# What's inside a node? Malicious IPFS nodes under the magnifying glass

Christos Karapapas[1], George C. Polyzos[1], and Constantinos Patsakis[2,3]

[1]Athens University of Economics and Business, Greece
[2]Department of Informatics, University of Piraeus, 80 Karaoli & Dimitriou str., 18534 Piraeus, Greece
[3]Information Management Systems Institute of Athena Research Centre, Greece

**Abstract**

InterPlanetary File System (IPFS) is one of the most promising decentralized off-chain storage mechanisms, particularly relevant for blockchains, aiming to store the content forever, thus it is crucial to understand its composition, deduce actor intent and investigate its operation and impact. Beyond the network functionality that IPFS offers, assessing the quality of nodes, i.e. analysing and categorising node software and data, is essential to mitigate possible risks and exploitation of IPFS. To this end, in this work we took three daily snapshots of IPFS nodes within a month and analysed each node (by IP address) individually, using threat intelligence feeds. The above enabled us to quantify the number of potentially malicious and/or abused nodes. The outcomes lead us to consider using a filter to isolate malicious nodes from the network, an approach we implemented as a prototype and used for assessment of effectiveness.

***Index terms***— InterPlanetary File System, Web3 Security

## 1 Introduction

Web 1.0 is known as the *read-only* web. For many years the Word Wide Web had an informative and educative role presented through static content. Few people generated content that was read by many people. From the constant interaction as well as the expanding familiarity that users had with the web, the number of users who also wanted to create content grew. Thus, the need for a more participative web arose, giving birth to Web 2.0. The latter, despite its shortcomings, is massively adopted. Single points of failure due to security incidents and control from a single organisation along with privacy violations for, e.g. marketing purposes, by centralised data storage facilities have been two of the most thorny issues in Web 2.0 for years.

Lately, there has been a lot of discussion around Web3. One of the pillars of Web3 is the decentralisation of the web to allow users to regain control over their data and selectively share and monetise the information they create. An integral part of attaining these goals are distributed ledger and blockchain technologies, and token-based economics [6, 13]. Web3 promises to offer decentralised services, meeting the needs of the Internet of Things (IoT) era and introducing a financial aspect of the web-user relationship through cryptocurrencies.

Web3 is considered to consist of different stacks, and these, in turn, of different protocols that cooperate with each other in order to provide services to the user. Some of them are data storage, domain name resolution, decentralised identities or, at a higher level, social media, gaming and marketplaces. All these different protocols, as well as the bridges between them, are still in their making; thus, their shortcomings may be exploited by attackers or utilised for malicious purposes. Indeed, ransomware and dark web marketplaces use cryptocurrencies to make siphon their payments, while blockchains and IPFS are used for the coordination of malware as C2 servers or to store malicious payloads.

In this work, we focus on distributed data storage and, more specifically, the InterPlanetary File System (IPFS), a cornerstone of the decentralised storage component of Web3. IPFS claims to have 2 million unique weekly users [1], and it has certainly caught the eye of the scientific community, as

---

[1]https://decrypt.co/resources/how-to-use-ipfs-the-backbone-of-web3

reflected by a total of more than 1160 papers found on Scopus with the search term "IPFS" in title and/or abstract. As IPFS is a collection of sub-protocols, it can be exploited by malicious users in a variety of ways. Immutability and decentralisation create a very dangerous mix that can be abused in various ways [5]. Karapapas et al. [9] illustrated how cybercriminals could exploit IPFS to set up an anonymous malware C2 facility alongside smart contracts. Patsakis et al. [14] showed that it could be abused to provide a robust malware C2 server infrastructure. Moreover, it is known to have been utilised by the Storm botnet[15] while new evidence has come to light linking IPFS to phishing [2].

To this end, we aim to unravel the structural elements of the IPFS network, and the nodes, focusing on suspicious activity. Initially, we crawl the IPFS network to enumerate it and make the first contact with the nodes. Following that, we collect intelligence from different sources regarding the aforementioned nodes. Moreover, we collect the exchanged data by nodes and analyse them to have a deeper understanding of the consistency of the network. Finally, we try to determine the extent of possible abuse of IPFS for copyright infringement.

The remainder of this paper is structured as follows: In Section 2, we present the required background information and overview of technologies. In Section 3, we present related research regarding IPFS monitoring. Section 4 focuses on nodes, presenting our data collection methodology and our findings regarding the nodes that constitute the IPFS network. Then, Section 5 goes a step higher in terms of abstraction, focusing on the content stored in IPFS. In Section 6, we discuss possible countermeasures to isolate malicious nodes. Finally, in Section 7, we summarise our findings and contributions, discussing possible future research directions.

## 2 Background

### 2.1 IPFS

InterPlanetary File System (IPFS) [3] is a peer-to-peer file-sharing system, consisting of many novel technologies, aiming to achieve decentralised data storage and low latency file distribution. Some of its main goals are to foster censorship circumvention and to avoid a single point of failure. E.g., in 2017 it was utilised to disseminate and store data regarding the Catalan independence referendum[3] when the Spanish government attempted to censor it.

