

Resurrecting Address Clustering in Bitcoin

Malte Möser
Princeton University

mmoeser@cs.princeton.edu

Arvind Narayanan
Princeton University

arvindn@cs.princeton.edu

Abstract

Blockchain analysis is essential for understanding how cryptocurrencies like Bitcoin are used in practice, and address clustering is a cornerstone of blockchain analysis. However, current techniques rely on heuristics that have not been rigorously evaluated or optimized. In this paper, we tackle several challenges of change address identification and clustering. First, we build a ground truth set of transactions with known change from the Bitcoin blockchain that can be used to validate the efficacy of individual change address detection heuristics. Equipped with this data set, we develop new techniques to predict change outputs with low false positive rates. After applying our prediction model to the Bitcoin blockchain, we analyze the resulting clustering and develop ways to detect and prevent cluster collapse. Finally, we assess the impact our enhanced clustering has on two exemplary applications.

1 Introduction

Motivation. Blockchain analysis techniques are essential for understanding how cryptocurrencies like Bitcoin are used in practice. A major challenge in analysing blockchains is grouping transactions belonging to the same user. Because users can create an unlimited amount of addresses, each of which can receive and send coins, their activity may be split among a multitude of such addresses.¹ Techniques to group activity of individual users are commonly referred to as *address clustering heuristics*, as they focus on identifying the addresses under an individual user’s control using heuristic assumptions about how their transactions are created. As the term *heuristic* suggests, address clustering today is more intuitive than rigorous; our overarching goal in this paper is to elevate it to a science.

There are at least four applications for which accurate address clustering is important. First, a law enforcement agency

¹Cryptocurrencies use digital signature schemes to authorize transactions. Public keys are used as account identifiers (called *addresses*), and the signatures created with the corresponding private key authorize transactions. Users can create and use an unlimited amount of key pairs.

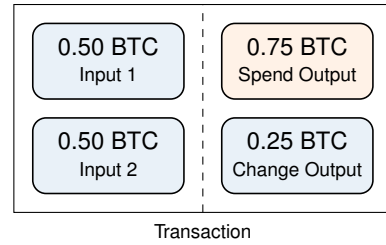


Figure 1: Schema of a typical Bitcoin transaction with two inputs and two outputs: the spend output is the intended payment, the change output returns the surplus coins to the sender. Each input and output is associated with an address.

may be interested in evaluating the transactions of a specific entity (e.g., a specific exchange, trader or gambling service). They may supplement their own investigation of the entity’s behavior with a set of reliable heuristics to identify relevant transactions. Second, and conversely, the ability to accurately determine a user’s transactions directly impacts their privacy. This tension between law enforcement needs and everyday users’ privacy is inherent to cryptocurrencies due to their transparency and pseudonymity. Advocates from one side push for greater privacy and from the other side for stronger regulation. To better understand this tug-of-war, it is important to quantify how reliable change address heuristics are in practice. Third, accurate grouping of transaction activity is important for aggregate analyses such as studying economic activity over time. This usually requires a full clustering of all addresses on the blockchain. Finally, the problem of address clustering itself may be interesting for researchers outside of cryptocurrencies. For example, it may pose as an application domain for machine learning models and could be used as a benchmarking application.²

Goals. The current state of address clustering techniques available to researchers is sub-optimal in multiple ways. The

²We plan to make our dataset public to facilitate this.

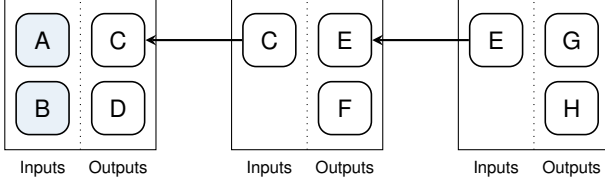


Figure 2: Clustering with only the multi-input heuristic misses addresses C and E that are not co-spent with other addresses.

most common heuristic, *multi-input*, groups addresses that are jointly used in inputs of a transaction (cf. Figure 1) [37, 38]. This heuristic is easy to apply, relatively effective in practice, [20] and widely used. However, it misses addresses that are never co-spent with other addresses (cf. Figure 2).

Many of these addresses can be identified using *change address* heuristics: as coins in Bitcoin cannot be spent partially, transactions need to return the surplus value back to the user who created the transaction. Identifying the change output thus allows grouping the associated address with the inputs’ addresses. While the effectiveness of change address detection has been demonstrated empirically and through simulation (e.g., [2, 30]), it remains difficult to assess how well it works in practice. As a result, clustering techniques are applied inconsistently across studies: many forgo change address clustering (e.g., [23, 24, 27, 39]), whereas some simply apply a single change heuristic (e.g., [11, 35]).

A major issue is the lack of ground truth data available to researchers. In the context of change output detection, ground truth consists of a set of transactions for which the change output is known. Such a dataset allows to assess the accuracy of individual heuristics aiming to identify the change output of a transaction. But because the Bitcoin blockchain is used for a variety of different use cases and by a diverse group of users, which may both change significantly over time, ground truth needs to reflect this diversity in order to allow for a reliable assessment. Such data is hard to collect, and, even if available, is unlikely to be made public, e.g., due to privacy concerns. We are only aware of one approach exploiting weaknesses in a specific type of lightweight client [34], which allowed to extract the addresses of 37585 wallets to assess four different clustering heuristics. Blockchain intelligence companies might have access to manually curated and refined data sets and clusterings, but their techniques and data aren’t generally available to researchers (or only shared in limited form, e.g., [19, 44]). As a result, analyses of clustering heuristics often fall short of quantifying their accuracy and instead resort to analyzing the resulting clusterings (e.g., [10, 45]).

Considering this state of affairs, our first goal in this paper is to address the lack of data and assessment methods. We build a ground truth set of transactions with known change from the Bitcoin blockchain that can be used to validate the efficacy of change address heuristics. Equipped with this data

set, our second goal is to develop new techniques that allow an analyst to combine different heuristics to predict change outputs with a high true positive and low false positive rate. Third, we aim to assess the quality of the enhanced clusterings that can be created using such techniques and develop ways to detect and prevent cluster collapse, which occurs when the clusters of two or more distinct entities are incorrectly merged together. Finally, we are interested in evaluating the potential impact of the enhanced clusterings on the results of typical blockchain analyses.

Contributions, methods and findings.

1. **A new ground truth method and dataset:** We put forward a procedure to select and filter transactions for which the change output has been revealed on the blockchain. Our approach exploits that future transactions of users can reveal change outputs in past transactions. We take specific care evaluating the data set with regards to potential issues, such as violations of our core assumption or existing cluster collapse. We extract a set of 30.05 million transactions with known change (carefully filtered down from 47 million candidate transactions) that can be used as ground truth for validation and prediction. Our method does not rely on interaction with intermediaries, though we use address tag data to assess the quality of our data set and methods. The data set can be continuously updated, improved and shared with other researchers. (Section 2)
2. **Evaluating existing heuristics:** We assess the true and false positive rates of various change address heuristics on our ground truth data set. We find that most change address heuristics have few false positives at low to medium true positive rates. Further, we find that Bitcoin’s protocol characteristics have become increasingly relevant due to their ability to provide a “fingerprint” for transactions. We also report the heuristics’ overall coverage (i.e., how often they return a result) for transactions with unknown change. (Section 3)
3. **Improved prediction:** We use a random forest classifier to identify change outputs and compare it against a baseline: the majority vote of individual heuristics. While machine learning has been used to classify the type of entity behind a transaction (e.g., [4, 19, 21, 23, 25, 42, 44]), to the best of our knowledge our work is the first to apply it to the problem of change identification. We find that a random forest model outperforms our baseline threshold voting mechanism, especially for low false positive rates (to prevent cluster collapse). For example, targeting a false positive rate below 0.1 % the random forest model correctly detects almost twice as many change outputs. (Section 4)

4. **Preventing cluster collapse:** We analyze the clustering that results from the predicted change outputs with regards to cluster collapse using multiple measures. We find that a naive clustering of predicted change outputs leads to cluster collapse, despite choosing a low threshold to prevent false positives. We then apply constraints to the union-find algorithm underlying our clustering, aiming to prevent cluster collapse that stems from frequent, repeated interaction between entities. This technique prevents large-scale cluster collapse while still enhancing a majority of the involved clusters. (Section 5)
5. **Assessing impact:** We assess the impact our enhanced clustering has on two exemplary applications: cash-out flows from darknet markets to exchanges and the velocity of bitcoins. We find that our enhanced clustering changes the outcomes of these longitudinal analyses by 11 % to 17 %. (Section 6)

Our process is summarized in Figure 3. We discuss our findings in Section 7 and conclude in Section 8

Limitations. Our results in this paper are limited by the availability of “real” (i.e. manually collected and validated) ground truth. As such, our analysis should be treated as a first step towards better understanding the feasibility of change address detection and clustering. However, we do not expect our high-level insights to change significantly in the light of minor corrections to our ground truth data set. We invite the research and blockchain community to evaluate our data set using their own ground truth data or analysis techniques.³ An interesting avenue for future research would be to build a privacy-preserving tool that allows to crowd-source the validation of both the ground truth data and prediction models, to which individual users could connect their wallet software.

Our extraction mechanism relies on change outputs revealed by the multi-input heuristic. This heuristic is effective in practice [20] and widely used, but vulnerable to false positives from techniques like CoinJoin and PayJoin transactions that are intentionally designed to break the heuristic (e.g., [12, 28, 29, 31]). While we take measures to detect CoinJoin transactions and pre-existing cluster collapse, some errors can remain. Furthermore, entities that more effectively prevent address reuse are less likely to be included in our data set.

We adopt the term *clustering* used in the blockchain community, however, our underlying implementation sequentially processes transactions and uses a union find algorithm that does not retain individual distance scores. In practice it may be desirable to give analysts a choice between different clustering outcomes (e.g., concerning the order in which clusters are merged or the tradeoff between false positive and true positive rates).

³We plan to make a list of transactions and output indexes available.

