



BRNO UNIVERSITY OF TECHNOLOGY

VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

FACULTY OF INFORMATION TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

DEPARTMENT OF INFORMATION SYSTEMS

ÚSTAV INFORMAČNÍCH SYSTÉMŮ

PROCESSING OF THE BLOCKCHAIN EMPLOYING IPFS

VYUŽITÍ IPFS PRO ZPRACOVÁNÍ BLOCKCHAINU

MASTER'S THESIS

DIPLOMOVÁ PRÁCE

AUTHOR

AUTOR PRÁCE

Bc. MATÚŠ MÚČKA

SUPERVISOR

VEDOUCÍ PRÁCE

Ing. VLADIMÍR VESELÝ, Ph.D.

BRNO 2020

Abstract

This work aims to design a system for processing blockchain data of selected cryptocurrencies for further exploring using IPFS. In this thesis, we created a proprietary decentralized and distributed database system that supports simple queries. The solution provides a well-arranged graphical user interface for data visualization as well as REST API, which makes it easy to connect to other systems. The benefit of this work is a new view of blockchain processing which opens up new possibilities in its exploring.

Abstrakt

Cieľom tejto práce je navrhnúť systém na spracovanie a preskúmavanie blockchainu vybraných kryptomien pri použití IPFS. Na riešenie tohoto problému bolo potrebné navrhnúť vlastný decentralizovaný a distribuovaný databázový systém, ktorý podporuje jednoduché dotazy. Vytvorené riešenie poskytuje prehľadné grafické užívateľské rozhranie, ktoré slúži na vizualizáciu dát a taktiež RestAPI, vďaka ktorému sa dá systém jednoducho napojiť na iné systémy. Prínosom tejto práce je nový pohľad na zpracovávanie blockchainu čo otvára nové možnosti v jeho prehľadávaní.

Keywords

IPFS, decentralization, blockchain, cryptocurrecny

Klíčové slová

IPFS, decentralizácia, blockchain, cryptocurrecny

Reference

MÚČKA, Matúš. *Processing of the Blockchain Employing IPFS*. Brno, 2020. Master's thesis. Brno University of Technology, Faculty of Information Technology. Supervisor Ing. Vladimír Veselý, Ph.D.

Processing of the Blockchain Employing IPFS

Declaration

I declare that this semester project was prepared as an original author's work under the supervision of Mr. Ing. Vladimír Veselý, PhD. All the relevant information sources, which were used during the preparation of this project, are cited and included in the list of references.

.....

Matúš Múčka

April 14, 2020

Contents

1	Introduction	3
2	Cryptocurrencies	4
2.1	Bitcoin	5
2.2	DigiByte	5
2.3	Ethereum	6
2.4	Decred	6
2.5	Monero	6
3	Decentralization (Web3.0)	7
4	IPFS	8
4.1	IPFS stack	9
4.2	libp2p	9
4.3	IPLD	10
4.4	IPFS	10
4.5	IPNS	11
4.6	Existing blockchain explorers in IPFS	11
5	Design	13
5.1	Database system	13
5.1.1	OrbitDB	13
5.1.2	Textile	13
5.1.3	Database system design	14
5.2	System components	17
5.3	Feeder	17
5.4	Explorer	18
5.4.1	ExplorerGUI	18
5.4.2	ExplorerAPI	19
6	Implementation	20
6.1	Explorer-core implementation	20
6.1.1	Indexes	20
6.1.2	Query system	21
6.1.3	Database	22
6.2	Feeder implementation	23
6.3	Explorer implementation	24
6.3.1	ExplorerCore	24

6.3.2	ExplorerGUI	26
6.3.3	ExplorerAPI	26
7	Benchmarking and testing	28
7.1	Benchmarking enviroment	28
7.1.1	Blockbook	28
7.1.2	Writing	28
7.1.3	Reading	28
7.1.4	Filtering	28
8	Conclusion	29
	Bibliography	30

Chapter 1

Introduction

HTTP is “good enough” for the most use cases of distributing files over the network. However, when we want to stream lots of data to multiple connected clients at once, we are starting to hit its limits. When two clients are requesting the same data, there is no mechanism in HTTP that would allow sending the data only once. Sending duplicate data has become a problem in large companies because of bandwidth capacity. Blizzard¹ started to distribute video game content via a distributed solution because it was cheaper for the company and faster for players[7]. Linux distributions use BitTorrent to transmit disk images².

The Bitcoin blockchain has now 242 gigabytes³. When blockchain is processed (all its data are parsed), the size can grow twice as much. If there are multiple blockchains, then data can have few terabytes. When we are sharing blockchains data from the server for several clients, there is a big chance that multiple clients want the same data. They may be working on the same case and investigating the same wallets. So in standard solution with relational database and some HTTP server, the server has to search in all data (that can have a size of few terabytes) and transmits selected data to the client for every request. This problem happens even if a different client asks for the same data in a few minutes ago. Behaviour mentioned above dramatically limits the scalability of the server.

Services that are using HTTP, often have client-server architecture, so there is also a problem with a single point of failure. If the server for some reason stops working, the client can not receive data. In distributed file system such as IPFS, there is no such problem as one point of failure because all data are duplicated on multiple clients.

This Semestral project is divided into six chapters. Chapter 2 describes the differences between cryptocurrencies used in this project. Chapter 4 describes IPFS and all its layers. Design of the system created in this thesis is in Chapter 5. The implementation of all applications that are provided with this project is in chapter 6. Finally, summarization of results achieved in this work is in chapter 8.

¹Game company - <https://www.blizzard.com/>

²Image of Debian downloadable by BitTorrent <https://www.debian.org/CD/torrent-cd/>

³Current size of bitcoin blockchain can be seen at <https://www.statista.com/statistics/647523/worldwide-bitcoin-blockchain-size/>

Chapter 2

Cryptocurrencies

There were hundreds of failed attempts of creating cryptographic payment systems before cryptocurrencies like Bitcoin and Ethereum came into existence. Some of these systems are listed in Figure 2.1. All of them were created before Bitcoin. Despite that, only a few of them survived to these days. Some of these attempts were only academic proposals while others were deployed and tested systems. One of the survival is PayPal. It is only because it quickly gave up its original idea of hand-held devices for cryptographic payments.[11]

So there is a question, what makes cryptocurrencies successful nowadays? It may be easy to use principle and no need for external hardware. Another critical component of cryptocurrencies discussed in this work is Blockchain. Generally, it is a ledger in which all transactions are securely stored. The idea behind blockchains is pretty old, and it was originally used for timestamping digital documents.[5]

ACC	CyberCents	iKP	MPTP	Proton
Agora	CyberCoin	IMB-MP	Net900	Redi-Charge
AlIMP	CyberGold	InterCoin	NetBill	S/PAY
Allopass	DigiGold	Ipin	NetCard	Sandia Lab E-Cash
b-money	Digital Silk Road	Javien	NetCash	Secure Courier
BankNet	e-Comm	Karma	NetCheque	Semopo
Bitbit	E-Gold	LotteryTickets	NetFare	SET
Bitgold	Ecash	Lucre	No3rd	SET2Go
Bitpass	eCharge	MagicMoney	One Click Charge	SubScrip
C-SET	eCoin	Mandate	PayMe	Trivnet
CAFÉ	Edd	MicroMint	PayNet	TUB
CheckFree	eVend	Micromoney	PayPal	Twitpay
ClickandBuy	First Virtual	MilliCent	PaySafeCard	VeriFone
ClickShare	FSTC Electronic Check	Mini-Pay	PayTrust	VisaCash
CommerceNet	Geldkarte	Minitix	PayWord	Wallie
CommercePOINT	Globe Left	MobileMoney	Peppercoin	Way2Pay
CommerceSTAGE	Hashcash	Mojo	PhoneTicks	WorldPay
Cybank	HINDE	Mollie	Playspan	X-Pay
CyberCash	iBill	Mondex	Polling	

Figure 2.1: Electronic payment systems before cryptocurrencies[9]

2.1 Bitcoin

Bitcoin is probably the most famous cryptocurrency. On 3 January 2009, Satoshi Nakamoto (an alias for a person or group persons authored the bitcoin white paper) mined the genesis block of bitcoin (block with height 0). Satoshi gets the reward of 50 bitcoins (half a million US dollars in the time of writing). This text was embedded in the block (see Figure 2.2): *The Times 03/Jan/2009 Chancellor on brink of second bailout for banks.*¹[3].