Contrary to traditional file systems, in IPFS files are addressed by their content and each one is assigned a unique content ID (CID). One of the main IPFS components is `libp2p`[10], an umbrella term for many underlying network protocols. IPFS uses Distributed Hash Table (DHT), a highly scalable coordinator of data lookup among the different nodes. libp2p provides IPFS with the `KAD-DHT`, a Kademlia [12] variant. The latter is responsible for storing three types of mappings: 1. Provider Records, i.e., what content is hosted by whom, 2. Peer Records, i.e., who (PeerID) has what address and finally, 3. InterPlanetary Name System (IPNS) records, i.e., static names pointing to varying data. Another noteworthy component is BitSwap, which is a data-exchanging protocol based on `want-have content` and `have content` messages[7]. Moreover, Merkle DAG, a combination of Merkle Tree and Directed Acyclic Graph (DAG), is used to certify that the data exchanged are unique and IPFS does not store any duplicates. Finally, users can have access to files stored on IPFS, through HTTPs, by visiting public gateways. Public gateways have been provided not only by Protocol Labs, which is the main developer of IPFS but also by various companies embracing Web3, like Cloudflare, Pinata, etc (`https://ipfs.github.io/public-gateway-checker/`). As of July 2022, i.e., `v0.14`, the implementation of IPFS is known as `Kubo` [4].

## 3 Related Work

P2P networks have been of interest to the scientific community for many years, and while their popularity fluctuates, they have never been outdone. In recent years, the advent of cryptocurrencies and blockchain technology has brought them back into the limelight. Thus, while P2P node profiling has been extensively studied in the past, research in the context of Web3 is minimal. Web3 is in a very early phase and its decentralised components are still under heavy development. Hence, the current research regarding its nodes is still in its infancy. Henningsen et al. in [8] make one of the first attempts to explore the IPFS

---

[2]`https://www.trustwave.com/en-us/resources/blogs/spiderlabs-blog/ipfs-the-new-hotbed-of-phishing`
[3]`https://edri.org/our-work/no-justification-for-internet-censorship-during-catalan-referendum/`
[4]`https://github.com/ipfs/kubo/blob/master/docs/changelogs/v0.14.md`

network. Adopting a hybrid design, passively and actively, they aim to enumerate the IPFS network and profile its nodes. The authors note that the overlay network outperforms the overlay induced by buckets. Furthermore, they observe that an overwhelming percentage of nodes, i.e. 94%, did not react to the authors' attempt to connect to them. The reason this happens is twofold. The first is because many nodes are behind NAT and thus advertise their local IP address. The second is that a large portion of users uses IPFS in an opportunistic way, therefore their footprint remains in buckets for longer than they remain online and connected.

Recently, researchers discovered a botnet hiding in the IPFS ecosystem [15]. The latter, named InterPlanetary Strom (IPStorm) and estimated size of 9000 devices, utilises IPFS at multiple levels. Initially, the researchers found that it uses the libp2p DHT to discover nodes. Bots identify each other with the attribute `Agent Version`: "storm". In addition, the botnet utilises the Pub/Sub protocol as a communication channel over specific topics. Finally, the botnet uses IPFS to share files so that it can be updated to a newer version.

Trautwein et al. [16] further to providing a basic guide of IPFS' design, they collected data from three different sources to shed light on various metrics related to IPFS performance. Initially, they crawled the IPFS network to gather information about peers. Among the conclusions drawn is that IPFS nodes are geographically distributed in 152 countries, yet more than 50% are located in just two countries, US and China. Furthermore, more than 50% of the IPs are covered by five automated systems, yet only 2.3% of the nodes are in some cloud infrastructure. The last insight extracted from this dataset is that the IPFS network suffers from high rates of churn, with 87.6% of peers having an uptime of less than 8 hours. Finally, the authors wanted to study the time performance in downloading data. To this end, they experimented with different AWS regions and recorded the download duration from the data they produce each time. In 50% of the cases, the download took less than 3s, and in 90% of the cases, less than 4.5s.

# 4 Profiling IPFS nodes

## 4.1 Data collection methodology

To enumerate the IPFS network we used the IPFS Crawler [8]. The IPFS crawler is a tool written in Go and is based on libp2p (v0.11.0). Acting as a Kademlia node the crawler uses precomputed keys to extract all the entries from most buckets for every node it encounters. In essence, it invokes FINDNODE actions repeatedly using the appropriate precomputed keys. Finally, the crawler produces two files: (i) a JSON file storing the tuple <`PeerID, multiaddress, agent, reachability`> for every distinct node met, and (ii) a CSV file containing all the pairs of connected nodes.

