

Predicting future botnet addresses with uncleanliness

Michael Collins, Timothy J. Shimeall, Sidney Faber, Jeff Janies, Rhiannon Weaver, Markus De Shon

CERT Network Situational Awareness Group
Software Engineering Institute
Pittsburgh, PA 15213

{mcollins,tjs,sfaber,janies,rweaver,mdeshon}@cert.org

May 9, 2007

Abstract

The increased use of botnets as an attack tool and the awareness attackers have of blocking lists leads to the question of whether we can effectively predict future bot locations. To that end, we introduce a network quality that we term uncleanliness: an indicator of the propensity for hosts in a network to be compromised by outside parties.

We hypothesize that unclean networks will demonstrate two properties: spatial and temporal uncleanliness. Spatial uncleanliness is the tendency for compromised hosts to cluster more densely within unclean networks. Temporal uncleanliness is the tendency for unclean networks to contain compromised hosts for extended periods.

We test for these properties by collating data from multiple indicators (spamming, phishing, scanning and botnet IRC log monitoring). We demonstrate evidence for both spatial and temporal uncleanliness. We further show evidence for cross-relationship between the various datasets, showing that botnet activity predicts spamming and scanning, while phishing activity appears to be unrelated to the other indicators.

1 Introduction

Botnets are a common attack tool due to the anonymity and flexibility that they provide attackers. Modern bots can be used for DDoS, spamming, infiltration of local networks, key-logging and other criminal acts [5, 15]. Past research, notably by Mirkovic *et al.* [18], has shown that botnet based attacks can be divided into distinct phases of acquisition and use.

We expect that bot acquisition is effectively oppor-

tunistic [2]: while attackers may avoid certain networks [24], in the majority of cases, attackers have no interest or knowledge about targets except that the target is vulnerable. With automatically propagating attack tools, an attacker may not know about the existence of a target until after he compromises it.

As bot software has become more sophisticated and flexible, it is now reasonable to expect that any publicly accessible host on the Internet will be attacked by every common method within a short period (for example, specific variants of Gaobot can spread themselves using network shares, AOL Instant Messenger, and multiple Windows vulnerabilities¹).

We therefore expect that within a short time, a host will be attacked by every possible means of compromise². If we assume that an attacker cannot distinguish between the hosts within a network, then he has an equal chance of attacking any of them. In addition, with no advance knowledge of what the target is vulnerable to, an attacker will use all attacks available to him. Consequently, the probability that a machine will be compromised during some period is not a function of that host's attacker, but of its *defenders*.

We hypothesize that networks have a property, which we term *uncleanliness* which is an indicator of the propensity that hosts within a network will be compromised. Our intuition is as follows: consider two institutions with different defensive postures. Institution A maintains an aggressive firewall policy, disables all

¹<http://www.symantec.com/enterprise/security.response/writeup.jsp?docid=2006-052712-1744-99&tabid=2>

²A report of the expected time between attacks for specific vulnerabilities is available at <http://isc.sans.org/survivaltime.html>; the interval between attacks for the average address is on the order of 20 minutes

email attachments, maintains all files on a central server and restores all hosts on the network from a ghosted state each night. Institution B has no central inventory of machines, runs a variety of hardware and software installations that administrators might not even know about, has a large number of self-administered machines and no firewall. We would expect that institution A would be less vulnerable to attacks, and that if a machine was compromised, it would be restored to its original state quickly. Conversely, machines in institution B will be reached by larger number of attacks, and when a machine is compromised, it may not be noticed or repaired until long after the compromise has taken place.

We can estimate the uncleanness of a network by examining its result: once an attacker has compromised hosts, he will use them for criminal activities. If a host is compromised, we expect that the attacker will use it to, for example, spam, scan and DDoS networks. If uncleanness is a network-specific property, we expect that compromised hosts will congregate in specific networks, which we quantify via the phenomena of *spatial* and *temporal* uncleanness. Note that uncleanness is a network level property: hosts are compromised, networks are unclean.

