

# Technology and Cryptocurrency Valuation: Evidence from Machine Learning\*

Yukun Liu    Jinfei Sheng    Wanyi Wang

January 2021

## Abstract

This paper studies the role of technological sophistication in Initial Coin Offering (ICO) successes and valuations. Using various machine learning methods, we construct technology indexes from ICO whitepapers to capture technological sophistication for all cryptocurrencies. We find that the cryptocurrencies with high technology indexes are more likely to succeed and less likely to be delisted subsequently. Moreover, the technology indexes strongly and positively predict the long-run performances of the ICOs. Overall, the results suggest that technological sophistication is an important determinant of cryptocurrency valuations.

*Keywords:* Cryptocurrency, Technological Sophistication, Machine Learning, Textual Analysis, FinTech, Blockchain.

---

\*Yukun Liu is with Simon School of Business at University of Rochester. Jinfei Sheng (Corresponding Author) is with Merage School of Business at University of California Irvine (Email: jinfei.sheng@uci.edu). Wanyi Wang is with Merage School of Business at University of California Irvine. For helpful comments, we thank David Hirshleifer, Chong Huang, Arthur Inuma, Alan Kwan, Ye Li, Chuchu Liang, Fangzhou Lu (discussant), Feng Mai, Asaf Manela, Amin Shams, Yushui Shi, Siew Hong Teoh, Aleh Tsyvinski, David Yang, Lu Zheng, Chenqi Zhu, and conference and seminar participants at University of California Irvine, 2019 Conference on Financial Economics and Accounting at New York University, 2020 American Finance Association Conference Poster Session, and 2020 CAFR Research Workshop on FinTech. All errors are our own.

# 1 Introduction

The rise of FinTech is one of the most critical developments in finance in the past decade. One important FinTech solution is initial coin offerings (ICOs), where investors can purchase blockchain-based cryptocurrencies directly from entrepreneurs. While ICO provides a new way of fundraising, there are extensive debates among practitioners and researchers about how to understand cryptocurrencies. On the one hand, there are growing concerns about whether speculations fuel the development of the market. For example, Satis—a security token advisory firm—claims that over 80 percent of ICOs in 2017 were scams.<sup>1</sup> In addition, there is evidence of price manipulation in Bitcoin and other cryptocurrencies (Griffin and Shams, 2020). On the other hand, many cryptocurrencies, such as Bitcoin and Ethereum, are highly valued on the market. The blockchain technology behind them is often referred to as the “*new internet*”, and some investors believe that it will bring revolutionary changes to every aspect of our lives. One central question of these debates is whether technological sophistication of cryptocurrencies affects their valuations.

Cryptocurrencies are a special asset class. Unlike stocks, cryptocurrencies do not distribute dividends, and there is no traditional accounting information for them. Cryptocurrencies are also different from fiat currencies in the sense that their value is not backed by any government. Many investors invest in cryptocurrencies because they trust the blockchain technology embodied in these digital coins. They believe that the blockchain technology is an important innovation and that at least some coins are assets that represent a stake in the future of this technology. Recent theoretical papers in the cryptocurrency literature echo with this viewpoint and emphasize the importance of technology in determining the viability and valuation of coins. For example, Fanti et al. (2019) show the pricing implications of different technologies used in setting up the cryptocurrency platforms. Because of these distinct features of cryptocurrencies, we focus on the technological sophistication of cryptocurrencies.

Measuring the technological sophistication of cryptocurrencies is challenging because of the

---

<sup>1</sup>For the full report, please see: [https://research.bloomberg.com/pub/res/d28giW28tf6G7T\\_Wr77aU0gDgFQ](https://research.bloomberg.com/pub/res/d28giW28tf6G7T_Wr77aU0gDgFQ).

limited information available. To overcome the challenge, we utilize the disclosure of the coins—their whitepapers—to measure the technological components employed in the cryptocurrencies. We use textual analysis to analyze the content of whitepapers. In particular, we use both supervised machine learning and unsupervised machine learning methods (i.e., word embedding and Latent Dirichlet Allocation (LDA)) to construct three measures of technological sophistication (Technology indexes) from a comprehensive database of ICO whitepapers. The supervised machine learning method we employ is a top-down approach that closely imitates the way investors assess ICOs. The unsupervised machine learning methods are bottom-up approaches to study the textual elements of whitepapers. The advantage of the unsupervised methods is that they require little human input. We also construct a composite index as the fourth technology index, which is the average of the three indexes mentioned above. We study the determinants of the tech indexes using various cryptocurrency characteristics. We find that cryptocurrencies that just use the Ethereum blockchain, have lower GitHub activities, have ambiguous whitepapers, and have less reliable teams tend to have lower tech indexes. However, the R-squared is only 0.136 when we use all the cryptocurrency characteristics, suggesting that the majority of the variation in the tech indexes are not captured by these characteristics.

To understand the role of technological sophistication in cryptocurrency pricing, we start by studying the relationship between the technology indexes and ICO successes. We first examine whether the technology indexes are related to ICO fundraising. If the entrepreneurs cannot raise any funding, the ICO is not likely to succeed, so the ability to raise funding is one of the most important steps in a successful ICO. If ICO performances are fully driven by speculations, investors would not care about the technology associated with the ICOs. Under this hypothesis, the technology indexes would not predict ICO successes. However, we find that ICOs with high technology indexes are more likely to raise capital and more likely to be traded in the secondary market subsequently. The economic magnitude of the effect is significant. For instance, a one standard deviation increase in the composite technology index is associated with a 10.4 percent increase in the listed probability, which is a 40.1 percent increase of the average. The results suggest that investors take

the underlying technology of the ICOs into consideration.

Next, we investigate whether the underlying technology of ICOs is associated with subsequent performances. The process to fully incorporate technology-related information may take months due to the complexity of blockchain technology. To test this conjecture, we examine the relationship between the technology indexes and the long-run performance of ICOs. We measure long-run performance using cumulative post-ICO returns, abnormal returns, and liquidity measures. We find that the ICOs with higher technology indexes tend to have better performance in the long run compared to other ICOs. A one standard deviation increase in the composite index is associated with a 23.9 percent increase in cumulative returns at the 300-day horizon.

We also investigate whether our indexes help understand ICO failure measured by delisting. We find that the ICOs with higher technology indexes are less likely to be delisted subsequently. The economic magnitude of the effect is also large. For instance, a one standard deviation increase in the composite technology index leads to a 2.52 percent decrease in delisting probability.

So far, we have shown that the technology indexes strongly and positively predict ICO successes and subsequent performances. We argue that the results are consistent with the notion that investors care about the technological sophistication of the cryptocurrencies, but it takes time for the market to incorporate the information, leading to predictable returns. We present additional evidence in support of the delayed reaction mechanism and attempt to rule out potential alternative explanations. An implication of the delayed reaction mechanism is that investors should be able to quickly incorporate the fundamental information if the whitepapers are written clearly. Consistent with the implication, we show that among the whitepapers with better readability, the long-horizon predictive power of the technology indexes is weaker. We also find that there is no return reversal phenomenon, suggesting that the return predictability results are unlikely to be driven by investor overreactions.

Overall, these results suggest that the underlying technology is an important determinant of cryptocurrency prices, and support the argument that investors do take the technological components in the ICO whitepapers into their consideration. However, it takes time for investors to

differentiate the fundamentally sound ICOs from the others fully. The delayed reaction from investors may be caused by investor inattention and the complex nature of the technologies, both of which necessitate more time to process related information.

## Related Literature

This paper contributes to the fast-growing literature on the economics of ICOs and digital assets in general. [Yermack \(2017\)](#) is the first paper to explore the financial implications of blockchain. [Liu and Tsyvinski \(2018\)](#) provide one of the first comprehensive analyses of the risk-return tradeoff of cryptocurrencies. [Liu et al. \(2019\)](#) examine the cross-section of cryptocurrency and establish a cryptocurrency three-factor model. Recently, several theoretical papers examine the rationale and mechanisms of ICOs and cryptocurrencies ([Cong and He, 2019](#); [Cong et al., 2019](#); [Catalini and Gans, 2018](#); [Sockin and Xiong, 2018](#)). Our paper is closely related to [Cong et al. \(2019\)](#) and [Sockin and Xiong \(2018\)](#), which argue that the value of cryptocurrency is fundamentally anchored by the underlying utility value. In other words, their models predict that coins have fundamental values and the fundamental values are crucial for performance. However, there is little evidence showing the importance of the fundamental values of coins because it is hard to measure that empirically. A set of empirical papers study factors that contribute to ICO success, including [Howell et al. \(2020\)](#), [Deng et al. \(2018\)](#), and [Lee et al. \(2019\)](#). [Lyandres et al. \(2020\)](#) studies the determinants of ICO successes and performances, and overturn some existing findings in the literature. In general, they find social media and team play a significant role in ICO success and performance. Although some prior papers touch about whitepapers (e.g., [Dittmar and Wu, 2019](#) and [Florysiak and Schandlbauer, 2019](#)), our paper is the first paper that tries to measure the technological sophistication of cryptocurrencies using various machine learning methods and account for the relationship between whitepapers with ICO short and long-run performance.<sup>2</sup> Our tech indexes appear to play a significant role in explaining ICO success, short-, and long-horizon performance, all of which are not well understood in the literature.

---

<sup>2</sup>Some studies also look at the text of social media about cryptocurrencies. For example, [Shams \(2019\)](#) use text from Reddit to measure the connectivity among cryptocurrencies.

This paper provides support to the theoretical literature that links the technological advances of blockchain to cryptocurrency valuations. [Budish \(2018\)](#), [Abadi and Brunnermeier \(2018\)](#), and [Hinzen et al. \(2019\)](#) discuss the limitations of proof-of-work technologies and the pricing implications of them. [Fanti et al. \(2019\)](#) show that the pricing implications of proof-of-stake. Consistent with the theoretical implications of the literature, our paper shows that the technological components affect the valuations of cryptocurrencies.

Our study also adds to the literature on machine learning and textual analysis in finance.<sup>3</sup> The application of machine learning in finance is a new and growing literature. Existing studies use machine learning methods to construct text-based uncertainty ([Manela and Moreira, 2017](#)), predict stock returns ([Gu et al., 2020](#)), measure corporate culture ([Li et al., 2019](#)), and analyze online reviews (e.g., [Sheng, 2019](#)). Recently, [Bybee et al. \(2020\)](#) use machine learning to measure the state of the economy via textual analysis of business news. To our best knowledge, this paper is the first paper to use machine learning methods to conduct textual analyses of cryptocurrencies. Our paper employs both supervised and unsupervised machine learning methods, which allows us to draw reliable references.

The rest of the paper is organized as follows. Section 2 explains the background of ICOs and the data we use. Section 3 introduces the construction and validation of the technology indexes. Section 4 describes our main empirical results and Section 5 documents the subsample and additional results. We conclude and discuss implications for policy in Section 6.

## 2 Background and Data

### 2.1 Initial Coin Offering (ICO) basics

In a typical ICO, entrepreneurs issue digital assets (“tokens”) that are implemented on a blockchain or a contract to deliver such tokens in the future (e.g., a Simple Agreement for Future Tokens, or

---

<sup>3</sup>See [Tetlock, 2014](#) and [Gentzkow et al., 2019](#) for reviews on textual analysis. Textual analysis includes both machine learning methods and other methods, such as word count. For recent studies using the word count method, please see [Liu and Matthies \(2018\)](#) and [Fisher et al., 2020](#).

SAFT). Entrepreneurs then use the raised capital to create an online platform or ecosystem where the native token can be used.

In general, these tokens can be classified into three types based on their purposes. The first type is called “utility token” because its purpose is to redeem a product or service in the future—this is the largest group of tokens. The second type is called “security token”, which is similar to conventional securities but recorded and exchanged on a blockchain to reduce transaction costs and create a record of ownership. This type of token gives holders the rights for associated cash flows, such as dividends. The third type is called “asset token”, which serves as a general-purpose medium of exchange and store of value. These are often termed “coins”, such as Bitcoin.

Initial Coin Offering is appealing to both start-up companies and investors. The start-up companies that choose to issue ICOs are usually those that "conventionally finance themselves with angel or venture capital (VC) investment" (see [Howell et al., 2020](#)). ICOs are attractive to these start-up companies because ICOs allow them to avoid regulations from SEC and intermediaries such as venture capitalists and banks, leading to lower financing costs and easier access to capital. Investors participate in ICOs for various reasons. Some investors may believe in the intrinsic value of the project and are optimistic about the technological innovations embedded therein. Other investors may be speculators who are attracted by the quick cash-out ability.

The first ICO was issued by Mastercoin in July 2013. In 2014, Ethereum also launched a token sale and raised over \$15 million to support its development. In 2017, ICOs have become popular, and 875 startups successfully raised capital using token sales during the year. As of February 2019, ICOs have raised over 25 billion USD.<sup>4</sup>

As a new source for seed and early-stage funding, ICOs raise money from many small investors over the Internet. In that sense, the ICOs are similar to crowdfunding, where investors get future rewards or deals on products and get securities for exchange. However, ICOs are different from crowdfunding in that they are blockchain-based and involve more advanced technology for their products and services. ICOs are also similar to initial public offerings (IPOs) in the sense that

---

<sup>4</sup>Source: <https://icobench.com/>.

tokens can be listed on one or more cryptocurrency exchanges, so investors can benefit from the price appreciation of a listed token even before the project launches. This process is usually much faster than that of IPOs. The whole process ranges from several days to several months, but there is no guarantee of listing.

## 2.2 Data on ICOs

Our dataset consists of three different components: ICO characteristics from trackico.com, daily trading data from coinmarketcap.com, and textual measures from ICO whitepapers. There are over 4,100 ICOs on trackico.com, with 2,452 closed, 575 trading, 264 ongoing, 82 pre-sale, 307 upcoming and 422 unknown. We focus on ICOs between January 2017 and December 2018. The final sample consists of 2,916 ICOs, which raised more than \$17 billion in total. For each ICO, we collect the following information: ICO start and end date, ICO price, total capital raised, trading status, pre-ICO, bonus, platform, accepted currency, the founder team, country, industry, links of whitepapers, official website, GitHub and Twitter.

We define two measures of ICO success. The first one is “Trading”, a self-reported indicator variable by fundraisers to trackico.com, indicating whether the token is trading on cryptocurrency exchanges. The second one is “Success”, which equals to 1 if an ICO successfully raised any capital (Benedetti and Kostovetsky, 2018). Other ICO characteristics serve as control variables. “ICO length” is the number of days between the start and end of an ICO. “ICO price” is the cost per token in US dollars. “Total Raised” is the amount of money raised in millions of US dollars. “Pre ICO”, “Bonus”, “Ethereum Based” and “Accept BTC” are indicator variables about whether the ICO has a pre-ICO, offers bonus to investors, is built on Ethereum platform and accepts Bitcoin as a payment currency, respectively. “Team size” is calculated as the number of team members. We define “Has GitHub” and “Has Twitter” to be indicator variables of whether the fundraiser has a GitHub or a Twitter homepage. We further control for Bitcoin price on the ICO start date or the coin’s listing day as a proxy for the market sentiment. Finally, we control for quarterly, categorical and geographical (continent-level) fixed effects.



Next, we merge ICO data with information from [coinmarketcap.com](https://coinmarketcap.com), the leading information source of cryptocurrency trading data, which is also a primary information source in the ICO literature. By the end of 2018, [coinmarketcap.com](https://coinmarketcap.com) has provided data for over 3,600 cryptocurrencies, among which 2,070 are active while 1,583 are delisted. We collect daily opening price and 24h dollar trading volume on all coins from August 2013 to December 2018. We then use token names, ticker symbols, and website slugs to merge these variables with our ICO data. Since many coins on [coinmarketcap.com](https://coinmarketcap.com) were not issued through ICO, and many ICOs do not list their coins on any exchange, we get a merged sample of 765 observations.

With the merged sample, we first define a third ICO success measure, “CMC Trading”, which equals to one if the coin has ever appeared on [coinmarketcap.com](https://coinmarketcap.com). This measure also aims at characterizing the same fact (i.e. whether the coin is traded on an exchange) as the self-reported measure “Trading”, but is more comprehensive.<sup>5</sup> Therefore, we use “CMC Trading” in our main analysis and consider the other measures in the robustness tests. We define “First Open/ICO Price” to measure the premium on the listing day and “Delist” to characterize whether the coin is delisted from cryptocurrency exchanges. We also calculate the cumulative rate of return, Bitcoin-adjusted rate of return and 24h trading volume after the coin has been listed for 7 days, 30 days, 90 days, 180 days, 240 days and 300 days. These measures capture the short- and long-term performance and liquidity of cryptocurrencies.

The last set of variables comes from textual analysis of ICO whitepapers, which are downloaded from [trackico.com](https://trackico.com). We obtained 1,629 valid whitepapers in PDF format. In Table OA.4, we list all other variations of whitepaper status. Next, we convert PDF files into TXT format, which can be used as the raw input for textual analysis.

Using this whitepaper corpus, we first construct our main measures of technology indexes, which we explain in detail in Section 3. Moreover, we consider three well-known textual measures as control variables: Readability, Tone, and Uncertainty. “Readability” is characterized by the Fog Index, a widely adopted measure in finance and accounting literature. Developed by Robert

---

<sup>5</sup>The correlation between CMC Trading and Trading is 75.8%. Trading is highly accurate if it equals to 1, but is not comprehensive, as we identified approximately 200 more trading tokens on [coinmarketcap.com](https://coinmarketcap.com).

Gunning in 1952, Fog Index is a linear combination of the percentage of complex words and the average number of words per sentence.<sup>6</sup> “Tone” is the difference between positive and negative words divided by the total number of words, and “Uncertainty” is the percentage of uncertainty words among all words used in a whitepaper. All lexical categories are defined in [Loughran and McDonald \(2011\)](#).

## 2.3 Summary Statistics

We report the summary statistics of the sample characteristics in Table 1. Panel A of Table 1 presents summary statistics on variables related to ICO characteristics. On average, it takes 51 days to complete an ICO with a team of 11 people. 18% of the ICOs are self-reported as trading and 38% have non-zero values of capital raised. Moreover, 60% have a GitHub homepage for their project and over 90% have set up their Twitter accounts.

Panel B of Table 1 presents summary statistics on the merged sample. Consistent with the literature, we identify that 26% of ICOs have listed tokens on an exchange at some point in time. Among these listed cryptocurrencies, only 10% are delisted while the remaining 90% are still active. On average, investing in a cryptocurrency during an ICO can earn a premium of 120% on the first trading day, indicating a large amount of first-day price reaction. Moreover, the return of cryptocurrency investment increases as time goes by, from 19% during a 7-day holding period to 151% during a 300-day holding period. The 24h trading volume fluctuates with different time spans, varying from \$1.5 million to \$2.78 million. ICO characteristics with respect to the merged subsample are also reported in this panel.

## 3 Measuring Technological Sophistication

In this section, we discuss how we measure the technological sophistication for the cryptocurrencies based on their whitepapers. We first present how we construct the technology indexes using

---

<sup>6</sup>The complete formula of Fog Index is:  $\text{Fog Index} = 0.4[(\text{words/sentences}) + 100((\text{complex words})/\text{words})]$ . “Complex words” are words consisting of three or more syllables.

different machine learning methods, and then we validate these measures.

### **3.1 Measure Construction**

We use several machine learning techniques to capture the technological components of the whitepapers, including both supervised and unsupervised methods. We first construct a supervised machine learning index. We mimic the way investors evaluate whitepapers and manually assign scores to 200 whitepapers, which we use as the training set. Then, we use a supervised machine learning algorithm, train it on the training sample, and extrapolate scores to the remaining whitepapers. The supervised machine learning method we employ is a top-down approach that closely imitates the way investors assess ICOs.

Additionally, we use two different techniques in the unsupervised machine learning literature to measure the technological aspect of the whitepapers: word embedding along with K-means clustering and Latent Dirichlet Allocation topic modeling approach. The unsupervised machine learning methods are bottom-up approaches to study the textual elements of whitepapers. One important advantage of unsupervised machine learning methods is that they require little human input. In other words, they do not require researchers to have good prior knowledge about what type of words they are looking for in the texts. Below, we briefly summarize the three methods and their estimations, and refer interested readers to the original paper for details.

#### **Supervised Machine Learning**

First, we use supervised machine learning methods to construct a technology index. Supervised machine learning methods learn from a training set in which both the input and the output are known. To construct the training sample, we read through 200 randomly selected whitepapers and give a score from 1 to 4 based on their technical sophistication. The process closely imitates the way investors evaluate the whitepapers. All the whitepapers emphasize on using blockchain and related technologies. Thus, these projects either employ more advanced blockchain technology or apply existing blockchain technology to different areas. The readers assign a high score (e.g.,

3 or 4) to a whitepaper when they think the ICO project involves more advanced and convincing technology. For example, Filecoin uses a novel class of Proof-of-Storage schemes called Proof-of-Replication, and receives an average score of 4. Then, we conduct preprocessing to the training set. We form all two-word phrases in the corpus, remove unigrams and bigrams that appear in less than ten documents, and convert the corpus to a document-term matrix. The final training set consists of 200 documents and 20,586 unique terms.

