

Designing Proof of Human-work Puzzles for Cryptocurrency and Beyond

Jeremiah Blocki¹ and Hong-Sheng Zhou²

¹ Purdue University, jblocki@purdue.edu

² Virginia Commonwealth University, hszhou@vcu.edu

Abstract. We introduce the novel notion of a Proof of Human-work (PoH) and present the first distributed consensus protocol from hard Artificial Intelligence problems. As the name suggests, a PoH is a proof that a *human* invested a moderate amount of effort to solve some challenge. A PoH puzzle should be moderately hard for a human to solve. However, a PoH puzzle must be hard for a computer to solve, including the computer that generated the puzzle, without sufficient assistance from a human. By contrast, CAPTCHAs are only difficult for other computers to solve — not for the computer that generated the puzzle. We also require that a PoH be publicly verifiable by a computer without any human assistance and without ever interacting with the agent who generated the proof of human-work. We show how to construct PoH puzzles from indistinguishability obfuscation and from CAPTCHAs. We motivate our ideas with two applications: HumanCoin and passwords. We use PoH puzzles to construct HumanCoin, the first cryptocurrency system with human miners. Second, we use proofs of human work to develop a password authentication scheme which provably protects users against offline attacks.

1 Introduction

The emergence of decentralized cryptocurrencies like Bitcoin [45] has the potential to significantly reshape the future of distributed interaction. These recent cryptocurrencies offer several advantages over traditional currencies, which rely on a centralized authority. At the heart of Bitcoin-like cryptocurrencies is an efficient distributed consensus protocol that allows for all users to agree on the same public ledger. When combined with other cryptographic tools like digital signatures the distributed consensus protocol prevents users from engaging in dishonest behavior like “double spending” their money or spending another user’s money. Fundamentally, the applications of a tamper-proof blockchain like the one in Bitcoin are not limited to cryptocurrency. For example, a tamper proof blockchain could help us construct secure and fair multiparty computation protocols [1,7,36,38], develop smart contracts [53,38], and build distributed autonomous agents, to name a few applications. In this paper we propose a fundamentally new technique, Proofs of Human-work (PoH), for constructing a secure blockchain, and we show that our techniques have several other valuable applications like password protection and non-interactive bot detection.

At its core, Bitcoin’s distributed consensus protocol is based on moderately hard Proofs of Work (PoW) [23]. In Bitcoin the Hashcash [3] PoW puzzles are used to extend the blockchain, a cryptographic data-structure in which the public ledger is recorded. A PoW puzzle should be moderately hard for a computer to solve, but the PoW solution should be easy for a computer to verify. Cryptocurrencies like Bitcoin require that this hardness parameter of PoW puzzles be tunable. An adversary would need to control 51% of the computational power in the Bitcoin network to be able to alter the blockchain and prevent users from reaching the correct consensus³. While Bitcoin cleverly avoids the Sybil attack by using PoW puzzles, there are still many undesirable features of this distributed consensus protocol. For example, constructing the proofs of work is energy intensive making the mining process in this distributed consensus protocol environmentally unfriendly. Furthermore, the mining process is dominated by a smaller number of professional miners with customized hardware making it unprofitable for others to join — this raises the natural concern that a few professional miners might collude to alter the public ledger [46]. Indeed, the mining pool GHash.io⁴ recently exceeded 50% of the computational power in Bitcoin. While other techniques like Proofs of Space [47,25] or Proofs of Stake [8] have been proposed to build the blockchain in a distributed consensus protocol each of these techniques has its own drawbacks. It is clearly desirable to find new techniques for reaching a stable distributed consensus. In this paper we ask the following question:

Is it possible to design proof of human-work puzzles that are suitable for a decentralized cryptocurrency?

We believe that a cryptocurrency based on Proof of Human-work might offer many advantages over other approaches. First, the mining process would be eco-friendly. Second, instead of wasting ‘human cycles,’ it might be possible to base the proofs of human work on activities that are fun [34], educational [33] or even beneficial to society [56,35]. Third, proofs of human work are fair by nature in the sense that two individuals will generally perform a comparable amount of work to produce a proof of human work. Thus, professional or rich miners would not have an significant advantage over regular users. By contrast, in Bitcoin the cost of computing the SHA256 hash function on customized hardware is dramatically less than the cost of computing SHA256 on personal computing⁵. Finally, we believe that the cryptocurrency would be less-vulnerable to 51% attacks by nation states or by a few professional miners. However, we stress that our purpose is not to enumerate all of the possible social consequences of

³ Technically byzantine agreement is only possible when the adversary has less than 50% of the hashing power and the network has high synchronicity — otherwise we need to ensure that the adversary has at most 33.3% of the hashing power [29].

⁴ See <http://arstechnica.com/security/2014/06/bitcoin-security-guarantee-shattered-by-anonymous-miner-with-51-network-power/>.

⁵ See <https://bitcoinmagazine.liberty.me/bitmain-announces-launch-of-next-generation-antminer-s7-bitcoin-miner/> (Retrieved 5/4/2016).

a cryptocurrency based on Proofs of Human-work. As with any new technology HumanCoin could potentially be used for good or for evil. See the full version [12] for additional discussions.

1.1 Cryptocurrencies meet AI: Proof of Human-work Puzzles

In this work we introduce the novel notion of Proofs of Human-work (PoH) which would be suitable for cryptocurrencies. Proofs of Human-work are fundamentally different from standard Proofs of Work. Informally, a PoH puzzle should be moderately hard for a human to solve meaning that it should require modest effort for a human to produce a valid proof of human work — again we require that this hardness parameter should be tunable. Furthermore, the puzzles should be easy for a computer to generate, but they need to be difficult for a computer to solve without sufficient human assistance — even for the computer that generated the puzzle. Finally, the puzzles need to be publicly verifiable meaning that it should be easy for a computer to verify the solution to the puzzle without any human assistance — even if the computer did not generate the puzzle. We stress that there is no interaction during the puzzle generation or during the puzzle verification process, and there is no trusted server in our distributed setting. Thus, a computer will need to validate proofs of human-work that were generated and solved by agents with whom it has never interacted.

Our description of a PoH puzzle might remind the reader of a CAPTCHA (Completely Automated Public Turing-Test to tell Computers and Humans Apart) [55]. CAPTCHAs have been widely deployed on the Internet to fight spam and protect against sybil attacks. Informally, a CAPTCHA is a puzzle that is easy for a human to solve, but difficult for a computer. CAPTCHAs are based on the assumption that some underlying artificial intelligence (AI) problem is hard for computers, but easy for humans (e.g., reading distorted letters).

While we do use CAPTCHAs to construct proofs of human work, we stress that a CAPTCHA itself *cannot* achieve our notion of proofs of human-work. Let (Z, σ) be a CAPTCHA puzzle-solution pair. Verifiers who receive the pair (Z, σ) would not necessarily be able to check that σ is the correct solution without interacting with a human. More importantly, the computer that generates the puzzle Z could produce the solution σ *without any human effort* because CAPTCHA generation algorithms start by randomly selecting a target solution σ and then outputting a randomly generated puzzle Z with the solution σ . Thus, a pair (Z, σ) does not constitute a proof of human work. The PoH verifier would need to ensure, without interacting with any other human agent or any other computer agent, that the challenge generator did not already have the answer σ to the puzzle Z .

We believe that our Proof of Human-work puzzles could also have applications in many other contexts. For example, to limit spam or prevent phishing attacks it might useful to verify that some human effort went into producing a message. When a human user is busy it would be convenient if the computer could validate this proof of human effort automatically without needing to interact with the sender who may no longer be available when the message

is received. Similarly, proofs of human-work might be a useful tool for honest preference elicitation — a challenging problem in mechanism design. A human could demonstrate that a particular issue or outcome is truly important to him by producing a proof of human-work.

1.2 AI meets Obfuscation: Constructing Proof of Human-work Puzzles

It is not immediately clear how to construct PoH puzzles. CAPTCHAs allow a computer to generate puzzles that other computers cannot solve, but how could a computer generate a puzzle that is meaningful to a human without learning the answer itself? Even if this were possible how could a puzzle verifier be convinced that the puzzle(s) was generated honestly (e.g., in a way that does not reveal the answer) without any interaction? How could the verifier be convinced that the answer is correct without help from a human? Building PoH puzzles is a challenging problem.

To address these issues, we need to have a way to generate CAPTCHAs *obliviously* in the sense that a computer is able to generate a well-formed puzzle instance Z without learning the corresponding solution σ . This is feasible by leveraging recent breakthroughs in indistinguishability obfuscation [30]. At an intuitive level, we can have a CAPTCHA puzzle Z generated inside an obfuscator, and now the corresponding answer σ remains hidden inside the obfuscated program. We note that the puzzle solution verification can also take place inside an obfuscated program, even without having human effort involved.

Once we have the idea of generating a CAPTCHA puzzle obliviously as mentioned above, we then can mimic the steps of constructing Proof of Work puzzle in Bitcoin to get a PoH scheme. In PoW, a prover/miner is given a puzzle instance x . The prover will compute the cryptographic hash $H(x, s)$ for many distinct witness s until the value $H(x, s)$ is smaller than a target value. In PoH, the miner uses (x, s) as the input for an obfuscated program, and inside the obfuscated program, a pseudorandom string r is generated from the input (x, s) , and this r will be used for generating the solution σ and the puzzle instance Z . The miner obtains Z but has no access to the internal state r and σ .

A human miner is now able to obtain the solution σ from the puzzle Z . As in PoW, the miner will repeat this process until he finds a witness s so that $H(x, s, \sigma, Z)$ is smaller than a target value. We note that, once a successful miner publishes a valid tuple (x, s, σ, Z) , any verifier is able to verify it without interaction with human: The verifier can reproduce Z inside the obfuscated program along with a verification tag, tag . While the verification tag allows the verifier to check whether a given solution σ is correct this value will not expose the solution σ (e.g., tag might be an obfuscated point function which outputs 1 on input $x = \sigma$ and 0 on all other inputs).

Our PoH scheme maintains many of the same desirable properties as a PoW. For example, we can tune the hardness of our PoH puzzle generator by having the verifier reject a valid triple (x, s, σ, Z) with probability $1 - 2^{-\omega}$ so that a human would need to generate and solve 2^ω on average to produce a valid proof of

human-work. Thus, the hardness of the PoH puzzles could be tuned by adjusting ω .

While the conceptual understanding of our PoH construction is quite simple, the security analysis is a bit tricky. In the PoW, we sample from a uniform distribution via random oracle, here we need to sample from a more sophisticated distribution. We rely on a newly developed tool *universal samplers* by Hofheinz et al. [32], which is based on the existence of indistinguishability obfuscation and one-way functions in the random oracle model. As discussed in [32], we stress that the random oracle is only used outside of obfuscated programs. There has been tidal wave of new cryptographic constructions using indistinguishability obfuscation since the groundbreaking results of Garg et al. [30]. However, to the best of our knowledge we are the first rigorous paper to explore the connection between AI and program obfuscation⁶. We believe that obfuscation is a powerful new tool that has the potential to fundamentally shape the nature of human-computer interaction. Could program obfuscation allow for a human to interact with a computer in fundamentally new ways? We view our work as a first step towards answering this question.

Remark 1. We view our Proof of Human-work construction as a novel proof of concept that is not yet practical due to the use of indistinguishability obfuscation. Since the work of Garg et al. [30] several other candidate indistinguishability obfuscation schemes have been proposed, but a practical obfuscation scheme would still be a major breakthrough. We note that PoH puzzles do not necessarily require general purpose indistinguishability obfuscation. It would be sufficient to obfuscate a few very simple programs (e.g., a CAPTCHA puzzle generator and a pseudorandom function). Constructing PoH puzzles without obfuscation (or without general purpose obfuscation) is an interesting open problem.

