

# BlockSci

Design and applications of a blockchain analysis platform

Harry Kalodner, **Malte Möser**, Kevin Lee, Steven Goldfeder, Martin Plattner, Alishah Chator, Arvind Narayanan

# Why build a blockchain analysis tool?

## Blockchain data

- Unprecedented research corpus of financial transactions
- Interesting for scientific analyses and commercial applications

## Tools to analyze blockchain data

- Commercial tools are often tailored towards specific use cases
- Lack of general-purpose tools

# Blockchain data is “small”

50 GB

Bitcoin transaction graph (efficient)

120 GB

EC2 instance < \$1000 per month

260 GB

Bitcoin blockchain

Max RAM available on an EC2 instance: 24 TB

# Goals for BlockSci

## Performance

---

in-memory database

domain-specific optimizations

## Capabilities

---

supports different blockchains

analytic tools included

## Usability

---

C++ / Python

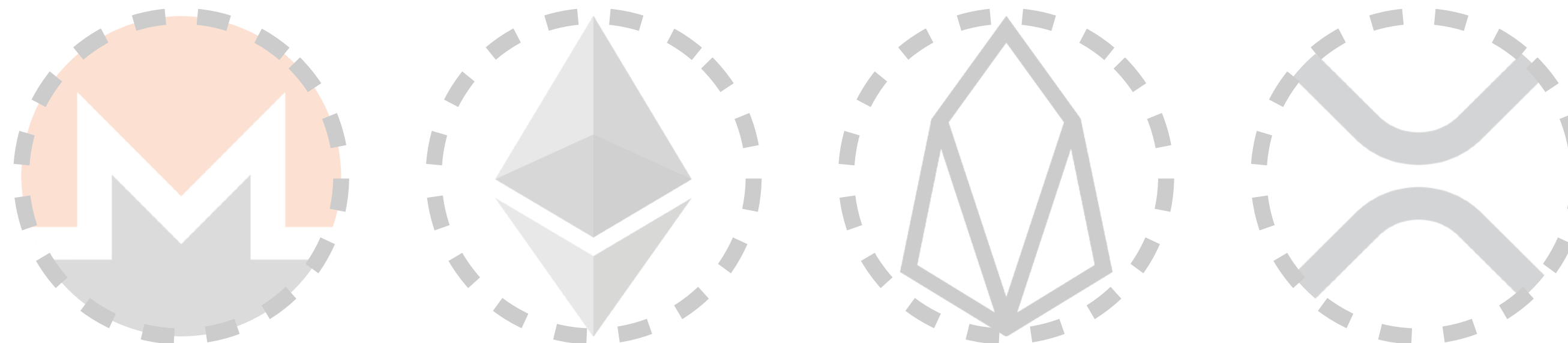
fluent interface

# Which blockchains should we support?

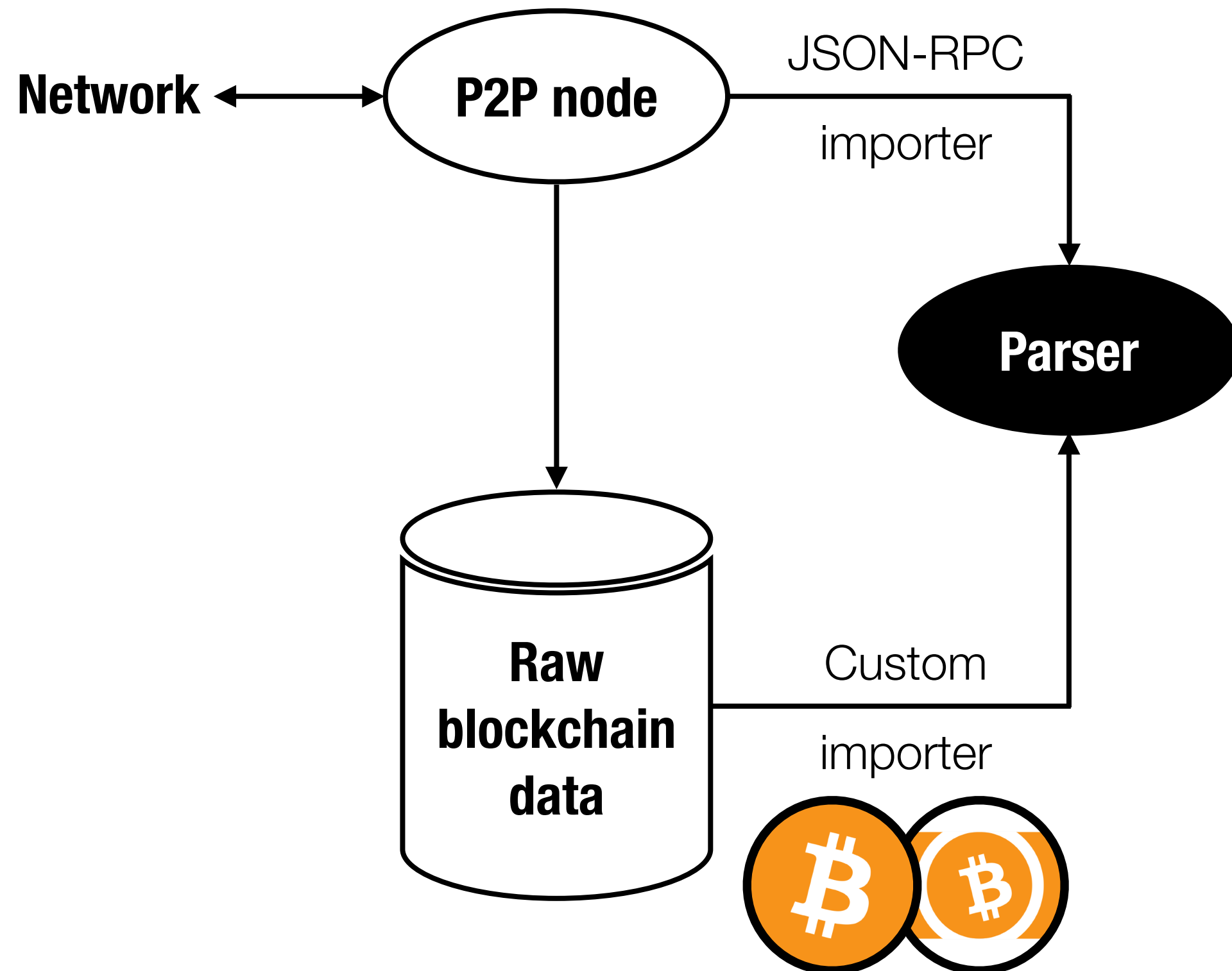
**Supported:** design similar to Bitcoin (to varying degree)



