

Applying Private Information Retrieval to Lightweight Bitcoin Clients

Kaihua Qin*, Henryk Hadass*, Arthur Gervais* and Joel Reardon†

*Department of Computing, Imperial College London

Email: kaihua.qin@imperial.ac.uk, henryk.hadass14@imperial.ac.uk, a.gervais@imperial.ac.uk

†Department of Computer Science, University of Calgary

Email: joel.reardon@ucalgary.ca

Abstract—Lightweight Bitcoin clients execute a Simple Payment Verification (SPV) protocol to verify the validity of transactions related to a particular user. Currently, lightweight clients use Bloom filters to significantly reduce the amount of bandwidth required to validate a particular transaction. This is despite the fact that research has shown that Bloom filters are insufficient at preserving the privacy of clients' queries.

In this paper we describe our design of an SPV protocol that leverages Private Information Retrieval (PIR) to create fully private and performant queries. We show that our protocol has a low bandwidth and latency cost; properties that make our protocol a viable alternative for lightweight Bitcoin clients and other cryptocurrencies with a similar SPV model. In contrast to Bloom filters, our PIR-based approach offers deterministic privacy to the user.

Among our results, we show that in the worst case, clients would like to verify 100 transactions occurring in the past week incurs a bandwidth cost of 33.54 MB with an associated latency of approximately 4.8 minutes, when using our protocol. The same query executed using the Bloom-filter-based SPV protocol incurs a bandwidth cost of 12.85 MB; this is a modest overhead considering the privacy guarantees it provides.

Index Terms—Private Information Retrieval; Bitcoin; Simple Payment Verification;

I. INTRODUCTION

Bitcoin [1], a pseudonymous cryptocurrency created in 2008, enables users to perform irreversible electronic transactions. Once a transaction occurs, the Bitcoin peer-to-peer network records this fact in a replicated and public database called the blockchain. The clients that verify and store every transaction in the network are called full nodes. This verification process is resource intensive, requiring hardware and storage nearly exceeding personal computers. At the time of writing, the Bitcoin blockchain was approximately 210 GB in size, monotonically growing every day.

Resource constrained devices, such as mobile phones, are unable to store the entire contents of the Bitcoin blockchain given its prohibitive size¹. Moreover, bandwidth usage on mobile phones is usually metered, which is another factor preventing full nodes from being run on such devices. These limitations, combined with the scale of the blockchain, prevents participants in the system from being able to validate their

own transactions on their mobile phones. Bitcoin, however, was aimed at completely removing such trust assumptions: no central third-party should be required for it to work properly. This suggests that users need to run their own full nodes at home in order to participate in the network.

This drawback was already anticipated when Bitcoin was first introduced [1]. Nakamoto suggests that a lightweight version of the full node protocol should be used in resource constrained environments, which was called Simple Payment Verification (SPV). SPV simply verifies transactions which are relevant to a specific user, in contrast to a full node which verifies every single transaction that has ever occurred.

A simple implementation of SPV requires the client to download the whole Bitcoin blockchain, much like a full node, in order to verify a single transaction. This approach preserves the privacy of the client since no other network participant can determine which transaction the client is interested in verifying. However, this approach wastes bandwidth since the majority of the downloaded data is discarded. We call this solution Naive SPV.

SPV as described in the Bitcoin white-paper was only realised in 2012 when a Bloom-filter-based SPV protocol was proposed in Bitcoin Improvement Proposal number 37 (BIP-37) [2]. This proposal allowed the SPV protocol to significantly reduce the amount of bandwidth used in verifying a particular transaction, when compared with Naive SPV. This resulted in BIP-37 becoming the de-facto standard in the Bitcoin community for lightweight clients, with mobile cryptocurrency wallet applications such as BitcoinJ [3] implementing it. We call this solution BIP-37 SPV.

Using Bloom filters, however, does not maintain the privacy of a user's queries. This has been acknowledged not only by the authors of BIP-37 [4], but also by the Bitcoin community [5] and by researchers in this field such as Gervais et al. [6]. Bloom filters present a trade-off between bandwidth and privacy. In the existing Bitcoin wallet implementations, almost no privacy is guaranteed because a low false positive rate is chosen for the sake of bandwidth. Bloom filters are moreover vulnerable to intersection attacks even though sacrificing bandwidth and increasing the false positive rate, which means that the adversary can read the user's interest by intersecting the results from different Bloom filters of the same wallet. Gervais et al. [6] offer solutions to improve the

¹At the time of writing, an iPhone 8 with 256 GB of storage costs \$749, a Samsung Galaxy Note 9 with 512 GB of storage costs \$1,249.99 and a SanDisk Ultra 400 GB microSD card alone costs \$175.13.

privacy provisions of Bloom filters. Nevertheless, we would like to provide the user with full and deterministic privacy.

In a public payments network such as Bitcoin, privacy is a fundamental requirement. Users who do not adequately preserve their privacy allow full node peers to discover the Bitcoin addresses they own and which transactions they have been involved in. Upon discovering this information, adversarial peers can engage in a denial of service to users or addresses which they disfavor. In addition, these peers can engage in a trivial process of tracing the funds of a particular user, given that the Bitcoin blockchain is public. This tracing could expose an individual to personal harm if it was discovered that they have access to a large quantity of Bitcoin. Due to these risks, it is highly desirable that Bitcoin's users preserve their privacy.

Summary: In this paper we describe our design of a Private Information Retrieval (PIR) based SPV client to demonstrate that a fully private, low bandwidth and low latency lightweight Bitcoin client is feasible. The contributions of this work are the following:

- **First PIR system for blockchain:** To the best of our knowledge, we are the first to design a PIR scheme to increase the privacy properties of cryptocurrency based systems.
- **Different PIR schemes:** We show that in the single-server setting our protocol is as efficient as BIP-37 SPV, while in the multi-server setting it has a comparable bandwidth cost.
- **Bandwidth comparison:** We provide a detailed analysis of the bandwidth cost of our novel protocol and compare its cost to the Naive and BIP-37 SPV protocols. We show that in both instances our protocol is far more bandwidth efficient than the Naive SPV protocol.
- **Practical Latency:** We also show that the latency of executing queries in our protocol is tolerable, with clients being able to perform fully private SPV in far less than a minute when verifying recent transactions.

This work is organized as follows. Section II covers the necessary background. Section III provides an overview of our system and Section IV outlines the details. In Section V, we present our evaluation. Section VI covers future work which improves our implementation and in Section VII, we discuss related work in this field. Finally, we conclude in Section VIII.

II. BACKGROUND

A. Bitcoin

Other than full nodes and lightweight clients, Bitcoin has another type of network participant called miners. The role of miners is to secure the Bitcoin blockchain through the execution of a Proof-of-Work (PoW) mechanism.

Blocks, which contain one or more transactions, have the PoW mechanism executed on them by miners. Blocks are said to be mined once a solution to the PoW mechanism has been found. The first mined block is called the Genesis block. Each block contains a block header, which is a lightweight summary of the contents of the block and is 80 bytes in size. Each

Bitcoin transaction consist of one or more inputs and outputs, with the inputs describing where the Bitcoin is coming from and the outputs where it is going to. Each transaction has a unique identifier called a transaction ID (TXID) which is obtained by applying a hash function to the transaction data.

The TXIDs of all of the transactions contained in a block are used to construct a Merkle [7] tree, whose root is included in the block header. This Merkle tree root serves as a unique identifier for all of the transactions included in a particular block. Since the outputs of a transaction state where the Bitcoin is sent to, spending Bitcoin means using the outputs of some transaction (for which you were the beneficiary) as the inputs to a new transaction. Since Bitcoin prevents double spending of the currency, this divides outputs into two categories: those that have been spent—i.e., used as inputs to a subsequent transaction—and those that are unspent. The latter are called Unspent Transaction Outputs (UTXOs) and they sum to the total amount of Bitcoin in existence.

