

Working Paper presented at the

# Peer-to-Peer Financial Systems 2016 Workshop

July, 2016

The Evolution of the Bitcoin  
Economy: Extracting and  
Analyzing the Network of Payment  
Relationships

**Prof. Paolo Tasca**

UCL

**Shaowen Liu**

Deutsche Bundesbank

**Adam S. Hayes**

University of Wisconsin-  
Madison

Powered by



**P2P Financial Systems**



# The Evolution of the Bitcoin Economy: Extracting and Analyzing the Network of Payment Relationships

Paolo Tasca<sup>\*1</sup>, Shaowen Liu<sup>2</sup> and Adam S. Hayes<sup>3</sup>

<sup>1</sup>University College London, Centre for Blockchain Technologies

<sup>2</sup>Deutsche Bundesbank

<sup>3</sup>University of Wisconsin-Madison

July, 2016

## Abstract

In this paper, we gather together the minimum units of Bitcoin identity (the individual addresses), and group them into approximations of business entities, what we call “super clusters”. While these clusters can remain largely anonymous, we are able to ascribe many of them to particular business categories by analyzing some of their specific transaction patterns, as observed during the period from 2009-2015. We are then able to extract and create a map of the network of payment relationships among them, and analyze transaction behavior found in each business category. We conclude by identifying three marked regimes that have evolved as the Bitcoin economy has grown and matured: from an early prototype stage; to a second growth stage populated in large part with “sin” enterprise (*i.e.*, gambling, black markets); to a third stage marked by a sharp progression away from “sin” and toward legitimate enterprises.

**Keywords:** Blockchain, Bitcoin, Digital Currencies, Business Analysis, Network Theory

**JEL Classification:** E42, L14, O12, O35, P40

---

<sup>\*</sup>Correspondence to: Paolo Tasca, UCL Centre for Blockchain Technologies, P.Tasca@ucl.ac.uk

# 1 Introduction

The data reveal that the number of transactions on the Bitcoin blockchain rose exponentially from around 1,000 per day in 2011 to around 300,000 per day at the moment of writing.<sup>1</sup> At the current exchange rate the notional value of daily transaction volume ranges between \$200 and \$300 million.<sup>2</sup> Thus, it becomes appropriate to explore how the Bitcoin economy is populated and extract the map of payment relationships, as well as to trace the evolution of those relationships over time. This paper takes that direction by identifying the interconnection between the economic agents that use the Bitcoin payment network to transfer the digital currency among each other. In particular, we start by clustering the minimum units of Bitcoin identity, which are the individual “addresses”, into what we call “super clusters”, and then we tag those clusters using de-anonymized addresses. A super cluster can be thought of as an approximation of a business entity in that it describes a group of addresses that are owned or controlled collectively for some particular economic purpose by the same entity.<sup>3</sup> Although the exact identities of super clusters can remain unknown, we are able to allocate many of them to specific business categories (namely to exchange, mining pool, online gambling, black market, or composite) by analyzing some of their specific transaction patterns, as observed during the period from 2009-2015.<sup>4</sup> With that information, we unveil and study the Bitcoin network of payment relationships both among super clusters and between them with other clusters in broader business categories (traders, gamblers, black market user-dealers, others, individuals or unknown). In doing so, we are subsequently able to identify three distinct regimes that have existed in the Bitcoin economy’s growth and development. Namely, a “proof of concept” or “mining-dominated” phase, a “sin” or “gambling/black market-dominated” phase, and a “maturation” or “exchange-dominated” phase.

---

<sup>1</sup>By convention we use Bitcoin with a capital ‘B’ to denote the protocol, network, and community, while bitcoin with a small ‘b’ denotes the digital currency and units of that currency.

<sup>2</sup>See Blockchain (2015a).

<sup>3</sup>In principle, a single entity may have control over more than one distinct super cluster if the common ownership of some of their addresses is not evident from the data.

<sup>4</sup>The category “composite” represents clusters that are identified as having a mix of more than one defined business activity.

It is possible to map some of the activity and interaction among Bitcoin users because *pseudonymity*, rather than strict anonymity, is a remarkable characteristic of the Bitcoin network, see *e.g.*, Reid and Harrigan (2013). As such, the identities of users are hidden behind their addresses that work as a pseudonyms, but which may be revealed upon transacting with somebody else.<sup>5</sup> In other words, if Alice remits payment to Bob, then their identities will be revealed to one another by virtue of exchanging addresses to send or receive bitcoin.

The novelty of our study is to elaborate and propose a general de-anonymization method that allows us to link some large clusters, including groups of addresses, to particular business categories. Furthermore, we map the network of payment relationships of these clusters and subsequently describe three distinct regimes within the Bitcoin economy.

There are two general procedures that have been proposed thus far to solve the problem of de-anonymizing Bitcoin addresses: “clustering” and “labeling”. Clustering consists of grouping together in one unique cluster all the addresses that belong to the same beneficial owner (*i.e.*, legal or individual person). This approach requires one to apply either the “input address heuristic” and/or the “change address heuristic.”<sup>6</sup> After clustering, one can apply labeling which consists of either: 1) manually tagging Bitcoin addresses to specific entities by directly participating in Bitcoin transactions with those entities;<sup>7</sup> or

---

<sup>5</sup>A Bitcoin address is an identifier of 26-35 alphanumeric characters that is derived from the public key through the use of one-way cryptographic hashing. The algorithms used to make a Bitcoin address from a public key are the Secure Hash Algorithm (SHA) and the RACE Integrity Primitives Evaluation Message Digest (RIPEMD), specifically SHA256 and RIPEMD160, see *e.g.*, Antonopoulos (2014).

<sup>6</sup>In the Bitcoin network, the output of a transaction is used as the input of another transaction. If the input is larger than the new transaction output the client generates a new Bitcoin address, and sends the difference back to this address. This is known as change. From the Bitcoin wiki: Take the case of the transaction 0a1c0b1ec0ac55a45b1555202daf2e08419648096f5bcc4267898d420dffef87, a 10.89 BTC previously unspent output was spent by the client. 10 BTC was the payment amount, and 0.89 BTC was the amount of change returned. The client can’t spend just 10.00 BTC out of a 10.89 BTC payment anymore than a person can spend \$1 out of a \$20 bill. The entire 10.89 BTC unspent output became the input of this new transaction and in the process produced are two new unspent outputs which have a combined value of 10.89 BTC. The 10.89 BTC is now “spent” and effectively destroyed because the network will prevent it from ever being spent again. Those unspent outputs can now become inputs for future transactions

<sup>7</sup>This is a very time consuming and inefficient activity, *e.g.*, Meiklejohn et al. (2013) by participating in 344 transactions was able to manually tag 1,017.



2) extracting information from specific web pages in which, for any reason, the identity of Bitcoin addresses' holders is public and can be extracted.

According to the *input address heuristic*, addresses used as inputs either synchronously in the same multi-input transactions or asynchronously in different multi-input transactions (when at least one input address is in common), are grouped together in clusters. In other words, if address  $x$  and address  $y$  are both inputs of a transaction, then we assume address  $x$  and  $y$  belong to the same cluster. Furthermore, if both address  $y$  and address  $z$  belong to another transaction, we would extrapolate that address  $x$ ,  $y$  and  $z$  are all belonging to the same cluster. Already Nakamoto (2008) indirectly recognize the power of the input address method by saying that, “[...] *some linking is still unavoidable with multi-input transactions, which necessarily reveal that their inputs were owned by the same owner. The risk is that if the owner of a key is revealed, linking could reveal other transactions that belonged to the same owner*”. Later, Ron and Shamir (2013) extensively discuss the input address method and applied it via the “Union-Find graph algorithm” in a study of the Bitcoin network through the 13th of May, 2012. Within the group of this heuristic, other works including Spagnuolo (2013) and Doll et al. (2014) try to provide some practical applications of the theory by developing front-end web services to show in real-time the cluster *id* correspondent to a specific address query. In particular, Spagnuolo (2013) proposes the BitIodine tool which parses the blockchain, clusters addresses that are likely to belong to a same user or group of users, classifies such users and labels them, and finally visualizes complex information extracted from the Bitcoin network.

With regard to the *change address heuristic*, a cluster is composed of the input addresses plus the output addresses that are predicted to be change addresses. A first proposal of this heuristic comes from Androulaki et al. (2013) which naively assumes that “[...] *in the current Bitcoin implementation, users rarely issue transactions to two different users*”. Presumably this assumption did once hold in the past but it is not holding true any longer. Therefore, this first version of the change address heuristic is relatively fragile compared to the input address heuristic, and aggressive implementa-

tions require large amounts of hand tuning to prevent false positives.<sup>8</sup> Meiklejohn et al. (2013) reevaluates the change address heuristic and applies it cautiously by identifying only *one-time* change addresses under the following conditions: 1) the transaction is not a coin generation; 2) it is the first appearance of the address and this is not the case for other output addresses (*i.e.*, all the other addresses has been previously used); and 3) there is no other address in the output that is the same as the input address (*i.e.*, no self-change address). The assumption behind this version of the change address heuristic is that the change address is one newly generated by the user’s wallet; even the owner may not acknowledge its existence. In contrast, the receiver’s address is usually determined in advance and notified to the sender. Thus, also the *one-time* change address requires significant human adjustment to avoid excessive false positives when: 1) the receiver is a new user or creates a new address never used before; 2) the transaction output has two receivers’ addresses without change address; 3) the sender uses an old address to receive change, or there is no change transaction at all. Despite some studies that rely on this version of the change address heuristic, *e.g.*, Garcia et al. (2014), for the purpose of our study we opt for the input address heuristic. Although this method is subject to some false negatives because it considers eligible for clustering only the addresses being used as transaction inputs, it is nonetheless robust because it is not subject to false positives. In fact, any false positive would compromise the result of the pattern analysis applied to ascribe the clusters to particular business categories; in such a case, the clusters themselves would be composed of wrong addresses that would likely follow incorrect behavioral patterns.

However, false positives could become a problem with the introduction, since 2013, of the *coinjoin* practice, see *e.g.*, Kristov Atlas (2015) and Tasca (2015). Coinjoin transactions are one example of a tool used to further anonymize transactions within a distributed ledger. The principle behind this method is quite simple: If for example, Alice wants to send one bitcoin to Bob, and Carla wants to send one bitcoin to David, a coin-

---

<sup>8</sup>A false positive exists when an address is wrongly included in a cluster (*i.e.*, all addresses are not controlled by the same entity), and a false negative when an address should be in a cluster but is not.

join transaction could be established whereby the addresses of Alice and Carla are both listed as inputs, and the addresses of Bob and David are listed as outputs in one unique transaction. Thus, when inspecting the 2-to-2 transaction from outside it is impossible to discern who is the sender and who the recipient. In our example, we cannot tell if it is Bob or David who is the recipient of Alice. If this simplified principle were the actual implementation of coinjoin, then obviously the input address method could mistakenly cluster together the addresses of Alice and Carla as if they belonged to the same entity.

However, in practice a coinjoin transaction works in a slightly different way. Indeed, the coinjoin service providers not only shuffle the addresses of the users but they also create new batches of addresses that are then added to the users' addresses and mixed together in the transactions. In our observations, coinjoin addresses could be reused for several times with many other addresses from different users. Thus, new unknown large clusters are created that do not belong to any precise business category because their addresses are very likely linked to more than two distinct entities or directly to coinjoin service providers. Those clusters are similar to "black holes" as they "absorb" addresses that should have been enclosed in other clusters with a clear business profile. To conclude, even in the presence of coinjoin transactions, our method is robust because the likelihood of encountering false positives is practically negligible.

The result of applying the input address heuristic returns more than 30 million clusters, which reduces to a more manageable 2,850 when considering only those composed of at least one hundred addresses and that have received at least one thousand bitcoins from January 2009 through May 2015. As mentioned at the beginning, we label such clusters as super clusters because they represent big agents with strong presence and intense activity in the Bitcoin economy. All together, those super clusters received tens of \$ billions over the study period, at the current exchange rate.

We acknowledge that it is practically impossible to correctly identify all the 2,850 super clusters in the sample. Moreover, our study has the less ambitious aim to ascribe all super clusters to a given broader business category and explore their network of

business relationships. Therefore, from a list of publicly available pre-identified addresses obtained from the Internet we successfully identify 209 super clusters out of 2,850. This subset of known clusters, called the “known group”, is used as the benchmark to identify the business category of the remaining 2,641 clusters in the “unknown group”.

We conclude our study by unveiling the network of payment relationship between these 2,850 super clusters and by exploring the relative interdependence among business categories, as well as their evolution over the study period.

The paper proceeds as follows: Section [2] introduces some preliminary definitions; Section [3] describes the data set we draw upon; Section [4] introduces pure user group (PUG) analysis to classify super clusters in the unknown group; Section [5] presents a transaction pattern (TP) analysis, examining inflows and outflows among those super clusters in known group; Section [6] backtests the results from the PUG analysis; and Section [7] describes the network of entities on the Bitcoin network as discerned from the PUG and TP analyzes and develops the progression of three regimes of the Bitcoin economy.

## 2 Preliminary Definitions

As explained in Section 1, the building block of our analysis is the concept of clustering Bitcoin addresses. In this section we provide a formal definition of clustering by omitting unnecessary technical information which may turn out to be redundant and therefore not useful for the scope of our analysis.

We define the set  $Tx$  of all the Bitcoin transactions, occurred during the period of our analysis, as  $Tx = (tx_1, \dots, tx_i, \dots, tx_z)$ . To each element  $tx_i$  of  $Tx$  corresponds the cluster set  $c_i = (a_1, a_2, \dots, a_n)_i$  containing all the input addresses  $(a_1, a_2, \dots)$  used in the transaction  $tx_i$ . By using a variant of a Union-Find graph algorithm (Cormen et al., 2009), if two or more clusters directly or indirectly (via other clusters) have at least one address in common, we merge those clusters into a single unique one. At the end of

the merging process, we get  $C = \{c_1, \dots, c_x, \dots, c_y, \dots, c_z\}$  which is the set of all disjoint clusters such that  $c_x \cap c_y = \emptyset$  for all  $c_x, c_y \in C$ . Let  $W(C)$  be a finite set  $W(C) = \{w_{xy}(c_x, c_y) \mid c_x, c_y \in C, c_x \neq c_y\} \cup \{w_{xx}(c_x, c_x) \mid \forall c_x \in C\}$ . Then,  $W \subseteq W(C)$  is the set of all (direct) transaction (with loops<sup>9</sup>) between clusters, where  $w_{xy}$  is the total quantity of bitcoins transferred from cluster  $c_x$  to cluster  $c_y$ :

$$w_{xy} = \begin{cases} w_{xy} & \text{if there is a transaction from } c_x \text{ to } c_y. \\ 0 & \text{otherwise.} \end{cases}$$

We define a super cluster,  $\hat{c}_x$ , as any special cluster that belongs to the partition  $\hat{C} \subset C$ :

$$\hat{C} = \left\{ c_x \in C \mid \sum_{h=1}^z w_{hx}(c_h, c_x) \geq 1,000 \text{ BTC} \wedge n(c_x) \geq 100 \right\}. \quad (1)$$

where  $n(c_x)$  denotes the number of addresses in cluster  $x$ .