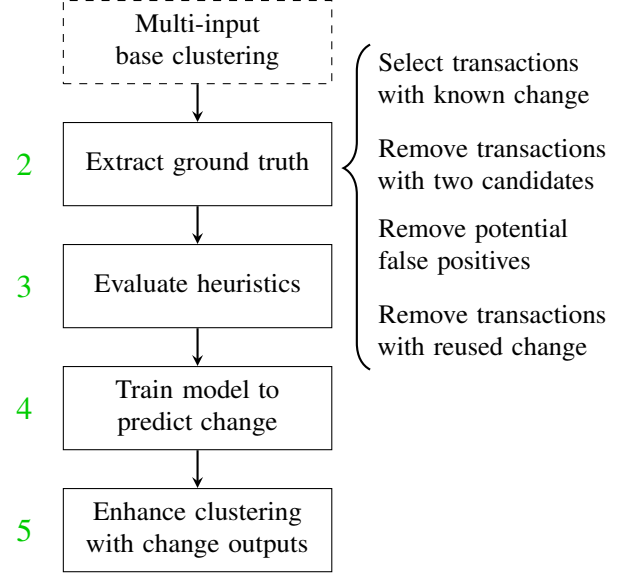


Figure 3: Our process in this paper

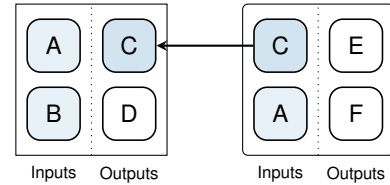


Figure 4: Address C is merged into the same cluster as addresses A and B by the multi-input heuristic, thereby revealed as the change address in the first transaction.

Finally, in this paper we work with the Bitcoin blockchain, currently the most popular cryptocurrency by market volume. Our methods are applicable to similar cryptocurrencies, potentially with a different set of heuristics, but may be less effective if transactions are more homogeneous.

2 Building a Ground Truth Data Set

Core assumption. In this paper we focus on the feasibility of detecting the change output in Bitcoin transactions with exactly two spendable outputs, by far the most common type of transaction as of June 2020 (72.6 % of all transactions, see Figure 5). Our core assumption is that one of these outputs is a payment, and the other output receives the change. We call this type of transaction a *standard* transaction, as they are created by typical end-user wallet software.⁴

For transactions with only one output there is no good indicator to directly and reliably determine whether the output

⁴There exist a separate notion of a standard transaction, namely those that pass the `isStandard` test of the Bitcoin reference implementation that checks whether a transaction uses one of a handful of default script types.

belongs to the same user. The transaction may correspond to a user sweeping the balance of their wallet, but the destination address may not be under the same user’s control (e.g., they might transfer all the bitcoins from a wallet on their personal computer to an online wallet, where private keys are managed by and in control of an exchange).

Transactions with more than two outputs are less likely to originate from an ordinary wallet. They may belong to an exchange that batches payouts to multiple users, or correspond to a restructuring of their internal hot and cold wallets. As a result, our assumption that exactly one of the outputs receives change may not hold. Large numbers of outputs could also indicate mixing services or CoinJoin transactions, where the funds of multiple users are mixed. Determining change in CoinJoin transactions requires solving a subset-sum problem (e.g., [16, 29, 31]) and is outside the scope of this paper.

Method. Our approach leverages the phenomenon that change outputs are sometimes revealed by the multi-input heuristic at a later point in time due to address reuse. Figure 4 shows an example of how such disclosure may unintentionally happen on the blockchain: a user spends coins at addresses *A* and *B*, their wallet directs the change to a new address *C*. Later, they spend the change at address *C* along with other coins at address *A*. At this point, the multi-input heuristic reveals that *A*, *B* and *C* belong to the same user, thereby revealing *C* as the change address in the first transaction.

Comparison to interactive collection. We briefly discuss the advantages and disadvantages of our approach to collecting ground truth interactively. We could download a Bitcoin wallet and send bitcoins to a number of different addresses, thereby creating a corpus of transactions for which both the spend and the change output are known. However, in order to be able to learn and generalize from our ground truth data it should capture the expected heterogeneity of implementations and use cases over Bitcoin’s entire history. An interactive collection would likely yield ground truth inferior in three dimensions: variety, scale, and the collection time frame.

Heterogeneous ground truth requires transactions from a *variety of different use cases and entities*. Compared to prior deanonymization studies that identified address clusters of specific entities by interacting with them (revealing some of the intermediary’s addresses, e.g., [30]), we are interested in the change of transactions made by those intermediaries. Purchasing an item from an online merchant can reveal one of their addresses, but we do not learn about change in any of the merchant’s transactions, only about our own transaction that pays them. The only conceivable way to learn about change in an intermediary’s transactions is to induce them to make a transaction to an address under our control. This may be possible with exchanges, where we could first deposit and then withdraw funds, but is not applicable to many other

intermediaries.⁵ Our method, instead, is not limited to a small set of intermediaries of our choosing.

Second, interactive collection would be hard to *scale* beyond a few hundred transactions, as we would need to individually engage with a variety of intermediaries and wallet implementations. This cannot easily be automated. While Nick [34] was able to collect data on a larger scale from exploiting a vulnerability in a specific type of lightweight client, his method is not transferable or generalizable to other types of wallets. Our approach, instead, yields a data set of millions of transactions.

A third issue of collecting data interactively is the *time frame*. Interactive collection, as described above, cannot be done retroactively and is therefore limited to a short, current time window. This limits its utility as the resulting data set wouldn’t capture shifting patterns over different epochs of Bitcoin’s history. Our non-interactive approach, on the other hand, is applicable to transactions from Bitcoin’s entire history.

Our method has a few important limitations. First, because we extract ground truth data non-interactively from the blockchain, we are not able to fully verify its correctness. Second, our core assumption that exactly one of the outputs is a change output may not hold in every scenario. When users transfer funds to an address under their control, determining the change and spend based on the multi-input heuristic is ambiguous. Similarly, there could be instances where none of the outputs is a change output because a user made a payment to two different entities using a perfectly matching set of inputs that does not require change to be returned. Third, as our method relies on address reuse, the resulting transaction corpus could be biased towards entities or wallet implementations that are more prone to reuse and merge addresses. It might contain fewer instances of transactions created with wallets that more effectively discourage address reuse.

2.1 Data collection and overview

We use and build upon BlockSci v0.7 [24], an open-source blockchain analysis framework that provides fast access to blockchain data upon which we can implement custom heuristics and extraction procedures. We parse the Bitcoin blockchain until the end of June 2020 (block height 637090) and create a *base clustering* using the multi-input heuristic (where we heuristically exclude CoinJoin transactions).

As of June 2020, the blockchain contains 66 million transactions with one output, 422 million with two outputs, and 56 million with three or more outputs (Figure 5). We divide the transactions into mutually exclusive categories. Transactions containing unspendable OP_RETURN outputs often signal the use of an overlay application that stores metadata in the blockchain [5]. Such transactions may have specific

⁵It could also raise ethical and legal questions, e.g., when interacting with gambling services or intermediaries in other countries.

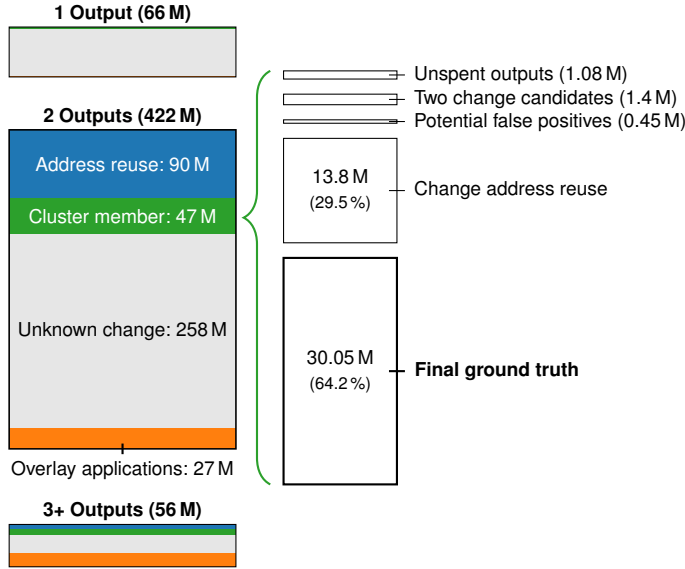


Figure 5: Distribution of different types of transactions in the Bitcoin blockchain until June 2020. Transactions with two outputs and change revealed through cluster membership form the basis of our ground truth data, which we further refine down to a ground truth data set of 30.05 million transactions.

rules for how they are constructed, potentially making change detection unreliable.

Both transactions reusing an input address as well as transactions where cluster membership (i.e. the multi-input heuristic) reveals a change output have their change output identified. Direct address reuse however makes change identification trivial and applying further change heuristics is never necessary. We thus only use transactions where the change has been revealed by multi-input clustering as the basis to construct our ground truth data set. For the remaining transactions, i.e. those with yet unknown change, we will later try to predict their change output.

2.2 Refining the candidate set of ground truth transactions

Our candidate set of ground truth transactions consists of transactions with two outputs (ignoring overlay transactions) where no input address is reused for change and where at least one output is in the same base cluster as the inputs. This yields a total of 46.79 million transactions. We further filter the transactions as follows (see Figure 5 for a visual breakdown and Appendix C.1 for an extended description):

1. We remove 1.08 million transactions with unspent outputs, as our subsequent analyses rely upon the spending transactions being known.
2. For 0.92 million transactions both outputs are in the

same base cluster as the inputs, violating our core assumption. We remove these transactions. We also find that some base clusters are more likely to produce such transactions. We thus exclude transactions from base clusters where more than 10 % of transactions exhibit such behavior. This removes an additional 0.47 million transactions in 8635 base clusters.

3. We check our base clustering for existing cluster collapse, which could produce false positives. First, we remove 0.37 million transactions belonging to the Mt.Gox supercluster. Second, we remove a possible instance of cluster collapse detected using address tags from Wallet-Explorer, further removing 0.09 million transactions.
4. We find many instances where the change address did not appear in the inputs, but had been seen before and was known to be the change at the time the transaction was created. In these instances, applying change address heuristics is unnecessary. We remove 13.81 million transactions where the change output was already known at the time the transaction was created.

2.3 Assessing the final set of ground truth transactions

Next, we assess the composition of our ground truth data set and compare it to transactions in the blockchain overall. This is useful to ensure that it contains heterogeneity with regards to scale, time frame and variety (see discussion above), as well as to spot potential biases in the ground truth that could result from our selection process.

Scale and time frame. Our ground truth of 30.05 million transactions makes up about 7.7 % of standard transactions and about 5.6 % of all transactions. Figure 6 shows that those percentages are relatively stable over time.

Variety of included clusters. Our ground truth includes transactions from 2.7 million base clusters. We consider two measures to assess their variety: the number of addresses those base clusters contain overall, as well as the number of transactions that originate from each (i.e., how ground truth transactions are distributed across base clusters).

