

# Malleability of the blockchain’s entropy

Cécile Pierrot<sup>1</sup> and Benjamin Wesolowski<sup>2</sup>

<sup>1</sup> Sorbonne Universités, UPMC Univ Paris 06, LIP6, 4 place Jussieu, 75005 Paris, France

<sup>2</sup> EPFL IC LACAL, Station 14, CH-1015 Lausanne, Switzerland

**Abstract.** Trustworthy generation of public random numbers is necessary for the security of a number of cryptographic applications. It was suggested to use the inherent unpredictability of blockchains as a source of public randomness. Entropy from the Bitcoin blockchain in particular has been used in lotteries and has been suggested for a number of other applications ranging from smart contracts to election auditing. In this *Article*, we analyse this idea and show how an adversary could manipulate these random numbers, even with limited computational power and financial budget.

**Keywords:** Random number generation, Blockchain, Random Beacon, Bitcoin, Dyck language.

## 1 Introduction

Randomness is a key concept in cryptography. Even though random numbers are often meant to be kept secret in cryptographic protocols, a number of applications require a form of publicly available randomness which cannot be predicted or manipulated. Some straightforward examples include draws of national lotteries, sampling of assemblies or citizen juries, or tie-breaking in sports or elections. This concept also appears in more technical settings, where a secure public source of randomness can be used to provide trust between communicating parties while implementing contract signing, confidential disclosures, or certified mail in an electronic mail system. These examples were first proposed by Rabin in [Rab83], where he defines the source of randomness as a *random beacon*: an online service that broadcasts allegedly unpredictable numbers at regular intervals.

In these situations the randomness should be generated in such a way that no one can willfully bias the outcome to anyone’s advantage or disadvantage. Simple methods hardly provide any strong guarantee, like a skilled prestidigitator can fool entire crowds while tossing a coin. The existing online services that make available fresh random numbers at regular intervals such as [NIS11,RAN98] can at best be considered as trusted third parties: despite extensive documentation and audits, they do not provide users with any mechanism that allows to verify the freshness and correct generation of the published numbers.

A line of research aims at developing and analysing methods to provide trustworthy public randomness. Clark and Hengartner proposed in [CH10] a beacon based on stock market prices, whose security relies on a presumed unmalleability of published financial data – hardly sufficient to convince the most skeptical users. In [LW15], Lenstra and Wesolowski constructed and proved the security of a verifiable beacon that anyone can influence, but no one can bias or predict, by feeding some public input data to a slow hash function.

In another attempt to build a decentralised beacon with no trusted parties, it has been suggested to exploit the inherent unpredictability of blockchains, in particular that of the Bitcoin protocol [Nak09]. In addition to the applications mentioned above, a source of randomness internal to Bitcoin – or any other blockchain based currency – allows that randomness to be used by the currency’s scripting language. This extends the range of implementable smart contracts considerably. There is currently no possibility in the Bitcoin scripting language to allow any form of random execution. It is still not the case in cryptocurrencies with more powerful scripting languages

such as Ethereum [Woo14]. A few Bitcoin protocols [BCG15,BNM<sup>+</sup>14,GGMR14,MGGR13] are however already relying on the blockchain as a source of randomness.

This idea becoming widespread together with the lack of a security model for the blockchain’s unpredictability motivated Bonneau, Clark and Goldfeder [BCG15] to propose an analysis of the method. They consider a model where the attacker controls all the miners, and suffers a monetary penalty whenever a miner finds a valid block that the attacker does not want to see included in the blockchain (see Paragraph 2.1 for a presentation of the notions of miners and valid blocks). This penalty is meant to compensate for the loss of the reward going to a miner finding a valid block. This reward currently consists of  $R_M = 12.5 \text{ ₿}$  (12.5 bitcoins<sup>3</sup>) per block plus some transaction fees in the case of the Bitcoin protocol. The amount of this reward decays over time: it started at 50 ₿, and is divided by two every 210 000 blocks. In [BCG15], the authors derive that if the next block is used to define a single random bit – presumably unbiased –, an attacker can force the bit to 1 by having the miners drop the valid blocks that give a 0, and this attacker would suffer an expected cost of 50 ₿ at the moment the paper was published.

We aim in the present paper at overcoming the limitations of this model. First, we measure the adversary in terms of computational resources rather than monetary. An associated cost in bitcoins as in [BCG15] can then be derived, but is not the primary focus. It makes sense indeed to investigate computational limitations when not only money is at stake in the random number generation. Also, it allows us to analyse adversaries that do not have control over the entire pool of miners: it seems unrealistic to assume that the adversary can bribe any miner, who would systematically and silently cooperate (this is discussed and nuanced in [Bon16] with models considering the possibility of renting mining power for a limited amount of time via bribery). We also account for adversaries with limited resources: whereas in the original analysis the adversary keeps spending money in bribes as long as its goal is not achieved (e.g., obtaining a bit 1 in the example above), our analysis can measure the extent to which the adversary can bias the outcome if they have a limited budget. It also accounts for more general strategies, where the adversary can take advantage of non canonical phenomena on the blockchain such as forking.

Whereas the authors of [BCG15] concluded from their analysis that the method is secure enough to be used in practice, we show that an attacker, even when not controlling the entire pool of miners and when limited by a budget, can still significantly bias the probability distribution of the outcome. For instance, again when the blockchain is applied to generate a single unbiased bit, we show that an adversary controlling only a quarter of the total mining power can increase the probability of having a 1 from 0.5 to 0.6 approximately, for an expected cost of around 12.5 ₿, i.e. US\$ 53 000. With the same parameters we also propose a cheaper strategy that permits to increase the probability from 0.5 to about 0.57, with an expected cost of 1.8 ₿  $\approx$  US\$ 7 600. We also show how this probability can be further improved by preparing the attack in advance: a head start of a day increases the probability to 0.74.

Also, the security properties that we define and analyse might be relevant beyond the problem of generating public random numbers, as a more general measure of the security of the blockchain construction. Namely, we study the malleability of the blockchain, a formalisation of the capacity of a party to influence the probability distribution of some bits in the chain. Even though unmalleability is not the primary goal of a secure public ledger, it comes as a desirable property of an ideal instantiation of the blockchain.

The details of the model, together with a brief introduction to blockchain protocols, are presented in Section 2. In Section 3, we study the simple case where the adversary simply wants to influence the next published block with immediate effect, as an example and a warm-up. The more complex setting is analysed in Section 4. Adversaries limited by a financial budget

---

<sup>3</sup> Note that 1 ₿ roughly equals US\$ 4 260 at the time of writing this article, keeping in mind that the Bitcoin exchange rate still fluctuates a lot.

are analysed in Section 5 whereas adversaries preparing themselves in advance are studied in Section 6.

## 2 A model for the network and the miners

### 2.1 Blockchains

A blockchain is a form of distributed database, introduced in [Nak09] as the backbone of the Bitcoin protocol. The protocol involves a set of active participants, called miners, distributed over a network, whose job consists in maintaining the blockchain: a global ledger of the history of the system. In the case of Bitcoin, this history is the set of all transactions that happened. At regular time intervals, a new block is added to the chain, testifying for the latest transactions or events. In other words, a block consists of the concatenation of various pieces of information, including a reference to the hash of the previous block, a list of transactions the miner is willing to account for, and finally a nonce, chosen by the miner.

A block is called *valid* if its hash is smaller than a target number. This constraint is controlled by a difficulty parameter  $d$  measuring approximately the number of leading zeros of the hash. The difficulty parameter is continuously adjusted so that it takes an expected time of 10 minutes for the worldwide mining effort to find a valid block. For instance, at the time of writing this article, the value of  $d$  for Bitcoin was approximately 68. Only valid blocks can be added to the blockchain. The process of finding a valid block is called mining, and it simply consists in hashing a block a tremendous amount of times with slight modifications (for instance, by varying the nonce) until a valid version is found.

In addition to the list of new transactions, a block also contains the hash of the previous valid block. The set of all broadcast blocks thus forms a tree, whose longest chain is the valid blockchain on which there is a consensus. In the case of Bitcoin, each block on that longest chain comes with a reward (some bitcoins) to the miner who mined it. This mechanism encourages the miners to always work on increasing the longest chain they are aware of, in order to keep a mostly linear tree: a tree that consists in the longest chain itself and a few small branches, that quickly stop growing. When such a branching occurs at the tail of the blockchain, it is called a fork. Miners might work on a different branch of the fork at a given moment, but if they keep agreeing to work on the longest one, only the branch with a majority of miners will survive, and consensus is quickly reached again.

It is justified in [BCG15, Paragraph 3.1] that assuming the underlying hash function is secure, each valid block found via the mining process contains at least  $d$  bits of computational min-entropy, and therefore  $\lfloor d/2 \rfloor$  near-uniform bits can securely be extracted by a cryptographic extractor [DGH<sup>+</sup>04]. An *extractor* is a function that transforms a random variable with sufficient entropy into an almost uniformly distributed variable in a smaller domain, which we call the *extract*. In this paper we consider an extractor  $\text{Ext}$  that maps valid blocks to elements of  $\mathcal{E}$ , the set of binary strings of length  $\lfloor d/2 \rfloor$ . Assuming the hash function used for mining is secure, we can consider the extract  $\text{Ext}(B)$  of a freshly generated valid block  $B$  to be uniformly distributed in  $\mathcal{E}$ .

### 2.2 Game model

We consider a single instance of a blockchain protocol executed in isolation. The quite detailed abstraction proposed in [GKL15, Section 2] fits our needs, yet we will not detail it here. Indeed, such a precise network diffusion model is not necessary as we will not be dealing with network adversaries: they will not be able to tamper with messages on transit in the network, or with the network's structure itself. This weaker adversary receives and broadcasts messages over the network, and its power will only be measured in terms of its computational resources.

