# Bitcoin User Identification Research & Transaction Data Visualisation

CS-8903-L03 Special Problem with Prof. Ling Liu

Shubhi Agarwal

# Bitcoin

| Market Price (USD) | Average Block Size | Transactions per Day | Mempool Size |
|---|---|---|---|
| **$7,202.72** | **1.09** | **310,205** | **595,712** |
| USD | Megabytes | Transactions | Bytes |
| Average USD market price across major bitcoin exchanges. | The 24 hour average block size in MB. | The aggregate number of confirmed Bitcoin transactions in the past 24 hours. | The aggregate size of transactions waiting to be confirmed. |

# Motivation

- 480 million transactions in last 8 years and around 3 million transactions per day
- Widely considered anonymous due to use of pseudonyms
- Makes fraud detection hard for regulators
- All transaction data publicly available
- Dire need for analysis of the data for link prediction and classification
- Analysis of privacy by user identification
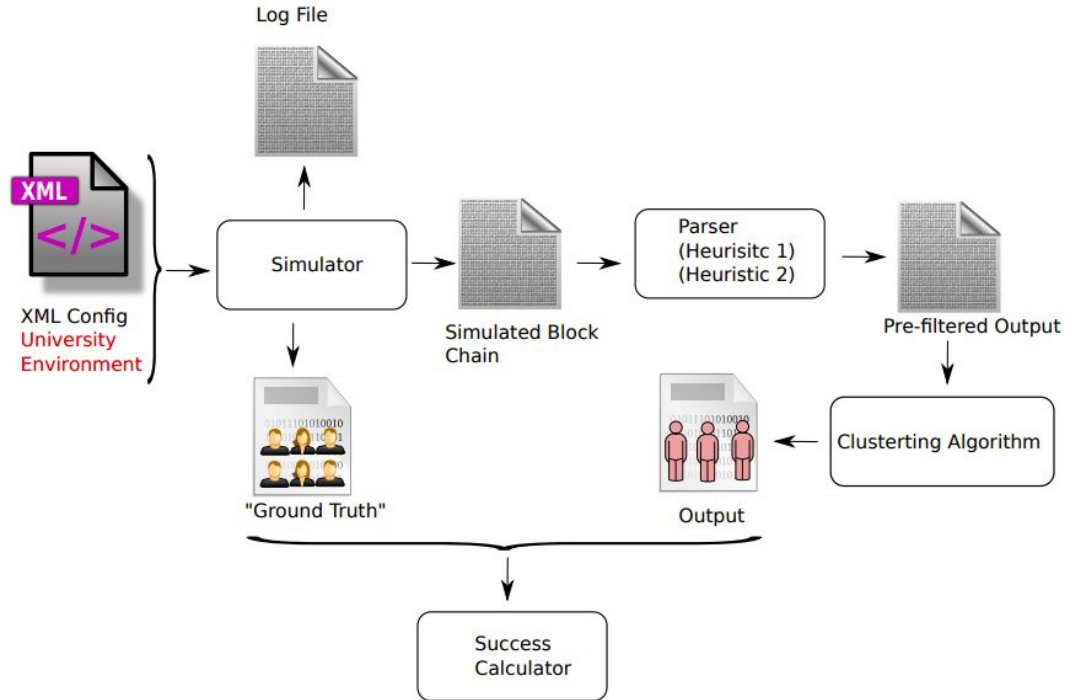
# Part - I

**Research on user identification problem**

# Related Work

- Works related to statistical characterization of Bitcoin transaction logs
- Mainly focusing on the problem of user identification

