



OSTBAYERISCHE
TECHNISCHE HOCHSCHULE
REGENSBURG



MASTERTHESIS

Valentin Brandl

Collaborative Monitoring of Fully Distributed Botnets

9th April 2022

Faculty:	Informatik und Mathematik
Study Programme:	Master Informatik
Supervisor:	Prof. Dr. Christoph Skornia
Secondary Supervisor:	Prof. Dr. Thomas Waas

Botnets pose a huge risk on general internet infrastructure and services. Distributed P2P topologies make it harder to detect, monitor and take those botnets offline. This work explores ways to make monitoring of fully distributed botnets more efficient, resilient and harder to detect, by using a collaborative, coordinated approach. [do me](#)

Keywords— P2P, botnet, monitoring, collaboration

Contents

1	Introduction	6
2	Background	7
2.1	Formal Model of P2P Botnets	9
2.2	Monitoring Techniques for P2P Botnets	10
2.2.1	Passive Monitoring	10
2.2.2	Active Monitoring	11
2.2.3	Monitoring Prevention Techniques	12
3	Methodology	13
3.1	Protocol Primitives	14
4	Coordination Strategies	16
4.1	Load Balancing	16
4.1.1	Round Robin Distribution	17
4.1.2	IP-based Partitioning	17
4.2	Reduction of Request Frequency	20
4.3	Creating Edges for Crawlers and Sensors	21
4.3.1	Use Other Known Sensors	26
4.3.2	Use Churned Peers After IP Rotation	28
4.3.3	Peers Behind Carrier-Grade NAT	28
5	Implementation	31
6	Conclusion, Lessons Learned	34
7	Further Work	35
	List of Figures	37
	List of Tables	38
	List of Listings	39

Contents

List of Acronyms	40
References	41

Todo list

do me	2
wording	7
explain structured and unstructured	9
more details	10
no context	11
good title	12
remove?	16
1 – 2 sentences about naive rr?	17
page numbers for forward refs?	17
better wording	17
don't use substrings, bit.int for 128 bit modulo, argumentation why this works	17
remove this?	19
better numbers for example?	20
timeline with peers per bucket	25
use better data?	25
repeat analysis, actual number	26
upper limit for NL size as impl detail	26
übergang	28
what is an AS	28
decide	34

1 Introduction

The internet has become an irreplaceable part of our day-to-day lives. We are always connected via numerous “smart” and internet of things (IoT) devices. We use the internet to communicate, shop, handle financial transactions, and much more. Many personal and professional workflows are so dependent on the internet, that they won’t work when being offline, and with the pandemic we are living through, this dependency grew even stronger.

In 2021 there were around 10 billion internet connected IoT devices and this number is estimated to more than double over the next years up to 25 billion in 2030 [17]. Many of these devices run on outdated software, don’t receive regular updates, and don’t follow general security best practices. While in 2016 only 77% of German households had a broadband connection with a bandwidth of 50 MBit/s or more, in 2020 it was already 95% with more than 50 MBit/s and 59% with at least 1000 MBit/s [4]. Their nature as small, always online devices—often without any direct user interaction—behind internet connections that are getting faster and faster makes them a desirable target for botnet operators. In recent years, IoT botnets have been responsible for some of the biggest distributed denial of service (DDoS) attacks ever recorded—creating up to 1 TBit/s of traffic [10].

2 Background

Botnets consist of infected computers, so called *bots*, controlled by a *botmaster*. *Centralized* and *decentralized botnets* use one or more coordinating hosts called *command and control (C2) servers* respectively. These C2 servers can use any protocol from internet relay chat (IRC) over hypertext transfer protocol to Twitter [19] as communication channel with the infected hosts. The abuse of infected systems includes several activities—DDoS attacks, banking fraud, proxies to hide the attacker's identity, sending of spam emails. . .

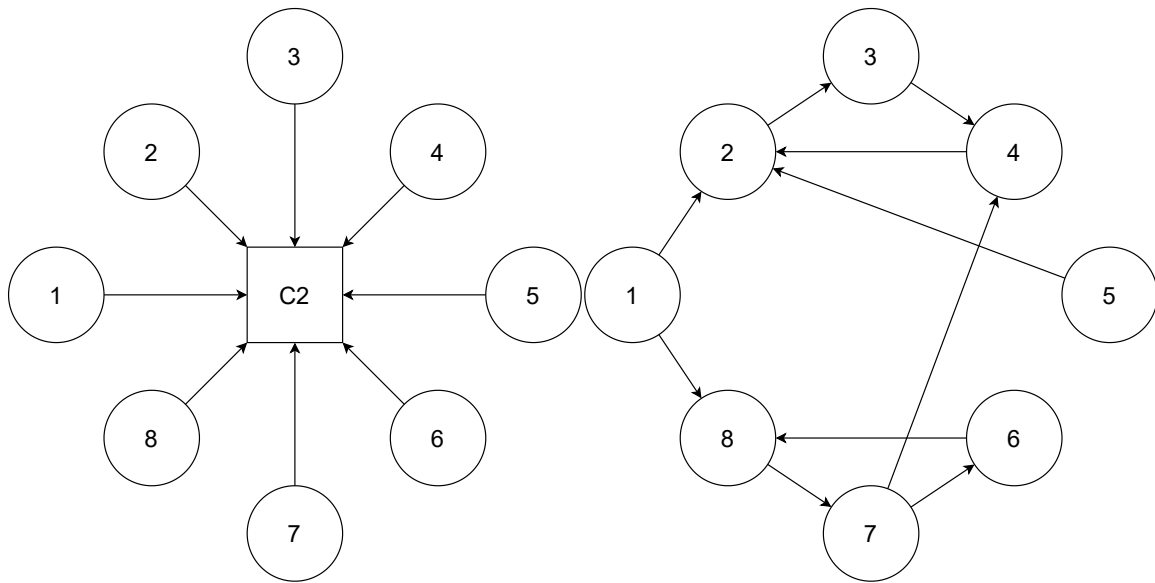
wording

Analyzing and shutting down a centralized or decentralized botnet is comparatively easy since the central means of communication (the C2 IP addresses or domain names, Twitter handles or IRC channels), can be extracted from the malicious binaries or determined by analyzing network traffic and can therefore be considered publicly known.

A coordinated operation with help from law enforcement, hosting providers, domain registrars, and platform providers could shut down or take over the operation by changing how requests are routed or simply shutting down the controlling servers/accounts.

To complicate take-down attempts, botnet operators came up with a number of ideas: domain generation algorithms use pseudorandomly generated domain names to render simple domain blacklist-based approaches ineffective [3] or fast-flux domain name system, where a large pool of IP addresses is assigned randomly to the C2 domains to prevent IP based blacklisting [16].

A number of botnet operations were shut down like this [14] and as the defenders upped their game, so did attackers—the concept of peer-to-peer (P2P) botnets emerged. The idea is to build a distributed network without single points of failure (SPOF) in the form of C2 servers as shown in Figure 1b. In a P2P botnet, each node in the network knows a number of its neighbors and connects to those, each of these neighbors has a list of neighbors on its own, and so on. The botmaster only needs to join the network to send new commands or receive stolen data. Any of the nodes in Figure 1b could be the botmaster but they don't even have to be online all the time since the peers will stay connected autonomously. In fact there have been arrests of operators of P2P botnets but due to the autonomy offered by the distributed approach, the botnet keeps communicating [23]. Especially worm-like



(a) Topology of a C2 controlled botnet (b) Topology of a peer-to-peer (P2P) botnet

Figure 1: Communication paths in different types of botnets

botnets, where each peer tries to find and infect other systems, can keep lingering for many years.

This lack of a SPOF makes P2P botnets more resilient to take-down attempts since the communication is not stopped and botmasters can easily rejoin the network and send commands.

Bots in a P2P botnet can be split into two distinct groups according to their reachability: peers that are not publicly reachable (e.g. because they are behind a Network Access Translation (NAT) router or firewall) and those, that are publicly reachable, also known as *superpeers*. In contrast to centralized botnets with a fixed set of C2 servers, in a P2P botnet, every superpeer might take the roll of a C2 server and *non-superpeers* will connect to those superpeers when joining the network.

As there is no well known server in a P2P botnet, they have to coordinate autonomously. This is achieved by connecting the bots among each other. Bot *B* is considered a *neighbor* of bot *A*, if *A* knows and connects to *B*. Since bots can become unavailable, they have to

permanently update their neighbor lists to avoid losing their connection into the botnet. This is achieved by periodically querying their neighbor's neighbors. This process is known as *Membership Management (MM)*.

MM comes in two forms: structured and unstructured [5]. Structured P2P botnets often use a Distributed Hash Table and strict rules for a bot's neighbors based on its unique ID. In unstructured botnets on the other hand, bots ask any peer they know for new peers to connect to, in a process called *peer discovery*. To enable peers to connect to unstructured botnets, the malware binaries include hardcoded lists of superpeers for the newly infected systems to connect to.

explain
struc-
tured
and
un-
struc-
tured

The concept of *churn* describes when a bot becomes unavailable. There are two types of churn:

- *IP churn*: A bot becomes unreachable because it got assigned a new IP address. The bot is still available but under another address.
- *Device churn*: The device is actually offline, e.g. because the infection was cleaned, it got shut down or lost its internet connection.