UTXOs represent the concept of “having bitcoin”—insofar that one can sign value transactions that include the UTXOs as inputs. Each transaction output includes a Bitcoin address of which there are a number of different types, the most common of which are Pay-to-Public-Key-Hash (P2PKH) and Pay-to-Script-Hash (P2SH) [8]. P2PKH are addresses owned by a single entity whereas P2SH can be multi-signature addresses. These addresses are usually encoded in the hexadecimal base58 format. Base58 removes characters from its encoding set which are similar to each other, such as 0, O, I and l, to reduce the risk of users mixing up distinct characters with similar glyphs.

B. Simple Payment Verification (SPV)

Three items are required to perform SPV on a transaction. The transaction itself, the list of TXIDs from the corresponding block and a list of all the block headers from the Genesis block, up to and including the block in which the transaction of interest is included. The transaction is used to calculate its TXID, to check that this TXID exists in the list of TXIDs from the relevant block and to confirm that an addresses belonging to the client is referenced in an output. The list of TXIDs is used to construct the Merkle tree and subsequently calculate the Merkle tree root and check if it matches the value stored in the corresponding block header. The list of block headers is used to prove that the header belonging to the block of interest can be placed in a valid location in the Bitcoin blockchain by recursively hashing them and confirming their corresponding Proofs-of-Work. Once the validity of a transaction has been determined, using the three items outlined above, confidence in the irreversibility of that transaction needs to be established. This is done by ensuring that the block containing the transaction of interest is embedded by normally 1 to 6 subsequently mined blocks depending on the transaction value [9].

A naive implementation of an SPV client would first establish one or more connections to a set of full nodes. Next, in order to verify a transaction, the client would request for every single block of transactions from the Genesis block up

to the block in which the transaction is included, from the peers it is connected to. For each received block the client would then select the data detailed above and discard the rest. The advantage of this approach is that the full node peers are not able to determine which particular transaction the client is interested in verifying, thus preserving the user's privacy. This is because the SPV client needs to download essentially the whole blockchain for each transaction which it would like to verify. However, this means that this approach is bandwidth intensive since the majority of the downloaded data is not used and is discarded.

C. BIP-37: Bloom filters & Merkle blocks

In order to improve bandwidth efficiency over the Naive SPV protocol, the BIP-37 SPV protocol introduced Bloom filters and a corresponding Merkle block to significantly reduce the bandwidth required to verify a particular transaction. We now explain both of these constructions.

A Bloom filter is a space efficient, probabilistic data structure that allows for the testing of an elements' membership in a set. Bloom filters work probabilistically and, by design, only have false positives (FP) but not false negatives (FN). This means that Bloom filters reliably tell when an element is *definitely* not in a set, but only offer probabilistic guidance as to whether an element may be in a set. Adding an element to a set and testing for membership are both constant time operations. Bloom filters that have too many elements contained within them are said to be "too full", which renders them unusable due to their excessively high FP rate.

Bloom filters work by using a small number of (non-cryptographic) hash functions that map the set element to a random bit position in the filter. So if you have five hash functions, there are five positions in the filter for each element, and this set of positions is characteristic to each element. To insert an element into the set, one simply makes all the element's characteristic positions to one (i.e., "on"). Thus, to see if an element is in the set, one simply checks if all the characteristic positions are "on": if any are "off" (i.e., not set to 1 but instead still 0) then we are certain that the element was never inserted. If all are on, however, it may be due to a coincidence from the other elements that were inserted. Therefore, there are no false negatives but there are false positives.

When using Bloom filters in the SPV setting, the FP rate is used as a proxy for the privacy level of the Bloom filter. An SPV client with ample bandwidth may choose to have a high FP rate. This means that the full node peer from which data is downloaded to verify a transaction, cannot accurately determine which transaction the SPV client is interested in verifying. Bloom filters with an exceedingly high FP rate will download as much data, and have the same privacy provisions, as Naive SPV. This is because such a Bloom filter would match all addresses and in effect download the entire blockchain.

In contrast, an SPV client with access to a minimal amount of bandwidth would create an accurate Bloom filter, by setting

a low FP rate. This would mean that the full node peer would know exactly which transactions the SPV client is interested in verifying, since only a specific set of data would be downloaded. As such, Bloom filters represent a trade-off, which is configurable by the user, between privacy—the precision of the data returned from a full node peer—and bandwidth.

Once a Bloom filter has been created, it is used in the construction of the corresponding Merkle block. This consists of a block header and a partial Merkle tree, called a Merkle branch. The Merkle branch consists of all the TXIDs for whom set membership passed in the Bloom filter—and which therefore includes false positives—as well as intermediate Merkle tree hashes, including the Merkle tree root. Altogether, this allows the SPV client to connect these TXIDs together and to verify the value of the Merkle tree root located in the block header. A Merkle block is followed by a list of transactions which were set in the Bloom filter, and whose TXIDs were already included in the Merkle branch.

In practice, users tend to use low FP rates when constructing Bloom filters [5], thus maintaining virtually no privacy. The reason low FP rates are used is because they allow faster synchronization for SPV clients. SPV synchronisation is a process by which a users' wallet is brought up to date with the current state of the blockchain by reflecting in the wallets' balance the final result of the transactions in which the user was involved. Increasing the FP rate leads to longer synchronization times, since more data needs to be downloaded. This has been commonly observed [4] to have a negative effect on users' satisfaction with SPV client implementations such as BitcoinJ [3]. This occurs because users are more concerned about the overall performance of their wallet application, rather than their privacy.

D. Private Information Retrieval

Private Information Retrieval (PIR) allows users to query a database, such that the database learns nothing about the users' query. The trivial solution to this would be for the client to download the entire database and perform the query offline. Over the years more efficient PIR protocols have emerged which have significantly improved on this upper bound. We provide a brief overview below.

There exist two classes of PIR: Information Theoretic PIR (IT-PIR) and Computational PIR (C-PIR). IT-PIR protocols involve replicating the database among a set of servers. The client makes different queries to different servers and determines the result of their query from the set of server responses. IT-PIR is information theoretically secure, provided that the different servers do not collude with each other, up to a certain threshold. C-PIR protocols, in contrast, require only one server to store the database and query privacy is guaranteed by cryptographic means.

Each class of PIR has its own advantages and disadvantages. The advantages of multi-server IT-PIR are that it generally incurs smaller communication and computation costs. This is because IT-PIR treats queries as vectors, and databases

as matrices, on which linear algebra operations of vector-by-matrix multiplication are performed. As such, the size of the database should be square in order for the communication cost between client and server to be optimal. IT-PIR is also robust to missing or incorrect database server responses.

The disadvantage of IT-PIR is that the threshold of non-colluding servers needs to be maintained and it is not clear how this requirement can be enforced in practice, particularly because a database server can mount a Sybil attack.

The advantage of single-server C-PIR is that it relies only on a single server, however, due to this property, C-PIR is not robust since it cannot overcome missing or incorrect database server responses.

C-PIR is generally computationally slower when compared to IT-PIR. This is because greater computation is required in encoding clients' queries and because the server is performing matrix-by-matrix multiplication—in contrast to vector-by-matrix multiplication in some IT-PIR schemes. Another disadvantage of some C-PIR protocols is that they are based on lattice problems whose security provisions are not well understood. Because of this, clients may have doubts about the privacy of their queries [10].

In our implementation we use a hybrid PIR scheme by Devet and Goldberg [10] which combines the advantages of IT-PIR and C-PIR, while minimising the disadvantages of either one. The underlying IT-PIR protocol in this hybrid scheme is by Devet et al. [11] which leverages Shamir Secret Sharing and Reed-Solomon decoding to create a robust and Byzantine-fault-tolerant IT-PIR protocol which can withstand up to t colluding servers (t is called the *privacy level*). The C-PIR scheme used is by Aguilar-Melchor and Gaborit [12].

The hybrid PIR scheme treats the underlying data as elements of the finite field $GF(2^8)$. This protocol can be found in the Percy++ [13] open source library which implements a number of state-of-the-art PIR protocols.

