while $B$ sends $s_B \equiv m\beta y$ spies to $A$. The process of spying is as follows. Manager $A$ generates $s_A$ miners and sends them to $B$ where they are being assigned $s_A$ tasks by manager of $B$. Then the manager of $A$ reassigns these $s_A$ tasks to *non-spy* $s_A$ miners in his pool. Note that the manager never calculates tasks directly as this is delegated to *non-spy* miners who report to him. Finally manager $A$ reports selected results to manager $B$ to receive mining rewards. Manager $B$ mirrors the behavior of manager $A$.

It is each pool manager's choice which results to forward to the other pool, therefore with malicious intent no fPoW and all pPow solutions are submitted. This allows to receive mining rewards corresponding to pPoW solutions. The manager of $B$ cannot receive any fPoW solution of $s_A$ tasks which reduces his the probability to win. Such behavior is defined as a BWA by $A$ to $B$.

Eyal (2015) showed that BWA bears the miner's dilemma. Individually the attack increases the relative winning probability by decreasing the victim pool's chances. However, if both managers attack each other both winning probabilities decrease which corresponds to the prisoners' dilemma.

In our approach several versions of a model and associated metrics were developed to measure the pool's performance and to quantify the manager's and miner's incentives.

### 2.2   Model 1: Confirming the Miner's Dilemma

Following Eyal (2015) we use the three metrics direct revenue, revenue density and efficiency. The direct revenue indicates the winning probability of a mining race in the whole network and the revenue density indicates the expected reward of each miner. With $m$ as total network task number the revenue density is $1/m$ if there are no attacks. The efficiency defining the revenue density multiplied by $m$ reflects the attractiveness of a pool as miners have stronger incentives to join a pool when the efficiency is high. An efficiency of one is considered to be normal attractive for miners to join. Therefore a pool manager aims to maximize its efficiency.

**Definition 1.** *The **direct revenues** of pool A and B denoted by $D_A$ and $D_B$ respectively are given as $D_A = \frac{m\alpha - m\alpha x}{m - m\alpha x - m\beta y}$ and $D_B = \frac{m\beta - m\beta y}{m - m\alpha x - m\beta y}$.*

**Definition 2.** *The **revenue density** of Pool A and B denoted by $R_A$ and $R_B$ respectively are given as $R_A = \frac{D_A + m\alpha x R_B}{m\alpha + m\beta y}$ and $R_B = \frac{D_B + m\beta y R_A}{m\beta + m\alpha x}$.*

**Definition 3.** *The **efficiency** of Pool A and B denoted by $E_A$ and $E_B$ respectively are given as $E_A = mR_A$ and $E_B = mR_B$.*