00000000	01 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00
00000010	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00
00000020	00 00 00 00 3B A3 ED FD	7A 7B 12 B2 7A C7 2C 3E;Łıýz{.²zÇ,>
00000030	67 76 8F 61 7F C8 1B C3	88 8A 51 32 3A 9F B8 AA	gv.a.È.Ã^ŠQ2:Ÿ,ª
00000040	4B 1E 5E 4A 29 AB 5F 49	FF FF 00 1D 1D AC 2B 7C	K.^J)«_IŸŸ...¬+
00000050	01 01 00 00 00 01 00 00	00 00 00 00 00 00 00 00
00000060	00 00 00 00 00 00 00 00	00 00 00 00 00 00 00 00
00000070	00 00 00 00 00 00 FF FF	FF FF 4D 04 FF FF 00 1DŸŸŸŸM.ŸŸ..
00000080	01 04 45 54 68 65 20 54	69 6D 65 73 20 30 33 2F	..E The Times 03/
00000090	4A 61 6E 2F 32 30 30 39	20 43 68 61 6E 63 65 6C	Jan/2009 Chancel
000000A0	6C 6F 72 20 6F 6E 20 62	72 69 6E 6B 20 6F 66 20	lor on brink of
000000B0	73 65 63 6F 6E 64 20 62	61 69 6C 6F 75 74 20 66	second bailout f
000000C0	6F 72 20 62 61 6E 6B 73	FF FF FF FF 01 00 F2 05	or banksŸŸŸŸ..ò.
000000D0	2A 01 00 00 00 43 41 04	67 8A FD B0 FE 55 48 27	*....CA.gŠŸ°pUH'
000000E0	19 67 F1 A6 71 30 B7 10	5C D6 A8 28 E0 39 09 A6	.gñ q0..\"Ö" (à9.!
000000F0	79 62 E0 EA 1F 61 DE B6	49 F6 BC 3F 4C EF 38 C4	ybaê.aB¶IÖ¼?Li8Ã
00000100	F3 55 04 E5 1E C1 12 DE	5C 38 4D F7 BA 0B 8D 57	óU.â.Á.Ð\8M÷°..W
00000110	8A 4C 70 2B 6B F1 1D 5F	AC 00 00 00 00	ŠLp+kñ._¬....

Figure 2.2: Genesis bitcoin block

First documents purchase happened in 22nd May 2010, when Laszlo Hanyecz bought two pizzas for 10 000 bitcoins (\$41 then, now about \$80 000 000)². This transactions hash is a1075db55d416d3ca199f55b6084e2115b9345e16c5cf302fc80e9d5fbf5d48d and will be stored in bitcoin blockchain forever³.

2.2 DigiByte

DigiByte was developed and released in 2013. It is based on Bitcoin with some adjustment in the code to improving functionality. In late 2017 there was 200 000 pending transaction in Bitcoin. Miners preferred transaction with a bigger fee, so to confirm a transaction, user needed to pay \$50. Digibyte solved this problem by adding a new block every 15 seconds (new block in Bitcoin is mined every 10 minutes). Average transaction occupies 570 bytes of data. One block can contain approximately 3 500 transactions given the 2 MB limit. This means that in DigiByte, 230 transactions can be confirmed in one second compared

¹https://en.bitcoin.it/wiki/Genesis_block

²<https://bitcointalk.org/index.php?topic=137.0>

³Actual photos of \$80 000 000 pizza are on <http://heliacal.net/~solar/bitcoin/pizza/>

to Bitcoin 4-7 transaction per second. DigiByte also has 1000:1 DigiByte to Bitcoin ratio, so for every Bitcoin there will be 1000 DigiByte.[2]

2.3 Ethereum

While both the Bitcoin and Ethereum networks are powered by the principle of distributed ledgers and cryptography, the two differ technically in many ways. For example, transactions on the Ethereum network may contain executable code, while data affixed to Bitcoin network transactions are generally only for keeping notes. Other differences include block time (an ether transaction is confirmed in seconds compared to minutes for bitcoin) and the algorithms that they run on (Ethereum uses ethash while Bitcoin uses SHA-256).[12]

2.4 Decred

Decred is cryptocurrency build from Bitcoin. Main difference from Bitcoin is the rewarding system from mining. In Bitcoin, the miner gets full reward for a mined block. Sometimes, the Bitcoin blockchain splits when two or more miners found a block at nearly the same time. The fork is resolved when the subsequent block(s) are added, and one of the chains becomes longer than the alternative(s). In Decred the chance of blockchain forks is minimized by hybrid proof of work and proof of state system. Each time a block is created (by a miner), it is not automatically part of the blockchain. Block needs to be approved by ticket holders. Then miner receives a block reward (newly created DCR). If the block is rejected by ticket holders, the miner does not receive a reward. Tickets holders for a new block are chosen randomly. The ticket validates the previous block. A block needs at least three of the five votes chosen to approve it for it to be validated. This hybrid system has many implications, including making a 51% hashpower attack very difficult, and making a minority fork very difficult as well.[6]

2.5 Monero

Three years after Bitcoin, in 2012, the competing Bytecoin cryptocurrency entered the market. The problem with this cryptocurrency, however, was that 80% of all coins were mined in advance by its authors. The chances of mining were, therefore, not balanced. This led to the decision that this cryptocurrency would start again. New cryptocurrency starts on 18 April 2014 and was called BitMonero, a composite of the word coin in Esperanto (Monero) and Bitcoin according to Bitcoin. However, after five days, the community decided to use only Monero for short. Monero's significant advantage is the dynamic size of the mined blocks. Bitcoin has one block size limited to 1 MB, while Monero adapts the block size to the network load. If the number of transactions increases, so does the block size to accommodate all transactions. Thus, unlike Bitcoin, the more transactions users make, the lower the transaction fee. Monero's main benefit is its full anonymity and interchangeability. Monero hides recipient and sender addresses.[10]

Chapter 3

Decentralization (Web3.0)

Chapter 4

IPFS

IPFS¹ stands for InterPlanetary File System and is a peer-to-peer distributed filesystem designed to make the Web faster, safer, and more open. In contrast with a standard filesystem, objects in IPFS are content-addressed by the cryptographic hash of their contents. In the case of the standard Web, when user wants some file, he needs to know on which server is a file located and the full path to the file (see Figure 4.1). In IPFS user needs only to know the hash of the requested file. He does not care about the location of the file (see Figure 4.2). Let us take an MIT license text, and add it to IPFS. If somebody tries to add this license as a file to IPFS, it will return `QmWpvK4bYR7k9b1feM48fskt2XsZfMaPfNnFxdbhJHw7QJ` every time. That is now and will be in the future, the *content address* of that file. Later, when user try to get this file by its hash, he can get it from a random person that added it into IPFS in the past.

IPFS can easily represent a file system consisting of files and directories. A small file (less than 256 kB) is represented by an IPFS object with data being the file contents (plus a small header and footer). Note that the file name is not part of the IPFS object, so two files with different names and the same content will have the same IPFS object representation and hence the same hash. A large file (more than 256 kB) is represented by a list of links to file chunks that are less than 256 kB, and only minimal information specifying that this object represents a large file. Currently, there is no known size limitations uploaded file or directory. There are already some big datasets hosted on IPFS such as Geocities archive² (704TB) or Project Apollo Archives³ (61GB).

IPFS provides API for manipulation with objects in several levels of abstraction. Two most used APIs are Files and Graph API. The files API⁴ enables users to use the File System abstraction of IPFS. Currently it has `add`, `cat`, `get` and `ls` methods. In the future, there should be MFS (Mutable File System), that would add methods `chmod`, `cp`, `mv`, `rm`, `stat` and `touch`⁵. Graph API⁶ is used in lower level of abstraction for manipulating with any IPFS object and creating links between objects. There are `put` and `get` methods available.