**Unsupported:** different design and/or interface

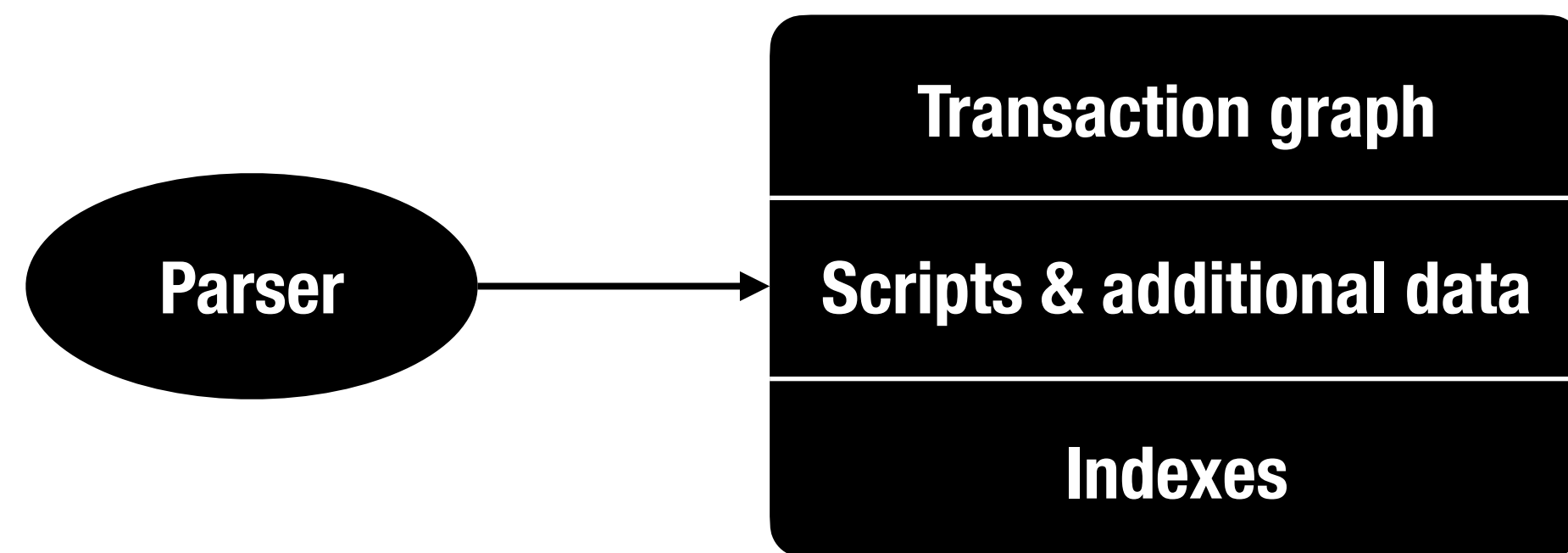


# Building an efficient parser



- Incremental updates
- Parallelized pipeline to resolve arbitrary transaction ordering in blocks
- Optimized address lookups
- Multi-chain mode

# Data layout: tradeoff between memory efficiency and performance



*BlockSci Data*

- Core transaction graph data optimized for sequential analyses
- Additional data stored in hybrid format, loaded on-demand
- Address data is deduplicated across different address types
- Indexes for common lookups