We conducted a series of consecutive crawls. Initially, the crawls were performed iteratively, every ten days during the period from March to April 2022. Each crawl series spanned over a day (24h) totalling about 360 crawls in a row per day. From the data in the JSON file, for each PeerID we extracted the IP addresses. Each IPFS node maintains an address book retaining information for the nodes it encounters. If any of the encountered nodes advertises a new address, then it is appended in the address book for reachability purposes. As a result, a single PeerID may correspond to more than one IP address. We studied each different address considering it as a unique node. Moreover, nodes behind a firewall or NAT use `p2p-circuit`, a libp2p relay transport protocol, to avoid connectivity barriers. In essence, these nodes advertise addresses through relay nodes. As a consequence, they do not reveal their real IP address but the IP address of the relay. The aforementioned peers as well as those which advertise only local IP addresses are excluded from our analysis. Clearly, the absence of such IP addresses prevents us from studying or fingerprinting the corresponding hosts.

## 4.2 Node Profiling

In this section, we present general information regarding the IPFS network and its nodes. We should mention that in the following findings, every different IP address is considered a different node. Although, we found that unique Peer IDs advertise multiple IP addresses since our study focuses on the "fabric" of the IPFS network. Thus, we want to enumerate and analyse every different IP address.

Figure 1 illustrates the nodes per crawl and the count of malicious nodes for which we collected intelligence. In Figure 2a, the exact results of IP addresses per crawl can be found. Moreover, from the same figure, we can observe that 16783 were found online in all three crawls. We can assume that the aforementioned nodes were found online at least once a day in the span of the whole month. Given the
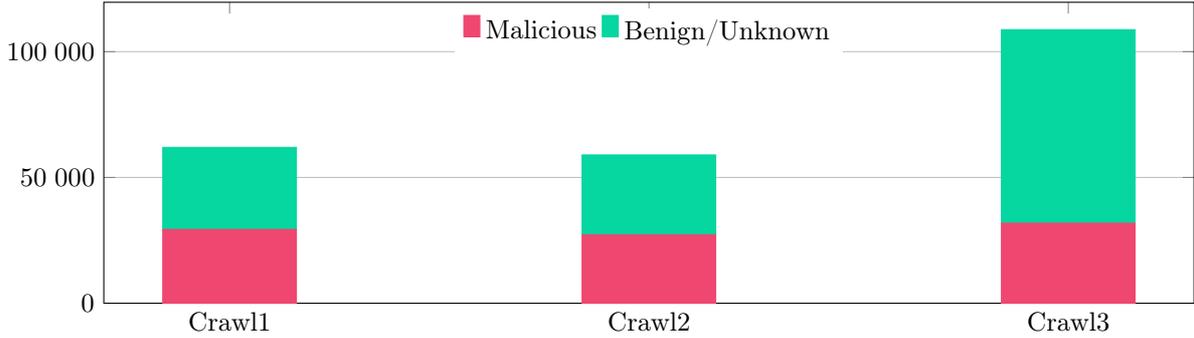
Figure 1: Malicious nodes per crawl.

periodic changes of IPs, we can assume that most of these IPs belong to some infrastructure that has been devoted to constantly working with IPFS.
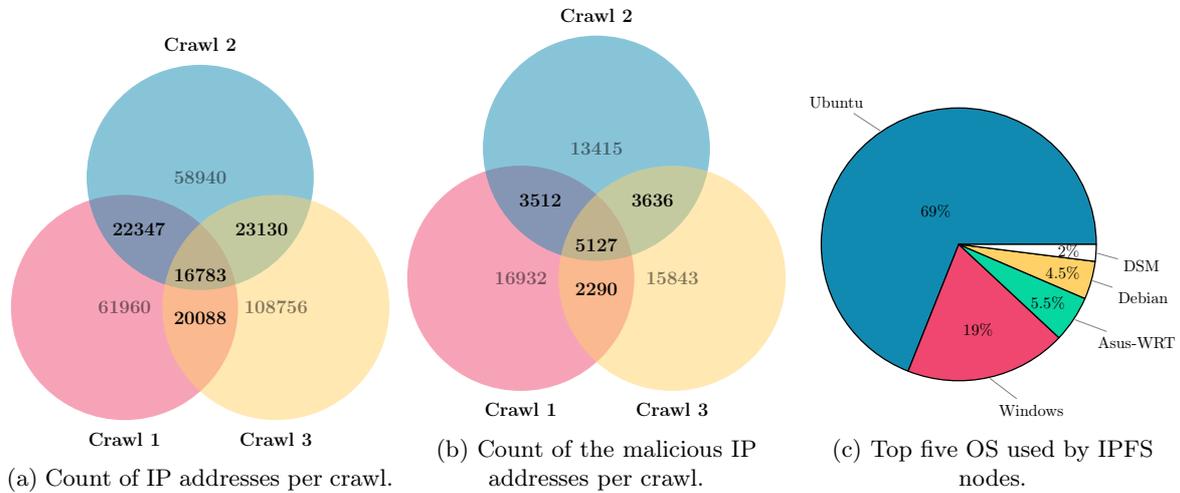


(a) Count of IP addresses per crawl.