We consider the following supervised machine learning approaches as potential candidates: panelized linear methods (LASSO, ridge, and elastic net), dimension reduction methods (PCR and PLS), decision tree boosting methods (random forest, gradient boosting), and neural networks. In the Online Appendix, we provide a brief introduction for each supervised method. In order to tune the hyperparameters of the supervised learning models and find the best model for constructing our supervised technology index, we need to quantify the performance of the model. We evaluate the model performance based on out-of-sample  $R^2$ . We use 5-fold cross-validation to build a validation set whose labels are known but are not used for training. Specifically, we divide the training set into five subsets, each of which contains 40 observations. Following that, each subset will be used as the validation set to evaluate the model based on  $R^2$ , while the remaining four subsets are used as the training set. The average out-of-sample  $R^2$  is the simple average of  $R^2$  on the five subsamples. Table 3 shows the best out-of-sample R-square ( $R_{OOS}^2$ ) for each supervised method and their corresponding hyperparameters. For our sample, partial least square (PLS) performs the best and has a  $R_{OOS}^2$  of 45.88%. Hence, we use the predicted technology score from PLS as our supervised technology index.

## Word Embedding

Word embedding is one of the most popular word representation methods in natural language processing (NLP) in recent years. Developed by Mikolov et al. (2013), its goal is to map words to numerical vectors, such that the semantic similarity between words is captured by the geometric distance in the vector space. How to construct such vectors? The intuition comes from the

famous quotation of [Firth \(1957\)](#)—“You shall know a word by the company it keeps.” In other words, the meaning of a word can be inferred from the context, so words appearing in similar contexts should have similar meanings.<sup>7</sup> Word embedding has two main advantages over traditional “bag-of-words” methods. First, it greatly reduces the number of dimensions. Word embedding vectors usually have only a few hundred dimensions, while bag-of-words models are typically sparse vectors of thousands of dimensions. Hence, it is a more efficient representation of the raw text. Second, word embedding maps synonyms to adjacent vectors, so we can use clustering methods on the vector space to divide words into different topics. We use K-means as the clustering algorithm. It is one of the simplest and most popular unsupervised machine learning methods. Given a fixed number of clusters ( $k$ ), K-means seeks a partition of the dataset, such that the within-cluster sum of squared distances between each observation and its closest centroid is minimized. In the Online Appendix, we provide details on the theoretical background of word embedding and k-means clustering and how to choose the optimal number of topics.

We find that the optimal number of topics detected by the algorithm is 20. Hence, we use K-means to cluster word embedding vectors into 20 topics. To interpret the embedding and clustering results, we give each topic a label. For clustering methods, topics are mutually exclusive, so each word can only be grouped into one topic. We name the topics based on the most frequent words in each cluster. Table [OA.1](#) lists the top 15 most frequent terms of each topic.

To further understand the relationship between topics, we apply two machine learning techniques. The first one is hierarchical agglomerative clustering ([Murtagh and Legendre, 2014](#)), which can be used to construct a taxonomy of our topic model. Following [Bybee et al. \(2020\)](#), we agglomerate topics recursively according to the semantic similarity between topics, as captured by the distance between cluster centroids. Figure [1](#) displays the result and shows that three topics (“blockchain”, “system”, and “algorithm”) belong to the same cluster. Another technique we use is multidimensional scaling (MDS, [Torgerson, 1958](#)), which is a non-linear dimensionality reduction algorithm such that the two-dimensional representation best preserves the distance between topics

---

<sup>7</sup>[Li et al. \(2019\)](#) provide a good example in the Appendix to illustrate the intuition.

in the original space. “Blockchain”, “system” and “algorithm” are combined into a broader topic in the taxonomy. They are also adjacent to each other in the inter-cluster distance map. Therefore, we consider these three topics as technology-related topics. For each whitepaper, we calculate the percentage of words that belong to the “blockchain”, “information” or “algorithm” topics, normalize it to zero mean and unit standard deviation, and define it as our embedding-based tech index.

### **Latent Dirichlet Allocation (LDA)**

The second unsupervised machine learning method we use is Latent Dirichlet Allocation (LDA). LDA is a popular method in the finance and economic literature. It has been used to analyze the structure of economic news (Bybee et al., 2020) and to detect latent topics among employee reviews (Sheng, 2019). The basic idea is that each document can be represented as a probability distribution over various topics, where each topic is a probability distribution over the vocabulary of a corpus. Similar to other textual analysis methods, LDA methods involve a step to remove useless information (i.e., stop words) and then represent the text as data. In the Online Appendix, we introduce the LDA model and describe the preprocessing procedures and the choice of topics in more detail. We find that 20 topics is optimal based on the selection process.

To understand the LDA output with 20 topics, we interpret these topics by looking at top words associated with each topic. This is a common approach adopted in most finance and economics literature (e.g., Hansen et al., 2018; Sheng, 2019). Table OA.2 displays the top 15 most relevant terms of each LDA topic (see Online appendix for details on how to find top words for each topic). We assign a label to each topic based on these key terms. Similar to word embedding, we also use hierarchical agglomerative clustering and multidimensional scaling (MDS) to understand the correlation between LDA topics. Panel A of Figure 2 shows the tree structure of LDA topics and Panel B shows the MDS results. These results suggest that we should group “information”, “blockchain” and “system” together and define the normalized proportional attention allocated to the three topics as our LDA-based tech index.

## **Composite Index**

Finally, we create a composite index to aggregate the information from the above three indexes. This is done by taking the simple average of the supervised, embedding-based and LDA-based technology indexes. Both supervised and unsupervised machine learning methods have pros and cons. The composite index can potentially reduce the noise of each index, resulting in a useful proxy. For most of the empirical analysis, we show the results for all of them, including the composite index.

## **3.2 Measure Determinant**

Given that the construction of the technology index is one of the key components of this paper, we use different machine learning methods to capture the technological sophistication of cryptocurrencies. These methods have different advantages. Supervised machine learning methods are relatively easy to interpret. Unsupervised machine learning methods require little human input and do not require prior knowledge about the subject from the researchers.

In our paper, we employ multiple methods to construct the technology indexes and find evidence that the measures capture meaningful information. First, the measures from different machine learning methods are highly correlated with each other, suggesting that the measures are just driven by noise. Second, we construct a composite index to reduce the noise of each measure. It is possible that each method capture only some aspect of the true technological components of cryptocurrencies. Then, the composite index would provide a better proxy because it aggregates the information from the three individual indexes. Third, we find consistent results based on all four measures of technology indexes.

To better understand the indexes, we study the determinants of the technology indexes. We utilize cryptocurrency characteristics from several dimensions, including whether they use Ethereum blockchain, GitHub data, whitepaper information, and other characteristics. GitHub is an open-source online platform that provides repository hosting service for developers. Using the API pro-

vided by GitHub, we obtain the number of (1) users subscribing updates of the repository (*watch*), (2) “likes” received by the repository (*star*), (3) copies made by other developers (*fork*), (4) code revisions (*commit*), (5) pointers to specific versions (*branch*) and (6) developers who have contributed to the source code (*contributor*). These measures are often used by researchers to proxy for product quality and post-ICO technology development (Deng et al., 2018; Dittmar and Wu, 2019). For the determinant results, we use *commit* as the GitHub measure. In the Online Appendix, we also use other GitHub measures as robustness checks and obtain qualitatively similar results.

Table 2 documents the results that relate the technology indexes to these cryptocurrency characteristics. We use the composite index as the dependent variable and we show qualitatively similar results using the other indexes in the Online Appendix. Each of columns (1)–(4) reports the determinant models based on a dimension of coin characteristics. Column (1) shows that cryptocurrencies that use Ethereum blockchain tend to have lower tech indexes, confirming the prior that cryptocurrencies that build their own blockchain on average have higher tech indexes. Column (2) shows that cryptocurrencies with more code revisions in GitHub have higher tech indexes. In Column (3), we find that cryptocurrencies with ambiguous whitepapers tend to have lower tech indexes. In Column (4), we find that cryptocurrencies with more reliable and supportive teams have higher tech indexes. For example, team size and the Twitter account dummy positively predict tech indexes. Column (5) combines all the cryptocurrency characteristics and delivers consistent messages. However, the R-squared of Model (5) is 0.136, suggesting that the majority of the variation in the tech indexes are not captured by the cryptocurrency characteristics.

## 4 Main Results

In this section, we examine whether the technological component of ICOs is associated with ICO success, short-run, and long-run performances. We capture the technological component of ICOs using the four technology indexes we defined above, and we evaluate an ICO using both its



fund-raising stage information and its subsequent performance data.

## 4.1 ICO Success

First, we study the set of characteristics in ICO whitepapers that are most related to ICO success. We use two ways to measure ICO success. The first measure of ICO success is based on whether the cryptocurrency is listed on the coinmarketcap.com (CMC trading) and the second measure is based on whether the ICO successfully raised capital. If the entrepreneur cannot raise any funding, the ICO is not likely to succeed. Therefore, the ability to raise funding is one of the most important steps in a successful ICO. If investors care about the technological components of ICOs, we should expect that it is easier for ICOs with more sophisticated technologies to raise funding. Companies voluntarily disclose whitepapers to communicate with investors in the fund-raising stage, and one of the primary ways that investors evaluate coins is through whitepapers. If whitepapers indeed inform investors about the different aspects of the ICOs, we would be able to extract information from the whitepapers. As discussed in Section 3, we form four measures to summarize the technological component of ICOs: (1) an index based on word embedding, (2) an index based on LDA, (3) an index based on supervised machine learning, and (4) a composite index.

Table 4 documents the results that relate ICO whitepapers' characteristics to ICO successes. Panel A of Table 4 presents results based on CMC trading and Panel B presents results based on whether the cryptocurrency successfully raised capital. We report coefficient estimates for each of the four whitepaper indexes as well as the control variables. Time, categorical, and geographic fixed effects are included in the specifications when indicated.

Panel A shows that the CMC trading indicator positively loads on all four tech indexes, suggesting that when the tech indexes are high, the cryptocurrencies are more likely to be listed on coinmarketcap.com. The coefficient estimates are 0.070, 0.107, 0.086, and 0.124 for the four indexes, respectively. The relationships are highly significant at the 1 percent level for all four cases. The economic magnitudes are large. For example, the coefficient estimate on the composite tech-

nology index is 0.124 in the univariate specification and the standard deviation of the tech\_comp index is 0.84. In other words, a one standard deviation increase in the composite index leads to an increase of the listed probability by 10.4 percent—a 40.1 percent increase of the sample average of the listed probability. In the multivariate specification with controls and fixed effects, the coefficient estimate on the composite technology index is 0.066. That is, a one standard deviation increase in the composite index is associated with an increase of the listed probability by 5.54 percent under the multivariate specification.

Panel B measures ICO success based on whether the ICO raised capital (Success indicator). The coefficient estimates are largely consistent with those in Panel A—the coefficients on the four indexes are 0.061, 0.077, 0.056, and 0.091. The loadings on the four technology indexes remain highly statistically and positively significant at the 1 percent level for all the specifications. The coefficient estimate on the composite technology index is 0.091 in the univariate regression, which suggests that a one standard deviation increase in the composite index is associated with a 7.64 percent increase of the probability that the ICO raised capital—a 20.1 percent increase of the sample average. In the multivariate specification with controls and fixed effects, the coefficient estimate on the composite technology index is 0.060. That is, a one standard deviation increases in the composite index is associated with a 5.04 percent increase in the probability that the ICO raised capital.

Further evidence that the technological component serves as an important factor for ICO success is the  $R^2$ . For example, Panel A of Table 4 shows that a single variable of each of the technology indexes already explains between 3 percent and 6 percent of the variation of CMC trading. Interestingly, in untabulated results, we find that the Quarterly Fixed Effects seem to be the most important factor—they explain 14 percent of the variation of the CMC Trading variable. In other words, the timing of the ICOs is important in determining whether they can successfully raise capital. Nevertheless, our technology indexes are still some of the most important factors that contribute to the success of an ICO. Overall, the results show that when an ICO whitepaper contains more discussion on technology-related topics as captured by our indexes, the ICO is more likely to

be successful.

## Industry Subsample

In this subsection, we test whether the technology indexes are stronger predictors of ICO successes in industries that technological components are deemed more important. In certain industries (e.g., platform; trading), investors may scrutinize the technological components of the ICOs, while in other industries (e.g., gaming; charity), this is not the case. We categorize “platform”, “cryptocurrency”, and “trading” as the technology-related industry, and construct an indicator variable (“industry”) to denote the technology-related industries. We test whether the technology indexes strongly predict ICO successes for coins in the technology-related industries.

We present the subsample results based on industries in Table 5. Consistent with the baseline results, the technology indexes positively predict ICO successes. The cross-terms between the technology indexes and the “industry” indicator are all positive and largely significant. The economic magnitudes are large. For example, judging from the composite index, the coefficient estimates almost double for the coins in the technology-related industries relative to the rest of the coins. These results also support the view that investors value the technological components of the cryptocurrencies, especially for the coins in the technology-related industries.

## 4.2 Long-Horizon Performance

In this section, we investigate whether the technology indexes help forecast the medium- to long-horizon ICO returns. In the equity market, initial public offerings tend to underperform in the long run (see Ritter, 1991; Loughran and Ritter, 1995). In sharp contrast, initial coin offerings perform well in the medium- to long-horizon (see Benedetti and Kostovetsky, 2018). In order to study the speed of information acquisition of the investors, we ask whether the long-horizon performance of the ICOs is related to the technology indexes.

We track the subsequent returns of the ICOs over different horizons—from 7-days ahead to 300-days ahead. Shumway (1997) documents that stock delisting is associated with a negative 10

percent return on average. In the robustness test section, we experiment with alternative assumptions and show that the results are consistent.

First, we look at how the technology indexes predict the subsequent performance of initial coin offerings. The results are documented in Table 6. Panel A, B, C, and D of Table 6 document the results for the index based on supervised machine learning, index based on word embedding, index based on LDA, and the composite index, respectively. We regress the cumulative ICO returns on current technology indexes, controls, and fixed effects. In general, we find that the technology indexes positively predict the subsequent performances of the ICOs. For example, based on the composite technology index, the point estimates are positive across all horizons. The point estimates steadily increase but are insignificant at short horizons. The point estimates start to become significant in longer horizons. At the 240-day horizon, the point estimate increases to 0.280, indicating a 23.5 percent increase in cumulative returns at this horizon for one standard deviation increase in the composite technology index. At the 300-day horizon, a one standard deviation increase in the composite technology index leads to a statistically significant 23.9 percent increase in cumulative returns.

The ICOs took place at different times, and our return measures do not take the time component information into consideration. A common factor that is important for the ICO market is Bitcoin returns. Thus, we also conduct a similar exercise with abnormal returns that are adjusted to Bitcoin returns. Table 7 reports the results of this test and shows similar results in terms of statistical significance and economic magnitude as in Table 6. For example, based on the composite technology index, the point estimates remain positive across all horizons. The point estimates become significant at the 180-day horizon. At the 240-day horizon, the point estimate increases to 0.300, indicating a 25.2 percent increase in cumulative returns at this horizon for one standard deviation increase in the composite technology index. At the 300-day horizon, a one standard deviation increase in the composite technology index leads to a statistically significant 29.7 percent increase in cumulative returns.

Overall, the medium- and long-horizon results are consistent with the idea that it takes time

for the market to fully incorporate information about technological sophistication. Although coins with high technology scores have a high probability of raising funds, investors undervalue these high-tech coins on average.

### 4.3 Other Measures of Performance

In this section, we use two additional measures to evaluate ICO performances. The first one is the liquidity measure and the second one is the delisting probability measure.

We measure coins' liquidity as the log transformation of the 24-hour trading volume. On average, we find that liquidities are higher for older coins, consistent with [Howell et al. \(2020\)](#). We examine the relationships between characteristics of whitepapers and coins' liquidity measures. We report the results in [Table 8](#). In our model specifications, we include quarterly, categorical, and geographic fixed effects. We find that the four technology indexes are positively associated with coin liquidity. These results are always statistically significant across the different horizons since inception.

We then investigate the relationships between coins' delisting probability and the characteristics of the whitepapers. We define Delist as an indicator variable, which is equal to 1 if a token is delisted from CMC. The results are reported in [Table 9](#). The results show that coins with high technology scores are less likely to be delisted subsequently. The economic magnitude of the effect is large. For instance, in the standalone specification, a one standard deviation increase in the composite technology index leads to a 2.52 percent decrease in delisting probability.

The results in this section highlight that coins with high technology scores are intrinsically superior. The results provide supports to our argument that the investors in the coin market take technical aspects of the ICOs into consideration. However, as we have shown above, it takes a considerable amount of time for the market to reach the proper pricing of the ICOs eventually.

## 5 Discussion

In the previous section, we show that the technology indexes strongly and positively predict ICO successes and subsequent performances. We argue that the results are consistent with the notion that investors care about the technological sophistication of the cryptocurrencies, but it takes time for the market to incorporate the information leading to predictable returns. In the first two parts of this section, we present additional evidence in support of the delayed reaction mechanism and attempt to rule out potential alternative explanations. In the last two parts of the section, we present additional robustness checks.

### 5.1 Comparison of Tech Index and Other Measures

In this subsection, we compare our tech indexes with other measures that may contain information on the technological sophistication of cryptocurrencies, including a GitHub measure and a simple word count measure.

One candidate that potentially captures some information of the technological sophistication of cryptocurrencies is the GitHub measures. However, the GitHub measures are *ex post* measures that capture information about the successes of the ICOs. Moreover, these measures may contain information such as the hype around cryptocurrencies.

In addition, there are multiple methods to conduct textual analysis. For example, the word-count method where we can just count the number of words that belong to a dictionary is well-accepted in the finance and economic literature (e.g., [Manela and Moreira, 2017](#); [Liu and Matthies, 2018](#); [Fisher et al., 2020](#)). The word-count method is particularly useful when researchers have good prior knowledge about what they are looking for and the list of words is straightforward. However, cryptocurrency and blockchain are new phenomena and researchers have limited knowledge about what should be a good list of words to describe the technology involved. In this case, unsupervised machine learning methods, such as LDA, are more proper and can overcome this issue. One important advantage of machine learning methods, especially the unsupervised machine

learning methods such as word embedding and LDA, is that they do not require researchers to have good prior knowledge about what type of words they are looking for in the texts.

With that being said, we construct technology measures from GitHub and from a simple word count method to compare with our tech indexes. The Github measure we use is *commits*, the number of code revisions of a project on GitHub. The simple word count measure captures the percentage of technology words in a whitepaper, where the technology words are defined by a blockchain dictionary.<sup>8</sup> The complete word list can be found in Table OA.3. In Table 10, we present results using the tech indexes to predict CMC trading, controlling these two types of measures. Columns (1)–(3) report results using the composite tech index, the GitHub *commits* measure, and the simple word count measure, and the point estimates on the variables are all positive and significant at the 5 percent level. Columns (4) and (5) report results using the composite tech index controlling for the GitHub *commits* measure and the simple word count measure, respectively. When both the composite tech index and the GitHub measure are included, the coefficient estimates on both measures remain positive and highly statistically significant. However, the tech index completely subsumes the explanatory power of the simple word count measure. When all three variables are included, the coefficient estimate on the composite tech index remains positive and statistically significant at the 1 percent level, and the point estimate on the simple word count measure is insignificant.

## 5.2 Subsample on Whitepaper Readability

In the main result section, we argue that the findings are consistent with the investor delayed reaction to technological sophistication due to thcan be found in Table OA.3. In Table 10, we present results using the tech indexes to predict CMC trading, controlling these two types of measures. Columns (1)–(3) report results using the composite tech index, the GitHub *commits* measure, and the simple wore complex nature of the cryptocurrency whitepapers. An implication of this argument is investors should be able to quickly incorporate the fundamental information if the whitepapers are written clearly. Therefore, among the whitepapers with high readability, we

---

<sup>8</sup>See <https://consensys.net/knowledge-base/a-blockchain-glossary-for-beginners/>; <https://blockgeeks.com/guides/blockchain-glossary-from-a-z/>; <https://www.blockchaintechnologies.com/glossary/>.

should expect weaker results on long-horizon performances.