Using these metrics we prove two theorems that allow us to conclude that there are strong incentives to attack competing pools as a pool manager. This holds in both cases when the counter party attacks or not as an attack always increases the efficiency regardless of the other managers strategy. We show that the efficiencies of both pools are equal to one if there are no attacks. In case a pool attacks the other pool has the incentive to retaliate to increase its relative efficiency. Proofs of Theorems are provided in the Appendix.
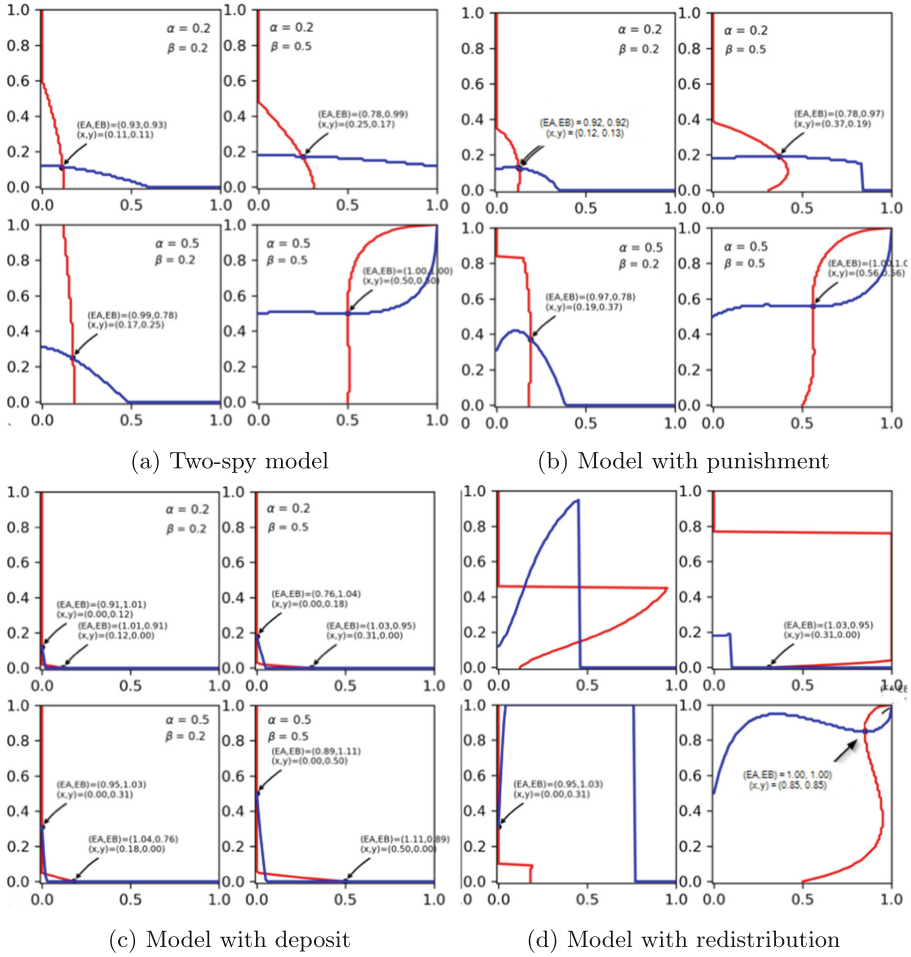
**Theorem 1.** *In the case that there are no attacks in the network it follows that the efficiency of both pools equals to one: $x = y = 0 \Rightarrow E_A = E_B = 1$*

**Theorem 2.** *In the case that one pool attacks the other it follows that the pool being attacked has an incentive to attack back to prevent having a lower efficiency: $x > 0 \land y = 0$ implies $E_A > E_B \land E_B < 1$. The efficiency of pool A is greater than one if the fraction $x$ of spies is below a threshold: $E_A > 1 \iff x < \frac{\beta}{1-\alpha}$*

This theorem implies pools efficiency suffers from attacks causing it to lose miners and an attacker can increase their attractiveness above 1. Since these implications influence each individual pool managers in aggregate the miner's dilemma will follow as proposed by Eyal (2015).

In a situation where both pool managers send spies to each other $x > 0$ and $y > 0$. To be more attractive for miners each pool manager wants to maximize their efficiency. Therefore w.l.o.g. manager A is looking for the optimal $x$ which maximizes $E_A$ for given $\alpha$, $\beta$, $y$ and $m$ and manager B makes analogous decisions. This idea is based on the best response analysis where $x^{BR}(y) = x_{E_A(y|\alpha,\beta,m)}$ is the best response for manager A while $y^{BR}(x) = y_{E_B(x|\alpha,\beta,m)}$ is the best response for manager B. If $(x,y) = (x^*,y^*)$ satisfies $(x^*,y^*) = (x^{BR}(y^*),y^{BR}(x^*))$ then $(x^*,y^*)$ is a Nash equilibrium. We defer from analytical treatment in this case and show the numerical results of best responses in Fig. 1 for all four versions of the model: standard two spy model, model with punishment, model with deposit and model with deposit redistribution.