III. SYSTEM OVERVIEW

In this section we describe our system design, its properties and considered adversarial models. We denote the IT-PIR protocol as Φ and the C-PIR protocol as Θ in our hybrid scheme.

A. System Model

Our system features the following entities:

- **Ledger:** We assume the existence of an immutable ledger, acting as coarse-grained timestamping service. The ledger's purpose is to discern not yet spent electronic coins (commonly referred to as unspent transaction outputs or short UTXO).
- **User:** A user possesses one or more private keys in the ledger and intends to validate the transactions towards the accounts (e.g., Bitcoin addresses) controlled by the private keys through SPV.
- **PIR server:** The PIR server responds to the PIR queries from the SPV clients and provides the data required to perform the SPV.

We assume users can communicate with PIR servers through the underlying network.

B. High-Level Operation

In this section, we outline the high-level operations of our system (cf. Figure 1).

- PIR servers download the whole blockchain and construct PIR databases. For each database, the PIR server creates a description file called *manifest file*.
- As outlined in Section II, the user collects all available block headers from e.g., full node peers.
- The user fetches the manifest files from the PIR servers to efficiently query the PIR database afterwards.
- The user executes the PIR-SPV protocol. PIR queries can be run recursively (i.e., in one PIR query, the user encodes and exchanges several query messages before receiving the final PIR response). The user decodes the responses and then performs SPV validation (verifying the partial Merkle tree of the received transactions given the downloaded block headers).

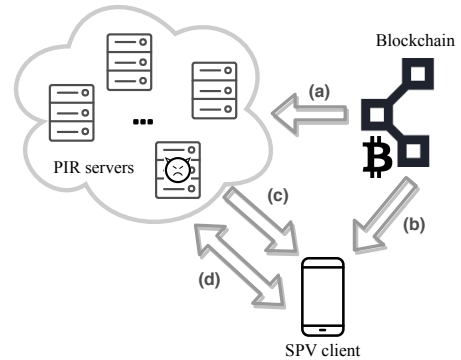


Fig. 1: High-level operations of our system: (a) PIR servers download the blockchain; (b) User keeps up to date with all available block headers; (c) User fetches the manifest files from the PIR servers; (d) User executes the PIR-SPV protocol.

C. Main Properties

Our system provides the following properties:

- **User:** Given our system, users can query their transactions from other blockchain nodes, without disclosing the transactions that they are interested in. We call this property *transaction query privacy*.
- **PIR server:** We assume that different PIR servers create equal databases and manifest files, given the same blockchain source. While serving user requests, the PIR servers do not learn any user-centric transaction data.

D. Adversarial Model

The privacy offered by the hybrid PIR scheme relies on the fact that both schemes Φ and Θ preserve the query privacy. For Φ , we assume that k out of ℓ servers respond per PIR query. We then define the privacy level t and the number of

Byzantine servers v , s.t. $v < k - t - 1$. We assume that no more than t servers are colluding to discover the contents of a users query, which ensures the privacy of Φ . This implies that we assume not all PIR servers are colluding in the case that all servers are responsive and there is no Byzantine server. For Θ , we assume that malicious PIR servers are computationally bound, which guarantees that disclosing user’s interest from a Θ PIR query is infeasible.

IV. SYSTEM DETAILS

In this section, we detail the architecture of our system. We begin by describing how PIR servers build the databases to facilitate PIR-based queries. Then, we describe the manifest files, which are used by clients to construct PIR-based queries. Finally, we depict our PIR-based SPV protocol.

A. Database Construction

1) *Database Content Structure*: We split the data required to perform SVP into three distinct components (i.e., the address PIR DB, the Merkle tree PIR DB and the transaction PIR DB), the formats of which are respectively detailed below.

a) *Address PIR DB*: In our implementation, we take a user-centric, address-first approach, which means that a user needs to select a Bitcoin address that belongs to them with which they would like to execute our protocol. Figure 2 illustrates the structure of a single entry in a row in the Address PIR database. A row is made up of one or more of these entries.

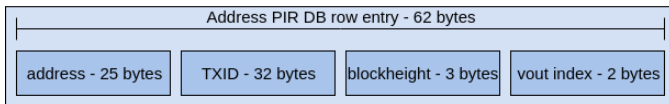


Fig. 2: Structure of a single Address PIR DB row

An explanation of each field follows:

- **Address**: This field holds a particular Bitcoin address can be involved in more than one transaction. As such, there can be multiple entries with the same address. This field serves as the entry point to our PIR SPV protocol.
- **TXID**: This field identifies the transaction whose output address the address field refers to. This field is used to query the Transaction PIR database.
- **Block height**: This field holds the block height of the block which includes the transaction the TXID field refers to. This field is used to query the Merkle Tree PIR database.
- **Vout Index**: This field holds the index value of the output of the transaction referenced by the TXID field, in which the address referenced by the address field was used. This field serves more of a convenience than a necessity, by allowing clients to quickly identify the location of their output in the relevant transaction, without having to look through every transaction output.

b) *Merkle Tree PIR DB*: This database contains the TXIDs which are used to calculate the Merkle tree root in the SPV protocol. The structure of a single row in the Merkle Tree PIR database is shown in Figure 3, where each TXID is 32 bytes in length.

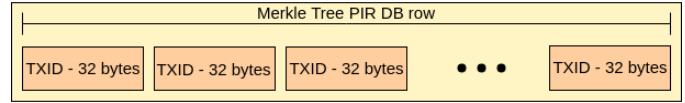


Fig. 3: Structure of a single Merkle Tree PIR DB row

c) *Transaction PIR DB*: This database contains the hexadecimal format of transactions in the UTXO set. Figure 4 illustrates the structure of a single row in the Transaction PIR database, where each row consists of transaction bytes.

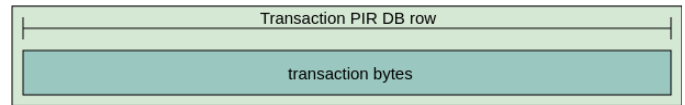


Fig. 4: Structure of a single Transaction PIR DB row

2) *Temporal Partitioning of Databases*: We explain our rationale for partitioning the databases across several time periods in this section.

The collected data are split into three distinct time periods, with associated manifest files: weekly, monthly and all-time. Since a new block is generated approximately every 10 minutes, weekly data consists of the most-recent 1 008 collected blocks, monthly data the preceding 4 032 blocks, and all-time data consisted of all remaining blocks—up to and including the Genesis block.

The temporal partitioning of the collected data means that enormous amounts of data do not usually need to be searched using PIR. This leaks the fact that a client is interested in, for example, more-recent data whenever the weekly databases are queried. This can be assumed to be true, however, given the knowledge that an entity is using Bitcoin; more importantly, simply accessing recent transactions does not imply one has recent transactions. The important piece of privacy to maintain is the mapping of a client to the set of addresses that they control. This is maintained implicitly since the client can only query the Address PIR DB under PIR.

Clients who are interested in verifying a transaction for a particular range of time simply issue their queries to the corresponding PIR database. This structure allows our static implementation to be easily transformed into one which is periodically updated whenever a new block is generated. This is outlined in Appendix C.

The other advantage of this partitioning is that it reduces the bandwidth and latency cost for clients who are interested in verifying more recent transactions.

3) *Database Dimensions*: In this section, we describe how the height and width of the databases are determined.

Several factors influence the appropriateness of the database dimensions when using PIR. PIR seeks to minimize the

communication cost, which is the bandwidth required to obtain one row of data from any of our databases. The PIR system that we use has an optimal communication cost when the “height” and “width” of the database are approximately equal. That is, the number of rows is proportional to the size of each row. This is because requests in PIR have a constant cost per row while the PIR replies are the size of a single row. As such, we tried to ensure that our databases were square-like when possible. Larger sets of data were placed into rectangular databases to avoid *bandwidth bloat*. Bandwidth bloat is when individual rows in a database contain a greater proportion of irrelevant information with respect to the clients query.