Other Applications. The applications of our techniques are not limited to cryptocurrency. In Section 5 we use our ideas to build a password authentication scheme that provably resists offline attacks even if the adversary breaches the authentication server. The basic idea is to require a proof of human-work during the authentication process so that it is not economically feasible for the adversary to check millions of password guesses. We also show how to develop a non-interactive bot detection protocol which allows Alice to send a message m to Bob along with a proof of human-work. Bob is able to verify that human-effort was used in the production/transmission of the message m without ever interacting with Alice.

⁶ Several existing altcoins (e.g., Bytecent, CaptchaCoin) do involve CAPTCHAs, but they rely on a trusted third party to generate the CAPTCHAs. There has also been informal discussion on the Bitcoin research chatroom about using obfuscation to base cryptocurrency on proofs of human labor. For example, see <https://download.wpsoftware.net/bitcoin/wizards/2014-05-29.html> or <http://vitalik.ca/files/problems.pdf>.

1.3 Related Work

While there are many variations of CAPTCHAs [55], they are all based on the fundamental assumption that some underlying AI problem is hard (e.g., reading garbled text [56], voice recognition with distorted audio [52], image recognition [26] or even motion recognition). While several CAPTCHAs have been broken (e.g., [43,54,16]) there is still a clear gap between human intelligence and artificial intelligence. We conjecture that in the foreseeable future we will continue to have viable CAPTCHA candidates suitable for proofs of human work. CAPTCHAs have many applications in security: fighting spam [55], mitigating Sybil attacks [20], preventing denial of service attacks [57] and even preventing fully automated man-in-the-middle attackers [24]. As we noted earlier CAPTCHAs alone are not suitable as PoH puzzles. Kumarasubramanian et al. [39] introduced the notion of human-extractable CAPTCHAs, and used them to construct concurrent non-malleable zero-knowledge protocols.

Canneti et al. [18] proposed a slight modification of notion of CAPTCHAs that they called HOSPs (Human Only Solvable Puzzles) as a defense against offline attacks on passwords. HOSPs are similar to PoHs in that the puzzles must be difficult even for the computer that generates them, but HOSP puzzles are not publicly verifiable by a computer and their construction assumes the existence of a large centralized storage server filled with unsolved CAPTCHA challenges. This makes their protocol vulnerable to pre-computation attacks⁷. By contrast, in Section 5 we present a protocol for password storage that provably protects users against offline attacks, does not require a large centralized storage server and is not vulnerable to pre-computation attacks. Blocki et al. [9] introduced GOTCHAs (Generating panOptic Turing Tests to Tell Computers and Humans Apart) as a defense against offline dictionary attacks on passwords. However, GOTCHAs have a high usability cost and are not suitable for cryptocurrency because the puzzle generation protocol requires interaction with a human and the solutions are not publicly verifiable by a computer. We refer an interested reader to the full version [12] for more details about CAPTCHAs and HOSPs.

The problem of designing distributed consensus protocols that work in the presence of an adversarial (Byzantine) parties has been around for decades [40,22,2]. Typically distributed consensus requires that $2/3$ of the parties are honest [40]. On the Internet this assumption is typically not valid because it is often possible for a malicious user to register for multiple fake accounts — a Sybil attack [20]. However, amazing ideas have been proposed in the original Bitcoin white paper [45] under a pseudo identity ‘Nakamoto’. At its core Bitcoin is based on an elegant distributed consensus protocol which in turn is based on Proof of Work puzzles [23] to allow users to agree on a common blockchain. Bitcoin uses the Hashcash Proof of Work algorithm due to Back [3]. Very recently, the underlying consensus protocol in the Bitcoin system have been rigorously analyzed in

⁷ In particular, the adversary might pay to solve every CAPTCHA challenge on the server. While expensive, this one-time cost would amortize over the number of users being attacked.

the cryptographic setting [29,48]; intensive analysis has also been given in the rational setting (e.g. [27,51]).

Since the breakthrough result of Garg et al. [30], demonstrating the first candidate of indistinguishability obfuscation for all circuits, a myriad of uses for indistinguishability obfuscation in cryptography have been found. Among these results, the puncturing methodology by Sahai and Waters [50] has been found very useful. Hofheinz et al explored the puncturing technique further introducing and constructing universal samplers in the random oracle model [32]. Their universal sampler is one of the key building blocks in our construction of proof of human-work puzzles. We remark that our work is distinct from previous applications in that we are using obfuscation to develop a new way for *humans* to interact with computers.

2 Preliminaries

We adopt the following notational conventions: Given a randomized algorithm \mathcal{A} we use $y \leftarrow \mathcal{A}(x)$ to denote a random sample from the distribution induced by an input x . If we fix the random bits r then we will use $y := \mathcal{A}(x; r)$ to denote the deterministic result.

We will consider two types of users: machine-only users and human-machine users. A machine-only user is a probabilistic polynomial time (PPT) algorithm who does not interact with a human. In general, when we say “human” user we mean a “human user equipped with a PPT machine.” Accordingly, we also consider two types of adversaries: a machine-only adversary \mathcal{A} , and a human-machine adversary $\mathcal{B}^{\mathcal{H}}$. The machine-only adversary is a PPT algorithm that does not get to query a human. The human-machine adversary $\mathcal{B}^{\mathcal{H}}$ is a PPT algorithm that gets to interact with a human oracle \mathcal{H} which could, for example, solve CAPTCHA puzzles. We typically restrict the total number of queries that human-machine adversary can make to the human oracle. We say that a human-machine adversary $\mathcal{B}^{\mathcal{H}}$ has m human-work units if it is allowed to query \mathcal{H} at most m times. We intentionally under-specify the behavior of the human oracle \mathcal{H} . At minimum we assume \mathcal{H} is capable of solving a CAPTCHA puzzle for one human-work unit (one query to the oracle). However, the human-machine adversary may use his queries to ask the human oracle to perform arbitrary tasks \mathcal{H} (e.g., solve basic arithmetic problems, write poetry) so long as each task takes (approximately) the same amount of human-effort as a single CAPTCHA puzzle.

2.1 CAPTCHAs

CAPTCHAs are a fundamental building block in our construction of Proof of Human-work puzzles. Traditionally, a CAPTCHA generator \mathbf{G} is defined as a randomized PPT algorithm that outputs a puzzle Z and a solution σ . In every CAPTCHA generator that we are aware of the program \mathbf{G} first generates a random target solution σ and then produces a random puzzle Z with solution

σ (e.g., by distorting the string σ). Given public parameters PP for CAPTCHA puzzle generation we adopt the syntax $(Z, \text{tag}) \leftarrow \mathbf{G}(\text{PP}, \sigma)$ to emphasize that the target puzzle Z is generated with complete knowledge of the CAPTCHA solution. In traditional CAPTCHA applications it is desirable for the agent who generates a puzzle Z to have knowledge of the corresponding answer σ so that he can verify another agent’s response to the challenge Z . However, in our setting this property is problematic since the agent who generates the puzzle Z is trying to produce a convincing proof of human-work. Thus, we will need additional tools to obtain proof of human-work puzzles from CAPTCHAs. Formally, a CAPTCHA puzzle-system is defined as follows.

Definition 1 (CAPTCHA). *A CAPTCHA puzzle system consists of a tuple of algorithms $(\text{Setup}, \mathbf{W}, \mathbf{G}, \mathbf{C}^{\mathcal{H}}, \text{Verify})$, where*

- **Setup** is a randomized system setup algorithm that takes as input 1^λ (λ is the security parameter), and outputs a system public parameter $\text{PP} \leftarrow \text{Setup}(1^\lambda)$, which includes a puzzle size parameter $\ell = \text{poly}(\lambda)$;
- **W** is a randomized sampling algorithm that takes as input the public parameter PP and outputs a target solution $\sigma \leftarrow \mathbf{W}(\text{PP})$ (e.g., a witness) of length ℓ ;
- **G** is a randomized puzzle generation algorithm that takes as input the public parameter PP and a solution σ , and outputs $(Z, \text{tag}) \leftarrow \mathbf{G}(\text{PP}, \sigma)$ where Z is a CAPTCHA puzzle and tag is a string that may be used to help verify a solution to Z ;
- **Verify** is a verification algorithm that takes as input the public parameters PP , a puzzle Z along with the associated tag and a proposed solution σ' outputs a bit $b := \text{Verify}(\text{PP}, Z, \text{tag}, \sigma')$. We require that $b = 1$ whenever $(Z, \text{tag}) \leftarrow \mathbf{G}(\text{PP}, \sigma)$ and $\sigma' = \sigma$;
- $\mathbf{C}^{\mathcal{H}}$ is a solution finding algorithm (i.e., human-machine solver) that takes as input the public parameter PP and a puzzle Z , and outputs a value $a \leftarrow \mathbf{C}^{\mathcal{H}(\cdot)}(\text{PP}, Z)$ as the solution to the puzzle Z . Here, $\mathcal{H}(\cdot)$ denotes the human oracle which takes intermediate human-efficient objects (such as images) as inputs, and returns machine-efficient values as outputs.

We typically require that **Setup**, **W**, **G** are probabilistic polynomial time algorithms, and **Verify** a deterministic polynomial time algorithm. \mathbf{C} should be a probabilistic polynomial time oracle machine.

For example, if we are defining a text based CAPTCHA puzzle-system the public parameters PP might specify the set of characters Σ , the set of fonts and a set of font sizes/colors. The public parameters PP would also describe the length $\ell = |\sigma|$ of the target solution (e.g., the number of characters in the CAPTCHA). In general, larger security parameters λ would imply longer puzzles. **W** is a randomized algorithm that outputs a random string $\sigma \in \Sigma^*$ (the target solution), and **G** is the randomized algorithm that produces a puzzle Z along with a tag which may be used for public verification of a potential solution σ' . We view the solution function $\mathbf{C}^{\mathcal{H}}$ as a human equipped with a PPT computer. Typically the computer would just be used to display the challenge to

the user, but it could also apply a more sophisticated algorithm to post-process the user’s answer.

Fixing the security parameter λ we define one human work unit to be the amount of time/energy that it takes a human to solve one honestly generated CAPTCHA puzzle $Z \leftarrow \mathbf{G}(\text{PP}, \sigma)$. Any CAPTCHA puzzle-system should be human usable, meaning that a typical human can consistently solve randomly generated CAPTCHA puzzles. While we recognize that solving a CAPTCHA puzzle may require more effort for some people than for others we will use the term human-work unit to denote the amount of human effort necessary to solve one CAPTCHA puzzle with security parameter λ .⁸

Definition 2 (Honest Human Solvability). *We say that a human-machine solver $\mathcal{C}^{\mathcal{H}}$ controls m human-work units if the machine \mathcal{C} can query the human oracle $\mathcal{H}(\cdot)$ at least m times. We say a CAPTCHA puzzle-system $(\text{Setup}, \mathbf{W}, \mathbf{G}, \mathcal{C}^{\mathcal{H}}, \text{Verify})$ is honest human solvable if for every polynomial $m = m(\lambda)$ and for any human $\mathcal{C}^{\mathcal{H}}$ who controls m human-work units, it holds that*

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda); \forall i \in [m] (\sigma_i^* \leftarrow \mathbf{W}(\text{PP})); \\ \forall i \in [m] ((Z_i^*, \text{tag}_i^*) \leftarrow \mathbf{G}(\text{PP}, \sigma_i^*)) \quad : \\ (\sigma_1^*, \dots, \sigma_m^*) \leftarrow \mathcal{C}^{\mathcal{H}(\cdot)}(\text{PP}, Z_1^*, \dots, Z_m^*) \end{array} \right] \geq 1 - \text{negl}(\lambda)$$

Finally, we require that CAPTCHAs are hard for computers to invert. More concretely, no (known) PPT adversarial machine should be able to find the solutions to $m + 1$ honestly-generated puzzles given only m -human work units. We introduce two similar notions of computer uncrackable CAPTCHAs. The first version states that an adversary with m human-work units cannot find the solution to $m + 1$ CAPTCHAs with non-negligible probability when he is only given the puzzles Z_1^*, \dots, Z_n^* ($n > m$).