Figure 1a shows that due to the miner's dilemma the efficiency of a pool manager is always below 1 even if best responses are selected. Furthermore the figure shows that there is an incentive to conduct a BWA, especially as a small pool - therefore managers of big pools have to consider suitable defences against BWA. If there is a pool with a notable share of mining power in the network the losses incurred by all pools due to the miner's dilemma are significant - and this effect is stronger the bigger the largest pool is. Consequently all managers observe the biggest pools of the network and aim to decrease their power.

(a) Two-spy model    (b) Model with punishment

(c) Model with deposit    (d) Model with redistribution

**Fig. 1.** Results of numerical analysis for best responses $x$ of pool manager A and $y$ for pool manager B are shown in red and blue respectively where $x$ is depicted on the horizontal axis and $y$ on the vertical axis in each panel. The Nash equilibrium is the intersection of best response lines where both efficiency values and both efficiency values $E_A$ and $E_B$ are shown. We set $m = 10,000,000$ and $f = 10$. (Color figure online)

## 2.3    Model 2: Punishment System

To prevent the BWA, Lee and Kim (2019) proposed sending spies. In this strategy the manager of $A$ send $s_A$ miners as spies to $B$, collects some open tasks and assigns those $s_A$ tasks to miners joining $A$. However these tasks are branded with information of the Coinbase Transaction which exposes their real affiliation to the respective pool and they can be therefore identified and correctly categorized by any miner. If spies are sent to $A$ by manager $B$ they would be able to

detect the attack - if these individuals are assigned the tasks with an inappropriate Coinbase transaction. As a consequence manager $B$ can sanction pool $A$ by not awarding any rewards for pPoW solutions if he discovers to be attacked. The probability to discover the attack can be calculated by an application of a combination problem: For the expected value of $d_{AB}$ of $A$'s spies detected by the manager of pool $B$ $d_{AB}$ it holds $d_{AB} = \sum_{k=0}^{m\beta y} k \frac{m\alpha x C_k \times m\beta C_{m\beta y-k}}{m\beta + m\alpha x C_{m\beta y}}$. The solution to this equation can be approximated by $d_{AB} = \frac{m\alpha\beta xy}{\alpha x + \beta}$ and in similar fashion $d_{BA} = \frac{m\alpha\beta xy}{\alpha + \beta y}$ to meet the requirements for this paper.

(Lee and Kim (2019)) proposed an approach which does provide improvements to previous models as it describes the one-attacker case but it does not cover the situation where both pools send spies to each other. It might very well be the case that pool manager $A$ himself considers being attacked if his $s_A$ miners find a task consisting of information on Coinbase Transactions including Pool $B$. The contribution of this paper is to extend Lee and Kim (2019)'s model to contain the two-attacker case and further to allow to assess both manager's strategies.

In comparison to the first model presented in the previous section the extended model of this section only differs in the definitions of equations for the revenue densities $R_A$ and $R_B$. Their numerators and denominators now include the terms for the expected detection rates $d_{AB}$ and $d_{BA}$ respectively which are subtracted from the mining tasks attributed to spying $m\alpha x$ and $m\beta y$ accounting for the fact that the revenue density of one pool is reduced by these terms if tasks are assigned to spies: $R_A = \frac{D_A + (m\alpha x - d_{AB})R_B}{m\alpha + m\beta y - d_{BA}}$ and $R_B = \frac{D_B + (m\beta y - d_{BA})R_A}{m\beta + m\alpha x - d_{AB}}$

Figure 1b shows the results of numerical analysis of best responses and Nash equilibria replicating the experiments presented in Fig. 1a. A comparison of both figures allows to conclude that a punishment scheme as defined by the second model performs worse when confronted with a BWA. Contrasting the numerical results for e.g. parameter values of $(\alpha, \beta) = (20\%, 20\%)$ reveals that in the former model pool efficiencies $E_A$ and $E_B$ are reduced to 93% when a fraction of 11% spies are sent while in the latter model both managers send 13% spies reducing the efficiency to 92%. The effects of the miner's dilemma are magnified in the model with punishment even more when considering pools of different sizes so that for e.g. $(\alpha, \beta) = (20\%, 50\%)$ the fraction of spies increases from 25% to 37% when comparing both models. Consequently the punishment system of mutually sending spies does not provide a satisfactory countermeasure to prevent BWH attacks.