(b) Count of the malicious IP addresses per crawl.

(c) Top five OS used by IPFS nodes.

Figure 2: Crawl statistics.

A node's agent version can be an indication of malicious activity. Nodes' agent version is public and advertised, thus, it can act as an identifier for malicious nodes to discover and track each other. The latter is a technique already implemented by "storm" agents. Figure 3 illustrates the ten most used agent versions we found in each crawl. We should highlight that the counts depicted correspond to the agents from the nodes we managed to connect to. In each crawl we found 50%, 61%, 49% respectively, unreachable peers, i.e., we found their address stored in the DHT but they were offline. Moreover, IPFS is open-source software; therefore, it is at the user's discretion whether to display the agent version. The latter results are aligned with the ones in [16]. In the third crawl we observe that there is an increase in nodes using the agent called `Hydra Booster` [5]. Hydra Booster is a node having many different Peer IDs over a common routing table. It is designed to accelerate IPFS' processes carried out through DHT-like content resolution, routing and discoverability. The existence, as well as the operation of these nodes, brought about an increase in the number of nodes of the third crawl. One of the features of open software, which has been hotly debated lately, is that upgrading to a newer version is at the user's discretion. Observing the crawling results of Figure 3, one can observe that there are many different software versions running and communicating simultaneously. For example, `go-ipfs 7.0` was released in July of 2020 while `go-ipfs 11.0` in August of 2021. Moreover, although the measurements were made in mid-2022, and version `go-ipfs 12.0` had already been released, we can conclude from the bar charts that the versions which are more widely used are the older ones. In addition, we must mention that agent storm, which has been found in all three crawls with a non-negligible number, is characteristic of the nodes belonging to the IPStorm botnet we have already mentioned.

In what follows, we study the maliciousness of nodes, so we used Virus Total to assess the correspond-
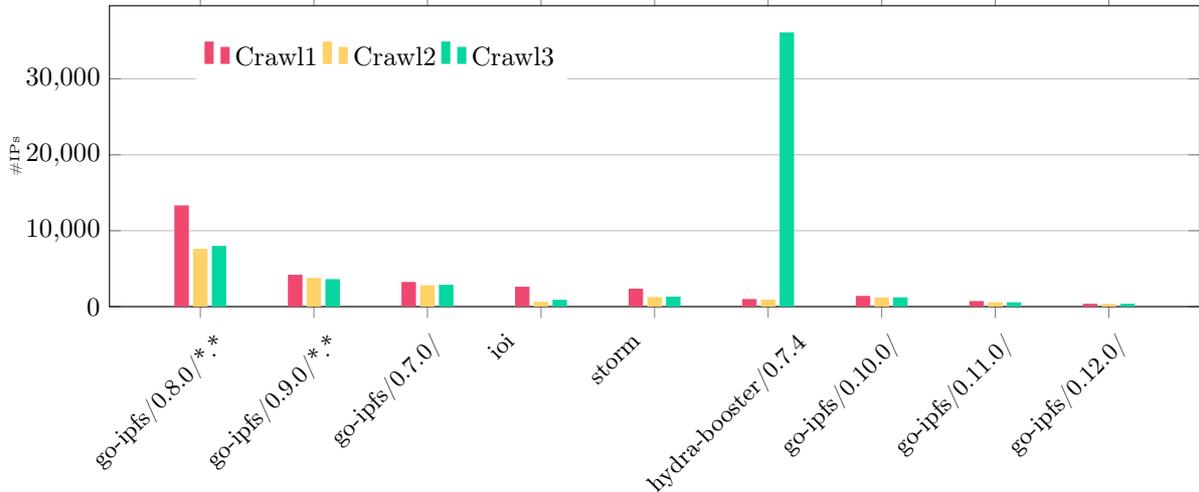
---

[5]`https://github.com/libp2p/hydra-booster/`

Figure 3: The ten most commonly used agent versions in each crawl. The `*.*` denotes varying subversions combined.

ing IPs. Nevertheless, Virus Total also provides valuable insights regarding the geographic distribution of the various nodes, regardless of whether they are malicious or not. The vast majority of the nodes are located in two countries, namely the United States and China. We notice that our results are aligned with [16].