# BlockSci is fast

Iterating over transaction inputs and outputs in C++

---

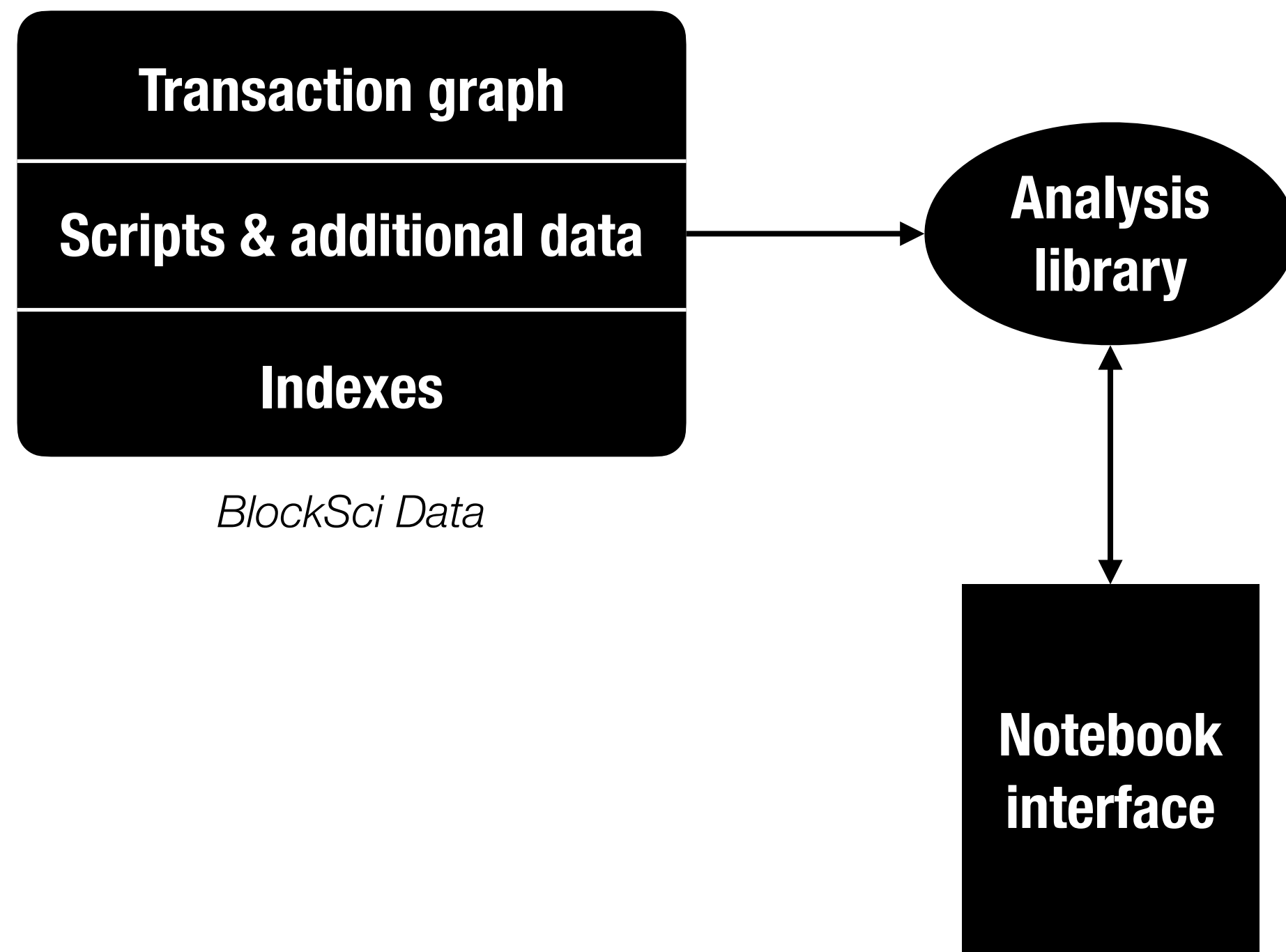
**11.3** seconds      single-threaded

**0.9** seconds      multithreaded (16 vCPUs)

Iterating over 610,695 Bitcoin blocks on an r5.4xlarge EC2 instance (16 vCPUs, 2.5 GHz Intel Xeon