

Snap Shot Samplings of the Bitcoin Transaction Network

Lambert Leong

University of Hawaii

www.lambertleong.com

1. Introduction

The invention of the internet has helped to provide a new platform for the exchange of money by providing the framework to connect individuals from all around the world. Most monetary exchanges over the internet occur with the help of third party entities like a bank or a credit card company. These third party entities usually have oversight over all transactions and are responsible for verifying the integrity of transactions. This exchange protocol, which requires a verifier, has become the conventional method of exchanging money over the internet. While conventional exchange methods account for a good amount of the daily transactions around the world, there are a few drawbacks that present problems to many users. Third party verifiers constantly receive a high volume of transactions which need to be verified before money is permanently moved from one individual to another. Different companies and banks have different methods of prioritizing the order in which transactions are verified. In any case, a high number of transactions being verified by a few third party entities leads to a bottleneck effect which can leave recipients waiting a while for payments to clear and be received. In an effort to reduce the volume of transactions many third party verifiers require minimum transactions amounts and, in some cases, a small transaction fee. While this may help to dampen the volume of transactions that need to be verified, it limits the financial freedom of those who participate in conventional exchange methods.

Bitcoin is a cryptocurrency that offered a solution to the shortcomings present in conventional exchange methods. Invented by Satoshi Nakamoto, bitcoin is a peer-to-peer payment system in which payments are validated by math rather than trust [5]. In other words, instead of trusting a potentially nefarious third party entity, transactions are verified by a proof of work algorithm known as block hashing [5 - 7]. A public ledger, known as the blockchain, is created as a result of the hashing algorithm and it contains every transaction in the history of the network. New blocks are hashed about every 10 minutes and when this happens, transactions are verified and added to the ledger indicating a successful transfer of bitcoin from one individual or entity to another [5]. Anonymity is preserved by keeping user addresses separate from personal information. In addition, there are some instances where a particular user's wallet can generate new sending and receiving addresses to further maintain a level of anonymity [11, 12].

Since its genesis, bitcoin has gained popularity and traction. As a result, the number of bitcoin addresses and the number of transactions have continued to increase at an exponential rate as seen in figures 1 and 2. In this paper, we analyze different snap shots of the growing bitcoin network. Since the bitcoin network is currently large and continuously growing, we are interested in investigating appropriate network sampling timeframes. Determining a sampling criterion would allow analysis to be done on a smaller sub graph from which, conclusions can be drawn and extrapolated to the network as a whole. Running a network analysis on a representative sub graph would be more efficient and computationally favorable than the whole bitcoin network.

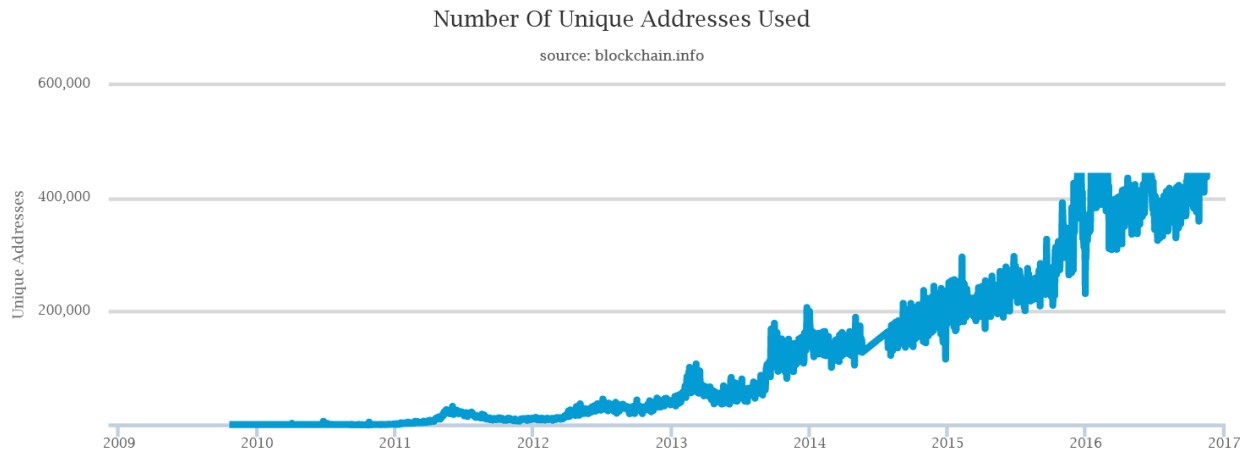


Figure 1. Graph showing the growth in the number of Bitcoin addresses over the lifespan of bitcoin [10]