To conduct a more in-depth analysis, we passed the crawling results to intelligence services. Namely, we used Shodan, a network monitoring tool, to fingerprint each node. Shodan returned intelligence for approximately 40960 unique nodes. Figure 4 illustrates the ten most commonly used ports by the total of nodes we examined. Port 22, the most widely used port by IPs related to IPFS, is typically used for Secure Shell (SSH) connections, which allow users to log in to a host and execute commands remotely. Port 80 is used as the default port for HTTP (Hypertext Transfer Protocol) traffic, port 8080 is an alternative to port 80 and moreover the default port of the IPFS gateway, and port 443 for HTTPS. Port 3389 is typically used by hosts running Microsoft Remote Desktop Protocol (RDP) to allow remote access to the host's desktop. Finally, port 4001 is used by default for IPFS traffic, but users can also set up a custom port. Regarding the operating system running on IPFS nodes, Shodan's results, depicted in Figure 2c, indicate that the lion's share uses Ubuntu Linux. The next runner-up is Microsoft Windows 10, followed by Debian Linux. The latter is also exhibited by the most used services, Figure 5, where most hosts appear to be using SSH as opposed to RDP. Moreover, most of them seem to have a web server (nginx and then Apache).
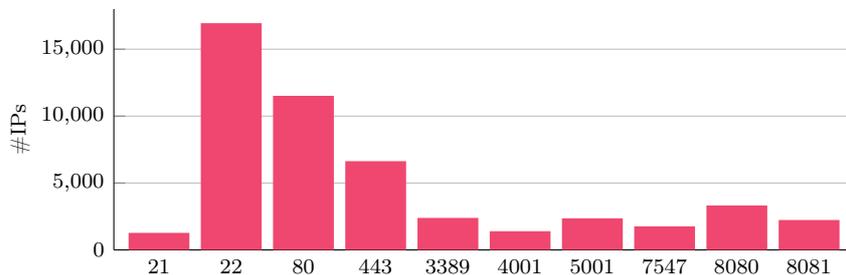


Figure 4: The ten most common ports.

JARM [1] is an open-source fingerprinting tool that generates a string based on the response of the host to ten TLS packets. JARM is used by the community as a software-wise host clustering tool, therefore it is also eligible to detect malware Command & Control (C2). We use JARM strings, extracted from Shodan and Virus Total, to detect any similarities among the different nodes. Finally, we combined them since for the same IP different services can provide varying information. For 1002 IP addresses, we found information in both services, so we considered both records. The JARMs indicate that there are several clusters of IPs in which servers have the same TLS configuration, which implies that the same
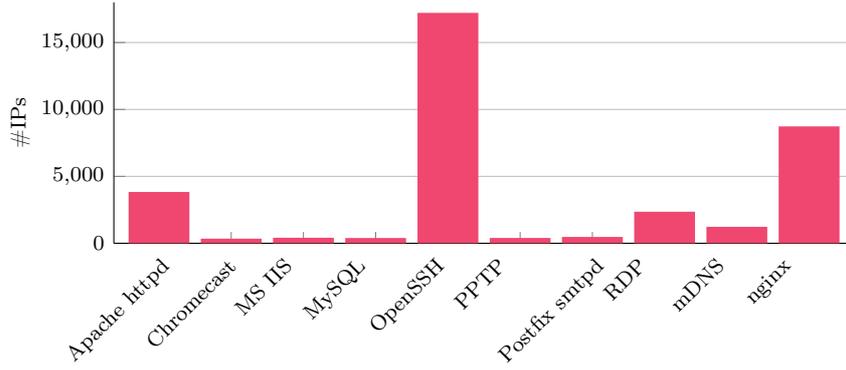
Figure 5: The ten most commonly used services.

entity is behind them. The most common ones are illustrated in Table 1.

| JARM | # IPs |
|------|-------|
| 2ad2ad0002ad2ad00042d42d0000008aec5bb03750a1d7eddfa29fb2d1deea | 2070 |
| 2ad2ad16d2ad2ad22c2ad2ad2ad2adfd9c9d14e4f4f67f94f0359f8b28f532 | 1378 |
| 15d3fd16d29d29d00042d43d000000fe02290512647416dcf0a400ccbc0b6b | 577 |
| 15d3fd16d29d29d00042d43d0000009ec686233a4398bea334ba5e62e34a01 | 562 |
| 15d3fd16d21d21d00042d43d000000fe02290512647416dcf0a400ccbc0b6b | 489 |

Table 1: Most common JARMs.

## 4.3   Malicious Activity

In this section, we investigate the moral character of IPFS nodes, i.e., we examine whether and to what extent there are malicious nodes. To this end, we collect and leverage existing intelligence to create and present their profile. Our goal is to assess the network structure, keeping IPFS users and the related community alert to the existence of malicious activity in the IPFS network. Due to the current IPFS rules, every node maintains several active connections varying from 600 to 900 peers. Thus, we argue that it is very important for each node to know what kind of alignment, i.e. neutral or malicious, the node it interacts with has.

Initially, we leveraged the intelligence provided by two popular services, namely Virus Total (`https://virustotal.com`) and SpamHaus (`https://www.spamhaus.org/`), to get a baseline for the reputation and past activity of nodes. SpamHaus uses several methods to find information about an internet resource. It uses sensors in large networks, i.e. a data-sharing community, from which it collects data about network traffic. In addition, SpamHaus deploys honeypots to attract malicious users. Along the same lines as SpamHaus and VT, in addition to monitoring more than 70 anti-malware and IP blocking services, it relies on data generated and shared by an already large community. Both the aforementioned services provide APIs to interact with their knowledge base and generate a JSON formatted output for each request. We combine the extracted output information with the SpamHaus output and we consider malicious those nodes with at least one record in one of the aforementioned services.