As in Garay et al.’s formal analysis of the Bitcoin backbone protocol [GKL15], the problem is made tractable by assuming that communications over the network are instantaneous. Honest miners are always working with the longest blockchain that they know of, and the broadcast messages are always received by all the miners. Therefore forks can appear when a valid block is found simultaneously by two miners, but the fork will disappear as soon as one of the branches will be lengthened before the other. This matches the real world situation, where the notion of simultaneity is loosened to take into account communication and processing delays: a fork only survives as long as new blocks are added simultaneously to both branches, which cannot happen for very long. We argue that accidental forks among the honest miners can only help the adversary, and therefore we assume for the rest of the paper that such forks do not occur. Even without accidental forks, malicious forks are still possible: an adversary can secretly work on a branch longer than the public one, and once this branch is revealed, it will completely replace the old one as the new longest chain. Other models such as [SSZ16] propose to look at adversaries that are highly connected to the rest of the network, giving them a timing advantage when sending or receiving new blocks. Again, this would only make the adversary stronger, so we do not consider this possibility.

The game involves a set of honest miners and an adversary  $\mathcal{A}$ . Honest miners simply aim at growing the blockchain, by broadcasting valid blocks to get the associated reward. Meanwhile, the goal of the adversary is to bias the probability distribution of the extract of an upcoming block in the chain, which we call the *decisive block*. For illustration, say that the extract  $x$  of that decisive block is a lottery draw: the adversary wins whenever  $x$  falls in a given subset of  $\mathcal{E}$ , the “winning tickets”. Let  $\chi : \mathcal{E} \rightarrow \{0, 1\}$  be the characteristic function of that subset of the set of possible extract values. The adversary  $\mathcal{A}$  wins if, at the end of the game,  $\chi(x) = 1$  where  $x$  is the extract of the decisive block, which we call *pleasant* in this case. Note that it does not matter whether that decisive block was broadcast by the adversary or by a honest miner. We denote by  $\mathcal{P}_{\mathcal{A}}$  the probability that  $\mathcal{A}$  wins.

The game starts at the round where a fixed *initial block* indexed by 0 is received and ends when the blockchain reaches a length of  $n + \Delta$ , for some predefined integers  $n$  and  $\Delta$ . The  $n$ -th block is the decisive block, whose extract determines whether or not the adversary wins. Since that block determines the outcome of the game, it might first seem surprising to continue for  $\Delta$  additional blocks. But the value of any specific block can be changed via forking: in blockchain protocols, consensus does not go to the version of a block that arrived first, but to the version of the block that belongs to the longest chain. An adversary could try to take advantage of that by attempting to create a fork to modify an unpleasant decisive block. The positive parameter  $\Delta$  takes into account the possibility that the decisive block can still be modified until the blockchain has been lengthened by  $\Delta$  blocks.

Finally we denote by  $\mu$  the probability for a uniformly distributed extract value  $x \in \mathcal{E}$  to satisfy  $\chi(x) = 1$ . We clearly have  $\mu = |\{\chi(x) = 1 \mid x \in \mathcal{E}\}|/|\mathcal{E}|$ .

### 2.3 Adversary model

We consider a weak adversary that cannot change the content of the messages traveling around the network nor prevent them from being delivered. The adversary can mine on the blockchain and, unlike honest miners, they are not assumed to always work on extending the longest chain, nor to broadcast all the valid blocks they find. Thus, in this context, the critical property of the adversary is their computing power, or more precisely, their relative power with respect to the total power of the miners. Concretely, this power is measured as a number  $s_{\mathcal{A}}$  (respectively  $s$ ) of hashes per second that the adversary (respectively the total set of miners, including  $\mathcal{A}$ ), is able to compute. Let

$$\alpha = \frac{s_{\mathcal{A}}}{s}$$

measure the relative power of  $\mathcal{A}$ . Such hash computations can be modelled as calls to a random oracle. The number of random oracle calls to be done before finding a valid block follows a Poisson distribution, and at any point in time,  $\alpha$  can be interpreted as the probability that the next oracle request giving a valid block has been done by the adversary.

*Remark 1.* The adversary could be a single entity as well as a set of colluding, corrupted miners. It therefore extends the model of [BCG15], where all the miners are ready to accept bribes to collaborate: this corresponds to  $\alpha = 1$  in our setting. We can still account for the money it takes to corrupt the miners, via the notion of *virtual cost* (see Paragraph 3.4), yet it is not the only metric available anymore.

### 3 A warm-up: manipulating the next published block

#### 3.1 Delay-free games

As a first example, let us look at the simple case where  $n = 1$  and  $\Delta = 0$ . The adversary aims at manipulating the next broadcast block, without bothering whether it will survive potential future forks. This case occurs whenever the extracted randomness is meant to be used immediately, such as a random beacon that broadcasts the random numbers as soon as they are available, or a simple national lottery. Indeed, a random beacon is meant to be unpredictable, and its output cannot be changed in the future – unlike a blockchain where forks can “change the past”. The beacon must broadcast a number as soon as it is available, which does not permit to wait for a safety delay  $\Delta > 0$  before outputting the numbers.

#### 3.2 Strategy

The adversary’s strategy here is straightforward: they will focus all their computational power on mining the next block. Whenever they find an unpleasant valid block  $B$  such that  $\chi(\text{Ext}(B)) = 0$ , they throw it away and continue to mine. The game ends whenever a honest miner finds a valid block or the adversary finds one that satisfies  $\chi(\text{Ext}(B)) = 1$ , which they immediately broadcast. In this context, it is clearly an optimal strategy as valid blocks can only be found by mining (which is a classical assumption for blockchains, and can be enforced by the random oracle model), and an adversary who would themselves publish a block such that  $\chi(\text{Ext}(B)) = 0$  would only decrease their winning probability, as well as an adversary who would not publish instantly a valid block such that  $\chi(\text{Ext}(B)) = 1$ . Note that this strategy has a certain cost, as a dropped block is a lost reward. A rational adversary with a financial incentive will want to compare the cost of dropping a block with the expected reward of winning the game. We analyse this cost in Section 3.4.

#### 3.3 The winning probability $\mathcal{P}_{\mathcal{A}}$

Let us compute the probability  $\mathcal{P}_{\mathcal{A}}$  that the adversary wins that game. The event that they win the game is denoted  $\mathfrak{w}$ . The possible outcomes can be split into the infinite family of disjoint events  $\epsilon_a$ , for any  $a \geq 0$ , that the adversary found and threw away exactly  $a$  valid blocks. The event  $\epsilon_a \wedge \mathfrak{w}$  further requires that the final block  $B$ , i.e. the  $(a+1)$ -th block that is found, satisfies  $\chi(h_B) = 1$  (and it does not matter who found that last one). The winning probability can now be expressed as:

$$\mathcal{P}_{\mathcal{A}} = \sum_{a \geq 0} \Pr(\epsilon_a \wedge \mathfrak{w}) = \sum_{a \geq 0} (\alpha(1 - \mu))^a \mu = \frac{\mu}{1 - \alpha(1 - \mu)}.$$

Since  $\alpha \geq 0$  and  $\mu \leq 1$  this is clearly better than  $\mu$ . For instance, when  $\mu = 1/2$ , and  $\alpha = 1/4$ , the probability that  $\mathcal{A}$  wins is increased from 0.5 to  $4/7 \approx 0.57$ . Note that  $\alpha = 1/4$  is a realistic assumption: for instance, the relative computational power  $\alpha$  of the four actual most influent mining pools Antpool, F2Pool, BitFury, and BTCChinaPool vary between 13% and 27% each, depending on sources and dates.

### 3.4 Virtual cost expectation for the adversary

The *virtual cost* refers to the number of blocks the adversary has to drop, thereby giving up the associated reward. We emphasise that this virtual cost does not take into account potential profits related to the outcome of the game. It is only a measure of the amount of mining power they sacrificed to play the game: namely, the difference between the expected reward they would have had by behaving as a honest miner, and the reward resulting from playing the described strategy. For each event  $\epsilon_a$ , exactly  $a$  blocks are dropped, therefore the expected virtual cost  $E$  is:

$$E = \sum_{a \geq 0} a \Pr(\epsilon_a) = \sum_{a \geq 0} a (\alpha(1 - \mu))^a (1 - \alpha(1 - \mu)) = \frac{\alpha(1 - \mu)}{1 - \alpha(1 - \mu)}.$$

Again, to illustrate this virtual cost we look at the example where  $\mu = 1/2$  and  $\alpha = 1/4$ . It gives  $1/7$  expected lost blocks. In terms of money, the actual mining reward at the time of writing is  $R_M = 12.5 \text{ ₿}$ . Thus this strategy can be approximately evaluated to have an expected cost of US\$ 7600. The case where  $\alpha = 1$  ensures the victory as long as  $\mu > 0$ , for an expected cost of  $(1 - \mu)/\mu$ . This is the situation in [BCG15] where all the miners are assumed to collude with the adversary, and they indeed found the same expected cost.

### 3.5 Financial game

In some situations, the goal of the adversary can be modelled as a monetary reward  $R_W$  associated to winning the game. In this situation, whenever they find a unpleasant block, a rational adversary will be comparing their expected earnings if they stick to the strategy (i.e., the expected reward minus the virtual cost,  $\mathcal{P}_A \cdot R_W - E \cdot R_M$ ) with the mining reward they would get by publishing the block (i.e.,  $R_M$ ). The adversary can take the decision of following the strategy or not once and for all at the beginning of the game since these quantities do not evolve until the decisive block finally gets published.

## 4 Reaching consensus

### 4.1 Games with delay

We now deal with the case where  $n = 1$  and  $\Delta > 0$ . Delaying the end of the game later than the broadcast of the decisive block is necessary to measure the security of protocols involving several parties that would base an agreement on the extract value of this decisive block. Smart contracts are a good example for such worries. Indeed, unlike in our first naive delay-free games model, forking can interfere and change the initial execution of the contract, thereby modifying the outcome of the contract itself after its execution. Here, the random value is only considered valid when it is deep enough in the chain. More precisely, choosing a sufficiently large  $\Delta$  helps to reach consensus, when it is considered impossible to create a fork of length strictly larger than  $\Delta$ . Then the parties involved in the contract gain confidence that the outcome is definitive.