According to our definition, a super cluster is any cluster that satisfies the following two thresholds: 1) having received at least 1,000 bitcoins during our research window; and 2) is composed of at least 100 unique addresses.

The first threshold is necessary in order to increase the likelihood of excluding inactive entities from the analysis. The second threshold is necessary to exclude as many private individuals as possible, who typically own only one or a few addresses. Together, these thresholds increase the robustness of the transaction pattern analysis of the clusters in Section [5], which is based on statistics requiring big enough data. In fact, smaller clusters composed of some tens, or even some hundreds of addresses are only able to generate a trivial amount of transaction data, giving us insufficient information to perform a meaningful analysis.

---

<sup>9</sup>Indeed, fork-merge patterns and self loops represent a frequent scenario in the Bitcoin economy, *e.g.*, Tasca (2015) and Ron and Shamir (2013).

### 3 Data Set

In our study we parsed data from the Bitcoin Core over the period of the 3rd of January 2009 (block 0) through the 8th of May 2015 (block 355551).<sup>10</sup> Over this interval, the Bitcoin network proliferates both in terms of number of addresses and in terms of number of transactions. See Table 1 for a summary of our data set. All the data related to Bitcoin transactions are imported into and managed via a MySQL database designed to run high-performing application codes (see the diagram in Figure 14 in Appendix A).

| <b>Bitcoin Core parsed from the 3rd of January 2009 until 8th of May 2015</b> |             |
|---|-------------|
| Max block height  | 355,551     |
| Total number of transactions  | 68,030,042  |
| Total number of input   | 172,743,139 |
| Total number of output  | 194,476,567 |
| Total addresses identified  | 75,191,953  |
| Total clusters identified(include at least one address)                       | 30,708,660  |
| Number of clusters with at least 2 addresses                                  | 9,847,999   |
| Total transactions between clusters   | 88,950,021  |

Table 1: Blockchain Database Facts. Source: Bitcoin Core.

By applying the input address heuristic, 75,191,953 unique Bitcoin addresses are grouped into 30,708,660 clusters, of which, about two-thirds are clusters composed of only a single address, as shown in Table 2.

---

<sup>10</sup>Bitcoin Core was, by far, the most dominant version of the Bitcoin blockchain over the study period.

| <b>Input Address Heuristic: Clustering Result</b> |                               |
|---|-------------------------------|
| Number of addresses included in each cluster      | Number of clusters identified |
| > 10001   | 194                           |
| 1001~10000  | 1,145                         |
| 101 ~ 1000  | 12,185                        |
| 11~100  | 436,093                       |
| 2 ~ 10  | 9,398,382                     |
| =1  | 20,860,661                    |
| Total number of clusters                          | 30,708,660                    |

Table 2: The clusters identified with the input address heuristic are grouped per number of addresses composing them. Source: Bitcoin Core.

Then, by applying the criteria defined in Equation (1), 2,850 super clusters are filtered out. Figure 1 shows the network of super clusters  $\hat{C}$  and their transactions among each other, as well as with all the remaining clusters in  $C \setminus \hat{C}$ .

By gathering publicly available address information, we are able to link part of super clusters  $\hat{c}_x \in \hat{C}$  to real world entities (*e.g.*, BTCChina, Kraken, Xapo) which belong to different business categories. Specifically, we gathered 359,776 deciphered addresses from Walletexplorer (2015) and Blockchain (2015b). According to their entity information, we could compose a group of deciphered sets of addresses,  $P = \{p_1, p_2, \dots, p_\gamma, \dots\} = \bigcup_{\gamma \in \Gamma} p_\gamma$ . Precisely,  $p_\gamma = \{a \mid a \text{ belong to the known beneficial owner } \gamma\}$  is the set of addresses that belong to the beneficial owner  $\gamma$  whose identity is publicly available from the Internet.<sup>11</sup>

<sup>11</sup>For example,  $p_{Huobi}$  is the set of addresses associated to Huobi with  $n(p_{Huobi}) = 37,756$ .

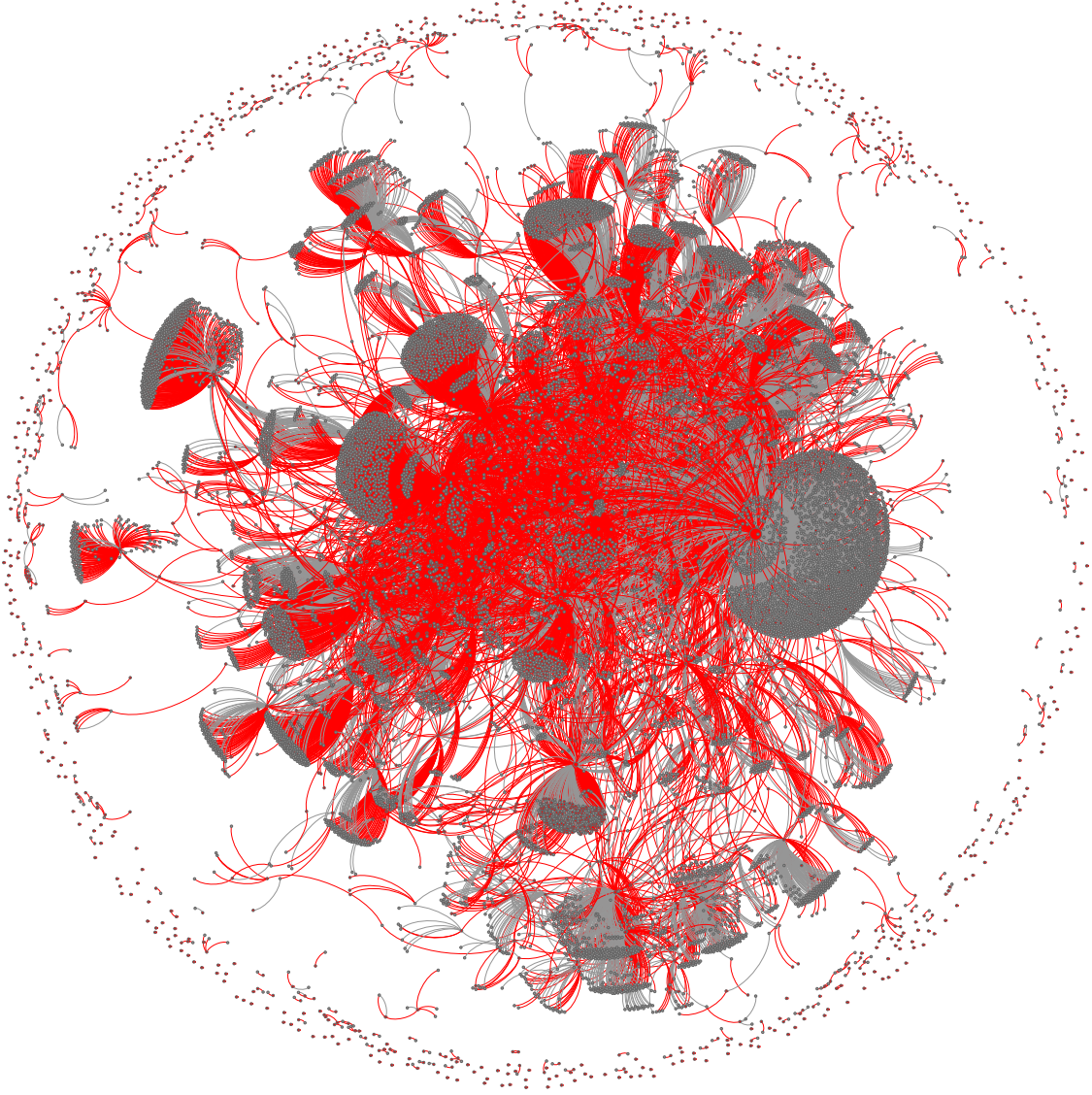


Figure 1: Network visualization of the interactions of the super clusters in  $\hat{C}$  with each other and also with all the remaining clusters in  $C \setminus \hat{C}$ . Every red node represents a single super cluster  $\hat{c}_x \in \hat{C}$  and every grey node represents a counterpart cluster. For visualization purposes, we set a threshold of at least 1,000 BTC transferred between a super cluster and its counterpart. Therefore, the plot shows 1,957 super clusters out of 2,850 in  $\hat{C}$ . One may observe that some of the clusters are highly connected to each other, although many are isolated. These isolated entities could be individuals, some highly self-contained businesses, or some clusters that belong to active business entities but which are kept untied from the others, *i.e.*, used independently for purposes different from the main business activity.

Thus, depending on whether a super cluster hold at least one address belonging to any  $p_\gamma \in P$  or not,  $\hat{C}$  is then decomposed into either a known group,  $\hat{C}^K$ , or a unknown group,  $\hat{C}^U$ . Formally,

$$\hat{C} = \hat{C}^K \cup \hat{C}^U \quad (2)$$



with  $\hat{C}^K \cap \hat{C}^U = \emptyset$  by definition and

$$\hat{c}_x \in \begin{cases} \hat{C}^K & \text{if } \hat{c}_x \cap p_\gamma \neq \emptyset \wedge \hat{c}_x \cap (P \setminus p_\gamma) = \emptyset, \forall \gamma \in \Gamma \\ \hat{C}^U & \text{if } \begin{cases} \hat{c}_x \cap p_\gamma \neq \emptyset \wedge \hat{c}_x \cap (P \setminus p_\gamma) \neq \emptyset, \forall \gamma \in \Gamma \\ \hat{c}_x \cap p_\gamma = \emptyset, \forall \gamma \in \Gamma. \end{cases} \end{cases} \quad (3)$$

The matching exercise turns out the following result:  $n(\hat{C}^K)=209$  and  $n(\hat{C}^U)=2,641$  such that  $n(\hat{C}^K) + n(\hat{C}^U) = n(\hat{C})= 2,850$ .<sup>12</sup> As a side note, we remark that Equation (3) follows a prudential principle that aims to avoid false positives. Namely, any cluster in  $\hat{C}$  that has addresses linked to more than one set  $p_\gamma \in P$ , is considered unknown and confined to the set  $\hat{C}^U$ .<sup>13</sup>

Then, according to their business model, each identified super cluster is allocated into one of the following primary business categories: *exchange*  $\hat{C}^{KX}$ , *mining pool*  $\hat{C}^{KP}$ , *online gambling*  $\hat{C}^{KG}$ , *black market*,  $\hat{C}^{KB}$ . Besides these big four business categories which are populated by economic entities with a clear business profile, there are also few other economic entities with a business models (*e.g.*, bitcoin wallets) ethereogeneous among them and disparate from the previous ones. Then, we classify them into the category *others*,  $\hat{C}^{KO}$ .

Table 8 in the Appendix B shows us the results, namely,  $n(\hat{C}^{KX}) = 104$ ,  $n(\hat{C}^{KP}) = 18$ ,  $n(\hat{C}^{KG}) = 45$ ,  $n(\hat{C}^{KB}) = 13$  and  $n(\hat{C}^{KO}) = 29$  such that  $n(\hat{C}^{KX}) + n(\hat{C}^{KP}) + n(\hat{C}^{KG}) + n(\hat{C}^{KB}) + n(\hat{C}^{KO}) = n(\hat{C}^K) = 209$ .

<sup>12</sup>See Figure 15 in Appendix A for a visualization of the problem we aim at solving.

<sup>13</sup>As an example, one of the biggest clusters holding about 6 million addresses which probably should have been included in  $\hat{C}^K$  is instead included in  $\hat{C}^U$  because although it has 2 million addresses linked to the MtGox exchange, it has one address linked to bitcoin-24.

