

Why is a Ravencoin Like a TokenDesk? An Exploration of Code Diversity in the Cryptocurrency Landscape

Pierre Reibel¹, Haaron Yousaf¹, and Sarah Meiklejohn¹

University College London
{pierre.reibel.16,h.yousaf,s.meiklejohn}@ucl.ac.uk

Abstract. Interest in cryptocurrencies has skyrocketed since their introduction a decade ago, with hundreds of billions of dollars now invested across a landscape of thousands of different cryptocurrencies. While there is significant diversity, there is also a significant number of scams as people seek to exploit the current popularity. In this paper, we seek to identify the extent of innovation in the cryptocurrency landscape using the open-source repositories associated with each one. Among other findings, we observe that while many cryptocurrencies are largely unchanged copies of Bitcoin, the use of Ethereum as a platform has enabled the deployment of cryptocurrencies with more diverse functionalities.

1 Introduction

Since the introduction of Bitcoin in 2008 [22] and its deployment in January 2009, cryptocurrencies have become increasingly popular and subject to increasing amounts of hype and speculation. Initially, the promise behind cryptocurrencies like Bitcoin was the ability to send frictionless global payments: anyone in the world could act as a peer in Bitcoin’s peer-to-peer network and broadcast a transaction that — without having to pay exorbitant fees — would send money to anyone else in the world, regardless of their location, citizenship, or what bank they used. This is achieved by the decentralization inherent in the open consensus protocol, known as proof-of-work, that allows any peer to not only broadcast transactions but also act to seal them into the official ledger.

While the realities of Bitcoin have shifted in the ensuing years, the landscape of cryptocurrencies has also shifted considerably. There are now thousands of alternative cryptocurrencies, supporting more exotic functionalities than the simple atomic transfer of money supported by Bitcoin. Ethereum, for example, promises to act as a distributed consensus computer (the Ethereum Virtual Machine, or EVM for short) by enabling arbitrary stateful programs to be executed by transactions, while Monero and Zcash promise to improve on the anonymity achieved by Bitcoin transactions. Others don’t promise new functionalities but instead aim to support the same functionality as Bitcoin in more cost-effective ways; e.g., Zilliqa [15,8,29,18,16] and Cardano [14,6] incorporate respective ideas

from the academic literature about achieving consensus without relying entirely on proof-of-work.

Alongside this rapid expansion in the functionality of cryptocurrencies (or indeed the general applicability of the underlying concept of a blockchain), there has also been a genuine explosion of investment into these technologies. In July 2013, for example, there were 42 cryptocurrencies listed on the popular data tracker CoinMarketCap,¹ and the collective market capitalization was just over 1 billion USD. In July 2018, in contrast, there were 1664 cryptocurrencies, and the collective market capitalization was close to 1 trillion USD. While comprehensive in terms of deployed cryptocurrencies, this list does not even include many of the recent “initial coin offerings” (ICOs) that have similarly attracted millions in investment despite there having been many documented scams.²³ Against this backdrop of hype and investment, it is thus crucial to gain some insight into the different types of functionalities offered by these many different cryptocurrencies, to understand which coins offer truly novel features and are backed by genuine development efforts, and which ones are merely hoping to cash in on the hype.

This paper takes a first step in this direction, by examining the entire landscape of cryptocurrencies in terms of the publicly available source code used to support each one. While source code may not be the most accurate representation of a cryptocurrency (as, for example, the actual client may use a different codebase), it does reflect the best practices of the open-source software community, so we believe it to be a reasonable proxy for how a cryptocurrency does (or should) represent itself. In particular, we begin by describing our collected dataset in Section 4, where we can already observe that many listed cryptocurrencies are in fact tokens based on the Ethereum blockchain. We then move on in Section 5 to a general identification of the ways in which cryptocurrencies copy their ideas from those of others, and observe that Bitcoin is by far the most popular in this respect. Due to the dominance of Bitcoin and Ethereum, we explore them in more detail in Sections 6 and 7 before concluding in Section 8.

2 Related Work

We treat as related research that measures either general properties of open-source software, or research that measures properties of cryptocurrencies. In terms of the former, there have been numerous papers measuring GitHub repositories. For example, Hu et al. [11] and Thung et al. [30] measured the influence of software projects according to their position of their repositories and developers in the GitHub social graph, and others have taken advantage of the volume of source code available on GitHub to analyze common coding practices [35] or how bugs vary across different programming languages [25].

In terms of the latter, there are by now many papers that have focused on measuring properties of both the peer-to-peer networks [17,7,3,1] and the

¹ <https://coinmarketcap.com/historical/20130721/>

² <https://deadcoins.com/>

³ <https://magoo.github.io/Blockchain-Graveyard/>

blockchain data associated with cryptocurrencies [26,27,19,28,21,13,5,31], as well as their broader ecosystem of participants [20,34,32,33]. Given the volume of research, we focus only on those papers most related to our own, in that they analyze properties across multiple cryptocurrencies, rather than within a single one like Bitcoin. In terms of comparing Bitcoin and Ethereum, Gencer et al. [9] compared the level of decentralization in their peer-to-peer networks and found, for example, that Ethereum mining was more centralized than it was in Bitcoin, but that Bitcoin nodes formed more geographic clusters. Azouvi et al. [2] also compared their level of decentralization, in terms of the discussions on and contributions to their GitHub repositories, and found that Ethereum was more centralized in terms of code contribution and both were fairly centralized in terms of the discussions. Gervais et al. [10] introduced a framework for identifying the tradeoff between security and performance in any cryptocurrency based on proof-of-work, and found that the same level of resilience to double-spending attacks was achieved by 37 blocks in Ethereum as by 6 blocks in Bitcoin. Finally, Huang et al. [12] compared the effectiveness of different mining and speculation activities for 18 cryptocurrencies, and found that the profitability of both was affected by when a cryptocurrency was listed on an exchange.

3 Background

Here we provide some brief background on the functionality of Bitcoin and Ethereum, which are the two biggest cryptocurrencies in terms of their market capitalization. As general terminology, we use *cryptocurrency* to refer to anything with an exchangeable unit of value, *coin* to refer to a cryptocurrency with its own dedicated blockchain, and *token* to refer to a cryptocurrency that operates instead using the blockchain of another cryptocurrency (e.g., Ethereum).