We measure the whitepaper readability using the Fog index. We construct an indicator variable (“Easy”) that equals to 1 if the whitepaper has a below-median Fog index and 0 otherwise. We present the results in Table 11. Panel A of Table 11 presents results based on the rate of returns. Consistent with the baseline results, we find that the technology indexes positively and significantly predict the long-horizon returns. The cross-terms between the technology indexes and the indicator variable (“Easy”) are negative and significant at the long-horizons, suggesting that the long-horizon return predictability of the technology indexes concentrate among the cryptocurrencies with low readability. For example, the coefficient estimate on the composite index at the 300-day horizon is 1.217, while the cross term between the composite index and the indicator variable at the same horizon is -1.160. That is, the long-horizon return predictive power of the composite index entirely concentrate on the cryptocurrencies with low readability.

Panel B of Table 11 shows results based on the Bitcoin-adjusted rate of returns. Similar to the results in Panel A, we find that the cross terms between the technology indexes and the indicator variables are negative and significant at the long-horizons across all the specifications. Overall, we confirm the implication of the investor delayed reaction mechanism: the long-horizon return predictability results are weaker for cryptocurrencies with high readability.

### **5.3 Return Reversal**

In the previous section, we find that the technology indexes positively and significantly predict cumulative ICO returns over the long horizons. We argue that the findings are consistent with investors’ delayed reaction to the technical aspects of the cryptocurrencies. An important alternative interpretation of the findings is that investors may overreact to technological sophistication of the cryptocurrencies, leading to results of ICO return predictability. [Barberis et al. \(1998\)](#) theoretically demonstrate that investor overreaction to fundamentals can lead to overvaluation of asset values. [Pastor and Veronesi \(2003\)](#) show that investor learning about uncertain fundamentals can lead to a bubble-like phenomenon. A common implication of the models based on investor overreaction or



learning of asset fundamentals is that the asset values would eventually reverse to the fundamental values. Technology is one important aspect of fundamental of cryptocurrencies. Therefore, we would expect return reversal if investors over-react to technological fundamentals of cryptocurrencies.

To test this alternative mechanism, in this section, we test whether there is a long-horizon return reversal phenomenon for the ICOs with high technology indexes. To detect any return reversal effect, we use the technology indexes to predict ICO returns from 180 days onward. The results are documented in Table 12. Panel A and Panel B of the table document results for the rate of returns and Bitcoin-adjusted rate of returns, respectively. Overall, we do not find evidence of more subsequent return reversal for coins with high technology indexes.

## 5.4 ICO First-Day Price

Extensive research has shown that there is a substantial amount of first-day performance in initial public offerings in the equity market.<sup>9</sup> Recently, [Benedetti and Kostovetsky \(2018\)](#) document a similar first-day price reaction in the initial coin offering market. In this section, we study whether the technology indexes help predict not only the long-horizon phenomenon but also the first-day price reaction of ICOs.

Our measure of first-day price reaction is defined as the natural logarithm of the ratio between the first opening price and the ICO offer price. By definition, the sample only includes coins with trading records. Table 13 reports the results for ICO first-day price. Quarterly, categorical and geographic fixed effects are included in the specifications when indicated. We find that the coefficient estimates of the four technology indexes are all positive and significant at the 1 percent level. In other words, the technology indexes positively and significantly predict the first-day price reaction. The coefficient estimates are 0.337, 0.414, 0.310, and 0.458 for the four indexes, respectively. The economic magnitudes of the coefficient estimates are large. The coefficient estimate remains stable in the multivariate specification with controls and fixed effects, where the

---

<sup>9</sup>For a survey paper, see [Beatty and Ritter \(1986\)](#).

coefficient estimate on the composite index is 0.417.

Overall, the technology indexes are strongly and positively predict both short-horizon and long-horizon ICO performances. These two sets of results suggest that, although coin market investors take the technical aspects of coins into consideration, they fail to incorporate the information fully.

## 5.5 Robustness

In this section, we conduct several robustness tests. First, we use an alternative measure of success, Trading, which indicates whether the token is traded on a cryptocurrency exchange. We examine whether the technology indexes predict ICO success under this measure and run a similar regression as in Table 4. Table 14 Panel A reports the result. The coefficients on the technology indexes are positive and significant and support the same conclusion as in Table 4.

Second, we use linear regression in Table 4 where the dependent variable is a binary variable. Alternatively, we can use a Logit or Probit model. Table 14 Panel B reports the results from a Logit regression and finds similar results as in Table 4. In the untabulated results, we show that the results under the Probit model are qualitatively similar.

Third, it is well-documented that we have to impute delisted returns for equity to avoid delisting bias in the data (Shumway, 1997). The equity return data from CRSP automatically contain imputed returns for delisted stocks. For the same reason, we may need to impute returns for delisted ICOs. We set a large negative value -99% as their returns after listed for all delisted ICOs. We then redo the tests on whether the technology indexes affect short-run and long-run returns with and without adjusting Bitcoin returns as in Table 6 and 7. Table 14 Panels C and D report the results. Similar to the results in Tables 6 and 7, ICOs with higher technology indexes tend to outperform in the long-run. The economic magnitudes are also close.

## 6 Conclusion

There are two views about cryptocurrency and blockchain technology. The first view is that the cryptocurrency market represents bubbles and fraud. The second one believes that the value of the cryptocurrency market comes from the innovative technologies and that a stake in cryptocurrencies is an investment in the future of the technology. This study contributes to this debate by providing novel measures of technological sophistication of cryptocurrencies via textual analysis of ICO whitepapers. We construct a set of four text-based technology indexes from a comprehensive sample of ICOs' whitepapers. We find that the ICOs with higher Tech-Index are more likely to succeed and less likely to be delisted subsequently. Although the Tech-index does not statistically significantly affect the short-run returns of ICOs, it has a positive impact on their long-run performance. In short, our findings suggest that technological sophistication is an important driving force for the performances and valuations of ICOs.

Our findings have important policy implications. Although SEC has launched several initiatives on regulating ICOs, there are no clear disclosure requirements. Our results show that disclosures such as whitepapers are potentially important for the long-term development of the cryptocurrency market. Thus, it might be useful to set up a requirement or guideline for formats and necessary components in the whitepaper, which is a natural analogy for disclosure requirements for public firms (e.g., 10K) and financial firms (e.g., 497K for mutual funds).

## References

- Abadi J, Brunnermeier M. 2018. Blockchain economics. *Working Paper, National Bureau of Economic Research* .
- Barberis N, Shleifer A, Vishny R. 1998. A model of investor sentiment. *Journal of Financial Economics* **49**: 307–343.
- Beatty RP, Ritter JR. 1986. Investment banking, reputation, and the underpricing of initial public offerings. *Journal of Financial Economics* **15**: 213–232.

- Benedetti H, Kostovetsky L. 2018. Digital tulips? returns to investors in initial coin offerings. *Working Paper, Boston College* .
- Blei DM, Ng AY, Jordan MI. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research* **3**: 993–1022.
- Budish E. 2018. The economic limits of bitcoin and the blockchain. *Working paper, University of Chicago and NBER* .
- Bybee L, Kelly BT, Manela A, Xiu D. 2020. The structure of economic news. *Working paper, Yale University* .
- Catalini C, Gans JS. 2018. Initial coin offerings and the value of crypto tokens. *Working paper, University of Toronto and NBER* .
- Cong LW, He Z. 2019. Blockchain disruption and smart contracts. *Review of Financial Studies* **32**: 1754–1797.
- Cong LW, Li Y, Wang N. 2019. Tokenomics: Dynamic adoption and valuation. *Working paper, University of Chicago* .
- Deng X, Lee YT, Zhong Z. 2018. Decrypting coin winners: Disclosure quality, governance mechanism and team networks. *Working paper, Shanghai University of Finance and Economics* .
- Dittmar RF, Wu DA. 2019. Initial coin offerings hyped and dehyped: An empirical examination. *Working paper, University of Michigan* .
- Fanti G, Kogan L, Viswanath P. 2019. Economics of proof-of-stake payment systems. *Working paper, MIT* .
- Firth JR. 1957. A synopsis of linguistic theory, 1930-1955. *Studies in Linguistic Analysis* .
- Fisher AJ, Martineau C, Sheng J. 2020. Macroeconomic attention and announcement risk premia. *Working paper, University of British Columbia* .
- Florysiak D, Schandlbauer A. 2019. The information content of ico white papers. *Working paper, Available at SSRN 3265007* .
- Gentzkow M, Kelly B, Taddy M. 2019. Text as data. *Journal of Economic Literature* **57**: 535–574.
- Griffin JM, Shams A. 2020. Is bitcoin really un-tethered? *Journal of Finance, forthcoming* .
- Griffiths TL, Steyvers M. 2004. Finding scientific topics. *Proceedings of the National Academy of Sciences* **101**: 5228–5235.
- Gu S, Kelly B, Xiu D. 2020. Empirical asset pricing via machine learning. *Review of Financial Studies, forthcoming* .
- Hansen S, McMahon M, Prat A. 2018. Transparency and deliberation within the fomc: a computational linguistics approach. *Quarterly Journal of Economics* **133**: 801–870.

- Hinzen FJ, John K, Saleh F. 2019. Proof-of-work's limited adoption problem. *Working Paper, New York University* .
- Howell ST, Niessner M, Yermack D. 2020. Initial coin offerings: Financing growth with cryptocurrency token sales. *Review of Financial Studies, forthcoming* .
- Lee J, Li T, Shin D. 2019. The wisdom of crowds in fintech: Evidence from initial coin offerings. *Working paper, University of Florida* .
- Li K, Mai F, Shen R, Yan X. 2019. Measuring corporate culture using machine learning. *Available at SSRN 3256608* .
- Liu Y, Matthies B. 2018. Long run risk: Is it there? *Working paper, Yale University* .
- Liu Y, Tsyvinski A. 2018. Risks and returns of cryptocurrency. *Working paper, Yale University and NBER* .
- Liu Y, Tsyvinski A, Wu X. 2019. Common risk factors in cryptocurrency. *Working paper, Yale University and NBER* .
- Loughran T, McDonald B. 2011. When is a liability not a liability? textual analysis, dictionaries, and 10-ks. *Journal of Finance* **66**: 35–65.
- Loughran T, Ritter JR. 1995. The new issues puzzle. *Journal of Finance* **50**: 23–51.
- Lyandres E, Palazzo B, Rabetti D. 2020. Ico success and post-ico performance. *Working Paper* .
- Manela A, Moreira A. 2017. News implied volatility and disaster concerns. *Journal of Financial Economics*. **123**: 137–162.
- Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- Murtagh F, Legendre P. 2014. Ward's hierarchical agglomerative clustering method: which algorithms implement ward's criterion? *Journal of Classification* **31**: 274–295.
- Pastor L, Veronesi P. 2003. Stock valuation and learning about profitability. *Journal of Finance* **58**: 1749–1789.
- Ritter JR. 1991. The long-run performance of initial public offerings. *Journal of Finance* **46**: 3–27.
- Röder M, Both A, Hinneburg A. 2015. Exploring the space of topic coherence measures. In *Proceedings of the eighth ACM international conference on Web search and data mining*. 399–408.
- Russell SJ, Norvig P. 2010. *Artificial Intelligence-A Modern Approach (3rd internat. edn.)*. Pearson Education.

- Satopaa V, Albrecht J, Irwin D, Raghavan B. 2011. Finding a "kneedle" in a haystack: Detecting knee points in system behavior. In *2011 31st international conference on distributed computing systems workshops*. IEEE, 166–171.
- Shams A. 2019. What drives the covariation of cryptocurrency returns? *Working paper, Ohio State University* .
- Sheng J. 2019. Asset pricing in the information age: Employee expectations and stock returns. *Working paper, University of California Irvine* .
- Shumway T. 1997. The delisting bias in crsp data. *Journal of Finance* **52**: 327–340.
- Sievert C, Shirley K. 2014. Ldavis: A method for visualizing and interpreting topics. In *Proceedings of the workshop on interactive language learning, visualization, and interfaces*. 63–70.
- Sockin M, Xiong W. 2018. A model of cryptocurrencies. *Working paper, Princeton University* .
- Taddy M. 2012. On estimation and selection for topic models. In *Artificial Intelligence and Statistics*. 1184–1193.
- Tetlock PC. 2014. Information transmission in finance. *Annual Review Financial Economics* **6**: 365–384.
- Torgerson WS. 1958. *Theory and Methods of Scaling*. Wiley.
- Yermack D. 2017. Corporate governance and blockchains. *Review of Finance* **21**: 7–31.

## Appendix: Variable Definition

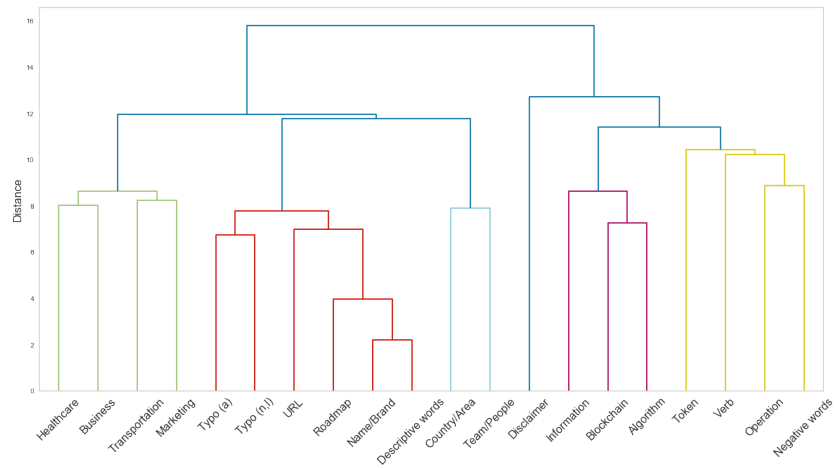
Variable	Definition
<b>ICO Success Measures:</b>	
CMC Trading	A dummy variable that equals to one if a cryptocurrency is shown as listed on coinmarketcap.com (CMC).
Trading	A self-reported dummy by ICO fundraisers about whether the cryptocurrency is traded on an exchange.
Success	A dummy variable indicating whether the ICO raises any capital.
<b>Trading Variables:</b>	
First Open/ICO Price	The ratio between the first day's opening price and the ICO price.
Delist	An indicator about whether a token is delisted from CMC.
Rate of Return	The rate of return that investors earn if they buy cryptocurrency at the opening price on the first listing day and sell them after a certain holding period.
Trading Volume	The 24-hour trading volume in millions of USD after they have been listed on CMC for a certain period of time.
<b>Whitepaper Measures:</b>	
Tech_sup	The normalized predicted technology score from partial least squares (PLS), a supervised machine learning approach.
Tech_embed	The normalized percentage of words in the "blockchain", "information" or "algorithm" topics of the word embedding and clustering approach.
Tech_lda	The normalized proportional attention allocated to the "information", "blockchain" and "system" topics of the LDA topic modelling approach.
Tech_comp	The simple average of the Tech_sup, Tech_embed and Tech_lda.
Fog Index	A readability measure defined as $0.4[(words/sentences) + 100((complexwords)/words)]$ , where "complex words" are words with three or more syllables.
Tone	The difference between number of positive and negative words defined in <a href="#">Loughran and McDonald (2011)</a> divided by the total number of words in a whitepaper.
Uncertainty	The number of uncertainty words defined in <a href="#">Loughran and McDonald (2011)</a> divided by the total number of words in a whitepaper.
<b>ICO characteristics:</b>	
Has GitHub	A dummy variable that equals to one if the ICO project has a GitHub homepage.
Has Twitter	A dummy variable that equals to one if the ICO project has a Twitter account.
ICO Length	The number of days from the start to the end of an ICO campaign.
Team Size	The number of ICO team members.
Pre ICO	A dummy variable indicating whether if a pre-ICO exists.
Bonus	A dummy variable indicating whether the fundraiser offers bonus to investors.
Ethereum Based	A dummy variable indicating whether the ICO project is built on Ethereum.
Accept BTC	A dummy variable indicating whether the ICO accepts Bitcoin as a currency of payment.
BTC Price (ICO)	The price of Bitcoin in thousands of US dollars on the day an ICO initiates.
BTC Price (List)	The price of Bitcoin in thousands of US dollars on the day an ICO is shown as listed on CMC.

# Figures & Tables

Figure 1: **Embedding Visualization**

This figure plots the relationship between embedding-based topics. Panel (a) displays the taxonomy generated by hierarchical agglomerative clustering. Panel (b) shows the similarity between topics in a two-dimensional space. The size of the circle represents the percentage of terms belonging to the topic. “Information”, “blockchain” and “algorithm” are used to construct the embedding-based tech index.

(a) Taxonomy



(b) Multidimensional scaling (MDS)

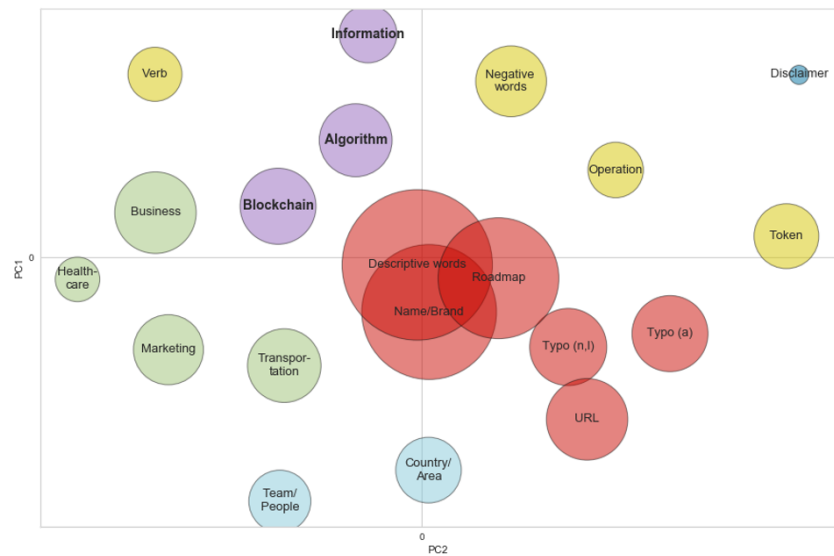
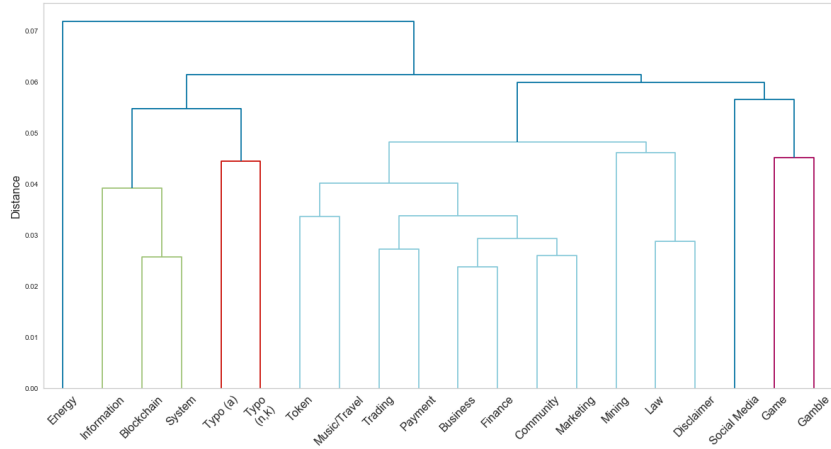




Figure 2: LDA Visualization

This figure plots the relationship between LDA-based topics. Panel (a) displays the taxonomy generated by hierarchical agglomerative clustering. Panel (b) shows the similarity between topics in a two-dimensional space. The size of the circle represents the relative topic prevalence in the corpus. “Information”, “blockchain” and “system” are used to construct the LDA-based tech index.

(a) Taxonomy



(b) Multidimensional scaling (MDS)

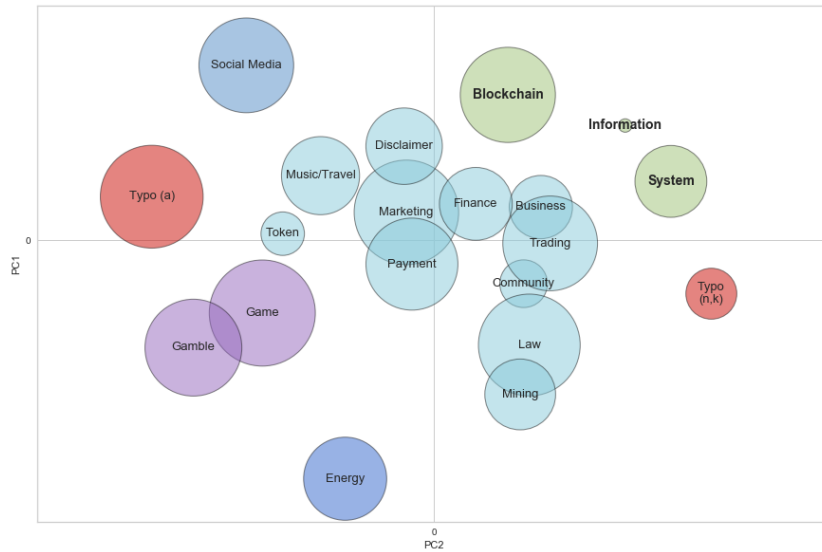


Table 1: Summary Statistics

This table presents summary statistics on variables related to ICO characteristics, outcomes and whitepaper measures. Panel A shows descriptive statistics for 2,916 ICOs completed before December 31st, 2018. Panel B summarizes a subsample of 765 ICOs listed on coinmarketcap.com. For each variable, we show the number of non-missing observations, the mean, the standard deviation and the 10th, 50th and 90th percentile values. Please refer to the “variable definition” in the Appendix for the definition of each variable.