2.1 Formal Model of P2P Botnets

A P2P botnet can be modelled as a digraph

$$G = (V, E)$$

With the set of vertices V describing the peers in the network and the set of edges E describing the communication flow between bots.

G is not required to be a connected graph but might consist of multiple disjoint components [20]. Components consisting of peers, that are infected by the same bot, are considered part of the same graph.

For a bot $v \in V$, the *predecessors* (neighbors) $\text{pred}(v)$ and *successors* $\text{succ}(v)$ are defined as:

$$\text{succ}(v) = \{u \in V \mid (u, v) \in E\}$$

$$\text{pred}(v) = \{u \in V \mid (v, u) \in E\}$$

The set of edges $\text{pred}(v)$ is also called the *peer list* of v . Those are the nodes, a peer will connect to, to request new commands and other peers.

For a vertex $v \in V$, the in and out degree deg^+ and deg^- describe how many bots know v or are known by v respectively.

$$\text{deg}^+(v) = |\text{pred}(v)|$$

$$\text{deg}^-(v) = |\text{succ}(v)|$$

more
de-
tails

2.2 Monitoring Techniques for P2P Botnets

There are two distinct methods to map and get an overview of the network topology of a P2P botnet:

2.2.1 Passive Monitoring

For passive detection, traffic flows are analysed in large amounts of collected network traffic (e.g. from internet service providers). This has some advantages in that it is not possible for botmasters to detect or prevent data collection of that kind, but it is not trivial to distinguish valid P2P application traffic (e.g. BitTorrent, Skype, cryptocurrencies, ...) from P2P bots. Zhang et al. propose a system of statistical analysis to solve some of these

problems in [24]. Also getting access to the required datasets might not be possible for everyone.

As most botnet detection mechanisms, also the passive ones work by building communication graphs and finding tightly coupled subgraphs that might be indicative of a botnet [15]. An advantage of passive detection is, that it is independent of protocol details, specific binaries or the structure of the network (P2P vs. centralized/decentralized) [11].

- Large scale network analysis (hard to differentiate from legitimate P2P traffic (e.g. BitTorrent), hard to get data, knowledge of some known bots required) [24]
- Heuristics: Same traffic patterns, same malicious behaviour

Passive monitoring is only mentioned for completeness and not a topic for this thesis.

no
con-
text

2.2.2 Active Monitoring

For active detection, a subset of the botnet protocol and behavior is reimplemented to take part in the network. To do so, samples of the malware are reverse engineered to understand and recreate the protocol. This partial implementation includes the communication part of the botnet but ignores the malicious functionality as to not support and take part in illicit activity.

There are two subtypes of active detection: *sensors* wait to be contacted by other peers, while *crawlers* actively query known bots and recursively ask for their neighbors [12]. Obviously crawlers can only detect superpeers and therefore only see a small subset of the network, while sensors are also contacted by peers in private networks and behind firewalls. To accurately monitor a P2P botnet, a hybrid approach of crawlers and sensors is required.

A crawler starts with a predefined list of known bots, connects to those and uses the peer exchange mechanism to request other bots. Each found bot is crawled again, slowly building the graph of superpeers on the way. Every entry E in the peer exchange response received from bot A represents an edge from A to E in the graph.

A sensor implements the passive part of the botnet's MM. They cannot be used to create the botnet graph (only edges into the sensor node) but are the only way to enumerate the whole network.

2.2.3 Monitoring Prevention Techniques

good
title

The constantly growing damage produced by botnets has many researchers and law enforcement agencies trying to shut down these operations [14, 13, 8, 7]. The monetary value of these botnets directly correlates with the amount of effort botmasters are willing to put into implementing defense mechanisms against take-down attempts.

Some of these countermeasures are explored by Andriess, Rossow, and Bos in “Reliable Recon in Adversarial Peer-to-Peer Botnets” and include deterrence, which limits the number of allowed bots per IP address or subnet to 1; blacklisting, where known crawlers and sensors are blocked from communicating with other bots in the network (mostly IP based); disinformation, when fake bots are placed in the peer lists, which invalidates the data collected by crawlers; and active retaliation like DDoS attacks against sensors or crawlers [1].

Successful take-downs of a P2P botnet requires intricate knowledge over the network topology, protocol characteristics and participating peers. In this work we try to find ways to make the monitoring and information gathering phase more efficient and resilient to detection.

3 Methodology

The implementation of the concepts of this work will be done as part of Botnet Monitoring System (BMS)¹, a monitoring platform for P2P botnets described by Böck et al. in “Challenges of Accurately Measuring Churn in P2P Botnets”. BMS is intended for a hybrid active approach of crawlers and sensors (reimplementations of the P2P protocol of a botnet, that won’t perform malicious actions) to collect live data from active botnets.

In an earlier project, we implemented different node ranking algorithms (among others “PageRank” [18]) to detect sensor candidates in a botnet, as described in “SensorBuster”. Both ranking algorithms exploit the differences in deg^+ and deg^- for sensors to weight the nodes. Another way to enumerate candidates for sensors in a P2P botnet is to find Weakly Connected Components (WCCs) in the graph. Sensors will have few to none outgoing edges, since they don’t participate actively in the botnet, while crawlers have only outgoing edges.

The goal of this work is to complicate detection mechanisms like this for botmasters by centralizing the coordination of the system’s crawlers and sensors, thereby reducing the node’s rank for specific graph metrics. The coordinated work distribution also helps in efficiently monitoring large botnets where one crawler is not enough to track all peers. The changes should allow the current crawlers and sensors to use the new abstraction with as few changes as possible to the existing code.

The goal of this work is to show how cooperative monitoring of a P2P botnet can help with the following problems:

- Impede detection of monitoring attempts by reducing the impact of aforementioned graph metrics
- Circumvent anti-monitoring techniques
- Make crawling more efficient

¹<https://github.com/Telecooperation/BMS>

The final results should be as general as possible and not depend on any botnet's specific behaviour (except for the mentioned anti-monitoring techniques which might be unique to some botnets), but we assume, that every P2P botnet has some way of determining a bot's neighbors.

In the current implementation, each crawler will itself visit and monitor each new node it finds. The general idea for the implementation of the ideas in this thesis is to report newfound nodes back to the BMS backend first, where the graph of the known network is created, and a fitting worker is selected to achieve the goal of the according coordination strategy. That worker will be responsible to monitor the new node.

If it is not possible, to select a specific sensor so that the monitoring activity stays inconspicuous, the coordinator can do a complete shuffle of all nodes between the sensors to restore the wanted graph properties or warn if more sensors are required to stay undetected.

The improved crawler system should allow new crawlers to register themselves and their capabilities (e.g. bandwidth, geolocation), so the amount of work can be scaled accordingly between hosts.

3.1 Protocol Primitives

The coordination protocol must allow the following operations:

Register Worker Register a new worker with capabilities (which botnet, available bandwidth and processing power, ...). This is called periodically and used to determine which worker is still active, when assigning new tasks.

Report Peer Report found peers. Both successful and failed attempts are reported, to detect churned peers, and blacklisted crawlers as soon as possible.

Report Edge Report found edges. Edges are created by querying the peer list of a bot. This is how new peers are detected.

Request Tasks Receive a batch of crawl tasks from the coordinator. The tasks consist of the target peer, if the worker should start or stop monitoring the peer, when the monitoring should start and stop and at which frequency the peer should be contacted.

Request Neighbors Sensors can request a list of candidate peers to return when their peer list is queried.

```
type Peer struct {
    BotID string
    IP     string
    Port  uint16
}
type PeerTask struct {
    Peer          Peer
    StartAt      *Time
    StopAt       *Time
    Frequency    uint
    StopCrawling bool
}
```

Listing 1: Relevant Fields for Peers and Tasks

Listing 1 shows the Go structures used for crawl tasks.

4 Coordination Strategies

Let C be the set of available crawlers. Without loss of generality, if not stated otherwise, we assume that C is known when BMS is started and will not change afterward. There will be no joining or leaving crawlers. This assumption greatly simplifies the implementation due to the lack of changing state that has to be tracked while still exploring the described strategies. A production-ready implementation of the described techniques can drop this assumption but might have to recalculate the work distribution once a crawler joins or leaves. The protocol primitives described in Section 3.1 already allow for this to be implemented by first creating tasks with the `StopCrawling` flag set to true for all active tasks, run the strategy again and create the according tasks to start crawling again.

4.1 Load Balancing

Depending on a botnet's size, a single crawler is not enough to monitor all superpeers. While it is possible to run multiple, uncoordinated crawlers, multiple crawlers can find and monitor the same peer, making the approach inefficient with regard to the computing resources at hand.