Platinum 8175M, 128 GiB memory, 800 GiB EBS volume)

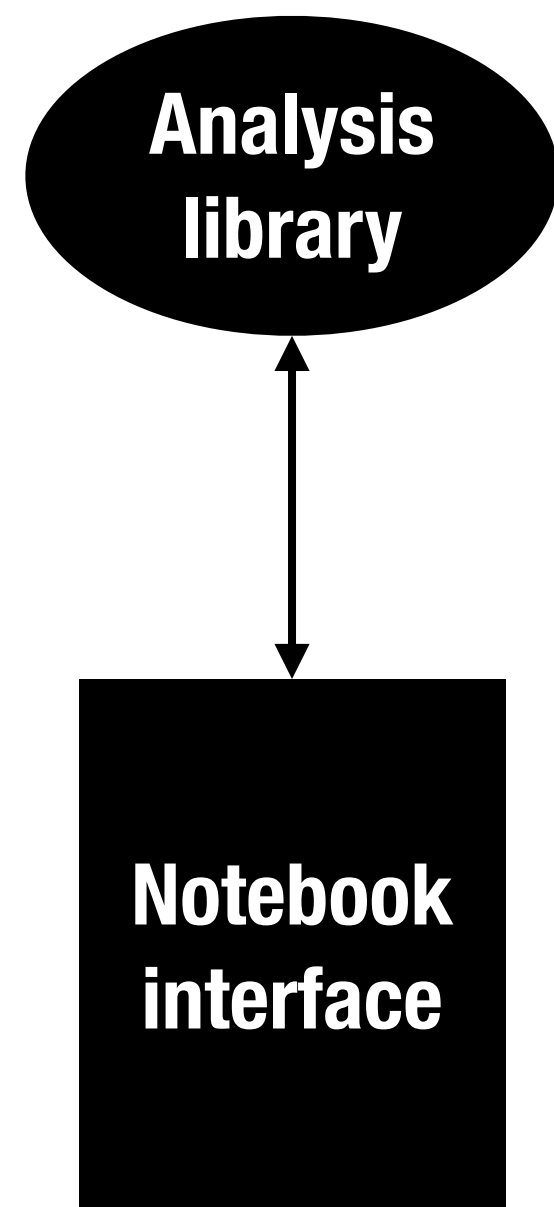


# Analysis library



- Provides static view on data that can be updated on disk
- Map-reduce abstractions
- Forensic capabilities (clustering, change address heuristics)
- Can collect P2P network timestamps