Philosophical Remark. There are two philosophical positions that one could take regarding CAPTCHA puzzles, the human oracle \mathcal{H} and Artificial Intelligence in general. The first view is that for any class of problems that a human oracle \mathcal{H} can solve there exists a (possibly unknown to mankind) PPT computer algorithm to solve the same class of problems. The second philosophical view is that there are some tasks that humans can solve that computers will never be able to solve (i.e., no PPT computer algorithm can consistently/accurately solve the task).

We will implicitly follow view 1 in our CAPTCHA security definitions. However, we do not advocate for either view and we stress that our construction would also work under view 2. Under this second view the class of PPT machine-human hybrid adversaries is *strictly more powerful* than the class of PPT adversaries. Thus, one would need to make the assumption that the cryptographic

⁸ In the same way some computers (ASICs) are much faster at evaluating the SHA256 hash function than others. However, we expect this difference to be less extreme for human users.

primitives used in our construction (e.g., $i\mathcal{O}$, OWF) are secure against machine-human hybrids. This assumption is highly plausible⁹, but also non-standard.

Following view 1 we can avoid such non-standard assumptions about cryptographic primitives. In particular, we assume that the behavior of the human oracle \mathcal{H} is fully described by some (unknown) PPT algorithm. We note that because there exists a PPT algorithm specifying the behavior of \mathcal{H} the class $\text{PPT}^{\mathcal{H}}$ (the class of PPT algorithms with oracle access to \mathcal{H}) is no more powerful than the class PPT. Thus, we do not need to rely on non-standard cryptographic assumptions (e.g., $i\mathcal{O}$ is secure against adversaries in $\text{PPT}^{\mathcal{H}}$). How can a CAPTCHA scheme be secure if there exists some PPT algorithm that accurately solves challenges without human assistance? We will use the set Discoverable to denote a subset containing all known turing machines and all turing machines that mankind might plausibly discover in the near future (e.g., 10–20 years). More specifically, $\text{Discoverable}_X = \{M \mid M \text{ is a turing machine that mankind will build within the next } X \text{ years}\}$. The security of a CAPTCHA scheme relies on the assumption that no PPT algorithm in Discoverable_X will be able to accurately solve CAPTCHA puzzles for some reasonably large value of X (e.g., 10–20 years).

Stating that no PPT algorithm $A \in \text{Discoverable}_X$ breaks CAPTCHAs is a statement about human ignorance. While the meaning of this statement is clear at an intuitive level it is vague in a formal mathematical sense. As Rogaway observed the same issue arises in the definition of (keyless) collision resistant hash functions [49]. There *is* an efficient algorithm to find collisions, but it is not known to mankind and the hope is that no such algorithm will be known to mankind for a long time in the future. We will follow the same approach taken by Rogaway [49] when making security statements about constructions (e.g., PoH) that rely on CAPTCHAs. For example, we prove that there is an explicit PPT blackbox reduction (blackbox-constructive form [49]) transforming an adversary who breaks Proof of Human-work security to an adversary who breaks CAPTCHAs.

Definition 3 (CAPTCHA Break v1). *We say that a PPT adversary \mathcal{A} who has at most m human-work units breaks security of a CAPTCHA puzzle-system $(\text{Setup}, \mathbf{W}, \mathbf{G}, \mathcal{H}, \text{Verify})$ if for some polynomials $m = m(\lambda)$, $n = \text{poly}(\lambda)$ and $\mu(\lambda)$ when \mathcal{A} controls at most m human-work units, it holds that*

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda); \forall i \in [n] (\sigma_i^* \leftarrow \mathbf{W}(\text{PP})); \\ \forall i \in [n] ((Z_i^*, \text{tag}_i^*) \leftarrow \mathbf{G}(\text{PP}, \sigma_i^*)); \\ S \leftarrow \mathcal{H}^{(\cdot)}(\text{PP}, Z_1^*, \dots, Z_n^*); \\ \forall i \in [n] (b_i \leftarrow \max_{\sigma \in S} \text{Verify}(\text{PP}, Z_i^*, \text{tag}_i^*, \sigma)) : \\ \sum_{i \in [n]} b_i \geq m + 1 \end{array} \right] \geq \frac{1}{\mu(\lambda)}$$

⁹ If a cryptographic primitives like $i\mathcal{O}$ or one-way functions were not secure against machine-human hybrids then these primitives would have to be considered broken in practice.

We say that the CAPTCHA puzzle-system is *computer uncrackable* for the next X years if for any PPT adversary $\mathcal{A} \in \text{Discoverable}_X$, \mathcal{A} does not break security of the CAPTCHA puzzle system.

Our second formulation of CAPTCHA security is slightly non-standard due to the fact that the adversary is given a tag tag_i along with each challenge Z_i . In particular, the value tag_i allows the adversary to run $\text{Verify}(\text{PP}, Z_i, tag_i, \sigma'_i)$ to test different candidate CAPTCHA solutions. While this formulation is non-standard we argue that we would expect that any CAPTCHA that is secure under definition 3 can be transformed into a CAPTCHA that is secure under definition 4. For example, tag_i might be the cryptographic hash of the solution σ_i or we might set $tag_i = \text{iO}(I_{Z_i, \sigma_i})$ to be the indistinguishability obfuscation of a point function $I_{Z_i, \sigma_i}(x) = 1$ if $x = (Z_i, \sigma_i)$; otherwise $I_{Z_i, \sigma_i}(x) = 0$ ¹⁰.

It is reasonable to believe that \mathbf{G} could produce a tag tag_i , which allows us verify whether or not a solution σ' is correct without revealing σ_i . For example, we might set $tag_i = \text{iO}(I_{Z_i, \sigma_i})$ to be the indistinguishability obfuscation of a point function $I_{Z_i, \sigma_i}(x) = 1$ if $x = (Z_i, \sigma_i)$; otherwise $I_{Z_i, \sigma_i}(x) = 0$. In this case $\text{Verify}(\text{PP}, Z_i, tag_i, \sigma')$ would simply output $tag_i(Z_i, \sigma')$.

Definition 4 (CAPTCHA Break v2). *We say that a PPT adversary \mathcal{A} unites breaks security of a CAPTCHA puzzle-system $(\text{Setup}, \mathbf{W}, \mathbf{G}, \mathbf{C}^{\mathcal{H}}, \text{Verify})$ if for some polynomials $m = m(\lambda)$, $n = \text{poly}(\lambda)$ and $\mu(\lambda)$ when \mathcal{A} controls at most m human-work units, it holds that*

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda); \forall i \in [n] (\sigma_i^* \leftarrow \mathbf{W}(\text{PP})); \\ \forall i \in [n] ((Z_i^*, tag_i^*) \leftarrow \mathbf{G}(\text{PP}, \sigma_i^*)); \\ S \leftarrow \mathcal{A}^{\mathcal{H}(\cdot)}(\text{PP}, (Z_1^*, tag_1^*), \dots, (Z_n^*, tag_n^*)); \\ \forall i \in [n] (b_i \leftarrow \max_{\sigma \in S} \text{Verify}(\text{PP}, Z_i^*, tag_i^*, \sigma)) \quad : \\ \sum_{i \in [n]} b_i \geq m + 1 \end{array} \right] \geq \frac{1}{\mu(\lambda)}$$

We say that the CAPTCHA puzzle-system is *computer uncrackable* for the next X years if for any PPT adversary $\mathcal{A} \in \text{Discoverable}_X$ \mathcal{A} does not break security of the CAPTCHA puzzle-system.

We will require λ to be large enough that a computer cannot reasonably find a solution by brute force. As Von Ahn et al. [55] observed we can always increase λ by composing CAPTCHA puzzles. Of course this will increase the amount of time that it would take to solve a puzzle. Bursztein et al. [17] conducted a large scale experiment on Amazon's Mechanical Turk to evaluate human performance on a variety of different CAPTCHAs. Based on these results we estimate that, if we define one human work unit to be about two minutes of human effort, it is plausible to believe that security could be amplified to the extent that that

¹⁰ Indistinguishability obfuscation provides 'best case' obfuscation [31] so it would be highly surprising if an adversary could use tag_i to extract σ_i as this would immediately imply that a host of alternative cryptographic techniques (e.g., one way functions, collision resistance hash functions) fail to hide σ_i . A recent result of Barak et al [4] provides evidence that evasive circuit families (e.g., point functions) can be obfuscated.

adversary’s odds of solving the long CAPTCHA challenge correctly (and without human assistance) is negligible (e.g., 2^{-100})¹¹. For traditional CAPTCHA applications like bot detection this would make the solution impracticable due to the high usability costs. However, for our applications such a delay can be acceptable (e.g., in Bitcoin the parameters are tuned so that a new block is mined every 10 minutes).

While some spammers have paid human workers to solve CAPTCHAs in bulk [44] we do not consider this an attack on our definition because human effort was involved to find the solution. A HumanCoin miner could pay users to solve CAPTCHAs for him, but human users would have incentive to mine their own HumanCoins if compensation was unfair.

2.2 Universal Samplers

In [32], Hofheinz et al introduce the notion of universal samplers. The essential property of a universal sampler scheme is that given the sampler parameters U , and given any program d that generates samples from randomness, it should be possible for any party to use the sampler parameters U and the description of d to obtain induced samples that look like the samples that d would have generated given uniform and independent randomness.

Definition 5. *A universal sampler scheme consists of algorithms (**Setup**, **Sample**) where*

- $U \leftarrow \mathbf{Setup}(1^\lambda)$ is a randomized algorithm which takes as input a security parameter 1^λ and outputs sampler parameters U .
- $p_d \leftarrow \mathbf{Sample}(U, d)$ takes as input sampler parameters U and a circuit d of size at most $\ell = \text{poly}(\lambda)$, and outputs induced samples p_d .

In our construction in the next section, we will use a slightly extended version of universal sampler scheme which allows an additional input. Note that in the basic version of universal sampler scheme in Definition 5 above, the algorithm $\mathbf{Sample}(U, d)$ receives as input a program d which specifies certain distribution. In our application the program d will be fixed ahead of time, and \mathbf{Sample} takes an additional input β where β is an index for specifying randomness for the program to generate a CAPTCHA puzzle Z with tag. Thus, for the slightly extended version of universal sampler scheme with an additional input, we will use

¹¹ Some CAPTCHA candidates have already been “broken” by PPT algorithms (e.g. [54,43,19,16]). For example, [16] was able to solve reCAPTCHA with accuracy 33.34%. There are solid guidelines about generating CAPTCHAs that are harder for a computer to crack (e.g., see [58]). Furthermore, we stress that even apparently “broken” CAPTCHAs may still be useful in our proof of human work context because it is acceptable to use CAPTCHA puzzles that take a long time (e.g., 2 minutes) for a human to solve. By contrast, most deployed CAPTCHAs (e.g., reCAPTCHA) are meant to be solvable in a few seconds. As long as there is some gap between human intelligence and artificial intelligence we can use standard hardness amplification techniques (e.g., parallel repetition) to obtain stronger CAPTCHAs [55].

the notation $\mathbf{Sample}(U, d, \beta)$ instead of $\mathbf{Sample}(U, d)$. This allows us to provide alternative and flexible description for a circuit d without changing its functionality. We note that this slightly extended version has been explored in [32], and it is straightforward to extend \mathbf{Sample} , without requiring a new construction or security analysis.

The formal security definition of adaptive security for the slightly extended universal samplers with additional inputs can be found in Appendix A. We briefly overview the notion of adaptive security here. Intuitively, adaptive security guarantees that induced samples are indistinguishable from honestly generated samples to an arbitrary interactive system of adversarial and honest parties. In a universal sampler with additional inputs, the program d is fixed, and when an additional input β is provided, the induced sample can be computed as $p_\beta \leftarrow \mathbf{Sample}(U, d, \beta)$.

We first consider an “ideal world,” where a trusted party with a fixed program description d , on input β , simply outputs $d(r_\beta)$ where r_β is independently chosen true randomness, chosen once and for all for each given β . In other words, if F is a truly random function, then the trusted party outputs $d(F(\beta))$. In this way, if any party asks for samples corresponding to a specific value of β , they are all provided with the same honestly generated value.