The load balancing strategy solves this problem by systematically splitting the crawl tasks into chunks and distributes them among the available crawlers. The following load balancing strategies will be investigated:

- Round Robin. See Section 4.1.1
- Assuming IP addresses are evenly distributed and so are infections, take the IP address as an 32 Bit integer modulo $|C|$. See Section 4.1.2 Problem: reassignment if a crawler joins or leaves

Load balancing in itself does not help prevent the detection of crawlers but it allows better usage of available resources. It prevents unintentionally crawling the same peer with multiple crawlers and allows crawling of bigger botnets where the uncoordinated approach

remove?

would reach its limit and could only be worked around by scaling up the machine where the crawler is executed. Load balancing allows scaling out, which can be more cost-effective.

4.1.1 Round Robin Distribution

This strategy distributes work evenly among crawlers by either naively assigning tasks to the crawlers rotationally or weighted according to their capabilities. To keep the distribution as even as possible, we keep track of the last crawler a task was assigned to and start with the next in line in the subsequent round of assignments. For the sake of simplicity, only the bandwidth will be considered as capability but it can be extended by any shared property between the crawlers, e.g. available memory or processing power. For a given crawler $c_i \in C$ let $cap(c_i)$ be the capability of the crawler. The total available capability is $B = \sum_{c \in C} cap(c)$. With G being the greatest common divisor of all the crawler's capabilities, the weight $W(c_i) = \frac{cap(c_i)}{G} \cdot \frac{cap(c_i)}{B}$ gives us the percentage of the work a crawler is assigned. The algorithm in Listing 2 distributes the work according to the crawler's capabilities.

This creates a list of crawlers where a crawler can occur more than once, depending on its capabilities. To ensure better distribution, first every crawler is assigned one task, then, according to the capabilities, every crawler with a weight of 2 or more is assigned a task, and so on. The set of crawlers $\{a, b, c\}$ with the capabilities $cap(a) = 3$, $cap(b) = 2$, $cap(c) = 1$ would produce $\langle a, b, c, a, b, a \rangle$, allocating two and three times the work to crawlers b and a respectively.

4.1.2 IP-based Partitioning

The output of cryptographic hash functions is uniformly distributed—even substrings of the calculated hash hold this property. Calculating the hash of an IP address and distributing the work with regard to $H(IP) \bmod |C|$ creates about evenly sized buckets for each worker to handle. For any hash function H , this gives us the mapping $m(i) = H(i) \bmod |C|$ to sort peers into buckets.

1 – 2 sentences about naive rr?

page numbers for forward refs?

better wording

don't use substrings, bit.int for 128 bit modulo, argumentation why

```
func WeightCrawlers(crawlers ...Crawler) map[string]uint {
    weights := []int{}
    totalWeight := 0
    for _, crawler := range crawlers {
        totalWeight += crawler.Bandwith
        weights = append(weights, crawler.Bandwith)
    }
    gcd := Fold(Gcd, weights...)
    weightMap := map[string]uint{}
    for _, crawler := range crawlers {
        weightMap[crawler.ID] = uint(crawler.Bandwith / gcd)
    }
    return weightMap
}

func WeightedCrawlerList(crawlers ...Crawler) []string {
    weightMap := WeightCrawlers(crawlers...)
    didSomething := true
    crawlerIds := []string{}
    for didSomething {
        didSomething = false
        for k, v := range weightMap {
            if v != 0 {
                didSomething = true
                crawlerIds = append(crawlerIds, k)
                weightMap[k] -= 1
            }
        }
    }
    return crawlerIds
}
```

Listing 2: Pseudocode for weighted round robin

4 Coordination Strategies

Any hash function can be used but since it must be calculated often, a fast function should be used. While the Message-Digest Algorithm 5 (MD5) hash function must be considered broken for cryptographic use [21], it is faster to calculate than hash functions with longer output. For the use case at hand, only the uniform distribution property is required so MD5 can be used without sacrificing any kind of security.

This strategy can also be weighted using the crawlers capabilities by modifying the list of available workers so that a worker can appear multiple times according to its weight. The weighting algorithm from Listing 2 is used to create the weighted multiset of crawlers C_W and the mapping changes to $m(i) = H(i) \bmod |C_W|$.

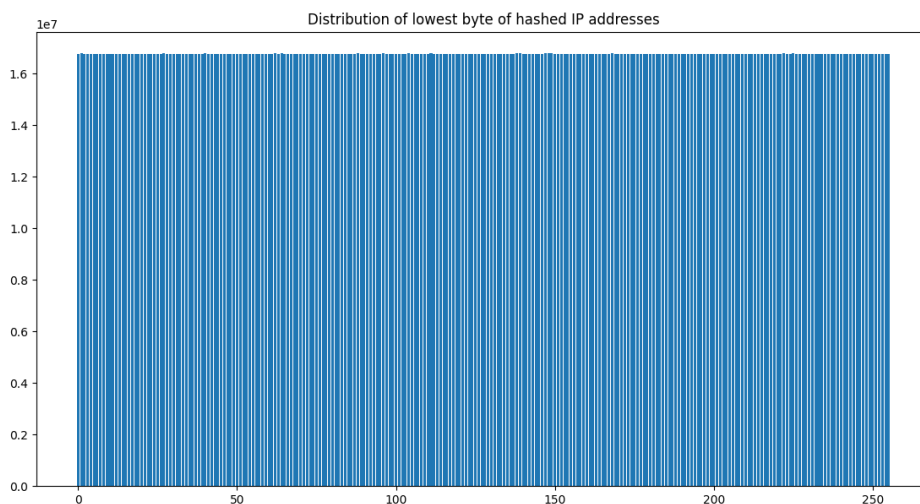


Figure 2: Distribution of the lowest byte of MD5 hashes over IPv4

MD5 returns a 128 Bit hash value. The Go standard library includes helpers for arbitrarily sized integers². This helps us in implementing the mapping m from above.

By exploiting the even distribution offered by hashing, the work of each crawler is also evenly distributed over all IP subnets, Autonomous System (AS) and geolocations. This ensures neighboring peers (e.g. in the same AS, geolocation or IP subnet) get visited by

²<https://pkg.go.dev/math/big#Int>

remove
this?

different crawlers. It also allows us to get rid of the state in our strategy since we don't have to keep track of the last crawler we assigned a task to, making it easier to implement and reason about.

4.2 Reduction of Request Frequency

The GameOver Zeus botnet deployed a blacklisting mechanism, where crawlers are blocked based in their request frequency [2]. In a single crawler approach, the crawler frequency has to be limited to prevent hitting the request limit.

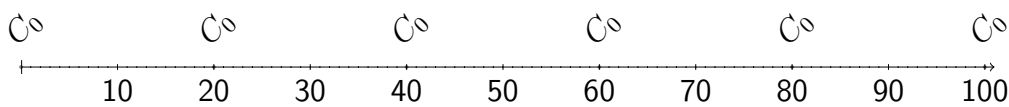


Figure 3: Timeline of crawler events as seen from a peer when crawled by a single crawler

Using collaborative crawlers, an arbitrarily fast frequency can be achieved without being blacklisted. With $L \in \mathbb{N}$ being the frequency limit at which a crawler will be blacklisted, $F \in \mathbb{N}$ being the crawl frequency that should be achieved. The amount of crawlers C required to achieve the frequency F without being blacklisted and the offset O between crawlers are defined as

$$C = \left\lceil \frac{F}{L} \right\rceil$$

$$O = \frac{1 \text{ req}}{F}$$

Taking advantage of the StartAt field from the PeerTask returned by the requestTasks primitive above, the crawlers can be scheduled offset by O at a frequency L to ensure, the overall requests to each peer are evenly distributed over time.

Given a limit $L = 5 \text{ req}/100\text{s}$, crawling a botnet at $F = 20 \text{ req}/100\text{s}$ requires $C =$

better numbers for example?

4 Coordination Strategies

$\lceil \frac{20 \text{ req}/100\text{s}}{5 \text{ req}/100\text{s}} \rceil = 4$ crawlers. Those crawlers must be scheduled $O = \frac{1 \text{ req}}{20 \text{ req}/100\text{s}} = 5 \text{ s}$ apart at a frequency of L for an even request distribution.

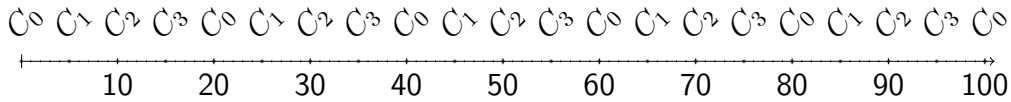
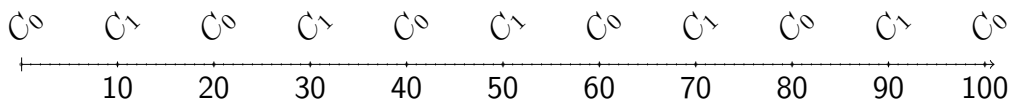


Figure 4: Timeline of crawler events as seen from a peer when crawled by multiple crawlers

As can be seen in Figure 4, each crawler C_0 to C_3 performs only $5 \text{ req}/100\text{s}$ while overall achieving $20 \text{ req}/100\text{s}$.

Vice versa given an amount of crawlers C and a request limit L , the effective frequency F can be maximized to $F = C \times L$ without hitting the limit L and being blocked.