¹<https://ipfs.io/>

²<https://ipfs.io/ipfs/QmVCjhoEFC9vwvaa8bKyJgwABYP4MXSogcyDGoZ4Lkc3ox>

³<https://ipfs.io/ipfs/QmSnuWmxptJZdLJpKRarxBMS2Ju2oANVrgbr2xWbie9b2D>

⁴<https://github.com/ipfs/interface-js-ipfs-core/blob/master/SPEC/FILES.md>

⁵<https://docs.ipfs.io/guides/concepts/mfs/>

⁶<https://github.com/ipfs/interface-js-ipfs-core/blob/master/SPEC/DAG.md>

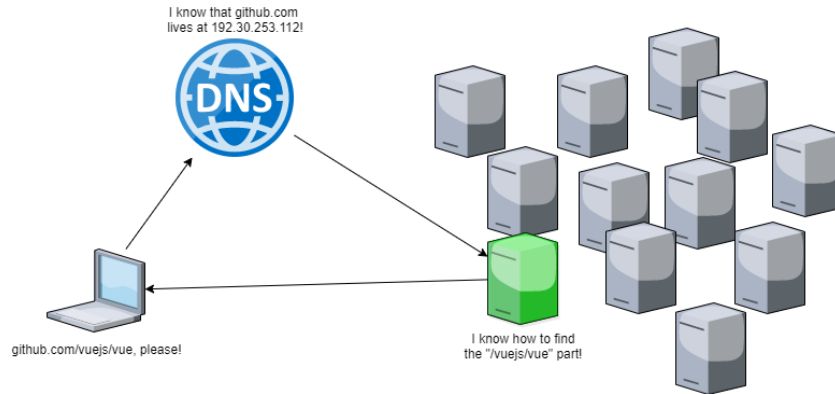


Figure 4.1: Classic web addressing

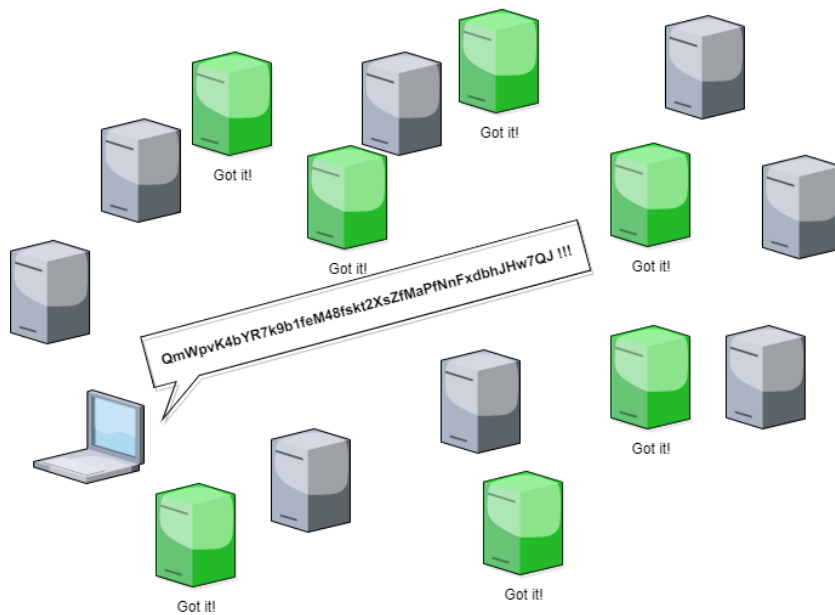


Figure 4.2: Content-based addressing

4.1 IPFS stack

We can split IPFS into layers (see Figure 4.3). *Libp2p*⁷ is at the bottom, which is a peer-to-peer networking module, that handles peer and content discovery, transport, security, identity, peer routing, and messaging. *IPDL* is the data model of the content-addressable web. It is providing linking between objects and multihash computing. On the top is *IPFS* which allows to publish and share files (or any data).[1]

4.2 libp2p

Libp2p is a modular system of protocols, specifications and libraries that enable the development of peer-to-peer network applications. It provides NAT Traversal, Peer Discov-

⁷<https://libp2p.io/>

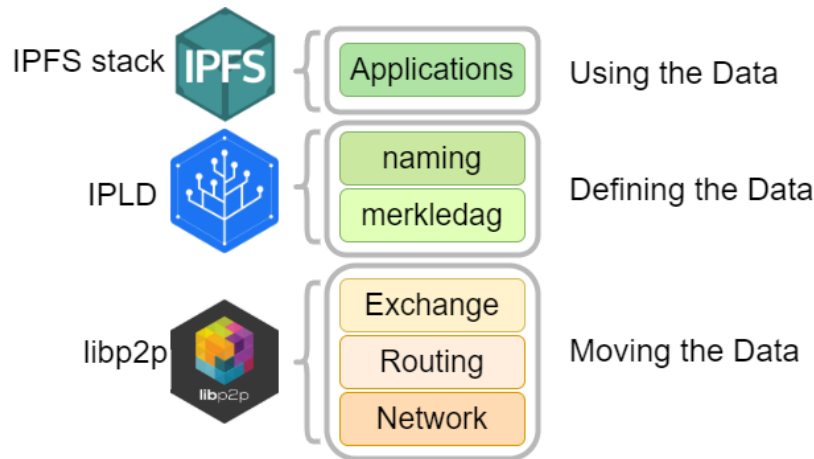


Figure 4.3: IPFS stack

ery, Routing, Stream Multiplexing, Protocol Multiplexing, Encryption, Authentication and more. It grew out of IPFS to solve networking problems in p2p networks, but now it does not require or depend on IPFS. Today many projects use libp2p as their network transporting layer, which is responsible for the actual transmission and receipt of data from one peer to another. For both content discovery and peer routing, libp2p uses Kademlia-based distributed hash table. With Kademlia, libp2p iteratively routes requests closer to the desired peer or content using Kademlia routing algorithm[8]. In the future, Kademlia might be changed easily to some other solutions that implements a simple interface for publishing and requesting data and finding a peer.[4]

4.3 IPLD

IPLD is providing linking and addressing objects with CID (Content ID). CID is hash-based self-describing content identifier (usually encoded to base58⁸ format) which includes codec and multihash. Multihash is then further composed of hashtype and hash value. Let us look closer on the MIT license file, that we add to IPFS at the beginning of this chapter (see Figure 4). It's CID is `QmWpvK4bYR7k9b1feM48fskt2XsZfMaPfNnFxdbhJHw7QJ`. It can be converted to human-readable format as can be seen in Figure 4.1, thanks to multicodec table⁹. We can see that this CID is encoded in base58 format and the file was stored using protobuf¹⁰ codec (this information is necessary to decode file correctly).

4.4 IPFS

IPFS is the top layer from the IPFS stack. It is used for pinning objects and files, naming system (see IPNS¹¹) and keys management. File or object is automatically pinned when a user adds it (but other IPFS commands do not include automatic pinning). Pinning a CID tells an IPFS server that the data is important and must not be thrown away. When

⁸<https://en.wikipedia.org/wiki/Base58>

⁹<https://github.com/multiformats/multicodec/blob/master/table.csv>

¹⁰https://en.wikipedia.org/wiki/Protocol_Buffers

¹¹<https://docs.ipfs.io/guides/concepts/ipns/>

Property		Value
Multibase		base58btc
Version		cidv0
Multicodec		dag-pb
Multihash	Hash Type	sha2-256
	Hash Length	256
	Hash	7e1b666c0327...3dc3022f

Table 4.1: Example of human redable version of CID

garbage collection is triggered on a node, any pinned content is automatically exempt from deletion. Non-pinned data may be deleted. The InterPlanetary Name System (IPNS) is a system for creating and updating mutable links to IPFS content. Since objects in IPFS are content addressed, an object address changes every time an object's content changes. A name in IPNS is the hash of a public key. It is associated with a record containing information about the hash it links to that is signed by the corresponding private key.

4.5 IPNS

Naming inside IPFS is governed by IPNS¹², the InterPlanetary Naming System. IPNS takes ideas from SFS¹³ to enable the creation of cryptographically signed mutable pointers, which can be used to the creation of name records inside the network.