In the real world, however, all parties would only have access to the trusted sampler parameters. Parties would use the sampler parameters to derive induced samples $d(r_\beta)$ for any specific inputs β . Now r_β is a pseudo random value corresponding to the randomness index β . We will require that for every real-world adversary \mathcal{A} , there exists a simulator \mathcal{S} that can provide simulated sampler parameters U to the adversary such that these simulated sampler parameters U actually induce the completely honestly generated samples $d(F(\beta))$ created by the trusted party: in other words, that $\mathbf{Sample}(U, d, \beta) = d(F(\beta))$. Note that since honest parties are instructed to simply compute induced samples, this ensures that honest parties in the ideal world would obtain these completely honestly generated samples $d(F(\beta))$.

3 Proof of Human-work Puzzles

In this section, we first define the syntax and security for proof of human-work puzzles; then we demonstrate a construction using universal samplers and CAPTCHAs.

3.1 Definitions

In a proof of work (PoW) puzzle, a party (i.e., prover) is allowed to prove to a bunch of verifiers that he completed some amount of computation/work. In general, those parties are machines. A typical PoW puzzle scheme consists of several algorithms: setup algorithm $\mathbf{Setup}()$ for generating the global system parameters and policies, puzzle instance generation algorithm $\mathbf{G}()$, puzzle solution

finding algorithm $\mathcal{C}()$, and solution verification algorithm $\mathcal{V}()$. To enable a consensus protocol, the PoW puzzle has to meet the following requirements: (i) it has to be moderately hard to compute (for machines), and no prover can create a proof of work in no time; (ii) it has to be easy to verify (for machines), and all verifiers can efficiently check if a proof is valid; (iii) the difficulty needed in order to solve the proof has to be adjustable in a linear way; and (iv) it has to be possible to ensure that proofs of work cannot be reused multiple times, and the proofs of work should be linked to some public information, e.g., the hash of the block header in a consensus protocol.

Proof of human-work puzzles are very similar to PoW puzzles, except that we intend to have the human in the loop for finding the solution. The key difference is that the prover (problem solver) should not be machine-only. In the above listed requirements, we therefore expect the PoH puzzle to be *moderately hard to compute for humans, and infeasible to compute for machines*. On the other hand, as in PoW, we expect the verification to be easy for machines¹². The syntax is as follows:

Definition 6 (Proof of Human-work Puzzle). *A proof of human-work puzzle-system consists of a tuple of algorithms $(\text{Setup}, \mathbf{G}, \mathcal{C}^{\mathcal{H}}, \mathbf{V})$, where*

- **Setup** is a randomized system setup algorithm that takes as input 1^λ (λ is the security parameter) and 1^ω (ω is the difficulty parameter), and outputs a system public parameter $\text{PP} \leftarrow \text{Setup}(1^\lambda, 1^\omega)$;
- **G** is a randomized puzzle generation algorithm that takes as input the public parameter PP , and outputs puzzle $x \leftarrow \mathbf{G}(\text{PP})$;
- $\mathcal{C}^{\mathcal{H}}$ is a solution finding algorithm (i.e., human-machine solver) that takes as input the public parameter PP and a puzzle x , and outputs value $a \leftarrow \mathcal{C}^{\mathcal{H}(\cdot)}(\text{PP}, x)$ as the solution to the puzzle x . Here, $\mathcal{H}(\cdot)$ denotes the human oracle which takes intermediate human-efficient objects (such as images) as inputs, and returns machine-efficient values as outputs.
- **V** is a deterministic puzzle-solution verification algorithm that takes as input the public parameter PP and a puzzle-solution pair (x, a) , and outputs bit $b := \mathbf{V}(\text{PP}, x, a)$ where $b = 1$ if a is a valid solution to the puzzle x , and $b = 0$ otherwise.

Following notation of Miller et al. [42] we will let $\zeta(m, \omega) \doteq 1 - (1 - 2^{-\omega})^m$. Intuitively, $\zeta(m, \omega)$ denotes the probability of finding a valid solution with m queries to the human-oracle.

Definition 7 (Honest Human Solvability). *A PoH puzzle system $(\text{Setup}, \mathbf{G}, \mathcal{C}^{\mathcal{H}}, \mathbf{V})$ is honest human solvable if for every polynomial $m = m(\lambda)$, and for any honest human-machine solver $\mathcal{C}^{\mathcal{H}(\cdot)}$ who controls m human-work units, it holds that*

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda, 1^\omega); x^* \leftarrow \mathbf{G}(\text{PP}); \\ a^* \leftarrow \mathcal{C}^{\mathcal{H}(\cdot)}(\text{PP}, x^*) \quad : \quad \mathbf{V}(\text{PP}, x^*, a^*) = 1 \end{array} \right] \geq \zeta(m, \omega) - \text{negl}(\lambda)$$

¹² We remark that, it might also be interesting to consider the variant in which verification is easy for human but not for machine-verifiers.

Definition 8 (Adversarial Human Unsolvability). We say that a PPT algorithm \mathcal{B} breaks security of the a PoH puzzle system $(\text{Setup}, \mathbf{G}, \mathbf{C}^{\mathcal{H}}, \mathbf{V})$ if for some polynomials $m = m(\lambda)$ and $\mu(\lambda)$ when \mathcal{B} controls at most m human-work units, it holds that

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda, 1^\omega); x^* \leftarrow \mathbf{G}(\text{PP}); \\ a^* \leftarrow \mathcal{B}^{\mathcal{H}(\cdot)}(\text{PP}, x^*) \quad : \quad \mathbf{V}(\text{PP}, x^*, a^*) = 1 \end{array} \right] \geq \zeta(m+1, \omega) + \frac{1}{\mu(\lambda)}$$

If no PPT human-machine adversary $\mathcal{B}^{\mathcal{H}(\cdot)} \in \text{Discoverable}_X$ breaks security we say that the PoH puzzle system is *adversarial human unsolvable* for the next X years.

Remark 2. We remark that the above definition can be strengthened by providing the adversarial \mathcal{B} additional access to a polynomial number of (x_i, a_i) pairs, where $x_i \leftarrow \mathbf{G}(\text{PP})$ and $\mathbf{V}(\text{PP}, x_i, a_i) = 1$. The definition can be further strengthened further by providing the adversarial \mathcal{B} multiple puzzle instances x_1^*, \dots, x_k^* , and asking \mathcal{B} to output a valid a_j^* for any $j \in [k]$. Our construction in next section can achieve these strengthened notions. For simplicity, we focus on the above simplified notion in this paper.

3.2 Construction

In this subsection, we show how to construct PoH puzzles for cryptocurrency. In Bitcoin each PoW puzzle instance is specified by the public ledger x . A motivated miner (i.e., the PoW prover) will produce a PoW by repeatedly querying a random oracle RO (e.g., the SHA256 hash function) to sample uniformly random elements in an attempt to produce a “small” output. More concretely, the miner computes random elements $y_i = \text{RO}(x, s_i)$ for different strings s_i 's. If there exist i so that $y_i < T_\omega$, then the corresponding s_i can be viewed as the PoW solution. Given a random oracle $\text{RO} : \{0, 1\}^* \rightarrow \{0, 1\}^n$ we will use the notation $T_\omega \doteq 2^{n-\omega}$. Intuitively, this ensures that $\text{RO}(x, s_i) < T_\omega$ with probability $2^{-\omega}$.

To have human in the loop, we need to first sample CAPTCHA instances for human solvers. Those instances are not in uniform distribution, and it is unclear if we can use a random oracle RO to generate such instances. We here use a cryptographic tool called “universal sampler” recently developed by Hofheinz et al. [32] to generate such CAPTCHA instances. Universal sampler can be viewed as an extended version of RO , which can generate elements in any efficiently samplable distributions. More concretely, we fix d to be a circuit for computing the CAPTCHA generation function CAPT.G . Thus, $d(r)$ generates a CAPTCHA puzzle Z_r and a tag tag_r from randomness r . Now, the miner begins by computing $(Z_i, \text{tag}_i) = \text{Sample}(U, d, \beta = (x, s_i))$; then the miner solves Z_i via human effort to get the corresponding CAPTCHA solution σ_i ; at this moment, we can adapt the strategy in the original PoW by computing $y_i = \text{RO}(x, s_i, \sigma_i)$ and if $y_i < T_\omega$ and if the CAPTCHA solution σ_i is correct, then the corresponding pair (s_i, σ_i) can be viewed as the PoH solution. We can verify that the solution is correct by re-sampling $(Z_i, \text{tag}_i) \leftarrow \text{Sample}(U, d, \beta = (x, s_i))$ and checking that $\text{Verify}(Z_i, \text{tag}_i, \sigma_i) = 1$ and that $\text{RO}(x, s_i, \sigma_i) < T_\omega$.

Construction Details. In our proof of human-work puzzle construction, we use a universal sampler scheme $\text{UNI} = \text{UNI}.\{\text{Setup}, \text{Sample}\}$, a CAPTCHA scheme $\text{CAPT} = \text{CAPT}.\{\text{Setup}, \text{W}, \text{G}, \text{C}^{\mathcal{H}}, \text{Verify}\}$, and a hash function \mathbf{G} . We will treat \mathbf{G} as a random oracle in our analysis. The constructed PoH puzzle scheme consists of algorithms $\text{POH}.\{\text{Setup}, \mathbf{G}, \text{C}^{\mathcal{H}}, \mathbf{V}\}$. Note that \mathcal{H} denotes a human oracle.

- The setup algorithm $\text{PP} \leftarrow \text{POH}.\text{Setup}(1^\lambda, 1^\omega)$: Compute $\tilde{\text{PP}} \leftarrow \text{CAPT}.\text{Setup}(1^\lambda)$; Compute $U \leftarrow \text{UNI}.\text{Setup}(1^\lambda)$; Define a program d as follows: On input randomness $r = (r_1, r_2)$, compute $\sigma := \text{CAPT}.\text{W}(\tilde{\text{PP}}; r_1)$, $(Z, \text{tag}) := \text{CAPT}.\mathbf{G}(\tilde{\text{PP}}, \sigma; r_2)$, and output (Z, tag) . Set $\text{PP} := (U, d, \tilde{\text{PP}}, T = T_\omega, \text{PARAM})$ where PARAM denotes the instructions of using the system.
- The puzzle generation algorithm $x \leftarrow \text{POH}.\mathbf{G}(\text{PP})$: Parse PP into $(U, d, \tilde{\text{PP}}, T, \text{PARAM})$; Based on the description of PARAM , sample x .
- The solution function $a \leftarrow \text{POH}.\text{C}^{\mathcal{H}}(\text{PP}, x)$: Upon receiving puzzle instance x , parse PP into $(U, d, \tilde{\text{PP}}, T, \text{PARAM})$; Randomly choose $s \leftarrow \{0, 1\}^\lambda$; Compute CAPTCHA puzzle instance $(Z, \text{tag}) \leftarrow \text{UNI}.\text{Sample}(U, d, \beta = (x, s))$; Use the human oracle \mathcal{H} to find a solution to CAPTCHA puzzle instance Z , i.e., $\sigma \leftarrow \text{CAPT}.\text{C}^{\mathcal{H}}(\tilde{\text{PP}}, Z)$; If $\mathbf{G}(x, s, \sigma) < T$, then set $a := (s, \sigma)$. Otherwise set $a := \perp$.
- The puzzle verification algorithm $b := \text{POH}.\mathbf{V}(\text{PP}, x, a)$: Parse a into (s, σ) ; Parse PP into $(U, d, \tilde{\text{PP}}, T, \text{PARAM})$; Compute $(Z, \text{tag}) \leftarrow \text{UNI}.\text{Sample}(U, d, \beta = (x, s))$; If $\text{CAPT}.\text{Verify}(\tilde{\text{PP}}, Z, \text{tag}, \sigma) = 1$ and $\mathbf{G}(x, s, \sigma) < T$, then set $b := 1$. Otherwise set $b := 0$.