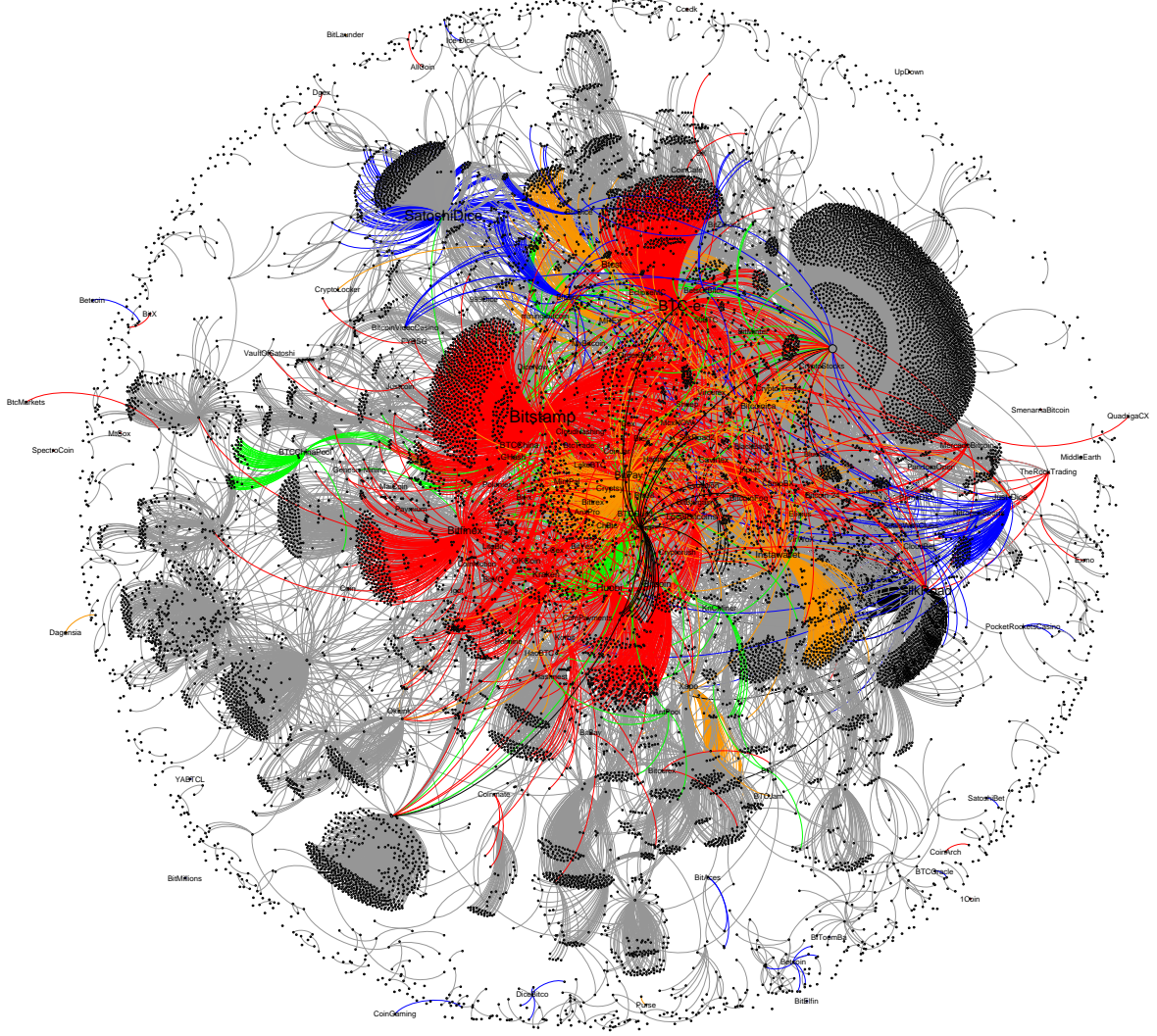


Figure 2: Network visualization of the 209 super clusters in the set  $\hat{C}^K$  that have been identified by cross linking the known addresses in the set  $P$  with the addresses in each  $\hat{c}_x \in \hat{C}$ . As it happens that more than one super cluster may belong to single entities in  $\Gamma$ , we combine them into one node in the network. For visualization purposes, we set a threshold of at least 5,000 BTC being transferred between a super cluster  $\hat{c}_x \in \hat{C}$  and its counterpart. Therefore, the plot shows only 94 super clusters out of 209 in  $\hat{C}^K$ . The grey nodes are the counterparts of each  $\hat{c}_x \in \hat{C}$ . Each super cluster is colored according to its business category: green for miners, red for exchange, blue for gambling, orange for others, black for black market, purple for composite category, and grey for the clusters which are the counterparts. The color of the edge is the same as the source market nodes. One may clearly observe some large entities with many counterparts, such as Silkroad (black market), SatoshiDice (online gambling), BitStamp (exchange) and BTC-e (exchange).

Figure 2 shows the payment network of the 209 identified super clusters<sup>14</sup> in  $\hat{C}^K$ . In the next sections, we will use the information on the super clusters in the set  $\hat{C}^K$  together with the information on their interactions with all the other clusters in  $C \setminus \hat{C}^K$  to derive

<sup>14</sup>Super clusters linked to the same real world entity are merged into one node in the network.

the business membership of each unknown super cluster in  $\hat{C}^U$ .

## 4 Pure User Group Analysis

The pure user group (PUG) analysis is carried out to classify (into specific business categories) super clusters in the unknown group and it is based on the definition and classification of “pure” users. By pure users we mean all those clusters populating the Bitcoin economy (except for those already in the known group) that had bilateral transactions with super clusters (in the known group) belonging to only one business category.

In other words, for each specific business category we build a correspondent PUG: 1) clusters having transactions only with exchanges in  $\hat{C}^{KX}$  are classified in the PUG *traders*; 2) clusters having transactions only with gambling services in  $\hat{C}^{KG}$  are classified in the PUG *gamblers*; and 3) clusters having transactions only with black market services in  $\hat{C}^{KB}$  are classified in the PUG black market *user-dealers*.

The classification of the clusters into different pure user groups is the first step of the PUG analysis. The second step consists of classifying the super clusters in the unknown group into a specific business category in the case they transact *only* with the corresponding specific PUG. For example, those super clusters in the unknown group that had transactions only with traders are classified as exchanges and so on also for the other categories. However, the clusters in the known group identified in the categories mining pools and others follow a peculiar business model. Thus, we do not create the set of PUGs having transactions only with mining pools in  $\hat{C}^{KP}$  because the mining pools in the unknown group will be identified via the coinbase analysis (see Section 4.2). Similarly, we do not create the set of PUGs having transactions only with others in  $\hat{C}^{KO}$  because those clusters do not have a clearly defined business profile. In other terms, to avoid false positives, we will not try to classify super clusters in the unknown group into the category others.

## 4.1 PUG Identification

In this first part of the analysis we consider only the following sets  $\hat{C}^{KX}$ ,  $\hat{C}^{KG}$  and  $\hat{C}^{KB}$ . Accordingly, we introduce the following set notation:  $U^X \subset C \setminus \hat{C}^K$  is the subset of pure traders that had transactions only with exchanges in  $\hat{C}^{KX}$ ;  $U^G \subset C \setminus \hat{C}^K$  is the subset of pure gamblers that had transactions only with gambling sites in  $\hat{C}^{KG}$ ;  $U^B \subset C \setminus \hat{C}^K$  is the subset of pure black market user-dealers that had transactions only with black markets in  $\hat{C}^{KB}$ . Formally:

$$U^X = \{c_x \in C \setminus \hat{C}^K \mid \exists \hat{c}_y \in \hat{C}^{KX} : (w_{xy}(c_x, \hat{c}_y) > 0 \vee w_{yx}(\hat{c}_y, c_x) > 0) \\ \wedge (w_{xj}(c_x, c_j) = 0 \vee w_{jx}(c_j, c_x) = 0), \forall c_j \in \hat{C}^K \setminus \hat{C}^{KX}\}. \quad (4)$$

$$U^G = \{c_x \in C \setminus \hat{C}^K \mid \exists \hat{c}_y \in \hat{C}^{KG} : (w_{xy}(c_x, \hat{c}_y) > 0 \vee w_{yx}(\hat{c}_y, c_x) > 0) \\ \wedge (w_{xj}(c_x, c_j) = 0 \vee w_{jx}(c_j, c_x) = 0), \forall c_j \in \hat{C}^K \setminus \hat{C}^{KG}\}. \quad (5)$$

$$U^B = \{c_x \in C \setminus \hat{C}^K \mid \exists \hat{c}_y \in \hat{C}^{KB} : (w_{xy}(c_x, \hat{c}_y) > 0 \vee w_{yx}(\hat{c}_y, c_x) > 0) \\ \wedge (w_{xj}(c_x, c_j) = 0 \vee w_{jx}(c_j, c_x) = 0), \forall c_j \in \hat{C}^K \setminus \hat{C}^{KB}\}. \quad (6)$$

This first part of the PUG analysis returns the following results:  $n(U^X) = 440,434$ ,  $n(U^G) = 415,528$  and  $n(U^B) = 74,233$ .

The statistics of the bitcoin transactions between pure users and clusters in the known group reveal that the average volume per transaction differs substantially with respect to each business category: The average volume per transaction from/to traders to/from exchanges is 20 BTC; the average volume per transaction from/to gamblers to/from gambling services is 0.5 BTC; and finally, the average volume per transaction from/to user-dealers to/from black market services is 3 BTC (see Table 3).

| Statistics for PUG Transaction |                 |                             |                           |                             |                           |
|--------------------------------|-----------------|-----------------------------|---------------------------|-----------------------------|---------------------------|
|                                |                 | PUG $\rightarrow \hat{C}^K$ |                           | $\hat{C}^K \rightarrow$ PUG |                           |
| PUG                            | Num of clusters | Avg tx volume (BTC)         | Avg tx interval (minutes) | Avg tx volume (BTC)         | Avg tx interval (minutes) |
| $U^X$                          | 440,434         | 23.4                        | 20,529.5                  | 17.6                        | 10,685.0                  |
| $U^G$                          | 415,528         | 0.5                         | 528.3                     | 0.5                         | 387.3                     |
| $U^B$                          | 74,233          | 2.7                         | 22,151.3                  | 3.4                         | 9,394.0                   |

Table 3: For PUG in different categories, the table summarizes the average transaction amount(BTC) and average transaction interval(minutes).

## 4.2 Identification of Mining Pools

The super clusters in the category mining pool are identified without utilizing the PUG analysis. Indeed, each newly generated Bitcoin block includes a reward to the successful miner: an amount equal to the sum of the block reward (or subsidy), *i.e.* newly available bitcoins, plus any accumulated fees paid by transactions included in that block. To allocate this sum, a new generation transaction is created whose input, called the “coinbase”, contains the reward for the miners. Thus, unlike all other transaction inputs, the coinbase is not linked to any previous output. This feature offers a simple and direct method to identify those clusters belonging to the mining category by filtering out the transactions with “null” input and only one output.

Let  $\hat{C}^{UCoinbase}$  be the set of clusters (in the unknown group) composed (also, but not only) of addresses with coinbase inputs,  $n(\hat{C}^{UCoinbase}) = 575$ . Not all these 575 clusters, however, can reliably be defined as mining pools; for some of them, mining is not their primary activity and rewards from coinbase transactions represent only a small percent of their activity. To make sure that the taxonomy is robust, we classify only clusters in the unknown group whose mining rewards occupy more than 80% of its total income, as mining pool,  $\hat{C}^{UP}$ . The remaining clusters  $\hat{C}^{UP} = \hat{C}^{UCoinbase} \setminus \hat{C}^{UP}$  that cannot be defined as mining pools according to our threshold are instead classified via the PUG analysis.

### 4.3 Classification of Unknown Super Clusters

The principle of PUG classification for unknown clusters is straightforward and works as follows: If one super cluster in the unknown group transacts only with one specific PUG, then we suspect that this cluster belongs to the business category correspondent to that specific PUG. For example, if during the period January 2009 - May 2015 one super cluster in  $\hat{C}^U$  records transactions with one or more traders in  $U^X$  but not with gamblers in  $U^G$  and user-dealers in  $U^B$ , it is classified as an exchange. One should note that this does not rule out the possibility for the exchange to transact with any other cluster in  $C$  beyond those in  $U^X$ . The clusters who transact with multiple PUGs are identified in the composite category,  $\hat{C}^{UM}$ , which implies those super clusters might have multi-business lines.

Let

$$\hat{C}^{U\bar{X}} = \{\hat{c}_x \in \hat{C}^U \mid \exists c_y \in U^X : (w_{yx}(c_y, \hat{c}_x) > 0 \wedge w_{xy}(\hat{c}_x, c_y) > 0)\} \quad (7)$$

be a broad subset of exchanges in  $\hat{C}^U$  that have transactions *not only* with traders.

Let

$$\hat{C}^{U\bar{G}} = \{\hat{c}_x \in \hat{C}^U \mid \exists c_y \in U^G : (w_{yx}(c_y, \hat{c}_x) > 0 \wedge w_{xy}(\hat{c}_x, c_y) > 0)\} \quad (8)$$

be a broad subset of gambling services in  $\hat{C}^U$  that have transactions *not only* with gamblers.

Let

$$\hat{C}^{U\bar{B}} = \{\hat{c}_x \in \hat{C}^U \mid \exists c_y \in U^B : (w_{yx}(c_y, \hat{c}_x) > 0 \wedge w_{xy}(\hat{c}_x, c_y) > 0)\} \quad (9)$$

be a broad subset of black market services in  $\hat{C}^U$  that have transactions *not only* with user-dealers.

Then, the subset of exchanges in  $\hat{C}^U$  that have transactions only with traders in  $U^X$

is:

$$\hat{C}^{UX} = \{\hat{c}_x \in \hat{C}^{U\bar{X}} \setminus (\hat{C}^{UG} \cup \hat{C}^{UB} \cup \hat{C}^{UP})\}. \quad (10)$$

Similarly, the subset of gambling services in  $\hat{C}^U$  that have transactions only with gamblers in  $U^G$  is:

$$\hat{C}^{UG} = \{\hat{c}_x \in \hat{C}^{U\bar{G}} \setminus (\hat{C}^{U\bar{X}} \cup \hat{C}^{UB} \cup \hat{C}^{UP})\}. \quad (11)$$

The subset of black market services in  $\hat{C}^U$  that have transactions only with user-dealers in  $U^B$  is:

$$\hat{C}^{UB} = \{\hat{c}_x \in \hat{C}^{U\bar{B}} \setminus (\hat{C}^{U\bar{X}} \cup \hat{C}^{UG} \cup \hat{C}^{UP})\}. \quad (12)$$

Finally, the subset of multi-business clusters in  $\hat{C}^U$  that have transactions with more than one user group is:

$$\hat{C}^{UM} = \{\hat{c}_x \in ((\hat{C}^{U\bar{X}} \cup \hat{C}^{UB} \cup \hat{C}^{UG} \cup \hat{C}^{UP}) \setminus (\hat{C}^{UX} \cup \hat{C}^{UG} \cup \hat{C}^{UB} \cup \hat{C}^{UP}))\}. \quad (13)$$

Table 4 shows that  $n(\hat{C}^{UX}) = 310$ ,  $\hat{C}^{UG} = 755$ ,  $\hat{C}^{UB} = 41$ ,  $\hat{C}^{UP} = 57$  and  $\hat{C}^{UM} = 630$ .

| Tagged Cluster in Unknown Group |                    |
|---------------------------------|--------------------|
| Category                        | Number of clusters |
| $n(\hat{C}^{UX})$               | 310                |
| $n(\hat{C}^{UG})$               | 755                |
| $n(\hat{C}^{UP})$               | 57                 |
| $n(\hat{C}^{UB})$               | 41                 |
| $n(\hat{C}^{UM})$               | 630                |

Table 4: The number of clusters tagged with the PUG method. To give the reader a complete view,  $\hat{C}^{UP}$  is also listed here, which is identified from coinbase transactions.

## 5 Transaction Pattern Analysis

In this section we introduce a transaction pattern (TP) analysis to study the different transaction patterns of the super clusters in set  $\hat{C}^K$  (listed in Table 8). The TP analysis is used to garner more insights into stylized facts characterizing the distinct business behaviors of the super clusters. Moreover, the TP analysis is used in Section 6 to measure the accuracy of the PUG analysis by testing the pattern similarity between the clusters in  $\hat{C}^K$  and those in  $\hat{C}^U$ . In the following we divide the TP analysis in *inflow* and *outflow* analysis.