## 4.2 Relevant values for $\Delta$

Relevant values for the delay  $\Delta$  can be deduced from the analysis performed in [GKL15] where the authors evaluate the probability to have a fork of two branches of a given length. From a practical point of view, a look back to the history of broadcast Bitcoin blocks permits to find all the accidental forks that have ever been observed. For instance, from [Blo16] we see that 540 forks were noticed since the block number<sup>4</sup> 142 257: two branches of length 4, four of length 3, fourteen of length 2, and what remains are orphan blocks. Thus, keeping a safety margin,  $\Delta = 6$  seems sufficient to make sure that no fork will alter the decisive block after the end of the game. This value of  $\Delta$  is already suggested in [Nak09].

## 4.3 Strategy

The adversary’s strategy begins as previously: they first focus all their computational power on mining the first block (that is also the decisive block since  $n = 1$ ) and throw away all valid blocks  $B$  they find as long as  $\chi(\text{Ext}(B)) = 0$ . At the first broadcast of the decisive block  $B_d$  two cases can occur. If it satisfies  $\chi(\text{Ext}(B_d)) = 1$  then  $\mathcal{A}$  starts to act as a honest miner. Namely, they continue to mine on this chain and the game ends when the  $(1 + \Delta)$ -th valid block is found, leading to the success of  $\mathcal{A}$ . Otherwise, if one honest miner broadcasts a decisive block satisfying  $\chi(\text{Ext}(B_d)) = 0$ ,  $\mathcal{A}$  continues to secretly mine on the previous block as if they received nothing. Again, they throw away all valid decisive blocks that would lead to their failure. If they find one such that  $\chi(\text{Ext}(B_d)) = 1$  they keep it secret<sup>5</sup>, create a fork and continue to mine on their new chain, hoping to develop a branch longer than the current public one containing the unpleasant decisive block. If they manage to do so, they broadcast their fork, and all miners switch to that new longest chain and proceed to mine on that one, thereby including the cheated pleasant decisive block in the new consensus chain. The game ends when a miner finds a valid final block, whether it is on the original chain, or on the branch later appeared.

## 4.4 Computing the probability $\mathcal{P}_{\mathcal{A}}$

To compute the probability that  $\mathcal{A}$  wins we consider a partition of our universe as follows. Let  $\mathfrak{w}_1, \mathfrak{f}_0, \mathfrak{f}_{01}$  and  $\mathfrak{w}_{01}$  ( $\mathfrak{w}$  implicitly stands for  $\mathcal{A}$  wins and  $\mathfrak{f}$  for  $\mathcal{A}$  fails) be the events defined as:

- $\mathfrak{w}_1$ : the first broadcast decisive block  $B_d$  is such that  $\chi(\text{Ext}(B_d)) = 1$ .
- $\mathfrak{f}_0$ : the first broadcast decisive block  $B_d$  is such that  $\chi(\text{Ext}(B_d)) = 0$  and  $\mathcal{A}$  fails to find any decisive block  $B'_d$  such that  $\chi(\text{Ext}(B'_d)) = 1$ .
- $\mathfrak{f}_{01}$ : the first broadcast decisive block  $B_d$  is such that  $\chi(\text{Ext}(B_d)) = 0$ ,  $\mathcal{A}$  manages to find a decisive block  $B'_d$  such that  $\chi(\text{Ext}(B'_d)) = 1$  but not to compute a branch sufficiently long to replace the main one.
- $\mathfrak{w}_{01}$ : the first broadcast decisive block  $B_d$  is such that  $\chi(\text{Ext}(B_d)) = 0$  but the adversary manages to find one  $B'_d$  verifying  $\chi(\text{Ext}(B'_d)) = 1$ , and to have it replace the original one.

So the event that  $\mathcal{A}$  wins, denoted  $\mathfrak{w}$ , is the union of  $\mathfrak{w}_1$  and  $\mathfrak{w}_{01}$ . Note that the probability of  $\mathfrak{w}_1$  has already been computed in Section 3, as:

$$\Pr(\mathfrak{w}_1) = \frac{\mu}{1 - \alpha(1 - \mu)}. \quad (1)$$

For convenience we also define  $\mathfrak{e}_0$  as  $\mathfrak{e}_0 = \mathfrak{f}_0 \vee \mathfrak{f}_{01} \vee \mathfrak{w}_{01}$ .

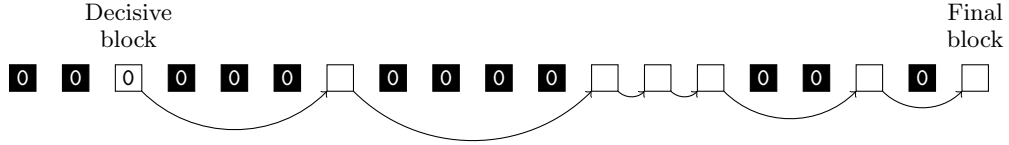
<sup>4</sup> At the time of writing the last blocks are around number 400 000.

<sup>5</sup> Note that keeping it secret is not critical in the strategy: the adversary could be completely public about this special fork, in the hope that some honest miners get confused and start unknowingly helping them (due to bad timing and propagation delays for instance).

**Probability of  $f_0$ .** Since  $f_0$  is a sub-event of  $e_0$ , its probability can be expressed as:

$$\Pr(f_0) = \Pr_{e_0}(f_0) \times \Pr(e_0) = \Pr_{e_0}(f_0) \times (1 - \Pr(w_1)),$$

where we denote by  $\Pr_B(A)$  the conditional probability of event  $A$  given that event  $B$  occurred.



**Fig. 1.** Secret and public blocks mined during the event  $f_0$ . White blocks represent the blockchain integrally mined by honest miners whereas black ones correspond to the  $k$  valid decisive blocks (secretly) found by  $\mathcal{A}$ , all such that  $\chi(\text{Ext}(B)) = 0$ . A bit inside a block indicates the value of  $\chi$  for its extract. It is only indicated on blocks that are meant to be decisive: here, the first broadcast block, and the adversary's failed attempts at finding a pleasant one. In this example,  $a = 12$  and  $\Delta = 6$ .

Figure 1 represents the situation for a fixed number  $a$  of valid blocks mined by  $\mathcal{A}$ . The probability that  $\mathcal{A}$  fails precisely  $a$  times knowing  $e_0$  is given by the product:

$$\binom{a + \Delta - 1}{a} (\alpha(1 - \mu))^a (1 - \alpha)^{\Delta - 1} (1 - \alpha).$$

Thus the formula for  $\Pr_{e_0}(f_0)$  is:

$$(1 - \alpha)^\Delta \sum_{a \in \mathbf{N}} \binom{k + \Delta - 1}{a} (\alpha(1 - \mu))^a.$$

The summation over  $a$  can be simplified using the fact that for any positive integer  $h$ , and any  $x \neq 1$ :

$$\sum_{a \in \mathbf{N}} \binom{a + h}{a} x^a = \frac{1}{(1 - x)^{h+1}}. \quad (2)$$

From Equation (2) we get the final value of the probability of  $f_0$  knowing  $e_0$ :

$$\Pr_{e_0}(f_0) = \left( \frac{1 - \alpha}{1 - \alpha(1 - \mu)} \right)^\Delta. \quad (3)$$

**Probability of  $f_{01}$ .** Again we compute  $\Pr_{e_0}(f_{01})$  the probability of  $f_{01}$  knowing  $e_0$  in order to find the probability of  $f_{01}$  that satisfies:

$$\Pr(f_{01}) = \Pr_{e_0}(f_{01}) \times (1 - \Pr(w_1)).$$

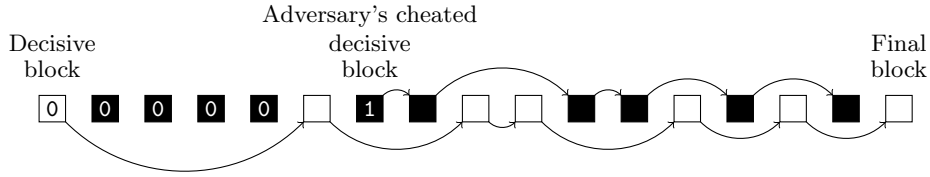
Figure 2 represents the situation and gives an indication of the various parameters we need to express this probability. We define:

- $a_2$  as the number of valid decisive blocks mined by  $\mathcal{A}$  after the first decisive block is broadcast by a honest miner and before  $\mathcal{A}$  finds a valid decisive block verifying  $\chi(\text{Ext}(B_d)) = 1$ .



- $h_2$  as the number of valid blocks mined by all honest miners in the blockchain in the meantime, namely after the first decisive block is broadcast by a honest miner, and before the adversary manages to find a decisive block verifying  $\chi(\text{Ext}(B_d)) = 1$ .
- $a_3$  as the number of valid blocks mined by  $\mathcal{A}$  on its chain after  $\mathcal{A}$  finds a valid decisive block verifying  $\chi(\text{Ext}(B_d)) = 1$  and before the last block is mined.

With this in hands, we see that the number of valid blocks mined by all honest miners, after  $\mathcal{A}$  finds a pleasant valid decisive block, and before the last block is mined, is given by  $\Delta - 1 - h_2$ . Besides, some parameters are upper bounded: we have  $h_2 < \Delta$  and  $a_3 < \Delta$ .



**Fig. 2.** Secret and public blocks mined during the event  $f_0$ . White blocks represent the blockchain integrally mined by honest miners whereas black ones correspond to the  $a_2 + 1 + a_3$  valid blocks (secretly) found by  $\mathcal{A}$ . Again a number inside a block corresponds to the value of  $\chi$  for its extract, when the block is meant to be a decisive one. In this example  $\Delta = 6$ ,  $a_2 = 4$ ,  $h_2 = 1$  and  $a_3 = 5$ .