Figure 2. Graph showing the growth in the number of Bitcoin transactions over the lifespan of bitcoin [10]

2. Method

2.1 Data Collection

Bitcoins are constantly being exchanged and moving throughout the network and as a result the numbers of transactions are ever growing. The number of users and addresses are also constantly increasing. Every transaction involving bitcoin is public and every transaction since bitcoins inception in 2009 is recorded on a public ledger. This is a unique feature of bitcoin and every node participating in the network has a copy of this public ledger. This ledger is constantly being updated as new transactions are confirmed. Updates occur on every ledger on every node participating in the network and thus every node is in consensus with other nodes on the current status of the ledger.

While we could have joined the network and downloaded a copy of the public ledger to build our network graph, we felt that it would not be an optimal method to evaluate the current status of the bitcoin network. As a result we decided to utilize a web socket API from BlockChain.Info. This API allowed us to stream live transactions as they were on their way to be confirmed. Each transaction came in as a JavaScript Object Notation (JSON) file and a simple script was written to extract and parse out the sender address, receiver address, and the bitcoin amount and store it in a comma separated file (.csv). The web socket was run for a total of six hours and separate csv files were taken at the one hour, two hour, and six hour time points. One of our objectives is to investigate if sample size and sampling time affect network characteristics. In other words, we are interested to see if any network metrics change as the network grows. It is important to note that the web socket continued to run and data collection was not interrupted when the one hour and two hour samples were taken. The resulting one hour and two hour samples are sub sets of the six hour sample with the one hour sample also being a subset of the two hour sample.

2.2 Graph Generation

The csv files for the one, two, and six hour samplings were imported into Gephi to generate a graph of the sampled bitcoin network. The csv files were imported as edge list where the sending address corresponded to the “source” and the receiving address corresponded to the “target”. The amount was converted from Satoshis to bitcoins, $1 \text{ bitcoin} = 1 \times 10^8 \text{ Satoshi}$, and stored as a float, an edge attribute. Graph files were exported as .graphml for further data processing.

2.3 Network Analysis

Graph metrics were analyzed with R utilizing the igraph, powerLaw, and linkcomm packages [3]. Transitivity, average degree, average distance, reciprocity, degree distribution, and maximal cliques were measured for each graph and compared to each other. The largest connected component was also extracted from each of the three sample graphs and metrics were computed with igraph. The six hour connected graph was further evaluated for the presence of communities. Using the linkcomm package and the getlinkcommunities by Alex T. Kalinka we were able to generate community’s nodes that linked to every node of a particular community [3]. Investigating community structures in the giant connected component may reveal bitcoin user financial behavior; in particular who exchanges money with whom.

3. Data & Analysis

There was a positive correlation between the sampling time and the amount of observed transactions. As mentioned above, the numbers of transactions are increasing at an exponential rate [10]. We plotted the number of edges, obtained in our sampling, over the duration of six hours. Transactions are indicated as edge in the resulting graph of the networks sampled for one, two, and six hours. Results are shown in figure 3 below.