Fundamentally, the main functionality offered by Bitcoin is the atomic transfer of money from a sender (or set of senders) to a recipient (or set of recipients). This is supported by a limited scripting language, known simply by the name Script, which defines how transactions are created and verified. Bitcoin is also, however, a standalone platform, and thus its codebase must do significantly more than support this so-called transaction layer [4]. In particular, it must define the peer-to-peer network, by which clients can find and communicate with each other, and the consensus protocol, by which they can come to agreement on the transactions that have taken place.

Beyond the relatively simple functionality offered by Bitcoin, Ethereum allows developers to create and deploy *smart contracts* onto the blockchain. These are stateful programs, typically written in a language called Solidity, that can be triggered by transactions and used to run (almost) arbitrary code. The only limitation is their complexity, as every operation they perform has an associated cost, and transactions have a maximum amount they are allowed to spend.

There are two special types of smart contracts, ERC20 and ERC721, that define *tokens*: ERC20 tokens⁴ are designed to be currency, and thus are fungible,

⁴ <https://github.com/ethereum/eips/issues/20>

Category	# coins	Representative examples
Animals	29	RabbitCoin, Birds
Computing	47	AI Doctor, Decentralized Machine Learning
Drugs	13	Cannation, KushCoin
Finance	22	iBank, Intelligent Investment Chain
Food	10	EggCoin, Olive
Gambling	17	CashBet Coin, CasinoCoin
Geography	12	Asiadigicoin, NewYorkCoin
Nation states	37	PutinCoin, Shekel, BritCoin
Outer space	17	Marscoin, SpaceChain
Metals & precious stones	37	GoldPieces, PlatinumBAR, Rubycoin
Religion	7	BiblePay, HalalChain

Table 1: Different categories of cryptocurrencies, based on the name of the coin, along with some representative examples.

whereas ERC721 tokens may be non-fungible and thus support collectibles such as CryptoKitties.⁵ At the most basic level, a token contract is a ledger mapping owners (identified by their Ethereum address) to the amount of tokens that they own, along with an associated set of rules determining how tokens are transferred between owners.

4 Data Collection

In order to collect the source code associated with each cryptocurrency, we started with the list maintained at CoinMarketCap, which is generally regarded as one of the most comprehensive resources for cryptocurrency market data. The site maintains not only market data for each cryptocurrency (its market capitalization, price, circulating supply, etc.), however, but also links to any websites, blockchain explorers, or — crucially for us — source code repositories. We last scraped the site on July 24 2018, at which point there were 1664 cryptocurrencies listed, with a cumulative market capitalization of 293B USD.

Of these cryptocurrencies, there were 866 categorized as a token rather than as a coin. There were 366 cryptocurrencies with a stated price of \$0.00, and in fact 1468 (88%) had a stated price of less than \$1.00. There were 276 cryptocurrencies with an unknown circulating supply (and thus an unknown market capitalization), and 924 (55.5%) had a market capitalization of over 1M USD. The names of the cryptocurrencies were typically designed to evoke a specific concept; e.g., wealth, computing, politicians, or existing fiat currencies. Based on their names, we manually sorted all of the listed cryptocurrencies into the categories shown in Table 1.

⁵ <https://www.cryptokitties.co/>

4.1 Source code repositories

Of the listed cryptocurrencies, 1123 had a link available to some source code repository. We performed manual spot checking to ensure that the links were legitimate, and in some cases replaced the links where the information was inaccurate (for Bitcoin Cash, for example, the provided link was for the repositories backing `bitcoincash.org` rather than the actual software code). Of these, 1108 pointed to GitHub (98.7%).

As should be expected, many of the cryptocurrencies had multiple software repositories available; indeed, the links provided were to the lists of repositories for a given GitHub organization, and in total there were 13,694 individual repositories available. The vast majority of these repositories had been created after October 2014, with a notable rise in frequency starting in April 2017. These repositories typically fell into one of three categories: (1) integral to the cryptocurrency itself, such as implementations of the reference client or supporting libraries; (2) irrelevant, such as a different project by the same organization; or (3) unchanged forks or mirrors of popular software projects, such as `llvm`. Given our goal of differentiating between different cryptocurrencies, we sought to isolate the first category of “meaningful” code.

Our initial hypothesis was that more integral repositories would be (1) better maintained, (2) more popular, and (3) re-used more frequently. There are several open-source tools for determining the “health” of the maintenance of a GitHub repository [24], which typically measure the activity, level of contribution, and responsiveness to pull requests and issues. In fact, there is even a tool called `CoinCheckup` that does this specifically for cryptocurrencies.⁶ We observed, however, that many of the repositories had a fairly low level of activity, so this was not on its own a good way to distinguish between different repositories. Another natural measure would be the total number of commits to a repository, but this information is not readily accessible via the GitHub API.⁷ We thus measured activity according to the gap in time (in weeks) between the current time and the last update of the repository, under the assumption that this would be shorter for more active repositories.

In terms of the second two criteria, the most natural representation of popularity and reusability is the number of stars and forks [23]. We again found, however, that stars were not a useful indication of whether or not a repository contained relevant source code. For example, important meta-information related to a cryptocurrency such as its whitepaper or its improvement proposals (IPs) was often contained in highly-starred repositories. We thus chose to ignore stars and look solely at forks.

Finally, in terms of relevance, we favored repositories with a name identical or close to that of the cryptocurrency, as well as ones that indicated they contained source code (e.g., with ‘core’ in the name) or code otherwise relevant to the operation of the cryptocurrency (e.g., ‘token’ or ‘contract’ for ERC20 tokens).

⁶ <https://coincheckup.com/analysis/github>

⁷ It is possible to derive it from information given for individual repositories, albeit in a relatively inefficient manner, but not from an organization’s list of repositories.

We also excluded repositories whose names contained terms that indicated they were not relevant; e.g., ‘website’ or ‘docs.’ A complete list of these excluded terms can be found in Table 3 in Appendix A.

In the end, we assigned a rating to each repository for a given cryptocurrency according to: (1) the gap between its last update and the current date, to capture activity (where this was subtracted from the rating, as a longer gap indicates less activity); (2) its number of forks, to capture popularity and reuse; and (3) information about the name of the repository, to capture relevance. For each cryptocurrency, we then cloned the top 20% of the list of repositories, sorted from high to low by these ratings (or cloned one repository, whichever was larger). Even after compiling this list, we made various manual adjustments in order to ensure that we had selected the “right” repositories. We cloned 2354 repositories in total, which comprised roughly 100 GB of data.