Figure 7 shows the distribution of address counts of base clusters that are represented with at least one transaction in our ground truth. Our ground truth contains transactions from base clusters of all sizes, giving us some confidence that it can be representative of the blockchain overall. The share of base clusters from which transactions are included is higher towards clusters with larger number of addresses. This is likely a side-effect of our selection process, as base clusters with more addresses may be more likely to combine addresses and thereby reveal change.

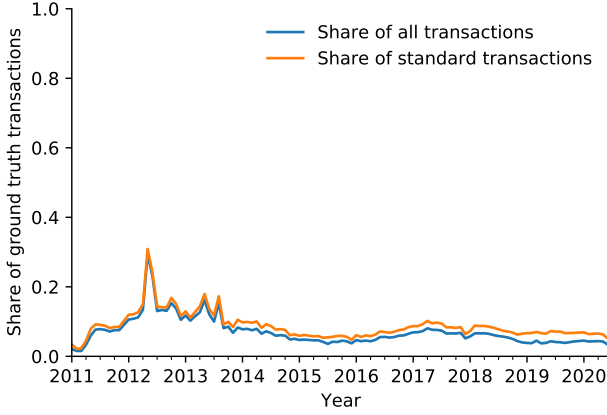


Figure 6: Share of ground truth transactions of all and standard transactions over time.

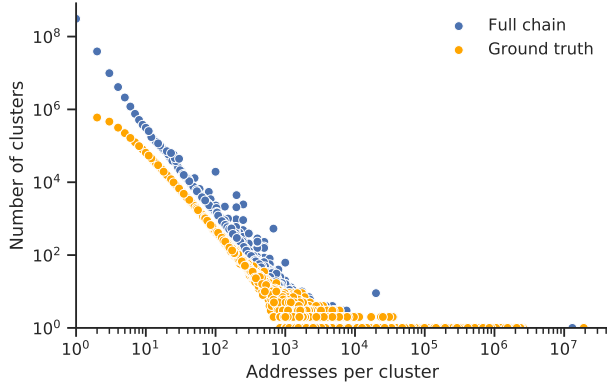


Figure 7: Number of base clusters of certain address counts contained in the ground truth data (with at least one transaction), as compared to the full blockchain.

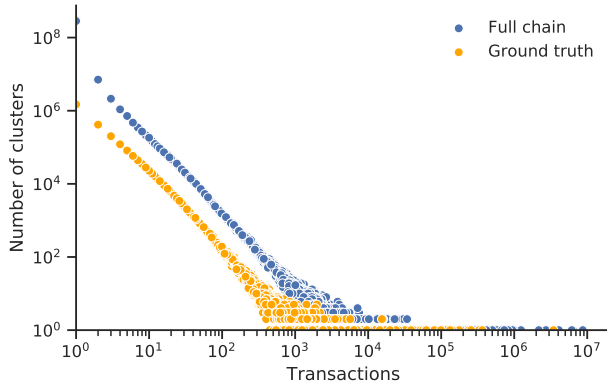


Figure 8: Number of transactions in ground truth and full blockchain, grouped by base cluster.

Table 1: Comparison of transaction characteristics between ground truth transactions and transactions with 2 outputs for which change is unknown.

Characteristic	Ground truth (%)	Remaining (%)
1 Input	38.96	78.28
2 Inputs	21.91	13.92
3+ Inputs	39.13	7.79
Version = 1	81.56	83.36
Locktime > 0	26.10	24.22
RBF	2.77	4.26
SegWit	14.89	23.36
n (in million)	30.05	258.68

Next, we inspect the distribution of transactions across base cluster that are included in the ground truth data (Figure 8). The highest number of transactions is 3.45 million, from a cluster that has 8.69 million transactions in total. It is not labeled by WalletExplorer.com. The second highest number of transactions is 366 109, again from an unlabeled cluster.

Transaction composition and use of protocol features. Table 1 compares characteristics of transactions in our ground truth data to those of standard transactions with yet unknown change, including the number of inputs as well as a number of important protocol features (an overview and description of the protocol characteristics used in this paper is available in Appendix A). Transactions in the ground truth data set notably tend to have more inputs than those in the set of transactions with unknown change. This is an artifact of our ground truth selection as it relies on transactions with more than one input to reveal change outputs. The share of transactions using SegWit serialization or allowing for fee bumping (RBF) is also higher in the set of remaining transactions.

3 Evaluating Individual Change Heuristics

3.1 Background on change address detection

The Bitcoin protocol does not explicitly distinguish between change and spend outputs. However, wallets create change outputs automatically to return surplus value when users make payments (cf. Figure 9). We briefly describe how this allows to identify change and present the heuristics that have been proposed in the literature (see Table 2 for details, limitations and references).

Spend output. For standard transactions, the user will typically be given a specific payment amount and an address to which the bitcoins should be sent (1). Payment amounts in Bitcoin are denominated in satoshi, with one bitcoin equal to

Table 2: Change heuristics proposed in the literature and used in this paper.

Heuristic	Notes and limitations	Used	Refs.
Optimal change: There should be no unnecessary inputs: if one output is smaller than any of the (2+) inputs, it is likely the change.	Only applies to transactions with 2+ inputs. We use two variants, one ignoring and one accounting for the fee.	✓	[34, 36]
Address type: The change output is likely to have the same address type as the inputs.	Wallets could use different address types to obfuscate the change output.	✓	[24, 36]
Power of ten: As purchase amounts may be rounded, and the change amount also depends on input values and the fee, it is more likely to have fewer trailing zeros.	We use six different variants, which are partially redundant.	✓	[24, 36]
Shadow address: Many clients automatically generate fresh change addresses, whereas spend addresses may be more easily reused.	Modern wallets discourage reuse of receiving addresses. We do not use the heuristic because our ground truth is filtered based on address freshness.	x	[2, 30]
Consistent fingerprint: The transaction spending a change output should share the same characteristics. We use 17 variants based on the following characteristics: <ul style="list-style-type: none"> • input/output counts and order • version • locktime • serialization format (SegWit) • replace-by-fee (RBF) • transaction fee • input coin age (zero-conf) • address and script types 	False positives are possible when a wallet implementation or the protocol change. We only consider characteristics after they are available in the protocol. Appendix A describes the characteristics we use in more detail.	✓	[8, 36]

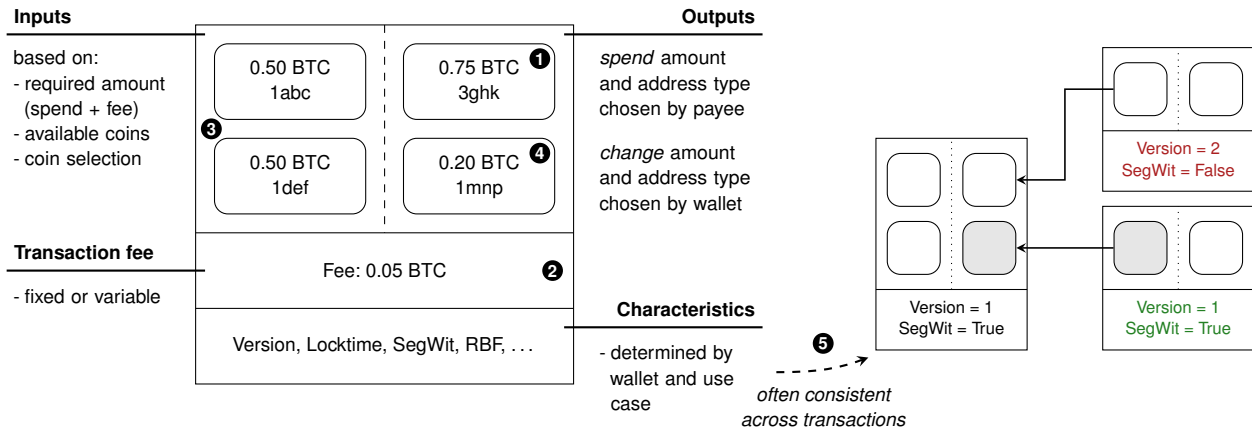


Figure 9: Schema of how transactions are created, and how consistency of a transaction's fingerprint allows to identify change.

10^8 satoshi. At current and historic exchange rates, a single satoshi is worth only a fraction of a cent. Merchants may thus round payment values to make it easier for users to enter the correct amounts in their wallet, and transfers initiated by users may use round values as well. Change can hence be distinguished from spends by potentially having fewer trailing zeros (*power of ten heuristic*).

Input selection. After the spend amount and a transaction fee have been determined, the wallet chooses a set of coins that covers the sum of both amounts (3). While this selection procedure is not standardized and can differ between wallet implementations (cf. [1, 13]), some behavior is common to many wallets, such as not including unnecessary inputs (*optimal change heuristic*).

Change output. The change output is automatically created by the wallet to return the surplus funds. Often, a fresh address is generated (i.e., one that has never received coins before). In the past it was common that payment addresses were reused, allowing to determine change based on this behavior (*shadow heuristic* / *one-time change*). We note that our ground truth data set is filtered based on address reuse (see Appendix C.1). Because this filtering effectively determines the performance of the heuristic, we decide not to use it in our subsequent analyses.

While the address type of the receiving address is determined by the payee, wallets usually use consistent address types. A single output with a different address type than the inputs might indicate a spend output, thus revealing the change (*address type heuristic*).

Consistency of transaction fingerprints. The heuristics above are based on expected behavior that should apply to many common wallet implementations. Furthermore, transactions can make use of protocol features that may allow to discriminate between different types of wallets. Such a fingerprint would allow to apply custom heuristics to subsets of transactions.

Transaction fingerprints can also be used for change identification. Assuming that users don’t change wallets very often, characteristics should be consistent across multiple transactions. When a wallet spends the change of a previous transaction, those two transactions’ fingerprints should be similar (*consistent fingerprint heuristic*, see (5) in Figure 9). We are not aware of any prior work that has evaluated this across the range of available protocol characteristics. We also note that this heuristic requires the outputs to be spent, whereas the previous heuristics are *universal* as they can be applied to all transactions, including those with unspent outputs. We selected our ground truth data such that it only contains transactions where both outputs are spent, so the consistent fingerprint heuristic can be applied to all of them.