4.6 Existing blockchain explorers in IPFS

There are already stored a few blockchains of cryptocurrencies in IPFS. For browsing them we can use dedicated applications^{14 15} or IPLD explorer¹⁶. Blockchains in IPFS are stored in raw binary format, so custom IPLD codec has to be created for every type of object. Using custom codecs for cryptocurrencies allows explorers to request blocks and transactions by its hash very fast, but there are also limitations.

- Existing IPLD codecs for cryptocurrencies are very limited. Only a few of cryptocurrency codecs are currently available in IPLD (namely Leofcoin, Ethereum, Bitcoin, Zcash, Steller, Decred and Dash)¹⁷.
- Addresses are not stored in IPFS, because they are not part of a blockchain. This means the explorer needs to go through the entire blockchain for computing address balance or to find address transactions.
- No additional information (for example transaction value in US dollars) can be stored with objects because it would change content, and therefore hash of the object.

¹²<https://docs.ipfs.io/guides/concepts/ipns/>

¹³https://en.wikipedia.org/wiki/Self-certifying_File_System

¹⁴<https://github.com/arcalineia/IPFS-Zcash-Explorer>

¹⁵<https://github.com/whyrusleeping/zcash-explorer>

¹⁶<https://explore.ipld.io/#/explore/z43AaGEvwdfzjrCZ3Sq7DKxdDHrwoaPQDtqF4jfdkNEVTiqGVFW>

¹⁷<https://github.com/multiformats/multicodec/blob/master/table.csv>

- There is no sorting or filtering. Explorer can only show the object (block or transaction) by its hash.

Chapter 5

Design

This chapter describes the proposed design of the whole system for storing and exploring blockchains in IPFS that is to be created as a result of this thesis. All parts of the system are described in this chapter.

5.1 Database system

We need a database in our system to store and index data. Database system needs to be decentralized and distributed. There are several databases build on top of IPFS already implemented. The two most known are OrbitDB and Textile.

5.1.1 OrbitDB

OrbitDB is a serverless, distributed, peer-to-peer database build on top of IPFS, developed by HAJA networks. OrbitDB is a decent solution for small user's databases, but it is still in the alpha stage of developing, and it is not well optimized to store hundreds of gigabytes of data. The biggest problem is that OrbitDB performs all queries locally. To perform a query that selects transactions that are more valuable than 1BTC, OrbitDB needs to load the whole database locally and then perform a cycle on all transactions to select only those transactions that meet the criteria. So every client ends up with a full copy of the database. This limitation is not usable for our case when we have a database that has hundreds of gigabytes of data.

5.1.2 Textile

Textile is a set of open-source tools that provide a distributed peer-to-peer database, remote storage, user management, and more, over the IPFS network. Textile already created applications for storing photos, notes or anything else (Anytype). Textile provides a high abstraction on top of the IPFS and provides simple API to store and index files securely. It uses Cafe peers to provide backups and indexing. Data are duplicated on several Cafe peers. When a client is performing some query, it contacts one of the Cafe peers to resolve the query for the client. Neither of these solutions fits our use case to store a large amount of data distributed on several nodes and performs queries that can be resolved by downloading only necessary parts of the database. Therefore, we need to create a new database system based on IPFS that would be decentralized and distributed.

5.1.3 Database system design

After some research, we concluded that currently for storing and indexing data in IPFS without large harddisk memory consumption, there is no solution. We created our own indexing system that currently supports three types of indexes. A database system that fits our needs is distributed and decentralized. That brings us lots of synchronization problems to solve. This database system consists of tables that contain records. For faster searching, tables have indexes. Relations between tables are represented via foreign keys. Also, this database system supports fluent query language, used for performing complex queries.

Record

Every record in our database system is stored in an append-only log that contains the whole history of the record. Every update of a record adds a new entry to its log which points to the previous entry.

In centralized systems such as git, conflicts are detected on write (for example, when two git users push changes to the server at the same time, one of them gets an error and needs to pull repository). This approach is impossible in a decentralized system. When we update some record in our database, we can not know if somebody does not update it before (and we do not receive changes yet). For this reason, we need to solve conflicts while reading. Record with conflict has more than one head in a record log, and users need to solve them in application logic. Look at the example in Figure 5.1. There is a record that is updated by every country in the World when they have new statistics about a pandemic. If two countries update data at the same time (their updates pointing to the same previous version of the entry) they create conflict. Luckily, this specific conflict is easy to solve. We just need to look at the previous entry and compute increment for both country. Then update the record with final increment.

Indexes

An index is B-tree optimized for IPFS (no cycled references and node size less than 256kB). Each table has at least one primary index. Primary Index - We use the primary index to reference record in foreign keys. A primary key is automatically created when a user does not specify it. Value of the primary key for an entity can not be changed (because we would need to scan all tables where the entity is referenced and change the referenced value to a new one). Also, when we execute a query without any condition or sorting, we use a primary key to obtain records. Every index has several components:

- **Comparator** - is a function that has two parameters (two keys of B-tree) and outputs a number. An index is using a comparator to search records and insert record to the right place in B-tree.
- **Key-getter** - is a function that returns a comparable object from the record (a comparable object is an object that can be compared using index comparator). Every index has key-getter that is using to obtain keys for B-tree.

Table

A table contains indexes and table name. Also it implements operations like **insert**, **update** and **delete**.

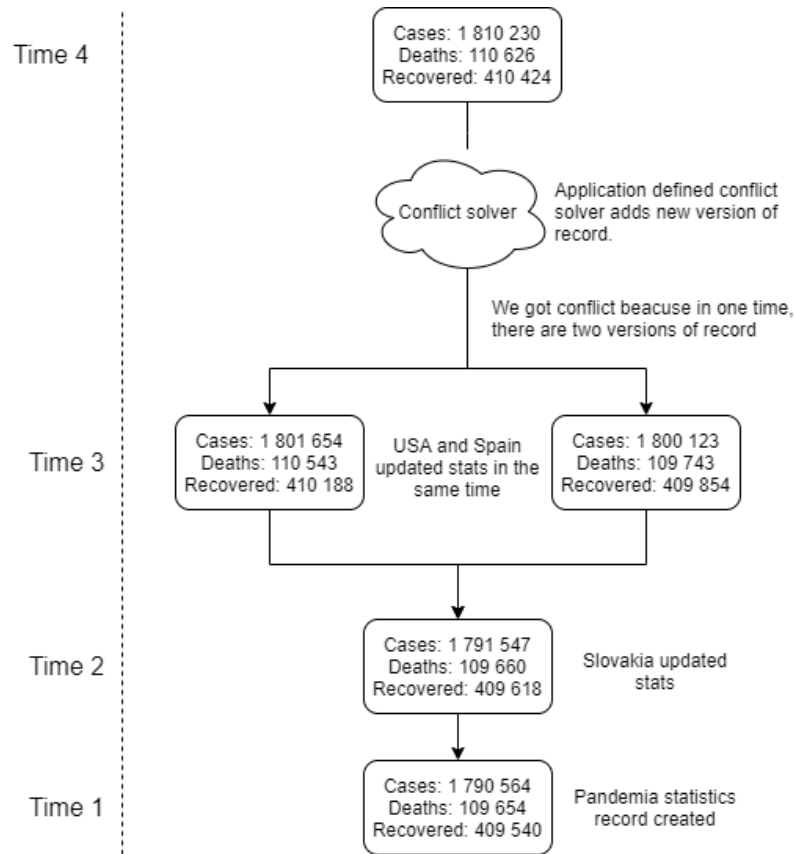


Figure 5.1: Record conflict

Foreign keys

We use foreign keys to represent relationships between tables. A foreign key is simply a table name and a value of a primary key of a referenced record.

Transactions

When we commit a transaction, it is inserted to the transactions queue. If there are more transactions in the queue, a database performs it one by one. We can not execute more transactions in one moment, because it could cause data inconsistency. There are five types of transactions:

- **read** - transaction for reading data from table. It is not logged in a database log (other peers don't have to know what we are reading from database).
- **sync** - when some peer publishes a new version of the database, all other peers have to migrate to it. This transaction has the biggest priority, and therefore it is inserted at the beginning of the transactions queue. After a migration is done, transactions queue can continue working normally.
- **insert, update, delete** - other types of transactions are written to database logs, because they are modifying the database.