# Related Work

- Simulator mimics the use of Bitcoin as a primary currency to support daily transactions in a university
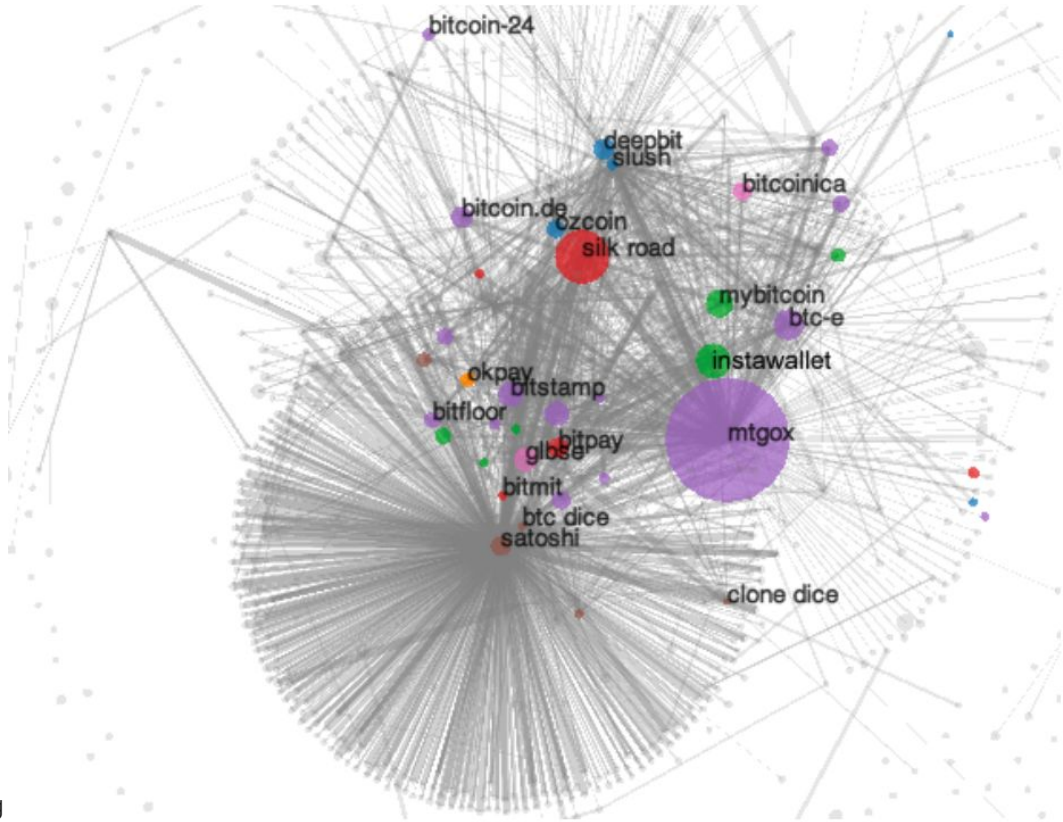- Results showed that almost 40% user profiles could be recovered



Evaluating User Privacy in Bitcoin, Elli Androulaki et al

# Heuristics

I - Common Spending

If two or more addresses are inputs of the same transaction with one output, then all these addresses are controlled by the same user. This heuristic is expected to be very accurate as common spending by different users should be based on high degree of trust between them

Evaluating User Privacy in Bitcoin, Elli Androulaki et al

# Related Work

- Heuristics to perform address clustering on tagged data
- Tagged users by simply transacting with them (e.g., depositing into and withdrawing bitcoins from Mt. Gox) and then observing the addresses they used
- Known (or assumed) addresses from various forums and websites,.

*A Fistful of Bitcoins: Characterizing Payments Among Men with No Names,* Sarah Meiklejohn et al, 2013

# Heuristics

II - One time Change (OTC)

It is based on the standard Bitcoin mechanism where the change from the transaction is returned to a new address. If the transaction satisfies the conditions of a one-time change transaction, then the OTC output and all the inputs of the transaction are controlled by the same user.

A Fistful of Bitcoins: Characterizing Payments Among Men with No Names, Sarah Meiklejohn et al, 2013
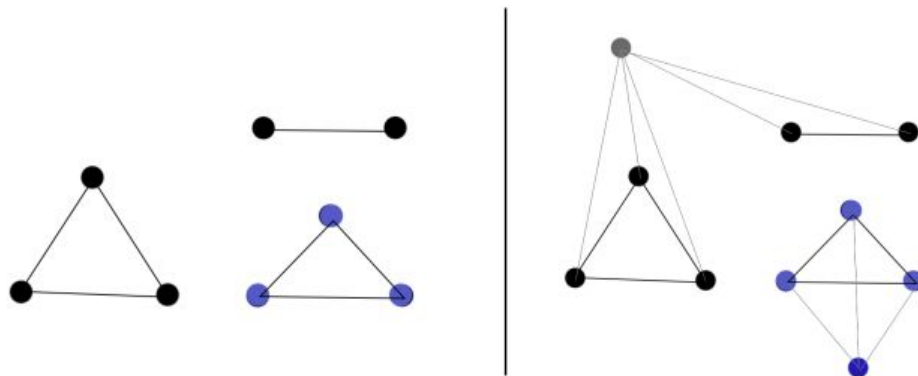
# Related Work