The probability, knowing  $\epsilon_0$ , that  $\mathcal{A}$  finds a pleasant valid decisive block exactly after  $a_2$  trials that failed, and to get  $h_2$  blocks mined by the honest miners before this new decisive block is:

$$\binom{a_2 + h_2}{a_2} \alpha^{a_2+1} (1 - \alpha)^{h_2} (1 - \mu)^{a_2} \mu.$$

To complete the computation of  $\text{Pr}_{\epsilon_0}(f_{01})$  we need to count the number of possible cases in which  $\mathcal{A}$  mines  $a_3$  valid blocks but fails to ever create a chain longer than the main one. This number in hand, it will only remain to sum over the possible values of the three variables  $a_2$ ,  $a_3$  and  $h_2$ . Fixing these parameters, an outcome where the adversary fails to replace the main blockchain with its own corresponds to a sequence where *at every point in time*, the adversary's chain is shorter than the main one. The insightful reader might have observed that this situation shares similarities with Dyck words. A word in two letters  $H$  and  $A$  is a *Dyck word* if no initial segment contains more  $A$ 's than  $H$ 's. Here the situation is slightly different as the honest main chain (whose blocks correspond to the letter  $H$ ) starts growing as soon as a honest miner finds a first valid block, while the adversary's chain (the  $A$ 's) starts growing only when the adversary finds a first valid block such that  $\chi(\text{Ext}(B)) = 1$ . Therefore we define *generalised Dyck words* as follows.

**Definition 1.** Let  $\Delta$ ,  $a$  and  $h$  be three non-negative integers such that  $a < \Delta$  and  $h < \Delta$ . We call generalised Dyck words (for the parameters  $\Delta$ ,  $a$ , and  $h$ ) words with  $a + h$  letters formed with exactly  $a$  letters  $A$  and  $h$  letters  $H$  that verify the following property: for any prefix of the word, we have  $\#A \leq \#H + \Delta - h - 1$  where  $\#H$  (respectively  $\#A$ ) represents the number of  $H$ 's (respectively  $A$ 's) in that prefix.

**Lemma 1.** Let  $\Delta$ ,  $a$  and  $h$  be three non-negative integers, and suppose  $a < \Delta$  and  $h < \Delta$ . Defining  $D(\Delta, a, h)$  as the number of generalised Dyck words, we exactly have:

$$D(\Delta, a, h) = \binom{a+h}{a} - \binom{a+h}{\Delta}.$$

*Proof.* The proof follows a classical idea developed to count Dyck words. We first look at words with exactly  $a$  letters  $A$  and  $h$  letters  $H$  that are not generalised Dyck words. For each such word, there exists a first  $A$  that breaks the property. Changing each  $A$  into an  $H$  and each  $H$  into an  $A$  strictly after this special  $A$  leads to write a word<sup>6</sup> with exactly  $a - \Delta + h$  letters  $H$  and  $\Delta$  letters  $A$ . Thus for a fixed  $\Delta$  this process describes an injection:

$$f : E \longrightarrow F$$

that sends  $E$  (the set of all the words with  $a$  letters  $A$ ,  $h$  letters  $H$  that are not generalised Dyck words) into  $F$  (the set of words with  $a - \Delta + h$  letters  $H$ ,  $\Delta$  letters  $A$ ). Besides, the map  $f$  is surjective. Indeed, let us consider a word  $w$  in  $F$ . From  $a < \Delta$  we see that  $\Delta > (a - \Delta + h) + \Delta - h - 1$ , which means that  $w$  breaks the condition  $\#A \leq \#H + \Delta - h - 1$ . Hence considering the first letter  $A$  in  $w$  that breaks the condition, and changing each  $A$  into an  $H$  and each  $H$  into an  $A$  strictly after this special  $A$ , we obtain a word in  $E$  such that the image of this word is exactly  $w$ . The bijection  $f$  yields  $\#E = \binom{a+h}{\Delta}$ , and at the end:

$$D(\Delta, a, h) = \binom{a+h}{a} - \binom{a+h}{\Delta}.$$

Summing over the possible values of the three variables  $a_2$ ,  $a_3$  and  $h_2$ , we obtain:

$$\Pr_{\epsilon_0}(f_{01}) = \alpha\mu(1-\alpha)^\Delta \sum_{\substack{0 \leq a_2 \\ 0 \leq a_3, h_2 < \Delta}} (\alpha(1-\mu))^{a_2} \binom{a_2+h_2}{a_2} \alpha^{a_3} D(\Delta, a_3, \Delta - 1 - h_2),$$

where  $D$  is the function in 3 variables defined in Lemma 1 that counts the number of possible cases in which the chain mined by  $\mathcal{A}$  never overtakes the honest one. The probability  $\Pr_{\epsilon_0}(f_{01})$  can be rewritten as:

$$\Pr_{\epsilon_0}(f_{01}) = \alpha\mu(1-\alpha)^\Delta \sum_{0 \leq a_3, h_2 < \Delta} \alpha^{a_3} D(\Delta, a_3, \Delta - 1 - h_2) \sum_{0 \leq a_2} \binom{a_2+h_2}{a_2} (\alpha(1-\mu))^{a_2}.$$

We can simplify the inner sum using Formula (2) and get:

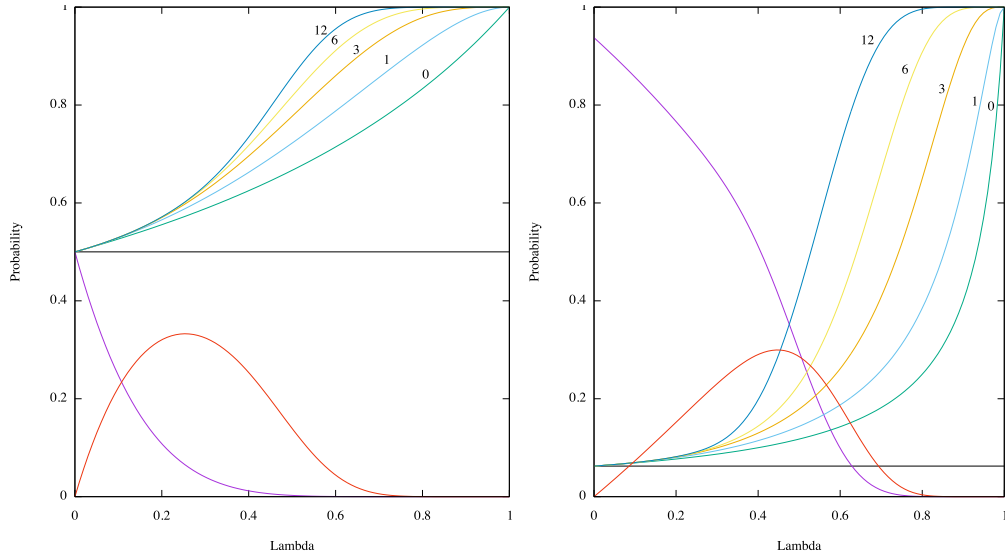
$$\Pr_{\epsilon_0}(f_{01}) = \frac{\alpha\mu(1-\alpha)^\Delta}{1-\alpha(1-\mu)} \sum_{0 \leq a_3, h_2 < \Delta} \frac{\alpha^{a_3}}{(1-\alpha(1-\mu))^{h_2}} D(\Delta, a_3, \Delta - 1 - h_2). \quad (4)$$

To conclude, we can put everything together using  $\mathcal{P}_{\mathcal{A}} = 1 - \Pr(f_0) - \Pr(f_{01})$ , with the formula

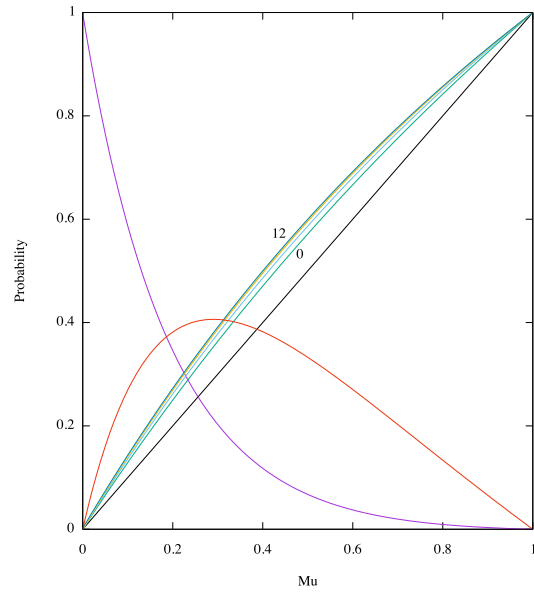
$$\mathcal{P}_{\mathcal{A}} = 1 - (1 - \Pr(\mathbf{w}_1)) \cdot (\Pr_{\epsilon_0}(f_0) + \Pr_{\epsilon_0}(f_{01})), \quad (5)$$

where  $\Pr(\mathbf{w}_1)$ ,  $\Pr_{\epsilon_0}(f_0)$  and  $\Pr_{\epsilon_0}(f_{01})$  are functions in  $\alpha, \mu$  and  $\Delta$  respectively given in Equations (1), (3) and (4). To compare again with  $\mu$  we plot several of these probabilities. Figure 3 gives the probability that  $\mathcal{A}$  wins as a function of their relative computational power  $\alpha$ , for a fixed value of  $\mu = 1/2$  on the left or  $\mu = 1/16$  on the right. For instance, if  $\mu = 1/2$ ,  $\alpha = 1/4$  and  $\Delta = 6$  – which are three arguably realistic parameters – then the probability that the adversary wins is increased from 50% to approximately 60%. Figure 4 takes the other point of view and considers how the probability  $\mathcal{P}_{\mathcal{A}}$  varies with  $\mu$  when  $\alpha$  is fixed to  $1/4$ .

<sup>6</sup> For instance, if  $h = 3$ ,  $a = 5$  and  $\Delta = 6$  the word HAAA**A**HHA becomes HAAA**A**AAH that has 6 letters  $A$  and  $5 - 6 + 3 = 2$  letters  $H$ .



**Fig. 3.** Probabilities  $\mathcal{P}_A$  as functions of  $\alpha$  when  $\mu = 1/2$  (first graph) and  $\mu = 1/16$  (second graph). The corresponding values for  $\Delta$ , ranging from 0 to 12, are indicated on the curves. For information purposes, the red curves correspond to  $\Pr(f_{01})$  and the purple ones to  $\Pr(f_0)$ , both with  $\Delta = 12$ .



**Fig. 4.** Probabilities  $\mathcal{P}_A$  as functions of  $\mu$  when  $\alpha = 1/4$ . Values for  $\Delta$  are indicated on the curves. Again the red curve corresponds to  $\Pr(f_{01})$  and the purple one to  $\Pr(f_0)$  both with  $\Delta = 12$ . The black line indicates the probability that a honest player wins.