Using the example from above with $L = 5 \text{ req}/100\text{s}$ but now only two crawlers $C = 2$, it is still possible to achieve an effective frequency of $F = 2 \times 5 \text{ req}/100\text{s} = 10 \text{ req}/100\text{s}$ and $O = \frac{1 \text{ req}}{10 \text{ req}/100\text{s}} = 10 \text{ s}$:



While the effective frequency of the whole system is halved compared to Figure 4, it is still possible to double the frequency over the limit.

4.3 Creating Edges for Crawlers and Sensors

“SensorBuster: On Identifying Sensor Nodes in P2P Botnets” describes different graph metrics to find sensors in P2P botnets. These metrics depend on the uneven ratio between incoming and outgoing edges for crawlers. One of those, “SensorBuster” uses WCCs since crawlers don’t have any edges back to the main network in the graph.

4 Coordination Strategies

Building a complete graph $G_C = K_{|C|}$ between the crawlers by making them return the other crawlers on peer list requests would still produce a disconnected component and while being bigger and maybe not as obvious at first glance, it is still easily detectable since there is no path from G_C back to the main network (see Figure 9b and Table 4).

With $v \in V$, $\text{succ}(v)$ being the set of successors of v and $\text{pred}(v)$ being the set of predecessors of v , PageRank is recursively defined as [18]:

$$\begin{aligned} \text{PR}_0(v) &= \text{initialRank} \\ \text{PR}_{n+1}(v) &= \text{dampingFactor} \times \sum_{p \in \text{pred}(v)} \frac{\text{PR}_n(p)}{|\text{succ}(p)|} + \frac{1 - \text{dampingFactor}}{|V|} \end{aligned}$$

For the first iteration, the PageRank of all nodes is set to the same initial value. Page et al. argue that when iterating often enough, any value can be chosen [18].

The dampingFactor describes the probability of a person visiting links on the web to continue doing so, when using PageRank to rank websites in search results. For simplicity—and since it is not required to model human behaviour for automated crawling and ranking—a dampingFactor of 1.0 will be used, which simplifies the formula to

$$\text{PR}_{n+1}(v) = \sum_{p \in \text{pred}(v)} \frac{\text{PR}_n(p)}{|\text{succ}(p)|}$$

Based on this, SensorRank is defined as

$$\text{SR}(v) = \frac{\text{PR}(v)}{|\text{succ}(v)|} \times \frac{|\text{pred}(v)|}{|V|}$$

In our experiments on a snapshot of the Sality [9] botnet obtained from BMS over the span of 21st to 28th April 2021 even 1 iteration were enough to get distinct enough values to detect sensors and crawlers.

4 Coordination Strategies

Iteration	Avg. PR	Crawler PR	Avg. SR	Crawler SR
1	0.24854932	0.63277194	0.15393478	0.56545578
2	0.24854932	0.63277194	0.15393478	0.56545578
3	0.24501068	0.46486353	0.13810930	0.41540997
4	0.24501068	0.46486353	0.13810930	0.41540997
5	0.24233737	0.50602884	0.14101354	0.45219598

Table 1: Values for PageRank iterations with initial rank $\forall v \in V : PR(v) = 0.25$

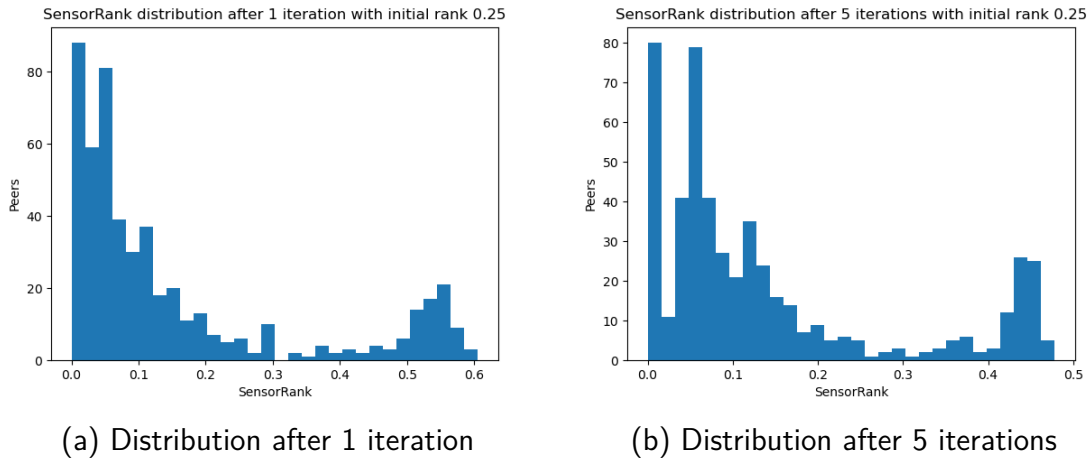


Figure 5: SensorRank distribution with initial rank $\forall v \in V : PR(v) = 0.25$

Iteration	Avg. PR	Crawler PR	Avg. SR	Crawler SR
1	0.49709865	1.26554389	0.30786955	1.13091156
2	0.49709865	1.26554389	0.30786955	1.13091156
3	0.49002136	0.92972707	0.27621861	0.83081993
4	0.49002136	0.92972707	0.27621861	0.83081993
5	0.48467474	1.01205767	0.28202708	0.90439196

Table 2: Values for PageRank iterations with initial rank $\forall v \in V : PR(v) = 0.5$

4 Coordination Strategies

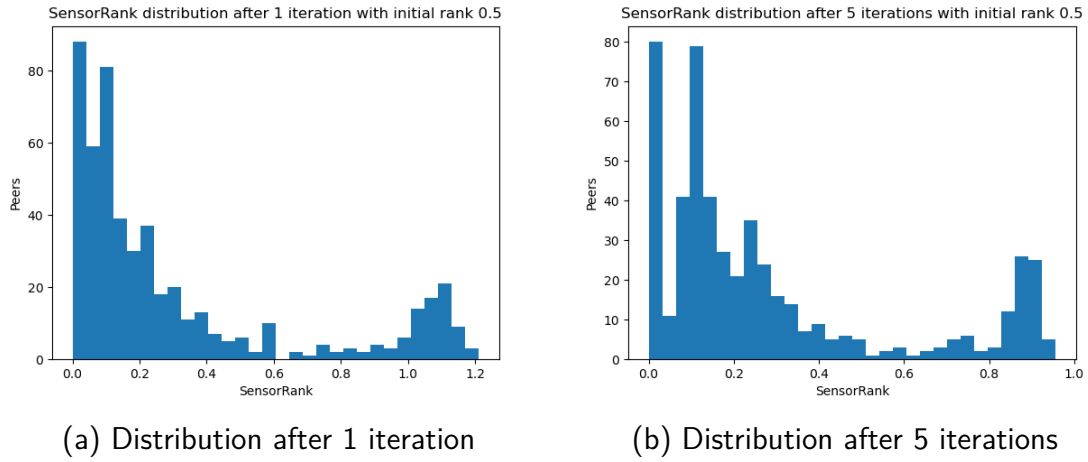


Figure 6: SensorRank distribution with initial rank $\forall v \in V : PR(v) = 0.5$

Iteration	Avg. PR	Crawler PR	Avg. SR	Crawler SR
1	0.74564797	1.89831583	0.46180433	1.69636734
2	0.74564797	1.89831583	0.46180433	1.69636734
3	0.73503203	1.39459060	0.41432791	1.24622990
4	0.73503203	1.39459060	0.41432791	1.24622990
5	0.72701212	1.51808651	0.42304062	1.35658794

Table 3: Values for PageRank iterations with initial rank $\forall v \in V : PR(v) = 0.75$

4 Coordination Strategies

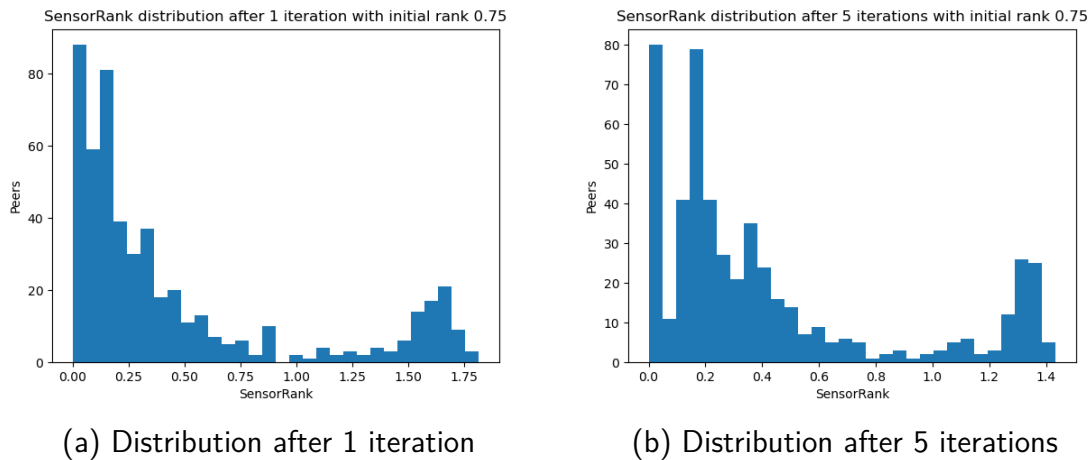


Figure 7: SensorRank distribution with initial rank $\forall v \in V : PR(v) = 0.75$

The distribution graphs in Figure 5, Figure 6 and Figure 7 show that the initial rank has no effect on the distribution, only on the actual numeric rank values.