Moreover, in Figure 2b, we notice that from the 27861 different IP addresses we encountered during the first crawl, 5126 of them, $\approx 18\%$ remained online throughout the whole month. The latter indicates that there is a number of nodes that constantly utilise the IPFS network for malicious purposes. Compared to the 16783 found online in all three crawls, as depicted in Figure 2a, a significant part of them, i.e., 30.5%, are known to be malicious. Based on SpamHaus' results, we conclude that the majority of malicious nodes were discovered using the `DNS Sinkhole` technique. According to this technique, security researchers create, at various levels, a DNS record of a known malicious URL pointing to an address they own, usually a sinkhole server. The gain from applying this technique is twofold: On the one hand, they prevent communication between bot and C2, and on the other hand, researchers can find which computers are infected, i.e. ask to connect to known malicious URLs.

In Table 2a, the five most commonly requested and sinkholed URLs in the number of unique IP addresses are illustrated. Note that several URLs such as `differentia.ru`, `atomictrivia.ru`,

`amnsreiuojy.ru` and `restlesz.su` are known to be leveraged as C2 by malware. `disorderstatus.ru` is a relatively newly created domain reported to be mostly used for spamming. To draw deeper conclusions about the URLs, we isolated the `Top Level Domain` (TLD) of the different requested URLs. To our surprise, while most requested URLs have a ".ru" TLD, this is not reflected among the unique TLDs. On the contrary, we notice that the most commonly encountered is ".xyz", a relatively new TLD offering many domains that would traditionally be registered by legitimate users. The fact that they are new and cheap and that traditional domain names are available has led `xyz` domains to be widely exploited[6]. Given that 11227 `xyz` domains are hosted by these addresses makes us conclude that some adversaries use nodes of IPFS for hosting malicious domains in addition to C2 infrastructure.Tinba a portmanteau of the words Tiny Banker, is a trojan that leverages packet sniffing to determine whether the user visits a bank's webpage. In that case, the trojan tries to steal the keystrokes and sends them to a C2. Nymaim and Ranbyus are well-known trojans, which steal information from the user and consequently send them to a C2. Some of their variants have been found to use domain fluxing to communicate with their orchestrator, and some have been found in DoS attacks. Mirai is used to infect Internet of Things (IoT) devices and turn them into bots that can be used to launch large-scale network attacks. The Mirai botnet was initially discovered in 2016 and was part of various high-profile cyberattacks, including distributed denial-of-service (DDoS) attacks that brought down popular websites and online services. The most frequently displayed campaigns are gathered in Table 2b.

| URL | Count |
|-----|-------|
| differentia.ru | 38681 |
| disorderstatus.ru | 15504 |
| atomictrivia.ru | 7049 |
| amnsreiuojy.ru | 5662 |
| restlesz.su | 2180 |

(a) The five most sinkholed URLs and the number of unique requests.

| Campaigns | Count |
|-----------|-------|
| tinba | 30019 |
| conflicker | 22650 |
| nymaim | 22228 |
| andromeda | 6403 |
| ranbyus | 4845 |
| mirai | 3750 |

(b) Malware campaigns with the largest participation from the encountered nodes.

Table 2: Extroversion of malicious nodes: Which groups do they belong to and what webpages they seek to visit.

Finally, we studied the JARMs of malicious nodes to better frame our research. As we have already mentioned, we combined knowledge from all intelligence services to produce the results. Notably, among them, we found a cluster of 68 nodes corresponding to the JARM fingerprint `15d3fd16d29d29d00042d43d00` `00009ec686233a4398bea334ba5e62e34a01` which is attributed to the notorious `emotet` botnet.

As already mentioned, the crawler we used, in addition to information about the nodes encountered, produces an edge list with each pair of connected nodes. Based on this, we constructed a mapping from one PeerID to the several PeerIDs we found connected during the second day. In essence, we built for each peer its buckets expanded to the span of a day. Consequently, we converted the aforementioned mapping to the corresponding IP addresses. This way, we can investigate whether there is a clique between the malicious nodes. The findings indicate that there is no such clique, as the median percentage of malicious nodes in the buckets of a malicious node is 7%, and the average is 9.5%. Along the same lines, the median percentage of nodes in the buckets of a benign node is also 7%, with the average being 9.2%.