It is easy to verify that the PoH scheme is honest human solvable if the underlying universal sampler is correct and the CAPTCHA scheme is honest-human solvable. Next we state a theorem for the security of our PoH scheme, and the proof can be found in the full version [12]. In our proof we give an explicit PPT reduction R from CAPTCHA to PoH security. Intuitively, this means that if mankind finds a PPT algorithm $\mathcal{B} \in \text{Discoverable}_X$ attacking PoH security then mankind will quickly find a PPT algorithm \mathcal{A} breaking CAPTCHA security. We take this to mean that if $\mathcal{B} \in \text{Discoverable}_X$ then $\mathcal{A} \in \text{Discoverable}_{X+\epsilon}$, where ϵ is the time necessary to apply a known, efficient blackbox reduction¹³.

¹³ In this (informal) line of reasoning we implicitly assume if there is an PPT algorithm breaking PoH is discovered then the reduction R will be quickly implemented by someone because it is publicly known and leads to a very useful result (breaking CAPTCHAs). We stress that we are not claiming that $R(\text{Discoverable}_X) \subset \text{Discoverable}_{X+\epsilon}$ for every known, explicit reduction R so that the set $\text{Discoverable}_{X+\epsilon}$ contains the result of applying every explicit, known reduction to every machine in the set Discoverable_X . If this was the case then we could claim that (at minimum) the set Discoverable_X contains all strings of length X/ϵ since the reductions $R_1 = \text{“append 1”}$ and $R_0 = \text{“append 0”}$ are explicit, known reductions. In our case we are assuming that the known, explicit reduction R from CAPTCHAs to PoHs would be implemented *because* it is known that the reduction leads to a useful result when we have an algorithm to break PoH security (breaking CAPTCHAs). By contrast, the reduction “append 1” is unlikely to lead to useful results when applied to most Turing Machines.

Thus, our POH construction is essentially as secure as the underlying construction. CAPT is a computer uncrackable CAPTCHA for the next $X + \epsilon$ years (definition 4), then the above proof of human-work scheme POH is adversarial human unsolvable for the next X years. We use ϵ to denote the time necessary (e.g., 1 day) to implement the reduction and build the resulting PPT CAPTCHA solver.

Theorem 1. *If UNI is an adaptively secure universal sampler then given any PPT algorithm \mathcal{B} that breaks POH security (Definition 8) there is an explicit PPT blackbox reduction producing a PPT algorithm \mathcal{A} that breaks CAPTCHA security (Definition 4).*

Proof (idea). The security of our PoH relies on the security of underlying building blocks, the universal sampler scheme UNI, and the CAPTCHA scheme CAPT. We start from the real security game. Based on the security of the universal sampler scheme UNI, we can modify the real security game into a hybrid world where CAPTCHA puzzle instances are generated independently and based on uniform randomness. Then we can use the security of CAPT to argue about the security of PoH. That is, we can construct a CAPT attacker $\mathcal{A}_{\text{CAPT}}$ based on a PoH attacker \mathcal{A}_{POH} . The CAPT attacker $\mathcal{A}_{\text{CAPT}}$ can simulate an internal copy of \mathcal{A}_{POH} , and embed his challenge into a simulated hybrid for \mathcal{A}_{POH} . If \mathcal{A}_{POH} wins with more than specified probability (i.e., $\zeta(m + 1, \omega)$) plus non-negligible probability, then $\mathcal{A}_{\text{CAPT}}$ can also win the computer-unbreakable game with non-negligible probability.

4 Application 1: HumanCoin

In this section we outline how a new cryptocurrency called *HumanCoin* could be built using Proofs of Human-work. At a high level HumanCoin closely follows the Bitcoin protocol, except that we use PoH puzzles to extend the blockchain instead of PoW puzzles. We will not attempt to describe HumanCoin in complete detail. Instead we will focus on the key modifications that would need to be made to an existing cryptocurrency like Bitcoin to use Proof of Human work puzzles. In our discussion we will use lowercase bitcoin (resp. humancoin) to denote the base unit of currency in the Bitcoin (resp. HumanCoin) protocol.

4.1 Bitcoin Background

We begin by highlighting several of the key features of Bitcoin. Our overview follows the systemization of knowledge paper by Bonneau et al. [14]. However, our discussion of Bitcoin is overly simplified and this choice is intentional. For example, we will completely ignore the use of Merkle Trees [41] in Bitcoin to compress the blockchain even though it is quite useful in practice. We make this choice so that we can focus on the key differences of HumanCoin (the use of Merkle Trees [41] in HumanCoin and Bitcoin would be identical). We do include additional discussion of Bitcoin in the full version, but even this discussion is

not intended to be complete. We refer interested readers to the excellent lectures by Narayanan et al. [46] for more details about Bitcoin or the original paper published under the pseudonym Nakamoto [45].

Blockchain. In Bitcoin all transactions (e.g., “Alice sends Bob 50 bitcoins”) are published on a public ledger. This public ledger is stored on a cryptographic data structure called a blockchain $b = B_0, \dots, B_t$. A blockchain b is valid if and only if all of the blocks B_i ($i \leq t$) are valid and an individual block $B_i = (tx_i, s_i, h_{i-1})$ is valid if and only if three key conditions are satisfied. First, all of the transactions recorded in the transcript tx_i must be valid (e.g., each transaction is signed by the sender and the spender has sufficient funds). Second, the block B_i must contain the cryptographic hash $h_{i-1} = \text{hash}(B_{i-1})$ of the previous block B_{i-1} ¹⁴. Finally, the block B_i should contain a nonce s_i which ensures that cryptographic hash $\text{hash}(B_i)$ begins with at least ω leading zeros, where ω is a hardness parameter that we will discuss later. Finding such a nonce s constitutes a proof of work in the Hashcash [3] puzzle system. The first property ensures that users cannot spend money they don’t have and that they cannot spend someone else’s money. The second property ensures that it is impossible to tamper with blocks B_i in the middle of the blockchain without creating an entirely new blockchain $b' = B_0, \dots, B_{i-1}, B'_i, B'_{i+1}, \dots, B'_t$. Finally, the third property ensures that it is moderately difficult to add new blocks to a blockchain. To incentivize miners to help validate transactions (i.e. extend the blockchain by finding a valid nonce s) the miner is allowed to add a special transaction (e.g., “I create 25 new bitcoins and give them to myself”) to the new block as a reward .

Distributed Consensus Protocol. Bitcoin’s distributed consensus protocol is simple, yet elegant. An agent should accept a transaction if and only if it is recorded on a block B_i of a valid blockchain $b = B_0, \dots, B_t$ and b is the longest valid that the agent has seen and $i \leq t - 6$. Unless a miner controls at least 25% of the hash power in the network the rational mining strategy is always to extend the longest blockchain because nobody will accept the Bitcoins they try to mine in a shorter blockchain (e.g., the special transaction in which a miner claims 25 bitcoins’ would only be recorded on a shorter blockchain which nobody accepts) [27]. Assuming that the network has high synchronicity [29] and that a malicious user controls at most 49% of the computational mining power he will never be able to tamper with any of the transactions in a block B_i from the middle of the blockchain because he would need to eventually produce a new blockchain $b' = B_0, \dots, B_{i-1}, B'_i, B'_{i+1}, \dots, B'_t$ that is at least as long as the true blockchain b and he will fail to accomplish this goal with high probability [45].

4.2 HumanCoin

Similar to Bitcoin all HumanCoin transactions (e.g., “Alice sends Bob 50 humancoins”) are recorded inside a blockchain $b = B_0, \dots, B_t$, where each block

¹⁴ Bitcoin uses the cryptographic hash function $\text{hash} = \text{SHA256}$. The function hash is typically treated as a random oracle in security analysis of Bitcoin.

$B_i = (tx_i, a_i, h_{i-1})$ contains three components: a list of transactions tx_i , a hash $h_{i-1} = \text{hash}(B_{i-1})$ of the previous block, and a Proof of Human-work which is encoded by a_i . As before all of the transactions in tx_i must be valid and the block must contain the hash $h_{i-1} = \text{hash}(B_{i-1})$ of the previous block. We additionally require that the PoH verifier accepts the Proof of Human-Work solution a_i . More formally, suppose that we are given a PoH puzzle system $(\text{Setup}, \mathbf{G}, \mathbf{C}^{\mathcal{H}}, \mathbf{V})$ and that we have already run $\text{Setup}(1^\lambda, 1^\omega)$ to obtain public parameters PP which are available to every miner. A valid block B_i must contain a value a_i such that the public verifier $\mathbf{V}(\text{PP}, x_i, a_i)$ outputs 1, where $x_i = \mathbf{G}(\text{PP}; r = \text{hash}(tx_i, h_{i-1}))$. Given a valid blockchain $b = B_0, \dots, B_t$ a miner can earn HumanCoins by finding a valid block $B_{t+1} = (tx_{t+1}, a_{t+1}, x_{t+1}, h_t)$ extending b . To find such a block the human-computer miner would first set $r = \text{hash}(tx_{t+1}, h_t)$ and then sample $x \leftarrow \mathbf{G}(\text{PP}; r)$. Finally, the human-computer miner can run $\mathbf{C}^{\mathcal{H}}(\text{PP}, x)$ to obtain a potential solution a . If $a = \perp$ then the miner will need to try again. Otherwise, the miner has found a valid proof of human-work and he can produce a valid new block $B_{t+1} = (tx_{t+1}, a, h_t)$ by adding inserting the PoH solution a into the block B_{t+1} . As before the miner is allowed to insert a special transaction into the new block (e.g., “I create 25 humancoins and give them to myself”) as a reward for extending the blockchain.

Parameter Selection. In Bitcoin ω is a public parameter is tuned to ensure that, on average, miners will add one new block to the blockchain every 10 minutes [46] — on average we need 2^ω hash evaluations to create one new block. The Bitcoin protocol would most likely work just fine with a shorter delay (e.g., 5 minutes) or a slightly longer delay (e.g., 20 minutes) between consecutive blocks — there is nothing magical about the specific target value of 10 minutes. However, it is clear that there needs to be some delay to promote stability. If multiple miners find a new block at the same time then we could end up with competing blockchains resulting in temporary confusion. Note that if the value of ω remains fixed then the average time to create one new block would begin to decrease as more miners join Bitcoin, or as existing miners upgrade their computational resources. Thus, the value of ω must be adjusted periodically. In Bitcoin the value of ω is adjusted every 2,016 blocks, which works out to two weeks on average (2 weeks = 2016×10 minutes), using the formula $\omega = \omega_{old} - \log\left(\frac{t_{elapsed}}{2016 \times 10 \text{ min}}\right)$, where $t_{elapsed}$ denotes the time span that it actually took to generate the last 2,016 blocks [46].

In HumanCoin we adjust ω in exactly the same way. Note that the PoH hardness parameter $T_\omega = 2^{n-\omega}$ in our PoH construction is a public parameter PP and can easily be modified as it is not embedded into any of the obfuscated programs. In HumanCoin we will need to select an initial value of ω that is *much* smaller than in Bitcoin if we want ensure that new block are discovered every 10 minutes. This is because computers can evaluate a hash function hash much faster than a human can solve a long CAPTCHA puzzle. However, we could still use the same basic formula to tune the hardness parameter ω of our proof of work puzzles in the event that many miners join/leave.

5 Application 2: Password Protection

An adversary who breaches an authentication server is able to mount an automated brute-force attack by comparing the cryptographic hash of each user's password with the cryptographic hashes of likely password guesses. These offline attacks have become increasingly prevalent and dangerous as password cracking resources has improved. In particular, the cost of computing a hash function H like SHA256 or MD5 on an Application Specific Integrated Circuit (ASIC) is orders of magnitude smaller than the cost of computing H on traditional hardware [21,46]. Similarly, data from previous breaches allow adversaries to improve their guessing strategies. Recent security breaches (e.g., Ashley Madison, LastPass, RockYou, LinkedIn and eBay to name a few ¹⁵), which have affected millions of users, highlight the importance of this problem.