We define *spatial uncleanness* as a tendency for compromised hosts to cluster in unclean networks. Spatial uncleanness implies that if we see a host engaged in hostile activity (such as scanning), we have a good chance of finding another IP address in the same network engaged in hostile activity. We will test for spatial uncleanness by examining the clustering of addresses within networks.

We define *temporal uncleanness* as a tendency for compromised hosts to repeatedly appear in unclean networks. Temporal uncleanness implies that if a host is compromised, then other hosts within that network will be compromised in the future. We will test for temporal uncleanness by examining the ability of unclean networks to predict future host compromises.

Figure 1 confirms our intuition for spatial uncleanness and temporal uncleanness. This figure shows two plots: the upper counts the number of unique hosts scanning a large network from January to April, 2006. The lower plot is a plot showing how many of these scanning addresses were also present in a botnet reported during the first week of March, 2006. This plot contains two lines: one counts the number of unique addresses from the bot report which were also identified scanning; the second counts the number of unique addresses from the bot report which were present in a 24-bit CIDR block where at least one address was also scanning.

First note that these reports resulted from two different detection methods: the bots were collected by observing IP addresses communicating on IRC channels, while scanning data was collected using a behavioral scan detection method deployed on an observed network [6]. Despite this, there is a strong intersection between the two sets: at its peak, 35% of the botnet's addresses are scanning the observed network.

Second, we observe that using the /24's comprising the botnet identifies more scanners than the botnet addresses alone; this value ranges between a 25% increase and 4 times as many addresses depending on the activity. We demonstrate in §4 that these results are statistically significant.

Finally, Figure 1 also explains our intuition for *temporal uncleanness*. As this figure shows, abnormal scanning (and therefore botnet compromise) occurs over several weeks. If bots take several weeks to be identified and removed, we expect that an unclean network will be unclean for some duration, and therefore we can predict future hostile activity from the same network.

In this paper, we examine four potential indicators of uncleanness: botnet data, scanning activity, spamming and phishing. We collect reports of unclean activity from multiple sources: public mailing lists and web sites, private studies, and by examining traffic crossing a large (multiple /8) network.

The primary contribution of this paper is an empirical study of uncleanness and its use as a predictive aid. We test for the existence of spatial and temporal uncleanness by comparing the traffic from various reports of hostile activity. We demonstrate that compromised hosts are both more densely clustered than normal traffic and predict future unclean activity. In addition, we show that scanning, spamming and botnet activity shows evidence of cross relationship, such as the scanning observed in Figure 1. We also show that while these phenomena do not predict future phishing sites, past phishing sites do, therefore demonstrating that temporal uncleanness holds for all four indicators. We then test the strength of this predictive mechanism by evaluating its suitability to block traffic crossing a large network. We demonstrate that limited predictive blocking is feasible, due to the impact of locality [17] evident in network traffic.

The remainder of this paper is structured as follows: §2 outlines relevant previous work in reputation management and identifying hostile groups by past history. §3 describes our model and the data sources we use in this paper. §4 examines the spatial uncleanness hypothesis, and §5 examines the temporal uncleanness hypothesis. §6 examines the impact of blocking unclean