# Programmer interface



## Querying the blockchain in Python

```
[tx for block in chain for tx in block if tx.fee > 10**7]
```

## Fluent interface: expressiveness & speed

- Queries are specified in Python, executed in C++
- Many operators: select, filter, group\_by, min/max, ...

```
chain.blocks.txes.where(lambda tx: tx.fee > 10**7).to_list()
```

# Example application

**Unfortunate effects of multisignature use**

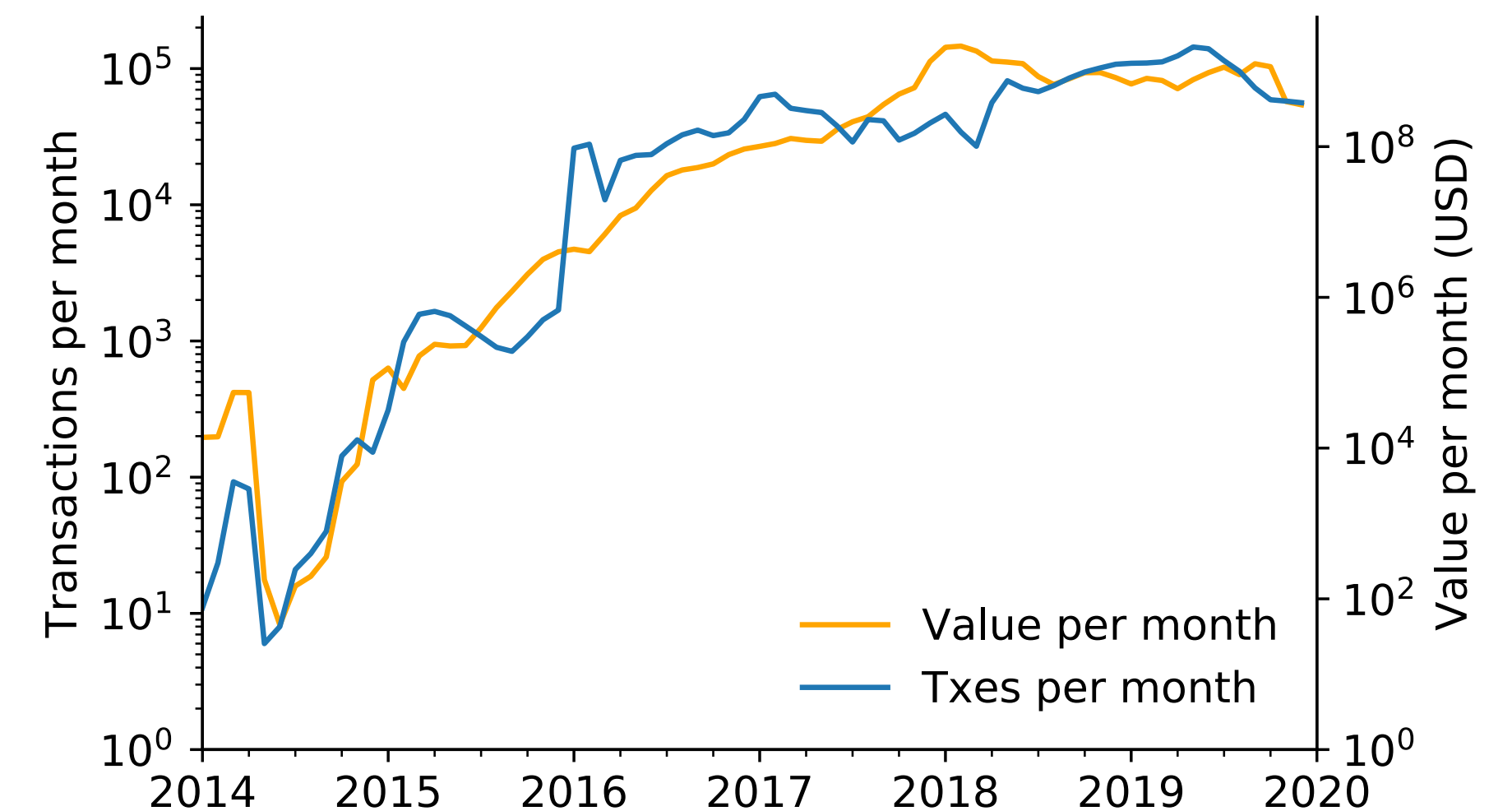
# Example: drawbacks of multisig

## Confidentiality

- Transactions specify list of keys and threshold
- Publicly exposes access control changes

Privacy

Security



# Example: drawbacks of multisig

## Confidentiality

## Privacy

- Different output type allows to distinguish multisig from single-key users
- May identify change output in up to 122 million transactions