## 2.4 Model 3: Deposit System

To overcome the shortcomings of the model with punishment we propose a deposit system. Each new miner who wants to join a pool is required by the pool manager to pledge a certain amount that lapses in the case of attack detection. Ceased amounts are attributed to the manager's income and are not redistributed to other pool members. If no misbehavior is ever discovered the deposit is fully returned to the initial miner when leaving the pool. The
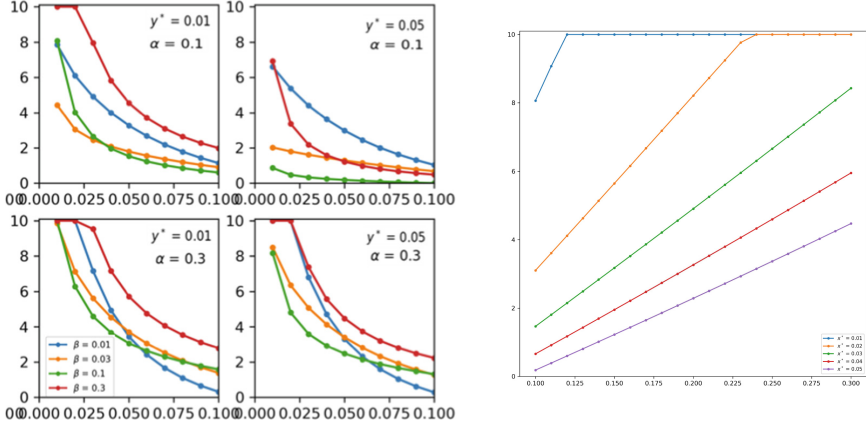
model of this system includes an additional parameter $f$ representing a deposit unit which is included in the definitions of revenue densities. This value satisfies that (deposit fee) $= f \times$ (expected reward unit) $= f/m$ and is equivalent to the expected reward of a task for $f = 1$ in the case of no attack: $R_A = \frac{D_A+(m\alpha x-d_{AB})R_B-d_{AB}f/m}{m\alpha+m\beta y-d_{BA}}$ and $R_B = \frac{D_B+(m\beta y-d_{BA})R_A-d_{BA}f/m}{m\beta+m\alpha x-d_{AB}}$

Numerical results of the model with deposit system are shown in Fig. 1c where the miner's dilemma appears to be resolved. A comparison of the panels for $(\alpha, \beta) = (50\%, 50\%)$ in Fig. 1a and Fig. 1c yields the conclusion that for the deposit system there no incentives to attack if the ratio of the other manager's spies exceeds a certain threshold. The best responses of both managers spy ratio indicated by the blue and red lines vanish almost everywhere in Fig. 1c in contrast to the model without punishment in Fig. 1a. Below a certain threshold of foreign spies however the incentive to attack becomes very large. Smaller pools are even less attractive to attack as the critical attacking threshold gets lower see e.g. $(\alpha, \beta) = (50\%, 20\%)$. This fact suggests that sending a small number of spies can serve as a form of insurance against potential attacks by the competing pool manager. In an environment with a deposit system mutual determent against BWA can be realized when all pools send a minimal amount of spies to each other. This result would not hold true without the deposit as this would even decrease the Pareto efficiency.

To find the optimal size of the required deposit for the insurance mechanism to work we refer to Fig. 2a where thresholds of the deposit unit are shown. A comparison of the various panels indicates that the threshold does not depend on the fraction of spies sent by the competing manager. This is consistent with the previous discussion because if a sufficient deposit unit is required, a pool manager who intends to do the BWA loses one's economic incentive to do so. When the calculating power of a pool is relatively large (the case of $\alpha = 30\%$), if the deposit unit sets $f = 6$, about $x = 5\%$ is needed for a pool with the same power. If the power of an opponent gets small, the threshold of $f$ drastically decreases.