- Off-chain information (tags) show contradictory tags in clusters created by common behavior pattern analysis (heuristics mentioned earlier)
- Tags used as votes against address union using a probabilistic clustering model

| Category | Number of tags | Number of common tags (size) | Examples of common tags |
|---|---|---|---|
| services | 33 | 5 (> 100K) | *Bitpay.com, Xapo.com* |
| gambling | 34 | 6 (> 50K) | *999Dice.com, primedice.com* |
| mixer | 3 | 1 (> 100K) | *BitcoinFog* |
| dnm | 14 | 5 (> 100K) | *SilkRoad Marketplace* |
| exchange | 64 | 12 (> 100K) | *BTC-e.com, Bittrex.com* |
| pool | 15 | 2 (> 50K) | *BTCChina, Hashnest.com* |

TABLE IV: Tags of the biggest cluster in case of clustering without constraints.

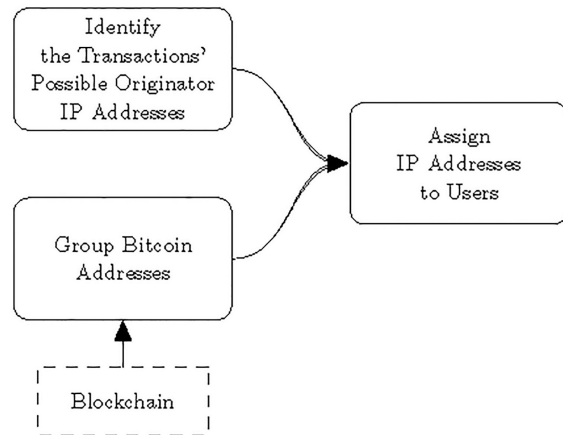Automatic Bitcoin Address Clustering, Dmitry Ermilov et al,

# Related Work

- Network attributes have information like script types, relaying IPs, timestamps, etc
- Constructed a hybrid structure-attribute graph
- Clustering using HDBSCAN
- Used chainalysis data as ground truth. Shows good results in the test cases analyzed in the paper



*Figures 5, 6: In the left image, we have the original graph. Beyond its structure, it has various attributional qualities that should also be taken into account (represented by the colors). On the right, we create points to represent the 'black' and 'blue' attributes and draw the appropriate edges into the graph. Now, the colors no longer matter – we have represented that information into the structure itself!*
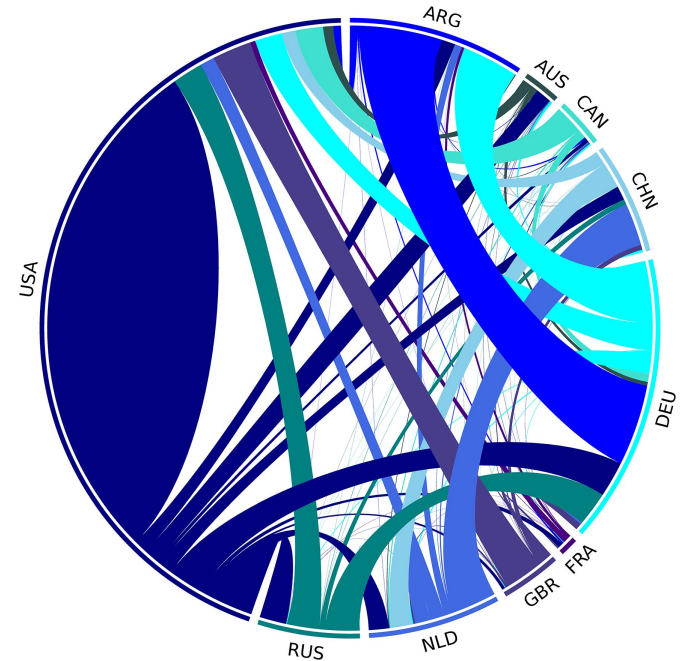
De-Anonymizing the Bitcoin Blockchain, Bharath Srivatsan et al. 2016

# Related Work

- Installed modified Bitcoin clients distributed in the network to observe as many bitcoin transaction messages as possible.
- Assigned probabilities of originator IP addresses to messages observed by the clients based on the relay time
- Combined these with address clustering heuristics to get final probability assignments for user identification



A Bayesian approach to identify Bitcoin users, Péter L. Juhász, József Stéger, 2018

# Related Work

- Traced geographic location from originating IP addresses
- Were able to identify several thousand Bitcoin clients and bind their transactions to geographical locations during a two month observation period.