4.2 Deployed source code

As evidenced by the 866 (52%) listed cryptocurrencies that were categorized as tokens (and the fact that 74 of these even had ‘token’ in their name), it is popular to launch new cryptocurrencies not as standalone coins, but as tokens that are supported by existing cryptocurrencies. Of these, by far the most popular type is an ERC20 token, supported by Ethereum. Of these listed tokens, 406 did not have any source code link available. For ERC20 tokens that have been deployed, however, it is often possible to obtain the contract code from another source: the version deployed on the Ethereum blockchain itself is compiled bytecode, but it is common practice to provide the Solidity code and display it on blockchain explorers such as Etherscan.⁸

For these tokens, we thus chose to use Etherscan as a data source (in addition to any provided repositories), in order to aid our Ethereum-based analysis in Section 7. At the time that we scraped Etherscan, there were 612 ERC20 tokens listed, identified by a name and a currency symbol (e.g., OmiseGO and OMG). Of these, we found 438 with a match on CoinMarketCap, where we defined a match as having (1) identical currency symbols, and (2) closely matching names. (We couldn’t also require the name to be identical because in some cases the name of the contract was somewhat altered from the name of the cryptocurrency; e.g., SPANK instead of SpankChain.) We scraped the available contract code for each of these tokens, which in all but 9 cases was Solidity code rather than just on-chain bytecode. We thus ended up with 429 deployed ERC20 contracts.

5 Detecting Code Reuse

In this section, we attempt to identify the extent to which cryptocurrencies reuse the codebases of others, in order to identify which cryptocurrencies incorporate genuinely novel ideas and which ones largely derive their ideas from others. We

⁸ <https://etherscan.io/>

do so using several different techniques, ranging from very simplistic (seeing if the name of one is a prefix of the name of another) to more complex. We refer to the cryptocurrencies borrowing ideas from others as *derivatives*.

Most of our methods label one cryptocurrency as a derivative of another if it has at least one repository that appears as derived from a repository of the second cryptocurrency. This means that, for example, if a cryptocurrency copied the Bitcoin repository but its actual platform consists largely of repositories written from scratch, we may unfairly label it as a derivative of Bitcoin. On the other hand, if we tried to do otherwise then we might unfairly reward cryptocurrencies that copy the Bitcoin repository and then create many other repositories that do not contribute to the functioning of the platform (e.g., web templates). Ultimately, without significantly more advanced analysis to determine how repositories are linked and which ones meaningfully support the platform, this is a potential limitation that we must accept in our analysis that follows.

5.1 Name derivations

As a first simple method, we observed that many cryptocurrencies attempt to profit from the name recognition of popular cryptocurrencies by using a similar name; e.g., Bitcoin Planet or Ethereum Gold. We thus decided to identify one cryptocurrency as a derivative of another cryptocurrency if the name of the second is a prefix of the name of the first (as in the examples above), manually excluding cryptocurrencies whose names are common words and thus might be prefixes anyway (e.g., Crypto). Using this method, we identified 28 cryptocurrencies with one derivative, and 8 with two derivatives (0x, Aurora, Dynamic, Hyper, Litecoin, Monero, Sentinel, and Waves). While they did not exactly match our prefix method, there were four cryptocurrencies whose names seem designed to evoke the ideas behind the popular “privacy coin” Zcash: ZClassic, ZCoin, Zoin, and Zero. Finally, the two most derived cryptocurrencies are — unsurprisingly — Bitcoin (with 17 derivatives) and Ethereum (with 6). While this method already yields some interesting insights into the extent to which the Bitcoin brand has been borrowed, the name of a cryptocurrency is not necessarily indicative of the contents of its underlying codebase. We thus continue to develop more code-specific methods for determining derivatives.

5.2 Git forks

Next, we considered the Git forks of a cryptocurrency, meaning the cryptocurrencies that were created as a fork of the GitHub repository of another cryptocurrency. To do this, we identified (by hash) every commit to every repository we cloned. We then mapped commits to the lists of the repositories in which they appeared, and considered the oldest repository containing that commit to be the original cryptocurrency and all other repositories in the list to be derivatives. To elevate this to the level of cryptocurrencies, we labelled a cryptocurrency as a derivative of another one if any of its repositories were forked from any repositories of the second one. The results are in Figure 1. The most immediate

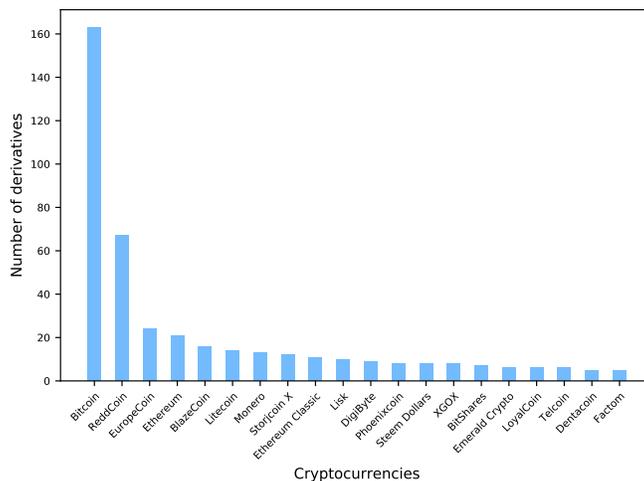


Fig. 1: The cryptocurrencies with three or more derivatives, where a derivative is defined as a cryptocurrency that has a repository initially created by the first one.

observation is that this method captures significantly more derivatives than the name-based one, and that it reinforces the popularity of Bitcoin and Ethereum (according to this method, Bitcoin has 163 forks and Ethereum has 21), as well as of Litecoin and Monero.

There are also, however, several cryptocurrencies with many Git forks that are less well-known; e.g., Reddcoin and EuropeCoin. These are due to being the earliest cryptocurrencies to incorporate independently popular repositories, and thus highlight the main limitation of this method; namely, that the earliest cryptocurrency to fork a popular general-purpose library or other integral software development tool will be incorrectly labelled as popular. In addition, it doesn't consider "intermediate" derivatives; e.g., Litecoin and Peercoin are both forks of Bitcoin that have themselves been forked many times.