Synchronization

Every time a transaction queue is empty (all pending transactions has been applied), a new version of a database is broadcasted to all connected users. When we receive information about the new database version, we load its root. Each database version has information from which version was created. Database versions create an append-only log called database log. Database log is an immutable, operation-based conflict-free replicated data structure (CRDT¹) for storing database versions. Every version in the log is saved in IPFS and each points to a hash of previous version(s) forming a graph. Database logs can be forked and joined back together.[?]

If more than one peer publishes a new version of the database that has been created from the same database version, other peers need to decide which version they would accept. Opposite to records conflicts, we need to solve database versions conflict fast and automatically. There are multiple ways to solve them. For example:

- **The biggest hash wins** - if there is more than one database version at the same (discrete) time, the one with the biggest hash wins. This strategy is present in Figure 5.2.
- **The longest connecting time wins** - a user who is connected to the database for the longest time wins.

Lots of more strategies can be implemented. But they need to be deterministic and as fair as possible. If we have access to geolocation, we can implement a strategy based on a distance to the North pole (the closest node to the North pole wins).

Queries

A Database system provides query language for performing selects. A query consists of:

- **conditions** - query may contain multiple conditions. There are logical operators (and - conjunction, or - disjunction) between conditions. When it is possible (an index is available on condition property), we use indexes for evaluating conditions. If there are more than one conditions, we use an intersection between and's conditions and union between or's a disjunction. These types of conditions are supported:
 - **equal** - record property equals to specified value,
 - **greater than** - record property is greater than specified value,
 - **less than** - record property is less than specified value,
 - **between** - record property is greater than specified minimal value, and less than specified maximal value.
- **filters** - filters are similar to conditions, but they can be more complicated. They are functions returning boolean that are being applied to query results. If any filter return **false** for the query result, the query result is ignored. A query can contain several filters.
- **offset** - offset saying how many query results should be ignored from the beginning.
- **evaluator** - we use evaluator for accessing results of the query.

¹https://en.wikipedia.org/wiki/Conflict-free_replicated_data_type

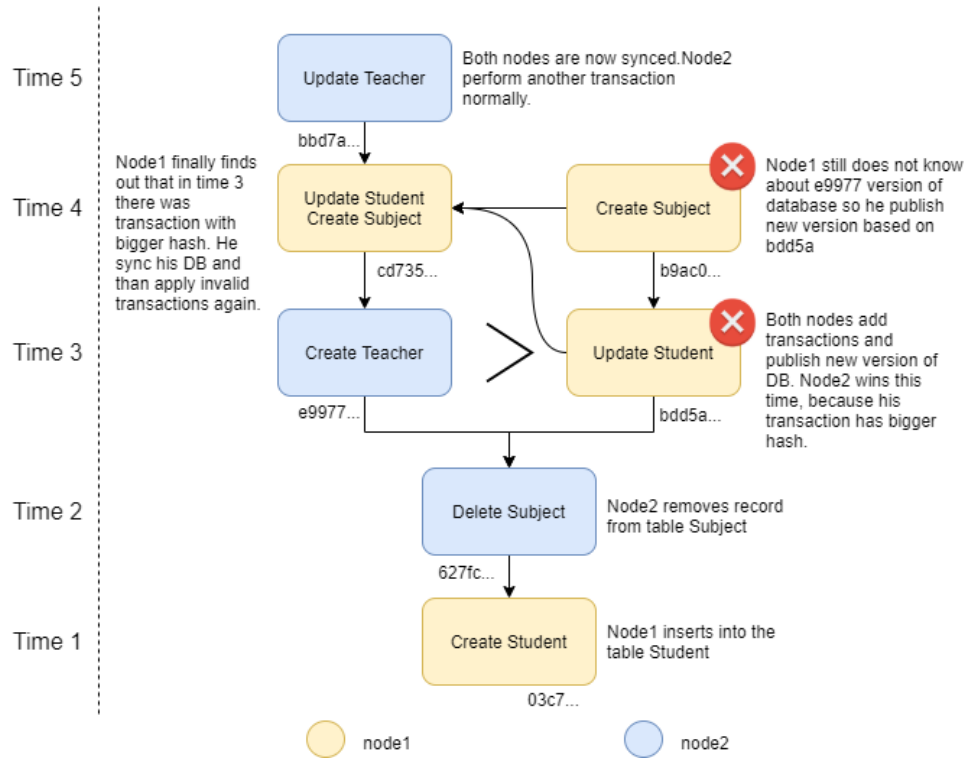


Figure 5.2: Synchronization of a database. Note that in time 4, Node1 add two transactions in one DB version.

- **all** - returns array of all results of the query,
- **first** - returns only first result,
- **take** - returns N number of results where N is an argument,
- **paginate** - returns page of results. Accept two arguments. **perPage** specifies the number of results in one page. **page** is number of requested page.
- **iterate** - returns iterator that can be used in cycles such as **for** or **while**.

5.2 System components

The system consists of one or more Feeders and Explorers. Feeders are connected to data sources and provide synchronization with cryptocurrency blockchains. Explorer can request data from the system network and display them to a user (see Figure 5.3).

5.3 Feeder

A Feeder is a service that stores data in IPFS and indexes them in our database system. Feeder stores data in an acyclic structure (see Figure ??). Every block has links to transactions that are confirmed in this block. An address also has links to all its transactions. This scheme allows storing blockchain data in IPFS in small objects with size less than 256kB (a limit for storing objects directly in DHT). Cyclic references (for example between next/previous blocks) are not possible, because links are changing CID of an object. Therefore if we store object A with a link to the object B, CID of A will change. Then if we add

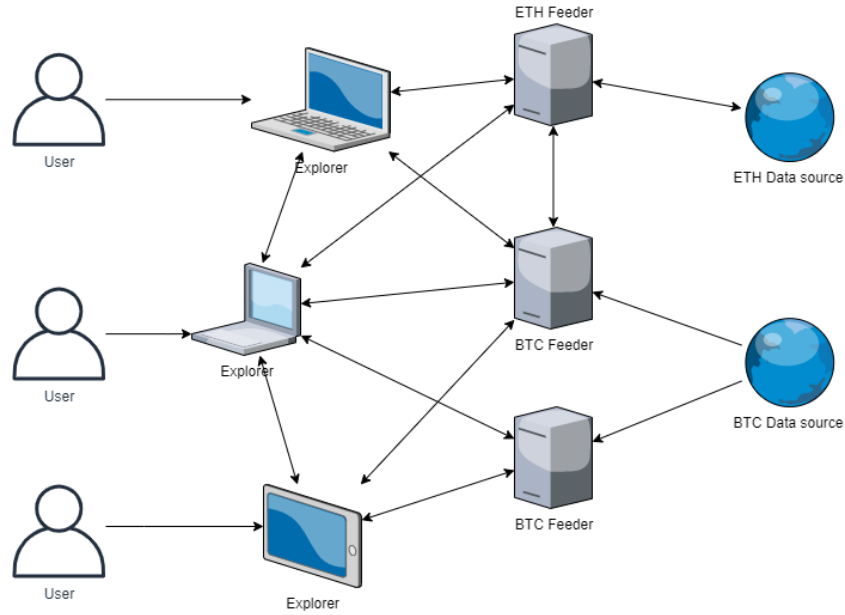


Figure 5.3: System architecture

backlink from B to A, CID of B will be changed and so A has now link to an old version of B. If we update link of object A to point on the new version of object B, CID of object A will change and therefore object B is now pointing to the old version of object A. From this example, it appears that there is no way of making cyclic references in IPFS.

A Feeder can use different sources for obtaining blockchains data. It can be connected directly to the blockchain as a full node or to some other data source providing blockchain data as Blockbook² or Insight³.

[[Describe used indexes and database structure that is present in feeder]]

5.4 Explorer

Explorer can perform basic queries like a search for block by its height or hash and search address and transaction by hash. Nevertheless, Explorer can also make more complex queries for example, get the first 20 transaction where the sum of inputs is more than some value, or get transactions between some time interval. Explorer is an application that can be used in two environments. In browser with graphical user interface, or in node.js⁴ as a RestAPI.

5.4.1 ExplorerGUI

ExplorerGUI is browser version of Explorer. It is implemented as a SPA⁵ to prevent restarting connection with IPFS after each time a user visit different page. It has separated views for block, transaction and address details. Also paginating, filtering and sorting objects by all its properties is supported.