**The flow of Bitcoin between countries**

# Part-II
**Data Visualization API for link prediction**

# Related Work

- Works that have successfully used network representation learning approaches on Bitcoin data for purposes like community detection and price prediction

# Related Work

- Refined address linking heuristics by also contracting output side addresses

- Following clustering techniques were used on the contracted transaction network
  - K-means on handpicked features
  - K-means on vector embeddings generated by Node2Vec
  - Spectral clustering by using the Fiedler vector method

- Different communities and user categories were identified from all the methods

User Categorization and Community Detection in Bitcoin Network, Dongyuan Mao and Yifei Zhang

# Related Work

- Representation learning approach based on the geometric structure and topological features of the Bitcoin transaction graph
- Betti numbers and sequences generated from the graph in a predictive learning model
- 40% improvement from previous baselines in 7-day ahead price prediction

TABLE I: Features used in Machine Learning models for a given day.

| Approach | Feature Set |
|---|---|
| Basic features | $Price, TotalTx, MeanDegree$ $MeanTxAmount, TotalTxAmount$ $NumNewAddress, ClusCoeff$ |
| Filtration (Sec III-A) | $Price, TotalTx, \mathcal{O}^{\epsilon_1} \ldots \mathcal{O}^{\epsilon_S}$ |
| Betti (Sec. III-B1) | $Price, TotalTx, \beta_0(\epsilon_1), \ldots, \beta_0(\epsilon_S)$ $\beta_1(\epsilon_1), \ldots, \beta_1(\epsilon_S)$ |
| Betti derivative (Sec. III-B2) | $Price, TotalTx$ $\beta_0(\epsilon_1), \ldots, \beta_0(\epsilon_S), \beta_1(\epsilon_1), \ldots, \beta_1(\epsilon_S),$ $\beta_0'(\epsilon_1), \ldots, \beta_0'(\epsilon_S), \beta_1'(\epsilon_1), \ldots, \beta_1'(\epsilon_S)$ |

ChainNet: Learning on Blockchain Graphs with Topological Features, Nazmiye Ceren Abay et al

# Proposed Methodology

Use network representation learning techniques on the bitcoin transaction graph to solve problems like -

- Link Prediction
- User Identification

# Proposed Methodology

Various models can be used according to the graph properties-

- Dynamic - Online Node2Vec
- Sparse - DeepWalk
- Behavioral/Structural similarity  - Struct2Vec

# Problem

- No real ground truth for Bitcoin user data
- All previous results are essentially unreliable
- Hard to explain the results of any study on Bitcoin analysis
- Need for a data visualization API that makes it possible to observe and analyze subgraphs and node activities for explanation purposes

# Revised Objective

Create bitcoin data visualizations that can help get intuitive explanations for the results from AI models.

Display

1. Node statistics - (Num transactions, In degree, Out degree, First Active, Last Active)
2. Most recent transactions (with ability to look at a particular subset of transaction ids)
3. Temporal activity of the node
4. Graphical visualization of the recent transactions
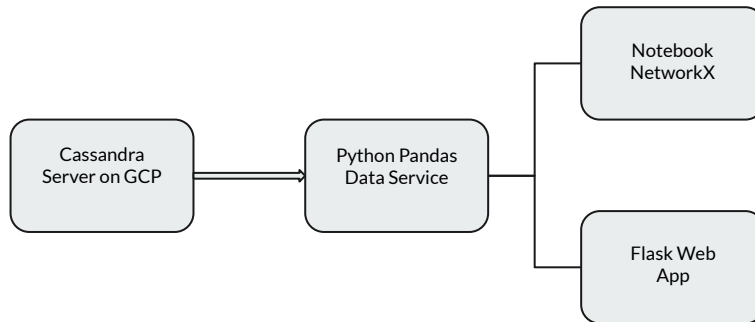5. Vector embedding of the node(s)

# Dataset

- Bitcoin network dataset - D Kondor et al
  https://senseable2015-6.mit.edu/bitcoin/
  - Bh.dat -- Block id, timestamp
  - Th.dat -- transaction id, block id
  - Txin.dat -- transaction id, input addresses
  - Txout.dat -- transaction id, output addresses