It should be noted that larger databases increase the latency of PIR queries, as they are executed over more data which is computationally intensive. In most cases the database size was slightly larger than the set of data being stored in it, resulting in the remaining space being padded with strings of zeros.

4) *Data Ordering*: Once the dimensions of the databases are determined, the collected data have to be ordered, as described below, before being placed into these databases.

For Address PIR DB data, each entry is lexicographically sorted by the address field. For Merkle Tree PIR DB data, each list of TXIDs is sorted in ascending order according to the associated block height. No ordering is imposed on Transaction PIR DB data.

After the sorting occurred of the respective sets of data, it is then arranged in row-oriented order in the databases. This is where a row is filled from left to right with entries before the next row is filled. This localises the entries, thus reducing the total number of queries that clients need to perform.

B. Database Manifest Files

In the PIR scheme we utilised, users are required to grasp some preliminary knowledge of the queried database (e.g., database dimensions and the position of the desired data) before making a request. Consequently, for each database, a manifest file is generated which provides sufficient information to bootstrap a PIR query. A manifest file consists of one or more records, each of which indicates the location of an individual piece of data. Users request from PIR servers to fetch the manifest files before making the queries. As manifest files are the global descriptions of the databases, the request does not reveal any user’s interest other than participating in Bitcoin. Note, in Section VI, we propose a solution that saves the trouble of downloading the whole manifest file replaced by iterative PIR queries to reduce bandwidth cost. We concretely describe the format of the manifest files in Appendix A.

C. PIR-Based SPV Protocol

We describe our PIR-based SPV protocol which facilitates private lightweight Bitcoin clients in Algorithm 1. The client is required to obtain the available block headers independently. Execution of the PIR-SPV protocol is divided into three rounds of queries to three PIR databases in the proper order (cf. Algorithm 1). Once completing the queries, the user can perform the SPV validation of the selected transaction with the transaction content, TXIDs and the block headers.

Algorithm 1: PIR-SPV protocol

Data: SPV client \mathcal{C} , one or multiple PIR servers \mathcal{S}

Result: \mathcal{C} obtains the necessary data for a SPV privately

Initialization: \mathcal{S} constructs the PIR databases and associated manifest files; \mathcal{C} downloads the manifest files from \mathcal{S} ;

- a.1 \mathcal{C} selects an address to fetch a record from the Address PIR DB manifest file and generates the PIR queries based on row indices of the selected record;
 - a.2 \mathcal{S} computes the result using the PIR queries on the Address PIR DB;
 - a.3 \mathcal{C} parses and decodes the result to obtain one or more Address PIR DB entries;
 - b.1 \mathcal{C} uses the value of the block height field of an entry to fetch the corresponding record from the Merkle Tree PIR DB manifest file and generates the PIR queries based on the row indices;
 - b.2 \mathcal{S} computes the result using the PIR queries on the Merkle Tree PIR DB;
 - b.3 \mathcal{C} parses and decodes the result to obtain the requested list of TXIDs;
 - c.1 \mathcal{C} uses the value of the TXID field from the same entry that was selected in step b.1, to fetch the corresponding record from the Transaction PIR DB manifest file and generates the PIR queries based on the row indices;
 - c.2 \mathcal{S} computes the result using the PIR queries on the Transaction PIR DB;
 - c.3 \mathcal{C} parses and decodes the result to obtain the requested transaction.
-

V. EVALUATION

We collected the necessary data between February and April 2018 from 9 virtual machines (VMs) running the core Bitcoin client. Each VM ran a Ubuntu 16.04 (64-bit) OS, had four 1 GHz CPUs, 8 GB of RAM and a 300 GB hard drive. For simplicity, in our implementation, we only considered the UTXO set of P2PKH addresses. Note the Bitcoin Core supports a descriptor scheme [14] to generally express different types of outputs. We therefore reasonably deduce the feasibility of extending our solution to cover various address types.

Table I summarises the sizes of the individual components of the system. We use *database width* to denote the number of records stored in a row. For example, the all-time data for the Address PIR DB states that a single row contains 906 entries, each of which is 62 bytes in size (cf. Section IV). Because the data fields vary across the different databases, the exact number of bytes that a database row consumes therefore varies. We detail the process of determining the dimensions in Appendix B.

A. Benchmark Methodology

We executed our PIR-based SPV protocol and compared its bandwidth cost with the Bloom-filter-based BIP-37 SPV protocol and with the fully private Naive SPV protocol. Since

TABLE I: The sizes and dimensions of the generated databases and their corresponding manifest files

		All-Time blocks 1 to 508462	Monthly blocks 508463 to 512494	Weekly blocks 512495 to 513502
Address PIR DBs	# entries per row (width)	906	214	124
	# rows (height)	56172	13268	7688
	Total size of DB	3.16 GB	176.04 MB	59.11 MB
	Size of manifest file	65.99 MB	175.33 MB	66.97 MB
Merkle Tree PIR DBs	# entries per row (width)	821	1196	1184
	# rows (height)	394080	38272	37888
	Total size of DB	10.18 GB	1.46 GB	1.44 GB
	Size of manifest file	40.90 MB	0.32 MB	78.62 KB
Transaction PIR DBs	Length of row in bytes (width)	758	848	876
	# rows (height)	20942782	1537424	512460
	Total size of DB	15.87 GB	1.30 GB	448.91 MB
	Size of manifest file	3.03 GB	218.68 MB	72.45 MB

we use Devet et al.’s [10] hybrid PIR scheme, the client can query one or more PIR servers when executing our protocol. The experiments were executed locally on a single machine, running one or more PIR server instances. This machine had a Ubuntu 16.04 (64-bit) OS, Intel Core i7 3.4 GHz CPU, 16 GB of RAM and a 512 GB hard drive. To ensure a fair and accurate comparison, we did not include the cost of downloading the block headers, since this is a common element of all three protocols. We did not include the bandwidth cost of client-side queries for BIP-37 and Naive SPV since these are contained in peer-to-peer layer messages. We did, however, include the bandwidth cost of client-side queries for PIR SPV since that is a fundamental element of performing PIR. The other steps that we followed when collecting our data were:

- **PIR SPV:** For every TXID located in each entry in the Address PIR DBs, we executed the PIR protocol as outlined in Section IV and used the total length of the returned rows, from one or more PIR servers, as part of the bandwidth cost. The downloading of manifest files is not included in our calculation of bandwidth cost because in Section VI we present a feasible solution which removes the need for clients to download them. Instead, clients perform interpolation search using iterative PIR queries to retrieve a particular record, which we imagine will have a small bandwidth cost.

Two sets of data were collected for PIR SPV. The first set of data used only one PIR server to simulate the case of single-server C-PIR being used by the hybrid PIR schema. The second set of data used three PIR servers to simulate the case of multi-server IT-PIR being used by the hybrid PIR schema. The bandwidth cost of using additional servers is linear, assuming that no Byzantine PIR servers exist. It should be noted that under C-PIR, a recursive depth of size 1 was used, as suggested by Aguilar-Melchor and Gaborit [12].

- **BIP-37 SPV:** For every TXID located in each entry in the Address PIR DBs, a Bloom filter was constructed and then used to send a Bitcoin peer-to-peer layer message to a full node hosted on one of our VMs. The length of the reply, which contained the Merkle block and the

associated transactions, was measured and used as the bandwidth cost. It should be noted that this reply also contained a network message header which was 24 bytes in length [15]. This value was deducted from the final bandwidth cost calculations since this does not exist in the other two protocols.

- **Naive SPV:** For every TXID located in each entry in the Address PIR DBs, we measured the total amount of bandwidth that had to be used to download the blocks between the Genesis block and the block in which the transaction of interest was included. We used a single full node hosted on one of our VMs.

B. Results and Analysis

The results are recorded in Figure 5, which shows three histograms displaying the bandwidth cost required to verify a single transaction for each of the different SPV protocols, across all-time, monthly and weekly data. These results are summarised in Table II, which shows the expectation and standard deviation of the collected data.