Canneti et al. [18] had a clever idea to deter an offline attacker that they called Human Only Solvable Puzzles. They proposed filling a hard drive with a dataset of unsolved CAPTCHA puzzles. When a user authenticates he will be challenged with a pseudorandom CAPTCHA puzzle from the dataset, and the server will append the solution to the user's password before computing the hash value. The choice of the pseudorandom CAPTCHA puzzle becomes deterministic once the user's password and username are fixed. Thus, if the user types in the same password he will receive the exact same CAPTCHA puzzle as a challenge. If the underlying CAPTCHA system is human usable, then the user will always be able to authenticate successfully provided that he can remember his password. If an offline adversary wants to verify a password guess he will need to find and solve the corresponding CAPTCHA puzzle. The key point is that each time the adversary tries a new guess he will need to solve a different CAPTCHA challenge.

Unfortunately, the Human Only Solvable Puzzles solution of [18] has one critical drawback. There are a finite number of CAPTCHAs on the hard drive, and the defense will break down once the adversary manages to solve all (or most) of them. Blocki et al. [9] estimated that it would cost about $\$10^6$ to solve all of the CAPTCHAs on an 8 TB hard drive. While this is certainly an expensive start-up cost it may not be sufficient to deter the adversary because these costs would amortize over all user accounts. Many password breaches affect millions of users, and each cracked password has significant value on the black market (e.g., $\$4$ – $\$30$). Blocki et al. [9] introduced their own scheme called GOTCHA based on inkblot images, but their protocol had higher usability costs and was based on newer untested AI assumptions.

In this section we introduce a provably secure password authentication scheme in the Random Oracle model using CAPTCHAs and program obfuscation. Unlike Blocki et al. [9] our solution can be based on standard CAPTCHA assump-

¹⁵ See <http://www.privacyrights.org/data-breach/> (Retrieved 9/1/2015).

tions. Unlike Canneti et al. [18] our solution is not vulnerable to pre-computation attacks¹⁶.

5.1 Password Authentication Scheme

We first formalize the notion of a password authentication scheme. Definition 9 formalizes the account creation and authentication algorithms from the perspective of an authentication server. We note that the server is allowed to interact with the human user \mathcal{H} during the account creation and authentication protocols.

Definition 9. *A password authentication scheme consists of a tuple of algorithms $(\text{Setup}, \text{CreateAccount}^{\mathcal{H}}, \text{Authenticate}^{\mathcal{H}})$ and a random oracle \mathbf{G} , where*

- *Setup is a randomized system setup algorithm that takes as input 1^λ (λ is the security parameter) and outputs a system public parameter $\text{PP} \leftarrow \text{Setup}(1^\lambda)$;*
- *CreateAccount ^{\mathcal{H}} is an account creation algorithm that takes as input the public parameter PP , a username u and a password pwd and outputs a tuple (h, s) . Here, s is typically a random bit string (salt) and h is a hash value produced by the random oracle. We note that CreateAccount ^{\mathcal{H}} is a human-machine algorithm and thus the hash value h may include the solution to CAPTCHAs that the human solves as well as the password pwd and salt s ;*
- *Authenticate ^{\mathcal{H}} is the algorithm that is invoked when a user wants to authenticate. The algorithm takes as input the public parameter PP , a username u , a password pwd , a hash h and a salt value s and outputs a bit $b \in \{0, 1\}$ indicating whether or not the authentication attempt was successful. We note that Authenticate ^{\mathcal{H}} is a human-machine algorithm and thus the human \mathcal{H} may be asked to solve CAPTCHAs as part of the authentication procedure.*

Our next definition says what it means for a password authentication scheme to be costly to crack. The game mimics an offline adversary who has breached the authentication server and stolen the record (u, h, s) indicating that user u has an account with salt value s and the salted hash of the user’s password needs to match h . In our definition we let \mathcal{P} denote a distribution over the passwords $\{\text{pwd}_1, \dots, \text{pwd}_n\}$ that the user u might select and let $p_i = \Pr_{\mathcal{P}}[\text{pwd}_i]$ denote the probability that the user selects password pwd_i . We assume that p_i and pwd_i are known to the adversary for all i and for convenience we assume that the passwords are ordered such that $p_1 \geq p_2 \geq \dots \geq p_n$. Informally, our definition states that an adversary with B units of human-work will succeed in cracking the user’s password with probability at most $p_1 + \dots + p_B + \text{negl}(\lambda)$.

Definition 10 (Costly to Crack). *We say a PPT adversary \mathcal{A} breaks security of a password authentication scheme $\{\text{Setup}, \text{CreateAccount}^{\mathcal{H}}, \text{Authenticate}^{\mathcal{H}}\}$,*

¹⁶ Of course the main downside to our approach is the dependence on indistinguishability obfuscation, which does not have practical solutions at this time.

\mathbf{G} } if for some polynomials $B = B(\lambda)$ and $\mu(\lambda)$ and user u , whenever \mathcal{A} has B human-work units it holds that

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda); \text{pwd} \leftarrow \mathcal{P}; \\ (h, s) \leftarrow \text{CreateAccount}^{\mathcal{H}}(\text{PP}, u, \text{pwd}); \\ \text{pwd}' \leftarrow \mathcal{A}^{\mathcal{H}(\cdot)}(\text{PP}, h, s) \quad : \\ \text{Authenticate}^{\mathcal{H}}(\text{PP}, u, \text{pwd}', h, s) = 1 \end{array} \right] \geq p_1 + \dots + p_B + \frac{1}{\mu(\lambda)}$$

We say that the password authentication scheme is *costly to crack* for the next X years if for any PPT adversary $\mathcal{A} \in \text{Discoverable}_X$ \mathcal{A} does not break security. We remark that we do not require the adversary's success probability to be negligibly small. Indeed, if the user selects passwords from a distribution with low entropy (and many users do [13]) then the adversary may have a good success rate. Thus, the problem is unavoidable as long as users are allowed to select low-entropy passwords¹⁷. We do not focus on helping users to select strong passwords [15,10], although this is indeed an important direction of research. Our goal is to provide the best possible protection for the passwords that users actually select.

The next definition quantifies human usability. Informally, the password authentication scheme is usable if an honest human user will always be able to authenticate if he remembers his password. We stress that our definition does not say anything about how easy it will be to remember the password. While this is certainly an important consideration it is orthogonal to our work. We are not focused on how to get users to choose stronger passwords, but rather how to more effectively protect the passwords that users actually choose. Our definition merely says that an honest user won't be locked out of his account as long as he remembers his password (e.g., because he cannot solve the CAPTCHAs).

Definition 11 (Human Usable). *We say that a password authentication scheme $\{\text{Setup}, \text{CreateAccount}^{\mathcal{H}}, \text{Authenticate}^{\mathcal{H}}, \mathbf{G}\}$ is human usable if for every human user \mathcal{H} who controls 1 human-work unit during authentication and 1 human work unit during account creation, it holds that*

$$\Pr \left[\begin{array}{l} \text{PP} \leftarrow \text{Setup}(1^\lambda); \text{pwd} \leftarrow \mathcal{P}; \\ (h, s) \leftarrow \text{CreateAccount}^{\mathcal{H}}(\text{PP}, u, \text{pwd}) \quad : \\ \text{Authenticate}^{\mathcal{H}}(\text{PP}, u, \text{pwd}, h, s) = 1 \end{array} \right] \geq 1 - \text{negl}(\lambda)$$

5.2 Construction

Construction Details. In our construction we use a universal sampler scheme $\text{UNI} = \text{UNI}.\{\text{Setup}, \text{Sample}\}$, a CAPTCHA scheme $\text{CAPT} = \text{CAPT}.\{\text{Setup}, \mathbf{W}, \mathbf{G}, \mathcal{C}^{\mathcal{H}}, \text{Verify}\}$, and a hash function \mathbf{G} . We will treat \mathbf{G} as a random oracle in our analysis. The constructed password authentication scheme consists of algorithms $\text{Password}.\{\text{Setup}, \text{CreateAccount}^{\mathcal{H}}, \text{Authenticate}\}$. Note that \mathcal{H} denotes a human oracle.

¹⁷ In addition to their high usability costs [28], policies aimed at forcing users to choose stronger passwords (e.g., requiring numbers and capital letters) can have the opposite affect on password strength [37,11].

- The setup algorithm $\text{PP} \leftarrow \text{Password.Setup}(1^\lambda)$: Compute $\tilde{\text{P}} \leftarrow \text{CAPT.Setup}(1^\lambda)$; Compute $U \leftarrow \text{UNI.Setup}(1^\lambda)$; Define a program d as follows: On input randomness $r = (r_1, r_2)$, compute $\sigma := \text{CAPT.W}(\tilde{\text{P}}; r_1)$, $(Z, \text{tag}) := \text{CAPT.G}(\tilde{\text{P}}, \sigma; r_2)$, and output Z . Set $\text{PP} := (U, d, \tilde{\text{P}}, \text{PARAM})$ where PARAM denotes the instructions of using the system.
- The account creation algorithm $(h, s) \leftarrow \text{Password.CreateAccount}^{\mathcal{H}}(\text{PP}, u, \text{pwd})$: Parse PP into $(U, d, \tilde{\text{P}}, \text{PARAM})$; randomly choose $s \leftarrow \{0, 1\}^\lambda$. Set $\beta = (u, \text{pwd}, s)$ and compute CAPTCHA puzzle instance $Z \leftarrow \text{UNI.Sample}(U, d, \beta = (u, \text{pwd}, s))$; Use the human oracle \mathcal{H} to find a solution to CAPTCHA puzzle instance Z , i.e., $\sigma \leftarrow \text{CAPT.C}^{\mathcal{H}}(\tilde{\text{P}}, Z)$; Compute $h \leftarrow \mathbf{G}(\text{pwd}|\sigma|s)$ and output (h, s) .
- The authentication algorithm $b \leftarrow \text{Password.Authenticate}^{\mathcal{H}}(\text{PP}, u, \text{pwd}, h, s)$: Parse PP into $(U, d, \tilde{\text{P}}, \text{PARAM})$, set $\beta = (u, \text{pwd}, s)$ and compute CAPTCHA puzzle instance $Z \leftarrow \text{UNI.Sample}(U, d, \beta)$; Use the human oracle \mathcal{H} to find a solution to CAPTCHA puzzle instance Z , i.e., $\sigma \leftarrow \text{CAPT.C}^{\mathcal{H}}(\tilde{\text{P}}, Z)$; Compute $h' \leftarrow \mathbf{G}(\text{pwd}|\sigma|s)$. If $h' = h$ then output $b = 1$; otherwise output 0.

It is easy to verify that **Password** is human usable if the underlying CAPTCHA scheme **CAPT** is honest human solvable. At an philosophical level we can interpret Theorem 2, our main technical result in this section, to say that the above password authentication scheme $\text{Password}.\{\text{Setup}, \text{CreateAccount}^{\mathcal{H}}, \text{Authenticate}^{\mathcal{H}}\}$ is costly to crack for the next X years as long as the underlying CAPTCHA scheme is computer uncrackable for the next $X + \epsilon$ years (Definition 10). Here, ϵ denotes the time it takes to implement an explicit (blackbox) PPT reduction from **CAPT** to the password scheme (e.g., one day). We stress that we only need to assume that the underling CAPTCHA scheme is computer uncrackable in the more traditional sense of Definition 3 (e.g., the adversary is only given the puzzles Z_1, \dots, Z_n and not the associated verification tags).