To get a sense of how much this method underestimated code reuse, we considered the file `addrman.cpp`, which is the way addresses are managed in Bitcoin's peer-to-peer network and is one of the most reused cryptocurrency-specific files we saw. It appeared in a repository for 536 other cryptocurrencies, meaning there are many potential derivatives of Bitcoin that this method does not capture.

5.3 Copyright derivations

Most open-source repositories are not written from scratch, but instead make use of established libraries or other code. This is expected (and in fact encouraged), as long as the authors acknowledge the original authors. We thus decided next to use this copyright information in order to identify potential derivatives.

More specifically, we looked in every repository to see if any `COPYING` files were available. If they were, then we scraped the comments concerning copyright

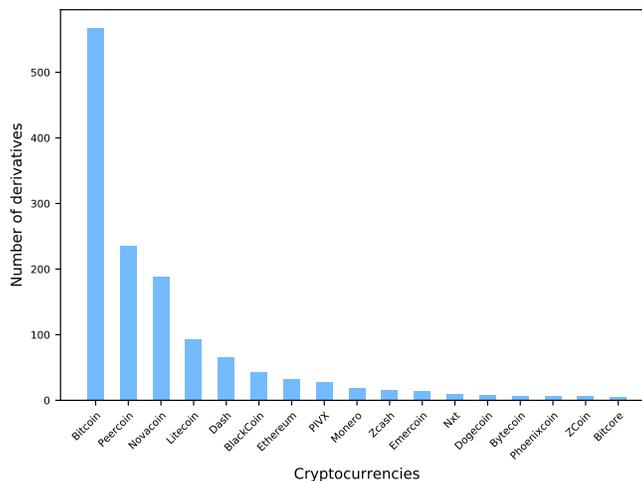


Fig. 2: The cryptocurrencies with five or more derivatives, where a derivative is defined as a cryptocurrency that contains information copyrighted by the first one.

from the beginning of these files. If not, then we went through every source code file in the repository and scraped the copyright information there instead, where we identified source code according to the file extensions maintained by the CLOC (Count Lines of Code) library.⁹ Either way, we then considered the collective authors of the repository to be the union of all of the individual authors identified by the copyright lines.

Given this set of collective authors, our next task was to assign them to a specific cryptocurrency. Luckily, many of the cryptocurrencies do not identify contributors by their individual names, but rather by the name of the cryptocurrency; e.g., “Bitcoin Developers” or “The go-ethereum Authors”. To handle the several prominent exceptions to this rule, we manually created a mapping from popular individual contributors to the coins with which they were associated; e.g., Pieter Wuille for Bitcoin. This also covered cases in which the name of the coin was altered slightly (e.g., PPCoin for Peercoin) and in which the copyright was given to the organization rather than the cryptocurrency (e.g., IOHK for Cardano). We did not include contributors to popular open-source libraries (e.g., Boost and LevelDB) in order to keep our focus on cryptocurrency-specific code rather than standard software development tools.

With this information, we then labelled one cryptocurrency as a derivative of another if the copyright lines in the first included the name of the second (or the name of one of its popular contributors). The results are in Figure 2. This graph again demonstrates the dominance of Bitcoin, as well as of several cryptocurrencies (most notably Litecoin, Peercoin, and Novacoin) that have also been very popular to fork. It also comes much closer to the expected number of derivatives given the occurrence of `addrman.cpp` discussed in Section 5.2.

⁹ <https://github.com/AlDanial/cloc>

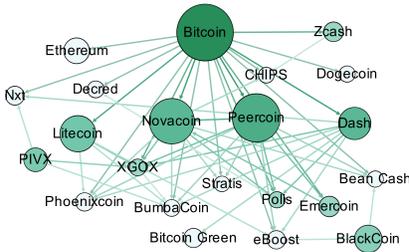


Fig. 3: In the largest connected component of the graph of copyright derivations, the cryptocurrencies with degree over 7, where the size (and color) of the node is proportional to its degree.

In addition to capturing a richer type of derivation, this approach also had the benefit over the previous ones that it did not flatten transitive relationships; i.e., if X forked from Y, which was itself a fork of Z, then X would be labelled as a derivative of both Y and Z, rather than just Z. It thus builds a tree of derivations, such as that seen on the (currently unmaintained) website Map of Coins.¹⁰ A version of this tree with the most derived (or derivative) cryptocurrencies can be seen in Figure 3. This graph reinforces the results from Figure 2, in terms of which cryptocurrencies have the most derivatives. It also highlights cryptocurrencies like BumbaCoin, which incorporate code from many different (10) sources.

While this method nicely captures code reuse, it goes perhaps a bit too far in labelling derivatives; again, code reuse is encouraged as a form of library support and using or modifying copyrighted code from another cryptocurrency does not indicate a lack of other novel ideas. Thus, our final method considers not just whether or not a cryptocurrency uses the codebase of another, but the actual proportion of code that goes unmodified.

5.4 File derivations

Finally, we looked for the most direct form of derivation: taking another repository and using it without any modification. To identify this, we computed and stored the hash of every source code file in our cloned repositories; as in Section 5.3, we identified source code file extensions using the CLOC library. We then computed a similarity score S_{hash} between a repository A and another one B by counting the number of files in A with an identical file in B (meaning the hash was the same), and then dividing by the total number of files in A . To elevate this to the level of cryptocurrencies C_1 and C_2 , we then computed $S_{\text{hash}}(C_1, C_2)$ as

$$S_{\text{hash}}(C_1, C_2) = \frac{\sum_{A \in C_1} S_{\text{hash}}(A, \cup_{B \in C_2} B)}{\sum_{A \in C_1} \# \text{ files in } A};$$