Figure 6 shows the cumulative bandwidth cost required to verify up to 100 transactions by the different SPV protocols, across the three time periods. The corresponding latency cost of executing our PIR SPV protocol in the single- and multi-server setting is also shown. A small sample of these results is shown in Table III. In Figure 6, to calculate the bandwidth cost of verifying one or more transactions through SPV, steps similar to those described in Section V-A were followed with some adjustments. For example, to calculate the total bandwidth cost of verifying 7 transactions, a random entry from the Address PIR DB was picked 7 times, and its TXID used to perform Naive, BIP-37 and PIR-based SPV. The bandwidth cost for each TXID was calculated as described in Section V-A and then the results of the 7 calculations were totaled. We did this five times and then an average was taken. This average is used as the final result.

It should be noted that when calculating the bandwidth cost of Bloom-filter-based SPV, a fresh filter was used on each TXID, rather than a single filter encoding the set of 7 TXIDs. This is because a Bloom filter gets full very quickly when it is used to match a large number of TXIDs. A full Bloom filter

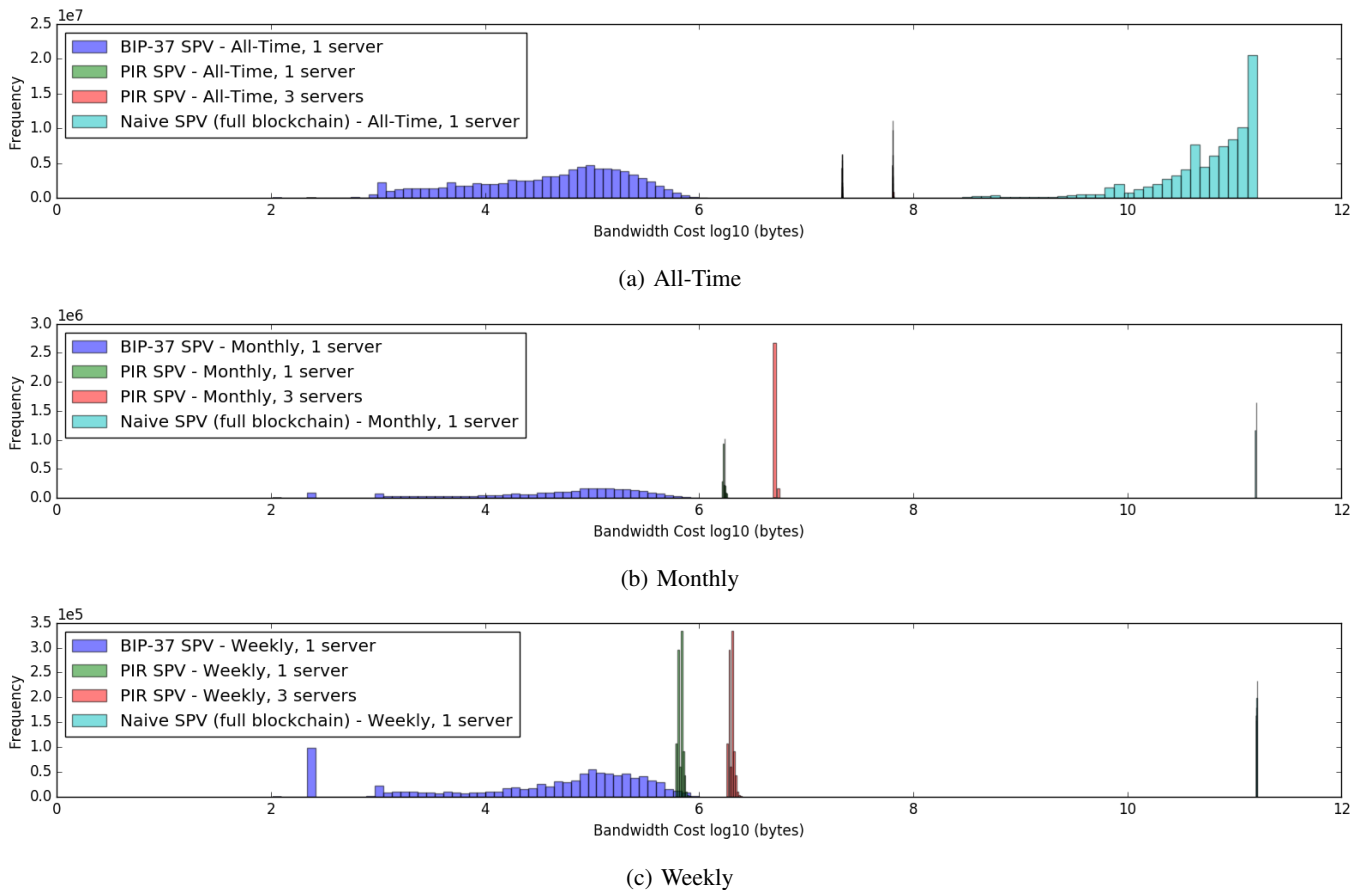


Fig. 5: Histograms showing the bandwidth cost of verifying a single TXID by BIP-37 SPV, single- and multi-server PIR SPV and Naive SPV, for the three time periods.

TABLE II: Summary of Figure 5, showing the expectation and standard deviation of the bandwidth cost for the different protocols, for the three time periods.

		BIP-37	PIR: 1 server	PIR: 3 servers	Naive
All-Time	expectation	102.80 KB	21.54 MB	64.61 MB	80.27 GB
	std. dev.	130.25 KB	35.43 KB	106.29 KB	50.24 GB
Monthly	expectation	132.43 KB	1.70 MB	5.11 MB	159.19 GB
	std. dev.	149.14 KB	34.00 KB	102.01 KB	1.12 GB
Weekly	expectation	128.52 KB	666.07 KB	2.00 MB	161.47 GB
	std. dev.	155.19 KB	34.00 KB	102.66 KB	292.08 MB

is one where most of the entries are set to “on”. Using a fresh filter for each TXID would mean that we would obtain a more accurate bandwidth cost.

The latency was measured by executing the hybrid PIR protocol from Percy++ [13] on the average number of rows needed to verify a set of TXIDs, for example 7, for each of the Address, Merkle Tree and Transaction PIR databases. Figure 6 shows the latency of a single server, for both the C-PIR and IT-PIR cases, because even though 3 servers were used to facilitate IT-PIR, the query was sent in parallel to each server rather than sequentially.

In the weekly and monthly case, Figure 6 shows that single-server PIR SPV has a similar bandwidth cost to BIP-37 SPV

when verifying between 20 and 100 transactions. In particular, Table III shows that if a client is interested in performing single-server PIR SPV on 100 transactions which occurred in the past week, it will take approximately 2 minutes for this query to be executed, with a bandwidth cost of 11.18 MB. In contrast, BIP-37 SPV will require 12.85 MB in the same scenario.