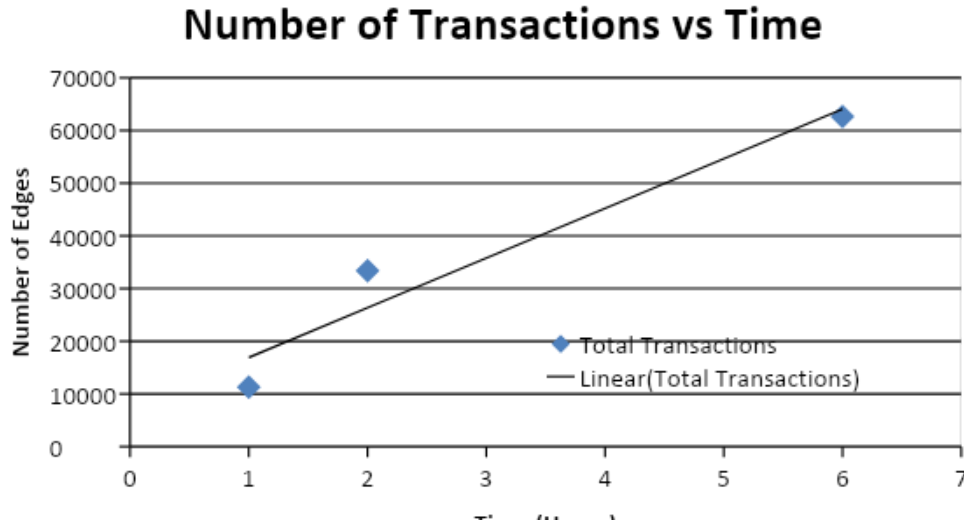


Figure 3. Graph representing the increase in the number of transactions obtained in our sampling.

The trend line would indicate that the number of transactions are growing at a logarithmic pace rather than an exponential one. This is contradictory to what others have reported however, this discrepancy can be explained by the amount of data and sample size [10]. Blockchain.info reports an exponential growth in the number of transactions over the whole life span of bitcoin. Our sampling window could be the limiting factor and perhaps we would start to see an exponential growth trend for a longer sampling duration.

The number of transactions since the inception of bitcoin has an overall general exponential growth trend but the transaction rate seems to fluctuate if observed from the perspective of weeks and months. There is a level of volatility when it comes to the prices or exchange rate of bitcoin [8]. Price volatility could correlate to these fluctuations and periods of lower transaction volumes. A six hour sampling window does not appear to be a long enough time frame to model the overall network growth. Since transaction rates can fluctuate for period of time ranging from days to months, it is best to look at the entire transaction history in order to draw conclusions about the overall growth rate of transactions. A smaller sampling window, like the 6 hours for which we took our samples, may be useful when analyzing how the network is growing in the short term.

Network analysis was performed on all three graphs with the purpose of investigating changes as the time and the size of the sample network increased as well as for comparison to previously reported findings. Analyzing the difference in resulting graph metrics would help in determining an appropriate sampling window that would yield a sample graph that is representative of the network as a whole. Table 1 contains the metrics calculated for all three graphs.

	1 Hour	2 Hour	6 Hour
Hours	1	2	6
Nodes	18654	52312	95209
Edges	11262	33408	62635
Reciprocity	1.82E-03	1.66E-03	1.99E-03
Transitivity (global)	2.08E-04	1.79E-04	2.19E-04
Mean Degree	1.207462	1.27726	1.315737
Maximal Cliques	3	3	3
Dyads	10943	32512	61122
Triads	2	8	34

Table 1. Graph metrics calculated using igraph [3]

Across all three sample graphs, reciprocity, global transitivity, and mean degree were fairly consistent. These metrics changed only slightly as the graph grew over the 6 hours. Transitivity was low for all three graphs, which was to be expected. Transitivity would increase when, for example, someone sent bitcoin to two separate people and those two people exchanged bitcoin with each other. All three graphs had a low mean degree with the largest mean degree value being 1.32 for the six hour sample graph. A low mean degree would help to explain a low transitivity as a low mean degree indicates that the majority of the nodes in a graph have one edge. Nodes need a minimum of two edges for transitivity. The mean degree values for each graph indicate that most addresses execute only one transaction. The number of dyads reflects that the graph is primarily composed of pairs of nodes that share only one edge. For each sample graph, dyads account for the majority of nodes in the network. Triads make up a significantly lesser portion of the graph. We did not expect to see many triads or cliques greater than that during our sampling. In the short term, it is unlikely that triads form because it would mean that some money comes back to the original sender implying that some portion of that transaction should not have been sent out in the first place. However, goods or services could have been exchanged amongst the triad which would explain its formation. The network is built only on bitcoin transactions and thus we do not know the exact motive behind each transaction. A longer sampling would probably have led to more triads and possibly greater maximal cliques. However, it is just as likely that the number of pairs of nodes would have increased at a similar rate.