### 5.1 Inflow Analysis

The inflow analysis consists of examining the properties of the transactions *toward* any super clusters in the known group  $\hat{C}^K$ . We select the transactions in the set:

$$\vec{W}^K \subset W = \{w_{yx}(c_y, \hat{c}_x) \in W \mid c_y \in C \setminus \hat{C}^K, \hat{c}_x \in \hat{C}^K, n_{yx} \geq 100\} \quad (14)$$

where  $n_{yx}$  denotes the number of transactions from cluster  $c_y$  to cluster  $\hat{c}_x$  during the period of the analysis. According to Equation (14), a pair  $(c_y, \hat{c}_x) \in \vec{W}^K$  is only considered if there has been at least 100 transactions from  $c_y$  to  $\hat{c}_x$ . This minimum transaction threshold is subjective and shall be set to any value able to ensure that the descriptive statistics calculated are robust. In our case, with  $n_{yx} \geq 100$  we obtain that  $n(\vec{W}^K) = 11,899$ , involving 148 super clusters in  $\hat{C}^K$ . After having defined the set of analysis, we calculate the median of transaction volume and the median of time interval in minutes for each pair  $(c_y, \hat{c}_x) \in \vec{W}^K$ .<sup>15</sup> Each dot in Figure 3 represents the measurement for one pair  $(c_y, \hat{c}_x) \in \vec{W}^K$ : red if  $\hat{c}_x \in \hat{C}^{KX}$ , green if  $\hat{c}_x \in \hat{C}^{KP}$ , blue if  $\hat{c}_x \in \hat{C}^{KG}$ , and black if  $\hat{c}_x \in \hat{C}^{KB}$ . The  $x$ -axis is the median of the transaction volume and the  $y$ -axis is the median of the time interval (in minutes) between inflow transactions for each pair

---

<sup>15</sup>For example, if the median value of intervals is 60 minutes, this means that counterparts tends to send to  $\hat{c}_x$  bitcoins every 60 minutes.



$$(c_y, \hat{c}_x) \in \vec{W}^K.$$

Figure 3 shows some clustering effects; we can see that for each of our four identified business categories there exists specific patterns of transaction behavior. For example, there are clearly plotted in blue, vertical lines at  $x=0.01$ ,  $0.02$  and so on. To capture this more clearly, we plot the kernel density in Figure 4.<sup>16</sup> This illustrates a notable characteristic for gambling behavior, that is gamblers tend to place bets with similar, round lot amounts (*i.e.* 0.1, 0.5, 1.0, etc.) again and again, with wagers of 0.01 BTC being placed most frequently. Gamblers may be accustomed to wagering in round amounts in traditional settings using casino or poker chips with specified round values (*e.g.*, \$1, \$5, or \$25), or online using virtual chips. Individuals may carry forward that behavior to bitcoin-based gambling even in instances where the size of bets are determined arbitrarily by the gambler placing bets.<sup>17</sup>

With respect to exchanges, although it is less obvious due to some overlap in the plot, we are still able to see some vertical lines in red at  $x=0.1$ ,  $0.5$ , and  $1.0$ , which indicates that the traders usually deposit into exchanges round amounts of bitcoins, rather than random amounts in order to presumably exchange them for fiat or alternative digital currency. In other words, it appears that traders may wait until they have accumulated some even amount, most commonly 1.0 BTC, before selling them.

Inflows to black markets show a wider variety of arbitrary transaction size, but still also show marked preference for round lots of bitcoin, notably at amounts of 0.1, 0.2, 0.3, 0.5, and 1.0 BTC. This may suggest that black market sellers explicitly place round lot prices on their items as a matter of doing business. Prescription and illegal drugs are notably sold on black markets, and this indicates that sellers will offer an amount of contraband that corresponds to a round price (say, 1.0 BTC), rather than determining what the price would be for a fixed quantity (say, for 1 ounce).<sup>18</sup>

---

<sup>16</sup>In order to capture exact density on point, a very precise width is needed. In this case, we set the width 0.00000001 BTC (or 1 satoshi).

<sup>17</sup>This is true, for example, in SatoshiDice, the largest bitcoin-based gambling service.

<sup>18</sup>This practice is common in transactions involving small amounts of street drugs where a "dime bag" is whatever quantity \$10 buys and a "nickel bag" whatever \$5 buys.

Mining pools exhibit a more or less random pattern of inflows, since a mining pool will only be credited with small amounts of bitcoin whenever it finds a new block of bitcoin. When this happens, the pool will generally extract a small profit consisting of either a nominal percentage of the block reward, or of the transaction fees associated with that block, or both.

In addition to studying patterns in the amounts of bitcoin inflows, we also consider transaction intervals. We observe a large density of dots, plotted in red, clustering horizontally just above  $y=1,000$  in Figure 3, specifically at 1,440 minutes, which is the number of minutes in one day. What this shows us is that there are a large number of traders who send small amounts of bitcoins to exchanges regularly each day. We suspect that these could be small miners who exchange mined bitcoins for cash on a daily basis, or “day traders” who are active daily but go home flat, having sold out any positions in bitcoin to avoid overnight price volatility. Figure 5 clarifies this effect, and shows the kernel density of the intervals between transactions (band=100 minutes). The 1,440 minute interval is prominent not only for traders to exchanges but also for the other business categories, suggesting that a “one-day” holding period for bitcoin transactions is somewhat typical; a one-day effect where traders, gamblers, black market participants, and miners tend to cash out on a daily basis.

We observe however, that gambling has, by far, the shortest interval as well as the highest transaction frequency. This is not difficult to understand, as gamblers can ante or re-bet many times in a matter of minutes.

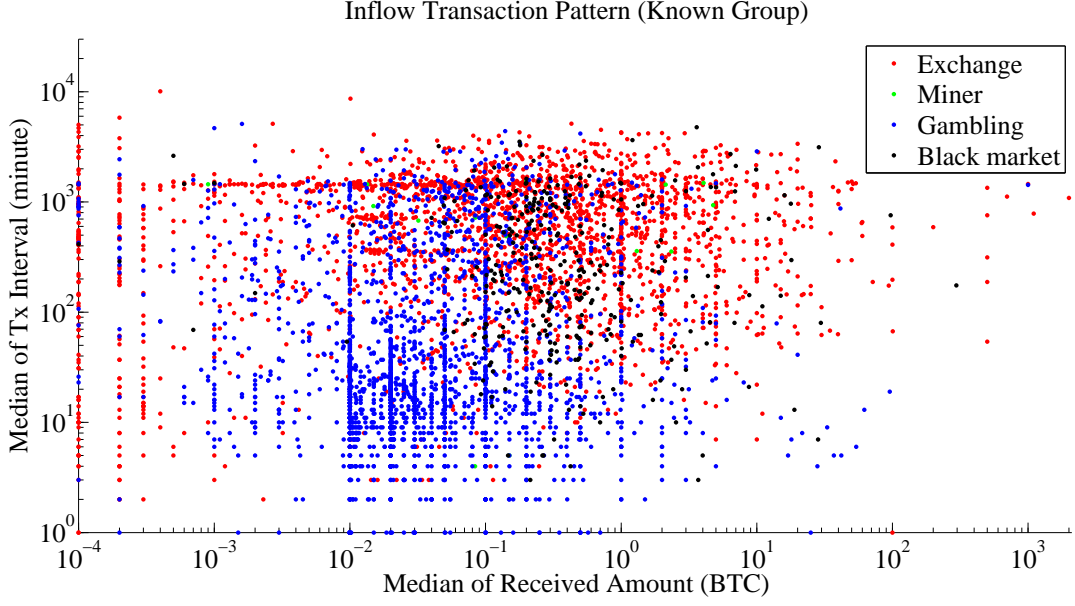


Figure 3: Inflow transaction pattern for the known group. Each dot characterizes one pair of clusters  $(c_y, \hat{c}_x) \in \vec{W}^K$ . The  $x$ -axis is the median transaction volume of all transactions between all the pairs of clusters  $\in \vec{W}^K$  during the period January 2009 – May 2015. The  $y$ -axis is the median transaction interval (in minutes) of the transactions between all the pairs of clusters  $\in \vec{W}^K$ .

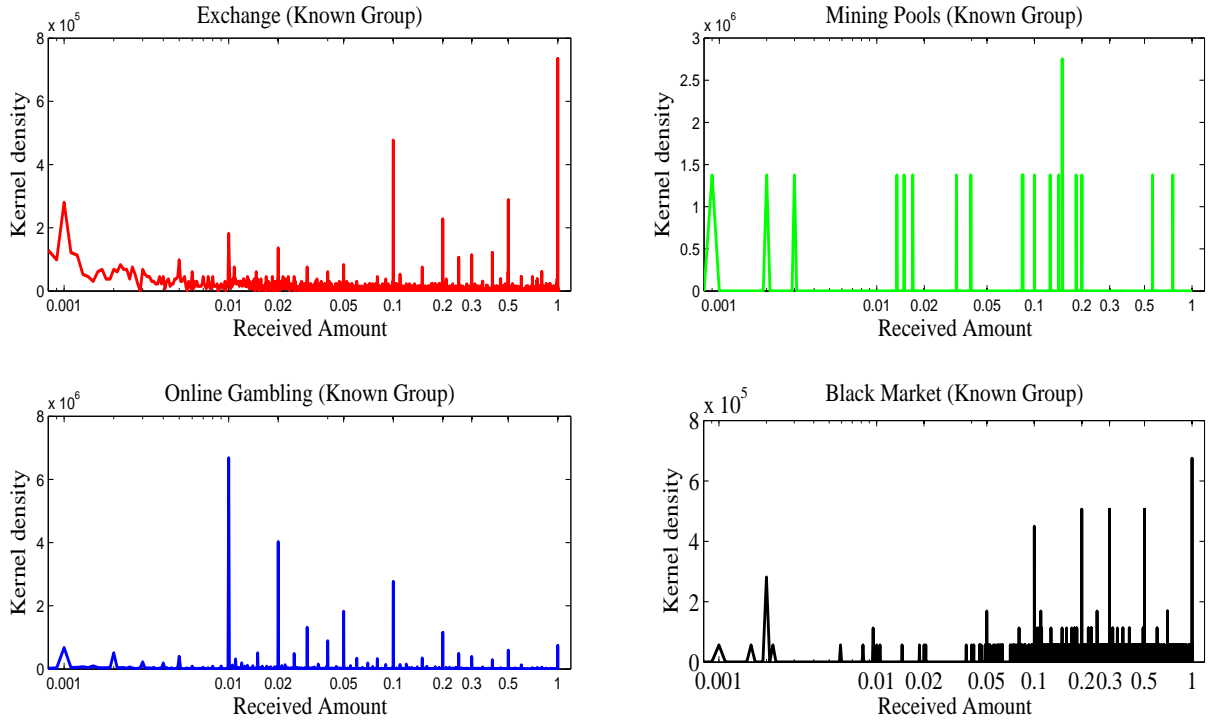


Figure 4: Kernel density of the inflow amount of bitcoins received by any  $c_y \in C \setminus \hat{C}^K$  during the period January 2009 – May 2015 by each cluster  $\hat{c}_x \in \hat{C}^k$  in the known group.

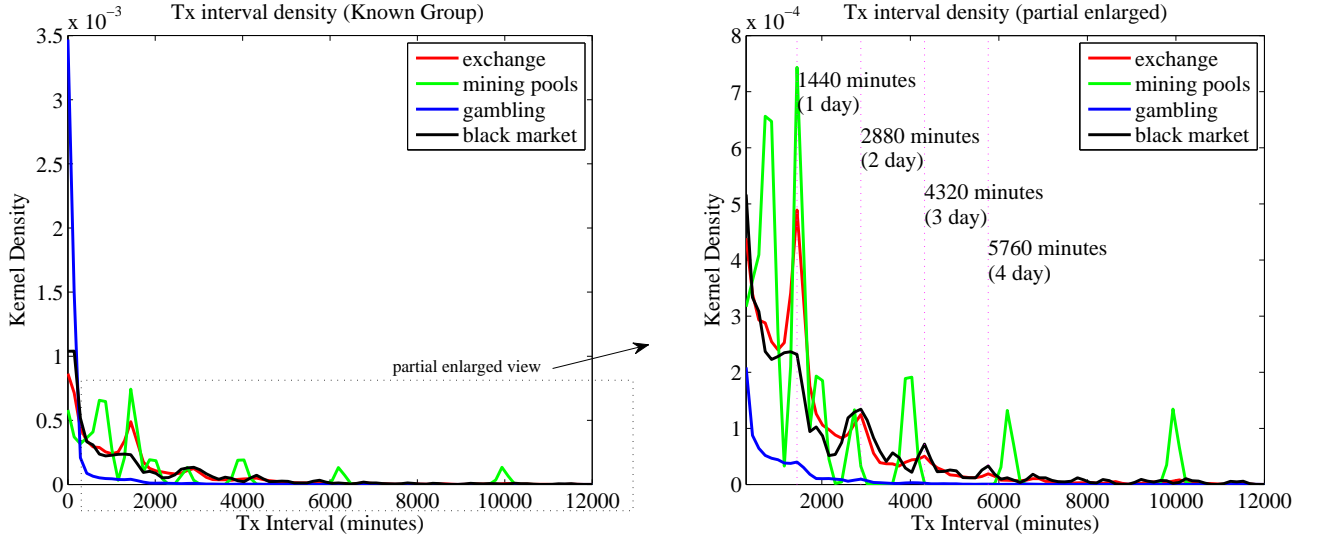


Figure 5: Kernel density of the inflow transaction intervals between subsequent transactions during the period January 2009 – May 2015 for each category in the known group.

## 5.2 Outflow Analysis

The outflow analysis consists of examining the properties of the transactions *from* the clusters in the known group  $\hat{C}^K$ . As for the inflow analysis we measure the median of transaction volume and the median of time interval in minutes. Additionally, we measure the median number of inputs and outputs in the transactions between each pair of clusters. To examine the transaction outflow from the super clusters  $\hat{c}_x \in \hat{C}^K$ , we select the transactions in the set:

$$\overleftarrow{W}^K \subset W = \{w_{xy}(\hat{c}_x, c_y) \in W \mid \hat{c}_x \in \hat{C}^K, c_y \in C \setminus \hat{C}^K, n_{xy} \geq 100\} \quad (15)$$

where  $n_{xy}$  denotes the number of transactions from cluster  $\hat{c}_x$  to cluster  $c_y$  during the period of the analysis. According to Equation (15), a pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$  is considered only if there has been at least 100 transactions from  $\hat{c}_x$  to  $c_y$ . From our database we obtain that  $n(\overleftarrow{W}^K) = 16,188$ , involving 148 super clusters in  $\hat{C}^K$ . To extract information about the number of inputs/outputs of the transactions in  $\overleftarrow{W}^K$ , we plot the median number

of inputs/outputs in all the transactions between each pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$ . Each dot in Figure 6 represents a value set for pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$ : red if  $\hat{c}_x \in \hat{C}^{KX}$ , green if  $\hat{c}_x \in \hat{C}^{KP}$ , blue if  $\hat{c}_x \in \hat{C}^{KG}$ , and black if  $\hat{c}_x \in \hat{C}^{KB}$ . The  $x$ -axis represents the median number of inputs, and the  $y$ -axis represents the median number of outputs.