<b>Panel A: Full Sample</b>						
	Obs.	Mean	SD	p10	p50	p90
<i>ICO Success Measures</i>						
CMC Trading	2916	0.26	0.44	0	0	1
Trading	2916	0.18	0.39	0	0	1
Success	2916	0.38	0.49	0	0	1
<i>Whitepaper Measures</i>						
Tech_sup	1629	0	1.00	-0.96	-0.19	1.44
Tech_embed	1629	0	1.00	-1.02	-0.23	1.35
Tech_lda	1629	0	1.00	-0.62	-0.49	1.63
Tech_comp	1629	0	0.84	-0.81	-0.25	1.20
Fog Index	1629	16.7	12.6	13.2	15.7	18.5
Tone	1629	0.28	0.73	-0.58	0.29	1.10
Uncertainty	1629	0.75	0.39	0.35	0.67	1.25
<i>ICO Characteristics</i>						
Has GitHub	2916	0.60	0.49	0	1	1
Has Twitter	2916	0.91	0.29	1	1	1
ICO Length	2683	50.7	45.8	14	32	100
Team Size	2916	11.0	7.05	3	10	20
Pre ICO	2916	0.51	0.50	0	1	1
Bonus	2916	0.20	0.40	0	0	1
Ethereum Based	2916	0.83	0.37	0	1	1
Accept BTC	2916	0.40	0.49	0	0	1
BTC Price (ICO)	2669	7.80	3.07	4.23	7.28	11.3
BTC Price (List)	710	7.59	3.70	2.73	7.03	13.5
ICO Price	1684	1.57	17.8	0.01	0.10	1

<b>Panel B: Listed Sample</b>						
	Obs.	Mean	SD	p10	p50	p90
<i>Trading Variables</i>						
First Open/ICO Price	413	2.20	4.66	0.16	0.97	3.79
Delist	765	0.10	0.30	0	0	1
<i>Rate of Return</i>						
7 Days	741	0.19	0.83	-0.45	-0.04	1.03
30 Days	730	0.30	1.84	-0.71	-0.28	1.60
90 Days	686	0.65	3.11	-0.87	-0.43	3.21
180 Days	566	1.00	5.14	-0.95	-0.64	4.00
210 Days	530	0.84	4.27	-0.96	-0.68	3.57
240 Days	486	0.69	3.89	-0.96	-0.70	3.20
270 Days	438	1.46	8.40	-0.96	-0.72	3.69
300 Days	397	1.51	8.91	-0.97	-0.74	3.85
330 Days	356	1.34	8.21	-0.98	-0.72	3.31
360 Days	289	1.60	7.74	-0.96	-0.67	4.20
<i>Trading Volume (\$ MIL)</i>						
Listing Days	751	2.40	11.9	0.0023	0.12	3.90
7 Days	739	1.63	5.58	0.0015	0.083	3.60
30 Days	725	1.50	5.62	0.0011	0.066	2.53
90 Days	680	1.60	5.58	0.00045	0.11	3.22
180 Days	564	2.78	13.5	0.00039	0.069	3.99
210 Days	526	2.24	8.30	0.00020	0.065	3.90
240 Days	482	2.60	12.4	0.00023	0.048	3.31
270 Days	436	1.73	5.91	0.00016	0.061	3.09
300 Days	393	2.60	11.5	0.00025	0.058	3.50
330 Days	352	2.22	8.48	0.00021	0.067	3.19
360 Days	285	2.45	9.65	0.000048	0.084	3.39
<i>Whitepaper Measures</i>						
Tech_sup	422	0.27	1.13	-0.92	0.0014	1.88
Tech_embed	422	0.41	1.17	-0.79	0.14	2.28
Tech_lda	422	0.33	1.23	-0.62	-0.26	2.51
Tech_comp	422	0.34	1.03	-0.73	0.033	1.87
Fog Index	422	17.2	18.9	13.3	15.5	18.3
Tone	422	0.20	0.72	-0.70	0.23	1.03
Uncertainty	422	0.79	0.40	0.35	0.71	1.30
<i>ICO Characteristics</i>						
Has GitHub	765	0.70	0.46	0	1	1
Has Twitter	765	0.96	0.19	1	1	1
ICO Length	656	34.9	42.0	2	30	63
Team Size	765	12.1	8.02	3	11	22
Pre ICO	765	0.26	0.44	0	0	1
Bonus	765	0.075	0.26	0	0	0
Ethereum Based	765	0.80	0.40	0	1	1
Accept BTC	765	0.30	0.46	0	0	1
BTC Price (ICO)	642	7.45	3.94	2.54	7.10	13.8
BTC Price (List)	710	7.59	3.70	2.73	7.03	13.5
ICO Price	420	2.37	19.9	0.01	0.12	1.22

Table 2: **Technology Indexes Determinant**

This table presents the determinants of our tech index. The dependent variable is the composite tech index (*Tech comp*). Column (1) links the tech index to whether an ICO uses Ethereum blockchain; column (2) presents the relation between the tech index and GitHub commits (the number of code revisions); column (3) considers other text-based measures of ICO whitepapers; column (4) presents estimates with ICO characteristics; column (5) includes all variables. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

	(1)	(2)	(3)	(4)	(5)
Ethereum Based	-0.219*** (0.066)				-0.122** (0.061)
Ln(commits)		0.104*** (0.011)			0.071*** (0.011)
Has GitHub		-0.115** (0.050)			-0.058 (0.050)
Fog Index			-0.003*** (0.001)		-0.002* (0.001)
Tone			-0.265*** (0.031)		-0.226*** (0.030)
Uncertainty			-0.186*** (0.054)		-0.217*** (0.050)
ICO Length				-0.002*** (0.001)	-0.001*** (0.000)
Team Size				0.015*** (0.003)	0.012*** (0.003)
Has Twitter				0.212** (0.085)	0.156* (0.084)
BTC Price (ICO)				-0.018** (0.008)	-0.010 (0.008)
Pre ICO				-0.040 (0.044)	-0.004 (0.042)
Bonus				-0.077* (0.047)	-0.059 (0.045)
Accept BTC				-0.095** (0.043)	-0.064 (0.041)
Constant	0.184*** (0.063)	-0.156*** (0.031)	0.260*** (0.055)	-0.094 (0.112)	0.127 (0.125)
$R^2$	0.009	0.098	0.047	0.048	0.136
Observations	1629	1629	1629	1483	1483

Table 3: **Technology Indexes**

This table presents results related to the construction of tech indexes. Panel A shows the correlation between the four tech indexes. Panel B compares various supervised machine learning methods with their out-of-sample (OOS)  $R^2$  and corresponding hyperparameters.

<b>Panel A: Correlation Matrix of Technology Indexes</b>									
	Tech_sup	Tech_embed	Tech_lda	Tech_comp					
Tech_sup	1.0000								
Tech_embed	0.5152	1.0000							
Tech_lda	0.4861	0.6838	1.0000						
Tech_comp	0.7929	0.8713	0.8597	1.0000					

<b>Panel B: Comparison of Different Supervised ML Methods</b>									
	OLS	LASSO	Ridge	ElNet	PCR	PLS	RF	GB	NN
Hyperparameter	—	$\lambda = 1.5$	$\lambda = 1.75$	$\alpha = 0.9$	$PC = 4$	$PC = 2$	$tree = 20$	$tree = 50$	$node = 50$
OOS $R^2$ (%)	27.14	30.02	27.15	30.08	35.40	<b>45.88</b>	37.91	32.53	-7.71

Table 4: ICO Success

This table examines the relationship between tech indexes and ICO success. The dependent variable is *CMC Trading* in Panel A and *Success* in Panel B. For each tech index, the first column presents the univariate result, and the second column displays estimates with control variables and fixed effects. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: CMC Trading</b>								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Supervised		Embedding		LDA		Composite	
Tech_sup	0.070*** (0.012)	0.039*** (0.012)						
Tech_embed			0.107*** (0.011)	0.048*** (0.012)				
Tech_lda					0.086*** (0.012)	0.047*** (0.013)		
Tech_comp							0.124*** (0.013)	0.066*** (0.015)
ICO Length		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)
Team Size		0.008*** (0.001)		0.009*** (0.001)		0.009*** (0.001)		0.008*** (0.001)
Has GitHub		0.053** (0.021)		0.044** (0.021)		0.048** (0.021)		0.045** (0.021)
Has Twitter		0.163*** (0.035)		0.172*** (0.036)		0.168*** (0.035)		0.166*** (0.035)
BTC Price (ICO)		0.010 (0.007)		0.010 (0.006)		0.010 (0.006)		0.011 (0.006)
Pre ICO		-0.035 (0.023)		-0.031 (0.023)		-0.031 (0.023)		-0.032 (0.023)
Bonus		0.008 (0.022)		0.010 (0.022)		0.007 (0.022)		0.009 (0.022)
Accept BTC		-0.013 (0.021)		-0.011 (0.021)		-0.010 (0.021)		-0.010 (0.021)
Ethereum Based		-0.017 (0.030)		-0.009 (0.030)		-0.012 (0.030)		-0.009 (0.030)
Fog Index		-0.000 (0.001)		0.000 (0.001)		-0.000 (0.001)		0.000 (0.001)
Tone		-0.000 (0.014)		0.003 (0.014)		0.004 (0.015)		0.006 (0.015)
Uncertainty		0.014 (0.028)		0.031 (0.028)		0.023 (0.028)		0.025 (0.028)
Constant	0.259*** (0.011)	0.610*** (0.094)	0.259*** (0.011)	0.521*** (0.099)	0.259*** (0.011)	0.503*** (0.099)	0.259*** (0.011)	0.511*** (0.097)
Fixed Effects	No	Yes	No	Yes	No	Yes	No	Yes
R <sup>2</sup>	0.026	0.322	0.060	0.324	0.038	0.323	0.057	0.327
Observations	1629	1382	1629	1382	1629	1382	1629	1382

<b>Panel B: Capital Raised &gt; 0</b>								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Supervised		Embedding		LDA		Composite	
Tech_sup	0.061*** (0.012)	0.041*** (0.013)						
Tech_embed			0.077*** (0.012)	0.043*** (0.014)				
Tech_lda					0.056*** (0.012)	0.037*** (0.014)		
Tech_comp							0.091*** (0.014)	0.060*** (0.016)
ICO Length		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)
Team Size		0.008*** (0.002)		0.009*** (0.002)		0.009*** (0.002)		0.008*** (0.002)
Has GitHub		0.089*** (0.025)		0.082*** (0.025)		0.086*** (0.025)		0.082*** (0.025)
Has Twitter		0.083 (0.053)		0.092* (0.053)		0.089* (0.054)		0.086 (0.053)
BTC Price (ICO)		-0.005 (0.007)		-0.005 (0.007)		-0.005 (0.007)		-0.004 (0.007)
Pre ICO		-0.012 (0.029)		-0.008 (0.028)		-0.009 (0.028)		-0.010 (0.028)
Bonus		0.103*** (0.031)		0.104*** (0.031)		0.102*** (0.031)		0.104*** (0.031)
Accept BTC		0.068*** (0.024)		0.070*** (0.024)		0.070*** (0.024)		0.071*** (0.024)
Ethereum Based		-0.014 (0.033)		-0.007 (0.034)		-0.011 (0.034)		-0.007 (0.033)
Fog Index		-0.001 (0.001)		-0.000 (0.001)		-0.001 (0.001)		-0.000 (0.001)
Tone		0.012 (0.018)		0.013 (0.018)		0.014 (0.018)		0.017 (0.018)
Uncertainty		0.066** (0.033)		0.081** (0.033)		0.073** (0.033)		0.076** (0.033)
Constant	0.377*** (0.012)	0.633*** (0.121)	0.377*** (0.012)	0.557*** (0.124)	0.377*** (0.012)	0.554*** (0.126)	0.377*** (0.012)	0.545*** (0.123)
Fixed Effects	No	Yes	No	Yes	No	Yes	No	Yes
R <sup>2</sup>	0.016	0.256	0.025	0.256	0.013	0.254	0.025	0.258
Observations	1629	1382	1629	1382	1629	1382	1629	1382

Table 5: ICO Success—Industry Subsample

This table examines the relationship between tech indexes and ICO success for different technology-related industries. The dependent variable is *CMC Trading. Industry* is a dummy that equals to 1 if the ICO belongs to “platform”, “cryptocurrency”, and “trading” industries. For each tech index, the first column presents the univariate result, and the second column displays estimates with control variables and fixed effects. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

	Supervised		Embedding		LDA		Composite	
Tech_sup	0.052*** (0.017)	0.030* (0.016)						
Tech_sup*Industry	0.036 (0.023)	0.017 (0.021)						
Tech_embed			0.085*** (0.017)	0.028* (0.016)				
Tech_embed*Industry			0.040* (0.022)	0.035* (0.021)				
Tech_lda					0.053*** (0.016)	0.012 (0.015)		
Tech_lda*Industry					0.069*** (0.023)	0.049** (0.022)		
Tech_comp							0.089*** (0.020)	0.034* (0.019)
Tech_comp*Industry							0.067** (0.027)	0.049* (0.026)
Industry	-0.008 (0.021)	0.010 (0.020)	-0.009 (0.021)	0.011 (0.020)	0.003 (0.021)	0.016 (0.020)	-0.004 (0.021)	0.014 (0.020)
ICO Length		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)
Team Size		0.008*** (0.001)		0.009*** (0.001)		0.009*** (0.001)		0.008*** (0.001)
Has GitHub		0.053*** (0.021)		0.044** (0.021)		0.048** (0.021)		0.045** (0.021)
Has Twitter		0.160*** (0.036)		0.169*** (0.036)		0.169*** (0.036)		0.163*** (0.035)
BTC Price (ICO)		0.010 (0.007)		0.010 (0.006)		0.010 (0.007)		0.010 (0.006)
Pre ICO		-0.036 (0.023)		-0.030 (0.023)		-0.029 (0.023)		-0.032 (0.023)
Bonus		0.009 (0.021)		0.009 (0.022)		0.007 (0.022)		0.009 (0.022)
Accept BTC		-0.017 (0.020)		-0.013 (0.020)		-0.015 (0.020)		-0.013 (0.020)
Ethereum Based		-0.012 (0.030)		-0.002 (0.030)		-0.007 (0.030)		-0.004 (0.030)
Fog Index		-0.000 (0.001)		-0.000 (0.002)		-0.000 (0.001)		-0.000 (0.001)
Tone		0.001 (0.014)		0.004 (0.014)		0.002 (0.014)		0.006 (0.014)
Uncertainty		0.025 (0.028)		0.043 (0.028)		0.033 (0.028)		0.037 (0.028)
Constant	0.263*** (0.016)	0.579*** (0.090)	0.264*** (0.016)	0.468*** (0.097)	0.259*** (0.016)	0.447*** (0.098)	0.262*** (0.016)	0.455*** (0.096)
Fixed effects	No	Yes	No	Yes	No	Yes	No	Yes
R2	0.028	0.309	0.062	0.313	0.044	0.310	0.061	0.314
Observations	1629	1382	1629	1382	1629	1382	1629	1382



Table 6: **Rate of Return**

This table presents the effects of tech indexes on cryptocurrency returns. The dependent variable is the log transformation of gross return over a given period. Panel A, B, C and D display the supervised, embedding-based, LDA-based and composite tech index respectively. Column (1)-(6) display results for six horizons: 7 days, 30 days, 90 days, 180 days, 240 days and 300 days. We include control variables related to ICO characteristics and whitepapers in all columns. Quarterly, categorical and geographical fixed effects are considered under all circumstances. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Supervised Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.013 (0.031)	0.036 (0.062)	-0.001 (0.090)	-0.031 (0.106)	0.019 (0.120)	-0.005 (0.142)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.003 (0.004)	-0.003 (0.005)	0.003 (0.008)
Team Size	-0.002 (0.004)	0.003 (0.007)	0.005 (0.010)	0.001 (0.013)	0.002 (0.016)	-0.001 (0.018)
Has GitHub	0.033 (0.076)	0.017 (0.135)	0.178 (0.176)	0.282 (0.239)	0.351 (0.259)	0.601 (0.363)
Has Twitter	0.066 (0.148)	0.028 (0.555)	0.025 (0.711)	-0.016 (0.840)	0.067 (0.755)	0.165 (0.820)
BTC Price (ICO)	-0.000 (0.013)	-0.030 (0.019)	-0.061** (0.026)	-0.098*** (0.029)	-0.084** (0.032)	-0.093** (0.044)
Pre ICO	-0.012 (0.076)	-0.169 (0.131)	-0.023 (0.176)	0.114 (0.284)	-0.118 (0.356)	0.471 (0.575)
Constant	0.114 (0.654)	-0.359 (1.292)	-1.298 (0.896)	-0.381 (1.218)	-0.400 (1.203)	-0.807 (1.425)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.076	0.162	0.238	0.351	0.396	0.362
Observations	316	310	286	218	184	140
<b>Panel B: Embedding Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_embed	0.032 (0.034)	0.080 (0.059)	0.123 (0.083)	0.195* (0.105)	0.331*** (0.117)	0.377** (0.150)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.003 (0.004)	-0.003 (0.004)	0.004 (0.007)
Team Size	-0.001 (0.004)	0.003 (0.007)	0.004 (0.010)	-0.003 (0.013)	-0.001 (0.015)	-0.003 (0.018)
Has GitHub	0.026 (0.075)	0.000 (0.134)	0.136 (0.172)	0.201 (0.235)	0.189 (0.249)	0.332 (0.374)
Has Twitter	0.063 (0.148)	0.019 (0.569)	0.031 (0.716)	-0.020 (0.847)	0.066 (0.784)	0.175 (0.812)
BTC Price (ICO)	0.001 (0.013)	-0.027 (0.019)	-0.056** (0.026)	-0.089*** (0.030)	-0.067** (0.034)	-0.074 (0.047)
Pre ICO	0.004 (0.079)	-0.127 (0.135)	0.038 (0.178)	0.248 (0.292)	0.080 (0.366)	0.776 (0.591)
Constant	0.059 (0.658)	-0.492 (1.292)	-1.530* (0.896)	-0.701 (1.215)	-0.722 (1.164)	-1.305 (1.319)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.078	0.167	0.246	0.364	0.430	0.403
Observations	316	310	286	218	184	140

<b>Panel C: LDA Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_lda	0.045 (0.034)	0.066 (0.056)	0.124 (0.076)	0.136 (0.086)	0.241** (0.098)	0.238* (0.133)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.003 (0.004)	-0.003 (0.005)	0.003 (0.008)
Team Size	-0.002 (0.004)	0.003 (0.007)	0.004 (0.010)	-0.003 (0.013)	-0.001 (0.016)	-0.003 (0.018)
Has GitHub	0.027 (0.077)	0.013 (0.136)	0.154 (0.176)	0.247 (0.241)	0.269 (0.259)	0.473 (0.376)
Has Twitter	0.063 (0.153)	0.019 (0.587)	0.022 (0.740)	-0.026 (0.862)	0.057 (0.818)	0.146 (0.854)
BTC Price (ICO)	-0.000 (0.012)	-0.029 (0.019)	-0.059** (0.026)	-0.096*** (0.029)	-0.082** (0.033)	-0.090** (0.044)
Pre ICO	0.004 (0.078)	-0.146 (0.134)	0.020 (0.178)	0.203 (0.291)	-0.007 (0.354)	0.534 (0.577)
Constant	-0.033 (0.665)	-0.555 (1.322)	-1.713* (0.938)	-0.830 (1.225)	-1.077 (1.179)	-1.429 (1.384)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.081	0.165	0.247	0.358	0.414	0.378
Observations	316	310	286	218	184	140
<b>Panel D: Composite Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_comp	0.041 (0.038)	0.083 (0.067)	0.113 (0.096)	0.144 (0.113)	0.280** (0.130)	0.284* (0.162)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.003 (0.004)	-0.003 (0.005)	0.003 (0.007)
Team Size	-0.002 (0.004)	0.003 (0.007)	0.003 (0.010)	-0.003 (0.013)	-0.003 (0.015)	-0.006 (0.018)
Has GitHub	0.027 (0.076)	0.006 (0.134)	0.153 (0.174)	0.241 (0.238)	0.259 (0.255)	0.467 (0.374)
Has Twitter	0.066 (0.146)	0.025 (0.563)	0.031 (0.708)	-0.013 (0.832)	0.076 (0.759)	0.163 (0.801)
BTC Price (ICO)	0.000 (0.012)	-0.029 (0.019)	-0.058** (0.026)	-0.095*** (0.029)	-0.079** (0.033)	-0.088* (0.045)
Pre ICO	-0.000 (0.078)	-0.145 (0.134)	0.010 (0.177)	0.187 (0.290)	-0.012 (0.359)	0.588 (0.587)
Constant	0.038 (0.657)	-0.501 (1.290)	-1.519* (0.891)	-0.627 (1.193)	-0.748 (1.143)	-1.199 (1.332)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.079	0.166	0.243	0.356	0.413	0.379
Observations	316	310	286	218	184	140