¹⁰ <http://mapofcoins.com>

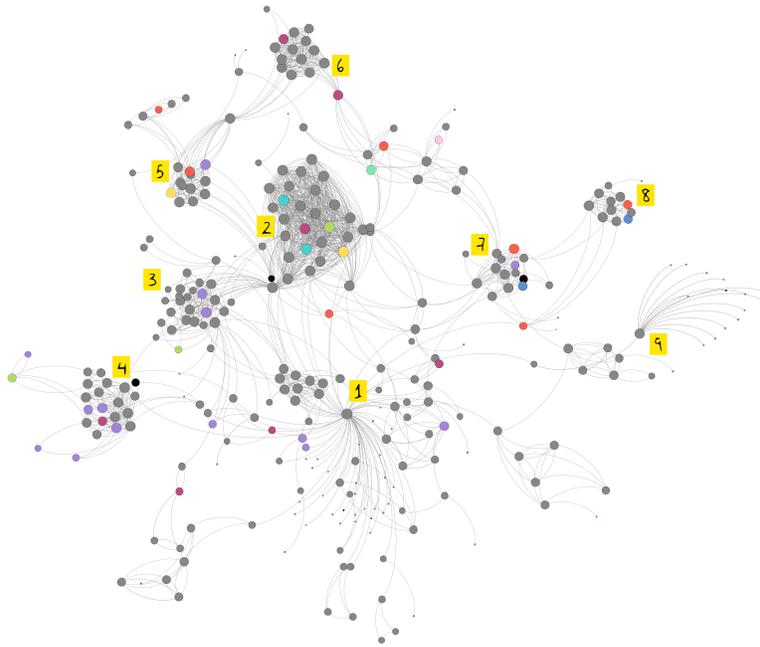


Fig. 4: The largest connected component of the graph formed by creating an edge from A to B if $S_{\text{hash}}(A, B) > 0.7$, along with labels for the most prominent clusters. The nodes are colored according to the categories from Table 1.

i.e., for each repository A contributing to C_1 we counted the number of files that were identical to a file in *any* repository contributing to C_2 , and then divided this by the total number of files across all repositories contributing to C_1 .

We ran this for every pair of cryptocurrencies A and B (for both $S_{\text{hash}}(A, B)$ and $S_{\text{hash}}(B, A)$, since they are not symmetric), and used the results to create a graph in which nodes represent cryptocurrencies and there is a directed edge from A to B if $S_{\text{hash}}(A, B) > 0.7$. This resulted in a graph with 445 nodes and 1854 edges, the largest connected component of which can be seen in Figure 4 (consisting of 302 nodes and 1599 edges).

Most of this component consists of Bitcoin forks, so we defer further discussion of these clusters to Section 6. The exception is cluster 9, which consists of one cryptocurrency (Zeepin) that is 100% similar to 16 other cryptocurrencies. The reason is simple: its repository consisted solely of an LGPL-3.0 license, so it matched other repositories with the same version of this license. At the time we scraped CoinMarketCap, Zeepin had a market capitalization of 23 million USD.

6 Bitcoin and Its Derivatives

In Section 5, we saw a clear dominance of Bitcoin in terms of the reuse of its code by other cryptocurrencies. We now explore these Bitcoin forks in more detail, by comparing them not only against one version of the Bitcoin codebase, but against

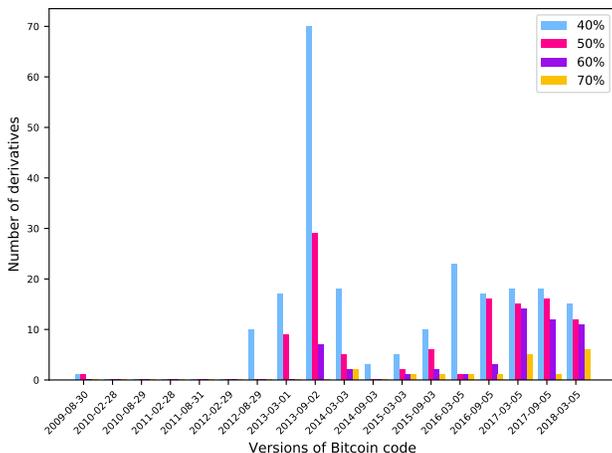


Fig. 5: For different percentage overlaps, the number of cryptocurrencies matching a given version of the Bitcoin codebase in terms of the hash similarity.

multiple versions reflecting its evolution over time. To do this, we first created a list of all commits to the Bitcoin repository, starting from August 2009, and collected one version of the codebase every six months after that. We labelled these versions from 1 to 18, but stress that these labels are not correlated with any official releases. We then re-ran the hash similarity code from Section 5.4 against each of these versions, assigning a cryptocurrency to a single version by picking the one with which it had the highest overlap. The results are in Figure 5, and demonstrate an increase in the number of derivatives over time, with a clear spike at version 9 (which represents the codebase in September 2013).

Hashing is a very brittle method of comparison, however, and what we saw copied far more often than the exact file contents was Bitcoin’s directory structure. We thus also compared the directory structure of two repositories; i.e., we computed a similarity score S_{dir} between a repository A and another one B by counting the number of files in A with an identical filename in B (meaning the name and path was the same), and then dividing by the total number of files in A . We elevated this to the level of cryptocurrencies in the same way as we did in Section 5.4 for S_{hash} . For completeness, the results are in Figure 7 in Appendix A, and demonstrate that the directory structure (and in particular the one associated with older versions of the codebase, prior to its re-organization due to SegWit adoption) is heavily borrowed by other cryptocurrencies.

Finally, we revisit the graph in Figure 4, as this connected component is almost entirely associated with Bitcoin forks. In particular, we can briefly explain clusters 1-8 as follows:

- **1.** The node at the center of this cluster, Akuya Coin, has a directory structure similar (63%) to a version of the Bitcoin codebase from 2013, but many (32%) of its files are empty and thus have the same hash, which makes it appear similar to 76 other Bitcoin forks.

- **2 and 3.** Both of these clusters also have a directory structure similar to older versions of the Bitcoin codebase (the average directory similarity was 89% for cluster 2 and 82% for cluster 3), and are similar to the same cryptocurrency (BumbaCoin). Many also incorporate the Zerocoin code:¹¹ 84% of the nodes in cluster 2 and 65% of the nodes in cluster 3. This is notable given that this code comes with the emphatic warning “THIS CODE IS UNMAINTAINED AND HAS KNOWN EXPLOITS. DO NOT USE IT.” In total it is included in repositories for 97 different cryptocurrencies.
- **4 and 5.** These clusters were the ones most similar to Bitcoin: on average we had $S_{\text{hash}} = 0.51$ and $S_{\text{dir}} = 0.80$ for cluster 4 and $S_{\text{hash}} = 0.37$ and $S_{\text{dir}} = 0.96$ for cluster 5. For cluster 4, the matching versions were also in quite a tight range from September 2013 to September 2014 (our versions 9 to 11), whereas most other clusters ranged more evenly across all 18 versions.
- **6 and 7.** These clusters consisted largely of forks from Litecoin: 100% of cluster 6 had the file `script.c`, which is unique to Litecoin, and indeed 100% identified as Litecoin derivatives using the copyright method from Section 5.3. 64% of cluster 7 had files with `script` in the name, although only 21% identified as copyright derivatives of anything other than Bitcoin.
- **8.** The nodes in this cluster were on average newer than the others (with the first repository created in June 2015), and indeed their directory structure is more consistent with newer versions of the Bitcoin codebase.