## Security

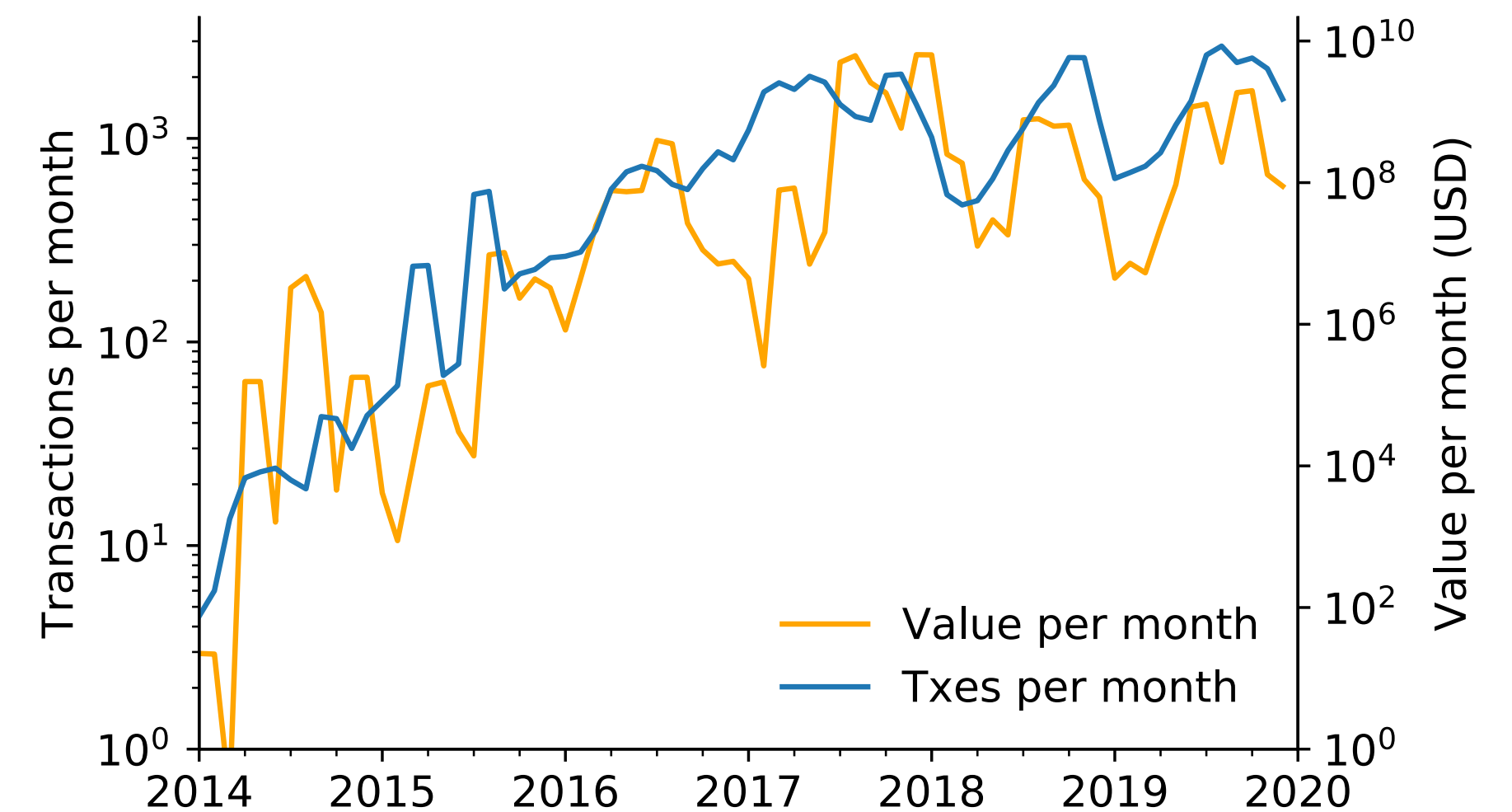
# Example: drawbacks of multisig

Confidentiality

Privacy

**Security**

- Pattern: Multisig → Regular → Multisig
- Temporary reduction in security



# BlockSci has been widely used

- At least 9 peer-reviewed articles and 6 preprints
- Educational tool for workshops and lectures
- Foundation for more specialized blockchain analysis tools
- Active on GitHub

When the cookie meets the blockchain:  
Privacy risks of web payments via cryptocurrencies

Steven Goldfeder\*, Harry Kalodner\*, Dillon Reisman<sup>†</sup>, Arvind Narayanan\*  
Princeton University  
\* {stevenag, kalodner, arvindn} @cs.princeton.edu  
<sup>†</sup> dillon@lonlon.io

2018 IEEE Symposium on Security and Privacy

Tracking Ransomware End-to-end

Danny Yuxing Huang<sup>1</sup>, Maxwell Matthaios Aliapoulos<sup>2</sup>, Vector Guo Li<sup>3</sup>  
Luca Invernizzi<sup>4</sup>, Kylie McRoberts<sup>4</sup>, Elie Bursztein<sup>4</sup>, Jonathan Levin<sup>5</sup>  
Kirill Levchenko<sup>3</sup>, Alex C. Snoeren<sup>3</sup>, Damon McCoy<sup>2</sup>

<sup>1</sup> Princeton University <sup>2</sup> New York University <sup>3</sup> University of California, San Diego <sup>4</sup> Google Inc <sup>5</sup> Chainalysis

Safeguarding the Evidential Value of Forensic Cryptocurrency  
Investigations<sup>\*,\*\*</sup>

Michael Fröwis<sup>a,\*</sup>, Thilo Gottschalk<sup>b</sup>, Bernhard Haslhofer<sup>c</sup>, Christian Rückert<sup>d</sup> and  
Paulina Pesch<sup>b</sup>

<sup>a</sup>University of Innsbruck, Technikerstr. 21a, 6020 Innsbruck, Austria