## 4.5 Virtual cost expectation for the adversary

Again, the virtual cost refers to the number of blocks the adversary has to drop, thereby giving up the associated reward. To compute the expected virtual cost  $E$  we consider three cases. In the event  $\mathfrak{w}_1$ , such blocks only appear before the first decisive block is mined. Thus the expected number of dropped blocks that correspond to this case is:

$$E_1 = \sum_{a \geq 0} a \Pr(\mathfrak{e}_a \wedge \mathfrak{w}) = \sum_{a \geq 0} a (\alpha(1-\mu))^a \mu = \frac{\mu \alpha (1-\mu)}{(1-\alpha(1-\mu))^2}. \quad (6)$$

Besides, considering the union of  $\mathfrak{f}_0$  and  $\mathfrak{f}_{01}$ , we see that in both cases the chain making the consensus at the end has entirely been mined by honest miners and consists in  $\Delta+1$  valid blocks. Since the relative computing power of  $\mathcal{A}$  is  $\alpha$ , it yields that  $\mathcal{A}$  dropped in average  $(\Delta+1)\alpha/(1-\alpha)$  valid blocks. So the expected number of lost blocks related to those two events is:

$$E_2 = (1 - \mathcal{P}_{\mathcal{A}}) \cdot \frac{(\Delta+1)\alpha}{1-\alpha}. \quad (7)$$

where  $\mathcal{P}_{\mathcal{A}}$  is given in Equation (5). Last but not least, we can see the event  $\mathfrak{w}_{01}$  as the union of all  $\mathfrak{e}_a$  for  $a \in \mathbf{N}$  where  $\mathfrak{e}_a$  is the event:  $\mathcal{A}$  finds precisely  $a$  valid blocks with an unpleasant extract, then mines a pleasant decisive block and manages to mine a chain longer than the current honest one. The expected number of dropped blocks related to this last case is:

$$E_3 = \sum_{a \geq 0} a \Pr(\mathfrak{e}_a).$$

To compute  $\Pr(\mathfrak{e}_a)$  we denote by  $h_2$  the number of honest valid blocks broadcast (strictly) between the two decisive blocks and by  $h_3$  the number of honest valid blocks broadcast after the pleasant decisive one is mined by  $\mathcal{A}$ . We get:

$$\Pr(\mathfrak{e}_a) = \alpha^a (1-\mu)^{a+1} \mu \sum_{0 \leq h_2 < \Delta} \binom{a+h_2}{h_2} (1-\alpha)^{h_2+1} \sum_{0 \leq h_3 \leq \Delta} \bar{D}(h_2+h_3+1, \Delta, h_3) \alpha^\Delta (1-\alpha)^{h_3}$$

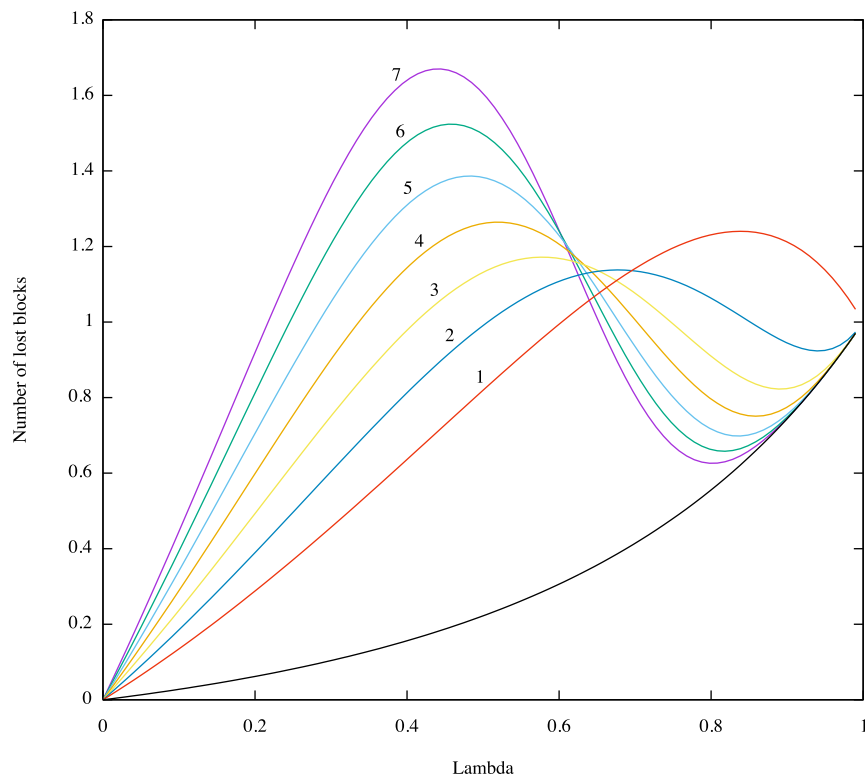
where  $\bar{D}(h_2+h_3+1, \Delta, h_3)$  is the number of words with  $\Delta$  letters  $A$  and  $h_3$  letters  $H$  that are *not* generalised Dyck words (with the corresponding parameters). Since we have the equality  $\bar{D}(h_2+h_3+1, \Delta, h_3) = \binom{h_2+h_3+1}{\Delta+h_3}$  we can write:

$$E_3 = \alpha^\Delta (1-\alpha)(1-\mu)\mu \times \sum_{a \geq 0} a (\alpha(1-\mu))^a \sum_{0 \leq h_2 < \Delta} \binom{a+h_2}{h_2} (1-\alpha)^{h_2} \sum_{0 \leq h_3 \leq \Delta} \binom{h_2+h_3+1}{\Delta+h_3} (1-\alpha)^{h_3} \quad (8)$$

Therefore the expected virtual cost  $E$  is given by:

$$E = E_1 + E_2 + E_3,$$

where  $E_1$ ,  $E_2$  and  $E_3$  are respectively given in Equations (6), (7) and (8). To illustrate this virtual cost we plot in Figure 5 the expected number of blocks that are dropped by the adversary as a function of their relative computational power  $\alpha$ , in the case where  $\mu = 1/2$ . For instance, when  $\alpha = 1/4$  and considering  $\Delta = 6$  it gives roughly 1.017 expected lost blocks. Thus this strategy can be approximately evaluated to have an expected cost of US\$ 53 000.



**Fig. 5.** Expected number of dropped blocks as functions of the relative computational power  $\alpha$ , when  $\mu = 1/2$ . The black curve corresponds to the basic case with no delay, i.e.  $\Delta = 0$ . The other curves correspond to values of  $\Delta$  that vary from 1 to 7.

## 4.6 Financial game

Consider again the situation where a monetary reward  $R_W$  is associated to winning the game. The above strategy is still valid, yet a rational adversary would adapt it a bit in order to maximize his earnings. We give here the basic ideas of this particular strategy, but do not provide a full analysis as the previous one already shows that the game is not resistant to malicious behaviors.

Dealing with a financial game,  $\mathcal{A}$  could compute at each step the expected cost of sticking to the previous strategy and decide whether or not to pursue their attack. More precisely, each time a block is found (either by  $\mathcal{A}$ , or broadcast by a honest miner),  $\mathcal{A}$  could reckon the two following expected earnings, in term of bitcoins:

- $G_1$ , their expected earnings at the end of the game, assuming they keep mining in secret on their chain, if they have any, or try to create a pleasant valid block
- $G_2$ , their expected earnings at the end of the game, assuming they withdraw from mining their current chain (if they managed to have one), and start mining as a honest miner.

Then  $\mathcal{A}$  keeps attacking the game as long as it is worth it, namely as long as:

$$G_2 < G_1.$$

Note that  $G_1$  and  $G_2$  are expressed in term of bitcoins, and not as dropped blocks anymore. Another situation might occur and prevent  $\mathcal{A}$  from carrying on their attack: if they have a budget, a limited amount of money they are ready to loose before aborting the strategy. We study this situation in more details in the next section.

## 5 Adversaries with a budget

A budget is the amount of money the adversary is ready to lose by not publishing valid blocks they found. In situations where winning the game cannot be easily modelled by a monetary reward, or the adversary's goal is not to maximize some earnings, considering a budget is a meaningful way to limit the power of an adversary. An adversary with a budget enforces a bound on the virtual cost, and thereby on the number of abandoned blocks. The analysis is similar to the previous case, but now, when the adversary runs out of budget, they must stop mining blocks that do not directly contribute to the longest chain.

Let  $b$  be the maximal number of blocks the adversary can drop, i.e., its budget. Then, the probability of  $\mathfrak{w}_1$  becomes:

$$\Pr(\mathfrak{w}_1) = \sum_{a=0}^b (\alpha(1-\mu))^a \mu = \frac{\mu(1 - (\alpha(1-\mu))^{b+1})}{1 - \alpha(1-\mu)}.$$

This is the probability to win when  $\Delta = 0$ , generalising the results of [BCG15] to adversaries with a budget. Let us compute the winning probability for  $\Delta > 0$ . Now, we cannot forget about what happened before the publication of the first decisive block as we did previously by conditioning over  $\epsilon_0$ , as we need to keep track of how many blocks the adversary dropped during that period.