Reciprocity was also low. This was expected for this type of network. It is unlikely that a node would reciprocate an edge because that would be like person A sending 1 bitcoin to person B in the first transaction then person B sending back, to person A, a given amount of bitcoin in a second transaction. The first and second transaction can be combined into one and the net amount, the difference in the amount between the first and second transaction, should be sent to the appropriate person which would eliminate an edge in one direction. The reciprocity value is small for all graphs but it is not zero which means that there are instance of reciprocating edges. The instances where a user sends the wrong amount to another user who, is kind enough to correct the transaction and send it back to the person it originated from may be an instance for reciprocation. Returning goods and getting a refund on something bought with bitcoin is another instance that may warrant reciprocation. Reciprocation can also be the result of gambling sites where, gamblers gamble in bitcoin and earn their winnings in bitcoin too. These instances are rare and small in comparison to the vast number of transactions constantly being received. The rarity of reciprocity is reflected in the value obtained for all sample graphs.

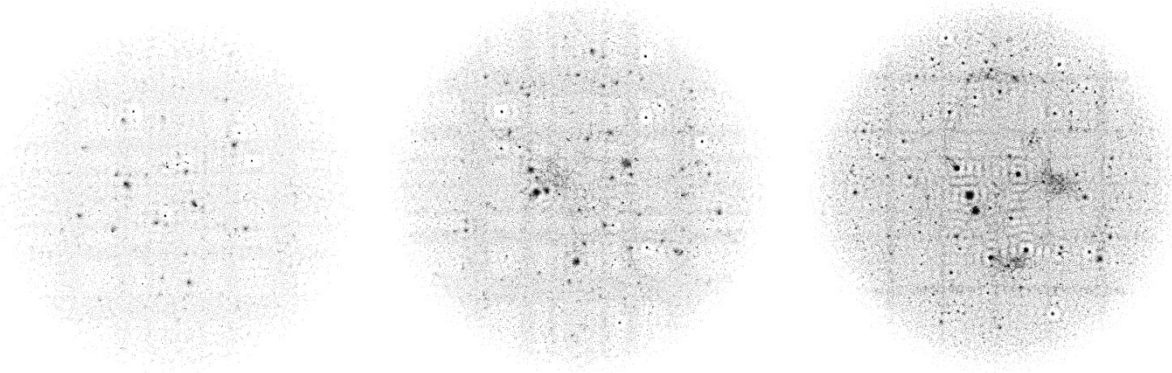


Figure 4. Visualizations of one, two, and six hour sample graphs, from left to right. Graph density increases as collection time increases

Visual representations are show in figure 4. The density of the graph increase as sampling time increases. Dark black spots are star like structures consisting of clusters of nodes around a single node with a high degree. As time goes on, more nodes form edges to these central nodes which increase the size of the black spots. However, it does not appear that links are formed from one high degree node to another. Nodes with high degree seem to increase their degree as more nodes join the network which exhibits characteristics of preferential attachment and implies a possible scale free network [1, 4].