<sup>b</sup>Karlsruhe Institute of Technology, Vincenz-Prißnitz-Str. 3, 76131 Karlsruhe, Germany

<sup>c</sup>AIT Austrian Institute of Technology, Giefinggasse 4, 1210 Vienna, Austria

<sup>d</sup>Friedrich-Alexander University of Erlangen-Nuremberg, Schillerstraße 1, 91054 Erlangen, Germany

20



# More details in our paper!

- Architecture and optimizations
- Performance measurements and comparison
- Multi-chain mode
- Three more applications
  - Cross-chain privacy
  - Effectiveness of transaction fee estimation
  - Velocity of cryptocurrencies

ARTIFACT  
EVALUATED



PASSED

BlockSci: Design and applications of a blockchain  
analysis platform

Harry Kalodner\*  
*Princeton University*

Malte Möser\*  
*Princeton University*

Kevin Lee  
*Princeton University*

Steven Goldfeder  
*Cornell Tech*

Martin Plattner  
*University of Innsbruck*

Alishah Chator  
*Johns Hopkins University*

Arvind Narayanan  
*Princeton University*

**Abstract**

Analysis of blockchain data is useful for both scientific re- search and commercial applications. We present BlockSci, an open-source software platform for blockchain analysis. BlockSci is versatile in its support for different blockchains and analysis tasks. It incorporates an in-memory, analytical (rather than transactional) database, making it orders of mag- nitudes faster than using general-purpose graph databases. We describe BlockSci’s design and present four analyses that il- lustrate its capabilities, shedding light on the security, privacy, and economics of cryptocurrencies.