# 5 File Investigation

Despite the processes and functionality IPFS offers through libp2p and its other components, its main purpose is undeniably storage-related. The largest NFT marketplaces use IPFS for the data storage and integrity it provides, while its widespread utilisation has already brought about the need for cooperation with other Web3 layers, such as ENS, which natively offers names corresponding to CIDs. No wonder the increasing popularity has also caught the eye of cyber criminals. A recent research [7] highlights that the volume of malware samples hosted in IPFS has increased during 2022. Moreover, researchers report the `Agent Tesla` malware, which using phishing techniques, leads to an IPFS public gateway,

---

[6]`https://www.spamhaus.com/resource-center/getting-the-low-down-from-xyz-registry-on-combating-domain-abuse/https://www.b`
[7]`https://blog.talosintelligence.com/ipfs-abuse/`

disguising the download of malicious content. To better frame our research into the storage of the IPFS ecosystem, we also researched the file side. Our research is twofold, in the first case, we eavesdropped on the files requested by IPFS users, while in the second, more actively, we searched for files we randomly downloaded from well-known torrent sites.

## 5.1 Bitswap Eavesdropping

According to the operating rules of IPFS, when a user searches for a file, a one-hop inquiry is first performed through Bitswarm, requesting it from nodes with an active connection to the initiator. If none of them responds, the query is then served by the DHT. To collect data, we tweaked our node so that it maintains active connections with around 4000 nodes; that is, according to our measurements, approximately 20% of the network's active nodes at that time. So when one of those nodes was looking for a file, thanks to Bitswap's functionality, that information would also go through us. This way, we could eavesdrop on about 20% of the network's requests and, in turn, request back to retrieve them. In total, we monitored the requests for 24 hours while we set each request to last no more than 15 seconds. This way, we avoided downloading very large files while, on the other hand, we cancelled the search in case it was routed through the DHT. In total, we collected 49155 files with a size of about 13.7 GB. To have a more complete picture of the type of files requested, we used the Python `mimetypes` module[8] to find the MIME type of each file. We shall mention that it managed to classify 13691 of the files. The latter can be attributed to Bitswap's design. When a user requests a file from Bitswap, the search is performed by the root CID of the file. The aforementioned file contains links to the chunks of which it is composed. Thus, when the requester receives the root CID and learns the CIDs of the chunks that make up the file, it requests through Bitswap consecutively all the chunks, which are essentially blocks of data. The file results illustrate that 3716 are image files with MIME types "image/png", "image/gif", "image/jpeg", and 9148 are JSON files, which is the most common format for NFT metadata. The latter clearly demonstrates and confirms our initial statement that IPFS is a cornerstone of NFT data storage and Web3 in general. Among others, we fetched 177 Javascript files and 27 videos of type "video/mp4". We then fed the image files to the Python `Not Suitable For Work (NSFW) Detector` module to determine whether IPFS is being used for inappropriate content. From the 1636 image files it examined successfully, it found 33 unsuitable [9]. The above indicates that some users leverage IPFS' anonymity to host inappropriate content that is difficult for LEAs to track and take down.

## 5.2 Torrent Files

Very often, inappropriate files are found in the form of torrent files disseminated through torrent search engines. We downloaded a sample from various widespread torrent sites, ten popular torrents in total. We computed their CIDs locally to determine whether they are shared on the IPFS. This way, not only did we not add any illegal files to the IPFS network, but we also limited the possibility of tampering with the results of our upcoming searches. The ten different torrent files yielded 72 different root CIDs. Each torrent file can contain a video file, a cover image for the video file, a text file with information about the file, etc. In turn, we made 72 requests to the DHT for providers of these CIDs. We found providers for seven of them, and in fact, for most of them, more than one. The latter implies that IPFS users may also share the same content in torrents and that intellectual infringement content is also distributed through IPFS.

# 6 Countermeasures

The amount of malicious nodes connected to IPFS is alarmingly high. Given the P2P nature of IPFS and its continuous exploitation, we believe that pruning nodes from the network might provide an initial measure of sanitising the network; otherwise, the benign peers facilitate the malicious ones. To this end, we opt for a periodical blacklist approach that is resolved through InterPlanetary Name System (IPNS). In essence, we propose the use of the proposed data crawling methodology to monitor the nodes on daily basis, the IPs are collected and using intelligence services, we determine whether the IP should be blocked or not. Each IP is four bytes long, so the expected size is rather small and easy to manage. For instance, using our experiments as a baseline, using the worst estimate of 32000 malicious nodes, the blocklist would be around 8KB if the IPs are directly stored. Given its size and possible optimisations (e.g.

---

[8] https://docs.python.org/3/library/mimetypes.html
[9] https://pypi.org/project/nsfw-detector/

use binary search over the sorted list), searching whether the connected peers are malicious can be very efficient. Moreover, since the amount of nodes is tolerable, the collection of data from intelligence services can be rather fast. Of course, one could mask the IPs with the use of Bloom filter-based approaches [4]. However, since the IPs are very small in size, the corresponding filters would for sure be significantly bigger. Moreover, due to the probabilistic nature of bloom filters, there is always a margin of considering a benign node as malicious. That said, for our experiments, with an error probability of only 0.001%, 93.6KB would be required, which is one scale of magnitude more.