For all combinations of initial value and PageRank iterations, the rank for a well known crawler is in the 95th percentile, so for our use case, those parameters do not matter.

On average, peers in the analyzed dataset have 223 successors over the whole week. Looking at the data in smaller buckets of one hour each, the average number of successors per peer is 90.

Since crawlers never respond to peer list requests, they will always be detectable by the described approach but sensors might benefit from the following technique.

By responding to peer list requests with plausible data, one can make those metrics less suspicious, because it produces valid outgoing edges from the sensors. The hard part is deciding which peers can be returned without actually supporting the network. The following candidates to place into the neighbor list will be investigated:

- Return the other known sensors, effectively building an complete graph $K_{|C|}$ containing all sensors

timeline
with
peers
per
bucket
use
better
data?

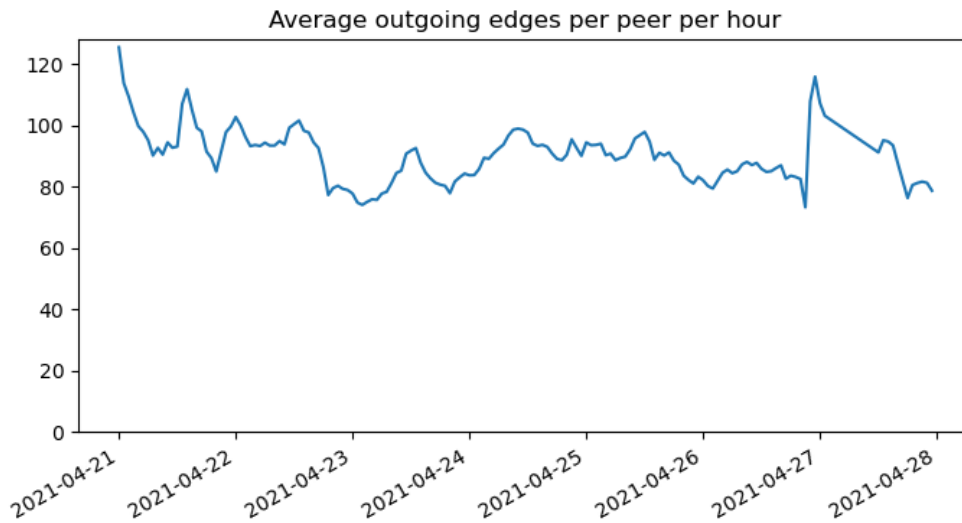


Figure 8: Average outgoing edges per peer per hour

- Detect churned peers from AS with dynamic IP allocation
- Detect peers behind carrier-grade NAT that rotate IP addresses very often and pick random IP addresses from the IP range

Knowledge of only 90 peers leaving due to IP rotation would be enough to make a crawler look average in Sality. This number will differ between different botnets, depending on implementation details and size of the network.

4.3.1 Use Other Known Sensors

By connecting the known sensors and effectively building a complete graph $K_{|C|}$ between them creates $|C| - 1$ outgoing edges per sensor. In most cases this won't be enough to reach the amount of edges that would be needed. Also this does not help against the WCC metric since this would create a bigger but still disconnected component.

repeat
anal-
ysis,
actual
num-
ber

upper
limit
for
NL
size
as
impl
detail

4 Coordination Strategies

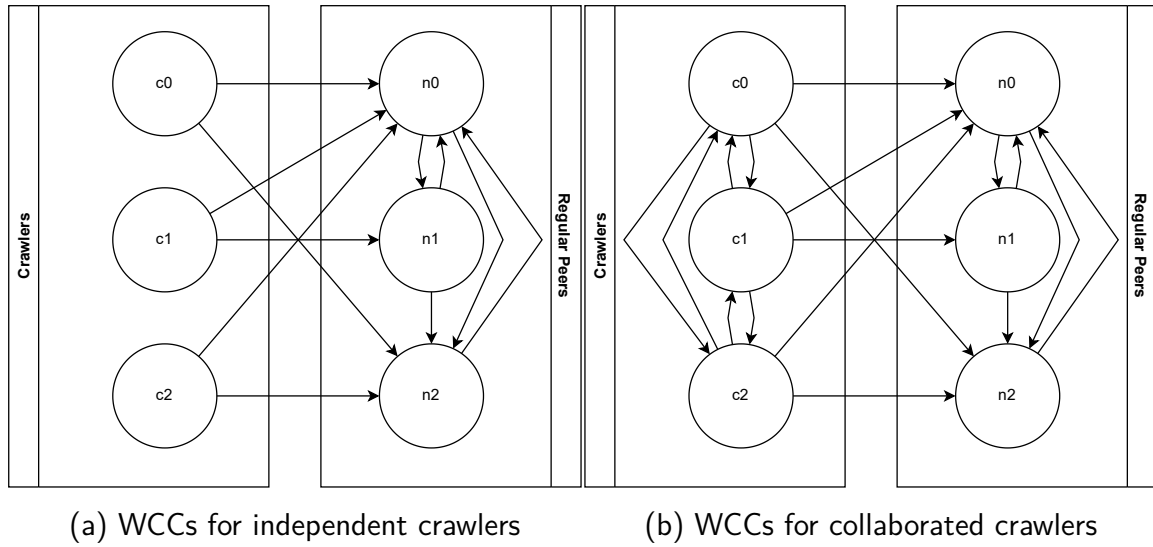


Figure 9: Differences in graph metrics

Applying PageRank once with an initial rank of 0.25 once on the example graphs in Figure 9 results in:

Node	deg ⁺	deg ⁻	In WCC?	PageRank	SensorRank
n0	0/0	4/4	no	0.75/0.5625	0.3125/0.2344
n1	1/1	3/3	no	0.25/0.1875	0.0417/0.0313
n2	2/2	2/2	no	0.5/0.375	0.3333/0.25
c0	3/5	0/2	yes (1/3)	0.0/0.125	0.0/0.0104
c1	1/3	0/2	yes (1/3)	0.0/0.125	0.0/0.0104
c2	2/4	0/2	yes (1/3)	0.0/0.125	0.0/0.0104

Table 4: Values for metrics from Figure 9 (a/b)

While this works for small networks, the crawlers must account for a significant amount of peers in the network for this change to be noticeable. The generated K_n needs to be at least as big as the smallest regular component in the botnet, which is not feasible.

4.3.2 Use Churned Peers After IP Rotation

Churn describes the dynamics of peer participation of P2P systems, e.g. join and leave events [22]. Detecting if a peer just left the system, in combination with knowledge about ASs, peers that just left and came from an AS with dynamic IP allocation (e.g. many consumer broadband providers in the US and Europe), can be placed into the crawler's peer list. If the timing of the churn event correlates with IP rotation in the AS, it can be assumed, that the peer left due to being assigned a new IP address—not due to connectivity issues or going offline—and will not return using the same IP address. These peers, when placed in the peer list of the crawlers, will introduce paths back into the main network and defeat the WCC metric. It also helps with the PageRank and SensorRank metrics since the crawlers start to look like regular peers without actually supporting the network by relaying messages or propagating active peers.

übergang

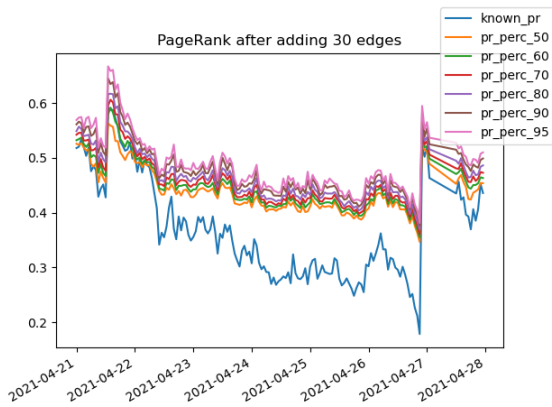
what
is an
AS

4.3.3 Peers Behind Carrier-Grade NAT

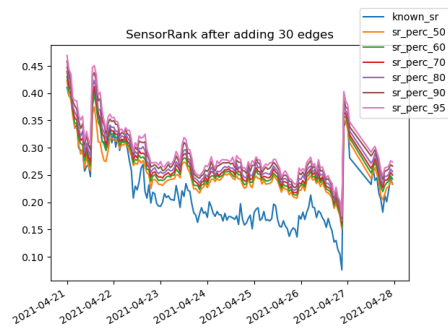
Some peers show behaviour, where their IP address changes almost after every request. Those peers can be used as fake neighbours and create valid looking outgoing edges for the sensor.