Table 7: **Bitcoin-Adjusted Rate of Return**

This table presents the effects of tech indexes on Bitcoin-adjusted returns. The dependent variable is the log transformation of gross return,  $\log(1+ROR)$ , minus the log transformation of Bitcoin gross return over the same period. Panel A, B, C and D display the supervised, embedding-based, LDA-based and composite tech index respectively. Column (1)-(6) display results for six different horizons: 7 days, 30 days, 90 days, 180 days, 240 days and 300 days. We include control variables related to ICO characteristics and whitepapers in all columns. Quarterly, categorical and geographical fixed effects are considered under all circumstances. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Supervised Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.012 (0.029)	0.037 (0.057)	0.011 (0.085)	0.061 (0.095)	0.091 (0.111)	0.058 (0.133)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.002 (0.003)	-0.002 (0.004)	0.004 (0.007)
Team Size	-0.001 (0.004)	0.002 (0.006)	0.008 (0.009)	0.004 (0.011)	0.006 (0.014)	-0.000 (0.017)
Has GitHub	0.053 (0.068)	-0.034 (0.126)	0.111 (0.162)	0.221 (0.199)	0.249 (0.241)	0.424 (0.331)
Has Twitter	0.191 (0.178)	-0.003 (0.650)	0.145 (0.912)	0.527 (0.881)	0.348 (0.788)	0.342 (0.864)
BTC Price (ICO)	0.009 (0.012)	-0.009 (0.018)	-0.034 (0.023)	-0.079*** (0.027)	-0.061* (0.031)	-0.060 (0.045)
Pre ICO	-0.017 (0.071)	-0.166 (0.113)	-0.005 (0.171)	0.062 (0.240)	-0.075 (0.310)	0.333 (0.492)
Constant	-0.056 (0.620)	-0.500 (1.206)	-1.873* (1.031)	-2.879** (1.227)	-3.431*** (1.135)	-2.894** (1.366)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.097	0.156	0.169	0.317	0.284	0.261
Observations	311	305	281	213	180	137
<b>Panel B: Embedding Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_embed	0.031 (0.033)	0.097* (0.056)	0.122 (0.080)	0.216** (0.094)	0.319*** (0.109)	0.407*** (0.144)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.000 (0.001)	-0.002 (0.003)	-0.002 (0.004)	0.005 (0.005)
Team Size	-0.001 (0.004)	0.002 (0.006)	0.007 (0.009)	0.003 (0.011)	0.006 (0.014)	0.001 (0.016)
Has GitHub	0.046 (0.067)	-0.057 (0.125)	0.072 (0.160)	0.143 (0.196)	0.098 (0.234)	0.138 (0.340)
Has Twitter	0.193 (0.178)	-0.001 (0.671)	0.168 (0.928)	0.548 (0.912)	0.400 (0.844)	0.483 (0.839)
BTC Price (ICO)	0.010 (0.012)	-0.005 (0.018)	-0.030 (0.024)	-0.069** (0.027)	-0.045 (0.031)	-0.037 (0.048)
Pre ICO	-0.002 (0.075)	-0.113 (0.116)	0.060 (0.176)	0.194 (0.250)	0.100 (0.324)	0.601 (0.515)
Constant	-0.114 (0.624)	-0.679 (1.207)	-2.119** (1.038)	-3.265*** (1.237)	-3.829*** (1.127)	-3.579*** (1.257)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.100	0.164	0.179	0.338	0.326	0.324
Observations	311	305	281	213	180	137

<b>Panel C: LDA Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_lda	0.040 (0.033)	0.066 (0.056)	0.119 (0.073)	0.136* (0.077)	0.226** (0.094)	0.290** (0.126)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.003 (0.003)	-0.002 (0.004)	0.004 (0.006)
Team Size	-0.001 (0.004)	0.002 (0.006)	0.007 (0.010)	0.004 (0.011)	0.005 (0.014)	0.000 (0.016)
Has GitHub	0.048 (0.068)	-0.039 (0.127)	0.090 (0.162)	0.199 (0.199)	0.175 (0.244)	0.274 (0.342)
Has Twitter	0.189 (0.179)	-0.015 (0.692)	0.143 (0.958)	0.511 (0.955)	0.357 (0.895)	0.420 (0.928)
BTC Price (ICO)	0.009 (0.012)	-0.008 (0.018)	-0.034 (0.024)	-0.078*** (0.027)	-0.059* (0.031)	-0.054 (0.046)
Pre ICO	-0.003 (0.073)	-0.143 (0.116)	0.041 (0.174)	0.138 (0.245)	0.020 (0.312)	0.360 (0.492)
Constant	-0.188 (0.630)	-0.693 (1.244)	-2.278** (1.082)	-3.331*** (1.271)	-4.117*** (1.166)	-3.765*** (1.352)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.102	0.159	0.179	0.325	0.304	0.292
Observations	311	305	281	213	180	137
<b>Panel D: Composite Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_comp	0.038 (0.037)	0.091 (0.065)	0.116 (0.092)	0.193* (0.102)	0.300** (0.122)	0.353** (0.154)
ICO Length	0.000 (0.001)	0.000 (0.001)	-0.001 (0.001)	-0.002 (0.003)	-0.002 (0.004)	0.004 (0.006)
Team Size	-0.001 (0.004)	0.001 (0.006)	0.006 (0.009)	0.002 (0.011)	0.002 (0.014)	-0.003 (0.016)
Has GitHub	0.047 (0.067)	-0.047 (0.125)	0.088 (0.160)	0.180 (0.198)	0.155 (0.239)	0.263 (0.340)
Has Twitter	0.195 (0.177)	0.001 (0.663)	0.163 (0.918)	0.540 (0.899)	0.396 (0.818)	0.443 (0.832)
BTC Price (ICO)	0.009 (0.012)	-0.007 (0.018)	-0.033 (0.024)	-0.075*** (0.027)	-0.056* (0.031)	-0.052 (0.047)
Pre ICO	-0.006 (0.073)	-0.138 (0.116)	0.033 (0.174)	0.143 (0.245)	0.025 (0.316)	0.428 (0.508)
Constant	-0.131 (0.624)	-0.666 (1.205)	-2.111** (1.029)	-3.226*** (1.209)	-3.875*** (1.098)	-3.497*** (1.254)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.100	0.161	0.176	0.329	0.309	0.295
Observations	311	305	281	213	180	137

Table 8: Trading Volume

This table presents the relationship between tech indexes and cryptocurrency liquidity. The dependent variable is the log transformation of 24-hour trading volume in USD. Column (1) displays results on the listing day. Column (2) to (7) display results for six time points: 7 days, 30 days, 90 days, 180 days, 240 days and 300 days. We include control variables in all columns. Quarterly, categorical and geographical fixed effects are considered under all circumstances. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Supervised Index</b>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Listing	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.334** (0.165)	0.333** (0.166)	0.346* (0.198)	0.501** (0.204)	0.491* (0.253)	0.464* (0.251)	0.344 (0.286)
ICO Length	-0.005 (0.005)	-0.006 (0.006)	-0.004 (0.005)	-0.006 (0.008)	-0.006 (0.010)	-0.001 (0.009)	-0.004 (0.018)
Team Size	0.014 (0.021)	0.017 (0.018)	0.027 (0.020)	0.022 (0.020)	0.074*** (0.027)	0.047 (0.029)	0.083** (0.036)
Has GitHub	0.263 (0.361)	0.510 (0.386)	0.444 (0.462)	0.805 (0.512)	0.889* (0.536)	1.272** (0.575)	1.580** (0.738)
Has Twitter	-0.487 (1.162)	-0.046 (1.073)	-0.676 (1.648)	-1.101 (1.614)	0.697 (2.080)	-0.431 (1.508)	-1.184 (1.557)
BTC Price (ICO)	0.073 (0.069)	0.096 (0.065)	-0.039 (0.072)	0.012 (0.072)	-0.049 (0.087)	-0.053 (0.087)	0.036 (0.108)
Pre ICO	-0.399 (0.411)	-0.487 (0.474)	-0.453 (0.501)	-1.005* (0.575)	-1.335* (0.749)	-2.061** (0.836)	-0.553 (1.238)
Constant	11.971*** (1.650)	14.015*** (1.722)	11.312*** (2.339)	9.960*** (2.125)	5.065 (3.157)	6.954*** (2.558)	9.246*** (2.992)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.195	0.221	0.172	0.210	0.342	0.433	0.428
Observations	323	316	308	283	217	183	139
<b>Panel B: Embedding Index</b>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Listing	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_embed	0.514*** (0.156)	0.490*** (0.153)	0.489*** (0.170)	0.729*** (0.172)	0.625** (0.262)	0.798*** (0.261)	0.780** (0.342)
ICO Length	-0.004 (0.004)	-0.005 (0.006)	-0.003 (0.005)	-0.005 (0.007)	-0.005 (0.009)	-0.000 (0.008)	-0.002 (0.016)
Team Size	0.019 (0.019)	0.023 (0.018)	0.033* (0.020)	0.030 (0.019)	0.085*** (0.026)	0.056** (0.026)	0.094** (0.037)
Has GitHub	0.154 (0.359)	0.417 (0.382)	0.360 (0.459)	0.626 (0.504)	0.736 (0.543)	0.954* (0.563)	1.068 (0.750)
Has Twitter	-0.545 (1.251)	-0.112 (1.207)	-0.753 (1.783)	-1.140 (1.799)	0.598 (2.187)	-0.500 (1.615)	-1.203 (1.473)
BTC Price (ICO)	0.090 (0.070)	0.113* (0.065)	-0.021 (0.073)	0.034 (0.073)	-0.020 (0.090)	-0.010 (0.092)	0.079 (0.118)
Pre ICO	-0.164 (0.429)	-0.246 (0.496)	-0.198 (0.528)	-0.652 (0.595)	-0.989 (0.753)	-1.602** (0.790)	-0.063 (1.178)
Constant	11.220*** (1.715)	13.273*** (1.785)	10.602*** (2.410)	8.893*** (2.197)	4.161 (3.229)	6.149** (2.547)	8.349*** (2.894)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.213	0.234	0.184	0.234	0.352	0.457	0.456
Observations	323	316	308	283	217	183	139

<b>Panel C: LDA Index</b>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Listing	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_lda	0.497*** (0.148)	0.429*** (0.159)	0.578*** (0.170)	0.725*** (0.171)	0.503** (0.209)	0.621*** (0.224)	0.600* (0.306)
ICO Length	-0.004 (0.004)	-0.006 (0.005)	-0.003 (0.005)	-0.006 (0.007)	-0.006 (0.009)	-0.001 (0.008)	-0.003 (0.017)
Team Size	0.018 (0.019)	0.022 (0.018)	0.030 (0.020)	0.029 (0.020)	0.083*** (0.026)	0.055** (0.026)	0.092** (0.037)
Has GitHub	0.223 (0.351)	0.487 (0.380)	0.415 (0.448)	0.738 (0.498)	0.854 (0.545)	1.123* (0.570)	1.310* (0.772)
Has Twitter	-0.552 (1.274)	-0.117 (1.236)	-0.762 (1.899)	-1.198 (1.902)	0.583 (2.229)	-0.519 (1.674)	-1.256 (1.518)
BTC Price (ICO)	0.075 (0.069)	0.098 (0.065)	-0.032 (0.071)	0.014 (0.073)	-0.046 (0.088)	-0.045 (0.089)	0.048 (0.108)
Pre ICO	-0.235 (0.412)	-0.339 (0.482)	-0.247 (0.499)	-0.752 (0.574)	-1.096 (0.752)	-1.826** (0.815)	-0.485 (1.218)
Constant	10.524*** (1.805)	12.772*** (1.902)	9.618*** (2.573)	7.839*** (2.330)	3.503 (3.302)	5.193* (2.664)	7.778** (3.003)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.212	0.229	0.194	0.235	0.345	0.443	0.442
Observations	323	316	308	283	217	183	139
<b>Panel D: Composite Index</b>							
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
	Listing	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_comp	0.610*** (0.183)	0.569*** (0.192)	0.640*** (0.209)	0.881*** (0.215)	0.744** (0.292)	0.885*** (0.295)	0.807** (0.361)
ICO Length	-0.004 (0.004)	-0.005 (0.005)	-0.003 (0.005)	-0.005 (0.007)	-0.006 (0.009)	-0.001 (0.008)	-0.003 (0.017)
Team Size	0.013 (0.020)	0.017 (0.018)	0.026 (0.020)	0.022 (0.020)	0.076*** (0.026)	0.046* (0.027)	0.083** (0.036)
Has GitHub	0.183 (0.355)	0.442 (0.380)	0.370 (0.452)	0.679 (0.497)	0.781 (0.530)	1.045* (0.558)	1.249* (0.740)
Has Twitter	-0.508 (1.233)	-0.072 (1.173)	-0.712 (1.744)	-1.127 (1.770)	0.627 (2.176)	-0.463 (1.580)	-1.217 (1.469)
BTC Price (ICO)	0.082 (0.069)	0.104 (0.064)	-0.028 (0.071)	0.024 (0.072)	-0.036 (0.087)	-0.034 (0.087)	0.053 (0.110)
Pre ICO	-0.238 (0.416)	-0.328 (0.483)	-0.263 (0.507)	-0.755 (0.572)	-1.051 (0.735)	-1.769** (0.799)	-0.323 (1.187)
Constant	10.978*** (1.723)	13.078*** (1.785)	10.259*** (2.366)	8.551*** (2.196)	3.906 (3.230)	5.835** (2.522)	8.237*** (2.921)
Other Controls	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.215	0.234	0.190	0.237	0.353	0.453	0.448
Observations	323	316	308	283	217	183	139

**Table 9: Delisting Probability**

This table presents OLS estimates of the relationship between tech indexes and ICO delisting probabilities. The dependent variable is *Delist*, a dummy variable that equals to 1 if a token was shown as “inactive” on coinmarketcap.com by the end of 2018. For each tech index, the first column presents univariate result, and the second column displays estimates with control variables and fixed effects. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Supervised		Embedding		LDA		Composite	
Tech_sup	-0.028*** (0.008)	-0.043*** (0.015)						
Tech_embed			-0.030*** (0.010)	-0.025* (0.014)				
Tech_lda					-0.013 (0.008)	-0.016 (0.012)		
Tech_comp							-0.030*** (0.009)	-0.037** (0.015)
ICO Length		-0.000 (0.000)		-0.000 (0.000)		-0.000 (0.000)		-0.000 (0.000)
Team Size		0.002 (0.002)		0.001 (0.002)		0.001 (0.002)		0.002 (0.002)
Has GitHub		-0.052 (0.040)		-0.050 (0.039)		-0.055 (0.040)		-0.050 (0.040)
Has Twitter		-0.203 (0.205)		-0.187 (0.203)		-0.189 (0.207)		-0.191 (0.203)
BTC Price (ICO)		-0.011** (0.005)		-0.012** (0.005)		-0.011** (0.005)		-0.011** (0.005)
Pre ICO		0.018 (0.042)		0.008 (0.045)		0.014 (0.043)		0.009 (0.043)
Bonus		0.013 (0.064)		0.021 (0.065)		0.023 (0.065)		0.019 (0.065)
Accept BTC		-0.003 (0.034)		-0.002 (0.034)		-0.005 (0.034)		-0.004 (0.034)
Ethereum Based		-0.021 (0.053)		-0.027 (0.052)		-0.027 (0.054)		-0.030 (0.053)
Fog Index		-0.001** (0.001)		-0.001** (0.001)		-0.001* (0.001)		-0.001** (0.001)
Tone		-0.025 (0.021)		-0.018 (0.021)		-0.017 (0.021)		-0.022 (0.021)
Uncertainty		0.010 (0.051)		0.006 (0.051)		0.012 (0.052)		0.006 (0.052)
Constant	0.079*** (0.014)	0.693** (0.336)	0.083*** (0.015)	0.684** (0.332)	0.075*** (0.013)	0.692** (0.340)	0.081*** (0.014)	0.718** (0.334)
Fixed Effects	No	Yes	No	Yes	No	Yes	No	Yes
R <sup>2</sup>	0.015	0.166	0.018	0.152	0.004	0.148	0.015	0.157
Observations	422	329	422	329	422	329	422	329

Table 10: Comparison with Other Measures

This table compares our tech index with other technology measures: GitHub commit and simple word counts.  $\ln(\text{commits})$  is the logarithm of the number of code revisions on GitHub.  $\text{Simple\_word\_count}$  measures the percentage of words in a whitepaper that belongs to a self-defined technology word list. The complete word list can be found in Table OA.3. The dependent variable is *CMC Trading*. We include control variables related to ICO characteristics and whitepapers in all columns. Quarterly, categorical and geographical fixed effects are considered under all circumstances. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

	(1)	(2)	(3)	(4)	(5)	(6)
Tech_comp	0.066*** (0.015)			0.059*** (0.015)	0.066*** (0.017)	0.059*** (0.017)
Ln(commits)		0.017*** (0.005)		0.013*** (0.005)		0.013*** (0.005)
Simple_word_count			0.024** (0.011)		0.001 (0.012)	0.001 (0.012)
Control Variables	Yes	Yes	Yes	Yes	Yes	Yes
Fixed Effects	Yes	Yes	Yes	Yes	Yes	Yes
$R^2$	0.327	0.322	0.318	0.331	0.327	0.331
Observations	1382	1382	1382	1382	1382	1382



Table 11: Long Horizon Performance—Subsample on Readability

This table examines the long horizon performance of cryptocurrencies for different readability subsamples. The dependent variable is rate of return in panel A and Bitcoin-adjusted return in panel B. *Easy* is a dummy that equals to 1 if the whitepaper has a below median Fog index. We include control variables related to ICO characteristics and whitepapers in all columns. Quarterly, categorical and geographical fixed effects are considered under all circumstances. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

Panel A: Rate of Return						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.015 (0.049)	0.065 (0.093)	0.068 (0.127)	0.202 (0.181)	0.323 (0.201)	0.489* (0.271)
Tech_sup*Easy	-0.008 (0.060)	-0.069 (0.115)	-0.143 (0.175)	-0.335* (0.195)	-0.445* (0.236)	-0.645** (0.314)
Easy	0.038 (0.069)	0.017 (0.122)	0.020 (0.168)	0.026 (0.222)	-0.127 (0.255)	-0.069 (0.328)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
Fixed effects	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.076	0.147	0.225	0.357	0.412	0.386
Observations	316	310	286	218	184	140
Tech_embed	0.024 (0.052)	0.107 (0.092)	0.174 (0.131)	0.402** (0.186)	0.673*** (0.223)	0.913*** (0.202)
Tech_embed*Easy	0.008 (0.061)	-0.082 (0.105)	-0.133 (0.155)	-0.300 (0.206)	-0.479** (0.231)	-0.698*** (0.221)
Easy	0.029 (0.071)	0.015 (0.120)	0.000 (0.168)	-0.018 (0.226)	-0.132 (0.249)	-0.110 (0.328)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FES	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.078	0.151	0.230	0.371	0.456	0.448
Observations	316	310	286	218	184	140
Tech_lda	0.017 (0.052)	0.060 (0.086)	0.193 (0.122)	0.437*** (0.167)	0.713*** (0.196)	0.946*** (0.231)
Tech_lda*Easy	0.040 (0.062)	-0.009 (0.109)	-0.136 (0.150)	-0.415** (0.186)	-0.632*** (0.220)	-0.865*** (0.248)
Easy	0.022 (0.068)	-0.008 (0.119)	-0.021 (0.166)	-0.030 (0.225)	-0.195 (0.253)	-0.150 (0.339)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FES	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.082	0.148	0.231	0.370	0.450	0.437
Observations	316	310	286	218	184	140
Tech_comp	0.027 (0.062)	0.111 (0.108)	0.206 (0.156)	0.517** (0.227)	0.886*** (0.266)	1.217*** (0.273)
Tech_comp*Easy	0.015 (0.072)	-0.078 (0.126)	-0.192 (0.191)	-0.511** (0.238)	-0.814*** (0.281)	-1.160*** (0.286)
Easy	0.027 (0.070)	0.009 (0.121)	0.007 (0.168)	0.024 (0.225)	-0.091 (0.250)	-0.034 (0.325)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FES	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.079	0.149	0.229	0.370	0.452	0.445
Observations	316	310	286	218	184	140