Clients who are interested in a smaller number of transactions for the same time period, such as 20, would incur a bandwidth cost of 8.31 MB in the multi-server setting, with a latency of approximately 62 seconds. In the monthly case, a similar query would incur a bandwidth cost of 12.20 MB with a latency of approximately 2.3 minutes, while in the all-time

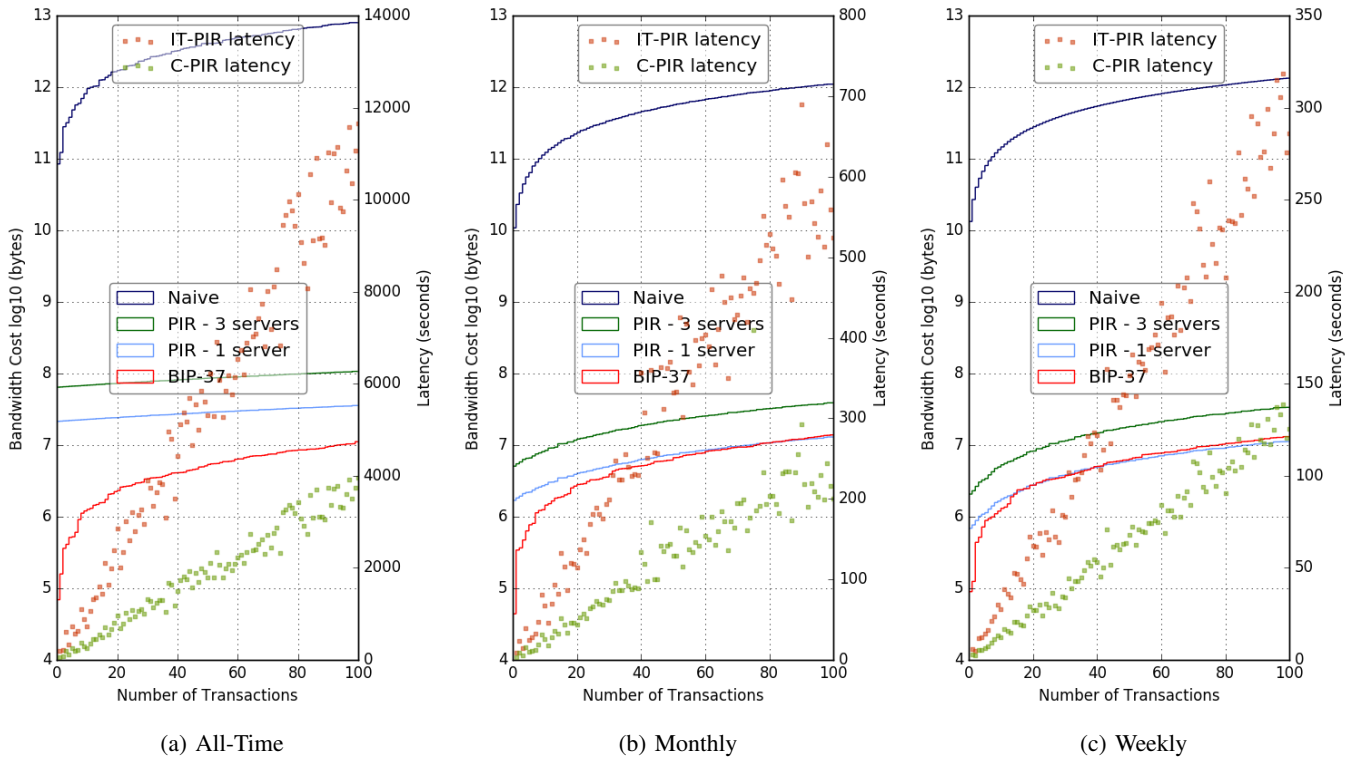


Fig. 6: Cumulative distribution of bandwidth cost for the verification of transactions under BIP-37 SPV, single- and multi-server PIR SPV and Naive SPV, across the three time periods. The latency of single- and multi-server PIR SPV is also displayed.

TABLE III: A small sample of results from Figure 6, showing the bandwidth and latency cost for a particular number of transactions under different SPV protocols, across the three time periods.

Number of Transactions		BIP-37	PIR: 1 server	PIR: 3 servers	Naive	C-PIR latency	IT-PIR latency
1	All-Time	69.32 KB	21.51 MB	64.53 MB	84.85 GB	68.44 s	203.52 s
	Monthly	44.08 KB	1.70 MB	5.09 MB	10.86 GB	3.89 s	9.39 s
	Weekly	89.78 KB	691.04 KB	2.07 MB	13.32 GB	2.84 s	6.02 s
20	All-Time	2.82 MB	23.41 MB	70.23 MB	1.72 TB	680.54 s	2008.16 s
	Monthly	2.87 MB	4.07 MB	12.20 MB	234.94 GB	53.07 s	137.28 s
	Weekly	2.76 MB	2.77 MB	8.31 MB	282.31 GB	26.85 s	61.78 s
40	All-Time	3.80 MB	25.46 MB	76.37 MB	3.16 TB	1736.52 s	5124.89 s
	Monthly	5.18 MB	6.32 MB	18.97 MB	461.65 GB	119.43 s	307.34 s
	Weekly	5.50 MB	4.89 MB	14.68 MB	549.78 GB	45.31 s	103.20 s
60	All-Time	6.38 MB	27.44 MB	82.32 MB	4.91 TB	2281.91 s	6732.03 s
	Monthly	8.20 MB	8.53 MB	25.60 MB	682.80 GB	149.21 s	389.73 s
	Weekly	8.58 MB	7.09 MB	21.27 MB	821.44 GB	68.89 s	156.38 s
80	All-Time	8.31 MB	29.61 MB	88.83 MB	6.44 TB	3079.08 s	9083.17 s
	Monthly	10.80 MB	10.81 MB	32.44 MB	898.70 GB	203.77 s	511.47 s
	Weekly	9.12 MB	9.26 MB	27.79 MB	1.09 TB	104.49 s	238.56 s
100	All-Time	10.09 MB	31.50 MB	94.49 MB	8.03 TB	3948.97 s	11650.12 s
	Monthly	13.50 MB	12.85 MB	38.56 MB	1.11 TB	200.78 s	523.98 s
	Weekly	12.85 MB	11.18 MB	33.54 MB	1.34 TB	125.31 s	286.26 s

case this would increase to 70.23 MB and approximately 34 minutes respectively. In the same scenario, BIP-37 SPV will require only 2.76, 2.87 and 2.82 MB respectively.

We expect that clients will query monthly and weekly data most often and that all-time data will only be queried when the client is synchronising with the Bitcoin blockchain for the

first time, and when the client has been disconnected from the Bitcoin blockchain for more than a month. We imagine that these scenarios will happen infrequently, and as such the user experience will not be adversely affected by the large latency of querying all-time data.

VI. FUTURE WORK

In this section, we mainly discuss the future work on the *manifest file trie* scheme, which attempts to spare the manifest file downloading bandwidth. We further discuss other extensions (i.e. database partitioning, dynamic protocol and integration with bitcoin) in Appendix C.

Since in some cases most of the PIR database entries hold distinct identities (e.g. the weekly address PIR database), which requires a single record in the manifest file for each entry in the database, manifest files have a nearly equivalent or even larger size compared with the database (cf. Table I). This may incur a poor performance in bandwidth because the client is required to download the whole manifest file.

Manifest files are currently structured as a list of entries which need to be sent to the client in order for the client to be able to query the corresponding databases. Sending the client an updated copy of these manifest files whenever a server-side change occurs (which would be every 10 minutes) is not bandwidth efficient. In order to mitigate this problem, we propose transforming the simple manifest files into a format suitable for PIR queries. These PIR compatible manifest files would then be stored on the PIR server and clients, in order to extract a particular record, would perform interpolation search on these manifest files, under PIR.

With this approach we remove the excessive bandwidth cost of sending full manifest files to the client, and ensure that the record of interest is extracted privately. The addition of PIR-based manifest files requires a slight update to our protocol, whereby the client would first need to perform interpolation search under PIR in order to select each record from each of the manifest files, before using that record in our protocol as outlined in Section IV.

VII. RELATED WORK

There exists a limited body of work which focuses on solving the problem associated with the privacy provisions of Bloom filters for the SPV protocol.

In 2017, Kanemura et al. [16] proposed a γ -deniability enabled Bloom filter to be used instead of the current implementation as specified in BIP-37. γ -deniability is a privacy metric that shows how many true positives are hidden by false positives. In this case, for the number of Bitcoin addresses or TXIDs set in a Bloom filter. Hence, a higher γ -deniability value would indicate a greater level of privacy since more true positives would be obscured by false positives. There exists an alternative proposal detailed in BIP-157 [17] which describes client-side block filtering. This approach is the inverse of BIP-37. Rather than having clients send a filter to a full node, full node peers generate deterministic filters on each block of transactions that they store and send these to the client. The client is then able to see if this deterministic filter matches any data it is interested in. In the case that there is a positive match, the client requests for the relevant full block of data to be downloaded. Compared to BIP-37, the privacy provision of this proposal is greatly improved. This is because full node peers are only able to determine which full block the client

is interested in, compared to a Bloom filter, which identifies individual addresses and transactions a client *may* be interested in.