**1 Introduction**

Public blockchains constitute an unprecedented research cor- pus of financial transactions. Bitcoin’s blockchain alone is 260 GB as of December 2019.<sup>1</sup> This data holds the key to measuring the privacy of cryptocurrencies in practice, study- ing user behavior with regards to security and economics, or understanding the non-currency applications that use the blockchain as a database.

We present BlockSci, a software platform that enables the science of blockchains. It addresses three pain points of ex- isting tools: poor performance, limited capabilities, and a cumbersome programming interface. Compared to the use of general-purpose graph databases, BlockSci is hundreds of times faster for sequential queries and substantially faster for all queries, including graph traversal queries. It comes bun- dled with analytic modules such as address clustering, exposes different blockchains through a common interface, collects “mempool” state and imports exchange rate data, and gives the programmer a choice of interfaces: a Jupyter notebook for intuitive exploration and C++ for performance-critical tasks. In contrast to commercial tools, BlockSci is not tailored to specific use cases such as criminal investigations or insights

for cryptocurrency traders. Instead, by providing efficient and convenient programmatic access to the full blockchain data, it enables a wide range of reproducible, scientific analyses.

BlockSci’s design starts with the observation that blockchains are append-only databases; further, the snapshots used for research are static. Thus, the ACID properties of transactional databases are unnecessary. This makes an in- memory analytical database the natural choice. On top of the obvious speed gains of memory, we apply a number of tricks such as converting hash pointers to actual pointers and deduplicating address data, which further greatly increase speed and decrease the size of the data. We plan to scale vertically as blockchains grow, and we expect that this will be straightforward for the foreseeable future, as commodity cloud instances currently offer up to a *hundred times* more memory than required for loading and analyzing Bitcoin’s blockchain. Avoiding distributed processing is further moti- vated by the fact that blockchain data is graph-structured, and thus hard to partition effectively. In fact, we conjecture that the use of a traditional, distributed transactional database for blockchain analysis has infinite COST (Configuration that Outperforms a Single Thread) [1], in the sense that no level of parallelism can outperform an optimized single-threaded implementation.

BlockSci comes with batteries included. First, it is not limited to Bitcoin: a parsing step converts a variety of blockchains into a common, compact format. Currently sup- ported blockchains include Bitcoin, Bitcoin Cash, Bitcoin SV, Litecoin, and Zcash (Section 2.1). A multi-chain mode optimizes for user-friendly and memory-efficient analyses of forked blockchains together with their parent chain. Smart contract platforms such as Ethereum are outside our scope.

Second, BlockSci includes a library of useful analytic tools, such as identifying special transactions (e.g., CoinJoin) and linking addresses to each other based on well-known heuris- tics, including across forked chains (Section 2.4). Third, BlockSci can record the time of transaction broadcasts on the peer-to-peer network and expose them through the same interface. Similarly, we make (historical and current) data on

<sup>\*</sup>These authors contributed equally to this work.

<sup>1</sup>All numbers in this paper are current as of December 2019, and analyses of the Bitcoin blockchain as of block height 610,695, unless stated otherwise.

16



# Thank you

Code & questions

<https://github.com/citp/BlockSci>

Documentation

<https://citp.github.io/BlockSci>

Contact

[blocksci@lists.cs.princeton.edu](mailto:blocksci@lists.cs.princeton.edu)

*Technical & how-to questions*



## BlockSci

Design and applications of a blockchain analysis platform

Harry Kalodner, **Malte Möser**, Kevin Lee, Steven Goldfeder, Martin Plattner, Alishah Chator, Arvind Narayanan