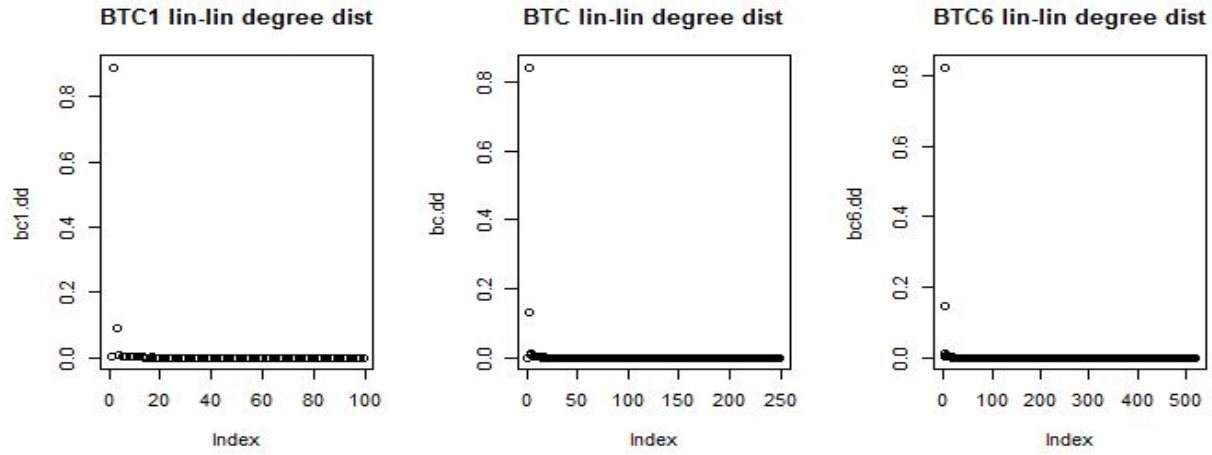


Figure 5. Degree distribution of one, two, and six hour networks, from left to right. Constructed with igraph [3]

Scale free networks are characterized by having a power-law degree distribution. We investigated the degree distribution in figure 5. In figure 5 we see that the degree distribution for all sample networks is heavily tailed to the right which is characteristic of power-law distributions.

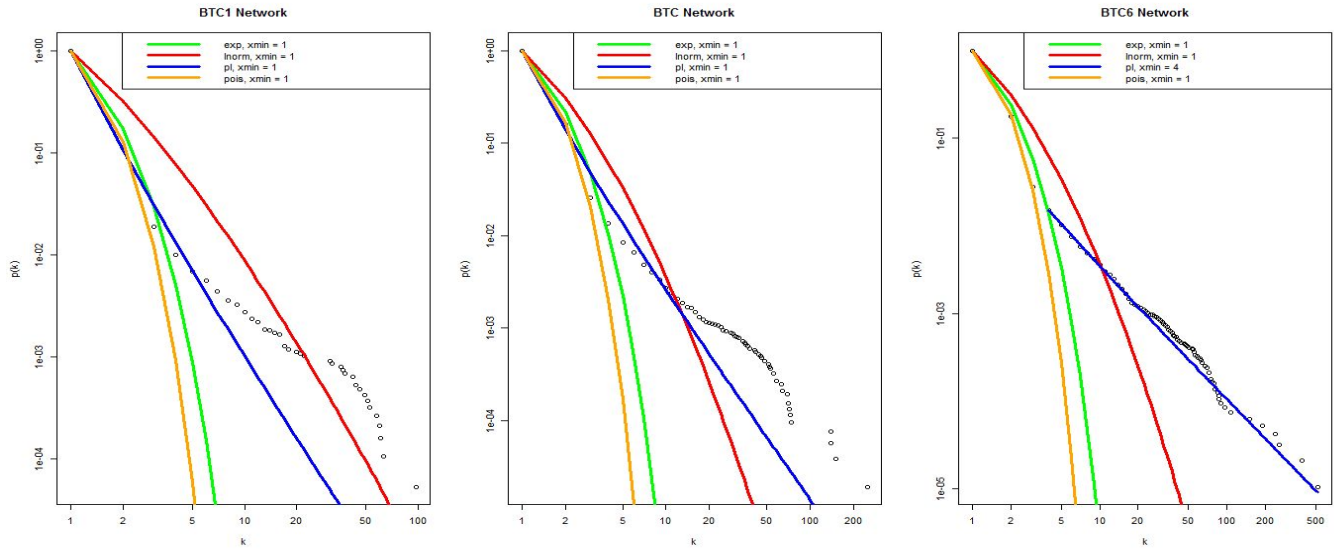


Figure 6. From left to right, plot of one, two, and six hour network degree distributions with exponential(Red), log normal(Green), power-law(Blue), and Poisson(Orange) best fit lines. Constructed with igraph [3]

The degree distributions for each graph was plotted in conjunction with exponential, log normal, power law, and Poisson best fit lines. Results in figure 6 would indicate that the one hour and two hour graphs do not fit any particular trend. However, when Xmin was set at 4, the six hour distribution

seemed to fit the power law distribution. Kolmogorov-Smirnov test results in a p value of 0.46 which indicates that we cannot reject the hypothesis that the six hour distribution fits a power law [3].