Only the mining pools show a significantly distinct transaction pattern from the others. Specifically, the outflow transactions for most of the mining pools are characterized by no more than ten inputs, but at the same time by a large amount of outputs, ranging from tens to thousands. This is consistent with the business model of mining pools: After successfully mining bitcoins, the mining pools will distribute the reward to all the small miners who have contributed some mining effort. So, the number of outputs is much larger than the number of inputs. One could speculate on the size of these mining pools according to the number of outputs in each outflow transaction.

As done for the inflow analysis, for each pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$ , we calculate the median transaction volume and the median time interval in minutes. The  $x$ -axis in Figure 7 represents the median transaction volume for each pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$ , while the  $y$ -axis represents the median time interval (in minutes) between outflow transactions for each pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$ . The blue dots scattered around the bottom-left area of the plot imply that gambling clusters send relatively small amounts of bitcoin but at a high-frequency to their counterparts.<sup>19</sup>

The outflow transaction interval is plotted in Figure 8, which also shows the one-day effect for mining pools. The combination of the results from Figures 5 and 8 reveal a clear stylized fact characterized by many small miners receiving daily rewards from mining pools and then exchanging those rewards for fiat currency on exchange platforms.

---

<sup>19</sup>This feature is consistent with the results in the former inflow analysis.

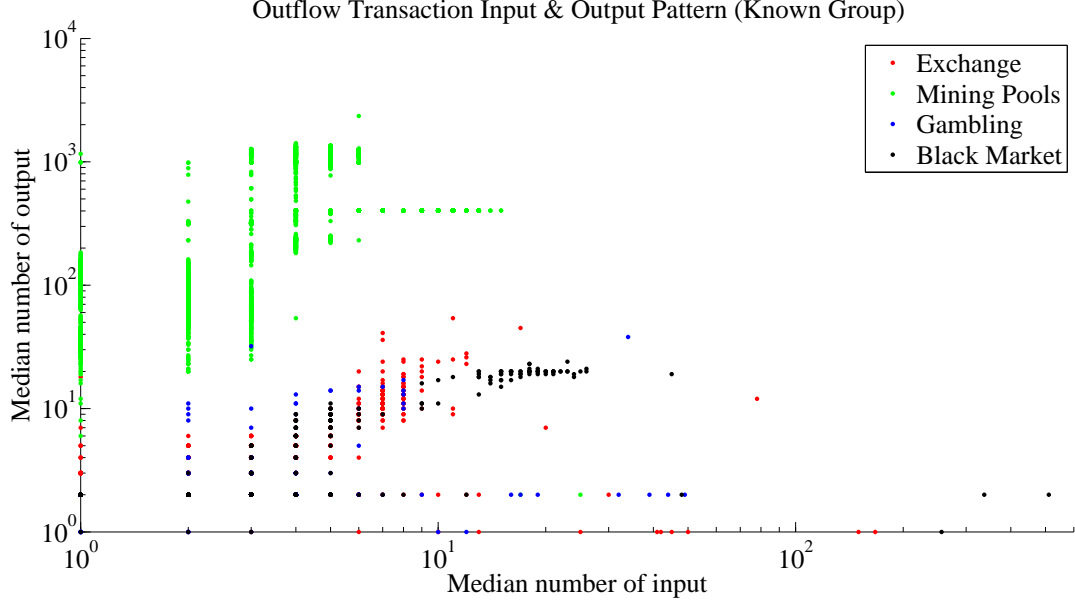


Figure 6: Median number of inputs and outputs for all the transactions among each pair  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$  during the period January 2009 – May 2015. Each dot characterizes one pair of clusters. The  $x$ -axis measures the median number of inputs and the  $y$ -axis measures the median number of outputs.

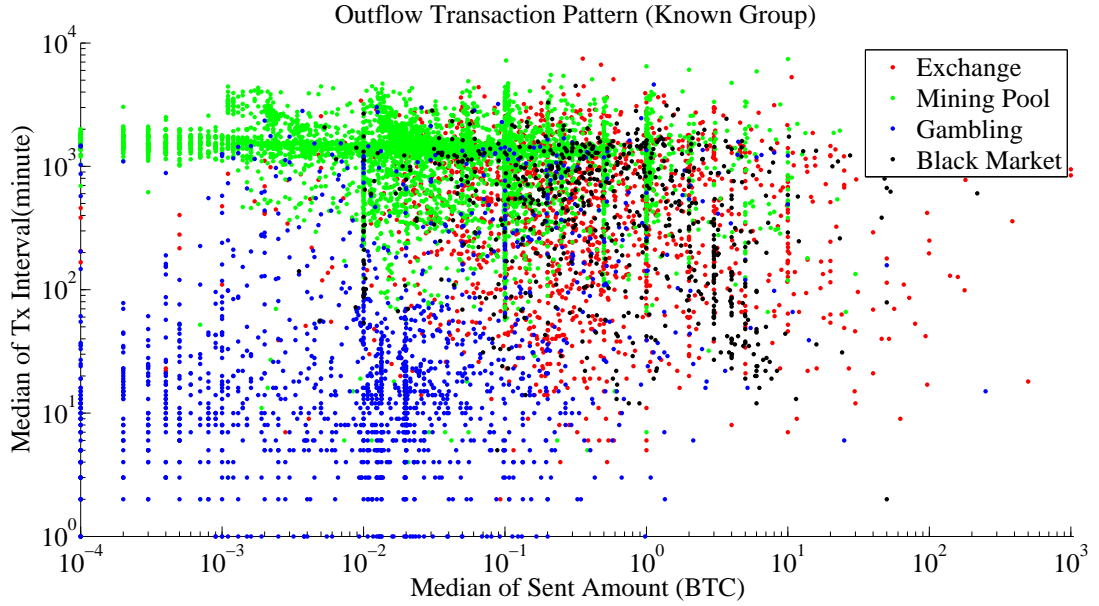


Figure 7: Outflow transaction pattern for the known group. Each dot characterizes one pair of clusters  $(\hat{c}_x, c_y) \in \overleftarrow{W}^K$ . The  $x$ -axis is the median transaction volume of all transactions between all the pairs of clusters  $\in \overleftarrow{W}^K$  during the period January 2009 – May 2015. The  $y$ -axis is the median transaction interval (in minutes) of the transactions between all the pairs of clusters  $\in \overleftarrow{W}^K$ .

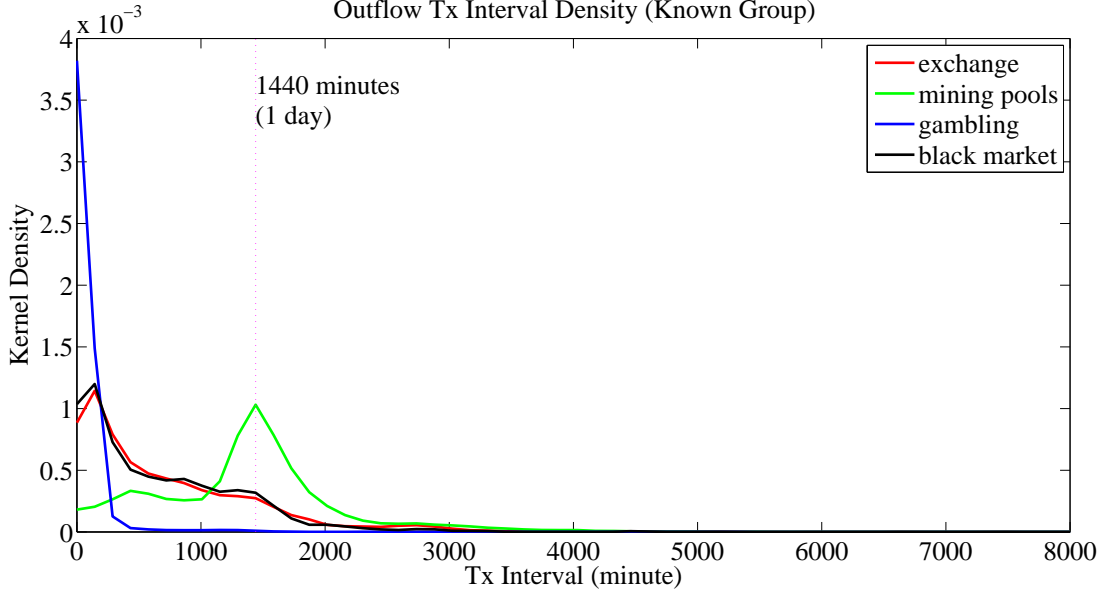


Figure 8: Kernel density of the outflow amount of bitcoins sent to any  $c_y \in C \setminus \hat{C}^K$  during the period January 2009 – May 2015 by each cluster  $\hat{c}_x \in \hat{C}^k$  in the known group.

## 6 PUG Control Test with TP Analysis

In this section, we test the results of the PUG classification conducted in Section 4 for clusters in the unknown group by analyzing whether they exhibit pattern similarities with the clusters in the known group. The test is twofold and is based on the translation in matrix form of the inflow transaction patterns and the outflow transaction patterns involving super clusters in the known group. Each matrix is then compared, via a 2-D correlation<sup>20</sup> analysis, with the correspondent one related to the inflow and outflow transaction patterns involving super clusters in the unknown group.

To start, we translate the patterns depicted in Figure 3 into a matrix of transaction volumes and time intervals for all cluster pairs  $(c_y, \hat{c}_x) \in \vec{W}^K$  with  $\hat{c}_x \in \hat{C}^K$ . We then create a matrix of transaction volumes and time intervals for all cluster pairs  $(c_y, \hat{c}_x) \in \vec{W}^U$  with  $\hat{c}_x \in \hat{C}^U$  where:

$$\vec{W}^U \subset W = \{w_{yx}(c_y, \hat{c}_x) \in W \mid c_y \in C \setminus \hat{C}^U, \hat{c}_x \in \hat{C}^U, n_{yx} \geq 100\} \quad (16)$$

<sup>20</sup>For deeper insight into the detail algorithm please see Barton et al. (1992).

which means that a pair  $(c_y, \hat{c}_x) \in \vec{W}^U$  is considered only if there has been at least 100 transactions from  $c_y$  to  $\hat{c}_x$ . For each pair  $(c_y, \hat{c}_x) \in \vec{W}^U$  we calculate the median transaction volume and the median of time interval (in minutes) between inflow transactions.

Let  $\vec{W}^{KX} \subset \vec{W}^K$ ,  $\vec{W}^{KG} \subset \vec{W}^K$ ,  $\vec{W}^{KB} \subset \vec{W}^K$  be the subsets of inflow transactions towards exchanges, gambling, and black markets in the known group, respectively. Similarly, let  $\vec{W}^{UX} \subset \vec{W}^U$ ,  $\vec{W}^{UG} \subset \vec{W}^U$ ,  $\vec{W}^{UB} \subset \vec{W}^U$  be the subsets of inflow transactions towards exchanges, gambling, and black markets in the unknown group, respectively.

Then, the 2-D correlations of the inflow transaction patterns between super clusters (in the known and unknown group) and clusters outside the groups are defined as follows:  $corr2D(\vec{W}^{KX}, \vec{W}^{UX})$  is the 2-D correlation between the inflow transaction patterns for the exchanges in the known and unknown groups;  $corr2D(\vec{W}^{KG}, \vec{W}^{UG})$  is the 2-D correlation between the inflow transaction patterns for the gamblers in the known and unknown groups;  $corr2D(\vec{W}^{KB}, \vec{W}^{UB})$  is the 2-D correlation between the inflow transaction patterns for the black markets in the known and unknown groups.

The correlation matrix in Table 5 shows that the classification of the super clusters according to the PUG analysis is consistent with the results of the TP analysis because the correlations along the main diagonal are greater than the values off-diagonal. Namely,

$$\begin{aligned} corr2D(\vec{W}^{KX}, \vec{W}^{UX}) &> corr2D(\vec{W}^{KX}, \vec{W}^{UG}), \\ &> corr2D(\vec{W}^{KX}, \vec{W}^{UB}) \end{aligned}$$

and

$$\begin{aligned} corr2D(\vec{W}^{KG}, \vec{W}^{UG}) &> corr2D(\vec{W}^{KG}, \vec{W}^{UX}), \\ &> corr2D(\vec{W}^{KG}, \vec{W}^{UB}) \end{aligned}$$



and

$$\begin{aligned} corr2D(\vec{W}^{KB}, \vec{W}^{UB}) &> corr2D(\vec{W}^{KB}, \vec{W}^{UX}), \\ &> corr2D(\vec{W}^{KB}, \vec{W}^{UG}) \end{aligned}$$

| Correlation Matrix - Inflow Transaction Volume/Interval Matrix |                |                |                |
|--|----------------|----------------|----------------|
|  | $\vec{W}^{UX}$ | $\vec{W}^{UG}$ | $\vec{W}^{UB}$ |
| $\vec{W}^{KX}$   | <b>0.8183</b>  | 0.2240         | 0.5090         |
| $\vec{W}^{KG}$   | -0.0412        | <b>0.8943</b>  | -0.0100        |
| $\vec{W}^{KB}$   | 0.6615         | 0.0389         | <b>0.6665</b>  |

Table 5: Correlation of category transaction(inflow) pattern between the known group and unknown group.

Finally, by following a reverse approach than the one adopted to build the 2-D correlation matrix for the inflow transaction patterns, we calculate also the 2-D correlation between pairs of outflow transactions involving clusters in the known and unknown group. Table 6 shows that also in this case the classification of the super clusters according to the PUG analysis is consistent with the results of the TP analysis because the correlations along the main diagonal are greater than the values off-diagonal.

| Correlation Matrix - Outflow Transaction Volume/Interval Matrix |                          |                          |                          |
|---|--------------------------|--------------------------|--------------------------|
|   | $\overleftarrow{W}^{UX}$ | $\overleftarrow{W}^{UG}$ | $\overleftarrow{W}^{UB}$ |
| $\overleftarrow{W}^{UX}$  | <b>0.4984</b>            | -0.0864                  | 0.3599                   |
| $\overleftarrow{W}^{UG}$  | 0.0378                   | <b>0.5933</b>            | -0.0705                  |
| $\overleftarrow{W}^{UB}$  | 0.4260                   | -0.0632                  | <b>0.4509</b>            |

Table 6: Correlation of category transaction(outflow) pattern between the known group and unknown group.