Compared to solutions based on probabilistic data structures, the only information leakage in our implementation is that of time since clients can choose which time period they would like to query. It can be easily assumed that all clients would like to query the most recent data found in the weekly databases, and the monthly and all-time set's of data are too large for an adversary to gain any meaningful information. In addition, our protocol uses less bandwidth since the PIR queries are precise. Since each block of transactions has a maximum size of 1 MB in Bitcoin, client-side block filtering will require the client to download at least 1 MB of data in the worst case. If a client would like to verify 20 transactions which occurred in the past week, they would incur a bandwidth cost of 8.31 MB by using our protocol, in the worst case. If these transactions are located in individual blocks, then the bandwidth cost of BIP-157 would be 20 MB in the worst case.

Matetic et al. [18] and Wüst et al. [19] proposed approaches to protect the privacy of lightweight clients, which however requires a trusted execution environment.

Other work has focused on increasing the privacy of Bitcoin transactions. Dandelion [20], which is specified in BIP-156, proposes privacy enhanced routing for transactions, but is yet to be integrated into the core Bitcoin client. Currently, when a transaction is submitted into the Bitcoin peer-to-peer network, an adversary can easily link an IP address to the node which broadcasts that particular transaction. Dandelion mitigates this class of attacks by sending transactions over a randomly selected path in the Bitcoin peer-to-peer network before broadcasting that transaction to other peers.

There are some other protocols attempting to mitigate the privacy issues (e.g., Monero [21], Zcash [22], Solidus [23], [24] and Cryptonote [25]). These techniques ensure that even though the blockchain of these cryptocurrencies is public, information regarding the source, amount and destination of a transaction is obfuscated and cannot be deduced. For further analysis of their effectiveness, we refer the reader to [26]–[28].

VIII. CONCLUSION

Bitcoin's lack of privacy with the use of Bloom filters for SPV clients has been a known issue for many years. Our work solves this problem by applying PIR to the SPV protocol and we develop an implementation to validate our design. We used this software to measure the bandwidth cost and latency of our protocol and found that its bandwidth cost was similar to that of BIP-37 SPV in the single-server setting of the hybrid PIR schema, and comparable in the multi-server setting. We also proposed some improvements to our protocol, one of which removed the bandwidth cost of downloading large manifest files. We believe we are the first to detail a fully private practical protocol for SPV in Bitcoin.

REFERENCES

- [1] S. Nakamoto *et al.*, “Bitcoin: A peer-to-peer electronic cash system,” 2008.
- [2] M. Hearn and M. Corallo, “BIP-37,” <https://github.com/bitcoin/bips/blob/master/bip-0037.mediawiki>.
- [3] “bitcoinj.” [Online]. Available: <https://bitcoinj.github.io/>
- [4] M. Hearn, “Bloom filter privacy and thoughts on a newer protocol.” [Online]. Available: <http://www.openbitcoinprivacyproject.org/2015/02/bloom-filter-privacy-and-thoughts-on-a-newer-protocol/>
- [5] J. Nick, “Privacy in BitcoinJ - nickler’s.” [Online]. Available: <https://jonasnick.github.io/blog/2015/02/12/privacy-in-bitcoinj/>
- [6] A. Gervais, S. Capkun, G. O. Karame, and D. Gruber, “On the privacy provisions of bloom filters in lightweight bitcoin clients,” in *Proceedings of the 30th Annual Computer Security Applications Conference*. ACM, 2014, pp. 326–335.
- [7] R. C. Merkle, “Protocols for public key cryptosystems,” in *1980 IEEE Symposium on Security and Privacy*. IEEE, 1980, pp. 122–122.
- [8] “BlockSci.” [Online]. Available: <https://citp.github.io/BlockSci/demo.html>
- [9] A. Gervais, G. O. Karame, K. Wüst, V. Glykantzis, H. Ritzdorf, and S. Capkun, “On the security and performance of proof of work blockchains,” in *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*. ACM, 2016, pp. 3–16.
- [10] C. Devet and I. Goldberg, “The best of both worlds: Combining information-theoretic and computational pir for communication efficiency,” in *International Symposium on Privacy Enhancing Technologies Symposium*. Springer, 2014, pp. 63–82.
- [11] C. Devet, I. Goldberg, and N. Heninger, “Optimally robust private information retrieval,” in *Presented as part of the 21st {USENIX} Security Symposium ({USENIX} Security 12)*, 2012, pp. 269–283.
- [12] C. Aguilar-Melchor and P. Gaborit, “A lattice-based computationally-efficient private information retrieval protocol,” in *Western European Workshop on Research in Cryptology*. Citeseer, 2007.
- [13] “Percy++ / PIR in C++.” [Online]. Available: <http://percy.sourceforge.net/>
- [14] “bitcoin/descriptors.md at master · bitcoin/bitcoin,” <https://github.com/bitcoin/bitcoin/blob/master/doc/descriptors.md>.
- [15] “Protocol documentation - Bitcoin Wiki.” [Online]. Available: https://en.bitcoin.it/wiki/Protocol_documentation#Message_structure
- [16] K. Kanemura, K. Toyoda, and T. Ohtsuki, “Design of privacy-preserving mobile bitcoin client based on γ -deniability enabled bloom filter,” in *2017 IEEE 28th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*. IEEE, 2017, pp. 1–6.
- [17] O. Osuntokun, A. Akselrod, and J. Posen, “BIP-157,” <https://github.com/bitcoin/bips/blob/master/bip-0157.mediawiki>.
- [18] S. Matic, K. Wüst, M. Schneider, K. Kostianen, G. Karame, and S. Capkun, “Bite: Bitcoin lightweight client privacy using trusted execution,” *IACR Cryptology ePrint Archive 2018, XXXX*, Tech. Rep., 2018.
- [19] K. Wüst, S. Matic, M. Schneider, I. Miers, K. Kostianen, and S. Capkun, “Zlite: Lightweight clients for shielded zcash transactions using trusted execution,” in *International Conference on Financial Cryptography and Data Security*. Springer, 2019.
- [20] B. Denby, A. Miller, G. Fanti, S. Bakshi, S. Venkatakrishnan, and P. Viswanath, “Dandyion.” [Online]. Available: <https://github.com/mablern8/bips/blob/master/bip-dandelion.mediawiki#Implementation>
- [21] “Home — Monero - secure, private, untraceable.” [Online]. Available: <https://getmonero.org/>
- [22] “Home Zcash.” [Online]. Available: <https://z.cash/>
- [23] I. Abraham, D. Malkhi, K. Nayak, L. Ren, and A. Spiegelman, “Solidus: An incentive-compatible cryptocurrency based on permissionless byzantine consensus,” *CoRR, abs/1612.02916*, 2016.
- [24] E. Cecchetti, F. Zhang, Y. Ji, A. Kosba, A. Juels, and E. Shi, “Solidus: Confidential distributed ledger transactions via pvorm,” in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2017, pp. 701–717.
- [25] N. Van Saberhagen, “Cryptonote v 2.0,” 2013.
- [26] M. Möser, K. Soska, E. Heilman, K. Lee, H. Heffan, S. Srivastava, K. Hogan, J. Hennessey, A. Miller, A. Narayanan *et al.*, “An empirical analysis of traceability in the monero blockchain,” *Proceedings on Privacy Enhancing Technologies*, vol. 2018, no. 3, pp. 143–163, 2018.
- [27] A. Kumar, C. Fischer, S. Tople, and P. Saxena, “A traceability analysis of moneros blockchain,” in *European Symposium on Research in Computer Security*. Springer, 2017, pp. 153–173.
- [28] G. Kappos, H. Yousaf, M. Maller, and S. Meiklejohn, “An empirical analysis of anonymity in zcash,” in *27th {USENIX} Security Symposium ({USENIX} Security 18)*, 2018, pp. 463–477.