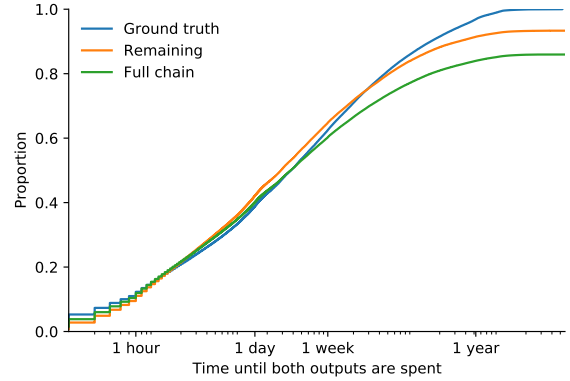


Figure 10: Time until both outputs are spent for transactions in our ground truth data, in the remaining standard transactions and the blockchain overall. Time until spent is set to infinity if not both outputs of a transaction are spent.

Figure 10 shows the distribution of time until both outputs are spent for transactions in the ground truth data set as well as remaining standard transactions with yet unknown change. Overall, the outputs of more than half of all transactions are spent in under a week, making the consistent fingerprint heuristic highly applicable to a majority of transactions, including recent ones.

3.2 Evaluating individual heuristics

We start by evaluating the universal and fingerprinting heuristics individually. We explicitly encode our constraint that only one of the outputs can be the change. Thus, applied to a transaction, the heuristic may either return an individual output if it is the only output determined to be change, or no output otherwise. Let h be a heuristic applied to transaction t which returns a set of potential change outputs, then our constrained heuristic h' returns:

$$h'(t) = \begin{cases} h(t) & |h(t)| = 1 \\ \emptyset & \text{otherwise} \end{cases}$$

This constraint is crucial to prevent cluster collapse: if a heuristic cannot determine a unique change output, we’d rather cluster none of the outputs than both, as clustering both would violate our core assumption and may lead to the merging of two clusters not belonging to the same user.

In Table 3 we report the individual heuristics’ true and false positive rate (TPR/FPR) for identifying change outputs in our ground truth. We also report the share of all standard transactions with unknown change for which the heuristic returns a unique output (denoted as “coverage”).

Most heuristics have a low FPR, with three fingerprinting heuristics being the exception: output count, input/output count and zero-confirmation spending. The power of ten

Table 3: True and false positive rates of heuristics applied to transactions in the ground truth data set.

Heuristic	Ground Truth		Remaining
	TPR	FPR	Coverage*
<i>Universal heuristics</i>			
Optimal change	0.299	0.027	0.134
• incl. fee	0.232	0.020	0.094
Address type	0.210	0.028	0.339
Power of ten			
• $n = 2$	0.489	0.012	0.405
• $n = 3$	0.444	0.006	0.335
• $n = 4$	0.400	0.005	0.277
• $n = 5$	0.326	0.006	0.191
• $n = 6$	0.229	0.005	0.115
• $n = 7$	0.115	0.001	0.053
<i>Consistent fingerprint</i>			
Output count	0.272	0.125	0.429
Input/output count	0.264	0.108	0.572
Version	0.224	0.003	0.297
Locktime	0.302	0.003	0.356
RBF	0.059	0.002	0.088
SegWit	0.155	0.019	0.224
SegWit-conform	0.023	0.001	0.029
Ordered ins/outs	0.241	0.053	0.430
Zero-conf	0.101	0.055	0.213
Absolute fee	0.134	0.029	0.335
Relative fee	0.045	0.009	0.215
Multisignature	0.137	0.000	0.151
Address type			
• P2PKH	0.214	0.014	0.287
• P2SH	0.236	0.012	0.305
• P2WPKH	0.146	0.017	0.219
• P2WSH	0.050	0.007	0.068
All address types	0.252	0.021	0.348

*Coverage denotes share of standard transactions with yet unidentified change where the heuristic returned exactly one output.

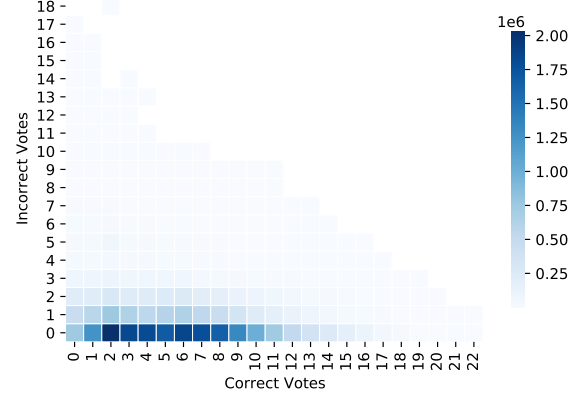


Figure 11: Number of votes from heuristics

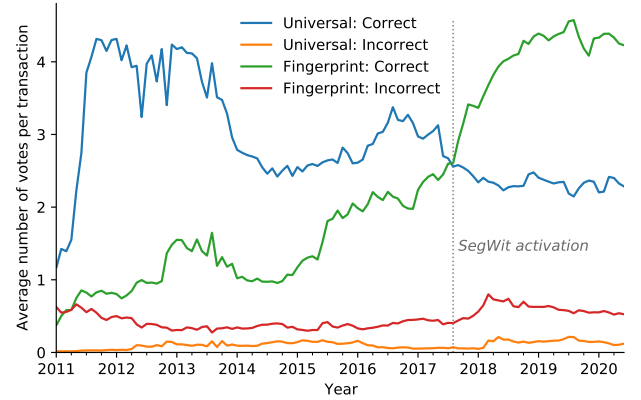


Figure 12: Average number of correct and uncorrect votes per transaction and type of heuristic, over time

heuristic has a notably high TPR compared to many other heuristics (its variants are in many cases redundant).

Comparing the TPR of the heuristics to their coverage, we see moderate positive correlation ($r = 0.55, p = 0.004$), meaning that heuristics that identify more change outputs in our ground truth also tend to be more applicable to the remaining standard transactions. One outlier is the optimal change heuristic, which only returns a result in 13.4 % of the remaining transactions (compared to a TPR of 29.9 % in the ground truth). This is due to the difference in the number of inputs for these two sets (see Table 1). At the same time, many of the fingerprinting heuristics appear to be more applicable to the remaining transactions, signaling higher heterogeneity among those.

23.79 million transactions have votes from a universal heuristic, and 27.56 million have votes from a fingerprint heuristic. 757937 transactions don't have any predictions. Among the 29.29 million transactions with at least one vote, 28.39 million (96.93 %) have at least one *correct* vote.

Figure 11 shows the number of (correct and incorrect) votes

received by each transaction. Figure 12 further breaks down the average number of correct and incorrect predictions per transaction over time, grouped by the type of heuristic. We notice three important trends: the universal heuristics drop over time, likely due to some of the power-of-ten variants becoming less useful (as they contain redundant information, we provide additional plots where the heuristics are aggregated in Appendix B). The consistent fingerprint heuristics instead see a steady uptick in the number of correct votes. This highlights how the increasing variety of protocol features also raises their utility for detecting change outputs. Finally, there’s an uptick in both correct and incorrect fingerprint votes in late 2017 and early 2018, when wallet implementations started to switch to SegWit transaction serialization and address formats (e.g., [6, 40]).

For the following analyses, we exclude the 757937 transactions that don’t have any predictions in order to reduce the chance of false positives. When we later apply the heuristics to the full blockchain, we will also skip transaction without votes from at least one heuristic.

4 Combining Heuristics

A clear disadvantage of choosing a single individual heuristic for change output detection is that they apply only to a subset of transactions (cf. Table 3). Furthermore, some heuristics may be more applicable during certain epochs of Bitcoin’s history than others. In contrast to prior work (e.g., an evaluation of three change heuristics [34]), we also have a larger variety of heuristics available that enable new ways of combining them. Here, we consider two approaches.

4.1 Threshold vote

Figures 11 and 12 suggest that, in general, a majority of cast heuristic votes should produce the correct outcome. However, the number of votes cast varies among transactions, and individual votes could be incorrect. We thus compute a threshold vote: if there are at least t more votes for output a than for output b , then output a is considered the change. Changing the threshold t thereby allows the analyst to require higher degrees of confidence.

We use all of the heuristics listed in Table 3 to compute the threshold vote on the full ground truth data set (as there is no training involved, we do not split the data set). The resulting ROC curve is shown in Figure 13 (for comparison, we also plot the FPR and TPR of the individual heuristics). We can achieve an ROC AUC of 0.9415, and, for example, a 35.2 % true positive rate (TPR) below a false positive rate (FPR) of 0.1 % with a threshold of $t = 7$.

Using a threshold vote may not be ideal as the individual heuristics have varying true positive and false positive rates, and some might be more or less reliable during different periods of Bitcoin’s history. Rather, a specific subset of heuristics

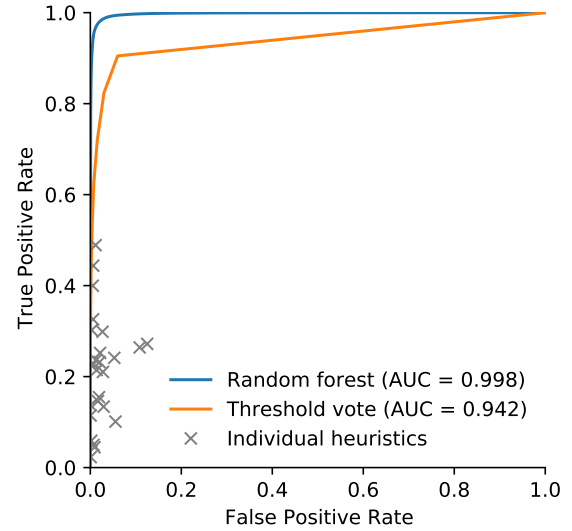


Figure 13: ROC curves for predicting change in the ground truth data set using the threshold vote and the random forest classifier, compared to individual heuristics. The curve of the threshold vote is based on the entire data set, whereas the curve of the random forest is based on the test set. The random forest model includes additional transaction and output characteristics.

may provide better classification accuracy. Instead of manually trying different combinations of heuristics, we opt to use a supervised learning classifier.

4.2 Random forest classifier

We use a random forest classifier to predict the change output of a transaction. A random forest is an ensemble classifier that trains and aggregates the results of individual decision trees. The first step is to transform our transaction-based predictions into an output-based binary classification problem. Every output is either a change (1) or spend (0) output. As before, an individual heuristic may produce one of three outcomes: vote for the output, against the output, or not be able to discern between the outputs.