Experiments were performed, in which a fixed amount of random outgoing edges were added to the known sensor and the data was plotted:

4 Coordination Strategies

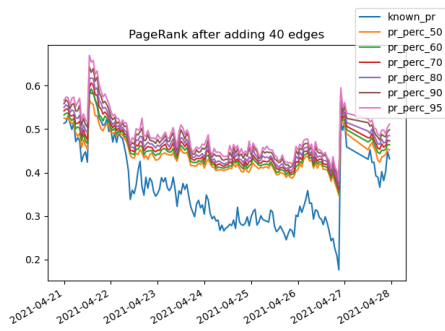


(a) PageRank with 30 additional edges

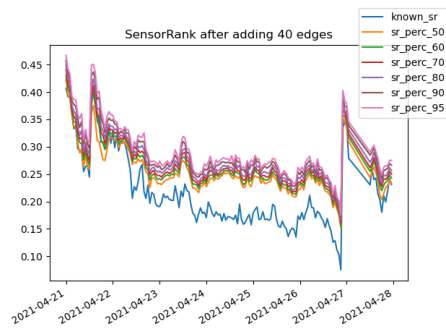


(b) SensorRank with 30 additional edges

Figure 10: Ranking with 30 additional edges



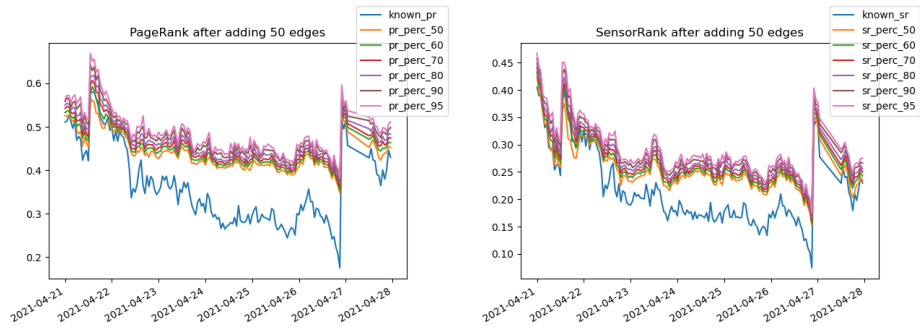
(a) PageRank with 40 additional edges



(b) SensorRank with 40 additional edges

Figure 11: Ranking with 40 additional edges

4 Coordination Strategies



(a) PageRank with 50 additional edges (b) SensorRank with 50 additional edges

Figure 12: Ranking with 50 additional edges

5 Implementation

Crawlers in BMS report to the backend using gRPC remote procedure calls (gRPCs)³. Both crawlers and the backend gRPC server are implemented using the Go⁴ programming language, so to make use of existing know-how and to allow others to use the implementation in the future, the coordinator backend and crawler abstraction were also implemented in Go.

BMS already has an existing abstraction for crawlers. This implementation is highly optimized but also tightly coupled and grown over time. The abstraction became leaky and extending it proved to be complicated. A new crawler abstraction was created with testability, extensibility and most features of the existing implementation in mind, which can be ported back to be used by the existing crawlers.

The new implementation consists of three main interfaces:

- `FindPeer`, to receive new crawl tasks from any source
- `ReportPeer`, to report newly found peers
- `Protocol`, the actual botnet protocol implementation used to ping a peer and request its peer list

Currently there are two sources `FindPeer` can use: read peers from a file on disk or request them from the gRPC BMS coordinator. The `ExactlyOnceFinder` delegate can wrap another `FindPeer` instance and ensures the source is only requested once. This is used to implement the bootstrapping mechanism of the old crawler, where once, when the crawler is started, the list of bootstrap nodes is loaded from a textfile. `CombinedFinder` can combine any amount of `FindPeer` instances and will return the sum of requesting all the sources.

The `PeerTask` instances returned by `FindPeer` contain the IP address and port of the peer, if the crawler should start or stop the operation, when to start and stop crawling and

³<https://www.grpc.io>

⁴<https://go.dev/>

5 Implementation

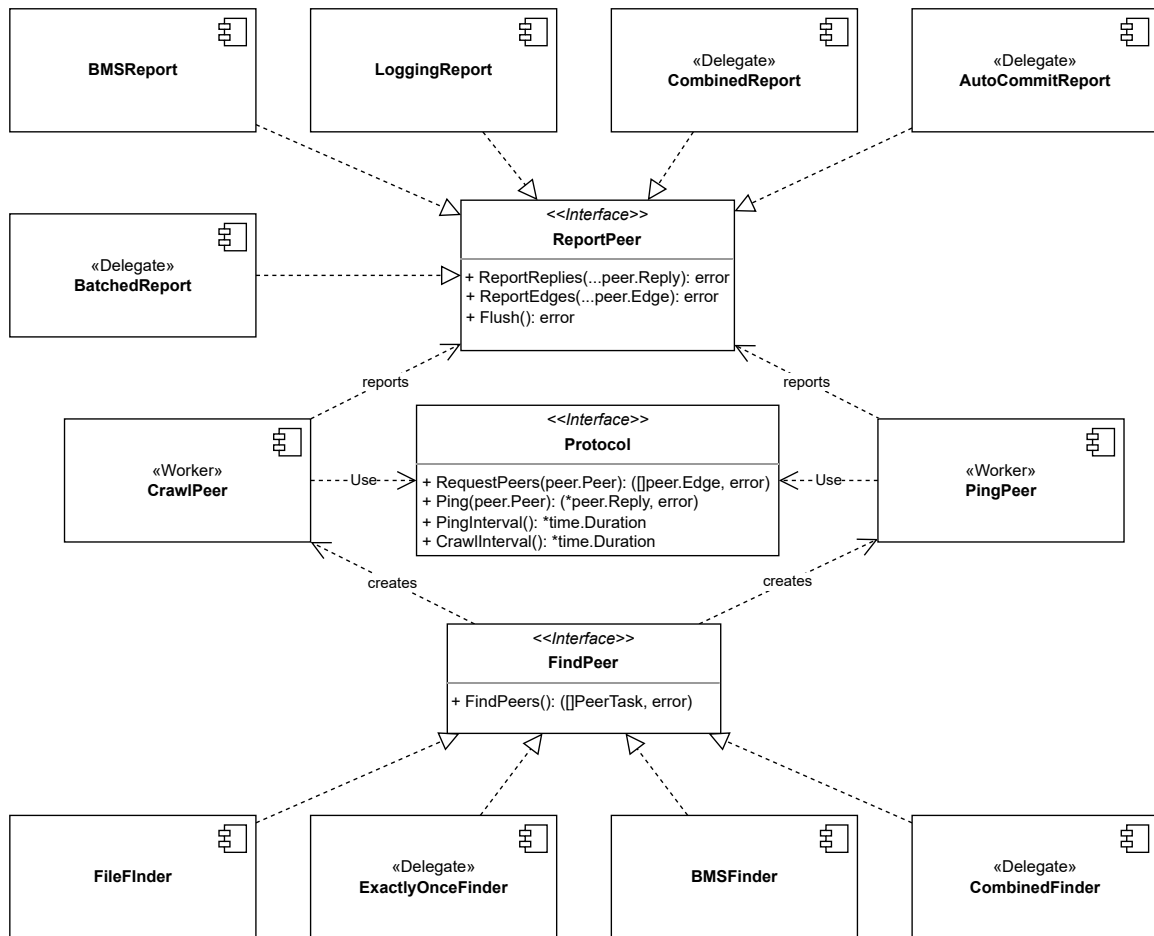


Figure 13: Architecture of the new crawler

in which interval the peer should be crawled. For each task, a `CrawlPeer` and `PingPeer` worker is started or stopped as specified in the received `PeerTask`. These tasks use the `ReportPeer` interface to report any new peer that is found.

Current report possibilities are `LoggingReport` to simply log new peers to get feedback from the crawler at runtime, and `BMSReport` which reports back to BMS. `BatchedReport` delegates a `ReportPeer` instance and batch newly found peers up to a specified batch size and only then flush and actually report. `AutoCommitReport` will automatically flush a delegated `ReportPeer` instance after a fixed amount of time and is used in combination with `BatchedReport` to ensure the batches are written regularly, even if the batch limit

5 Implementation

is not reached yet. CombinedReport works analogous to CombinedFinder and combines many ReportPeer instances into one.

PingPeer and CrawlPeer use the implementation of the botnet Protocol to perform the actual crawling in predefined intervals, which can be overwritten on a per PeerTask basis.

The server-side part of the system consists of a gRPC server to handle the client requests, a scheduler to assign new peers, and a Strategy interface for modularity over how work is assigned to crawlers.