7 Ethereum as a Platform

As discussed in Sections 3 and 4.2, it is increasingly popular to deploy cryptocurrencies as tokens on the Ethereum blockchain; indeed, over half of the cryptocurrencies listed on CoinMarketCap fell into this category. This section thus explores this type of cryptocurrency deployment, focusing again on the extent to which ERC20 tokens are similar to or different from each other. As an ERC20 token consists of just a single file, our methods from the previous sections do not apply here so we develop new methods for identifying similarities.

The basic functionality of an ERC20 token—allowing the transfer of tokens from one holder to another—defines a contract type called `Basic` (or `BasicToken`) or—with one slight functional difference—`ERC20`. There are, however, many additional types that ERC20 tokens can have. For example, if they want to allow for the creation of new tokens they can be `Mintable` and if they want to allow for the destruction of existing tokens they can be `Destructible` or `Burnable`. These types are not standardized, and in fact new types can be defined and used within the Solidity code for a contract.

To identify the types of a given token, we identified all lines in its contract of the form `contract X is Y {`, where `X` is the name of the contract and `Y` is its type. Some intermediate types themselves appear as names (e.g., `contract Mintable is Ownable`), which we exclude from our final results but carry over transitively to the higher-level contract names; e.g., if `X` is `Mintable`

¹¹ <https://github.com/Zerocoin/libzerocoin>

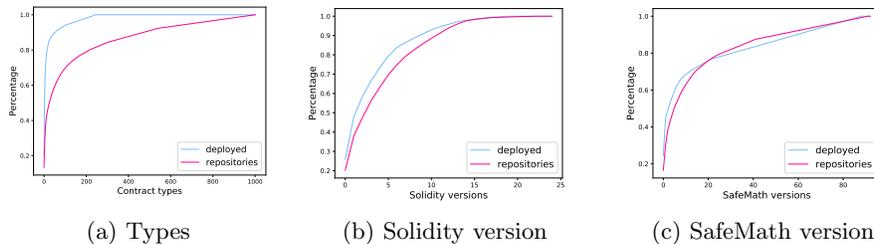


Fig. 6: When ranked from most to least popular, the cumulative percentage of contracts matching three different features, for both the set of deployed contracts and the ones found in repositories.

and `Mintable` is `Ownable` then `X` is both `Mintable` and `Ownable`. This resulted in a map from the higher-level token names to a list of all of their types.

Beyond these types, there are several ERC20 templates whose reuse we sought to identify. Based on a manual inspection of a random subset of contracts, we chose five points of distinction: first, which version the contract used of (1) Solidity and (2) the SafeMath library, which provides safe arithmetic operations. Second, we considered whether or not it used (3) the definition of `StandardToken` created by the FirstBlood token; (4) the `UpgradeableToken` type, created by the Golem token; or (5) a template by OpenZeppelin,¹² who also created SafeMath.

For the version of Solidity, we looked for lines starting with `pragma solidity` and extracted the version from what followed (typically of the form `0.4.X`). To determine the version of SafeMath, we first used CLOC to strip the comments from the `.sol` file. We then identified the lines of code that defined the SafeMath library (starting with either `contract SafeMath {` or `library SafeMath {` and ending with `}`), and hashed this substring to form a succinct representation. Finally, to identify the use of templates, we simply looked at whether or not the file explicitly mentioned the relevant keywords.

We extracted this information from all Solidity files, whether deployed on the Ethereum blockchain (and thus scraped from Etherscan, as described in Section 4.2) or contained in a repository.¹³ For the types, Solidity and SafeMath versions, we ordered them from most to least popular and plotted this as a CDF, as seen in Figure 6; i.e., we plotted the percentage y of all contracts that had one of the top x attributes. For the templates, the results are in Table 2.

The relatively long tails in all of the figures indicate a relatively high level of diversity among these features in both deployed contracts and those still under development. For example, the Solidity version most popular among deployed contracts (version 18) was still used in only 23% of them. Whereas Figures 6b and 6c show similar curves for both sets of contracts, Figure 6a shows a much longer tail for contracts contained in repositories, with 246 distinct types in

¹² <https://openzeppelin.org/>

¹³ Interestingly, these sets were non-intersecting; i.e., there was no contract in a repository that was identical to a deployed one.

	Deployed		Repositories	
	#	%	#	%
FirstBlood	134	30.8	190	2.5
Lunyr	22	5.1	19	0.2
OpenZeppelin	72	16.6	245	3.2
SafeMath	287	65.9	400	5.2

Table 2: The number of contracts of each type (deployed or in repositories) that are derived from one of the first three templates, or using the SafeMath library, as well as the percentage of all contracts this represents.

deployed contracts and 1002 in ones in repositories. This indicates — as should perhaps be expected — that (1) there are just many more possibilities for contract types than for versions, and (2) there is greater experimentation with types in contracts still under development. Even among deployed contracts, 129 out of 429 had a type that did not appear in any other deployed contracts, and 148 of the 246 distinct types appeared in only a single contract.

Table 2 also demonstrates the relatively high diversity across contracts, with no one template being used in a dominant way. Even though 65.9% of deployed contracts use SafeMath, Figure 6c demonstrates that there is quite a lot of variety within this category; indeed, there were 90 different versions of SafeMath used in deployed contracts (and 93 different versions in the repositories). Again, we see significantly more experimentation in contracts still under development.

Finally, we view the points of similarity that did exist as operating primarily in support of the safety of deployed contracts. For example, among the 20 most popular types across both deployed and repository contracts, five of them defined the basic ERC20 functionality, and six of them were related to safety in terms of either including a standard library or in defining an owner who could take action if something went wrong. The same is true of the usage of FirstBlood’s `StandardToken`, which was the first safe implementation of this type, or of the `SafeMath` library. We thus view these similarities as a sign of good development practices, rather than the copying of ideas.