Next, we add additional characteristics about each output and corresponding transaction that may help the classifier differentiate between distinct types of transactions, or wallets. Output characteristics include the ratio of an output’s value to the total output value of a transaction (the change output is the smaller output in 76.76 % of our ground truth transactions) and its index. Transaction characteristics include its total value in satoshi, the transaction fee paid per byte, its version number, whether it uses SegWit serialization and sets a non-zero locktime, as well as the number of inputs and the time of inclusion (as epochs of 1008 blocks, about one week).

We use the `RandomForestClassifier` implementation in scikit-learn 0.24.1. We add regularization by searching a small parameter grid using a successive halving strategy, optimizing two parameters: the number of features considered at each split, and the minimum number of samples required before each further split. As we consider an analyst that works with a static snapshot of the blockchain, we randomly split our data set into 80 % training and 20 % test set. We use the training set for the hyperparameter search using 4-fold cross-validation, optimizing the area under the curve (AUC) as our scoring metric.⁶ To account for the fact that transactions in the same base cluster may be highly similar, we explicitly ensure that all outputs of a base cluster remain in the same set and fold.

Applying the model to the test set, we receive an AUC of 0.9978 (Figure 13).⁷ We see that random forest model is able to detect a higher share of outputs, especially at low false positive rates, compared to the threshold vote.

In Figure 14 we show the ROC curves of both the threshold vote and the random forest on the same test set, log-transforming the x-axis to highlight the important difference in low false positive rates. The random forest achieves much higher true positive rates at low false positive rates, meaning that it correctly identifies the change output in a larger number of transactions. For example, if we target a false positive rate below 0.1 %, the threshold vote achieves a TPR of around 39 % at a FPR of 0.06 %. For the same FPR, the random forest achieves a TPR of 73 %, almost twice as high.

We train a second random forest model only based on transactions that contain predictions of the universal heuristics in order to later predict change in transactions that contain unspent outputs. Using a similar evaluation strategy as for the full model, the AUC of this model is 0.9978.

We note one caveat: our grouping based on the base cluster ID may not be fully effective at preventing homogeneous transactions from the same entity to appear in both sets because the base clustering is likely incomplete (e.g., an entity’s transactions can be split among multiple clusters, of which some end up in the training and some in the test set). Other researcher or companies with access to private, more heterogeneous ground truth may be able to evaluate this possibility.

4.3 Model validation

We use two data sets to assess the performance of the random forest model outside of our ground truth. First, we use the list of 16 764 transactions identified by Huang et al. [22] as ransom payments related to the Locky and Cerber ransomware. Those payments were identified through clustering, transaction graph analysis and known characteristics of the ransom amounts. The data set contains not only transaction hashes, but also the index (and amount) of the predicted payment output. Out of the 16 764 transactions in the data set, we exclude

⁶Additional details are available in Appendix C.2.

⁷The AUC on the training set is 0.9997

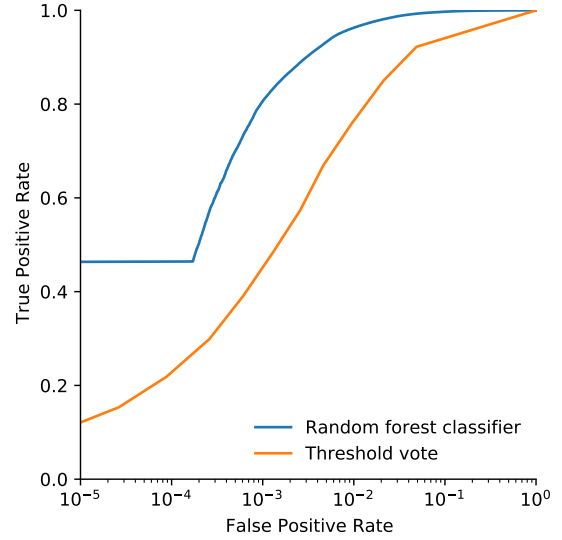


Figure 14: Comparison of the random forest classifier and the threshold vote on the test set. There is a notable difference between the two classifiers for low false positive rates.

3057 because they don’t match our definition of a standard transaction. For 24 transactions, none of the heuristics return any distinct votes. 1636 transactions directly reuse change, and in 850 transactions multi-input clustering already correctly revealed the change output. For the remaining 11 197 transactions we predict the change output using the random forest model and achieve an AUC of 0.996.

Our second data set is constructed using a GraphSense tag-pack [17] that contains 382 tags for addresses of 273 distinct entities (such as exchanges or gambling services). We identify each associated cluster and then extract transactions between the clusters, assuming that the output belonging to a different cluster is the spend output. In total, we extract 2 167 588 transactions between the clusters. As the data is highly skewed towards a few clusters, we limit the total number of transactions for each combination of interacting clusters to 1000 (sampled randomly), giving us 853 011 transactions. Out of these, 13 353 don’t have any predictions, 113 847 reuse an address for change and 457 818 change outputs have already been identified through cluster membership. For the remaining 267 993 transactions we predict the change output and achieve an AUC of 0.973.

5 Clustering Change Outputs

We now use the random forest model to enhance the base clustering. To this end, we predict the change output in 258 million standard transactions with yet unknown change. We exclude 8.9 million transactions where no individual heuristic identified a change output. We use our second model that

does not include the consistent fingerprint heuristics for 12.7 million transactions with unspent outputs.

Informed by the histogram of probabilities (cf. Figure 20 in the appendix), and in order to keep the likelihood of false positives low, we use a conservative probability threshold of 0.99.⁸ This gives us 119.86 million change outputs (for 46.34 % of transactions). In comparison, using the threshold vote to predict change with $t = 7$ returns a change output for 36.19 % of transactions.⁹ We then enhance the base clustering by merging the base cluster of the inputs with the base cluster of the change address in the order that the transactions appear on the blockchain.

5.1 Naive merging leads to cluster collapse

We inspect the enhanced clustering with regards to cluster collapse. First, we look at the size of the enhanced clusters. A typical measure for the size of the cluster is the number of addresses contained in it. However, this may not always be reliable: if entities reuse addresses, their clusters will appear small despite being responsible for a large volume of on-chain transactions. We therefore also inspect the number of transactions originating from the cluster, which is independent of address reuse. Without address reuse, these two measures should correlate as a fresh address is created for every outgoing transaction.

Clustering the identified change outputs reduces 142.4 million affected base clusters into 31.6 million enhanced clusters. However, it leads to severe cluster collapse: there is one large supercluster, of which the prior Mt. Gox supercluster is a part of, that consists of 172.3 million addresses (a 1214 % increase) and 98.6 million transactions (a 2313 % increase). Inspecting the 273 labeled clusters from the Graphsense tag pack, we find that 148 of them have been merged into this supercluster.

5.2 Constraints prevent cluster collapse

The majority of merges we observe involve address clusters that so far had only been seen in a single transaction. In these cases, the impact of a single misclassification is low unless a sequence of such merges collapses multiple clusters. However, in a small number of merges two large clusters are combined, leading to cluster collapse. For example, if there are two large exchanges that interact frequently with each other, a single misidentified change address would collapse their clusters.

Approach. We use this intuition to constrain which clusters we merge. While change outputs predicted by our model

⁸This corresponds to a false positive rate of 0.079 % for our full model. We use a threshold of 0.9955 for the model excluding fingerprint heuristics to match the FPR.

⁹The threshold $t = 7$ corresponds to a FPR of 0.06 % on the test set, the closest to match the FPR of the random forest.

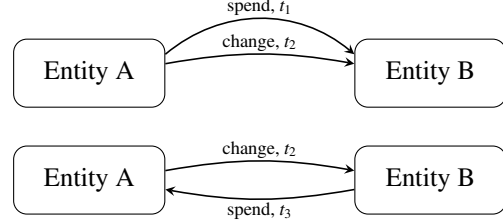


Figure 15: Our constrained clustering prevents the merging of clusters A and B due to conflicting types of payments between them.

should be clustered, we can use outputs predicted to be spends to prevent cluster merges: the input cluster should not be clustered into the cluster of the spend. Let p_i be the probability returned by the random forest model for output i , then we define two thresholds p_{change} and p_{spend} such that if $p_i > p_{change}$ the clusters should be merged, and if $p_i < p_{spend}$ then the clusters should not be merged. In many cases, these constraints prevent the spend and change output of a single transaction to end up in the same cluster (cf. Figure 15). It is therefore a stronger assumption than our core assumption, which only considered the change.

This approach is comparable to that by Ermilov, Panov, and Yanovich [14] to use address tags in combination with a probabilistic model to reduce the number of conflicting tags in the final clustering. However, most public sources of address tags only contain information on a limited amount of intermediaries and clusters. Our approach, instead, potentially covers all clusters appearing in the 258 million standard transactions, including those that may be hard to interact with (and thereby tag) manually. Due to the size of the blockchain, and the large number of predictions we have, here we only consider the binary case of preventing any potential conflict, accepting that we may prevent some valid merges in the process.

We implement this approach as a constrained union-find algorithm that prevents merging two clusters that are related by a predicted spend output. That is, for every spend from cluster c_i to cluster c_j predicted by the random forest with $p \leq p_{spend}$, we add a constraint to cluster c_i that it must not be merged with cluster c_j . When merging two clusters, we compare each cluster’s set of constraints to the set of members of the other cluster and skip the merge if the sets intersect.

Results. Using the same p_{change} as before and setting $p_{spend} = 0.01$, the constrained clustering skips 419641 merges that would have violated constraints and retains 247111 more individual clusters than the unconstrained clustering. 8.9 million of the predicted merges are redundant because the clusters had already been merged.

We find that the constrained merging prevents the previously observed severe cluster collapse. For example, the constrained clustering does not produce the large Mt. Gox

Table 4: Transaction count of smaller cluster being merged

Percentile	Transaction count	
	Naive	Constrained
90	1	1
99	6	6
99.9	27	26
99.99	143	114
99.999	2943	646

supercluster: the enhanced cluster contains 4.3 million transactions (a 5 % increase) and 14.4 million addresses (a 10 % increase). Assessing the 273 labeled clusters, there are only ten instances left where two labeled clusters are merged together. Three of those instances are pairs of gambling services (e.g., dadice.com and 999dice.com), the others also involve exchanges and darknet markets. We suspect that unusual types of payouts from these services might have triggered their collapse.

Table 4 shows percentiles of the number of transactions associated with each smaller cluster that was merged into a larger cluster. In at least 90 % of merges, the transaction count of the smaller cluster was equal to or below 1 (0 if unspent and not used otherwise). This highlights how change address clustering is useful to merge small clusters occurring in only a single transaction that are missed by multi-input clustering. We can also see that until the 99.99 percentile, transaction counts are very similar between the unconstrained and the constrained clustering, but that the constraints prevent large merges from taking place, preventing cluster collapse.