Theorem 2. *If **UNI** is an adaptively secure universal sampler then given a PPT algorithm \mathcal{B} that breaks security of our password authentication scheme there is a PPT, blackbox reduction which produces a PPT algorithm \mathcal{A} to break CAPTCHA security (Definition 3).*

Proof (idea). At a high level we show that we can construct an adversary that breaks CAPTCHAs (under Definition 3) from an adversary that breaks the password authentication scheme. To do this we embed challenge CAPTCHAs Z_1, \dots, Z_n inside the UniversalSampler **UNI** (we can do this by the security of the Universal Sampler scheme). Intuitively, in order to check that a password guess pwd_i is correct the adversary will need to query the random oracle \mathbf{G} with the value $\beta_i = (\text{pwd}_i|\sigma_i|u)$, where σ_i is the correct solution to CAPTCHA Z_i . If the adversary queries \mathbf{G} with $B + 1$ unique solutions then we can win the CAPTCHA challenge (Definition 3) by simply outputting these $B + 1$ solutions. If the adversary queries \mathbf{G} with at most B unique solutions then we can show that his success rate is at most $p_1 + \dots + p_B + \text{negl}(\lambda)$. A formal proof of Theorem 2 can be found in the full version [12].

Discussion. We believe that the construction of our secure password authentication scheme might lead to many other useful applications. For example, the scheme might allow us to use human memorable (i.e., lower entropy) secrets to secure highly confidential data like secret keys. Let pwd_i be the user’s password and let σ_i denote the solution to the corresponding CAPTCHA challenge. The random oracle value $R_i = \mathbf{G}(pwd_i, \sigma_i, s, 1|i)$ is completely uncorrelated with any information that the adversary can obtain without discovering the user’s password. The random values R_1, R_2, \dots could be used as a one-time pad to efficiently encrypt/decrypt information on a hard drive or to (re)derive private keys for a signature scheme.

We also note that our authentication scheme could potentially be modified to make the proof of human work *safely exportable* and that the amount of human work during authentication can easily be tuned. For example, suppose that Bob wants to protect his passwords, but that he is too busy to solve CAPTCHAs. During authentication, after Bob enters his password and receives the CAPTCHA challenge, Bob might like to pay other human(s) to solve the CAPTCHA challenge for him. However, he wants to make sure that his password is not exposed if these contracted workers are malicious. For example, we might replace the hash of the password with an obfuscation of two program $P_{K,pwd}$ and G_K . Here, $G_K(x, pwd)$ generates a CAPTCHA puzzle Z using randomness $(r_1, r_2) = \text{PRF}_K(x, pwd)$ in the procedures CAPT.W and CAPT.G respectively. $P_{K,pwd}(pwd', \sigma, x)$ outputs 1 if and only if $pwd' = pwd$ and σ is the correct solution to the puzzle Z output by $G_K(x, pwd)$. During authentication we can obtain the puzzle Z by running $G_K(x, pwd')$ with a uniformly random string $x \in \{0, 1\}^\lambda$ which should be discarded immediately after the authentication session finishes. As long as Bob keeps the value x secret he can safely share the puzzle Z with other users. However, this modified authentication protocol is merely a heuristic as we do not have any formal security proof that it is hard for a computer to solve CAPTCHA puzzles generated by G_K when given the obfuscated source code $i\mathcal{O}(G_K)$.

As another application we could use the same general framework as a way to detect bots *without* interaction! Suppose that we rename the algorithms `CreateAccountH` and `AuthenticateH` to algorithms `GenerateVerifiedMessageH` and `VerifyMessageH`. The algorithms have essentially the same functionality except for a few minor modifications: 1) the password field pwd is renamed to denote a message m that a user Alice wishes to send to Bob, 2) we replace the username u with a pair (u_1, u_2) where u_1 denotes the sender and u_2 denotes the intended receiver, and we fix the salt value $s = \mathbf{G}(u_1, u_2, m)$ for a given message m that a user u_1 wishes to send to u_2 . To send the message m to Bob Alice would first execute `GenerateVerifiedMessageH(PP, (Alice, Bob), m)` and solve the corresponding CAPTCHA to obtain a tuple (h, s) . Now Alice sends the tuple $(Alice, Bob, m, h, s)$ to Bob. At this point Alice is finished with the protocol. Bob runs `VerifyMessageH(PP, (Alice, Bob), m, h, s)` and solves the corresponding CAPTCHA to obtain a bit b . If $b = 1$ then Bob accepts that a human

(possibly Alice) spent time and energy to send the him the message m ¹⁸. If $b = 0$ then Bob may dismiss the message as potentially being produced by a bot.

6 Future Challenges

While we believe that Proofs of Human Work could have many benefits, we see three primary challenges for future research. First, because our construction of PoH puzzles is based on $i\mathcal{O}$ HumanCoin is not practical without a large breakthrough in the design of practical $i\mathcal{O}$ schemes. Could we design efficient targeted obfuscation schemes for specific programs like our PoH algorithms? Second, because our PoH puzzles rely on the assumption that some underlying AI problem is hard it is possible that a cryptocurrency like HumanCoin might have a shorter shelf life (e.g., if it takes 15 years for AI researchers to break the underlying CAPTCHA then HumanCoin would expire in at most 15 years). Would it be possible for HumanCoin participants to reach a consensus to change the underlying CAPTCHA in the event of an AI breakthrough? Finally, our Proof of Human Work construction, and by extension HumanCoin, requires an initial trusted setup phase for the Proof of Human Work construction. If the Proof of Human Work system is generated by a malicious party then that party might be able to insert a trapdoor which would allow him to mine HumanCoins without any human effort. We note that this concern is not unique to HumanCoin. Other cryptocurrencies like Zerocash [5] also require an initial trusted setup phase¹⁹. Ben-Sasson et al. [6] proposed to run this trusted setup phase using secure multiparty computation. As long as at least one of the parties in this computation are honest it would be impossible for a malicious adversary to insert a backdoor. Similar techniques could also be used to minimize risks during the HumanCoin setup phase.

In addition to cryptocurrency we also showed that our PoH techniques could be applied to protect passwords and to detect bots without interaction. What other applications are possible?

Acknowledgments: The authors thank paper shepherd Peter Gaži for his very constructive feedback which helped us to improve the quality of the paper. In particular, we are thankful for his suggestions about formalizing security statements involving hard AI problems.

The authors also thank Andrew Miller, and the PC of ITCS 2016 and TCC 2016B for their helpful comments.

¹⁸ If Bob wanted to additionally verify that Alice was the human that sent the message Alice and Bob would need to use other cryptographic tools like digital signatures.

¹⁹ Arguably, even Bitcoin does require some trust assumptions during setup. For example, we need to trust that the cryptographic hash function $h = \text{SHA256}$, which is modeled as a random oracle in the Bitcoin protocol, does not have any secret backdoors. A malicious miner with a secret backdoor could easily reverse old transactions.

References

1. M. Andrychowicz, S. Dziembowski, D. Malinowski, and L. Mazurek. Secure multi-party computations on bitcoin. In *2014 IEEE Symposium on Security and Privacy*, pages 443–458. IEEE Computer Society Press, May 2014.
2. J. Aspnes, C. Jackson, and A. Krishnamurthy. Exposing computationally-challenged Byzantine impostors. Technical Report YALEU/DCS/TR-1332, Yale University Department of Computer Science, July 2005.
3. A. Back. Hashcash — A denial of service counter-measure. 2002. <http://hashcash.org/papers/hashcash.pdf>.
4. B. Barak, N. Bitansky, R. Canetti, Y. T. Kalai, O. Paneth, and A. Sahai. Obfuscation for evasive functions. In Y. Lindell, editor, *TCC 2014*, volume 8349 of *LNCS*, pages 26–51. Springer, Heidelberg, Feb. 2014.
5. E. Ben-Sasson, A. Chiesa, C. Garman, M. Green, I. Miers, E. Tromer, and M. Virza. Zerocash: Decentralized anonymous payments from bitcoin. In *2014 IEEE Symposium on Security and Privacy*, pages 459–474. IEEE Computer Society Press, May 2014.
6. E. Ben-Sasson, A. Chiesa, M. Green, E. Tromer, and M. Virza. Secure sampling of public parameters for succinct zero knowledge proofs. In *2015 IEEE Symposium on Security and Privacy*, pages 287–304. IEEE Computer Society Press, May 2015.
7. I. Bentov and R. Kumaresan. How to use bitcoin to design fair protocols. In J. A. Garay and R. Gennaro, editors, *CRYPTO 2014, Part II*, volume 8617 of *LNCS*, pages 421–439. Springer, Heidelberg, Aug. 2014.
8. I. Bentov, C. Lee, A. Mizrahi, and M. Rosenfeld. Proof of activity: Extending bitcoins proof of work via proof of stake. In *Proceedings of the ACM SIGMETRICS 2014 Workshop on Economics of Networked Systems, NetEcon*, 2014.
9. J. Blocki, M. Blum, and A. Datta. GOTCHA password hackers! In *AISec’13, Proceedings of the 2013 ACM Workshop on Artificial Intelligence and Security*, pages 25–34, 2013. <http://www.cs.cmu.edu/~jblocki/papers/aisec2013-fullversion.pdf>.
10. J. Blocki, S. Komanduri, L. F. Cranor, and A. Datta. Spaced repetition and mnemonics enable recall of multiple strong passwords. In *NDSS 2015*. The Internet Society, Feb. 2015.
11. J. Blocki, S. Komanduri, A. Procaccia, and O. Sheffet. Optimizing password composition policies. In *Proceedings of the fourteenth ACM conference on Electronic commerce*, pages 105–122. ACM, 2013.
12. J. Blocki and H.-S. Zhou. Designing proof of human-work puzzles for cryptocurrency and beyond. In *IACR Cryptology ePrint Archive 2016/145*, 2016. Full version. <http://eprint.iacr.org/2016/145>.
13. J. Bonneau. The science of guessing: Analyzing an anonymized corpus of 70 million passwords. In *2012 IEEE Symposium on Security and Privacy*, pages 538–552. IEEE Computer Society Press, May 2012.
14. J. Bonneau, A. Miller, J. Clark, A. Narayanan, J. A. Kroll, and E. W. Felten. SoK: Research perspectives and challenges for bitcoin and cryptocurrencies. In *2015 IEEE Symposium on Security and Privacy*, pages 104–121. IEEE Computer Society Press, May 2015.
15. J. Bonneau and S. Schechter. Toward reliable storage of 56-bit keys in human memory. In *Proceedings of the 23rd USENIX Security Symposium*, August 2014.
16. E. Bursztein, J. Aigrain, A. Moscicki, and J. C. Mitchell. The end is nigh: Generic solving of text-based captchas. In *8th USENIX Workshop on Offensive Technologies (WOOT 14)*, San Diego, CA, Aug. 2014. USENIX Association.