²<https://github.com/trezor/blockbook>

³<https://insight.is>

⁴<https://nodejs.org/>

⁵https://en.wikipedia.org/wiki/Single-page_application

5.4.2 ExplorerAPI

ExplorerAPI is a node.js application provides several endpoints for obtaining the data:

- GET /tx - get collection of transaction,
- GET /tx/[:txHash] - get transaction by hash,
- GET /block - get block collection,
- GET /block/[:blockHeightOrHash] - get block by height or hash,
- GET /address - get address collection,
- GET /address/[:hash] - get address by hash.

All endpoints support more detailed path specification by providing multiple parameters at the end of the path (for example, it is possible to request only a block hash by calling /block/42/hash). Also, every endpoint supports query parameters

- **filter** where value of this parameter can be simple arithmetic expression like **height>50**,
- **limit** parameter for pagination with default value set to 20,
- **orderBy** sorts results of a sequence in ascending order.

Chapter 6

Implementation

This chapter describes the implementation details of the system, shows the internal architecture. As mentioned in the previous chapter, the program was created in Typescript, and they are using js-ipfs¹ implementation of IPFS node. This selection of tools enables code sharing between these two separated applications.

6.1 Explorer-core implementation

Explorer-core is the most complex module of a whole system with more than 5000 lines of code. The database is a main part of the explorer-core module. Database system consists of a query system, indexes and a database abstract layer.

6.1.1 Indexes

For indexing there is currently implemented only B-tree, but different structures such as tries² can be implemented easily. They only need to implement index's interface (function such as `insert`, `delete`, `update`, `find`). We can create indexes on database entities with decorators³, which are part of ECMAScript 6 standard. Every database entity has to have `PrimaryKey` decorator on property that is used as a primary index. Index has three parameters:

- **comparator** - is a function that accepts two arguments and returns number that is less than zero if first arguments is greater than second, zero if arguments are the same, more then zero if a second argument is greater than first. If a user has not set any custom comparator default one $((a, b) \Rightarrow a < b ? -1 : +(a > b))$ is used. This default comparator works on atomic keys such as string or numbers.
- **keyGetter** - is another function that accepts a whole entity as arguments and returns key, that is used in an index. A default key getter is function that returns property of the entity (for example default key getter for property `height` is $(e) \Rightarrow e["height"]$).
- **branching** - the branching factor is the number of children at each node, the out-degree. Default branching factor is 16.

¹<https://github.com/ipfs/js-ipfs>

²<https://en.wikipedia.org/wiki/Trie>

³<https://www.typescriptlang.org/docs/handbook/decorators.html>

The B-tree structure (see Figure 6.1) is optimized for IPFS. Unlike PostgreSQL B-tree⁴, lead nodes haven't got left and right sibling references (due to cyclic reference this is not possible in IPFS). And therefore we don't have to store data only in leaves, but we can store data in nodes itself. This approach lead to better performance when executing queries. In insert queries, statistically fewer nodes needs to be updated. Also in search queries, there is a chance that we found a key in some non-leaf node, and therefore we would need fewer node visits.

All search queries (equal, less, greater, between) has two steps:

- **Find subtree** - First, we need to find minimal subtree that contains all results. We start in a root node of the B-tree. Then, we use an index's **comparator** function to determine if we can visit some child node and still have all the result in it. If not, we found a minimal subtree for the query.
- **Traverse** - To get results of the query we can traverse minimal subtree in two direction. In-order for **greater than** and **between**. For **less than** reversed in-order. With **equal** it does not matter.

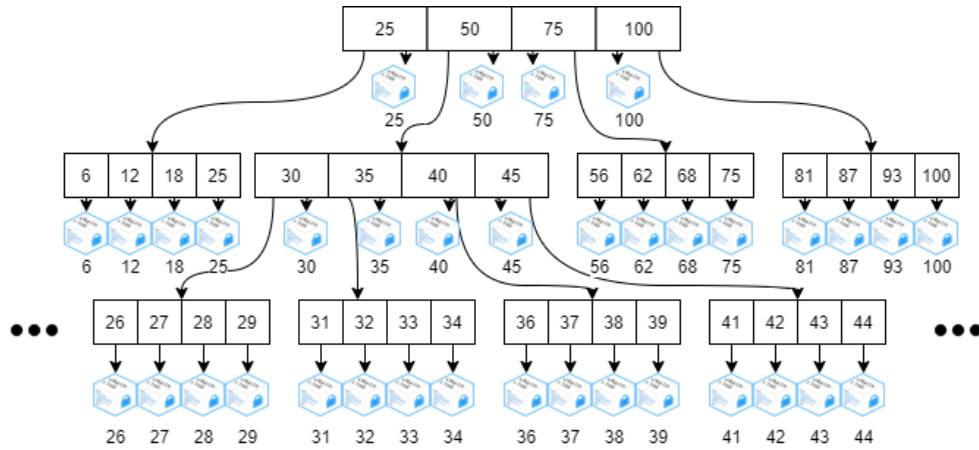


Figure 6.1: Btree indexing first 100 block by their height

6.1.2 Query system

Database system offers a complex query system. A query can consist of multiple conditions and the **Query Planner** is responsible for resolving them. It decides which indexes are used for query and choose a strategy. If there is no condition, a primary key is used for query execution. In the case of a single condition, **Query Planner** checks if there is an index on the condition's property. If yes, then this index is used to perform the query. Else condition is transformed to filter and a primary key is used to obtain results which are then filtered. For multiple conditions connected with logical operators **AND** or **OR**, **Query Planner** creates **OR-hashset** and **AND-hashset**. **AND-hashset** is initialized with the results of a condition which has the smallest number of results and uses **AND** operator. Then we check for each result hash for all other **AND** conditions if **AND-hashset** contains it. If no, a hash is deleted from **AND-hashset**. This creates intersection between all **AND** conditions. The **OR-hashset**

⁴<https://github.com/postgres/postgres/tree/master/src/backend/access/nbtree>

is created empty. Then, we add a hash of every result of all **OR** conditions to it. This creates union between **OR** conditions. At the end, we perform union between **AND-hashset** and **OR-hashset**. This creates final hashset which is used later to obtain actual data with resolvers such as **all**, **first**, **paginate** etc. If one of the **AND** conditions has great selectivity (it has less than a hundred results), **Query Planner** can decide to don't use other indexes, but transform condition to filter and cycle over the results. Example query returning blocks between 38 and 42 is shown in Figure 6.2. A peer needs to download only data that are highlighted in green. After downloading them, they are stored on peer's filesystem and are offered to other peers.

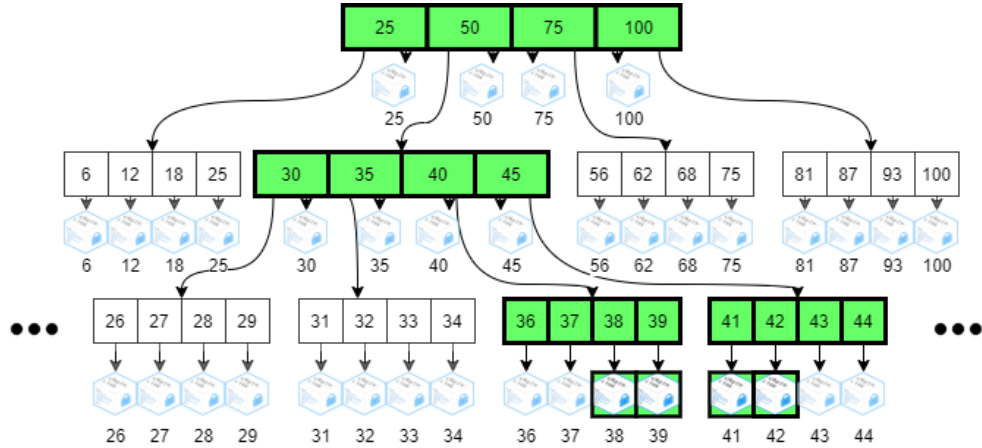


Figure 6.2: Data access for performing query that returns blocks that have a height between 38 and 42.