<b>Panel B: Adjusted Rate of Returns</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.005 (0.048)	0.026 (0.086)	0.025 (0.121)	0.256* (0.154)	0.383** (0.179)	0.508** (0.249)
Tech_sup*Easy	0.008 (0.058)	0.001 (0.107)	-0.046 (0.163)	-0.279 (0.173)	-0.427* (0.217)	-0.589** (0.286)
Easy	0.005 (0.066)	-0.024 (0.115)	0.000 (0.155)	0.050 (0.198)	-0.056 (0.236)	0.028 (0.302)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FEs	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.096	0.139	0.155	0.319	0.302	0.284
Observations	311	305	281	213	180	137
Tech_embed	0.027 (0.050)	0.101 (0.088)	0.159 (0.121)	0.342** (0.158)	0.605*** (0.190)	0.870*** (0.188)
Tech_embed*Easy	0.000 (0.059)	-0.038 (0.102)	-0.101 (0.145)	-0.173 (0.174)	-0.399* (0.202)	-0.600*** (0.200)
Easy	0.002 (0.068)	-0.029 (0.113)	-0.006 (0.154)	-0.016 (0.200)	-0.072 (0.234)	-0.035 (0.304)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FEs	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.098	0.146	0.164	0.338	0.349	0.366
Observations	311	305	281	213	180	137
Tech_lda	0.007 (0.052)	0.032 (0.085)	0.153 (0.115)	0.337** (0.142)	0.576*** (0.178)	0.913*** (0.210)
Tech_lda*Easy	0.049 (0.062)	0.037 (0.107)	-0.075 (0.140)	-0.276* (0.161)	-0.465** (0.200)	-0.754*** (0.220)
Easy	-0.008 (0.065)	-0.041 (0.112)	-0.027 (0.152)	0.001 (0.202)	-0.126 (0.235)	-0.071 (0.311)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FEs	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.102	0.142	0.164	0.328	0.329	0.348
Observations	311	305	281	213	180	137
Tech_comp	adjror7 0.020 (0.061)	adjror30 0.079 (0.105)	adjror90 0.163 (0.149)	adjror180 0.463** (0.190)	adjror240 0.819*** (0.230)	adjror300 1.193*** (0.253)
Tech_comp*Easy	0.023 (0.070)	-0.009 (0.122)	-0.110 (0.179)	-0.364* (0.202)	-0.692*** (0.249)	-1.033*** (0.257)
Easy	-0.003 (0.067)	-0.034 (0.114)	-0.009 (0.155)	0.029 (0.201)	-0.032 (0.233)	0.034 (0.299)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
FEs	Yes	Yes	Yes	Yes	Yes	Yes
R2	0.098	0.143	0.161	0.335	0.347	0.362
Observations	311	305	281	213	180	137

**Table 12: Is There Return Reversal?**

This table presents the effects of tech indexes on long term returns of cryptocurrencies. The dependent variable is  $\log(1 + ROR_{180 \rightarrow j})$ , the gross return from 180 listing days onward. Panel A displays results on rate of returns and panel B shows Bitcoin-adjusted rate of returns. Column (1)-(6) display results for six horizons from the listing day: 210 days, 240 days, 270 days, 300 days, 330 days and 360 days. We include control variables related to ICO characteristics and whitepapers in all columns. Quarterly, categorical and geographical fixed effects are considered under all circumstances. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Rate of Returns</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	210 Days	240 Days	270 Days	300 Days	330 Days	360 Days
Tech_sup	0.052 (0.036)	0.047 (0.041)	0.081* (0.043)	0.034 (0.045)	0.061 (0.041)	0.035 (0.082)
$R^2$	0.297	0.306	0.404	0.553	0.481	0.495
Tech_embed	0.048 (0.035)	0.086** (0.038)	0.080* (0.041)	0.089* (0.048)	0.070 (0.059)	0.008 (0.077)
$R^2$	0.296	0.319	0.403	0.563	0.482	0.493
Tech_lda	0.051 (0.034)	0.065* (0.039)	0.012 (0.041)	0.032 (0.050)	-0.049 (0.063)	-0.115 (0.071)
$R^2$	0.298	0.312	0.390	0.553	0.479	0.512
Tech_comp	0.071* (0.039)	0.093** (0.046)	0.082* (0.046)	0.073 (0.056)	0.040 (0.063)	-0.037 (0.085)
$R^2$	0.301	0.316	0.400	0.557	0.477	0.495
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
Observations	207	184	156	140	132	103
<b>Panel B: Bitcoin-Adjusted Rate of Returns</b>						
	210 Days	240 Days	270 Days	300 Days	330 Days	360 Days
Tech_sup	0.037 (0.030)	0.021 (0.038)	0.039 (0.041)	-0.001 (0.039)	0.031 (0.041)	0.024 (0.062)
$R^2$	0.287	0.343	0.419	0.529	0.448	0.564
Tech_embed	0.045 (0.030)	0.037 (0.035)	0.076* (0.042)	0.069 (0.046)	0.042 (0.059)	0.028 (0.071)
$R^2$	0.291	0.346	0.429	0.538	0.449	0.565
Tech_lda	0.047 (0.030)	0.017 (0.033)	0.014 (0.043)	0.025 (0.045)	-0.064 (0.059)	-0.082 (0.064)
$R^2$	0.292	0.342	0.416	0.530	0.453	0.574
Tech_comp	0.060* (0.034)	0.035 (0.041)	0.061 (0.048)	0.044 (0.051)	0.005 (0.062)	-0.016 (0.076)
$R^2$	0.293	0.344	0.421	0.531	0.447	0.564
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
Observations	202	180	153	137	129	101

Table 13: ICO First-Day Price

This table presents OLS estimates of the relationship between tech indexes and ICO first-day price. The dependent variable is the log transformation of the ratio between the first day's opening price and ICO price. For each tech index, the first column presents univariate result, and the second column displays estimates with control variables and fixed effects. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at 1%, 5%, and 10% respectively.

Ln(First Opening Price/ICO Price)								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Supervised		Embedding		LDA		Composite	
Tech_sup	0.337*** (0.090)	0.299*** (0.109)						
Tech_embed			0.414*** (0.076)	0.368*** (0.098)				
Tech_lda					0.310*** (0.074)	0.287*** (0.091)		
Tech_comp							0.458*** (0.092)	0.417*** (0.117)
ICO Length		-0.002 (0.002)		-0.001 (0.002)		-0.001 (0.002)		-0.001 (0.002)
Team Size		-0.001 (0.014)		0.007 (0.012)		0.008 (0.011)		0.003 (0.012)
Has GitHub		-0.059 (0.263)		-0.128 (0.253)		-0.046 (0.254)		-0.100 (0.254)
Has Twitter		-0.074 (0.368)		-0.162 (0.392)		-0.284 (0.431)		-0.151 (0.376)
BTC Price (ICO)		-0.012 (0.048)		0.010 (0.048)		-0.008 (0.049)		-0.002 (0.047)
Pre ICO		-0.363 (0.341)		-0.205 (0.348)		-0.278 (0.352)		-0.266 (0.346)
Bonus		0.152 (0.402)		0.128 (0.407)		0.137 (0.414)		0.158 (0.404)
Accept BTC		-0.014 (0.236)		-0.049 (0.229)		0.014 (0.236)		-0.015 (0.232)
Ethereum Based		-0.299 (0.385)		-0.206 (0.367)		-0.136 (0.384)		-0.183 (0.371)
Fog Index		-0.007 (0.007)		-0.005 (0.007)		-0.007 (0.007)		-0.005 (0.007)
Tone		0.103 (0.146)		0.105 (0.147)		0.102 (0.148)		0.133 (0.145)
Uncertainty		-0.036 (0.318)		0.121 (0.328)		0.021 (0.327)		0.070 (0.325)
Constant	-0.295*** (0.105)	-2.998*** (0.923)	-0.404*** (0.111)	-3.392*** (0.916)	-0.324*** (0.105)	-3.578*** (0.974)	-0.380*** (0.110)	-3.552*** (0.924)
Fixed Effects	No	Yes	No	Yes	No	Yes	No	Yes
R <sup>2</sup>	0.064	0.305	0.111	0.328	0.066	0.306	0.102	0.325
Observations	238	199	238	199	238	199	238	199

Table 14: **Robustness Tests**

This table displays several robustness tests. Panel A redoes Table 4 using Trading as the dependent variable. Panel B is the Logit regression version of Table 4. Besides, to mitigate the concern of survivorship bias, we impute -99% to returns of delisted cryptocurrencies and redo table 6 and table 7. Results are presented in Panel C and panel D. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Trading</b>								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Supervised		Embedding		LDA		Composite	
Tech_sup	0.058*** (0.011)	0.034*** (0.010)						
Tech_embed			0.101*** (0.010)	0.047*** (0.010)				
Tech_lda					0.082*** (0.011)	0.051*** (0.011)		
Tech_comp							0.114*** (0.013)	0.065*** (0.013)
ICO Length		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)		-0.001*** (0.000)
Team Size		0.004*** (0.001)		0.005*** (0.001)		0.005*** (0.001)		0.004*** (0.001)
Has GitHub		0.004 (0.016)		-0.005 (0.016)		-0.002 (0.016)		-0.005 (0.016)
Has Twitter		0.102*** (0.028)		0.110*** (0.028)		0.106*** (0.029)		0.104*** (0.028)
BTC Price (ICO)		0.007 (0.006)		0.007 (0.006)		0.007 (0.006)		0.007 (0.006)
Pre ICO		-0.029* (0.017)		-0.025 (0.017)		-0.025 (0.017)		-0.027 (0.017)
Bonus		0.021 (0.013)		0.022* (0.013)		0.020 (0.013)		0.022* (0.012)
Accept BTC		0.013 (0.015)		0.016 (0.015)		0.017 (0.015)		0.017 (0.015)
Ethereum Based		0.007 (0.023)		0.015 (0.023)		0.013 (0.022)		0.015 (0.022)
Fog Index		-0.000 (0.001)		-0.000 (0.001)		-0.000 (0.001)		-0.000 (0.001)
Tone		-0.005 (0.011)		-0.001 (0.011)		0.001 (0.011)		0.002 (0.011)
Uncertainty		0.011 (0.021)		0.027 (0.021)		0.020 (0.021)		0.022 (0.021)
Constant	0.167*** (0.009)	0.767*** (0.065)	0.167*** (0.009)	0.678*** (0.068)	0.167*** (0.009)	0.649*** (0.069)	0.167*** (0.009)	0.668*** (0.067)
Fixed Effects	No	Yes	No	Yes	No	Yes	No	Yes
R <sup>2</sup>	0.024	0.383	0.074	0.389	0.049	0.390	0.066	0.393
Observations	1629	1382	1629	1382	1629	1382	1629	1382

<b>Panel B: Logit regression</b>								
	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)
	Supervised		Embedding		LDA		Composite	
Tech_sup	0.343*** (0.054)	0.287*** (0.094)						
Tech_embed			0.527*** (0.055)	0.345*** (0.085)				
Tech_lda					0.398*** (0.052)	0.293*** (0.089)		
Tech_comp							0.596*** (0.065)	0.443*** (0.108)
ICO Length		-0.012*** (0.003)		-0.012*** (0.003)		-0.012*** (0.003)		-0.012*** (0.003)
Team Size		0.064*** (0.012)		0.069*** (0.012)		0.069*** (0.012)		0.066*** (0.012)
Has GitHub		0.426** (0.175)		0.374** (0.175)		0.402** (0.175)		0.381** (0.175)
Has Twitter		1.463*** (0.417)		1.491*** (0.432)		1.461*** (0.418)		1.459*** (0.425)
BTC Price (ICO)		0.051 (0.034)		0.055 (0.033)		0.051 (0.034)		0.055 (0.034)
Pre ICO		-0.263 (0.186)		-0.211 (0.186)		-0.236 (0.187)		-0.229 (0.187)
Bonus		0.122 (0.245)		0.129 (0.245)		0.108 (0.246)		0.126 (0.247)
Accept BTC		-0.186 (0.174)		-0.179 (0.174)		-0.167 (0.174)		-0.164 (0.175)
Ethereum Based		-0.148 (0.255)		-0.094 (0.257)		-0.097 (0.259)		-0.082 (0.261)
Fog Index		-0.000 (0.008)		0.000 (0.008)		-0.000 (0.008)		0.001 (0.008)
Tone		-0.007 (0.118)		0.023 (0.118)		0.014 (0.120)		0.036 (0.120)
Uncertainty		0.059 (0.239)		0.220 (0.237)		0.137 (0.233)		0.152 (0.239)
Constant	-1.077*** (0.058)	-6.461*** (0.986)	-1.108*** (0.059)	-6.693*** (1.010)	-1.083*** (0.058)	-6.538*** (0.987)	-1.101*** (0.059)	-6.593*** (0.999)
Fixed Effects	No	Yes	No	Yes	No	Yes	No	Yes
Pseudo R <sup>2</sup>	0.0214	0.322	0.050	0.325	0.031	0.322	0.046	0.326
Observations	1629	1351	1629	1351	1629	1351	1629	1351

<b>Panel C: Rate of return (-99% return for delisted coins)</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.044	0.081	0.037	0.006	0.094	0.071
	(0.049)	(0.070)	(0.093)	(0.104)	(0.117)	(0.132)
Constant	-1.905	-2.869*	-2.378**	-0.430	0.029	0.028
	(1.246)	(1.490)	(1.027)	(1.162)	(1.258)	(1.440)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.152	0.204	0.256	0.373	0.438	0.428
Observations	323	319	293	228	198	157
<hr/>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_embed	0.112**	0.156**	0.189**	0.241**	0.379***	0.396***
	(0.052)	(0.070)	(0.088)	(0.101)	(0.113)	(0.141)
Constant	-2.045	-3.035**	-2.658***	-0.747	-0.391	-0.495
	(1.246)	(1.472)	(1.015)	(1.162)	(1.252)	(1.373)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.164	0.215	0.271	0.389	0.473	0.462
Observations	323	319	293	228	198	157
<hr/>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_lda	0.120**	0.138**	0.182**	0.185**	0.300***	0.253*
	(0.046)	(0.062)	(0.078)	(0.088)	(0.096)	(0.129)
Constant	-2.240*	-3.208**	-2.910***	-0.954	-0.710	-0.528
	(1.273)	(1.517)	(1.070)	(1.194)	(1.276)	(1.455)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.167	0.212	0.270	0.383	0.459	0.441
Observations	323	319	293	228	198	157
<hr/>						
	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_comp	0.127**	0.171**	0.187*	0.205*	0.364***	0.342**
	(0.055)	(0.076)	(0.099)	(0.112)	(0.127)	(0.155)
Constant	-2.112*	-3.113**	-2.712***	-0.764	-0.473	-0.477
	(1.254)	(1.484)	(1.034)	(1.156)	(1.225)	(1.395)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.164	0.213	0.266	0.381	0.460	0.445
Observations	323	319	293	228	198	157

**Panel D: Bitcoin-adjusted rate of return (-99% return for delisted coins)**

	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_sup	0.011 (0.029)	0.050 (0.058)	0.019 (0.085)	0.077 (0.095)	0.152 (0.113)	0.096 (0.130)
Constant	-0.315 (0.658)	-1.844 (1.578)	-1.716* (1.014)	-2.876** (1.265)	-3.648*** (1.348)	-2.607* (1.466)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.091	0.156	0.163	0.311	0.302	0.287
Observations	312	308	282	217	186	144

	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_embed	0.052 (0.039)	0.116* (0.060)	0.137* (0.081)	0.236** (0.096)	0.354*** (0.111)	0.428*** (0.146)
Constant	-0.417 (0.676)	-2.038 (1.586)	-1.991* (1.021)	-3.252** (1.285)	-4.174*** (1.371)	-3.415** (1.409)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.097	0.166	0.175	0.333	0.341	0.342
Observations	312	308	282	217	186	144

	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_lda	0.058 (0.037)	0.085 (0.058)	0.129* (0.073)	0.156* (0.081)	0.255*** (0.094)	0.256* (0.132)
Constant	-0.508 (0.689)	-2.090 (1.614)	-2.150** (1.069)	-3.368** (1.329)	-4.465*** (1.413)	-3.414** (1.523)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.099	0.160	0.174	0.320	0.319	0.305
Observations	312	308	282	217	186	144

	(1)	(2)	(3)	(4)	(5)	(6)
	7 Days	30 Days	90 Days	180 Days	240 Days	300 Days
Tech_comp	0.056 (0.041)	0.114* (0.067)	0.131 (0.092)	0.219** (0.103)	0.358*** (0.125)	0.369** (0.156)
Constant	-0.428 (0.674)	-2.041 (1.581)	-1.982* (1.010)	-3.245** (1.261)	-4.234*** (1.341)	-3.283** (1.409)
Controls & FEs	Yes	Yes	Yes	Yes	Yes	Yes
R <sup>2</sup>	0.096	0.162	0.171	0.324	0.328	0.315
Observations	312	308	282	217	186	144



# Online Appendix: Technical Notes on Measure Construction

## 1 Supervised Machine Learning

### 1.1 Basics

Supervised learning is “the machine learning task of learning a function that maps an input to an output based on example input-output pairs.” (Russell and Norvig, 2010). Mathematically, this can be expressed as estimating a function  $f(\cdot)$  given input variables ( $X$ ) and output variables ( $Y$ ), such that the mapping function  $Y=f(X)$  is satisfied as much as possible. Various machine learning models impose different constraints on the function, resulting in different optimization results.

We use the simplest regression model, ordinary least squares (OLS), as our benchmark. The objective function of OLS is:

$$\min_{\beta} ||y - x\beta||^2$$

OLS works well when there are only a few predictors, but its performance deteriorates significantly as the dimension of predictors increases. Unfortunately, high dimensionality and sparseness are both common features of text data. Hence, we apply more advanced machine learning methods to avoid the "curse of dimensionality".

The first set of methods we use are penalized linear approaches. The idea is to add a penalty term in the objective function to reduce a model's fit on noise and hence enhances prediction accuracy. We consider LASSO, ridge regression and elastic net for this approach. Another common approach to deal with high-dimensional data is dimension reduction. While penalized linear methods select a subset of predictors that have strong predictive power, dimension reduction methods combine predictors into several main components while retaining as much information as possible. We apply principal component regression (PCR) and partial least squares (PLS) in this vein. All the methods above are linear regression models, but we are also interested in using non-linear ap-

proaches to get better prediction accuracy. We consider decision trees (random forest and gradient boosting) and neural network algorithms. Next, I briefly introduce each of these machine learning methods that we consider as candidates to construct our supervised tech index.

### 1.1.1 LASSO

LASSO (least absolute shrinkage and selection operator) is a common approach employed to deal with high-dimensional sparse data. The objective function for LASSO is:

$$\min_{\beta} \left\{ \frac{1}{N} \|y - X\beta\|^2 + \lambda \|\beta\| \right\}.$$

The first term is the same as OLS, while the second term is a penalty on non-zero coefficients, with  $\lambda$  representing the regularization strength. The effect of LASSO is to select only a subset of predictors by pushing other predictor coefficients to 0.

### 1.1.2 Ridge

Ridge regression (also known as Tikhonov regularization) is another useful method to mitigate the problem of dimensionality by adding a L2-norm regularization term as penalty. The objective function is:

$$\min_{\beta} \left\{ \frac{1}{N} \|y - X\beta\|^2 + \lambda \|\beta\|^2 \right\}.$$

The difference of Ridge regression from LASSO is that it shrinks the coefficients of unimportant predictors but do not set them to 0. Hence, ridge regression is a regularization approach, but not a variable selection approach.

### 1.1.3 Elastic Net

Elastic net is a combination of LASSO and ridge regression. It optimizes the following objective function:

$$\min_{\beta} \left\{ \frac{1}{N} \|y - X\beta\|^2 + \lambda\alpha \|\beta\| + \lambda(1 - \alpha) \|\beta\|^2 \right\}.$$

$\alpha$  controls the weight between L1 and L2 norm penalty. If  $\alpha = 1$ , it is the same as LASSO; if  $\alpha = 0$ , it becomes ridge regression. By averaging between LASSO and ridge regression, elastic net is expected to combine the advantages of both methods.

#### 1.1.4 Principal component regression (PCR)

Principal component regression combines standard linear regression with principal component analysis (PCA). Specifically, PCR regresses the dependent variable (Y) on principal components of independent variables (X), as opposed to regressing Y directly on X in OLS. Since the principal components are extracted based on their ability to explain the variation in X, the forecasting goal (Y) does not come into play until the final regression step.