8 Conclusions

This paper considered diversity in the cryptocurrency landscape, according to the source code available for each one, in order to identify the extent to which new cryptocurrencies provide meaningful innovation. This was done by examining the source code for over a thousand cryptocurrencies, and — in the case of ERC20 tokens — the deployed code of hundreds more. While more sophisticated static analysis of the source code would likely yield further insights, even our relatively coarse methods clearly indicated the dominance of Bitcoin and Ethereum, as well as the extent to which creating a standalone platform is a significantly greater undertaking (leading to the reuse of much of the Bitcoin codebase) than defining just the transaction semantics of an Ethereum-based token.

Acknowledgements

The authors were supported in part by EPSRC Grant EP/N028104/1 and in part by the EU H2020 TITANIUM project under grant agreement number 740558.

References

1. M. Apostolaki, A. Zohar, and L. Vanbever. Hijacking bitcoin: Routing attacks on cryptocurrencies. In *2017 IEEE Symposium on Security and Privacy*, pages 375–392, San Jose, CA, USA, May 22–26, 2017. IEEE Computer Society Press.
2. S. Azouvi, M. Maller, and S. Meiklejohn. Egalitarian society or benevolent dictatorship: The state of cryptocurrency governance. In *Proceedings of the 5th Workshop on Bitcoin and Blockchain Research*, 2018.
3. A. Biryukov, D. Khovratovich, and I. Pustogarov. Deanonymisation of clients in bitcoin P2P network. In G.-J. Ahn, M. Yung, and N. Li, editors, *ACM CCS 14*, pages 15–29, Scottsdale, AZ, USA, Nov. 3–7, 2014. ACM Press.
4. J. Bonneau, A. Miller, J. Clark, A. Narayanan, J. A. Kroll, and E. W. Felten. SoK: Research perspectives and challenges for bitcoin and cryptocurrencies. In *2015 IEEE Symposium on Security and Privacy*, pages 104–121, San Jose, CA, USA, May 17–21, 2015. IEEE Computer Society Press.
5. J. W. Bos, J. A. Halderman, N. Heninger, J. Moore, M. Naehrig, and E. Wustrow. Elliptic curve cryptography in practice. In N. Christin and R. Safavi-Naini, editors, *FC 2014*, volume 8437 of *LNCS*, pages 157–175, Christ Church, Barbados, Mar. 3–7, 2014. Springer, Heidelberg, Germany.
6. B. David, P. Gazi, A. Kiayias, and A. Russell. Ouroboros praos: An adaptively-secure, semi-synchronous proof-of-stake blockchain. In J. B. Nielsen and V. Rijmen, editors, *EUROCRYPT 2018, Part II*, volume 10821 of *LNCS*, pages 66–98, Tel Aviv, Israel, Apr. 29 – May 3, 2018. Springer, Heidelberg, Germany.
7. J. A. D. Donet, C. Pérez-Solà, and J. Herrera-Joancomartí. The bitcoin P2P network. In R. Böhme, M. Brenner, T. Moore, and M. Smith, editors, *FC 2014 Workshops*, volume 8438 of *LNCS*, pages 87–102, Christ Church, Barbados, Mar. 7, 2014. Springer, Heidelberg, Germany.
8. I. Eyal, A. E. Gencer, E. G. Sirer, and R. van Renesse. Bitcoin-ng: A scalable blockchain protocol. In *Proceedings of NSDI 2016*, 2016.
9. A. E. Gencer, S. Basu, I. Eyal, R. V. Renesse, and E. G. Sirer. Decentralization in Bitcoin and Ethereum networks. In *Proceedings of the 22nd International Conference on Financial Cryptography and Data Security (FC)*, 2018.
10. A. Gervais, G. O. Karame, K. Wüst, V. Glykantzis, H. Ritzdorf, and S. Capkun. On the security and performance of proof of work blockchains. In E. R. Weippl, S. Katzenbeisser, C. Kruegel, A. C. Myers, and S. Halevi, editors, *ACM CCS 16*, pages 3–16, Vienna, Austria, Oct. 24–28, 2016. ACM Press.
11. Y. Hu, J. Zhang, X. Bai, S. Yu, and Z. Yang. Influence analysis of GitHub repositories. *SpringerPlus*, 5(1), 2016.
12. D. Y. Huang, K. Levchenko, and A. C. Snoeren. Measuring profitability of alternative crypto-currencies. In *Proceedings of the 22nd International Conference on Financial Cryptography and Data Security (FC)*, 2018.
13. G. Kappos, H. Yousaf, M. Maller, and S. Meiklejohn. An empirical analysis of anonymity in Zcash. In *Proceedings of the USENIX Security Symposium*, 2018.