Figure 16 shows histograms of the increase in address count and transaction count for each affected cluster. Increase here is the difference between the resulting size of the enhanced cluster and the largest individual base cluster before change address clustering. The histograms are plotted on a log-log axis, highlighting the extreme skewness of the distribution. Most clusters increase only by a few addresses or transactions.¹⁰ The supercluster can also easily be spotted.

Finally, we inspect the clustering with regards to the maximum time difference between two transactions. We would expect that transactions in the merged clusters should, in general, not be too far apart. If two clusters are merged such that the maximum time gap between transactions increases substantially (say, from a few months to a few years), then this potentially indicates an incorrect clustering. For the majority (28 million) of the new clusters, each prior base cluster had fewer than 2 transactions. The median max time gap for those clusters amounts to just below 2 hours. For instances where

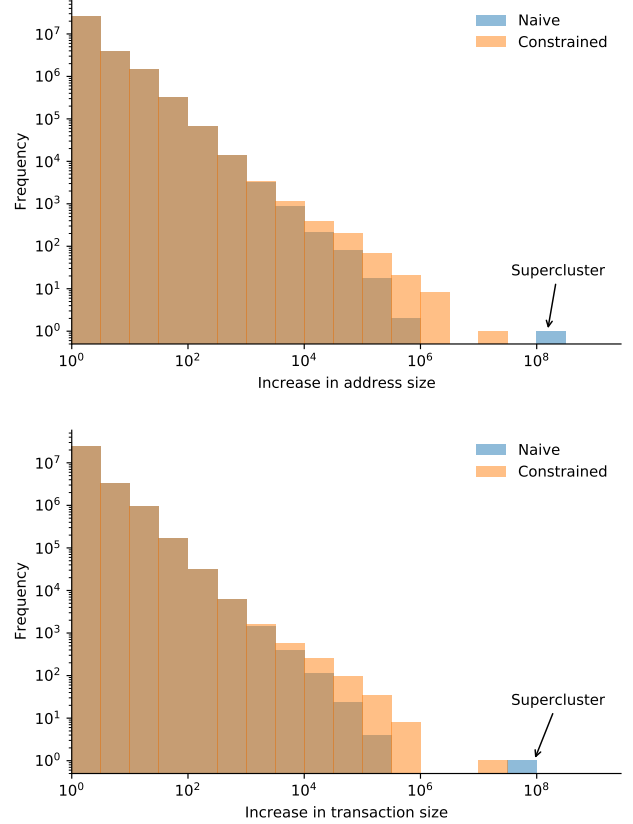


Figure 16: Absolute increase in address and transaction count per affected cluster, using either the naive or the constrained approach. Note the caveat in footnote 10.

¹⁰Due to the log-scale, the histogram showing the increase in transaction size excludes 2.4 million clusters that, using the naive approach, didn't increase in their transaction count, and 2.6 million using the constrained approach, respectively.

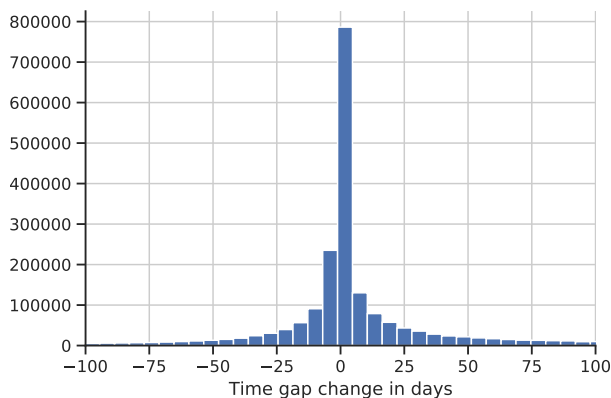


Figure 17: Change in maximum time gap for clusters where at least one merged base cluster had two or more transactions. The X-axis is cut off to highlight the most important area (accounting for 86 % of clusters).

at least one original cluster already consisted of 2 or more transactions, the median increase in the maximum time gap is only around half an hour. Overall, the distribution is highly centered (cf. Figure 17).

6 Impact on Blockchain Analyses

Change address clustering is generally used as a pre-processing step before analyzing activity of entities on the blockchain. Using different heuristics (or none at all) thus affects the outcome of these analyses. In this section we present two exemplary applications to highlight the impact change address clustering has on blockchain analyses. Furthermore, we compare our clustering to one created with a popular change address heuristic based on address reuse.

6.1 Cashout flows from darknet markets to exchanges

For our first example we evaluate the impact of our enhanced clustering on analysing payment flows from darknet markets to exchanges, using address tags from the GraphSense tag pack to identify relevant clusters. The tag pack contains address tags for 117 exchanges and 15 darknet markets.

We extract the value of all outputs in transactions initiated by a darknet market that are sending bitcoins to an exchange, comparing the transaction volume calculated using our base clustering to that of our enhanced clustering. The median increase in output value across all 15 exchanges amounts to 13.86 %. Overall, the total amount of bitcoins flowing from the darknet markets to exchanges increases from BTC 821 500 to BTC 961 519 (a 17 % increase). An overview of each individual market’s increase in transaction volumes is presented in Table 5 in the appendix.

6.2 Improved estimate of velocity

We replicate the analysis of velocity conducted by Kalodner et al. [24], an example for a longitudinal analysis of economic activity occurring on the Bitcoin blockchain. For this analysis, clustering is used to remove self-payments of users (such as change outputs), which would artificially inflate estimates of economic activity. The better and more complete our clustering, the more self-payments are removed and hence the lower the estimate will be.

Our refined clustering reduces their estimate of bitcoins moved per day between January 2017 to June 2020 by about 11.7 %. This number is quite similar to the impact on cash-out flows.

In general, much activity is generated by large intermediaries, which are more prone to be clustered by the multi-input heuristic. Our enhanced clustering predominantly merged small clusters, which should give a more realistic estimate for their activity. As we chose relatively conservative thresholds to enhance the clustering, loosening them could lead to larger effect sizes.

6.3 Comparison to the Meiklejohn heuristic

We compare our constrained clustering to one created using the address reuse-based heuristic presented by Meiklejohn et al. [30], which has subsequently been used in other studies (e.g., [11, 35]). While the authors highlight the need for manual intervention to prevent cluster collapse, this is likely infeasible for analysts without in-depth domain knowledge or the right set of tools. There are two possible ways to implement the heuristic, either such that an output is considered to be change if its address has not appeared in any previous transactions (we call this variant local), or such that over the entire blockchain it appears in only one output (global). The latter thus depends on the particular state of the blockchain.

Applying the local heuristic to the standard transactions with unknown change produces a large supercluster, including 106.3 million transactions and 216.1 million addresses. 166 of the tagged clusters appear in this supercluster. Similarly, the global variant produces a supercluster containing 113.3 million transactions and 230.4 million addresses, with 175 tagged clusters ending up in the supercluster.

To characterize the difference between the clusterings, we look at the probability that two random addresses are clustered in a particular clustering. The probability of two addresses being clustered together increases by a factor of 40 when using the Meiklejohn heuristic compared to our constrained clustering, further highlighting the cluster collapse (cf. Table 6 in the appendix).

Finally, we take closer look at the individual predictions. First, comparing the predictions of the local Meiklejohn heuristic to those of our constrained clustering, they overlap on 53.6 million transactions and differ for 0.9 million

of those. The global heuristic differs on 1.2 million transactions out of an overlapping 62.3 million. For those conflicting predictions, the difference in output value to our clustering amounts to BTC 1.6 million, or USD 8.8 billion of economic activity for the local heuristic. For the global heuristic, the difference amounts to BTC 3 million, or USD 16.3 billion.

7 Discussion

Our results confirm existing expectations (e.g., [36]) that in many instances change address detection is feasible with high precision. This may motivate users to adopt privacy-enhancing countermeasures. Techniques to avoid address reuse (such as one-time addresses [43]) would reduce the effectiveness of the multi-input heuristic and thus also of our ground truth extraction technique. Widespread use of cooperative obfuscation techniques like CoinJoin [28] and PayJoin [12] would make multi-input clustering unreliable (by causing cluster collapse), and ambiguity-based obfuscation such as the randomization of address types (e.g., [6]) could thwart the fingerprinting heuristics. If users adopted these techniques more widely (cf. [31]), it would also make law enforcement investigations into cryptocurrencies more challenging, increasing the need for complimentary approaches that do not rely on address clustering [33].

At the same time, the transparency blockchains provide shouldn't categorically be considered a disadvantage, and we think that our work is useful beyond law enforcement purposes. Researchers studying social or economic questions, using cryptocurrencies as a "social science lab" [9], benefit from increased transparency. Blockchains contain a valuable trove of financial data that traditionally hasn't been available to researchers [7]. In this context, our techniques can help to remove ambiguity about ownership of funds, improving data quality and leading to more realistic estimates.

8 Conclusion

Address clustering is an important cornerstone of many blockchain analyses. In this paper, we've taken a first step towards building better models that allow analysts to identify change outputs in transactions, enabled by a new ground truth data set extracted from the Bitcoin blockchain. We've further demonstrated how constraints based on our model's predictions can prevent cluster collapse. Finally, we've shown that the resulting, enhanced clustering changes the outcomes of economic analyses, in our examples by about 11 % to 17 %. Given this impact, we hope that our work will encourage further research into change address clustering.

Acknowledgements

We thank Rainer Böhme and Kevin Lee for their feedback on an earlier draft of this paper. This work is supported by NSF Award CNS-1651938 and a grant from the Ripple University Blockchain Research Initiative.