Let  $\mathfrak{f}$  be the event that  $\mathcal{A}$  fails the game,  $\mathfrak{b}$  the event that  $\mathcal{A}$  used its full budget during the game, and  $\mathfrak{o}$  be the event that the blocks found by the adversary within their budget are all unpleasant. The event  $\mathfrak{f}$  is partitioned into four disjoint events:

$$\Pr(\mathfrak{f}) = \Pr(\mathfrak{f} \wedge \mathfrak{b} \wedge \mathfrak{o}) + \Pr(\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \mathfrak{o}) + \Pr(\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \neg \mathfrak{o}) + \Pr(\mathfrak{f} \wedge \mathfrak{b} \wedge \neg \mathfrak{o}).$$

Let us first focus on the event  $\mathfrak{f} \wedge \mathfrak{b} \wedge \mathfrak{o}$ . Look at the prefix of the sequence of found blocks up to the point where the adversary went out of budget. If  $b = 0$ , that prefix is empty and the

probability of failure is  $1 - \mu$ ; otherwise, the last block of that prefix is the  $b$ th block found by  $\mathcal{A}$ . Let  $h$  denote the number of blocks honestly mined during that period. Then:

$$\Pr(\mathfrak{f} \wedge \mathfrak{b} \wedge \mathfrak{o}) = \begin{cases} 1 - \mu & \text{if } b = 0, \\ (1 - \mu)(\alpha(1 - \mu))^b \sum_{h=0}^{\Delta} \binom{b+h-1}{b-1} (1 - \alpha)^h & \text{otherwise.} \end{cases} \quad (9)$$

Now, let us compute the probability of  $\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \mathfrak{o}$ . Look at the full sequence of found blocks,  $\Delta + 1$  of which are honestly mined, and  $a < b$  are found by the adversary. The probability can then be written as:

$$\Pr(\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \mathfrak{o}) = (1 - \mu)(1 - \alpha)^{\Delta+1} \sum_{a=0}^{b-1} \binom{\Delta+a}{a} (\alpha(1 - \mu))^a. \quad (10)$$

We now look at  $\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \neg \mathfrak{o}$ . Split the sequence of found blocks into two parts: the initial segment up to the pleasant block found by the adversary, and the rest of the sequence until the end. Let  $a_1$  be the number of blocks found by  $\mathcal{A}$  during the initial part, excluding the last one (i.e., only the unpleasant blocks), and  $h_1$  the number of blocks honestly mined in the meantime, including the unpleasant broadcast decisive block. Let  $a_2$  be the number of blocks found by  $\mathcal{A}$  during the second phase. Since we suppose  $\neg \mathfrak{b}$ , we have  $a_1 + 1 + a_2 < b$ . The number of blocks honestly mined in the second phase is  $\Delta + 1 - h_1$ , including the final block. Observe that since  $\mathcal{A}$  loses, the second part is a generalised Dyck word. Then, the probability can be computed as:

$$\Pr(\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \neg \mathfrak{o}) = (1 - \mu)\alpha\mu(1 - \alpha)^{\Delta+1} \sum_{\substack{0 \leq a_1 < b-1 \\ 1 \leq h_1 \leq \Delta \\ 0 \leq a_2 < b-a_1}} \binom{a_1 + h_1}{a_1} (\alpha(1 - \mu))^{a_1} D(\Delta, a_2, \Delta - h_1) \alpha^{a_1}. \quad (11)$$

The last event  $\mathfrak{f} \wedge \mathfrak{b} \wedge \neg \mathfrak{o}$  needs to be further split into two events: the case  $\mathfrak{u}$  where the pleasant decisive block found by  $\mathcal{A}$  is the last block in their budget, and the case  $\mathfrak{v}$  where  $\mathcal{A}$  still has some budget after that pleasant block was found. Then,  $\Pr(\mathfrak{f} \wedge \mathfrak{b} \wedge \neg \mathfrak{o}) = \Pr(\mathfrak{u}) + \Pr(\mathfrak{v})$ . We first focus on  $\mathfrak{u}$ . Look at the prefix of the sequence of blocks up to the pleasant decisive block. That prefix contains exactly  $b$  blocks found by  $\mathcal{A}$  (including the last one), and a number  $h \leq \Delta$  of blocks honestly mined (including the unpleasant decisive block). The probability is then:

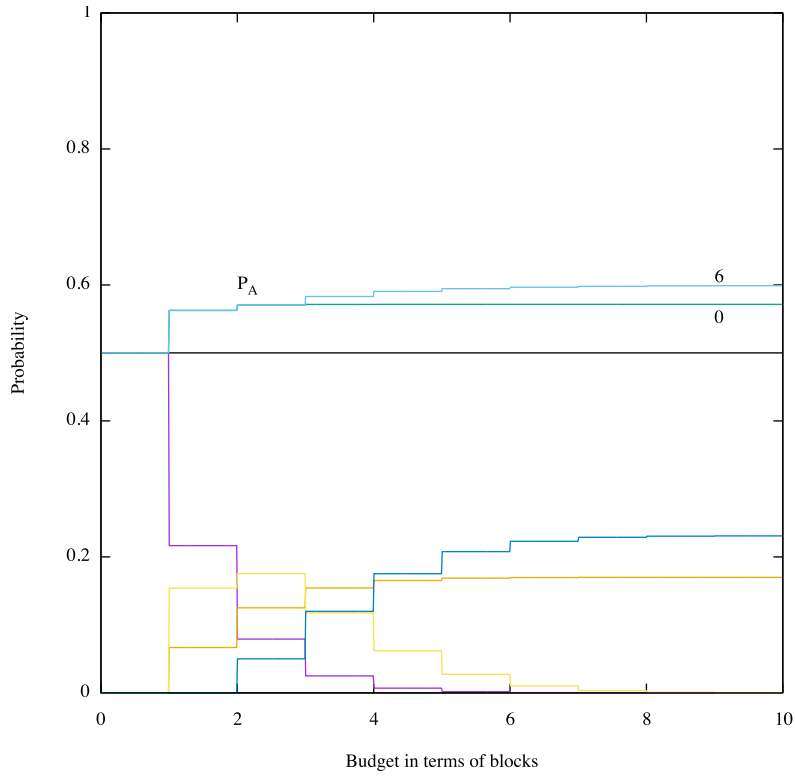
$$\Pr(\mathfrak{u}) = \alpha^b \mu (1 - \mu)^b \sum_{h=1}^{\Delta} \binom{b+h-1}{h} (1 - \alpha)^h. \quad (12)$$

It only remains to compute  $\Pr(\mathfrak{v})$ . Again, focus on the prefix of the sequence of found blocks which stops when the adversary runs out of budget, i.e. when they mine their  $b$ th block. Split that prefix into two parts: the initial part up to the pleasant decisive block found by the adversary, and the second part, from that pleasant block to the end, i.e., the  $b$ th block of  $\mathcal{A}$ . Let  $a_1$  be the number of blocks found by the adversary during the first part, excluding the last one (i.e., only the unpleasant blocks),  $h_1$  the number of blocks honestly mined in the meantime, including again the unpleasant decisive block, and  $h_2$  the number of blocks honestly mined in the second part. The number of blocks found by  $\mathcal{A}$  in that part is  $b - a_1 - 1$ . The probability is then:

$$\Pr(\mathfrak{v}) = (1 - \mu)\alpha^b \mu \sum_{\substack{0 \leq a_1 < b-1 \\ 1 \leq h_1 \leq \Delta \\ 0 \leq h_2 \leq \Delta - h_1}} \binom{a_1 + h_1}{a_1} (1 - \mu)^{a_1} (1 - \alpha)^{h_1 + h_2} D(h_1 + h_2, b - a_1 - 2, h_2). \quad (13)$$

Finally, the winning probability is computed by subtracting all the possible failure probabilities:

$$\mathcal{P}_{\mathcal{A}} = 1 - \Pr(\mathfrak{f} \wedge \mathfrak{b} \wedge \mathfrak{o}) - \Pr(\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \mathfrak{o}) - \Pr(\mathfrak{f} \wedge \neg \mathfrak{b} \wedge \neg \mathfrak{o}) - \Pr(\mathfrak{u}) - \Pr(\mathfrak{v}),$$



**Fig. 6.** Probabilities as functions of the budget (number of blocks that can be dropped), when  $\alpha = 1/4$  and  $\mu = 1/2$ . The two highest curves correspond to  $\mathcal{P}_A$  in the delay-free case ( $\Delta = 0$ ) and when  $\Delta = 6$ . The lowest curves are the four components that permit to compute  $\mathcal{P}_A$  when  $\Delta = 6$ : the purple curve corresponds to  $\Pr(f \wedge \mathbf{b} \wedge \mathbf{o})$ , the orange one to  $\Pr(f \wedge \neg \mathbf{b} \wedge \mathbf{o})$ , the blue one to  $\Pr(f \wedge \neg \mathbf{b} \wedge \neg \mathbf{o})$ , and the yellow one to  $\Pr(f \wedge \mathbf{b} \wedge \neg \mathbf{o})$ .



where  $\Pr(f \wedge \mathbf{b} \wedge \mathfrak{o})$ ,  $\Pr(f \wedge \neg \mathbf{b} \wedge \mathfrak{o})$ ,  $\Pr(f \wedge \neg \mathbf{b} \wedge \neg \mathfrak{o})$ ,  $\Pr(\mathbf{u})$  and  $\Pr(\mathbf{v})$  are respectively given in Equations (9), (10), (11), (12) and (13). Figure 6 details how  $\mathcal{P}_{\mathcal{A}}$  varies as a function of the number of blocks that can be dropped in the two cases  $\Delta = 0$  (no delay) and  $\Delta = 6$ . Note that these probabilities tend to 0.57 (resp. 0.6) which were the probabilities previously found in Section 3.3 (resp. Section 4.4), when the budget was unlimited. Moreover, we can see that the winning probability converges fast enough so that even a small budget is sufficient to significantly cheat the game.

## 6 Provident adversaries

### 6.1 Games with early start and delay

Let us consider the more general case where not only there is a final delay  $\Delta \geq 1$  as previously, but the game also starts early, i.e. the decisive block is not the next one, but the  $n$ -th one, where  $n > 1$ . This corresponds to the real world situation where an attacker knows in advance which future block will be decisive for the goal he tries to achieve. This case is relevant when a far-sighted adversary may get ahead of the game to manipulate an upcoming decisive block. In truth, nothing prevents adversaries that already cheat from preparing themselves several blocks in advance, without waiting for the beginning of the game to start mining in secret.

The analysis of this situation is essential to applications where some parties commit on an upcoming random number. In a lottery, the ticket seller should announce in advance which (future) block will determine the winner, giving the adversary a head start to manipulate that precise block. It is also crucial in electronic voting systems, where unpredictable, verifiable, public randomness is required for the security of the random auditing after the election [CCC<sup>+</sup>08].