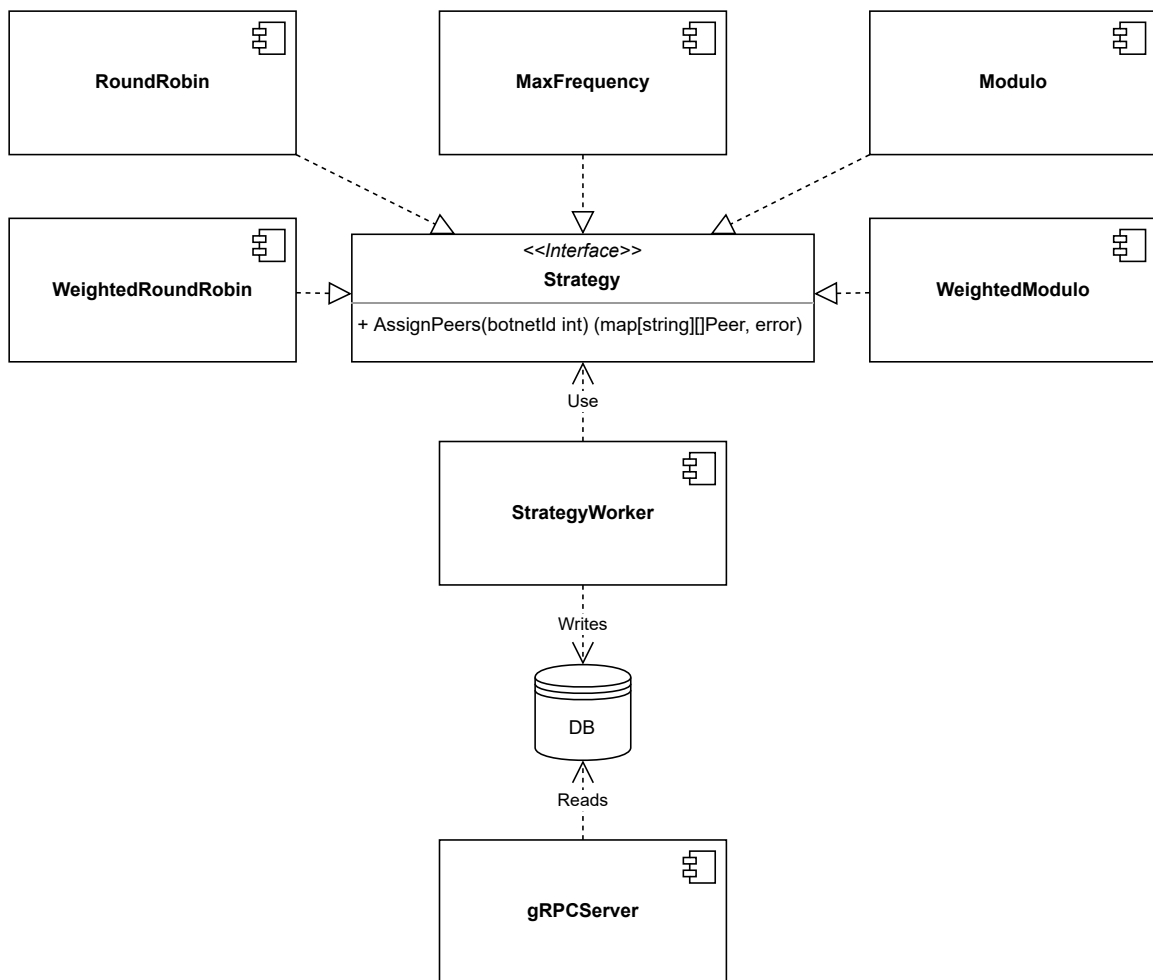


Figure 14: Architecture of the gRPC backend

6 Conclusion, Lessons Learned

decide

Collaborative monitoring of P2P botnets allows circumventing some anti-monitoring efforts. It also enables more effective monitoring systems for larger botnets, since each peer can be visited by only one crawler. The current concept of independent crawlers in BMS can also use multiple workers but there is no way to ensure a peer is not watched by multiple crawlers thereby using unnecessary resources.

7 Further Work

Following this work, it should be possible to rewrite the existing crawlers using the new abstraction. This might bring some performance issues to light which can be solved by investigating the optimizations from the old implementation and applying them to the new one.

Another way to expand on this work is automatically scaling the available crawlers up and down, depending on the botnet size and the number of concurrently online peers. Doing so would allow a constant crawl interval for even highly volatile botnets.

Placing churned peers or peers with suspicious network activity (those behind carrier-grade NATs) might just offer another characteristic to flag sensors in a botnet. This should be investigated and maybe there are ways to mitigate this problem.

Autoscaling features offered by many cloud-computing providers should be evaluated to automatically add or remove crawlers based on the monitoring load, a botnet's size and number of active peers. This should also allow create workers with new IP addresses in different geolocations fast and easy.

Acknowledgments

In the end, I would like to thank

- Prof. Dr. Christoph Skornia for being a helpful supervisor in this and many earlier works of mine
- Leon Böck for offering the possibility to work on this research project, regular feedback and technical expertise
- Valentin Sundermann for being available for insightful ad hoc discussions at any time of day for many years
- Friends and family who pushed me into continuing this path

List of Figures

1	Communication paths in different types of botnets	8
2	Distribution of the lowest byte of MD5 hashes over IPv4	19
3	Timeline of crawler events as seen from a peer when crawled by a single crawler	20
4	Timeline of crawler events as seen from a peer when crawled by multiple crawlers	21
5	SensorRank distribution with initial rank $\forall v \in V : PR(v) = 0.25$	23
6	SensorRank distribution with initial rank $\forall v \in V : PR(v) = 0.5$	24
7	SensorRank distribution with initial rank $\forall v \in V : PR(v) = 0.75$	25
8	Average outgoing edges per peer per hour	26
9	Differences in graph metrics	27
10	Ranking with 30 additional edges	29
11	Ranking with 40 additional edges	29
12	Ranking with 50 additional edges	30
13	Architecture of the new crawler	32
14	Architecture of the gRPC backend	33

List of Tables

1	Values for PageRank iterations with initial rank $\forall v \in V : PR(v) = 0.25$. . .	23
2	Values for PageRank iterations with initial rank $\forall v \in V : PR(v) = 0.5$. . .	23
3	Values for PageRank iterations with initial rank $\forall v \in V : PR(v) = 0.75$. . .	24
4	Values for metrics from Figure 9 (a/b)	27

List of Listings

1	Relevant Fields for Peers and Tasks	15
2	Pseudocode for weighted round robin	18

List of Acronyms

AS Autonomous System	19, 25, 27
BMS Botnet Monitoring System	13 f., 16, 22, 31 f., 34
C2 command and control	7 f.
DDoS distributed denial of service	6 f., 12
gRPC gRPC remote procedure call	31, 33
IoT internet of things	6
IRC internet relay chat	7
MD5 Message-Digest Algorithm 5	17, 19
MM Membership Management	8 f., 11
NAT Network Access Translation	8, 25, 35
P2P peer-to-peer	7–13, 21, 27, 34
SPOF single point of failure	7 f.
WCC Weakly Connected Component	13, 21, 26 f.

References

- [1] Dennis Andriessse, Christian Rossow, and Herbert Bos. “Reliable Recon in Adversarial Peer-to-Peer Botnets”. In: *Proceedings of the 2015 Internet Measurement Conference*. IMC '15: Internet Measurement Conference. Tokyo Japan: ACM, Oct. 28, 2015, pp. 129–140. ISBN: 978-1-4503-3848-6. DOI: 10.1145/2815675.2815682. URL: <https://dl.acm.org/doi/10.1145/2815675.2815682> (visited on 11/16/2021).
- [2] Dennis Andriessse et al. “Highly Resilient Peer-to-Peer Botnets Are Here: An Analysis of Gameover Zeus”. In: *2013 8th International Conference on Malicious and Unwanted Software: "The Americas" (MALWARE)*. 2013 8th International Conference on Malicious and Unwanted Software: "The Americas" (MALWARE). Fajardo, PR, USA: IEEE, Oct. 2013, pp. 116–123. ISBN: 978-1-4799-2534-6 978-1-4799-2535-3. DOI: 10.1109/MALWARE.2013.6703693. URL: <https://ieeexplore.ieee.org/document/6703693/> (visited on 02/27/2022).
- [3] Manos Antonakakis et al. “From Throw-Away Traffic to Bots: Detecting the Rise of DGA-Based Malware”. In: *21st USENIX Security Symposium (USENIX Security 12)*. Bellevue, WA: USENIX Association, Aug. 2012, pp. 491–506. ISBN: 978-931971-95-9. URL: <https://www.usenix.org/conference/usenixsecurity12/technical-sessions/presentation/antonakakis>.
- [4] *Availability of broadband internet to households in Germany from 2017 to 2020, by bandwidth class*. Statista Inc. Aug. 16, 2021. URL: <https://www.statista.com/statistics/460180/broadband-availability-by-bandwidth-class-germany/> (visited on 11/11/2021), archived at <https://web.archive.org/web/20210309010747/https://www.statista.com/statistics/460180/broadband-availability-by-bandwidth-class-germany/> on Mar. 9, 2021.