- Bitcoin Whos Who - Potential ground truth
  https://bitcoinwhoswho.com/

# Framework

- Google Cloud Platform
- Cassandra 3.x
- Python Pandas Dataframe
- NetworkX
- Flask

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│  Cassandra   │      │ Python Pandas│      │  Notebook    │
│ Server on GCP│─────▶│ Data Service │──────│  NetworkX    │
└──────────────┘      └──────────────┘   │  └──────────────┘
                                         │  ┌──────────────┐
                                         └──│  Flask Web   │
                                            │     App      │
                                            └──────────────┘
```

# Cassandra Tables

Node Info:

- add_id : Integer node id
- first_active : Timestamp of first transaction
- last_active : Timestamp of last transaction
- in_deg : Total number of incoming edges
- out_deg : Total number of outgoing edges
- num_txs : Number of unique transactions

| add_id | add_hash | first_active | in_deg | last_active | num_rows |
|--------|----------|--------------|--------|-------------|----------|
| 265689 | None | 2011-02-27 23:02:46 | 15 | 2011-03-15 22:01:02 | 26 |

| num_txs | out_deg |
|---------|---------|
| 12 | 11 |

# Cassandra Tables

Transactions:

- add_id : Integer node id
- block_id : block number for cassandra partitioning
- txid : Unique transaction id
- role : Role of the add_id in the given transaction
- nbr_id : Node ids on other end of transaction
- timestamp : Time of the transaction

| | add_id | block_id | txid | nbr_id | amount | role | timestamp |
|---|---|---|---|---|---|---|---|
| 25 | 265689 | 0 | 335300 | 305470 | 0.17 | source | 2011-03-15 22:01:02 |
| 24 | 265689 | 0 | 319797 | 205987 | 0.17 | target | 2011-03-09 12:28:19 |
| 23 | 265689 | 0 | 313276 | 283929 | 0.47 | source | 2011-03-06 12:21:03 |
| 22 | 265689 | 0 | 313276 | 282408 | 0.32 | source | 2011-03-06 12:21:03 |
| 21 | 265689 | 0 | 313225 | 283871 | 0.02 | source | 2011-03-06 12:00:53 |
| 20 | 265689 | 0 | 313225 | 282408 | 0.48 | source | 2011-03-06 12:00:53 |
| 19 | 265689 | 0 | 313224 | 282341 | 0.49 | target | 2011-03-06 12:00:53 |
| 18 | 265689 | 0 | 313224 | 263200 | 0.10 | target | 2011-03-06 12:00:53 |
| 17 | 265689 | 0 | 313224 | 262441 | 0.10 | target | 2011-03-06 12:00:53 |
| 16 | 265689 | 0 | 313224 | 205987 | 0.10 | target | 2011-03-06 12:00:53 |
| 15 | 265689 | 0 | 311960 | 269798 | 0.50 | target | 2011-03-05 18:00:10 |
| 14 | 265689 | 0 | 310542 | 280699 | 2.00 | source | 2011-03-05 01:23:18 |