14. A. Kiayias, A. Russell, B. David, and R. Oliynykov. Ouroboros: A provably secure proof-of-stake blockchain protocol. In J. Katz and H. Shacham, editors, *CRYPTO 2017, Part I*, volume 10401 of *LNCS*, pages 357–388, Santa Barbara, CA, USA, Aug. 20–24, 2017. Springer, Heidelberg, Germany.
15. E. K. Kogias, P. Jovanovic, N. Gailly, I. Khoffi, L. Gasser, and B. Ford. Enhancing Bitcoin security and performance with strong consistency via collective signing. In *Proceedings of USENIX Security 2016*, 2016.
16. E. Kokoris-Kogias, P. Jovanovic, L. Gasser, N. Gailly, and B. Ford. Omniledger: A secure, scale-out, decentralized ledger. In *Proceedings of the 39th IEEE Symposium on Security & Privacy*, 2018.
17. P. Koshy, D. Koshy, and P. McDaniel. An analysis of anonymity in bitcoin using P2P network traffic. In N. Christin and R. Safavi-Naini, editors, *FC 2014*, volume 8437 of *LNCS*, pages 469–485, Christ Church, Barbados, Mar. 3–7, 2014. Springer, Heidelberg, Germany.
18. L. Luu, V. Narayanan, C. Zheng, K. Baweja, S. Gilbert, and P. Saxena. A secure sharding protocol for open blockchains. In E. R. Weippl, S. Katzenbeisser, C. Kruegel, A. C. Myers, and S. Halevi, editors, *ACM CCS 16*, pages 17–30, Vienna, Austria, Oct. 24–28, 2016. ACM Press.
19. S. Meiklejohn, M. Pomarole, G. Jordan, K. Levchenko, D. McCoy, G. M. Voelker, and S. Savage. A fistful of bitcoins: characterizing payments among men with no names. In *Proceedings of the 2013 Internet Measurement Conference (IMC)*, pages 127–140, 2013.
20. T. Moore and N. Christin. Beware the middleman: Empirical analysis of Bitcoin-exchange risk. In A.-R. Sadeghi, editor, *FC 2013*, volume 7859 of *LNCS*, pages 25–33, Okinawa, Japan, Apr. 1–5, 2013. Springer, Heidelberg, Germany.
21. M. Möser, K. Soska, E. Heilman, K. Lee, H. Heffan, S. Srivastava, K. Hogan, J. Hennessey, A. Miller, A. Narayanan, and N. Christin. An empirical analysis of linkability in the Monero blockchain. *Proceedings on Privacy Enhancing Technologies*, 2016(3):143–163, 2018.
22. S. Nakamoto. Bitcoin: A peer-to-peer electronic cash system, 2008. bitcoin.org/bitcoin.pdf.
23. M. Papamichail, T. Diamantopoulos, and A. Symeonidis. User-perceived source code quality estimation based on static analysis metrics. In *2016 IEEE International Conference on Software Quality, Reliability, and Security (QRS)*, 2016.
24. pelson. Github repository health report. <https://github.com/pelson/repohealth.info>.
25. B. Ray, D. Posnett, V. Filkov, and P. Devanbu. A large scale study of programming languages and code quality in github. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE)*, pages 155–165, 2014.
26. F. Reid and M. Harrigan. An analysis of anonymity in the Bitcoin system. In *Security and privacy in social networks*, pages 197–223. Springer, 2013.
27. D. Ron and A. Shamir. Quantitative analysis of the full Bitcoin transaction graph. In A.-R. Sadeghi, editor, *FC 2013*, volume 7859 of *LNCS*, pages 6–24, Okinawa, Japan, Apr. 1–5, 2013. Springer, Heidelberg, Germany.
28. M. Spagnuolo, F. Maggi, and S. Zanero. BitIodine: Extracting intelligence from the bitcoin network. In N. Christin and R. Safavi-Naini, editors, *FC 2014*, volume 8437 of *LNCS*, pages 457–468, Christ Church, Barbados, Mar. 3–7, 2014. Springer, Heidelberg, Germany.

29. E. Syta, I. Tamas, D. Visher, D. I. Wolinsky, P. Jovanovic, L. Gasser, N. Gailly, I. Khoffi, and B. Ford. Keeping authorities “honest or bust” with decentralized witness cosigning. In *2016 IEEE Symposium on Security and Privacy*, pages 526–545, San Jose, CA, USA, May 22–26, 2016. IEEE Computer Society Press.
30. F. Thung, T. F. Bissyandé, D. Lo, and L. Jiang. Network structure of social coding in GitHub. In *Proceedings of the 17th European Conference on Software Maintenance and Reengineering*, 2013.
31. M. Vasek, J. Bonneau, R. Castellucci, C. Keith, and T. Moore. The bitcoin brain drain: Examining the use and abuse of bitcoin brain wallets. In J. Grossklags and B. Preneel, editors, *FC 2016*, volume 9603 of *LNCS*, pages 609–618, Christ Church, Barbados, Feb. 22–26, 2016. Springer, Heidelberg, Germany.
32. M. Vasek and T. Moore. There’s no free lunch, even using bitcoin: Tracking the popularity and profits of virtual currency scams. In R. Böhme and T. Okamoto, editors, *FC 2015*, volume 8975 of *LNCS*, pages 44–61, San Juan, Puerto Rico, Jan. 26–30, 2015. Springer, Heidelberg, Germany.
33. M. Vasek and T. Moore. Analyzing the Bitcoin Ponzi scheme ecosystem. In *Proceedings of the 5th Workshop on Bitcoin and Blockchain Research*, 2018.
34. M. Vasek, M. Thornton, and T. Moore. Empirical analysis of denial-of-service attacks in the bitcoin ecosystem. In R. Böhme, M. Brenner, T. Moore, and M. Smith, editors, *FC 2014 Workshops*, volume 8438 of *LNCS*, pages 57–71, Christ Church, Barbados, Mar. 7, 2014. Springer, Heidelberg, Germany.
35. J. Zhu, M. Zhou, and A. Mockus. Patterns of folder use and project popularity: A case study of Github repositories. In *Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, 2014.

A Extra Tables and Figures

Concepts	Excluded terms
Binaries	binaries
Documentation	docs, document, papers, whitepaper, wiki
Improvement proposals (IPs)	[a-z]1,2ips
Other	electrum, explorer, faucet, vanitygen
Packaging	docker, homebrew, install
Testing	example, test
Tools	kit, lib, plugin, sdk, service, tools
User interface	android, gui, ios, macos, mobile, window
Website	.com, .info, .io, .net, .org, -org, site, website, www

Table 3: The terms which, if we encountered them in the name of a repository, led to the exclusion of that repository from further consideration. We added manual exceptions to these where appropriate (e.g., allowing the ‘ips’ pattern for the CHIPS and Vipstar Coin cryptocurrencies).

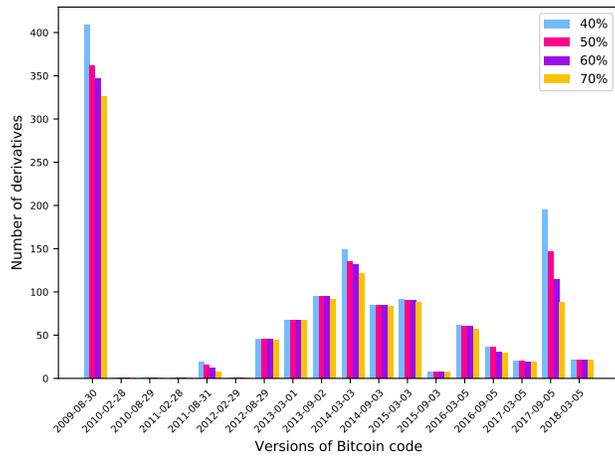


Fig. 7: For different percentage overlaps, the number of cryptocurrencies matching a given version of the Bitcoin codebase in terms of the directory similarity.