In reality a pool manager never knows where a joining miner comes from and thus cannot change the deposit unit depending on a miner, rather it must be set depending on the most powerful pool. Although the relative mining power shares fluctuate on a daily basis they never change drastically. Therefore the optimal deposit requirement should depend on the biggest pool share. With no deposit, Fig. 1b shows that two big pools have strong economic incentives to do the BWH attack with many spies. Figure 2b shows the suitable deposit for this case, e.g. when the biggest pool in the network has a relative share of 25% and the insurance spies are $x^* = 3\%$, then the suitable deposit requirement is about 4.

(a) $x$ represented on the horizontal axis while threshold of $f$ shown on the vertical axis. For given $(\alpha, \beta)$ and $x$ the best response $y$ is regarded as a function of $f$: $y(f)$. The threshold is the minimum value of $f$ that satisfies $y(f) < y^*$.

(b) Suitable deposit requirement: The horizontal axis represents a relative mining power of the biggest pool in the whole network while the vertical axis represents the threshold of $f$ where $x^*$ is the fraction of spies sent to each other pool.

**Fig. 2.** Analysis of deposits: thresholds shown in Fig. 2a while optimal deposits shown in Fig. 2b. Thresholds above 10 are set to 10 in both Subfigures. We set $m = 10,000,000$.

Although the deposit system resolves the second issue of the miner's dilemma, the first issue is not resolved. How should a pool manager owning a small pool do against the BWA by a big pool? If the deposit unit is set to an extreme large value (for example, $f = 100$), this issue can possibly be resolved. This approach is feasible in practice as the deposit system bears no costs for honest miners besides opportunity costs of their deposits.

## 2.5   Model 4: Distributing the Lost Deposit

In this version of the model we consider the case of redistribution of claimed deposits. In Model 3 they are attributed to the pool manager, and thus are never redistributed to the miners in the pool. Here, we consider two arguments when dealing with the deposit.

In the first, a BWA is done by a pool manager and not a miner and the losses are incurred by the manager. Therefore the manager should receive the income from deposits to compensate his losses. However a pool manager is also very interested in maximizing its efficiency to attract more miners. Thus the income from collected deposits from attacks should be redistributed to the joining miners to raise the pools efficiency. Model 4 considers this argument.

For a manager of Pool $A$ resp. $B$ the collected deposit is $d_{BA}f/m$ resp. $d_{AB}f/m$ and thus the equations of $R_A$ and $R_B$ are revised to

$R_A = \frac{D_A + (m\alpha x - d_{AB})R_B - d_{AB}f/m + d_{BA}f/m}{m\alpha + m\beta y - d_{BA}}$ and
$R_B = \frac{D_B + (m\beta y - d_{BA})R_A - d_{BA}f/m + d_{AB}f/m}{m\beta + m\alpha x - d_{AB}}$

The result is shown in Fig. 1d. We see that the incentives to send spies and attack the opponent reappear compared to Fig. 1c.

## 3   Conclusion

In our analysis we have shown that results for the Miner's dilemma introduced by (Eyal (2015)) can be replicated in a model extending the approach presented in (Lee and Kim (2019)) where both pool managers send spies to the other system. Our numerical results indicate that an effective countermeasure to the BWH attack can be posed by introducing a deposit system in mining pools. In this case even a small number of spies sent to competing pools can act as an insurance against being attacked. This analysis is supported with the derivation of optimal deposit units and thresholds for various scenarios. Contrary to previously proposed countermeasures, this method seems to be effective even without the requirement to change protocol code and can be therefore implemented directly into already operating systems. Our analysis further shows that slashed deposits from malicious miners cannot be redistributed to the mining pool to increase it's efficiency since this would again reintroduce the problem of the Miner's dilemma and render the deposit insurance ineffective.