APPENDIX

A. Manifest Files Formats

In this section, we describe the formats for the three databases’ manifest files.

1) *Address PIR DB Manifest File*: Each record in this manifest file consists of a mapping of an address to the location of the corresponding row entry(ies) in the Address PIR DB. Figure 7 illustrates the structure of this record.

```
{...
  (i)   "address" : [
  (ii)  "row_index_start",
  (iii) "row_index_end",
  (iv)  "column_index_start",
  (v)   "column_index_end",
  ]
...}
```

Fig. 7: Example record in Address PIR DB manifest file

- (i) The address contained in the address field of the entry referenced by this record.
- (ii) The row index where the first entry is located.
- (iii) The row index where the last entry is located.
- (iv) The column index where the first entry is located.
- (v) The column index where the last entry is located.

2) *Merkle Tree PIR DB Manifest File*: Each record in this manifest file consists of a mapping of a block height to the location of the corresponding list of TXIDs in the Merkle Tree PIR DB. Figure 8 illustrates the structure of this record.

```
{...
  (i)   "block-height" : [
  (ii)  "row_index_start",
  (iii) "row_index_end",
  (iv)  "column_index_start",
  (v)   "column_index_end",
  ]
...}
```

Fig. 8: Example record in Merkle Tree PIR DB manifest file

- (i) The block height of the TXIDs referenced by this record.
- (ii) The row index where the first TXID is located.
- (iii) The row index where the last TXID is located.
- (iv) The column index where the first TXID is located.
- (v) The column index where the last TXID is located.

3) *Transaction PIR DB Manifest File*: Each record in this manifest file consists of a mapping of a TXID to the location of the corresponding transaction in the Transaction PIR DB. Figure 9 illustrates the structure of this record.

```

{...
(i)    "txid" : [
(ii)   "row_index_start",
(iii)  "row_index_end",
(iv)   "column_index_start",
(v)    "column_index_end",
      ]
...}

```

Fig. 9: Example record in Transaction PIR DB manifest file

- (i) The TXID of the transaction referenced by this record.
- (ii) The row index where the first byte of the transaction is located.
- (iii) The row index where the last byte of the transaction is located.
- (iv) The column index where the first byte of the transaction is located.
- (v) The column index where the last byte of the transaction is located.

B. Database Dimensions Determination

Below we detail the methodology we followed for determining the dimensions of the databases.

a) Address PIR DB Dimensions: The size of the all-time, monthly and weekly data was measured and the height and width values were chosen such that the databases were square in shape, to minimise the communication cost.

b) Merkle Tree & Transaction PIR DB Dimensions: To determine the width of the rows for these two database types, two graphs were plotted. Figure 10 shows the number of transactions per block and Figure 11 shows how the length of transactions changes over time. For both database types, for the three time periods, the width was taken as the expected value of the respective data sets. The statistical measure of expectation was used for the width because this would reduce the number of times a query would have to be performed, given that in expectation, the row result given to the client would contain the complete set of data required. However, the width of the Merkle Tree PIR DB for the set of all-time data was taken as the running average instead, since Figure 10 shows an extremely skewed distribution. Monthly and weekly Merkle Tree DBs were square, while the rest, namely the all-time Merkle Tree DB and all Transaction DBs, were rectangular.

C. Extended Future Work

1) Database Partitioning: Since all-time databases are notably large, partitioning them would reduce the bandwidth and latency cost of queries, since they will be executed over a smaller data set. Database partitioning would require for two new fields, "pir_db_start" and "pir_db_end", to be added to the respective manifest files, to track the databases in which a particular entry is included. This would require the PIR based SPV protocol to be slightly adjusted, with the client selecting both the database and row indices when constructing queries. However, database partitioning would expose the client to a

statistical privacy attack. Here, a malicious PIR server can monitor for a pattern of querying across multiple databases by the same client and attempt to deduce the user's interests. As such, the partitioning should not be excessive. As a guideline, each partitioned all-time sub-database should not be smaller than the size of the corresponding database for the monthly set of data.

2) Dynamic Protocol: Currently our proof-of-concept implementation is static, which means that it does not reflect the latest state of the Bitcoin blockchain. A new block of transactions is mined and sent out into the Bitcoin peer-to-peer network approximately once every 10 minutes.

In order to create a dynamic implementation of our protocol, these new blocks would first need to be parsed into the format described in Sections IV. The new data would then be appended to the final rows of the appropriate weekly databases, with Address PIR DB data being subsequently sorted to maintain lexicographic ordering.

After collecting a week's worth of new data in the weekly databases, these databases are subsequently emptied by having their data migrated to the monthly databases. Similarly, once a month's worth of new data is collected in the monthly databases, these databases are also emptied by having their data migrated to the all-time databases, where it permanently remains. This migration simply involves the new data being appended to the final rows of the appropriate monthly or all-time databases, with Address PIR DB data being subsequently sorted to maintain lexicographic ordering. In addition, whenever such updates occur, either when new data arrives in the weekly databases or when data is migrated, these changes are reflected in the manifest files of the corresponding databases.

This process ensures that the freshest data exists in the set of weekly databases, followed by the monthly and all-time databases. In addition, PIR servers do not need to perform synchronisation with each other with regards to the current state of the blockchain as that is implicitly handled by Bitcoin.

Our protocol can be further extended by having finer-grained temporal slices of the Bitcoin blockchain, such as for the past day or hour. Having too many temporal partitions, however, can expose the client to a privacy attack. A malicious PIR server can monitor a client's query and determine that the client is interested in hourly data, thus reducing the set of data in which that client's transaction of interest is included, when compared to the set of weekly data. This increases the likelihood of the malicious PIR server determining which addresses the client controls.

3) Integration with Bitcoin: Our protocol as described in this paper can be easily integrated with Bitcoin and other cryptocurrencies with a similar SPV model. Our protocol acts as an alternative SPV implementation which is fully private, and as such would only require changes at the network protocol layer. At a minimum, new network message headers would need to be created which would allow clients to identify which full nodes support PIR SPV. Full nodes which do support PIR SPV would also need to account for extra storage

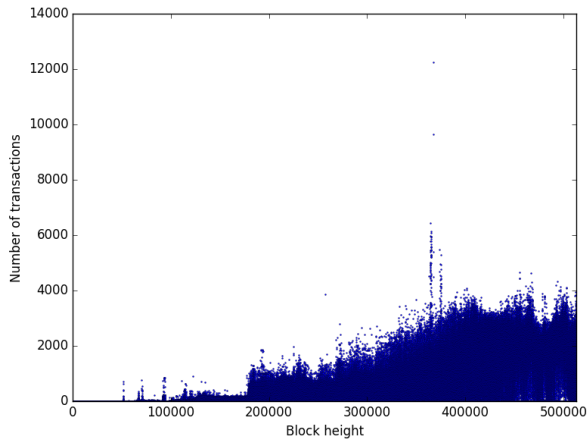


Fig. 10: Number of TXID's per block

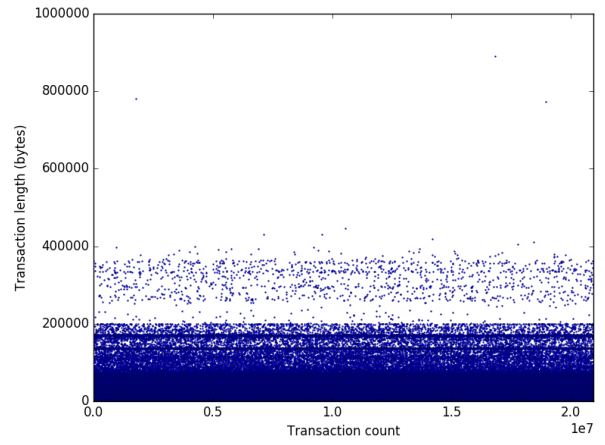


Fig. 11: Length of transactions over time

to accommodate the additional databases and corresponding manifest files that are necessary to facilitate PIR based queries.