6.1.3 Database

The most important part of the explorer-core module is a Database system which connects query system with indexes. When we make a query, it translates to the database, and data are obtained from the indexes.

Tables

A database contains tables. A table contains indexes and table name and implements operations that modifies its state:

- **insert** - creates history log for entity with first entity version. Then adds it to every table index.
- **update** - adds new version of the entity to its history log. Than update every index of the table.
- **delete** - removes entity from every table index. From now, an entity can not be find in this table.

Example of the class **User** is in Figure 6.3. It has two indexes. One primary on property **name** and one normal index on property **age**. The second index has also specified **comparator** that is a bit faster than default one, but works only on numbers, and **keygetter** that returns a year when a user has born.

```

class User extends Queriable<User> {
    @PrimaryKey()
    name: string;

    @Index(
        (a, b) => a - b,
        u => new Date().getFullYear() - u.age,
    )
    age: number;
}

```

Figure 6.3: Example of class.

Transactions

A database can be executing only a single transaction at the time to prevent data inconsistency. For that reason, we implemented a transaction queue where transactions are stored before executing in the order in which they came. Executing more transactions in a row is significantly more effective than executing them one by one. If a transactions queue has only one transaction, it will wait 100ms for more transactions to come.

Synchronization

There are two types of transactions. Those that changes the database state (**create**, **update**, **delete**) and those that don't (**read**). Every transaction that changes database state need to be synchronized with other peers. We created Database log for that. It is an append-only log with a discrete time. There can be only one valid database transaction in every point in time. There are several strategies to choose which transaction is globally accepted and which transactions need to rollbacks (described in design [\[\[add link\]\]](#)).

6.2 Feeder implementation

A Feeder is a simple program (less than a hundred lines of code) written in Typescript (see Algorithm 1). It strongly depends on the Explorer-core module. Currently, we support only Blockbook connector as a source of blockchain's data, but Feeder can be simply expanded to support more data source such as InsightAPI or direct connection to a blockchain as a full node. Each Feeder has configuration file (usually called `.env`) with Feeder's settings. Main Feeder settings are `URL` of the source for the blockchain's data and `DB_NAME` which is a name of the database where Feeder inserts new blocks. A Feeder can be connected to only one blockchain. We can create multiple Feeders for single blockchain, but we need to provide some deterministic algorithm, that insures that each block is parsed by only one Feeder. This can be done by providing `FeederId` and `FeedersCount` in `.env` config file. If

those values are provided, Feeder parses only blocks that have height modulo `FeedersCount` equals to `FeederId`.

Algorithm 1: Simplified Feeder algorithm

```

load configuration;
while there is new block do
    fetch block;
    for transaction in block do
        | save transaction;
    end
    save block;
end

```

6.3 Explorer implementation

We can split Explorer to several modules. The central part is ExplorerCore which communicates with IPFS and provides an interface to explore blockchain and to perform queries. There can be ExplorerGUI or ExplorerAPI on the top of the ExplorerCore. These modules are providing communication with user or other systems. This architecture is shown in Figure 6.4.

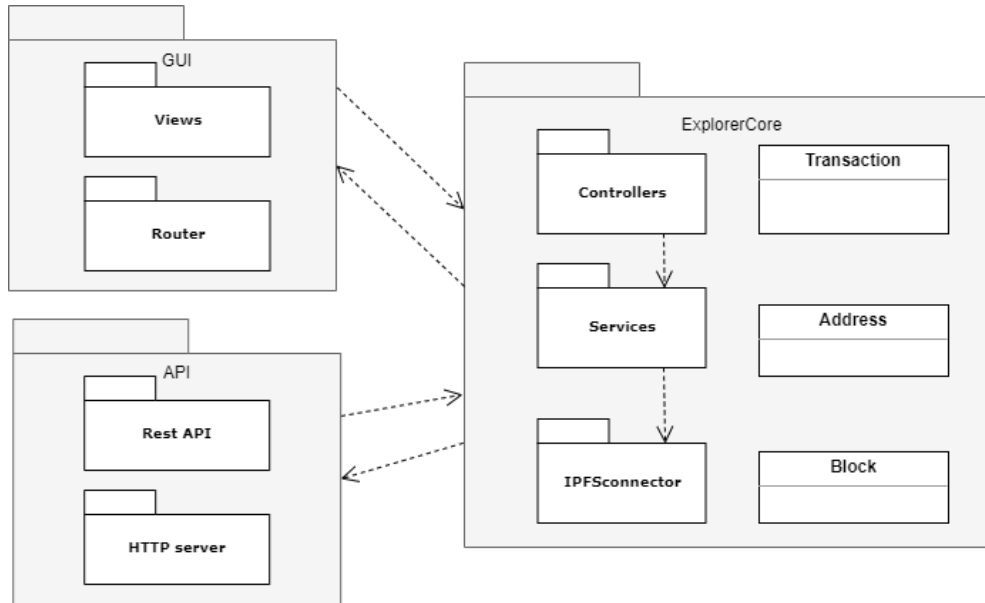


Figure 6.4: Explorer architecture

6.3.1 ExplorerCore

ExplorerCore is main module of the Explorer. It is responsible for obtaining and resolving data from IPFS and performing queries. It has simple interface:

- `getObject(CID[, path])` - gets object from IPFS with optional path,
- `getTx(txHash)` - gets transaction by its hash,

- `getBlock(blockHash)` - gets block by its hash,
- `getBlock(blockHeight)` - gets block by its height,
- `getAddress(addressHash)` - gets address by its hash.

All ExplorerCore methods except `getObject` are returning `QueryObject`.

Query system

Explorer supports a simple query system. Every queriable object (address, transaction, block) has methods:

- `where(propertyName)` - create a condition on the property. There can be multiple conditions in one query. A condition has to be followed by one of the functions:
 - `gt(value)` - property is greater or equal than `value`. The query finds in an index first object that has property (set by `propertyName` in condition) equal or greater than `value`, and traverse index to the right (to bigger objects).
 - `lt(value)` - property is less or equal than `value`. Similar as in the `gt` function, the query finds the closest object that has property equal or less than `value`. Then, the query traverses index to the smaller objects with smaller index value (to the left).
 - `between(min, max)` - property is greater or equal than `min` and less or equal than `max`.
- `offset(offsetValue)` - query will skip `offsetValue` number of results,
- `limit(limitValue)` - set maximum number of results. After query has `limitValue` count of matched objects it will stop browsing the index,
- `all()` - return all objects that matched query,
- `first()` - return the first object that matched a query,
- `and(childQuery)` - logical and between two queries. Parent query resolves `childQuery` (call `all()` function) and makes a hashtable from its results. Later, when parent query being resolved, query checks if the result exists in the hashtable for every result.
- `or(childQuery)` - logical or between two queries. `childQuery` will be resolved by parent query, and its results stores in hashtable. When the parent query is resolved, its results will be added to the hashtable. This will remove duplicates in the final results.
- `[Symbol.iterator]()` - iterator that is used in cycle `for (result of query)`.

The data obtaining will happen only when we access the data inside the `Query` object (for example, in for-loop) or by calling `all()` or `first()` method. Index is used only for the first condition in a query, that is on the property, that has created index. Other conditions are performed client-side.

6.3.2 ExplorerGUI

ExplorerGUI is a single page application with simple user interface that runs in browser implemented with Vue.js⁵. Browser's implementation of IndexedDB is used as a storage for IPFS as can be seen in Figure 6.5. Communication with other peers is provided though WebRTC⁶ or WebSockets. Every tab opened in the browser is the same IPFS node. Opening a new tab in Incognito mode or different browser will spawn different IPFS node.