References

- [1] Svetlana Abramova and Rainer Böhme. "Your Money or Your Privacy: A Systematic Approach to Coin Selection." In: *Cryptoeconomic Systems '20*. 2020.
- [2] Elli Androulaki, Ghassan O Karame, Marc Roeschlin, Tobias Scherer, and Srdjan Capkun. "Evaluating user privacy in Bitcoin". In: *International Conference on Financial Cryptography and Data Security*. Springer. 2013, pp. 34–51.
- [3] Kristov Atlas. *Lexicographical Indexing of Transaction Inputs and Outputs*. URL: <https://github.com/bitcoin/bips/blob/master/bip-0069.mediawiki> (visited on 01/20/2020).
- [4] Massimo Bartoletti, Barbara Pes, and Sergio Serusi. "Data mining for detecting Bitcoin Ponzi schemes". In: *2018 Crypto Valley Conference on Blockchain Technology (CVCBT)*. IEEE. 2018, pp. 75–84.
- [5] Massimo Bartoletti and Livio Pompianu. "An analysis of Bitcoin OP_RETURN metadata". In: *Financial Cryptography and Data Security: FC 2017 International Workshops, WAHC, BITCOIN, VOTING, WTSC, and TA, Sliema, Malta, April 7, 2017, Revised Selected Papers*. Springer. 2017, pp. 218–230.
- [6] *Bitcoin Core 0.16.0*. URL: <https://bitcoincore.org/en/releases/0.16.0/> (visited on 02/25/2021).
- [7] Desamparados Blazquez and Josep Domenech. "Big Data sources and methods for social and economic analyses". In: *Technological Forecasting and Social Change* 130 (2018), pp. 99–113.
- [8] *Blockchair.com API v.2.0.76 Documentation: Privacy-o-meter*. URL: https://blockchair.com/api/docs#link_M6 (visited on 02/22/2021).
- [9] Rainer Böhme, Nicolas Christin, Benjamin Edelman, and Tyler Moore. "Bitcoin: Economics, technology, and governance". In: *Journal of Economic Perspectives* 29.2 (2015), pp. 213–38.
- [10] Tao-Hung Chang and Davor Svetinovic. "Improving bitcoin ownership identification using transaction patterns analysis". In: *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 50.1 (2018), pp. 9–20.

- [11] Mauro Conti, Ankit Gangwal, and Sushmita Ruj. “On the economic significance of ransomware campaigns: A Bitcoin transactions perspective”. In: *Computers & Security* 79 (2018), pp. 162–189.
- [12] Nicolas Dorier. *A Simple Payjoin Proposal*. URL: <https://github.com/bitcoin/bips/blob/master/bip-0078.mediawiki> (visited on 01/20/2020).
- [13] Mark Erhardt. *An Evaluation of Coin Selection Strategies*. Master Thesis. 2016.
- [14] Dmitry Ermilov, Maxim Panov, and Yury Yanovich. “Automatic Bitcoin address clustering”. In: *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE. 2017, pp. 461–466.
- [15] Mark Friedenbach, BtcDrak, Nicolas Dorier, and kinoshitajona. *Relative lock-time using consensus-enforced sequence numbers*. URL: <https://github.com/bitcoin/bips/blob/master/bip-0068.mediawiki> (visited on 01/20/2020).
- [16] Steven Goldfeder, Harry Kalodner, Dillon Reisman, and Arvind Narayanan. “When the cookie meets the blockchain: Privacy risks of web payments via cryptocurrencies”. In: *Proceedings on Privacy Enhancing Technologies* 2018.4 (2018), pp. 179–199.
- [17] *GraphSense Public TagPacks*. URL: <https://github.com/graphsense/graphsense-tagpacks> (visited on 04/01/2021).
- [18] David A. Harding and Peter Todd. *Opt-in Full Replace-by-Fee Signaling*. URL: <https://github.com/bitcoin/bips/blob/master/bip-0125.mediawiki> (visited on 01/20/2020).
- [19] Mikkel Alexander Harlev, Haohua Sun Yin, Klaus Christian Langenheddt, Raghava Mukkamala, and Ravi Vatrpu. “Breaking bad: De-anonymising entity types on the Bitcoin blockchain using supervised machine learning”. In: *Proceedings of the 51st Hawaii International Conference on System Sciences*. 2018.
- [20] Martin Harrigan and Christoph Fretter. “The unreasonable effectiveness of address clustering”. In: *2016 Intl IEEE Conferences on Ubiquitous Intelligence & Computing, Advanced and Trusted Computing, Scalable Computing and Communications, Cloud and Big Data Computing, Internet of People, and Smart World Congress (UIC/ATC/ScalCom/CBDCoM/IoP/SmartWorld)*. IEEE. 2016, pp. 368–373.
- [21] Yining Hu, Suranga Seneviratne, Kanchana Thilakarathna, Kensuke Fukuda, and Aruna Seneviratne. “Characterizing and Detecting Money Laundering Activities on the Bitcoin Network”. In: *arXiv preprint arXiv:1912.12060* (2019).
- [22] Danny Yuxing Huang, Maxwell Matthaios Aliapoulos, Vector Guo Li, Luca Invernizzi, Elie Bursztein, Kylie McRoberts, Jonathan Levin, Kirill Levchenko, Alex C Snoeren, and Damon McCoy. “Tracking ransomware end-to-end”. In: *IEEE Symposium on Security and Privacy*. IEEE. 2018, pp. 618–631.
- [23] Marc Jourdan, Sebastien Blandin, Laura Wynter, and Pralhad Deshpande. “Characterizing entities in the Bitcoin blockchain”. In: *2018 IEEE International Conference on Data Mining Workshops (ICDMW)*. IEEE. 2018, pp. 55–62.
- [24] Harry Kalodner, Malte Möser, Kevin Lee, Steven Goldfeder, Martin Plattner, Alishah Chator, and Arvind Narayanan. “Blocksci: Design and applications of a blockchain analysis platform”. In: *29th USENIX Security Symposium*. 2020, pp. 2721–2738.
- [25] Yu-Jing Lin, Po-Wei Wu, Cheng-Han Hsu, I-Ping Tu, and Shih-wei Liao. “An evaluation of Bitcoin address classification based on transaction history summarization”. In: *2019 IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*. IEEE. 2019, pp. 302–310.
- [26] Eric Lombrozo, Johnson Lau, and Pieter Wuille. *Segregated Witness (Consensus layer)*. URL: <https://github.com/bitcoin/bips/blob/master/bip-0141.mediawiki> (visited on 01/20/2020).
- [27] Damiano Di Francesco Maesa, Andrea Marino, and Laura Ricci. “Data-driven analysis of bitcoin properties: exploiting the users graph”. In: *International Journal of Data Science and Analytics* 6.1 (2018), pp. 63–80.
- [28] Gregory Maxwell. *CoinJoin: Bitcoin Privacy for the Real World*. 2013. URL: <https://bitcointalk.org/index.php?topic=279249.0> (visited on 06/30/2019).
- [29] Sarah Meiklejohn and Claudio Orlandi. “Privacy-enhancing overlays in Bitcoin”. In: *International Conference on Financial Cryptography and Data Security*. Springer. 2015, pp. 127–141.
- [30] Sarah Meiklejohn, Marjori Pomarole, Grant Jordan, Kirill Levchenko, Damon McCoy, Geoffrey M Voelker, and Stefan Savage. “A fistful of bitcoins: characterizing payments among men with no names”. In: *Internet Measurement Conference*. ACM. 2013, pp. 127–140.
- [31] Malte Möser and Rainer Böhme. “The price of anonymity: empirical evidence from a market for Bitcoin anonymization”. In: *Journal of Cybersecurity* 3.2 (2017), pp. 127–135.

- [32] Malte Möser, Rainer Böhme, and Dominic Breuker. “An inquiry into money laundering tools in the Bitcoin ecosystem”. In: *eCrime Researchers Summit (eCRS)*, 2013. IEEE. 2013, pp. 1–14.
- [33] Malte Möser and Arvind Narayanan. *Effective cryptocurrency regulation through blacklisting*. Preprint. 2019.
- [34] Jonas David Nick. “Data-driven de-anonymization in Bitcoin”. MA thesis. ETH-Zürich, 2015.
- [35] Francesco Parino, Mariano G Beiró, and Laetitia Gauvin. “Analysis of the Bitcoin blockchain: socio-economic factors behind the adoption”. In: *EPJ Data Science* 7.1 (2018), p. 38.
- [36] *Privacy - Bitcoin Wiki*. URL: <https://en.bitcoin.it/Privacy> (visited on 12/15/2020).
- [37] Fergal Reid and Martin Harrigan. “An Analysis of Anonymity in the Bitcoin System”. In: *Security and Privacy in Social Networks*. Springer, 2013, pp. 197–223.
- [38] Dorit Ron and Adi Shamir. “Quantitative analysis of the full Bitcoin transaction graph”. In: *International Conference on Financial Cryptography and Data Security*. Springer. 2013, pp. 6–24.
- [39] Jürgen E Schatzmann and Bernhard Haslhofer. “Bitcoin Trading is Irrational! An Analysis of the Disposition Effect in Bitcoin”. In: *arXiv preprint arXiv:2010.12415* (2020).
- [40] *SegWit FAQ*. URL: <https://help.coinbase.com/en/pro/getting-started/general-crypto-education/segwit-faq> (visited on 04/27/2021).
- [41] Peter Todd. *Discourage fee sniping with nLockTime #2340*. 2014. URL: <https://github.com/bitcoin/bitcoin/pull/2340> (visited on 08/11/2015).
- [42] Kentaroh Toyoda, Tomoaki Ohtsuki, and P Takis Mathiopoulos. “Multi-class Bitcoin-enabled service identification based on transaction history summarization”. In: *2018 IEEE International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCoM) and IEEE Smart Data (SmartData)*. IEEE. 2018, pp. 1153–1160.
- [43] Nicolas Van Saberhagen. *CryptoNote v2.0*. 2013. URL: <https://cryptonote.org/whitepaper.pdf> (visited on 01/14/2018).
- [44] Mark Weber, Giacomo Domeniconi, Jie Chen, Daniel Karl I Weidele, Claudio Bellei, Tom Robinson, and Charles E Leiserson. “Anti-money laundering in Bitcoin: Experimenting with graph convolutional networks for financial forensics”. In: *arXiv preprint arXiv:1908.02591* (2019).
- [45] Yuhang Zhang, Jun Wang, and Jie Luo. “Heuristic-Based Address Clustering in Bitcoin”. In: *IEEE Access* 8 (2020), pp. 210582–210591.