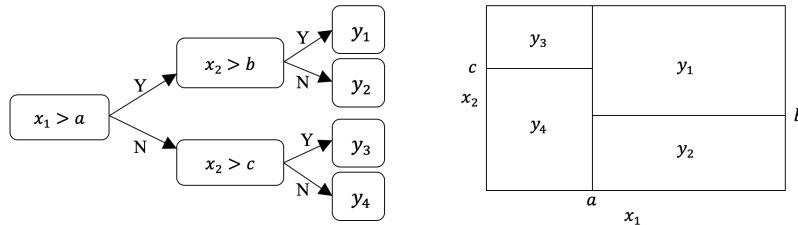
#### 1.1.5 Partial least square (PLS)

Partial least squares (PLS) regression shares some similarities with PCR, but it constructs the principal components of X with the goal to best explain the covariance between X and Y. It first projects both the independent variables (X) and dependent variables (Y) to a new space, in which the projection of the X-space that explains the most variation of the Y-space. It then runs a linear regression model in the new space. PLS is especially helpful when predictors are more than available observations or when predictors are highly collinear.

#### 1.1.6 Random Forest

Random forests come from decision trees. A decision tree is a set of logic conditions on input variables (X) that lead to predictions on the target output (Y). The following figure illustrates a regression tree example. The first condition used to determine y is whether  $x_1$  is greater than a. Conditional on the answer to this question, another logic condition will be raised. This process iterates until the value of y is determined. Different from linear regressions, the regression tree is a

non-linear and non-parametric method. A random forest is an ensemble of multiple decision trees. It outputs the average prediction of each individual tree. Although a single tree may be a weak prediction model, through combination the random forest can have a strong performance.



### 1.1.7 Gradient Boosting

Gradient boosting is another approach to ensemble regression trees. At each step, a new tree is fitted on the negative gradient of a given loss function. Hence, new trees aim at correcting the error of preceding trees. To avoid overfitting on residuals, following trees will be discounted at each step. This process is repeated until a total number of N trees is reached.

### 1.1.8 Neural Network (NN)

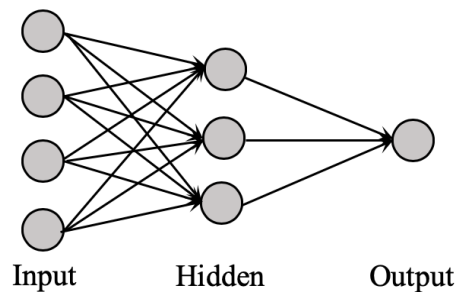
Artificial neural network is a broad set of machine learning algorithms inspired by the biological neural structure of human brains. It is a layer-by-layer structure, where each layer is composed of “neurons”, and the layers are connected by "edges". The following figure shows an example, the feedforward neural network. The input layer is the input variables (X), and the output layer is the outcome (Y). Each node of the hidden layer represents the following operation:

$$H_i = f(w_0 + XW),$$

where  $W$ , the linear weight matrix on the inputs, represents the “edges” connecting the input layer and the hidden layer. There are multiple choices for  $f(\cdot)$ , one of which is the sigmoid function:

$$f(x) = \frac{1}{1 + e^{-x}}$$

The output of the hidden layer ( $H_i$ ) can then be used as the input of the output layer or another concatenated hidden layer. This process continues until the output layer is arrived.



## 1.2 Hyperparameter search

We tune one hyperparameter for each of the supervised machine learning methods. For LASSO and ridge regression, we change the regularization strength ( $\lambda$ ); for elastic net, we alter the linear weight ( $\alpha$ ); for PCR and PLS, we vary the number of principal components; for random forest and gradient boosting, we adjust the number of trees; for neural network, we tune the number of nodes of the hidden layer. Figure OA.1 presents the hyperparameter search results. Table 3 shows the best out-of-sample R-square ( $R_{OOS}^2$ ) for each supervised method and their corresponding parameters. It may be surprising that the most popular and advanced neural network approach works the worst among all methods and even underperforms the most basic OLS. This is due to the mismatch between the high-dimensional predictors and the relatively small sample size. NN is a highly parameterized model, and we do not have enough observation to get all parameters well-tuned. This mismatch can also explain why dimension reduction methods (especially PLS) works particularly

well on our dataset. By limiting the predicting variables to only a few principal components, the number of parameters is manageable for our training set.

## 2 Word Embedding & K-Means Clustering

### 2.1 Model

#### 2.1.1 Word Embedding model

In practice, word embedding vectors are estimated with a two-layer neural network, the Skip-gram model. Given a sequence of words  $w_1, w_2, \dots, w_T$ , the inference problem is to maximize the average log probability of the context of  $w_t$ :

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

where  $c$  denotes the size of the context.  $p(w_{t+j} | w_t)$  is calculated as:

$$p(w_o | w_I) = \frac{\exp(v_{w_o}'^T v_{w_I})}{\sum_{w=1}^W \exp(v_{w_o}'^T v_{w_I})}$$

where  $v(w_I)$  and  $v(w_o)$  represent the input and output representation of word  $w_t$ , and  $W$  denotes vocabulary size. The embedding of  $w_t$  is the projection vector between the input and output layer.

#### 2.1.2 K-means clustering algorithm

Given a fixed number of clusters ( $k$ ), the objective function of k-means is to find a partition of the dataset, such that the within-cluster sum of squared distances between each observation and its closest centroid are minimized. Equivalently, this can be expressed as:

$$\arg \min_S \sum_{i=1}^K \sum_{x \in S_i} \|x - \mu_i\|^2$$

where  $\mu_i$  is the average of data points in  $S_i$ .

A k-means algorithm works as follows:

1. Specify the number of clusters  $k$ . Randomly select  $k$  data points as cluster centroids ( $\mu$ ).
2. For each data point, assign it to the nearest centroid:

$$label(i) = \arg \min_j \|x_i - \mu_j\|^2$$

3. Update each centroid as the average of data points in that cluster:

$$\mu_j = \frac{1}{||S_j||} \sum_{x \in S_j} x$$

4. Repeat 2) and 3), until the assignments of data points no longer change.

## 2.2 Preprocessing

Before estimation, we preprocess the raw text step by step to get a cleaner input. We first split all documents into words and convert them to lowercases. We then apply lemmatization to convert all words to its root form. Because word embedding uses contextual information, we do not remove individual words before estimating the vector representation, so as not to affect the sentence structure. After obtaining the word embedding vector, we drop stop-words and low-frequency words that appear less than 20 times in the vocabulary. Finally, we transform preprocessed text into numerical counts that we use as the input of word embedding estimation. The corpus is represented as a  $D \times V$  document-term matrix  $M$ , where  $M(d, v)$  indicates the count of the  $v$ -th word in the  $d$ -th document. This is the “bag-of-words” representation. The underlying assumption is that the order of words does not matter. Although this is an oversimplification of reality, it retains a large amount of information while keeping the algorithm simple. The final corpus consists of 2,262 documents and 20,145 unique terms.

## 2.3 Choice of topics

One important step of k-means is to find the optimal number of topics. We take a data-driven approach to select the best model. To be specific, we apply the “Elbow method” to the distortion score (the sum of squared distances between each point and its assigned centroid), which is a heuristic method to find the appropriate number of clusters on a dataset. “Elbow” refers to the point where adding another cluster does not give much improvement to the model.<sup>10</sup> To determine the “Elbow”, [Satopaa et al. \(2011\)](#) propose an algorithm detecting the point of maximum curvature as the elbow, where the curvature can be calculated as:

$$K_f(x) = \frac{f''(x)}{\left(1 + f'(x)^2\right)^{\frac{3}{2}}}$$

Figure [OA.2](#) presents the results on the optimal number of topics. We find that the optimal number of topics detected by the algorithm is 20.

## 3 Latent Dirichlet Allocation (LDA)

### 3.1 LDA Model

Latent Dirichlet Allocation (LDA), developed by [Blei et al. \(2003\)](#), is a generative probabilistic modeling approach. The basic idea is that each document can be represented as a probability distribution over various topics, where each topic is a probability distribution over the vocabulary of a corpus. Suppose there are  $K$  latent topics,  $D$  documents and  $V$  unique terms in the corpus. LDA assumes the following data generating process for each document  $d$ :

1. Draw  $\beta_k$  from a multinomial distribution, where  $\beta_k$  (a  $1 \times V$  vector) denotes the word distribution of topic  $k$  for each  $k = 1, 2, \dots, K$ .

---

<sup>10</sup>[Hansen et al. \(2018\)](#) use the method to select the number of topics for the FOMC transcripts.



2. Draw  $\theta_d$  from a Dirichlet distribution, where  $\theta_d$  (a  $1 \times K$  vector) denotes the topical distribution of document  $d$ .
3. For each word  $w$  in document  $d$ :
  - (a) Choose a topic  $k$  from  $\theta_d$ ;
  - (b) Choose a word  $w$  from  $\beta_k$ .

Intuitively, one can think of generating a document with  $N$  words as repeating the action of "generating a word" by  $N$  times, where each word is generated in two steps: first, roll a  $K$ -sided dice to select a topic; conditional on the topic being selected, roll another  $V$ -sided dice to choose a word. Note that the probability of obtaining each side is not equal. It corresponds to  $\theta_d$  and  $\beta_k$  respectively.

Given a corpus and a latent topic number  $K$ , the inference problem of LDA is to compute the posterior distribution of hidden variables  $\Theta = (\theta_1, \theta_2, \dots, \theta_D)$  and  $B = (\beta_1, \beta_2, \dots, \beta_K)$ , such that the generated distribution resembles the observed distribution of words of each document. Since the distribution is usually mathematically intractable, it is solved with Gibbs sampling algorithm ([Griffiths and Steyvers, 2004](#)) in practice.

## 3.2 Preprocessing

Similar to word embedding, we preprocess the raw text to get a cleaner input of the LDA model. First, we split all documents into words and convert to lowercases. We then remove common stop-words like "the", "a" and "I", as they appear frequently in text but convey little information. Second, we convert all words to its root form, so that words like "communicates", "communicating" all become "communicate". Third, we identify common two-word collocations which appears more than 20 times in the corpus. For example, "machine learning" conveys a specific meaning different from "machine" and "learning". Fourth, we drop infrequent unigrams and bigrams that appear in less than 10 documents. Finally, we convert the preprocessed text to a

document-term matrix, as what we do for word embedding analysis. The final corpus consists of 2,262 documents and 26,410 unique terms.

### 3.3 Choice of topics

An important yet challenging task of LDA is to find the optimal number of topics ( $K$ ). As discussed in Hansen et al. (2018), there is a trade-off between model interpretability and statistical goodness-of-fit. If  $K$  is too small, the model does not fit the data well, and the topics generated are often too general and mix multiple themes. However, if  $K$  is too large, the topics are too fine-grained, which impairs the interpretability of the model. To balance the two effects, we adopt a statistical measure—topic coherence—to select  $K$  (Röder et al., 2015). A topic is said to be coherent if its top words frequently co-occur with each other. In particular, we use normalized pointwise mutual information (NPMI) that has been proved to have the largest correlation to human topic coherence ratings to calculate co-occurrence:

$$NPMI(w_i, w_j) = \frac{\log\left(\frac{P(w_i, w_j) + \epsilon}{P(w_i)P(w_j)}\right)}{-\log(P(w_i, w_j) + \epsilon)},$$

where  $P(w_i)$ ,  $P(w_j)$  and  $P(w_i, w_j)$  denote the probability that  $w_i$  appears,  $w_j$  appears and  $w_i$  and  $w_j$  jointly appear in the corpus.  $\epsilon$  is added to avoid taking logarithm on zero.

We consider candidates of topic numbers ( $K$ ) ranging from 10 to 80 in increments of ten. Figure OA.3 shows the topic coherence of each LDA model with different specifications of  $K$ . It indicates that  $K = 20$  maximizes the coherence measure and produces the best results. To understand the LDA output with 20 topics, we need to interpret the estimated topics. Since each topic is a probability distribution over all unique terms in the vocabulary, a natural way to name each topic is to read the terms with the highest probabilities and manually assign a label. However, the most frequent terms often appear in multiple topics, making it difficult to distinguish between topics. An alternative way is to look for terms that exclusively appear in a given topic. This is defined as the ratio of a term's probability within a topic to its probability across all topics (Taddy, 2012).

Bybee et al. (2020) adopts this approach to analyze the structure of economics news from the Wall Street Journal. However, this measure may put too much weight on very rare terms, which can also be hard to interpret. Following Sievert and Shirley (2014), we use the relevance measure, which is defined as the weighted average of the two measures above:

$$Relevance(term\ w|topic\ t) = \lambda \times p(w|t) + (1 - \lambda) \times \frac{p(w|t)}{p(w)}$$

We find LDA topics with  $\lambda = 0.6$  yields the best topic interpretability.

Figure OA.1: **Supervised Learning Hyperparameter Search**

This figure plots the hyperparameter search results of the supervised method. For each subplot, the solid blue line indicates how  $R^2_{OOS}$  varies with different parameter choices, and the dashed red line indicates the parameter that gives the best performance.

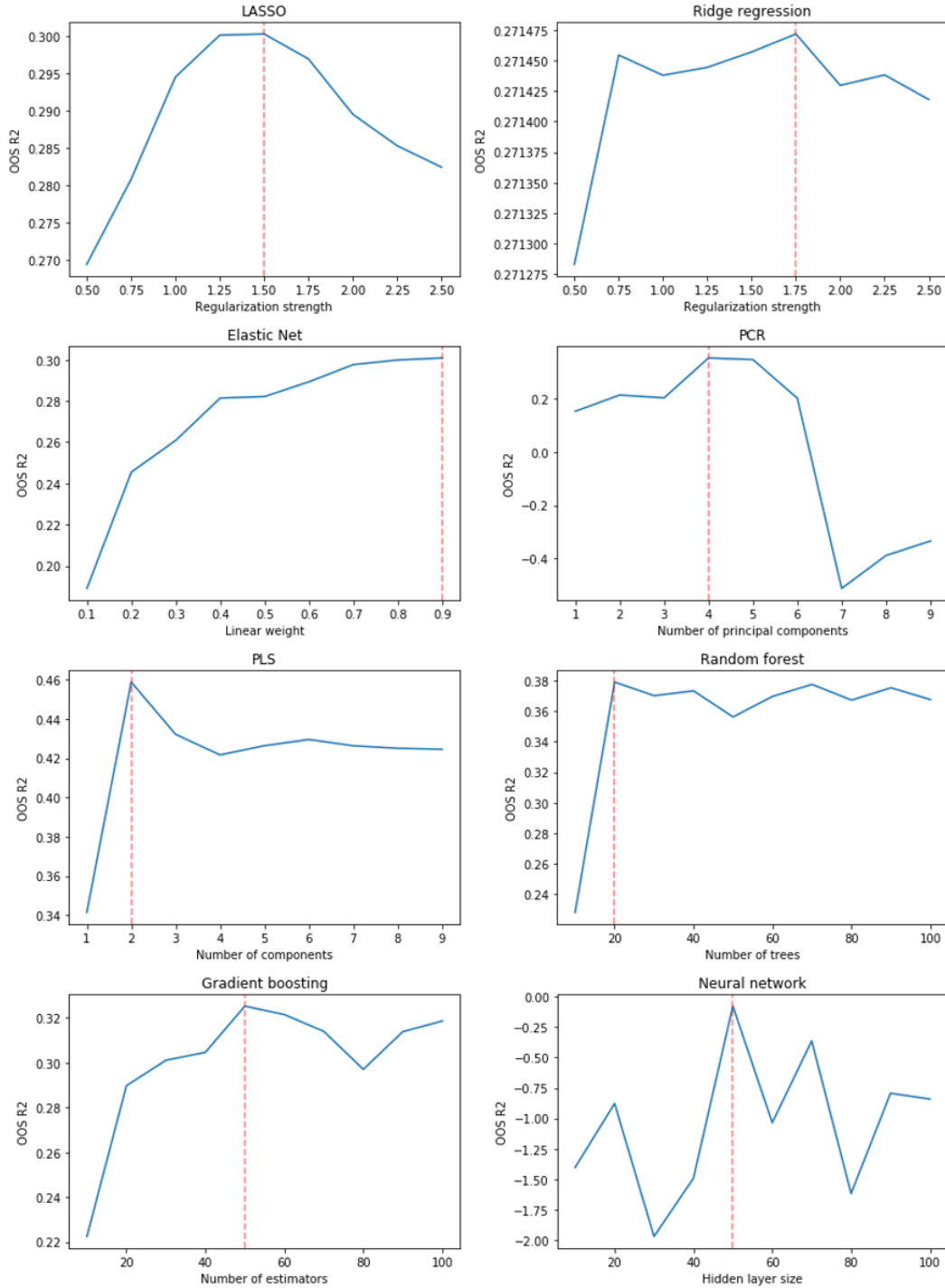


Figure OA.2: **Elbow Method**

This figure shows the elbow method used to select the most appropriate number of clusters. The blue solid line plots the elbow curve of the distortion score, and the red dashed line indicates the “elbow” detected by the algorithm.

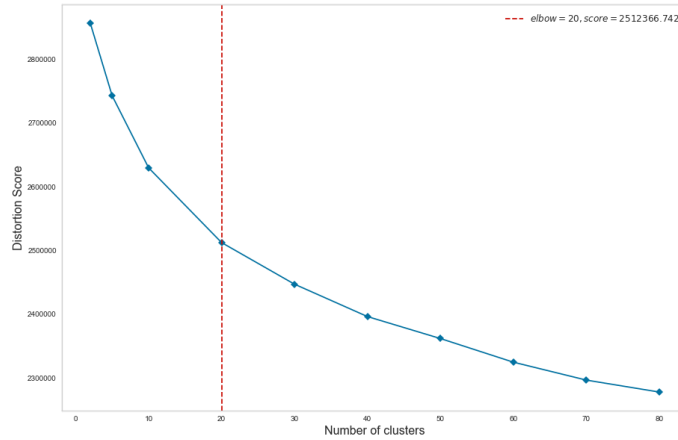


Figure OA.3: **Topic Coherence**

This figure plots the topic coherence measure with different specifications of LDA topic numbers.

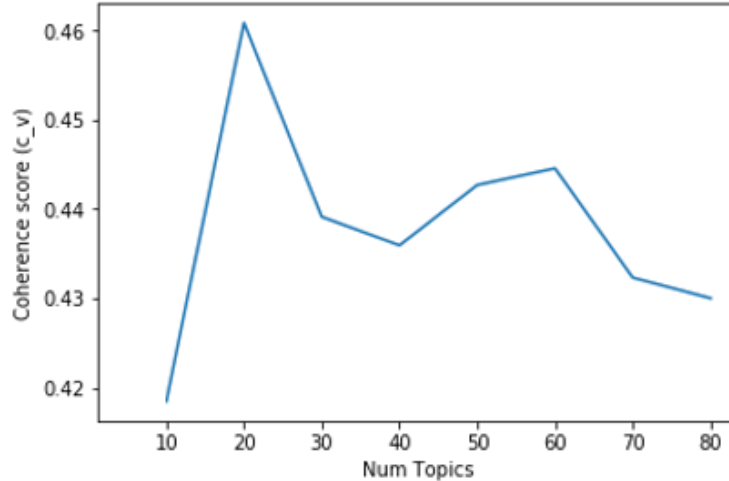
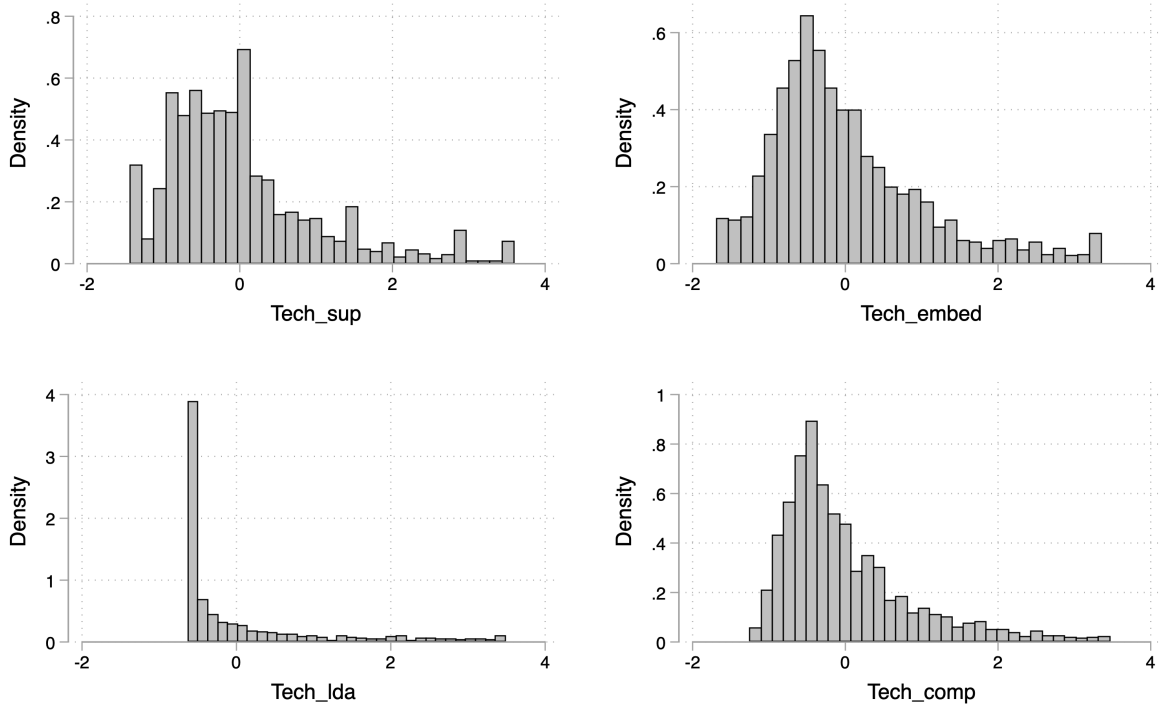


Figure OA.4: **Distribution of Technology Indexes**

This figure plots the distributions of the four technology indexes.