## 7 The Bitcoin Network

From the PUG analysis, we are able to classify some unknown super clusters into specific business categories. To illustrate the result, Figure 9 plots the payment network between the super clusters in  $\hat{C}$  and their counterparts. For the sake of visualisation purpose, two thresholds are set for plotting: First, we only plot for transactions (edges) with a volume larger than 1,000 BTC; second, the degree of the nodes must be larger than 2.

Figure 10 is a matrix of transactions between those super clusters in  $\hat{C}$  ascribed to the major business categories (exchange, mining pool, online gambling, black market and composite). The  $y$ -axis depicts the sending clusters (grouped by business category) and the  $x$ -axis depicts the receiving clusters (also grouped by business category). There is no transaction volume limit for plotting this matrix; a dot is plotted as long as a  $y$ -axis super cluster has ever sent (even once) bitcoins to an  $x$ -axis super cluster, no matter what the transaction volume is. All the dots are colored according to the category to which the source belong to. For example, all the transactions sent from exchanges are signified by red dots.

We observe that mining pools typically only send coins to other categories, and do not receive any. We also observe that black markets tend to interact most with exchanges and composite services. A more comprehensive analysis of the results shown in Figure 10 is offered by the inflow dependency matrix in Table 7. Table 7 (A) lists the bilateral transaction volume between all the pairs of business categories. The number in the cell  $(i, j)$  is the amount category  $i$  sent to  $j$ . For example,  $\text{cell}(6,1) = 6,003,342.66$  tells us the category traders sent around six million bitcoins to exchanges. Table 7 (B) calculates, for a given category, the percentage of bitcoins received from other categories. For example,  $\text{cell}(6,1) = 26.72\%$  shows us that the 26.72% of the total inflow for the category exchanges comes from traders.

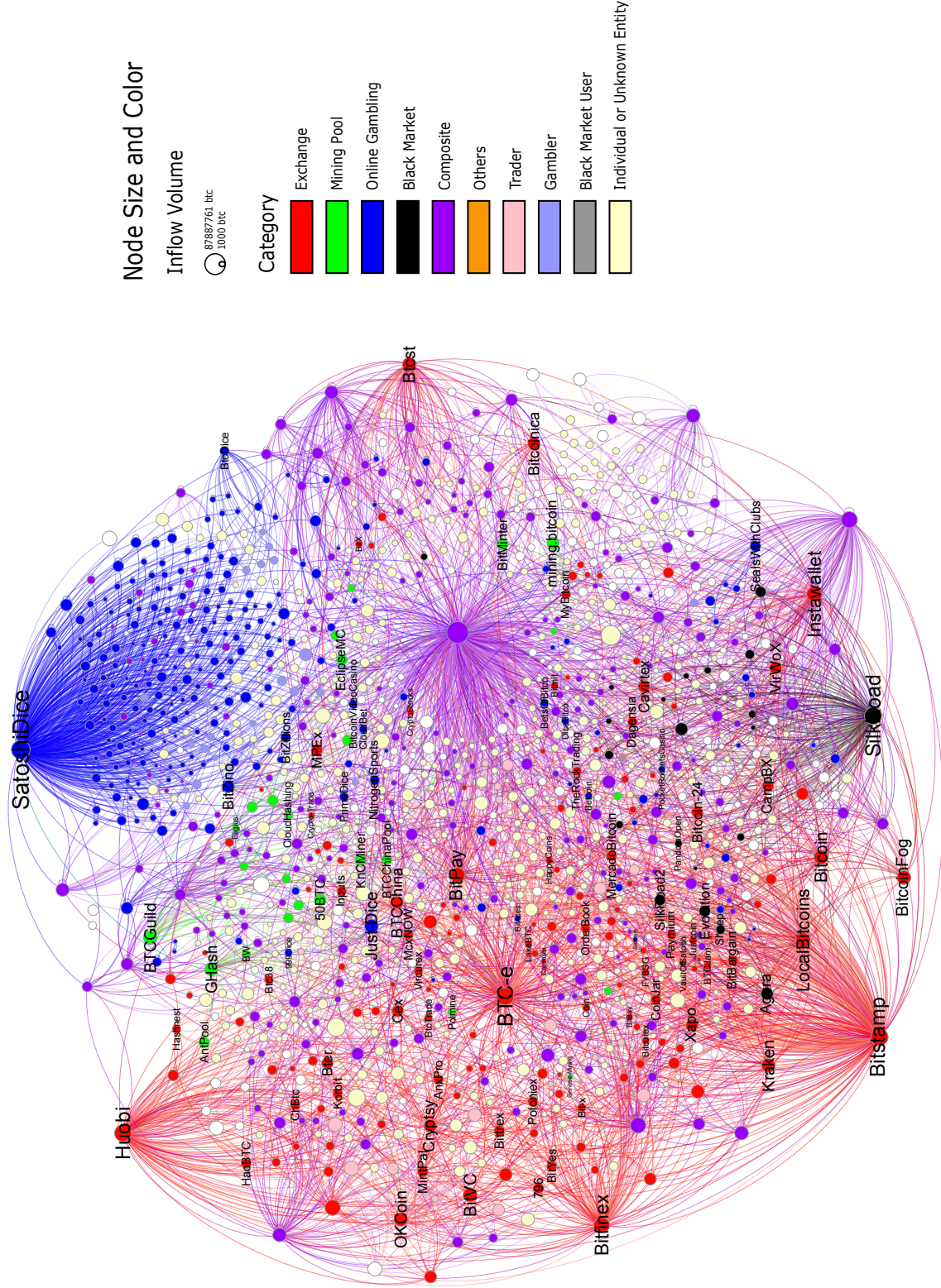


Figure 9: Payment network between super clusters in  $\hat{C}$  and their counterparts. An edge is traced as long as one party of the transaction belongs to our sample group. All clusters are colored according to the categories we explored. For clearer visualization, two thresholds are set here: First, transaction volume between each pair of nodes must be larger than 1,000 BTC. Second, the nodes' degree must be larger than 2.

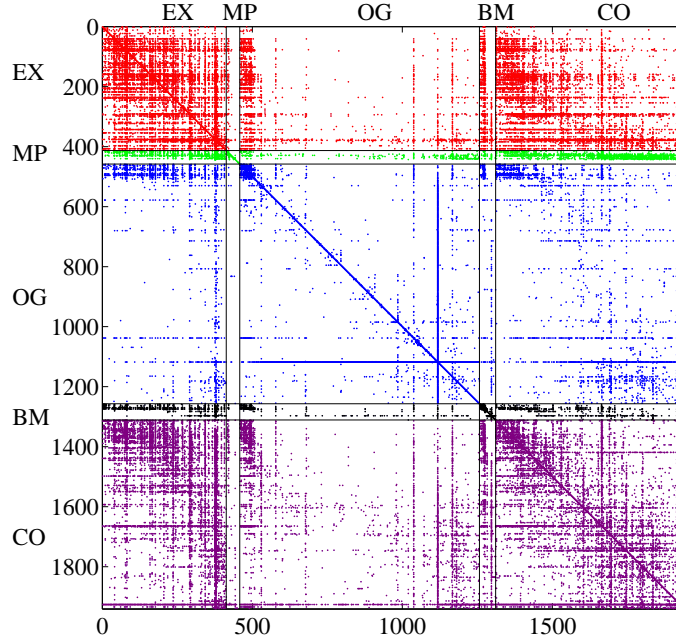


Figure 10: Transaction matrix between super clusters in different business groups: “EX” for exchange, “MP” for mining pool, “OG” for online gambling, “BM” for black market, “CO” for composite. The  $y$ -axis depicts the sending clusters (grouped by business category) and the  $x$ -axis depicts the receiving clusters (also grouped by business category). A dot is plotted if during the period January 2009 – May 2015 a  $y$ -axis cluster has ever sent bitcoins to an  $x$ -axis cluster, no matter what the transaction volume is. All the dots are colored according to the category of source clusters, which is in line with the color rule used in Figure 9.

## 7.1 Evolution of the Bitcoin Economy

Using the above analysis, we can measure the relative prevalence of each general business category (*i.e.* mining pool, exchange, online gambling, and black market) in our sample, and track their evolution over the study period and visualize shifts in their relative centrality. Figure 11 shows the income inflow from January 2009 through May 2015.

| Inflow Dependency Matrix                     |             |             |             |                  |             |            |            |           |  |
|--|-------------|-------------|-------------|------------------|-------------|------------|------------|-----------|--|
| A. Inflow Transaction Matrix (in volume)     |             |             |             |                  |             |            |            |           |  |
|  | Exchange    | Mining Pool | Gambling    | Black Market(BM) | Composite   | Trader     | Gambler    | BM User   |  |
| Exchange                                     | 12723006.59 | 5496.23     | 103072.23   | 377087.03        | 3026163.73  | 7311315.30 | 658.27     | 320.85    |  |
| Mining Pool                                  | 207146.47   | 152704.12   | 18717.09    | 1050.73          | 371430.37   | 1974.96    | 120.87     | 93.94     |  |
| Gambling                                     | 66925.21    | 4440.94     | 20670014.00 | 14443.62         | 668594.19   | 5632.81    | 1493367.85 | 322.65    |  |
| Black Market(BM)                             | 384240.90   | 1573.13     | 30706.69    | 605954.69        | 534554.26   | 458.29     | 0.00       | 527171.73 |  |
| Composite                                    | 3079138.63  | 19936.42    | 900872.93   | 858455.30        | 54573371.03 | 519874.09  | 145867.51  | 57825.19  |  |
| Trader                                       | 6003342.66  | 2.61        | 1881.18     | 279.75           | 534557.53   | 0.00       | 0.00       | 0.00      |  |
| Gambler                                      | 686.01      | 0.00        | 1296775.33  | 1093.50          | 210665.85   | 0.00       | 0.00       | 0.00      |  |
| BM User                                      | 1492.33     | 0.00        | 2220.96     | 319168.90        | 76486.66    | 0.00       | 0.00       | 0.00      |  |
| B. Inflow Transaction Matrix (in percentage) |             |             |             |                  |             |            |            |           |  |
|  | Exchange    | Mining Pool | Gambling    | Black Market(BM) | Composite   | Trader     | Gambler    | BM User   |  |
| Exchange                                     | 56.63       | 2.98        | 0.44        | 17.31            | 5.04        | 93.26      | 0.04       | 0.05      |  |
| Mining Pool                                  | 0.92        | 82.92       | 0.08        | 0.05             | 0.61        | 0.03       | 0.01       | 0.02      |  |
| Gambling                                     | 0.29        | 2.41        | 89.77       | 0.66             | 1.11        | 0.07       | 91.05      | 0.06      |  |
| Black Market(BM)                             | 1.71        | 0.85        | 0.13        | 27.82            | 0.89        | 0.01       | 0.00       | 90.00     |  |
| Composite                                    | 13.70       | 10.82       | 3.91        | 39.42            | 90.96       | 6.63       | 8.89       | 9.87      |  |
| Trader                                       | 26.72       | 0.00        | 0.01        | 0.01             | 0.89        | 0.00       | 0.00       | 0.00      |  |
| Gambler                                      | 0.00        | 0.00        | 5.63        | 0.05             | 0.35        | 0.00       | 0.00       | 0.00      |  |
| BM User                                      | 0.01        | 0.00        | 0.01        | 14.65            | 0.12        | 0.00       | 0.00       | 0.00      |  |

Table 7: This table shows the transaction relationship between categories. The transaction flow in subtable(A) is from row to column. For example, cell(2,1) means mining pools send 207,146 BTC to exchanges. In subtable(B), the percentage is calculated column-wise, such that the figure reflects the inflow ratio for each category. For example, cell(2,1) = 0.92 tells us 0.92% of total income for exchanges is from mining pools.

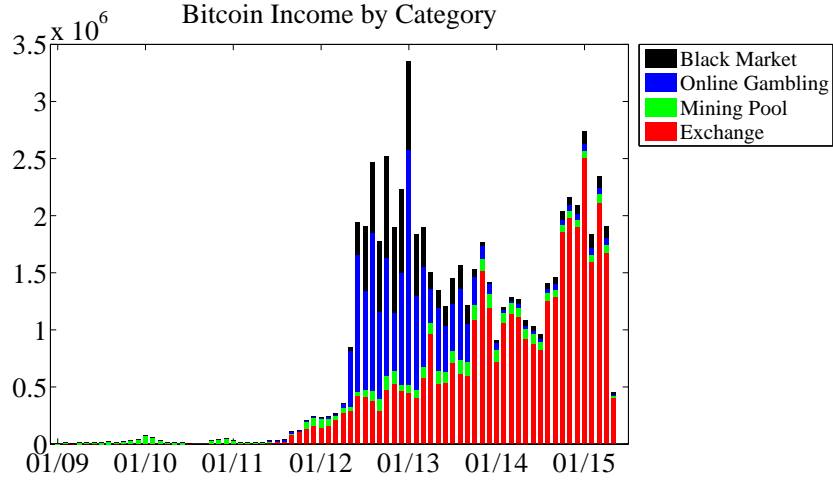


Figure 11: Stacked plot of the inflow income amount for each business category in our sample (i.e., sum of the bitcoin inflows across all the super clusters in  $\hat{C}$  belonging to the same major business category) over the Bitcoin network, monthly from January 2009 through May 2015.

We identify three distinguishable regimes that have occurred in the Bitcoin economy since its inception. The first period runs from approximately January 2009 through March 2012. This “proof-of-concept” period is characterized largely by test transactions among a small number of users, and with very little meaningful commercial activity.<sup>21</sup> Our analysis shows that this initial period is dominated almost entirely by mining, which is what we’d expect from a system still devoid of material economic activity.

Next, from approximately April 2012 through October 2013 a second period consisting of “early adopters” appears. This period is characterized by a sudden influx of gambling services and “darknet” black markets; due to the overwhelming prevalence of these arguably nefarious categories, another name for this phase could be the period of “sin.” These types of businesses initially responded to the unique features of Bitcoin such as its relative anonymity (pseudonymity), lack of regulatory and legal oversight, borderless transactions, and low transaction costs absent from taxation. This new form of secure digital cash was ideal for the purchase and sale of illicit drugs, stolen items, and other contraband that could not be easily traded elsewhere online, or for gambling from a lo-

<sup>21</sup>One notable exception is on the 22nd May 2010 in a purchase made by Laszlo Hanyecz, a software developer who paid a fellow BitcoinTalk online forum user 10,000 BTC for two Papa John’s pizzas. At today’s prices that is the equivalent of \$2.25 million per pizza!



cation where such a practice would be prohibited. Often, users of these “sin” sites would mask their internet traffic via services such as a virtual private network (VPN) or via the TOR network, encouraging usage growth where the probability of being caught would be minimal (Dingledine, 2004). In fact, our data show that in January of 2013, gambling and black markets together accounted for fully 51% of all transactional inflows on the Bitcoin blockchain (in our sample). Figure 12 shows the relative percentage of inflow transactions for each business category from January 2009 through May 2015.