17. E. Bursztein, S. Bethard, C. Fabry, J. C. Mitchell, and D. Jurafsky. How good are humans at solving CAPTCHAs? A large scale evaluation. In *2010 IEEE Symposium on Security and Privacy*, pages 399–413. IEEE Computer Society Press, May 2010.
18. R. Canetti, S. Halevi, and M. Steiner. Mitigating dictionary attacks on password-protected local storage. In C. Dwork, editor, *CRYPTO 2006*, volume 4117 of *LNCS*, pages 160–179. Springer, Heidelberg, Aug. 2006.
19. K. Chellapilla and P. Y. Simard. Using machine learning to break visual human interaction proofs (HIPs). In *Neural Information Processing Systems (NIPS)*, pages 265–272, 2004. <https://papers.nips.cc/paper/2571-using-machine-learning-to-break-visual-human-interaction-proofs-hips.pdf>.
20. J. R. Douceur. The sybil attack. In *Peer-to-peer Systems*, pages 251–260. Springer, 2002.
21. C. Dwork, A. Goldberg, and M. Naor. On memory-bound functions for fighting spam. In D. Boneh, editor, *CRYPTO 2003*, volume 2729 of *LNCS*, pages 426–444. Springer, Heidelberg, Aug. 2003.
22. C. Dwork, J. Y. Halpern, and O. Waarts. Performing work efficiently in the presence of faults. *SIAM Journal on Computing*, 27(5):1457–1491, 1998.
23. C. Dwork and M. Naor. Pricing via processing or combatting junk mail. In E. F. Brickell, editor, *CRYPTO’92*, volume 740 of *LNCS*, pages 139–147. Springer, Heidelberg, Aug. 1993.
24. S. Dziembowski. How to pair with a human. In J. A. Garay and R. D. Prisco, editors, *SCN 10*, volume 6280 of *LNCS*, pages 200–218. Springer, Heidelberg, Sept. 2010.
25. S. Dziembowski, S. Faust, V. Kolmogorov, and K. Pietrzak. Proofs of space. In R. Gennaro and M. J. B. Robshaw, editors, *CRYPTO 2015, Part II*, volume 9216 of *LNCS*, pages 585–605. Springer, Heidelberg, Aug. 2015.
26. J. Elson, J. R. Douceur, J. Howell, and J. Saul. Asirra: a CAPTCHA that exploits interest-aligned manual image categorization. In P. Ning, S. D. C. di Vimercati, and P. F. Syverson, editors, *ACM CCS 07*, pages 366–374. ACM Press, Oct. 2007.
27. I. Eyal and E. G. Sirer. Majority is not enough: Bitcoin mining is vulnerable. In N. Christin and R. Safavi-Naini, editors, *FC 2014*, volume 8437 of *LNCS*, pages 436–454. Springer, Heidelberg, Mar. 2014.
28. D. Florêncio and C. Herley. Where do security policies come from. In *Proc. of SOUPS*, page 10, 2010.
29. J. A. Garay, A. Kiayias, and N. Leonardos. The bitcoin backbone protocol: Analysis and applications. In E. Oswald and M. Fischlin, editors, *EUROCRYPT 2015, Part II*, volume 9057 of *LNCS*, pages 281–310. Springer, Heidelberg, Apr. 2015.
30. S. Garg, C. Gentry, S. Halevi, M. Raykova, A. Sahai, and B. Waters. Candidate indistinguishability obfuscation and functional encryption for all circuits. In *54th FOCS*, pages 40–49. IEEE Computer Society Press, Oct. 2013.
31. S. Goldwasser and G. N. Rothblum. On best-possible obfuscation. In S. P. Vadhan, editor, *TCC 2007*, volume 4392 of *LNCS*, pages 194–213. Springer, Heidelberg, Feb. 2007.
32. D. Hofheinz, T. Jager, D. Khurana, A. Sahai, B. Waters, and M. Zhandry. How to generate and use universal samplers. Cryptology ePrint Archive, Report 2014/507, 2014. <http://eprint.iacr.org/2014/507>.
33. K.-F. Hwang, C.-C. Huang, and G.-N. You. A spelling based CAPTCHA system by using click. In *Biometrics and Security Technologies (ISBAST), 2012 International Symposium on*, pages 1–8, March 2012.

34. J. Kani and M. Nishigaki. Gamified CAPTCHA. In *Proceedings of Human Aspects of Information Security, Privacy, and Trust*, pages 39–48, 2013.
35. R. A. Khot and K. Srinathan. iCAPTCHA: Image tagging for free. In *Proc. Conference on Usable Software and Interface Design*, 2009.
36. A. Kiayias, H.-S. Zhou, and V. Zikas. Fair and robust multi-party computation using a global transaction ledger. In *EUROCRYPT*, 2016. <http://eprint.iacr.org/2015/574>.
37. S. Komanduri, R. Shay, P. Kelley, M. Mazurek, L. Bauer, N. Christin, L. Cranor, and S. Egelman. Of passwords and people: measuring the effect of password-composition policies. In *Proceedings of the 2011 annual conference on Human factors in computing systems*, pages 2595–2604. ACM, 2011.
38. A. Kosba, A. Miller, E. Shi, Z. Wen, and C. Papamanthou. Hawk: The blockchain model of cryptography and privacy-preserving smart contracts. In *IEEE Symposium on Security and Privacy*, 2016.
39. A. Kumarasubramanian, R. Ostrovsky, O. Pandey, and A. Wadia. Cryptography using captcha puzzles. In K. Kurosawa and G. Hanaoka, editors, *PKC 2013*, volume 7778 of *LNCS*, pages 89–106. Springer, Heidelberg, Feb. / Mar. 2013.
40. L. Lamport, R. Shostak, and M. Pease. The byzantine generals problem. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 4(3):382–401, 1982.
41. R. C. Merkle. A digital signature based on a conventional encryption function. In C. Pomerance, editor, *CRYPTO'87*, volume 293 of *LNCS*, pages 369–378. Springer, Heidelberg, Aug. 1988.
42. A. Miller, A. E. Kosba, J. Katz, and E. Shi. Nonoutsourcable scratch-off puzzles to discourage bitcoin mining coalitions. In I. Ray, N. Li, and C. Kruegel, editors, *ACM CCS 15*, pages 680–691. ACM Press, Oct. 2015.
43. G. Mori and J. Malik. Recognizing objects in adversarial clutter: Breaking a visual CAPTCHA. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 134–144, 2003.
44. M. Motoyama, K. Levchenko, C. Kanich, D. McCoy, G. M. Voelker, and S. Savage. Re: CAPTCHAs-understanding CAPTCHA-solving services in an economic context. In *19th USENIX Security Symposium, Washington, DC, USA, August 11-13, 2010, Proceedings*, pages 435–462, 2010.
45. S. Nakamoto. Bitcoin: A peer-to-peer electronic cash system. 2008. <https://bitcoin.org/bitcoin.pdf>.
46. A. Narayanan, J. Bonneau, E. Felten, and A. Miller. *Bitcoin and Cryptocurrency Technology (online course)*. 2015. <https://piazza.com/princeton/spring2015/btctech/resources>.
47. S. Park, K. Pietrzak, A. Kwon, J. Alwen, G. Fuchsbauer, and P. Gazi. Spacemint: A cryptocurrency based on proofs of space. *Cryptology ePrint Archive, Report 2015/528*, 2015. <http://eprint.iacr.org/2015/528>.
48. R. Pass, L. Seeman, and abhi shelat. Analysis of the blockchain protocol in asynchronous networks. In *Cryptology ePrint Archive, Report 2016/454*, 2016. <http://eprint.iacr.org/2016/454>.
49. P. Rogaway. Formalizing human ignorance. In P. Q. Nguyen, editor, *Progress in Cryptology - VIETCRYPT 06*, volume 4341 of *LNCS*, pages 211–228. Springer, Heidelberg, Sept. 2006.
50. A. Sahai and B. Waters. How to use indistinguishability obfuscation: deniable encryption, and more. In D. B. Shmoys, editor, *46th ACM STOC*, pages 475–484. ACM Press, May / June 2014.

51. A. Sapirshstein, Y. Sompolinsky, and A. Zohar. Optimal selfish mining strategies in bitcoin. In *FC*, 2016. <http://arxiv.org/abs/1507.06183>.
52. G. Sauer, H. Hochheiser, J. Feng, and J. Lazar. Towards a universally usable CAPTCHA. In *Proceedings of the 4th Symposium on Usable Privacy and Security*, 2008.
53. N. Szabo. Formalizing and securing relationships on public networks. In *First Monday*, 1997. <http://firstmonday.org/ojs/index.php/fm/article/view/548/469>.
54. J. Tam, J. Simsa, S. Hyde, and L. Von Ahn. Breaking audio captchas. *Advances in Neural Information Processing Systems*, 1(4), 2008.
55. L. von Ahn, M. Blum, N. J. Hopper, and J. Langford. CAPTCHA: Using hard AI problems for security. In E. Biham, editor, *EUROCRYPT 2003*, volume 2656 of *LNCS*, pages 294–311. Springer, Heidelberg, May 2003.
56. L. Von Ahn, B. Maurer, C. McMillen, D. Abraham, and M. Blum. reCAPTCHA: Human-based character recognition via web security measures. *Science*, 321(5895):1465–1468, 2008.
57. B. Waters, A. Juels, J. A. Halderman, and E. W. Felten. New client puzzle outsourcing techniques for DoS resistance. In V. Atluri, B. Pfizmann, and P. McDaniel, editors, *ACM CCS 04*, pages 246–256. ACM Press, Oct. 2004.
58. J. Wilkins. Strong CAPTCHA guidelines v1.2. 2009. <http://bitland.net/captcha.pdf>.

A Universal Samplers: Security Definition

Definition 12. Consider efficient algorithms $(\text{Setup}, \text{Sample})$ where $U \leftarrow \text{Setup}^{\text{RO}}(1^\lambda)$, d is the fixed program supporting additional input, and $p_\beta \leftarrow \text{Sample}^{\text{RO}}(U, d, \beta)$. We say $(\text{Setup}, \text{Sample})$ is an adaptively-secure universal sampler scheme for a circuit d , if there exist efficient interactive Turing Machines SimSetup , SimRO such that for every efficient admissible adversary \mathcal{A} , there exists a negligible function $\text{negl}()$ such that the following two conditions hold:

$$\Pr[\text{Real}(1^\lambda) = 1] - \Pr[\text{Ideal}(1^\lambda) = 1] = \text{negl}() \text{ and } \Pr[\text{Ideal}(1^\lambda) = \text{aborts}] < \text{negl}()$$

where admissible adversaries, the experiments **Real** and **Ideal** and the notion of the **Ideal** experiment aborting, are described below

- An admissible adversary \mathcal{A} is an efficient interactive Turing Machine that outputs one bit, with the following input/output behavior:
 - \mathcal{A} initially takes input security parameter 1^λ and sampler parameters U , as well as the program d .
 - \mathcal{A} can send a message (RO, x) corresponding to a random oracle query. In response, \mathcal{A} expects to receive the output of the random oracle on input x .
 - \mathcal{A} can send a message (sample, β) . The adversary does not expect any response to this message. Instead, upon sending this message, \mathcal{A} is required to honestly compute $p_\beta = \text{Sample}(U, d, \beta)$, making use of any additional RO queries, and \mathcal{A} appends (β, p_β) to an auxiliary tape.

Remark. Intuitively, (sample, β) messages correspond to an honest party seeking a sample generated by the fixed program d on input β . Recall that \mathcal{A} is meant to internalize the behavior of honest parties.

- The experiment **Real**(1^λ) is as follows:
 - Throughout this experiment, a random oracle **RO** is implemented by assigning random outputs to each unique query made to **RO**.
 - $U \leftarrow \text{Setup}^{\text{RO}}(1^\lambda)$.
 - $\mathcal{A}(1^\lambda, U, d)$ is executed; when \mathcal{A} sends every message of the form (RO, x) , it receives the response $\text{RO}(x)$.
 - The output of the experiment is the final output of the execution of \mathcal{A} (which is a bit $b \in \{0, 1\}$).
- The experiment **Ideal**(1^λ) is as follows:
 - Throughout this experiment, a Samples Oracle \mathcal{O} is implemented as follows: On input β , \mathcal{O} outputs $d(F(\beta))$, where F is a truly random function.
 - $(U, \tau) \leftarrow \text{SimSetup}(1^\lambda)$. Here, **SimSetup** can make arbitrary queries to the Samples Oracle \mathcal{O} .
 - $\mathcal{A}(1^\lambda, U, d)$ and **SimRO**(τ) begin simultaneous execution. Messages for \mathcal{A} or **SimRO** are handled as:
 1. Whenever \mathcal{A} sends a message of the form (RO, x) , this is forwarded to **SimRO**, which produces a response to be sent back to \mathcal{A} .
 2. **SimRO** can make any number of queries to the Samples Oracle \mathcal{O} .
 3. In addition, after \mathcal{A} sends messages of the form (sample, β) , the auxiliary tape of \mathcal{A} is examined until \mathcal{A} adds entries of the form (β, p_β) to it. At this point, if $p_\beta \neq d(F(\beta))$, the experiment aborts and we say that an “Honest Sample Violation” has occurred. Note that this is the only way that the experiment **Ideal** can abort. In this case, if the adversary itself “aborts”, we consider this to be an output of zero by the adversary, not an abort of the experiment itself.
 - The output of the experiment is the final output of the execution of \mathcal{A} (which is a bit $b \in \{0, 1\}$).