A Protocol characteristics used for fingerprinting

- Input/output count: the number of inputs and/or outputs may indicate a wallet software’s behavior of creating transactions. While the number of inputs depends on the UTXOs available to the user, some commonly occurring patterns such as peeling chains ([30, 32]) have consistent input and output counts.
- Version: BIP 68 [15] introduced relative timelocks for transactions, which requires transaction to set the transaction version to 2.
- Locktime: Transactions can set a timelock such that they are valid only after the tip of the chain has passed a specific block height or timestamp. Some clients (e.g., Bitcoin Core) produce timelocked transactions by default to prevent fee-sniping [41].
- Replace-by-fee (RBF): Transactions opting into the replace-by-fee policy can be replaced by a similar transaction paying a higher fee [18].
- SegWit: Segregated witness [26] is a protocol update that enabled storing the inputs’ signatures outside of the transaction, thereby increasing available space in blocks. As the upgrade is backwards-compatible, not all wallets produce SegWit transactions. A wallet might also be able to produce SegWit transaction, but may be required to use non-SegWit serialization if none of the inputs use SegWit. We call this behavior SegWit-conform.
- Ordered inputs/outputs: BIP 69 [3] defines non-binding rules (i.e. not enforced by the consensus mechanism) for lexicographically sorting inputs and outputs in a transaction. (A limitation of our implementation is that it does not compare the raw `scriptPubKey` in case the output values are equal, as they are not available in BlockSci).
- Zero-conf: Bitcoin user’s are encouraged to wait for up to six confirmations (about an hour) before accepting a payment, as there is a risk that funds might be double-spent. A transaction spending inputs without any confirmations indicates willingness to accept the double-spending risk, which could be specific to certain intermediaries.
- Transaction fee: Bitcoin users pay transaction fees for their transactions to be included into the blockchain by miners. Some clients may pay the same exact fee (either absolute, or relative to the transaction’s size) for every transaction.

- **Multisignature:** Multisignature scripts allow to specify a list of public keys and a threshold m such that the redeemer must provide valid signatures for m out of n of these keys. They aren't typically used by normal end-user wallets.
- **Address types:** Bitcoin Core defines a number of standardized output scripts types including Pay-to-Pubkey-Hash (P2PKH), Pay-to-Script-Hash (P2SH) as well as their respective SegWit variants (P2WPKH and P2WSH). Often, a wallet consistently uses a specific address type. (Compared to the normal address type heuristic, the fingerprint checks for overlap with the address types of all inputs of the spending transaction).

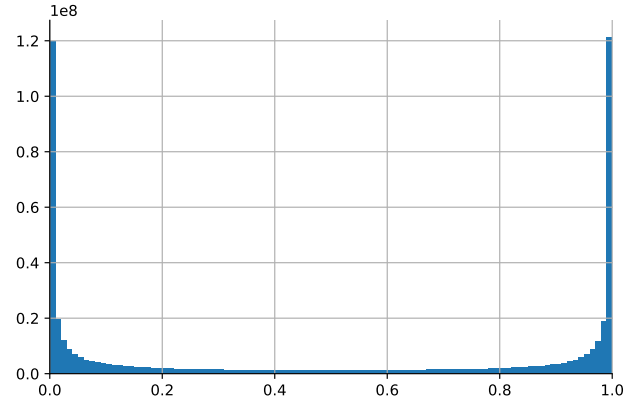


Figure 20: Probabilities returned by the random forest classifier for standard transactions with unknown change

B Additional plots and tables

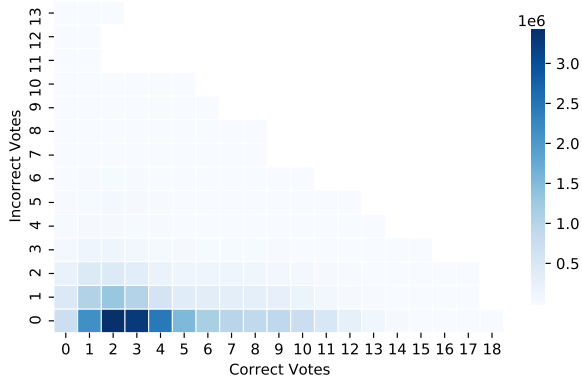


Figure 18: Number of votes from heuristics (with compressed power-of-ten heuristic)

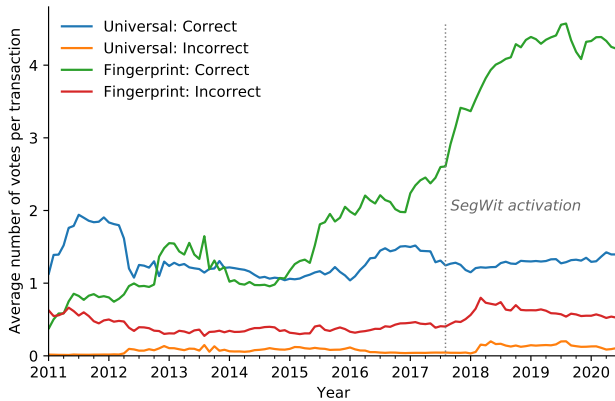


Figure 19: Average number of correct and incorrect votes per transaction and type of heuristic, over time (with compressed power-of-ten heuristic)

Table 5: Change in outgoing transaction volumes of darknet markets to exchanges using the base clustering (before) and our enhanced clustering (after).

Tag name	Volume (BTC)		
	Before	After	Change (%)
abraxasmarket	21 925	24 334	10.99
agoramarket	158 355	172 523	8.95
alphabaymarket	35 495	51 682	45.61
babylonmarket	222	287	29.21
blackbankmarket	8 292	9 942	19.89
blueskymarket	2 520	3 391	34.58
cannabisroadmarket	6	8	39.64
doctordmarket	224	426	90.52
evolutionmarket	49 891	85 919	72.21
middleearthmarket	11 793	12 389	5.06
nucleusmarket	45 265	50 289	11.1
pandoraopenmarket	8 708	9 539	9.54
sheepmarket	12 104	13 391	10.64
silkroad2market	47 292	49 847	5.4
silkroadmarket	419 409	477 551	13.86
<i>Total</i>	821 500	961 519	17.05

Table 6: Characteristics of clusterings created with the Meiklejohn heuristic in comparison to our constrained clustering

Characteristic	Meiklejohn heuristics	
	Local	Global
Coverage	47.8 %	54.7 %
Largest cluster		
• # addresses	216.1 M	230.4 M
• # transactions	106.3 M	113.3 M
Address pairs clustered		
• Neither	90.80 %	90.00 %
• Ours only	0.06 %	0.05 %
• Meiklejohn only	8.98 %	9.78 %
• Both	0.17 %	0.18 %
Predictions		
• Total	123.8 M	141.5 M
• Overlapping	54.6 M	62.3 M
Conflicting predictions		
• Count	0.9 M	1.2 M
• Difference in BTC	1.6 M	3.0 M
• Difference in USD	8.8 B	16.3 B

C Further insights and technical details

C.1 Filtering the ground truth data set

Selecting transactions with two outputs, no OP_RETURN outputs, where no input address has been directly reused in the outputs and where at least one output is in the same base cluster as the inputs yields a total of 46.79 million transactions. We first exclude 1.08 million transactions with unspent outputs, as our subsequent analyses rely upon the spending transactions being known.

Transactions with two change candidates. Out of the 45.7 million transactions with at least one change candidate, for 0.92 million transactions *both* outputs are in the same base cluster. This can happen when a user transfers funds to an address in their own wallet, an online service restructures their funds, or cluster collapse leads to merging of both outputs’ addresses. In a first step, we exclude all transactions with two change candidates.

However, it is possible that there are yet unidentified transactions where both outputs do belong to the same entity. This should occur only in rare cases, but there may be specific intermediaries that create such transactions more frequently (e.g., 28.67 % of these transactions originate from only two base clusters). We therefore exclude *all* transactions from base clusters where more than 10 % of transactions exhibit such

behavior. This removes an additional 472 607 transactions in 8635 base clusters from our ground truth.

Potential false positives. A risk of using the base clustering to extract ground truth is that the multi-input heuristic could already have produced false positives. For example, if a user Alice makes a payment to merchant Bob and their wallet addresses are incorrectly clustered together, her spend output would appear to be the change.

To this end, we first remove 365 010 transactions belonging to the Mt. Gox supercluster that resulted from users being able to import their private keys into the service [20]. Next, we spot-check our base clustering against the website WalletExplorer.¹¹ For the 100 largest base clusters in our ground truth we select 25 addresses at random and collect the tag (which is either explicitly named or pseudo-random) that WalletExplorer assigns to the address. In five instances, the addresses yield multiple tags. Four of these return only additional pseudo-random tags, which upon manual inspection we believe to be the result of a heuristic to not link addresses in transactions with large numbers of inputs. Only one base cluster contains addresses with two different named tags: “LocalBitcoins.com-old” and “AnxPro.com”. This could be a result of cluster collapse, or an instance of mislabeling on the side of WalletExplorer. We remove the 87 931 transactions from this base cluster from our ground truth. Overall, this check gives us some confidence that our base clustering does not already include wide-spread cluster collapse.

Change address reuse. Our initial selection removed transactions where the change address appeared in an input of the transaction. Yet, we find many instances where the change address did not appear in the inputs but had been seen before. For example, a base cluster labeled by WalletExplorer as the gambling service “SatoshiDice”, contains 5.77 million transactions that use only 50 different change addresses. Similarly, there are 1.27 million transactions from a base cluster tagged as “LuckyB.it” that all use a single change address.¹² In many of these cases, the change address could have already been revealed (before the transaction took place) through the multi-input heuristic.

If the change is known at the time the transaction is created, applying change heuristics is unnecessary. In contrast, whenever a transaction uses a fresh address for change, it cannot possibly be revealed as the change at the time the transaction was created. With this intuition, we remove transactions with change addresses that were not freshly generated if, at the time they were included in the blockchain, the change had already been revealed by the multi-input heuristic. This removes a total of 13.81 million transactions (92.80 % of transactions with reused change addresses). Table 7 provides an overview

¹¹<https://walleterexplorer.com>

¹²1NxaBCFQwejSZbQfWcYNwgqML5wWoE3rK4

Table 7: Number of transactions (in million) in our ground truth data set with fresh or reused spend and change outputs.

<i>Change</i>	<i>Spend</i>		Total
	Reused	Fresh	
Reused	0.52	0.55	1.07
Fresh	16.68	12.30	28.98
Total	17.20	12.85	30.05

of whether the change and spend addresses are fresh in our ground truth data.

C.2 Random forest classifier

Encoding. Due to the large size of the data set, we forgo one-hot encoding and instead use the following ordinal encoding for the heuristics:

- 1** the heuristic votes for the output
- 0** the heuristic votes neither for nor against the output
- 1** the heuristic votes against the output

Model parameter. Our hyperparameter search returns the following parameters:

- Full model
 - max_features: 5
 - min_samples_split: 50
- No fingerprinting model
 - max_features: 5
 - min_samples_split: 100