Evidence of a power law distribution was seen in our largest graph that was sampled for the longest period of time. It would appear that smaller samplings of the network do not resemble a scale-free network however, the largest sampling started to show evidence of a power law distribution which is indicative of a scale-free network. It is likely that a longer sampling window would yield a larger graph in which a power law is more apparent.

Although the graphs are mainly composed of pairs of nodes, we were interested in the presence of communities. We extracted the largest connected component from the sample graphs. This left us with graphs that excluded the many isolated pairs of nodes. Metrics were calculated and are shown in table 2.

	1 Hour Connected Component	2 Hour Connected Component	6 Hour Connected Component
Nodes	238	3028	9052
Edges	243	3193	9698
Reciprocity	0.00E+00	6.29E-04	1.04E-03
Transitivity (global)	0.00E+00	1.06E-04	2.20E-04
Mean Degree	2.042017	2.108983	2.142731
Maximal Cliques	2	3	3
Dyads	2	3172	9582
Triads	0	3	20

Table 2. Graph metrics for giant component subgraph calculated using igrap [3]h

Directedness was not taken into account when extracting the giant components from each graph. Each sub graph is significantly smaller than its original graph implying again that the graphs are primarily made up of isolated pairs of nodes. This means that, during the sampling window, addresses only sent one transaction. For the giant component subgraphs, the metrics were not consistent across all giant component subgraphs, as seen in the original graphs.

Our analysis showed that the six hour graph showed to be the best representative sample network. Therefore, we chose to run community detection on the six hour giant component subgraph. The resulting graphs can be viewed in the appendix. Figures 7 & 8 are sub graphs of the graph in the appendix.



Figure 7 & 8. Two community sub graphs extracted out of the six hour giant component sub graph in appendix. Nodes are sized based on degree. Pink nodes correspond to bitcoin addresses while green nodes indicate the community each node belongs to. Figure 7(left) is a community in which one address node has a high degree. Figure 8(right) represents a community in which all nodes have the same relative degree.

Both figures are two of the biggest communities contained in the six hour giant component graphs. While they are both two of the biggest communities, they have slightly different structure. Figure 7 illustrates a community in which many nodes share edges with one central node. This creates a star like structure seen in most of the other communities. The second community in figure 8 does not contain a single node with a high degree. Most of the nodes have the same degree and yet they form a community. This would imply that the nodes are more interconnected to each other, as there is no distinct hub or authority. Figure 8 represents more of a decentralized community.

Community structures provide clues as to what types of individual or entity may be contained within a network. For communities that resemble that of figure 7, it is possible that the high degree node may be a type of vendor or a spender depending whether or not there is a high in-degree or out-degree. A high in-degree would mean that one address is receiving a lot of transactions while a high out-degree would mean that an address is paying a lot of other individuals. Figure 8 resembles a different type of community. It is difficult to draw conclusions about the nature of the nodes in this community. Gambling is one possible instance in which many transactions occur amongst a community of individuals. However, it is also likely this community is a bunch of nodes doing business together. Given that it was a six hour window, it is unlikely that business would execute that many transactions with each other. Further analysis is needed to determine the nature of the individuals but community detection is a good starting point.

4. Conclusion

The six hour sample graph seemed to provide a sample graph that was most representative of the reported nature of the bitcoin transaction network. Modeling the network growth with the one, two and six hour subgraphs did not seem to be a reliable means of evaluating long term growth of the network. Growth fluctuations are experienced in the short term and thus any short term network growth sampling may not be indicative of the network's overall growth.

Sample graphs yielded metrics that we expected to see. Our six hour sample graph showed evidence of a possible power law degree distribution which suggest that we cannot dismiss the possibility that our network is scale-free. The emergence of such correlation was only seen in the six hour graph and not in the two hour or one hour graphs. Further investigations should consider a longer sampling window in order to get a better snap shot of the network. A longer sample window may be beneficial and provide a more definitive distribution.