### Figure OA.5: Technology Index Validation

This figure plots the relationship between the composite technology index and GitHub measures. In panel (a), the variable of interest is *subscriber*, which measures the number of users subscribing repository updates; in panel (b), *star* indicates the number of “likes” received by the repository; in panel (c), *fork* proxies for repository copies made by other developers; in panel (d), *commit* represents how many times the code has been revised; in panel (e), *branch* is the amount of pointers to specific versions of the repository; and in panel (f) *contributor* reflects how many developers have contributed to the source. The red solid line represents the linear fitting of GitHub measures on the composite tech index.

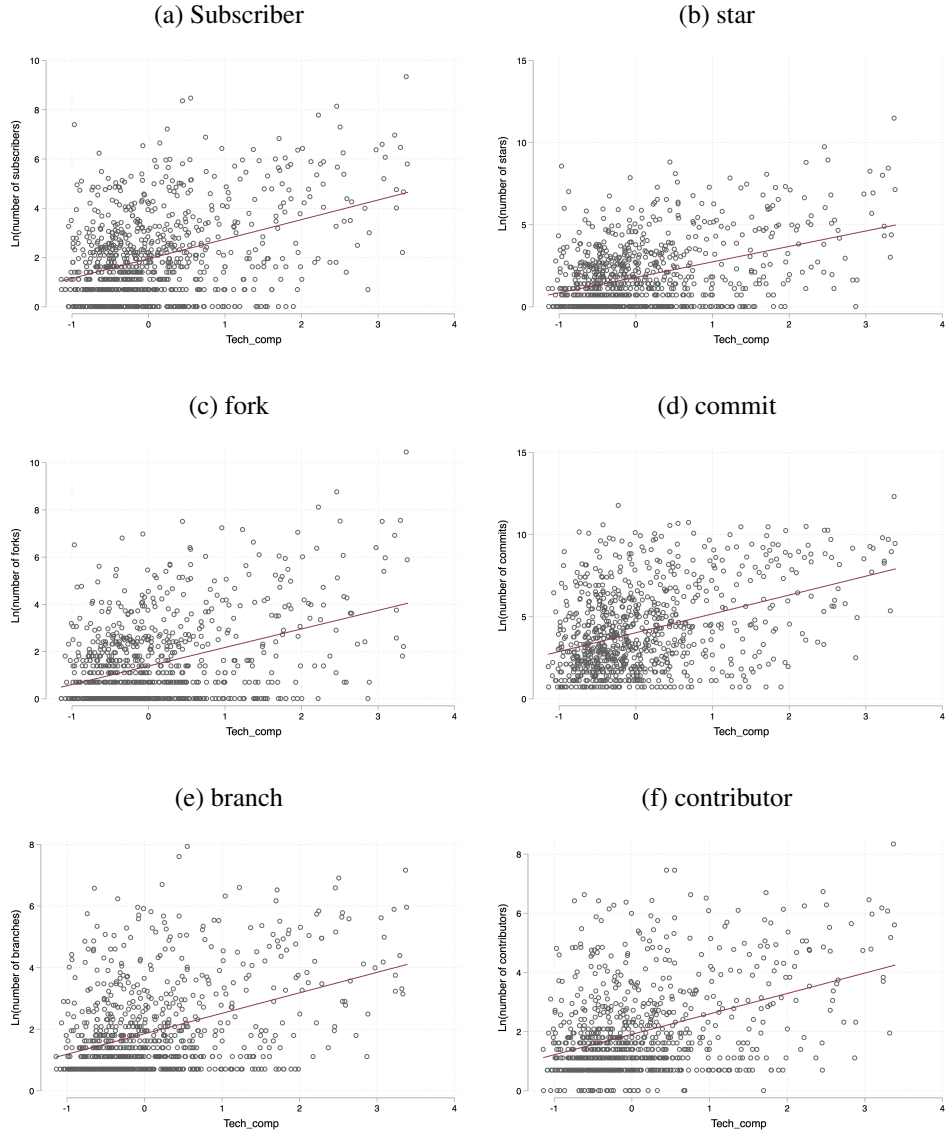


Table OA.1: **Embedding Key Terms**

This table displays the top 15 most frequent terms of each word embedding clustering and their topic labels. The number in parentheses indicates the percentage of terms belonging to the topic.

Topic Label	Most frequent terms
Information (1.8%)	user, datum, contract, transaction, information, process, access, wallet, order, node, would, account, public, key, store
Blockchain (2.9%)	platform, blockchain, system, network, base, smart, development, application, ethereum, chain, protocol, design, developer, open, software
Algorithm (2.7%)	one, model, block, follow, level, different, bitcoin, example, two, point, state, type, function, proof, algorithm
Healthcare (1.5%)	research, tool, health, professional, quality, analysis, report, knowledge, machine, patient, test, ai, medical, al, human
Business (3.5%)	market, business, technology, project, new, also, world, high, ecosystem, community, industry, crypto, solution, digital, cryptocurrency
Transportation (2.7%)	mining, energy, small, area, production, delivery, location, car, range, hardware, retail, gold, physical, home, producer
Marketing (2.5%)	product, content, customer, com, marketing, game, app, event, online, member, social, like, medium, program, marketplace
Token (2.1%)	token, exchange, sale, time, value, fund, payment, asset, coin, purchase, price, number, currency, cost, investment
Verb (1.7%)	use, provide, make, create, include, work, offer, allow, need, develop, increase, support, share, take, require
Operation (1.7%)	service, company, security, management, financial, legal, operation, partner, trust, bank, third, individual, activity, various, foundation
Negative words (2.5%)	result, change, case, without, problem, could, control, even, possible, low, however, due, loss, reduce, therefore
Country/Area (2.2%)	year, country, group, international, university, US, united, startup, partnership, estate, China, center, states, Singapore, Asia
Team/People (2.0%)	experience, advisor, founder, co, expert, manager, director, CEO, tech, entrepreneur, executive, degree, strategic, head, science
Disclaimer (1.3%)	may, risk, party, paper, right, future, whitepaper, term, white, part, form, law, document, person, limit
Typo (a) (2.9%)	ot, ond, dnd, con, ore, os, thot, doto, sen, ho, blockchoin, ds, hove, morket, ct
Typo (n,l) (3.0%)	th, ahd, tor, ih, wiii, hot, oh, ah, biockchain, tokehs, blockchaih, ts, aii, oi, tokeh
URL (3.5%)	https, mm, ii, en, http, de, iii, st, et, _, nd, er, pdf, ng, es
Roadmap (12.4 %)	team, www, launch, plan, page, main, io, utility, roadmap, ltd, usage, introduction, copyright, overview, disclaimer
Name/Brand (18.6%)	man, ago, litecoin, ltc, forex, paypal, anol, wp, sam, proj, nakamoto, hat, wire, cite, eight
Descriptive words (28.4%)	direct, org, flow, yes, late, successfully, old, previously, additionally, pro, soon, whose, ofa, maker, ten



Table OA.2: LDA Key Terms

This table displays the top 15 most relevant terms of each LDA topic. The number in parentheses indicates the relative prevalence of the topics in the corpus.

Topic Label	Most relevant terms ( $\lambda = 0.6$ )
Information (6.8%)	datum, health, data, patient, medical, healthcare, provider, identity, care, information, doctor, user, service, use, access
Blockchain (8.8%)	node, network, block, transaction, blockchain, proof, consensus, protocol, hash, use, chain, system, message, validator, contract
System (3.1%)	node, storage, cloud, file, quantum, datum, chain, compute, computing, de, system, application, blockchain, storage node, machine
Token (4.5%)	contract, order, exchange, trade, chain, network, asset, protocol, transaction, liquidity, user, token, fee, smart, dispute
Music/Travel (2.4%)	music, artist, travel, diamond, contract, driver, smart, smart contract, song, forest, industry, music industry, ride, cargo
Trading (6.7%)	trading, trader, market, platform, exchange, user, trade, ai, intelligence, crypto, strategy, system, service, use, development
Payment (9.8%)	payment, user, cryptocurrency, wallet, service, coin, merchant, card, exchange, transaction, use, currency, crypto, system
Business (9.3%)	business, product, token, sale, global, customer, consumer, year, blockchain, company, platform, technology, market, industry, experience
Finance (8.5%)	loan, asset, token, bank, credit, estate, platform, borrower, fund, financial, investor, lending, real estate, market, investment
Community (6.7%)	project, token, vote, community, platform, voting, fund, team, ico, user, member, reputation, bounty, crowdfunde, market
Marketing (7.2%)	token, platform, sale, user, marketing, token sale, team, online, use, social, tournament, player, service, advertising, development
Mining (3.2%)	mining, issuer, gold, der, currency, investment, die, mine, investor, EUR, crypto, holder, price, fund, coin
Law (6.1%)	may, token, company, purchaser, risk, law, include, car, purchase, regulation, party, platform, sale, jurisdiction, person
Disclaimer (3.8%)	whitepaper, distributor, statement, token, representation, forward, information, thereof, dissemination, look, constitute, risk uncertainty, person, warranty, uncertainty
Gamble (1.8%)	bet, ticket, gambling, player, casino, betting, sport, lottery, event, game, online gambling, jackpot, poker, online, gamble
Game (3.6%)	game, ad, advertiser, publisher, advertising, gamer, developer, gaming, AR, player, VR, game developer, virtual, user, games
Social Media (3%)	content, video, creator, ond, influencer, user, content creator, fan, medium, doto, viewer, social, tv, blockchoin, photo
Energy (1.6%)	energy, electricity, production, grid, water, solar, power, carbon, plant, renewable, green, waste, renewable energy, oil, fuel
typo (a) (1.4%)	dnd, ot, tor, ore, ond, sid, hove, tth, ds, cube, thot, con, hos, thdt, dny
typo (n,k) (1.7%)	ahd, tol, ih, insurance, wihi, oh, tokehs, ah, ens, blocl, ahy, to_lens, hot, blocl_chain, tol_en

Table OA.3: **Blockchain technology word list**

This table presents the complete word list that we use to count blockchain technology words as an alternative measure of technology sophistication.

accenture	DAPP	gigabyte	protocal
address	DDOS	halve	record
airdrop	DDOS attack	hard fork	relayer
altcoin	decentralize	harware wallet	reproduction
AML	decryption	hash	robustness
API	deposit	hashcash	Satoshi Nakamoto
ASIC	difficulty	hashrate	scalability
authentication	digital asset	hot wallet	scrypt
Bitcoin	digital identity	IBM	self execute
BTC	digital signature	immutable	serialization
block	distributed ledger	IPFS	server
block height	double spend	KYC	SHA-256
blockchain	EEA	ledger	shard
bounty	EIP	liquid democracy	smart contract
bug bounty	encryption	liquidity	soft fork
chain	ERC	mainnet	solidity
cipher	ETH	merkle tree	stable coin
client	Ether	multi signature	stablecoin
coin	Ethereum	NFT	testnet
cold storage	EVM	node	timestamp
cold wallet	exchange	oracle	transaction fee
collective	fiat	private key	validator
confirmation	fiat currency	public key	wallet
consensus	fork	proof	wallet address
cryptocurrency	gartner	proof of authority	workflow
cryptography	gas	proof of stake (PoS)	
DAO	genesis block	proof of work (PoW)	

**Table OA.4: Summary Statistics on Whitepaper Status**

This table lists all possible whitepaper status and their frequencies.

	<b>Frequency</b>	<b>Percent (%)</b>
Downloaded.	1629	55.90
URL response: client error.	535	18.36
URL response: server error.	104	3.57
Unable to get URL response.	403	13.83
Invalid PDF files.	155	5.32
Whitepaper not found.	54	1.85
Whitepaper is accessible but not downloadable.	27	0.93
Permission is required to access.	7	0.24
<b>Total</b>	<b>2914</b>	<b>100.00</b>

**Table OA.5: Summary Statistics on ICO Industries**

This table lists all ICO industries and their frequencies.

<b>Industry</b>	<b>Frequency</b>	<b>Percent</b>
Business	212	7.27
Charity	15	0.51
Connectivity	40	1.37
Cryptocurrency	384	13.17
Ecology	30	1.03
Finance	219	7.51
Games & Entertainment	174	5.97
Health & Medicine	83	2.85
Internet	56	1.92
Other	223	7.65
Platform	977	33.50
Production	30	1.03
Real Estate	65	2.23
Social Media	45	1.54
Software	94	3.22
Sports	17	0.58
Study	17	0.58
Trading	158	5.42
Transport	40	1.37
Travel	37	1.27
<b>Total</b>	<b>2916</b>	<b>100.00</b>

Table OA.6: **Technology Indexes Validation**

This table validates tech indexes with measures from GitHub. Panel A, B, C and D display the supervised, embedding-based, LDA-based and composite tech index respectively. For each column, *watch* measures the number of users subscribing repository updates; *star* indicates the number of “likes” received by the repository; *fork* proxies for the copies made by other developers; *commit* represents the number of times the code has been revised; *branch* is the amount of pointers to specific versions of the repository; and *contributor* reflects how many developers have contributed to the source code. All GitHub measures are in logarithmic forms. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Supervised Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	ln(watch)	ln(star)	ln(fork)	ln(commits)	ln(branch)	ln(contributor)
Tech_sup	0.561*** (0.064)	0.635*** (0.081)	0.530*** (0.073)	0.754*** (0.089)	0.429*** (0.052)	0.453*** (0.057)
Constant	1.985*** (0.056)	1.842*** (0.063)	1.428*** (0.055)	4.081*** (0.084)	1.887*** (0.047)	1.943*** (0.049)
Observations	861	861	861	861	861	861
R <sup>2</sup>	0.107	0.106	0.098	0.090	0.091	0.094
<b>Panel B: Embedding Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	ln(watch)	ln(star)	ln(fork)	ln(commits)	ln(branch)	ln(contributor)
Tech_embed	0.651*** (0.061)	0.792*** (0.073)	0.665*** (0.068)	0.995*** (0.078)	0.563*** (0.050)	0.580*** (0.051)
Constant	1.950*** (0.054)	1.795*** (0.060)	1.389*** (0.051)	4.018*** (0.080)	1.852*** (0.044)	1.908*** (0.046)
Observations	861	861	861	861	861	861
R <sup>2</sup>	0.148	0.170	0.158	0.160	0.161	0.158
<b>Panel C: LDA Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	ln(watch)	ln(star)	ln(fork)	ln(commits)	ln(branch)	ln(contributor)
Tech_lda	0.504*** (0.066)	0.600*** (0.079)	0.496*** (0.072)	0.726*** (0.088)	0.440*** (0.052)	0.454*** (0.056)
Constant	1.983*** (0.056)	1.836*** (0.062)	1.424*** (0.054)	4.072*** (0.083)	1.879*** (0.046)	1.936*** (0.048)
Observations	861	861	861	861	861	861
R <sup>2</sup>	0.095	0.105	0.095	0.092	0.106	0.104
<b>Panel D: Composite Index</b>						
	(1)	(2)	(3)	(4)	(5)	(6)
	ln(watch)	ln(star)	ln(fork)	ln(commits)	ln(branch)	ln(contributor)
Tech_comp	0.796*** (0.072)	0.940*** (0.089)	0.785*** (0.083)	1.148*** (0.094)	0.665*** (0.059)	0.691*** (0.062)
Constant	1.949*** (0.054)	1.796*** (0.060)	1.390*** (0.052)	4.023*** (0.080)	1.853*** (0.044)	1.908*** (0.046)
Observations	861	861	861	861	861	861
R <sup>2</sup>	0.161	0.174	0.161	0.155	0.164	0.163

Table OA.7: **Determinant of Technology Index**

This table presents the determinants of tech indexes. The dependent variable in panel A, B and C is the supervised, embedding-based and LDA-based tech index, respectively. Column (1) links the tech index to whether an ICO uses Ethereum blockchain; column (2) presents the relation between the tech index and GitHub commits (the number of code revisions); column (3) considers other text-based measures of ICO whitepapers; column (4) presents estimates with ICO characteristics; column (5) includes all variables. The reported t-statistics are based on robust standard errors. \*\*\*, \*\*, and \* indicate statistical significance at the 1%, 5%, and 10% levels respectively.

<b>Panel A: Supervised Index</b>					
	(1)	(2)	(3)	(4)	(5)
Ethereum Based	-0.102 (0.076)				-0.063 (0.073)
ln_commits		0.093*** (0.013)			0.066*** (0.013)
Has GitHub		-0.125** (0.062)			-0.136** (0.060)
Fog Index			-0.003*** (0.001)		-0.002 (0.002)
Tone			-0.211*** (0.035)		-0.193*** (0.033)
Uncertainty			0.028 (0.070)		-0.029 (0.062)
ICO Length				-0.002*** (0.001)	-0.002*** (0.000)
Team Size				0.031*** (0.004)	0.028*** (0.003)
Has Twitter				0.281*** (0.096)	0.246*** (0.095)
BTC Price (ICO)				-0.016* (0.009)	-0.009 (0.009)
Pre ICO				0.043 (0.053)	0.072 (0.052)
Bonus				-0.036 (0.059)	-0.007 (0.057)
Accept BTC				-0.074 (0.050)	-0.050 (0.049)
Constant	0.085 (0.071)	-0.124*** (0.039)	0.081 (0.065)	-0.414*** (0.123)	-0.348** (0.148)
$R^2$	0.001	0.052	0.026	0.075	0.119
Observations	1629	1629	1629	1483	1483

<b>Panel B: Embedding Index</b>					
	(1)	(2)	(3)	(4)	(5)
Ethereum Based	-0.294*** (0.077)				-0.172** (0.073)
ln_commits		0.122*** (0.013)			0.083*** (0.013)
Has GitHub		-0.119** (0.060)			-0.014 (0.061)
Fog Index			-0.004** (0.002)		-0.004 (0.003)
Tone			-0.282*** (0.040)		-0.227*** (0.039)
Uncertainty			-0.364*** (0.065)		-0.383*** (0.065)
ICO Length				-0.002*** (0.001)	-0.002*** (0.001)
Team Size				0.008** (0.004)	0.004 (0.004)
Has Twitter				0.150 (0.124)	0.075 (0.126)
BTC Price (ICO)				-0.027*** (0.010)	-0.016* (0.009)
Pre ICO				-0.115** (0.053)	-0.071 (0.051)
Bonus				-0.103* (0.055)	-0.094* (0.053)
Accept BTC				-0.134*** (0.051)	-0.095* (0.048)
Constant	0.247*** (0.072)	-0.193*** (0.037)	0.419*** (0.072)	0.202 (0.154)	0.569*** (0.175)
$R^2$	0.012	0.097	0.042	0.040	0.131
Observations	1629	1629	1629	1483	1483

<b>Panel C: LDA Index</b>					
	(1)	(2)	(3)	(4)	(5)
Ethereum Based	-0.262*** (0.078)				-0.132* (0.074)
ln_commits		0.098*** (0.014)			0.063*** (0.013)
Has GitHub		-0.100* (0.060)			-0.023 (0.062)
Fog Index			-0.002* (0.001)		-0.001 (0.002)
Tone			-0.304*** (0.036)		-0.257*** (0.034)
Uncertainty			-0.222*** (0.059)		-0.240*** (0.057)
ICO Length				-0.001* (0.001)	-0.001 (0.001)
Team Size				0.007* (0.004)	0.004 (0.004)
Has Twitter				0.205** (0.090)	0.148* (0.086)
BTC Price (ICO)				-0.012 (0.010)	-0.003 (0.010)
Pre ICO				-0.047 (0.054)	-0.015 (0.052)
Bonus				-0.092 (0.058)	-0.077 (0.056)
Accept BTC				-0.076 (0.053)	-0.047 (0.052)
Constant	0.220*** (0.073)	-0.151*** (0.036)	0.281*** (0.061)	-0.071 (0.130)	0.161 (0.140)
$R^2$	0.009	0.062	0.042	0.016	0.081
Observations	1629	1629	1629	1483	1483