| 13 | 265689 | 0 | 310542 | 276410 | 8.00 | source | 2011-03-05 01:23:18 |

# GUI

## Bitcoin Data Search
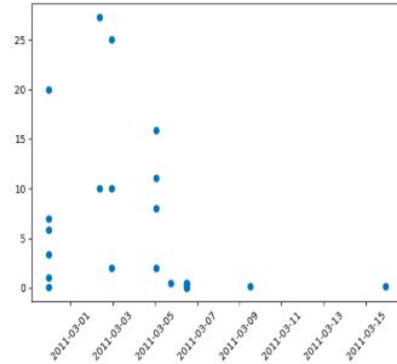
Node ID:

[                    ] [ Submit ]

# GUI

Node ID:
[        ] Submit

## Search results for node 265689!

| Node Id | Num transactions | First Active | Last Active | In Degree | Out Degree |
|---|---|---|---|---|---|
| 265689 | 12 | 2011-02-27 23:02:46 | 2011-03-15 22:01:02 | 15 | 11 |



| Serial Number | Time | Amount (BTC) | Role | Neighbor NodeID | TXID |
|---|---|---|---|---|---|
| 0 | 2011-03-15 22:01:02 | 0.17 | source | 305470 | 335300 |
| 1 | 2011-03-09 12:28:19 | 0.17 | target | 205987 | 319797 |
| 2 | 2011-03-06 12:21:03 | 0.47 | source | 283929 | 313276 |
| 3 | 2011-03-06 12:21:03 | 0.32 | source | 282408 | 313276 |
| 4 | 2011-03-06 12:00:53 | 0.02 | source | 283871 | 313225 |
| 5 | 2011-03-06 12:00:53 | 0.48 | source | 282408 | 313225 |
| 6 | 2011-03-06 12:00:53 | 0.49 | target | 282341 | 313224 |
| 7 | 2011-03-06 12:00:53 | 0.1 | target | 263200 | 313224 |
| 8 | 2011-03-06 12:00:53 | 0.1 | target | 262441 | 313224 |

# GUI
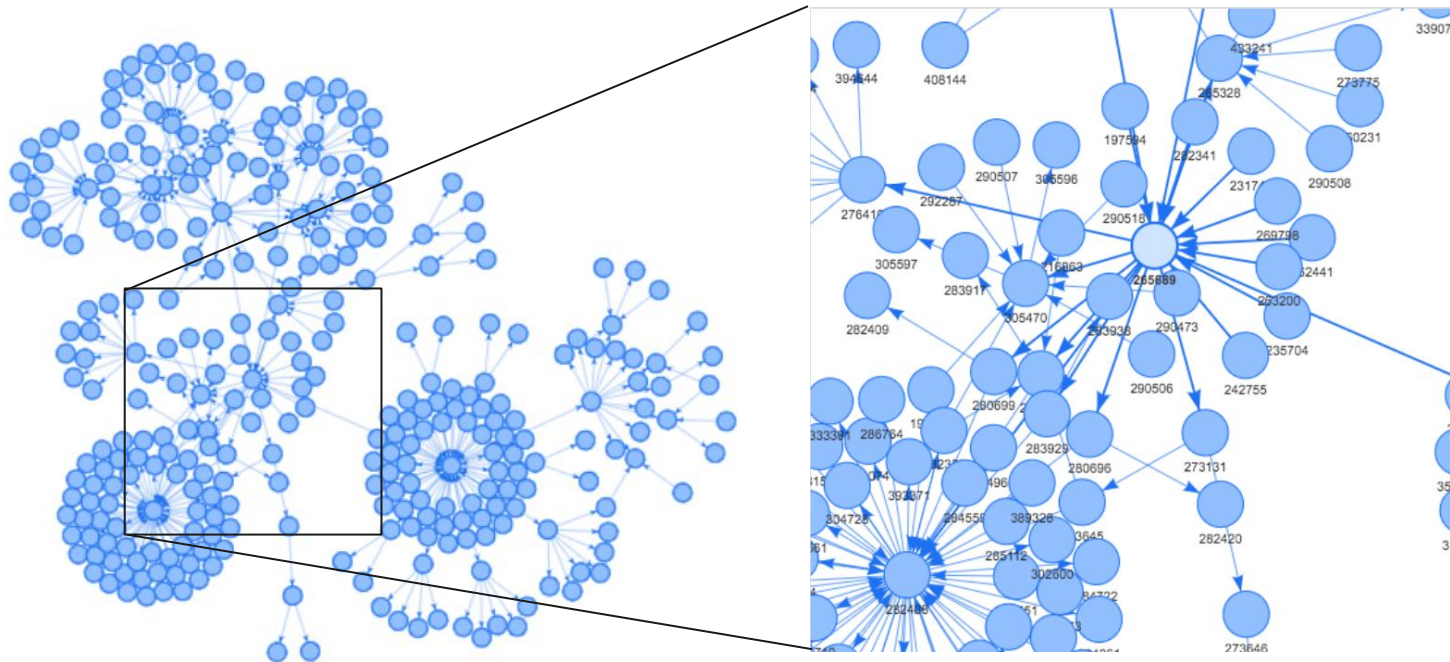


Graphical view of recent transactions

Temporal view of all transactions

# GUI

K-step
information
(k=3)

# GUI

Sample vector embedding of a node

744
[0.000615999975707382, 0.00035099995513477087, −0.0004220000118948519, −0.0001429999974789098, 0.0004030000127386302, 0.0003769999893847853, −0.0006849999772384763, 0.00015700000221841037, −3.9999998989515007e−05, 0.00016399999731220305, 0.00043799998820759356, 6.399999983841553e−05, 0.00022899999748915434, 0.0004579999949783087, 0.0005319999763742089, 0.0001440000509340316, 0.00011100000119768083, 6.900000153109431e−05, −9.200000204145908e−05, 0.000707999977748841, 0.0006430000066757202, 0.00026699995801597783, 0.00012099997730708078, 0.00025099990038502574, 2.4000000848900527e−05, −0.0001500000071246177, 0.00011300000187475234, −0.00011000000085914508, −9.999999974752427e−07, 0.0004879999905824661, −2.2000000171829015e−05, 5.500000042957254e−05, −3.600000127335079e−05, 0.0005789999850094318, −0.00036100001307204366, 0.00028199999360367656, 4.099999932805076e−05, 0.0002739999908953905, 0.00016500000492669642, −8.800000068731606e−05, −0.00036899998667649925, 5.0999999075429514e−05, 0.00026199998683296144, 0.00021699999342672527, 0.0001849999971454963, −0.0001429999974789098, −0.0002119999990100041, 4.70000013592653e−05, 0.00017299999308306724, 0.000297999999024865, 0.00012500000059371814, −3.199999991920777e−05, 0.00032500000088475645, 0.00013699999544769526, 