## A     Proof of Theorems

### A.1     Proof of Theorem 1

*Proof (Theorem 1).* If no attacks occur in the network both $x = y = 0$ by definition. Plugging both in into Definition 1 we get

$$D_A = \frac{m\alpha - m\alpha x}{m - m\alpha x - m\beta y} = \frac{m\alpha}{m} = \alpha$$

and

$$D_B = \frac{m\beta - m\beta y}{m - m\alpha x - m\beta y} = \frac{m\beta}{m} = \beta$$

which respecting Definition 2 yields

$$R_A = \frac{D_A + m\alpha x R_B}{m\alpha + m\beta y} = \frac{\alpha}{m\alpha} = \frac{1}{m}$$

and

$$R_B = \frac{D_B + m\beta y R_A}{m\beta + m\alpha x} = \frac{\beta}{m\beta} = \frac{1}{m}.$$

Due to Definition 3 we get $E_A = E_B = 1$.

## A.2    Proof of Theorem 2

*Proof (Theorem 2).* For $x > 0$ and $y = 0$ plugging in into 1 yields

$$D_A = \frac{m\alpha - m\alpha x}{m - m\alpha x - m\beta y} = \frac{m\alpha - m\alpha x}{m - m\alpha x} = \frac{m(\alpha - \alpha x)}{m(1 - \alpha x)} = \frac{\alpha(1 - x)}{1 - \alpha x}$$

and

$$D_B = \frac{m\beta - m\beta y}{m - m\alpha x - m\beta y} = \frac{m\beta}{m - m\alpha x} = \frac{m\beta}{m(1 - \alpha x)} = \frac{\beta}{1 - \alpha x}$$

which again can be inserted into 2 and thus since $y = 0$ then $R_B$ reduces to

$$R_B = \frac{D_B + m\beta y R_A}{m\beta + m\alpha x} = \frac{\frac{\beta}{1 - \alpha x}}{m(\beta + \alpha x)} = \frac{\beta}{1 - \alpha x} \cdot \frac{1}{m(\beta + \alpha x)} = \frac{\beta}{m\beta + m\alpha x - m\beta\alpha x - m\alpha^2 x^2}$$

which can be plugged into $R_A$

$$
\begin{aligned}
R_A &= \frac{D_A + m\alpha x R_B}{m\alpha + m\beta y} = \frac{\frac{\alpha(1-x)}{1-\alpha x} + m\alpha x \frac{\beta}{1-\alpha x} \cdot \frac{1}{m(\beta+\alpha x)}}{m\alpha} \\
&= \left( \frac{\alpha(1-x)}{1-\alpha x} + m\alpha x \frac{\beta}{1-\alpha x} \cdot \frac{1}{m(\beta+\alpha x)} \right) \cdot \frac{1}{m\alpha} \\
&= \left( \frac{\alpha(1-x)}{1-\alpha x} \cdot \frac{m(\beta+\alpha x)}{m(\beta+\alpha x)} + \frac{m\alpha x\beta}{1-\alpha x} \cdot \frac{1}{m(\beta+\alpha x)} \right) \cdot \frac{1}{m\alpha} \\
&= \left( \frac{\alpha(1-x) \cdot m(\beta+\alpha x) + m\alpha x\beta}{(1-\alpha x) \cdot m(\beta+\alpha x)} \right) \cdot \frac{1}{m\alpha} \\
&= \left( \frac{(\alpha - \alpha x) \cdot (m\beta + m\alpha x) + m\alpha x\beta}{(1-\alpha x) \cdot m(\beta+\alpha x)} \right) \cdot \frac{1}{m\alpha} \\
&= \left( \frac{m\beta\alpha - m\beta\alpha x - m\alpha^2 x^2 + m\alpha^2 x + m\alpha x\beta}{(1-\alpha x) \cdot m(\beta+\alpha x)} \right) \cdot \frac{1}{m\alpha} \\
&= \left( \frac{m\alpha(\beta - \alpha x^2 + \alpha x)}{(1-\alpha x) \cdot m(\beta+\alpha x)} \right) \cdot \frac{1}{m\alpha} \\
&= \frac{\beta - \alpha x^2 + \alpha x}{(1-\alpha x) \cdot m(\beta+\alpha x)}
\end{aligned}
$$