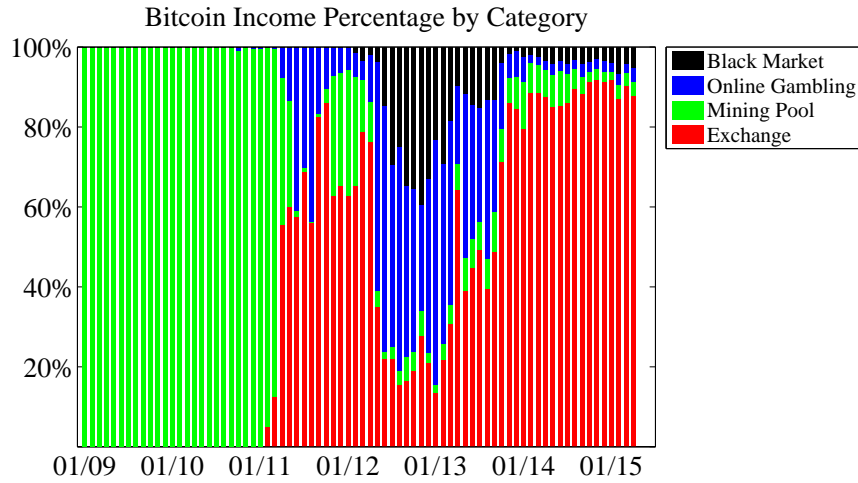


Figure 12: Stacked plot of the *relative* income for each business categories as a percentage of total income inflows, monthly from January 2009 through May 2015. Mining dominates initially, then “sin” categories (gambling in blue and black markets in black) rise, but recede over time in favor of exchanges.

The largest black market at the time was the Silk Road (see Figure 9). That service was famously raided and shut down by the FBI in October of 2013, which could help explain the sudden drop in black market activity that brought this period to a close, although this event cannot satisfactorily explain the concurrent drop in gambling activity. The drop in gambling as a percentage of overall bitcoin transactions may have been due instead to the increase in value of one bitcoin, from a few tens of dollars to a few hundreds of dollars. If a gambler tends to bet one hundred dollars per day, what used to translate into dozens of bitcoins steadily became fractions of one bitcoin. Indeed, even though the relative amount of gambling has declined, the absolute amount wagered in dollar terms has risen modestly. Over one hundred active gambling services currently exist that use

bitcoin.

It is also worth noting that while the overall amount of business being transacted on “sin” entities has fallen quite significantly, the actual number of black market sites available on the Bitcoin economy has grown, with at least four reboots of the Silk Road, and no less than fifty other (now defunct) marketplaces established since January 2014. There are still a dozen or more such marketplaces active at the time this paper is written (Branwen, 2015).

By November 2013, the amount of inflows attributable to “sin” entities had shrunk significantly to just 3% or less of total transactions. This third period, which we are still currently experiencing, is characterized by a maturation of the Bitcoin economy away from “sin” enterprise and diversifying into legitimate payments, commerce and services. This claim is moreover supported by the ascendancy of the centrality of exchanges in the Bitcoin network. Figure 13 takes the sum of the monthly betweenness centralities of the super clusters in each business category and it ranks them from January 2012 through May 2015.<sup>22</sup> Each cell is colored according to the category we have identified. Since January 2014 we see red cells outnumber all the others in each column, which tells us that exchanges are the center of transaction activity.

When a licit merchant or service provider enters the Bitcoin economy and accepts bitcoin as payment, we expect that they will cash out on a steady basis in order to cover business costs and to reduce exposure to bitcoin’s price volatility; in doing so they require the regular use of exchanges. At the same time, investors and other users who see bitcoin as a financial asset would increasingly require exchanges. It is also around this time that external venture capital investment grew in support of Bitcoin-related start-ups and infrastructure, further legitimizing it. According to Tasca (2015) startups in the Bitcoin space raised almost \$1 billion in three years (Q1 2012 – Q1 2015). In 2012, around \$2 million of VC money made its way to Bitcoin start-ups. In 2013 that number had

---

<sup>22</sup>For consistency we also check other centrality measures like weighted degree and closeness, but the result does not change. We refer the reader to Hanneman and Riddle (2014) for more details on network centrality measures.



grown to \$95 million, followed by \$361.5 million in 2014 and more than half a billion dollars in 2015. Mining pools have stayed out of the spotlight in terms of our analysis of inflows and outflows. This should not understate the value of miners and their role in the Bitcoin economy. Firstly, we would not expect mining pools to receive much in the way of income as those joining pools will only extract bitcoins and not send any to the pool. The pool generally earns income by taking a nominal percentage (1-2% or less) of the block reward, and/or by taking in the transaction fees associated with a found block. In terms of outflows, despite the amount of miners active on the network, the rate of unit formation for new bitcoins remains fixed at one block every ten minutes. At the beginning of our study period, the block reward was 50 BTC per block, from March 1, 2009 until November 28, 2012, so on any given day miners collectively produced just 7,200 BTC, a small fraction of total daily transaction volume. After November 28, 2012 and until approximately July 9th 2016, the block reward was reduced to 25 BTC, so that only 3,600 BTC were produced by miners daily, on average. After July of 2016, the block reward is again to be reduced by half to 12.5 BTC, or 1,800 BTC to be produced per day through mining. Therefore, even if all participants of mining pools cashed out daily, their contribution to the overall network of payments will always be very small, and in fact decrease over time. At the same time, the mining system is the *de facto* “central bank” of the Bitcoin economy, expanding the money supply and validating every transaction. Without a robust and “honest” segment of miners, the fidelity of all other payments in the network would be suspect. In fact, a weak network of miners would leave the Bitcoin economy prone to a so-called 51% attack, where a bad actor could begin to censor transactions by controlling a majority of power that validates transactions. Even though the relative centrality of miners is very small compared to the other business categories, the value they confer on to the network may instead be manifest via the price of bitcoin and miners’ profitability as described by Hayes (2016).

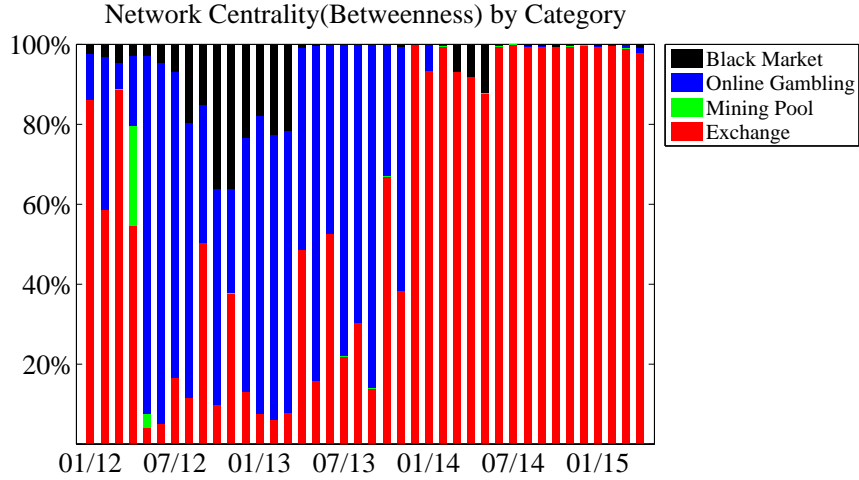


Figure 13: Monthly evolution, from January 2012 through May 2015, of the sum of the (betweenness) centrality measures across all the super clusters in  $\hat{C}$  belonging to the same major business category in the Bitcoin network.

## 8 Conclusions

As the Bitcoin economy grows in size and scope, it becomes important to understand the key components and players in that system. However, this task has largely proven cumbersome since many tens of millions of individual addresses exist, which are not obviously linked to any specific individual or business entity, simply representing nondescript public keys in a public-private key pair.

In this paper we start by analysing a database composed of millions of individual bitcoin addresses that we distill down to 2,850 super clusters, each comprised of more than 100 addresses and having received at least 1,000 BTC from January, 2009 to May, 2015. A super cluster can be thought of as an approximation of a business entity in that it describes a number of individual addresses that are owned or controlled collectively by the same beneficial owner for some particular economic purpose. These important clusters are, for the most part, initially unknown and uncategorized. However, we are able to ascribe most of them to one of four specific business categories – mining pools, exchanges, online gambling, and black markets – by mapping and analyzing the network of payments among those and a smaller, known set of clusters. In particular, we achieve

this using a Pure User Group (PUG) analysis that examines inflows and outflows to and from each of these clusters as well as Transaction Pattern (TP) analysis to confirm those findings. Our method of de-anonymizing otherwise pseudonymous clusters allows us to not only visualize the Bitcoin network of payments but also to extract stylized facts that describe its internal economy.

We find that there are, in fact, distinct patterns of transaction flows for each business category. For example, flows between traders and exchanges average just around 20 BTC, and traders buy or sell on average every 11 days. Meanwhile, gamblers wager just 0.5 BTC on average, but re-bet often within the same day. There seems to be a strong preference to do business within the bitcoin economy in round lot amounts (*e.g.*, 0.1, 0.2, 0.5, 1.0 BTC, etc.), whether it is traders exchanging for fiat money, gamblers placing bets, or black market goods being bought and sold.

In terms of transaction interval, there is an observable one-day effect for each business category. For instance, one stylized fact that emerges is that many small miners receive mining rewards and may subsequently sell those rewards on exchanges daily. This is interesting, as it could suggest most miners are operating for-profit and are not doing so in order to accumulate and hoard bitcoins. Whether or not this effect has any bearing on the price of bitcoin is open to further study. Transaction flows from miners in our sample, however are relatively small compared to the rest of the Bitcoin economy since miners in aggregate are only able to produce no more than either 3,600 or 7,200 BTC per day, on average (with a block reward of either 25 or 50 BTC), as the Bitcoin protocol enforces a controlled rate of new unit formation at one block every ten minutes.

We then trace the evolution of the prevalence of each business category over time, and identify three distinguishable regimes that have existed over the lifespan of the Bitcoin economy. First, a proof-of-concept phase made up largely of small test transactions and dominated by mining, with little substantive economic activity. Next, a period of rapid growth occurs as early adopters consisting mainly of “sin” enterprises (*i.e.* gambling services and black markets) who flock to the unique attributes that a cryptocurrency

typifies. In the third phase, “sin” enterprises are supplanted by legitimate merchants and a proliferation of exchange activity as those businesses convert digital currency into fiat to cover costs and avoid price volatility. We can thus refer to the first regime as the “proof of concept” or “mining-dominated” phase, the second as the “sin” or “gambling/black market-dominated” phase, and the third as the “maturation” or “exchange-dominated” phase.

The result of this work is to show that the Bitcoin economy, rather than being a fleeting and frivolous pursuit, has grown and matured over the few years that it has been operational, with distinct patterns of behavior among its most influential entities and participants. As the Bitcoin economy continues to expand and evolve, the type of de-anonymization and analysis employed in this paper can be used to ascribe unknown entrants to perhaps new, distinct business categories, as well as further update and refine the network of payments. Moreover, the methodology can in general be extended and applied to other cryptocurrency networks, for example to the Ethereum or Litecoin blockchains.

To conclude, the outcome of our study provides a quantitative assessment of the systematically important categories within the Bitcoin economy and their network of payment relationships. This information and analysis can be relevant to a broad audience of interested parties, including financial professionals, data scientists, and social scientists; as well as to policymakers, regulators and risk management practitioners. Finally, our results suggest that some recent concerns regarding the use of bitcoin for illegal transactions at the present time might be overstated, and that whatever such transactions may exist could further diminish as the Bitcoin economy continues to mature.

# Appendices

## A Figures

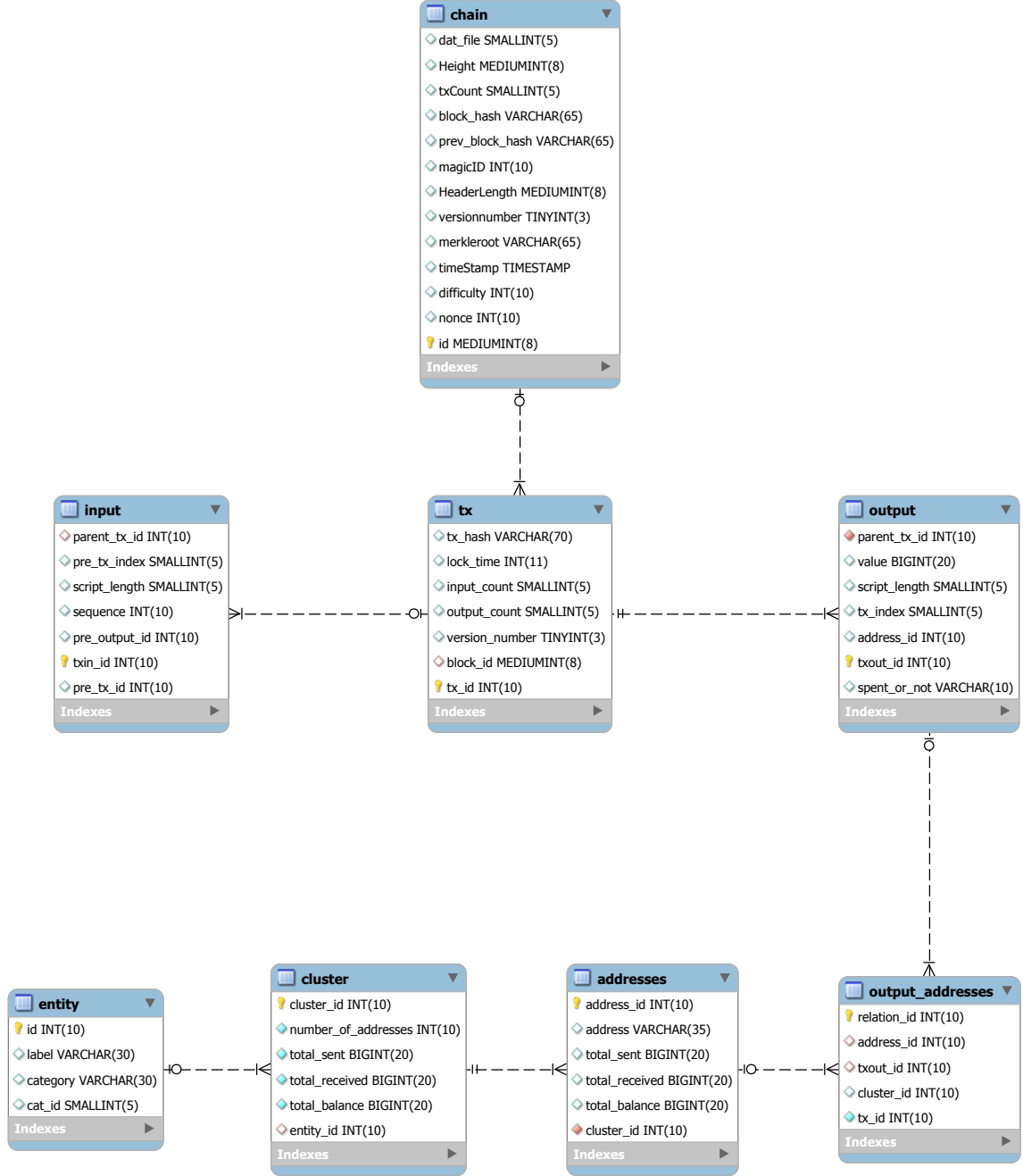


Figure 14: Structure of the MySQL database created from the data described in Table 1. We parse the transaction data from the Bitcoin Core, and then we populate them into a MySQL database. Our database carries all the data in 'blk00000.dat' files, from the 1st of Jan 2009 to the 7th of May 2015. The skeleton of the database is composed of 5 tables: chain, tx, input, output and addresses. Except for linkage between output and addresses, all the other linkage are 1:n relationship. Output and address have m:n ( $m > 1, n > 1$ ) relationship, as in case of multi-sig transactions, one output could contain several addresses, and at the same time, one address could also be used to receive bitcoins from time to time.