IPFS is becoming institutional, after all, there are many organisations participating in it and supporting it. Recent research efforts indicate that it could frame the existing banking system [11], while at the same time it constitutes a cornerstone of Decentralised Finance (DeFi). Our research does not intend to act as a brake on its use, on the contrary, it intends to inform, alert and promote its secure use. For instance, the network administrator of an organisation participating in the IPFS network can block the traffic towards and from a suspicious IP address by adding a rule to the firewall. Note that it can also remove alert fatigue from SOCs who might observe malicious IPs connected to the monitored infrastructure due to IPFS traffic. Finally, while IPFS provides the ability to disconnect from a node, it does not provide natively the option for the user to maintain a blacklist.

# 7 Conclusions

Open and decentralised systems are, by their very nature, prone to several attacks. However, given the crucial role of IPFS for Web3, it is essential to protect the ecosystem. Our measurements indicate that an alarming number of IPs reported as malicious through intelligence services are using IPFS. Rather than making it centralised, we opt for soft measures that allow nodes to isolate malicious ones selectively. We argue that this isolation can significantly benefit the network as the content of most of these nodes may be malicious, leading legitimate ones to facilitate nefarious acts and malicious campaigns. Therefore, their isolation, in the long run, may increase the robustness of the network and trust in it.

IPFS seems to have sacrificed part of the privacy to succeed in terms of performance, speed, and robustness [2]. This shortcoming can be exploited for malicious purposes, but it can also be leveraged by security analysts to monitor malicious nodes. Thus, apart from the fact that we can obtain critical information regarding a malicious node, such as its IP address, we can also monitor it from a content point of view, i.e., its requests as well as what it provides. Therefore, a future direction of this work is an extension of the implementation of the proposed filter so that it associates malicious nodes with the corresponding content.

# Acknowledgements

# References

[1] John Althouse. Easily Identify Malicious Servers on the Internet with JARM, 2020.

[2] Leonhard Balduf, Sebastian Henningsen, Martin Florian, Sebastian Rust, and Björn Scheuermann. Monitoring data requests in decentralized data storage systems: A case study of IPFS. In *2022 IEEE 42nd International Conference on Distributed Computing Systems (ICDCS)*, pages 658–668. IEEE, 2022.

[3] Juan Benet. Ipfs-content addressed, versioned, p2p file system. *arXiv preprint arXiv:1407.3561*, 2014.

[4] Burton H Bloom. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426, 1970.

[5] Fran Casino, Eugenia Politou, Efthimios Alepis, and Constantinos Patsakis. Immutability and decentralized storage: An analysis of emerging threats. *IEEE Access*, 8:4737–4744, 2019.

[6] Allan V Cook, Mike Bechtel, Siri Anderson, David R Novak, Nicole Nodi, and Jay Parekh. The Spatial Web and Web 3.0: What business leaders should know about the next era of computing. *Deloitte Insights*, 2020.

[7] Alfonso De la Rocha, David Dias, and Yiannis Psaras. Accelerating Content Routing with Bitswap: A multi-path file transfer protocol in IPFS and Filecoin, 2021.

[8] Sebastian Henningsen, Martin Florian, Sebastian Rust, and Björn Scheuermann. Mapping the Interplanetary Filesystem. In *2020 IFIP Networking Conference (Networking)*, pages 289–297, 2020.

[9] Christos Karapapas, Iakovos Pittaras, Nikos Fotiou, and George C Polyzos. Ransomware as a service using smart contracts and IPFS. In *2020 IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, pages 1–5. IEEE, 2020.

[10] Protocol Labs. Libp2p. `https://libp2p.io/`, 2021.

[11] Abdullah Al Mamun, Sheikh Riad Hasan, Md Salahuddin Bhuiyan, M. Shamim Kaiser, and Mohammad Abu Yousuf. Secure and Transparent KYC for Banking System Using IPFS and Blockchain Technology. In *2020 IEEE Region 10 Symposium (TENSYMP)*, pages 348–351, 2020.

[12] Petar Maymounkov and David Mazieres. Kademlia: A peer-to-peer information system based on the xor metric. In *International Workshop on Peer-to-Peer Systems*, pages 53–65. Springer, 2002.

[13] Alex Murray, Dennie Kim, and Jordan Combs. The promise of a decentralized Internet: What is web 3.0 and HOW can firms prepare? *Business Horizons*, 2022.

[14] Constantinos Patsakis and Fran Casino. Hydras and IPFS: a decentralised playground for malware. *International Journal of Information Security*, 18(6):787–799, 2019.

[15] Silvia Pripoae. Looking Into the Eye of the Interplanetary Storm, 2020.

[16] Dennis Trautwein, Aravindh Raman, Gareth Tyson, Ignacio Castro, Will Scott, Moritz Schubotz, Bela Gipp, and Yiannis Psaras. Design and evaluation of IPFS: a storage layer for the decentralized web. In *Proceedings of the ACM SIGCOMM 2022 Conference*, pages 739–752, 2022.