and therefore both efficiencies compute as:

$$E_A = m R_A = m \cdot \frac{\beta - \alpha x^2 + \alpha x}{(1-\alpha x) \cdot m(\beta+\alpha x)} = \frac{\beta - \alpha x^2 + \alpha x}{(1-\alpha x)(\beta+\alpha x)}$$

and

$$E_B = m R_B = m \cdot \frac{\beta}{(1-\alpha x) \cdot m(\beta+\alpha x)} = \frac{\beta}{(1-\alpha x)(\beta+\alpha x)}$$

This shows that $E_A > E_B$ because

$$E_A - E_B > 0$$

$$\frac{\beta - \alpha x^2 + \alpha x}{(1 - \alpha x)(\beta + \alpha x)} - \frac{\beta}{(1 - \alpha x)(\beta + \alpha x)} =$$

$$\frac{\alpha(x - x^2)}{(1 - \alpha x)(\beta + \alpha x)} > 0$$

This is true because $(x - x^2) > 0$ as $x < 1$. This further shows that $E_B < 1$ because:

$$E_B = \frac{\beta}{(1 - \alpha x)(\beta + \alpha x)}$$

$$= \frac{\beta}{\beta + \alpha x - \beta \alpha x + \alpha^2 x^2}$$

$$= \frac{\beta}{\beta + \underbrace{\alpha x}_{>0} \underbrace{(1 - \beta + \alpha x)}_{>0}} < 1$$

since the denominator is greater than the numerator. Further we show when $E_A > 1$:

$$E_A > 1$$

$$\frac{\beta - \alpha x^2 + \alpha x}{(1 - \alpha x)(\beta + \alpha x)} > 1$$

$$\frac{\beta - \alpha x^2 + \alpha x}{\beta - \beta \alpha x + \alpha x - \alpha^2 x^2} > 1$$

$$\beta - \alpha x^2 + \alpha x > \beta - \beta \alpha x + \alpha x - \alpha^2 x^2$$

$$-\alpha x^2 > -\beta \alpha x - \alpha^2 x^2$$

$$0 > -\beta \alpha x - \alpha^2 x^2 + \alpha x^2$$

$$0 > x(-\beta \alpha - \alpha^2 x + \alpha x)$$

which is the case when

$$0 > -\beta \alpha - \alpha^2 x + \alpha x$$

$$\beta \alpha > x(-\alpha^2 + \alpha)$$

$$\frac{\beta \alpha}{-\alpha^2 + \alpha} > x$$

$$x < \frac{\beta \alpha}{\alpha(1 - \alpha)} = \frac{\beta}{(1 - \alpha)}$$

# References

Eyal, I., Sirer, E.G.: Majority is not enough: bitcoin mining is vulnerable. In: Conference Paper (2013)

Eyal, I.: The Miner's dilemma. In: IEEE Symposium, Security and Privacy, pp. 89–103 (2015)

Eyal, I., Sirer, E.G.: Majority is not enough: bitcoin mining is vulnerable. Commun. ACM **61**, 95–102 (2018)

Lee, S., Kim, S. Countering block withholding attack efficiently. In: IEEE INFOCOM (2019)

Luu, L., Saha, R., Parameshwaran I., Saxena P., Hobor A.: On power splitting games in distributed computation: the case of bitcoin pooled mining. In: 2015 IEEE 28th Computer Security Foundations Symposium (2015)

Rosenfeld, M.: Analysis of bitcoin pooled mining reward systems. arXiv, 1112.4980 (2011)

Zhu, S., Li, W., Li, H., Hu, C., Cai, Z.: A survey: reward distribution mechanisms and withholding attacks in Bitcoin pool mining. Math. Found. Comput. **1**(4), 393–414 (2018)