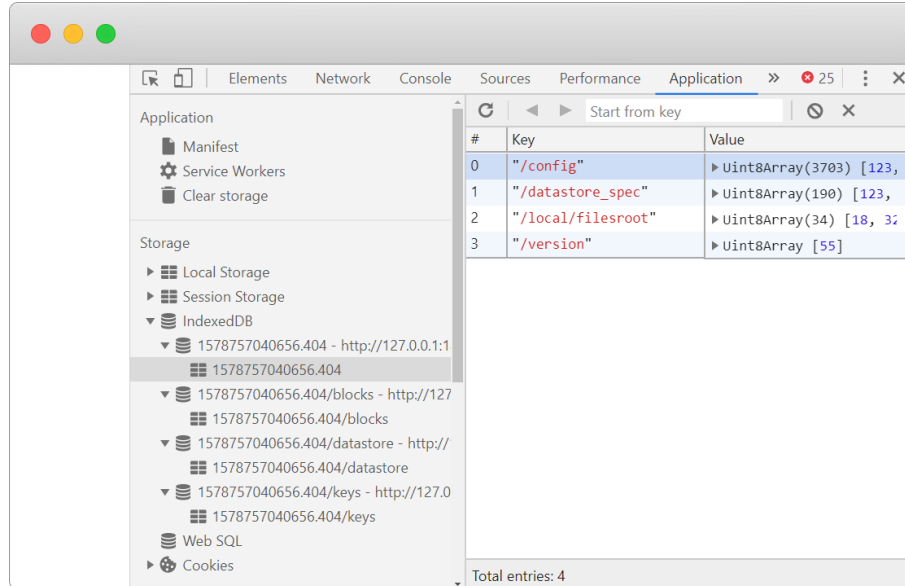


Figure 6.5: IPFS storage in browser

6.3.3 ExplorerAPI

Node.js⁷ implementation of IPFS uses filesystem to store data (see Figure 6.6). On the top of ExplorerCore, there is a simple web framework Express⁸, where are routes (endpoints) defined. Currently available routes are described in Section 5.4.2. Every route supports query parameters **filter** (used for filtering results) and **limit** (which reduces number of results). Pagination can be made by setting **filter** to be greater/smaller (depending on ordering) as key of the last displayed object and **limit** to page size. For example, if a user is looking at page of transactions ordered by time (ordered from the newest transactions to the oldest), the next page of transactions are first N transactions that happened before last displayed transaction (where N is page size).

Rest API supports optional path param **path** that is useful for traversing objects in IPFS. If we want to get fifth transaction of the block with height 1000 one of the way is request URL `/block/998/next_block/next_block/txs/5` (get block with height 998, get next block two times, get transaction, a get fifth item from array of transactions). This approach allows the user to explore IPFS storage as graph very quickly by objects links.

⁵<https://vuejs.org/>

⁶<https://webrtc.org/>

⁷<https://nodejs.org/>

⁸<http://expressjs.com/>

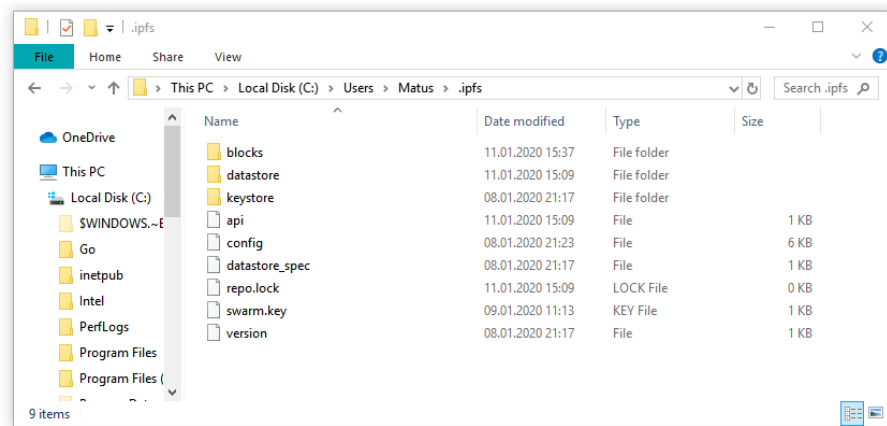


Figure 6.6: IPFS storage in node.js

Chapter 7

Benchmarking and testing

7.1 Benchmarking enviroment

7.1.1 Blockbook

Blockbook¹ is a blockchain indexer for Trezor Wallet², developed by SatoshiLabs³. It currently supports more than 30 coins (and the community implemented some others). For data storage Blockbook is using RocksDB⁴ developed by Facebook, which is a NoSQL database that stores only key-value pairs. Blockbook is providing fast API for accessing blocks, addresses and transactions. Main limitations of blockbook:

- **Not distributed** (client-server architecture) - problem with scaling for more users.
- **Not a SQL database** - it does not have a relational data model, it does not support SQL queries, and it has no support for indexes.
- **Single-Process** - only a single process (possibly multi-threaded) can access a particular database at a time.

7.1.2 Writing

7.1.3 Reading

7.1.4 Filtering

¹<https://github.com/trezor/blockbook>

²<https://wallet.trezor.io/>

³<https://satoshilabs.com/>

⁴<https://github.com/facebook/rocksdb/wiki>

Chapter 8

Conclusion

The goal of this project was to create a system that uses IPFS to explore blockchain. To achieve this, we needed to create our own decentralized and distributed database system that supports queries and indexes on top of the IPFS. Feeder and ExplorerCore communicate with this database system to store (Feeder) or retrieve data (ExplorerCore). To visualize data in GUI, ExplorerGUI can be used, and for obtaining data with RestAPI there is ExplorerAPI.

The research section of this project focuses mainly on the IPFS principles (content-based addressing, nodes linking, and others). Selected cryptocurrencies for exploring are also briefly discussed. Great effort was devoted to designing the whole system and creating a functional prototype.

In the next semester, we want to highly improve database system by support more operators (`like`, `groupBy`). We want to research and compare more indexing strategies and make performance testing and profiling. Also, comparison with Blockbook in the same conditions would be interesting. There is high potential by using IPFS in this system to create new use cases that no other system for exploring cryptocurrencies blockchain supports. For example, with some heuristic, we believe, that with this system, we can follow crypto coins after exchange for another crypto coins (thanks to linking between objects in the IPLD layer).

Bibliography

- [1] BENET, J. IPFS - Content Addressed, Versioned, P2P File System. *CoRR*. 2014, abs/1407.3561. Available at: <http://arxiv.org/abs/1407.3561>.
- [2] DASCANO, M. *Digibyte: An Easy Guide to Learning the Essentials*. 1stth ed. North Charleston, SC, USA: CreateSpace Independent Publishing Platform, 2018. ISBN 1725629178.
- [3] DAVIS, J. The Crypto-Currency. *The New Yorker*. october 2011. Available at: <https://www.newyorker.com/magazine/2011/10/10/the-crypto-currency>.
- [4] DIAS, D. and BENET, J. Distributed Web Applications with IPFS, Tutorial. In: BOZZON, A., CUDRE MAROUX, P. and PAUTASSO, C., ed. *Web Engineering*. Cham: Springer International Publishing, 2016, p. 616–619. ISBN 978-3-319-38791-8.
- [5] HABER, S. and STORNETTA, W. S. How to time-stamp a digital document. In: Springer. *Conference on the Theory and Application of Cryptography*. 1990, p. 437–455.
- [6] JEPSON, C. DTB001: Decred Technical Brief. Available at <https://coss.io/documents/white-papers/decred.pdf> Additional information available at <https://www.decred.org>. 2015.
- [7] KALLE, K. *Big data in video games*. Lappeenranta, FI, 2017. Bachalar Thesis. Lappeenranta University of Technology, School of Business and Management, Computer Science. Available at: <http://lutpub.lut.fi/handle/10024/147666>.
- [8] MAYMOUNKOV, P. and MAZIERES, D. Kademlia: A peer-to-peer information system based on the xor metric. In: Springer. *International Workshop on Peer-to-Peer Systems*. 2002, p. 53–65.
- [9] NARAYANAN, A., BONNEAU, J., FELTEN, E., MILLER, A. and GOLDFEDER, S. *Bitcoin and cryptocurrency technologies: a comprehensive introduction*. Princeton University Press, 2016.
- [10] NOETHER, S. Ring Signature Confidential Transactions for Monero. *IACR Cryptology ePrint Archive*. 2015, vol. 2015, p. 1098.
- [11] WAYNER, P. *Digital cash: Commerce on the net*. Academic Press Professional, Inc., 1997.
- [12] WOOD, G. et al. Ethereum: A secure decentralised generalised transaction ledger. *Ethereum project yellow paper*. 2014, vol. 151, no. 2014, p. 1–32.