### 6.2 Strategy

For the adversary, the goal will be to obtain a chain of length  $n - 1$  as soon as possible. If they can obtain such a branch before the honest miners – and keeps it secret –, they can start mining for a pleasant  $n$ -th block (the decisive one) while the honest miners are still working on building the initial segment of the chain. The adversary tries to build a chain longer than the honest one as follows: on one hand there is the public, honestly mined chain, and on the other hand there is the secret branch of the adversary. The adversary only mined on their own secret branch. At any point, if the honest chain becomes longer than the adversary’s branch, the latter is thrown away and replaced by the longer honest chain. With this strategy, the adversary’s branch is never shorter than the public one.

When an  $n - 1$ -th block has finally been found (either by  $\mathcal{A}$ , or by the honest miners if  $\mathcal{A}$  has not been fast enough) the adversary starts following the same strategy as in the previous game – see Section 4.3. Again, they focus all their computational power on mining the decisive block and throw away all valid but unpleasant blocks. At the first broadcast of the decisive block  $B_d$  two cases can occur. If it satisfies  $\chi(\mathbf{Ext}(B_d)) = 1$  then  $\mathcal{A}$  starts to act as a honest miner. Otherwise, if one honest miner broadcasts a decisive block satisfying  $\chi(\mathbf{Ext}(B_d)) = 0$ ,  $\mathcal{A}$  continues to secretly mine on the previous block as if they received nothing. Again, they throw away all valid decisive blocks that would lead to their failure. If they find a pleasant one, they keep it secret, create a fork and continue to mine on their new chain, hoping to develop a branch longer than the current public one containing the unpleasant decisive block. If they manage to do so, they broadcast their fork, and all miners switch to that new longest chain and proceed to mine on that one, thereby including the cheated pleasant decisive block in the new consensus chain. The game ends when a miner finds a valid final block, whether it is on the original chain, or on the branch later appeared.

### 6.3 Computing the probability $\mathcal{P}_{\mathcal{A}}$

To evaluate the probability that  $\mathcal{A}$  wins, we consider again the event  $\mathbf{e}_0$  that the first decisive broadcast block is an unpleasant one. As in Equation (5),  $\mathcal{P}_{\mathcal{A}}$  is given by:

$$\mathcal{P}_{\mathcal{A}} = 1 - \Pr(\mathbf{e}_0) \cdot (\Pr_{\mathbf{e}_0}(f_0) + \Pr_{\mathbf{e}_0}(f_{01})), \quad (14)$$

where  $\Pr_{\mathbf{e}_0}(f_0)$  and  $\Pr_{\mathbf{e}_0}(f_{01})$  are functions given in Equations (3) and (4). Yet in this far-sighted adversary case, the analysis differs when computing the probability of  $\mathbf{e}_0$ .

Let us consider the family of disjoint events  $\mathbf{a}_k$  that  $\mathcal{A}$  finds an  $n - 1$ -th block with an advantage of  $k$  blocks over the public branch, where  $k$  varies from 0 (meaning that the  $n - 1$ -th block is mined by a honest miner) to  $n - 1$  (when  $\mathcal{A}$  has computed the entire chain from blocks 1 to  $n - 1$  while the honest miners did not find a single valid block). This leads to:

$$\Pr(\mathbf{e}_0) = \sum_{k=0}^{n-1} \Pr_{\mathbf{a}_k}(\mathbf{e}_0) \cdot \Pr(\mathbf{a}_k).$$

If  $a$  denotes the number of blocks dropped by  $\mathcal{A}$  while the honest miners broadcast all the blocks needed to complete their chain up to the unpleasant decisive block,  $\Pr_{\mathbf{a}_k}(\mathbf{e}_0)$  is given by the sum  $(1 - \alpha)^{k+1} (1 - \mu) \sum_{a \geq 0} \binom{a+k}{a} (\alpha(1 - \mu))^a$ . Using Equation (2) we find:

$$\Pr_{\mathbf{a}_k}(\mathbf{e}_0) = (1 - \mu) \cdot \left( \frac{1 - \alpha}{1 - \alpha(1 - \mu)} \right)^{k+1}. \quad (15)$$

Thus it suffices now to evaluate the probabilities  $\Pr(\mathbf{a}_k)$ .

**Probability  $\Pr(\mathbf{a}_k)$  for a fixed  $k$ .** Let  $k$  be an integer between 0 and  $n - 1$ . The different outcomes that lead the adversary to obtain an advantage of  $k$  blocks can be modelled as paths on an  $(n - 1) \times (n - 1)$  grid of blocks from the point  $(0, 0)$  to the point  $(n - 1, n - 1 - k)$  as in Figure 7, in Appendix A. Following certain constraints.

Each point  $(a, h)$  of a path corresponds to the moment at which  $\mathcal{A}$  has a valid branch from block 0 to block  $a$  whereas the honest chain only goes up to the block  $h$ . The adversary's strategy enforces  $h \leq a$ , i.e. the path never overpasses the diagonal line. The short arrows forming the path represent the events that a new valid block has been found. The arrow is black if the block was found by  $\mathcal{A}$  and light gray if it was found by a honest miner. There are only three kinds of arrows:

1. Black arrows to the right: when  $\mathcal{A}$  finds a block, it only increases their own secret branch.
2. Gray arrows that go up, without overpassing the diagonal: when a honest miner finds a new block but  $\mathcal{A}$ 's chain is longer, this block only increases the honest branch.
3. Gray arrows on the diagonal: when a honest miner finds a new block and  $\mathcal{A}$ 's branch coincides with the honest chain, this block increases both chains.

Also, the path never follows the blue line: once the blue line has been reached, the next phase of  $\mathcal{A}$ 's strategy begins. To determine the probability of  $\mathbf{a}_k$ , the strategy will consist in computing the number of such paths to  $(n - 1, n - 1 - k)$  for any specified number of black arrows.

Let  $C_{x,y,r}$  denote the number of paths from  $(0, 0)$  to  $(x, y)$  that follow all the previous requirements and go through  $r$  black arrows precisely.  $C_{x,y,r}$  can be computed thanks to the following recursive formulae:

$$C_{x,y,r} = \begin{cases} 0 & \text{if } y = 0, r \neq x \\ 1 & \text{if } y = 0, r = x \\ C_{x-1,y,r-1} & \text{if } x = n - 1, y \neq n - 1 \\ C_{x-1,y-1,r} & \text{if } x = n - 1 = y \\ C_{x,y-1,r} + C_{x-1,y-1,r} & \text{if } x = y, y \neq n - 1 \\ C_{x,y-1,r} + C_{x-1,y,r-1} & \text{otherwise.} \end{cases}$$

At the end, since the number  $r$  of black arrows to go to  $(n-1, n-1-k)$  varies between  $k$  and  $n-1$ , the probability of  $\mathbf{a}_k$  is given by:

$$\Pr(\mathbf{a}_k) = \sum_{r=k}^{n-1} (1-\alpha)^{n-1-k} \alpha^r C_{n-1, n-1-k, r}.$$

Together with Equation (15), we find:

$$\Pr(\epsilon_0) = \frac{(1-\mu)(1-\alpha)^n}{1-\alpha(1-\mu)} \sum_{k=0}^{n-1} \left( \frac{1}{1-\alpha(1-\mu)} \right)^k \sum_{r=k}^{n-1} \alpha^r C_{n-1, n-1-k, r}. \quad (16)$$

The probability  $\mathcal{P}_{\mathcal{A}}$  that  $\mathcal{A}$  wins is so computed thanks to Equation (14), using Equations (3), (4) and (16) as explicit formulae for the inner probabilities  $\Pr_{\epsilon_0}(f_0)$ ,  $\Pr_{\epsilon_0}(f_{01})$  and  $\Pr(\epsilon_0)$ . Figure 8, in Appendix A, shows how  $\mathcal{P}_{\mathcal{A}}$  varies with  $\alpha$ , for  $\mu = 1/2$  and a fixed value of  $n = 144$  that roughly corresponds to an adversary that starts mining with a day in advance. Note that with a relative computational power  $\alpha = 1/4$ , the probability that a provident adversary wins takes off from 50% to 74% in the realistic case when  $\Delta = 6$ . Figure 8 also suggests that being too provident is not necessary. Indeed,  $\Pr(\epsilon_0)$  is the only parameter that governs the variations of the probability  $\mathcal{P}_{\mathcal{A}}$  when the adversary mines more or less in advance. We note that this probability stabilises very quickly, as mining with half-a-day, a day or a week in advance is roughly the same in terms of success.

## Conclusion

We have shown that a random beacon producing random numbers by extracting entropy from the last block of a blockchain does not provide a strong level of security. Their “random” output can easily be biased by large groups of miners.

We have been studying a beacon that outputs the hash of the very last block of the chain. One could naturally wonder if a stronger beacon could be build by hashing the last  $n$  blocks  $B_{N-n+1}, \dots, B_N$  for  $n > 1$ , say  $H(B_{N-n+1}, \dots, B_N)$  for a secure hash function  $H$ . Unfortunately, any bit of such a beacon could be manipulated in the very same way, with the same resulting biases, by still considering the last block as the decisive one. Indeed, once the blocks  $B_{N-n+1}, \dots, B_{N-1}$  are determined, the function  $H'(x) = H(B_{N-n+1}, \dots, B_{N-1}, x)$  is a hash function with security properties as good as  $H$  itself (assuming  $H$  is cryptographically secure).

Other games could use the randomness of the last  $n$  blocks without first hashing them together. A lot of naturally occurring games seem susceptible to a form of “last block attack”, yet each game would require a case-by-case analysis. In particular, we do not claim that it is impossible to generate secure random numbers from the blockchain. On the contrary, it seems to be possible by incorporating ideas from [LW15] in a blockchain protocol (see [But16] for an informal discussion).