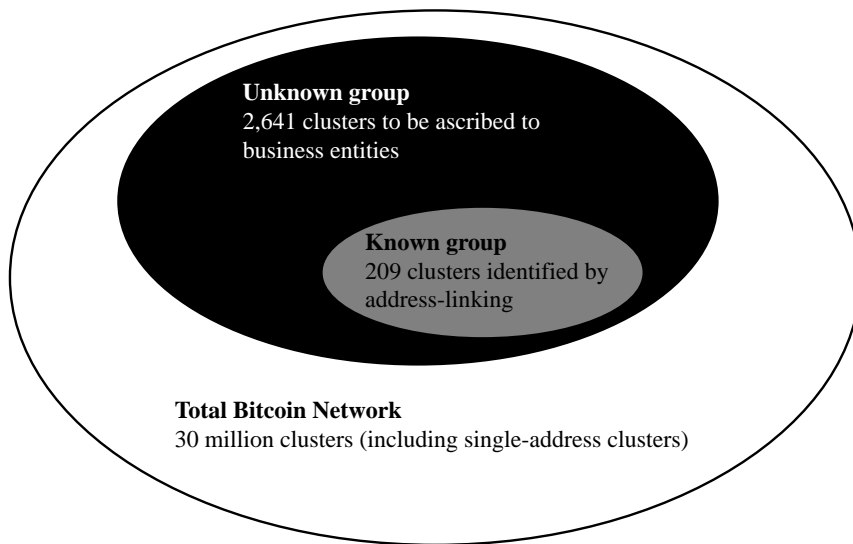


Figure 15: The total number of clusters in the Bitcoin network is about 30 million. Our research focuses on 2,850 large clusters that include at least 100 addresses and that also have received at least 1,000 bitcoins from the 3rd of January 2009 to the 8th of May 2015 (Known Group + Unknown Group).

## B Tables

| List of identified super clusters in the set $\hat{C}^K$ |                 |          |                |              |                |              |                     |                         |
|--|-----------------|----------|----------------|--------------|----------------|--------------|---------------------|-------------------------|
| Exchange   |                 | 24850005 | CoinSpot       | 45437770     | BitYes         | 9549829      | BitMillions         | 25855276 Evolution      |
| Cluster_id   | Entity_name     | 25671458 | CoinTrader     | 45936035     | Huobi          | 10352795     | Betcoins            | 28893773 BlackBank      |
| 414294   | VirWoX          | 25686021 | Poloniex       | 46523945     | CleverCoin     | 11592006     | BTCTOracle          | 32188483 CannabisRoad   |
| 725951   | Cavirtex        | 25764718 | LiteBit        | 47947669     | BtcTrade       | 12554372     | Coinroll            | 34450135 PandoraOpen    |
| 1152538  | CampBX          | 26526196 | AllCoin        | 48224085     | BitVC          | 14592351     | Just-Dice           | 39307422 MiddleEarth    |
| 1477742  | MercadoBitcoin  | 26768188 | VaultOfSatoshi | 49491730     | Coinmate       | 15935043     | BIToomBa            | 53431789 Nucleus        |
| 1538322  | BTCCChina       | 27058962 | MintPal        | 49497346     | LocalBitcoins  | 17858445     | YABTCL              | 58138309 Abraxas        |
| 1591210  | Bitcash         | 27430132 | C-Cex          | 51159703     | Bter           | 18104553     | BitZillions         |                         |
| 1640270  | BTC-e           | 27436957 | Indacoin       | 52772381     | ChBtc          | 18242583     | Ice-Dice            | Others                  |
| 1745068  | Bitstamp        | 27516236 | 1Coin          | 63125407     | BTCCChina      | 18844372     | SatoshiRoulette     | Cluster_id Entity_name  |
| 1872280  | Bitcoin         | 27883925 | Bittrex        | 64576584     | Exmo           | 19074264     | Peerbet             | 51591631 HaoBTC         |
| 2403973  | TheRockTrading  | 28195895 | Paymium        | 64714050     | HitBtc         | 20709464     | Betcoin             | 2481605 BitcoinFog      |
| 3160452  | OrderBook       | 29314526 | AllCrypt       | 67252210     | Bter           | 21368294     | AnoniBet            | 38948179 BitLauder      |
| 3169372  | BitBargain      | 29907632 | Dgex           | 74399779     | MtGox          | 22181037     | NitrogenSports      | 17685858 CryptoLocker   |
| 5128017  | LocalBitcoins   | 30131409 | CoinMotion     |              |                | 22815568     | CoinGaming          | 1400957 MPEx            |
| 5946497  | HappyCoins      | 30139877 | Bter           | Mining Pools |                | 23210454     | SatoshiBet          | 2062018 Bitcoinica      |
| 6299268  | Cryptonit       | 30804162 | CoinArch       | Cluster_id   | Entity_name    | 24545072     | 999Dice             | 3165186 Bitcoinica      |
| 6606601  | MtGox           | 30852641 | BTCCChina      | 177451       | Eligius        | 24857474     | BitcoinVideoCasino  | 32965397 UpDown         |
| 6960785  | Bitfinex        | 31778769 | Coin-Swap      | 2400970      | mining.bitcoin | 26783278     | PocketRocketsCasino | 1406234 Btctest         |
| 7522909  | Bitcoin-24      | 32344865 | BitBay         | 2440660      | BitMinter      | 28382823     | BitAces             | 17144983 Purse          |
| 8058186  | Justcoin        | 32394318 | Bter           | 4886325      | EclipseMC      | 32814149     | BitStarz            | 21601241 Bylls          |
| 8764670  | FYBSG           | 33419156 | CoinCafe       | 5272039      | GHash          | 33495508     | Betcoin             | 14543862 Bitbond        |
| 11025414   | BitX            | 34085743 | BX             | 7530073      | BTCTGuild      | 37042731     | CloudBet            | 36933042 BTCTJam        |
| 11196419   | SmenarnaBitcoin | 34277949 | BtcExchange    | 8388553      | 50BTC          | 38624871     | PrimeDice           | 39317993 BitLendingClub |
| 11749226   | Cryptorush      | 35226292 | MeXBT          | 11551066     | 50BTC          | 39363482     | DiceNow             | 17815289 BTCT           |
| 12637441   | McxNOW          | 35431781 | Zyado          | 12547187     | mining.bitcoin | 41129839     | DiceBitco           | 1075785 BitPay          |
| 12797521   | Korbit          | 35636277 | QuadrigeCX     | 13455133     | KnCMiner       | 43427199     | PrimeDice           | 16248472 CoinPayments   |
| 13228368   | Vircurex        | 36674288 | MaiCoin        | 18761724     | CloudHashing   | 44125199     | SatoshiMines        | 65645195 BitPay         |
| 13539065   | Crypto-Trade    | 36837273 | HitBtc         | 21224287     | BTCTChinaPool  | 45607266     | FortuneJack         | 1582623 Bitnit          |
| 13549778   | Cryptsy         | 37013580 | Matbea         | 23855294     | Polmine        | 48934666     | SecondsTrade        | 13255854 CryptoStocks   |
| 14777694   | Coins-e         | 37776533 | Btc38          | 34581906     | Genesis-Mining | 50523669     | Betcoin             | 454407 Instawallet      |
| 14833131   | AnxPro          | 38951758 | Ceedk          | 45656162     | AntPool        | 52248120     | SatoshiDice         | 869503 MyBitcoin        |
| 15004560   | BitKonan        | 39963036 | 796            | 48150806     | mining.bitcoin | 57476416     | BitcoinVideoCasino  | 8341192 Dagensia        |
| 16030982   | OKCoin          | 40161739 | LakeBTC        | 58048160     | AntPool        | 58900551     | PrimeDice           | 14011339 CoinJar        |
| 17494455   | Huobi           | 41193900 | Bitso          | 61166475     | BW             | 64148592     | BitAces             | 14359270 Xapo           |
| 17518823   | CoinMkt         | 41323542 | SpectroCoin    |              |                | 65420930     | SwCPoker            | 14773742 Inputs         |
| 17747783   | Kraken          | 41433875 | OKCoin         | Gambling     |                | 73161189     | PrimeDice           | 31631652 BitcoinWallet  |
| 18055670   | Cex             | 41555907 | BTC-e          | Cluster_id   | Entity_name    |              |                     | 51620287 OkLink         |
| 18847146   | BtcMarkets      | 41614840 | BTC-e          | 184867       | Just-Dice      | Black Market |                     | 59825409 GoCelery       |
| 19681395   | Bitcoin         | 41923963 | Hashnest       | 1687007      | BetsOfBitco    | Cluster_id   | Entity_name         |                         |
| 20789150   | Coinomat        | 42369494 | Cryptsy        | 2254800      | SealsWithClubs | 4401158      | SilkRoad            |                         |
| 21373812   | Bleutrade       | 42879690 | C-Cex          | 3486952      | SatoshiDice    | 9563241      | Sheep               |                         |
| 21653414   | Bitfinex        | 43277175 | Bit-x          | 4169604      | BitZino        | 19517829     | PandoraOpen         |                         |
| 23421684   | Coin            | 43970673 | Bter           | 4831753      | BtcDice        | 20627442     | SilkRoad2           |                         |
| 23672561   | Masterxchange   | 43974172 | Bter           | 8339663      | BitElfin       | 22735225     | Agora               |                         |
| 24089310   | Igot            | 45333046 | Bitcurex       | 9510403      | Playt          | 22917766     | BlueSky             |                         |

Table 8: List of the super clusters in  $\hat{C}^K$ . One entity could own and control more than one cluster. The cluster IDs are generated internally from MySQL database, and each cluster has one unique cluster ID. Entities are classified according to their business objective. We focus on the biggest four categories (exchange, mining pool, gambling, black market). Entities with exposure to more than one category, such as HaoBTC (both wallet and mining pools) and categorized as “composite”.

## References

- Androulaki, E., Karame, G. O., Roeschlin, M., Scherer, T., and Capkun, S. (2013). Evaluating User Privacy in Bitcoin. In *Financial Cryptography and Data Security*, pages 34–51. Springer.
- Antonopoulos, A. M. (2014). *Mastering Bitcoin: Unlocking Digital Cryptocurrencies*. ” O’Reilly Media, Inc.”.
- Barton, F., Himmelsbach, D., Duckworth, J., and Smith, M. (1992). Two-dimensional vibration spectroscopy: correlation of mid-and near-infrared regions. *Applied Spectroscopy*, 46(3):420–429.
- Blockchain (2015a). <http://blockchain.info>. (Date last accessed: 01-June-2015).
- Blockchain (2015b). [https://blockchain.info/tags?form\\_type=0](https://blockchain.info/tags?form_type=0). (Date last accessed: 01-June-2015).
- Branwen, G. (2015). Black market risks. <http://www.gwern.net/Black-market%20survival>. (Date last accessed: 15-Jun-2016).
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). Data structures for disjoint sets. *Introduction to Algorithms*, pages 561–585.
- Dingledine, e. a. (2004). The second-generation onion router. *Naval Research Lab*.
- Doll, A., Chagani, S., Kranch, M., and Murti, V. (2014). btctracker: Finding and displaying clusters in bitcoin.
- Garcia, D., Tessone, C. J., Mavrodiev, P., and Perony, N. (2014). The Digital Traces of Bubbles: Feedback Cycles Between Socio-Economic Signals in the Bitcoin Economy. *Journal of The Royal Society Interface*, 11(99):20140623.
- Hanneman, R. and Riddle, M. (2014). Introduction to social network methods, university of california, riverside, 2005. URL: <http://faculty.ucr.edu/hanneman/nettext>.



- Hayes, A. (2016). Cryptocurrency Value Formation: An empirical study leading to a cost of production model for valuing Bitcoin. *Telematics & Informatics*.
- Kristov Atlas (2015). Weak Privacy Guarantees for SharedCoin Mixing Service. <http://www.coinjoinsudoku.com/advisory>. (Date last accessed: 01-June-2015).
- Meiklejohn, S., Pomarole, M., Jordan, G., Levchenko, K., McCoy, D., Voelker, G. M., and Savage, S. (2013). A Fistful of Bitcoins: Characterizing Payments Among Men with No Names. In *Proceedings of the 2013 conference on Internet measurement conference*, pages 127–140. ACM.
- Nakamoto, S. (2008). Bitcoin: A peer-to-peer electronic cash system. *Consulted*, 1:2012.
- Reid, F. and Harrigan, M. (2013). *An analysis of Anonymity in the Bitcoin System*. Springer.
- Ron, D. and Shamir, A. (2013). Quantitative analysis of the full bitcoin transaction graph. In *Financial Cryptography and Data Security*, pages 6–24. Springer.
- Spagnuolo, M. (2013). Bitiodine: extracting intelligence from the bitcoin network.
- Tasca, P. (2015). Digital currencies: Principles, trends, opportunities, and risks. *ECUREX Research WP (September 7, 2015)*.
- Walletexplorer (2015). <https://www.walletexplorer.com/>. (Date last accessed: 01-June-2015).