References

- [5] Leon Böck et al. “Next Generation P2P Botnets: Monitoring Under Adverse Conditions”. In: *Research in Attacks, Intrusions, and Defenses*. Ed. by Michael Bailey et al. Vol. 11050. Series Title: Lecture Notes in Computer Science. Cham: Springer International Publishing, 2018, pp. 511–531. ISBN: 978-3-030-00469-9 978-3-030-00470-5. DOI: 10.1007/978-3-030-00470-5_24. URL: http://link.springer.com/10.1007/978-3-030-00470-5%5C_24 (visited on 04/08/2022).
- [6] Leon Böck et al. “Poster: Challenges of Accurately Measuring Churn in P2P Botnets”. In: *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. CCS '19: 2019 ACM SIGSAC Conference on Computer and Communications Security. London United Kingdom: ACM, Nov. 6, 2019, pp. 2661–2663. ISBN: 978-1-4503-6747-9. DOI: 10.1145/3319535.3363281. URL: <https://dl.acm.org/doi/10.1145/3319535.3363281> (visited on 11/12/2021).
- [7] Joseph Demarest. *Taking Down Botnets*. Federal Bureau of Investigation. July 15, 2014. URL: <https://www.fbi.gov/news/testimony/taking-down-botnets> (visited on 03/23/2022), archived at <https://web.archive.org/web/20220318082034/https://www.fbi.gov/news/testimony/taking-down-botnets>.
- [8] David Dittrich. “So You Want to Take over a Botnet”. In: *Proceedings of the 5th USENIX Conference on Large-Scale Exploits and Emergent Threats*. LEET'12. San Jose, CA: USENIX Association, 2012, p. 6. DOI: 10.5555/2228340.2228349.
- [9] Falliere, Nicolas. *Sality: Story of a Peer-to-Peer Viral Network*. July 2011. URL: <https://papers.vx-underground.org/archive/Symantec/sality-story-of-peer-to-peer-11-en.pdf> (visited on 03/16/2022), archived at https://web.archive.org/web/20161223003320/http://www.symantec.com/content/en/us/enterprise/media/security_

References

- response/whitepapers/sality_peer_to_peer_viral_network.pdf on Dec. 23, 2016.
- [10] Dan Goodin. *Brace yourselves — source code powering potent IoT DDoSes just went public*. Ars Technica. Oct. 2, 2016. URL: <https://arstechnica.com/information-technology/2016/10/brace-yourselves-source-code-powering-potent-iot-ddoses-just-went-public/> (visited on 11/11/2021), archived at <https://web.archive.org/web/20211022032617/https://arstechnica.com/information-technology/2016/10/brace-yourselves-source-code-powering-potent-iot-ddoses-just-went-public/> on Oct. 22, 2021.
- [11] Guofei Gu et al. “BotMiner: Clustering Analysis of Network Traffic for Protocol- and Structure-Independent Botnet Detection”. In: *Proceedings of the 17th Conference on Security Symposium*. SS’08. San Jose, CA: USENIX Association, 2008, pp. 139–154.
- [12] Shankar Karuppayah et al. “SensorBuster: On Identifying Sensor Nodes in P2P Botnets”. In: *Proceedings of the 12th International Conference on Availability, Reliability and Security*. ARES ’17. New York, NY, USA: Association for Computing Machinery, Aug. 29, 2017, pp. 1–6. ISBN: 978-1-4503-5257-4. DOI: 10.1145/3098954.3098991. URL: <https://doi.org/10.1145/3098954.3098991> (visited on 03/23/2021).
- [13] Yacin Nadji, Roberto Perdisci, and Manos Antonakakis. “Still Beheading Hydras: Botnet Takedowns Then and Now”. In: *IEEE Transactions on Dependable and Secure Computing* 14.5 (Sept. 1, 2017), pp. 535–549. ISSN: 1545-5971. DOI: 10.1109/TDSC.2015.2496176. URL: <http://ieeexplore.ieee.org/document/7312442/> (visited on 03/17/2022).
- [14] Yacin Nadji et al. “Beheading hydras: performing effective botnet takedowns”. In: *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security - CCS ’13*. the 2013 ACM SIGSAC conference.

References

- Berlin, Germany: ACM Press, 2013, pp. 121–132. ISBN: 978-1-4503-2477-9. DOI: 10.1145/2508859.2516749. URL: <http://dl.acm.org/citation.cfm?doid=2508859.2516749> (visited on 03/15/2022).
- [15] Shishir Nagaraja et al. “BotGrep: Finding P2P Bots with Structured Graph Analysis”. In: *Proceedings of the 19th USENIX Conference on Security*. USENIX Security’10. Washington, DC: USENIX Association, 2010, p. 7. ISBN: 8887666655554.
- [16] Jose Nazario and Thorsten Holz. “As the net churns: Fast-flux botnet observations”. In: *2008 3rd International Conference on Malicious and Unwanted Software (MALWARE)*. 2008 3rd International Conference on Malicious and Unwanted Software (MALWARE). Fairfax, VI: IEEE, Oct. 2008, pp. 24–31. ISBN: 978-1-4244-3288-2. DOI: 10.1109/MALWARE.2008.4690854. URL: <https://ieeexplore.ieee.org/document/4690854/> (visited on 03/15/2022).
- [17] *Number of Internet of Things (IoT) Connected Devices Worldwide from 2019 to 2030*. Statista Inc. Dec. 2020. URL: <https://www.statista.com/statistics/1183457/iot-connected-devices-worldwide/> (visited on 11/11/2021), archived at <https://web.archive.org/web/20211025185804/https://www.statista.com/statistics/1183457/iot-connected-devices-worldwide/> on Oct. 25, 2021.
- [18] Lawrence Page et al. *The PageRank Citation Ranking: Bringing Order to the Web*. Jan. 29, 1998. URL: <http://ilpubs.stanford.edu:8090/422/1/1999-66.pdf> (visited on 11/30/2021).
- [19] Nick Pantic and Mohammad I. Husain. “Covert Botnet Command and Control Using Twitter”. In: *Proceedings of the 31st Annual Computer Security Applications Conference on - ACSAC 2015*. the 31st Annual Computer Security Applications Conference. Los Angeles, CA, USA: ACM Press, 2015, pp. 171–180. ISBN: 978-1-4503-3682-6. DOI: 10.1145/2818000.2818047.

References

- URL: <http://dl.acm.org/citation.cfm?doid=2818000.2818047> (visited on 03/15/2022).
- [20] Christian Rossow et al. "SoK: P2PWNEED - Modeling and Evaluating the Resilience of Peer-to-Peer Botnets". In: *2013 IEEE Symposium on Security and Privacy*. 2013 IEEE Symposium on Security and Privacy (SP) Conference dates subject to change. Berkeley, CA, USA: IEEE, May 2013, pp. 97–111. ISBN: 978-1-4673-6166-8 978-0-7695-4977-4. DOI: 10.1109/SP.2013.17. URL: <https://ieeexplore.ieee.org/document/6547104/> (visited on 03/15/2022).
- [21] Marc Stevens. "Fast Collision Attack on MD5". In: (2006). <https://ia.cr/2006/104>.
- [22] Daniel Stutzbach and Reza Rejaie. "Understanding Churn in Peer-to-Peer Networks". In: *Proceedings of the 6th ACM SIGCOMM on Internet Measurement - IMC '06*. The 6th ACM SIGCOMM. Rio de Janeiro, Brazil: ACM Press, 2006, p. 189. ISBN: 978-1-59593-561-8. DOI: 10.1145/1177080.1177105. URL: <http://portal.acm.org/citation.cfm?doid=1177080.1177105> (visited on 03/08/2022).
- [23] Alex Turing, Hui Wang, and Genshen Ye. *The Mostly Dead Mozi and Its' Lingering Bots*. 360 Netlab. Aug. 30, 2021. URL: <https://blog.netlab.360.com/the-mostly-dead-mozi-and-its-lingering-bots/> (visited on 04/07/2022), archived at <https://web.archive.org/web/20220130162722/https://blog.netlab.360.com/the-mostly-dead-mozi-and-its-lingering-bots/> on Jan. 30, 2022.
- [24] Junjie Zhang et al. "Building a Scalable System for Stealthy P2P-Botnet Detection". In: *IEEE Transactions on Information Forensics and Security* 9.1 (Jan. 2014), pp. 27–38. ISSN: 1556-6013, 1556-6021. DOI: 10.1109/TIFS.2013.2290197. URL: <http://ieeexplore.ieee.org/document/6661360/> (visited on 11/09/2021).

Erklärung

1. Mir ist bekannt, dass dieses Exemplar der Masterthesis als Prüfungsleistung in das Eigentum der Ostbayerischen Technischen Hochschule Regensburg übergeht.
2. Ich erkläre hiermit, dass ich diese Masterthesis selbstständig verfasst, noch nicht anderweitig für Prüfungszwecke vorgelegt, keine anderen als die angegebenen Quellen und Hilfsmittel benutzt sowie wörtliche und sinngemäße Zitate als solche gekennzeichnet habe.

Ort, Datum und Unterschrift

Presented by: Valentin Brandl
Student ID: 3220018
Study Programme: Master Informatik
Supervisor: Prof. Dr. Christoph Skornia
Secondary Supervisor: Prof. Dr. Thomas Waas