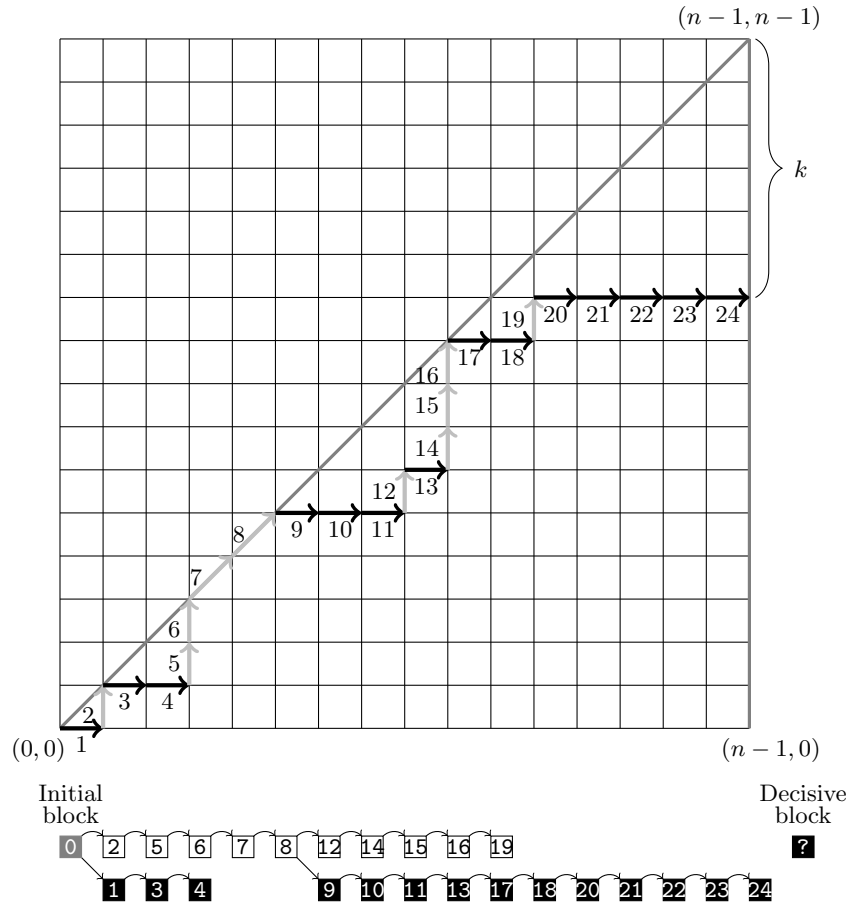
## Acknowledgements

We would like to thank Arjen Lenstra for his comments and Antoine Joux for both his technical and financial support. A special thanks goes to Direction Générale de l’Armement and CNRS for funding the first author, and to the Swiss National Science Foundation for supporting the second author via grant number 200021-156420.

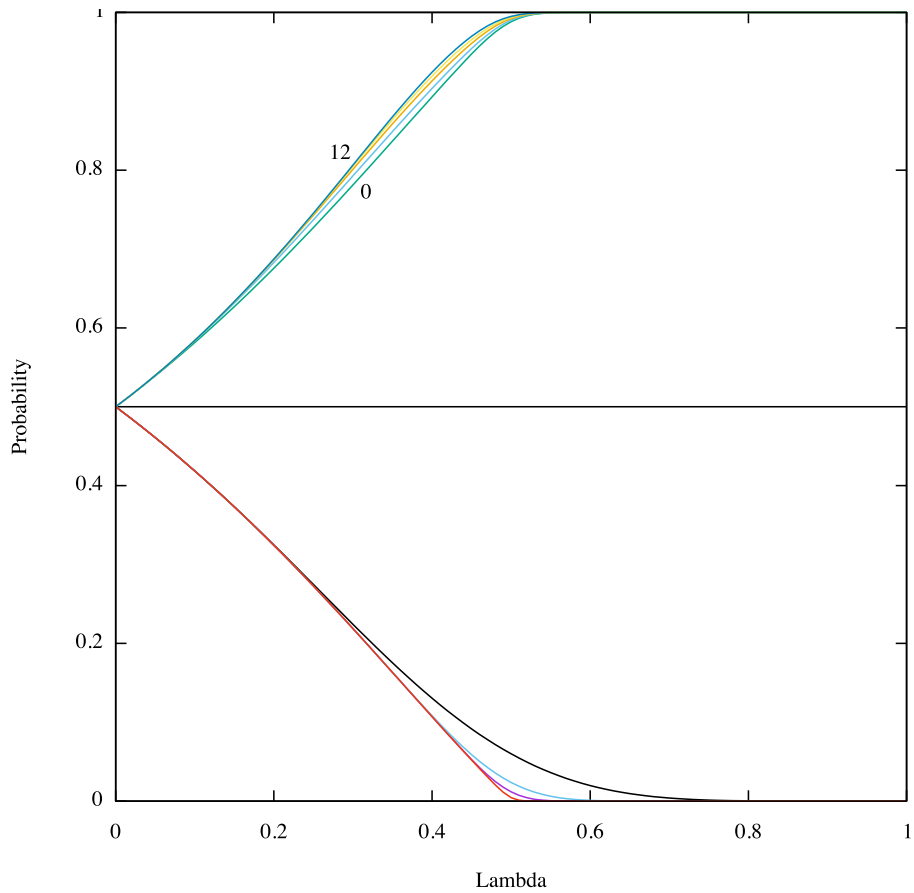
## References

- [BCG15] Joseph Bonneau, Jeremy Clark, and Steven Goldfeder. On bitcoin as a public randomness source. *IACR Cryptology ePrint Archive*, 2015:1015, 2015.
- [Blo16] Blockchain Luxembourg SARL. <https://blockchain.info>, 2016.
- [BNM<sup>+</sup>14] Joseph Bonneau, Arvind Narayanan, Andrew Miller, Jeremy Clark, Joshua A. Kroll, and Edward W. Felten. Mixcoin: Anonymity for bitcoin with accountable mixes. In *Financial Cryptography and Data Security: 18th International Conference, FC 2014*, pages 486–504, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg.
- [Bon16] Joseph Bonneau. Why buy when you can rent? - bribery attacks on bitcoin-style consensus. In *Financial Cryptography and Data Security - FC 2016 International Workshops, BITCOIN, VOTING, and WAHC, Christ Church, Barbados, February 26, 2016, Revised Selected Papers*, pages 19–26, 2016.
- [But16] Vitalik Buterin. Could Ethereum do this better? [Tor Project is working on a web-wide random number generator]. Reddit post (accessed on 2017-08-25), 2016.
- [CCC<sup>+</sup>08] David Chaum, Richard Carback, Jeremy Clark, Aleksander Essex, Stefan Popoveniuc, Ronald L. Rivest, Peter Y. A. Ryan, Emily Shen, and Alan T. Sherman. Scantegrity ii: End-to-end verifiability for optical scan election systems using invisible ink confirmation codes. In *USENIX/ACCURATE Electronic Voting Technology Workshop (EVT)*, 2008.
- [CH10] Jeremy Clark and Urs Hengartner. On the use of financial data as a random beacon. In *USENIX EVT/WOTE*. USENIX Association, 2010.
- [DGH<sup>+</sup>04] Yevgeniy Dodis, Rosario Gennaro, Johan Håstad, Hugo Krawczyk, and Tal Rabin. Randomness extraction and key derivation using the cbc, cascade and hmac modes. In *Advances in Cryptology - CRYPTO 2004: 24th Annual International Cryptology Conference*, pages 494–510, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
- [GGMR14] Christina Garman, Matthew Green, Ian Miers, and Aviel D. Rubin. Rational zero: Economic security for zerocoin with everlasting anonymity. In *Financial Cryptography and Data Security: BITCOIN*, pages 140–155. Springer Berlin Heidelberg, 2014.
- [GKL15] Juan A. Garay, Aggelos Kiayias, and Nikos Leonardos. The bitcoin backbone protocol: Analysis and applications. In *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Sofia, Bulgaria, 2015, Proceedings, Part II*, pages 281–310, 2015.
- [LW15] Arjen K. Lenstra and Benjamin Wesolowski. A random zoo: sloth, unicorn, and trx. Cryptology ePrint Archive, Report 2015/366, 2015. <http://eprint.iacr.org/2015/366>.
- [MGGR13] Ian Miers, Christina Garman, Matthew Green, and Aviel D. Rubin. Zerocoin: Anonymous distributed e-cash from bitcoin. In *Security and Privacy (SP), 2013 IEEE Symposium on*, pages 397–411, May 2013.
- [Nak09] Satoshi Nakamoto. Bitcoin: A peer-to-peer electronic cash system, <http://bitcoin.org/bitcoin.pdf>, 2009.
- [NIS11] NIST randomness beacon. <https://beacon.nist.gov>, 2011.
- [Rab83] Michael O. Rabin. Transaction protection by beacons. *Journal of Computer and System Sciences*, 27(2):256 – 267, 1983.
- [RAN98] RANDOM.ORG. <https://www.random.org>, 1998.
- [SSZ16] Ayelet Sapirshstein, Yonatan Sompolinsky, and Aviv Zohar. Optimal selfish mining strategies in bitcoin. In *Financial Cryptography and Data Security - 20th International Conference, FC 2016, Christ Church, Barbados, February 22-26, 2016, Revised Selected Papers*, pages 515–532, 2016.
- [Woo14] Gavin Wood. Ethereum: A secure decentralized transaction ledger. <http://gavwood.com/paper.pdf>, 2014.

## A Figures for the provident adversary



**Fig. 7.** Secret and public blocks mined during the event  $\alpha_k$ . White blocks (or light gray arrows) represent the blockchain integrally mined by honest miners whereas black blocks (or black arrows) are those secretly mined by  $\mathcal{A}$ . On this example,  $n = 17$  and  $\mathcal{A}$  finds the block number  $n - 1$  with a head start of  $k = 6$  blocks. The number inside each block (or near each arrow) indicates the order in which the corresponding block was mined, irrespective of whether it was broadcast or kept secret.



**Fig. 8.** Probability  $\mathcal{P}_A$  as a function of  $\alpha$  when  $\mu = 1/2$  and  $n = 144$  (in the case of Bitcoin, 144 blocks are found in an expected time of one day). The corresponding values of  $\Delta$  are 0, 1, 3, 6 and 12 as suggested on the curves. The lower curves represent the inner probability  $\Pr(\epsilon_0)$  as a function of  $\alpha$  for  $\mu = 1/2$  and various choices of  $n$ . The black one corresponds to  $\Pr(\epsilon_0)$  when  $n = 6$  (corresponding, in Bitcoin, to an expected time of one hour), the blue one when  $n = 36$  (six hours), the purple one when  $n = 144$  (one day) and the red one when  $n = 1008$  (one week).