A deeper investigation into the owners of each address would aid in a more comprehensive analysis of the bitcoin transaction network. Certain security feature most likely lead to some level of skewness in the data set. The feature which allows users to change the sending and receiving address is a powerful security feature but it could lead to a misrepresentation of the amount of actual bitcoin users participating in the network. Since all transactions are kept on a public ledger, addresses associated with a transaction involving a lot of bitcoin may become targets of malicious attacks. Some bitcoin wallets allow users to generate new public addresses in order to maintain a level of anonymity on the network. This feature may explain the large number of pairs of nodes in our network samples. It could be that individuals change their address for every transaction and a new address is used for subsequent transactions. Situations are present where an address cannot be changed or it would be inconvenient and this is why we see nodes with high degrees despite the fact that an address can be changed. For example, if a vendor accepts bitcoins as payment it would be inconvenient for the vendor and patrons if the vendor constantly changed addresses because patrons, especially "regulars", would have to constantly check that the address is current and that they are sending money to the right address.

Individuals are allowed to have multiple bitcoin wallets just as someone may have multiple bank accounts or credit cards. Multiple wallets per individual would lead to a similar phenomenon to generating new addresses however; one more thing needs to be considered. Money can be sent from one wallet to another wallet but one individual may have ownership to both wallets; it is like transferring funds between two bank accounts you own. To acknowledge this transfer the, the network needs to verify it and thus it would show up as a transaction in the network.

Community detection may help to investigate ownership of a particular address because multiple addresses in a particular community may belong to one individual. Combining these addresses would lead to a better resolution of the network. The way it stands now, our network primarily resembles the movement of bitcoin. Building a transaction network in which the nodes represent a single individual or entity rather than a single address would result in a network that would resemble how people spend

bitcoin. Separating the identity from the bitcoin address was in the original design of the network to protect individual

5. References

- [1]
B. Albert-László, *Barabási Albert-László - Books*. .
- [2]
M. Fleder, M. S. Kester, and S. Pillai, "Bitcoin Transaction Graph Analysis," *arXiv:1502.01657 [cs]*, Feb. 2015.
- [3]
C. Gabor and N. Tamas, "The igraph software package for complex network research," *InterJournal*, vol. Complex Systems, p. 1695, 2006.
- [4]
D. Kondor, M. Pósfai, I. Csabai, and G. Vattay, "Do the Rich Get Richer? An Empirical Analysis of the Bitcoin Transaction Network," *PLOS ONE*, vol. 9, no. 2, p. e86197, Feb. 2014.
- [5]
S. Nakamoto, "Bitcoin: A peer-to-peer electronic cash system,," 2008. [Online]. Available: <https://bitcoin.org/bitcoin.pdf>. [Accessed: 12-Dec-2016].
- [6]
F. Reid and M. Harrigan, "An Analysis of Anonymity in the Bitcoin System," in *Security and Privacy in Social Networks*, Y. Althuler, Y. Elovici, A. B. Cremers, N. Aharony, and A. Pentland, Eds. Springer New York, 2013, pp. 197–223.
- [7]
D. Ron and A. Shamir, "Quantitative Analysis of the Full Bitcoin Transaction Graph," in *Financial Cryptography and Data Security*, 2013, pp. 6–24.
- [8]
D. Yermack, "Is Bitcoin a Real Currency? An economic appraisal," National Bureau of Economic Research, Working Paper 19747, Dec. 2013.
- [9]
A. Yu and B. Bunz, "Community_Detection_and_Analysis_in_the_Bitcoin_Network.pdf," 09-Dec-2015. [Online]. Available: http://snap.stanford.edu/class/cs224w-2015/projects_2015/Community_Detection_and_Analysis_in_the_Bitcoin_Network.pdf. [Accessed: 12-Dec-2016].
- [10]
"Bitcoin Charts & Graphs - Blockchain," *Blockchain.info*. [Online]. Available: <https://blockchain.info/charts>. [Accessed: 12-Dec-2016].
- [11]
"How It Works - My Wallet - blockchain.info." [Online]. Available: <https://blockchain.info/wallet/how-it-works>. [Accessed: 12-Dec-2016].
- [12]
"Why did my wallet address change?," *Coinbase*. [Online]. Available: <https://support.coinbase.com/customer/en/portal/articles/2276500>. [Accessed: 12-Dec-2016].