6.000000212225132e−06, −0.0002730000123847276, 0.00011100000119768083, −0.00015100000018719584, 0.00010299999848939478, −3.199999991920777e−05, 0.00014200000441633165, 9.7999996796716e−05, 0.0005849999724887311, 0.00030499999411404133, −0.00010599999950500205, −0.00022499996135011325, −0.00041700000292621553, −0.0001630000042496249, 7.400000322377309e−05, −4.600000102072954e−05, −0.0005370000144466758, −0.00022800000442657625, −0.0002570000069681555, 4.999999873689376e−05, 0.0007660000119358301, −0.0003910000086762011, 0.00021100000594742596, 0.00017800000205170363, 0.00042399999802000082, −4.099999932805076e−05, 0.00038899993447214375, −0.0005530000198632479, −0.00022800000442657625, 0.00025499991739516876, 0.0001069999984353781, 0.0004360000020824373, 0.00025499991739516876, −0.00012700000661425292, −0.00011999999696854502, 0.0003110000106971711, 0.00011500000255182385, −0.0004579999949783087, −4.400000034365803e−05, −0.00011100000119768083, 6.70000008540228e−05, −0.0002610000083222985, −3.099999958067201e−05, 0.00010499999916664663, 9.7999996796716e−05, 0.00035499999648891139]

# Future Work

- Utilize data visualization API to explain results of the Link Prediction model
- Expand the database with a dynamic data ingestion pipeline

# Conclusion

- Hard to solve the problem of bitcoin user identification as there is no real ground truth. Most previous results are unreliable.

- Good data visualization plays a significant role in making sense of AI models.

# References

- Nakamoto S. Bitcoin: A peer-to-peer electronic cash system," http://bitcoin.org/bitcoin.pdf; 2009
- Androulaki E, Karame GO, Roeschlin M, Scherer T, Capkun S. Evaluating User Privacy in Bitcoin. In: Sadeghi AR, editor. Financial Cryptography and Data Security. Berlin, Heidelberg: Springer Berlin Heidelberg; 2013. p. 34–51
- Meiklejohn S, Pomarole M, Jordan G, Levchenko K, McCoy D, Voelker GM, et al. A Fistful of Bitcoins: Characterizing Payments Among Men with No Names. Commun ACM. 2016;59(4):86–93
- D. Ermilov, M. Panov, Y. Yanovich, "Automatic Bitcoin address clustering", Proc. Int. Conf. Mach. Learn. Appl. (ICMLA), pp. 461-466, Dec. 2017
- Bharath Srivatsan. 2016. De-Anonymizing the Bitcoin Blockchain. (2016)
- Juhász PL, Stéger J, Kondor D, Vattay G (2018) A Bayesian approach to identify Bitcoin users. PLoS ONE 13(12): e0207000
- D. Mao, Y. Zhang, User Categorization and Community Detection in Bitcoin Network, snap.stanford.edu
- Nazmiye C. Abay et al. "ChainNet: Learning on Blockchain Graphs with Topological Features". In: (2019)