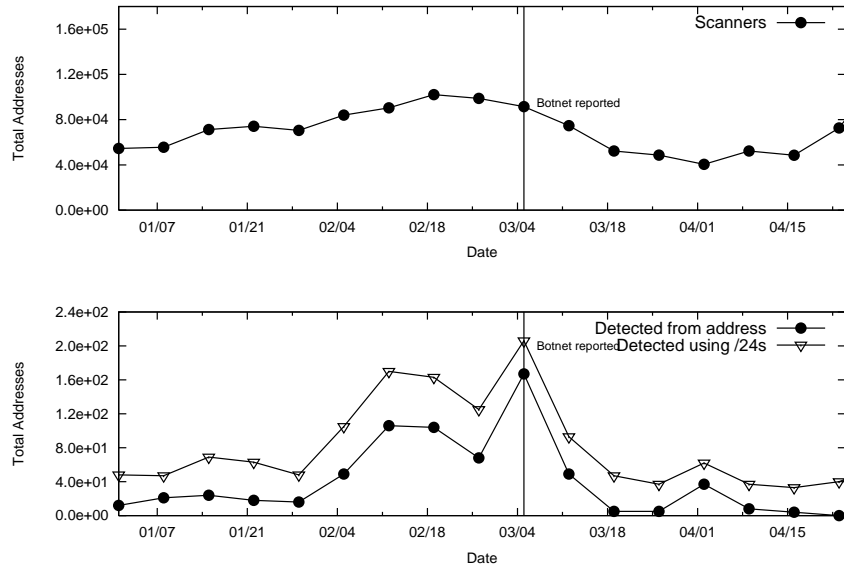


Figure 1: Relationship between scanning and botnet population

networks, and §7 discusses the results.

2 Previous Work

Researchers initially studied botnets due to their use in DDoS attacks; in this domain, Mirkovic *et al.*[18] defined a DDoS attack as a two-phase process: acquiring hosts to use for the DDoS and then using those hosts to conduct an attack. Freiling *et al.*[5] identify a variety of other attacks that botnets can conduct efficiently, Collins *et al.* [2] define bot occupation attacks as conducted by opportunistic attackers: that is, the attacker has no interest or knowledge of the target except that the target is exploitable. Our work uses these concepts to study the impact of largely automated acquisition and its impact on network defense.

Botnet demographics have been studied using Honeybots and by actively probing bot networks [8, 9, 21]. Rajand *et al.*'s [21] analysis is particularly relevant due to the extended period during which they observed network traffic, allowing them to identify not only botnet demographics but activity. Our work differs from these analyses by comparing multiple observed phenomena

and using this information to predict future activity.

In operational security, blacklists are commonly used to identify and block hosts that are already assumed to be hostile. Examples of such blacklists include Spamhaus' ZEN list [20] and the Bleeding Snort rule set [23]. Researchers such as Levy [16] note that spammers increasingly rely on the use of occupied hosts to generate spam messages - these approaches are more attractive to spammers because they offload processing requirements from the spammer (as noted by Laurie *et al.*[15]) and because they hide the attacker's identity[4].

In addition, researchers have studied the impact of blacklists on spamming and other hostile activity Jung *et al.*[12] compare spamming blacklists against spam traffic to MIT in 2000 and 2004, finding that in 2004, 80% of spammers were identified by blacklists. Ramachandran *et al.*[22], examine Blacklist abuse by botnet owners. Ramachandran notes that botnet owners appear to place a higher premium on addresses not present on block lists. Since uncleanliness is intended to predict future hostile addresses, this may impact the costs noted by Ramachandran.

McHugh *et al.* use locality to characterize normal network behavior and differentiate attacks. Krishna-

murthy *et al.* [14] use netblocks to characterize target audiences for networks, and demonstrate that many sites have common audiences. This leads to a generalized netblock-level approach developed by Jung *et al.* [10] for DDoS defense. These methods of blocking are predicated on the assumption that attack traffic differs from normal traffic due to a limited and clustered audience for any normal service. Our filtering approach differs from the past history used in these cases by developing a set of explicitly untrusted networks.

3 Source Data

We demonstrate evidence of uncleanliness by showing that address distributions from unclean data sets show specific qualities; in order to do so, we must collate information from various sources with different collection methods. In this section we describe a simple taxonomy and notation scheme for managing our data; in the following sections we use this data to demonstrate significance. This section is divided as follows, §3.1 explains the taxonomy and notation for reports, and §3.2 describes the individual reports.

3.1 Model

In order to estimate the uncleanliness of a network, we must compare data from multiple sources. For example, an attacker may initially use a bot for scanning, then for spamming. We call these sources *reports*, each of which consists of a set of IP addresses describing a particular phenomenon over some period. Reports differ by the *class* of data reported, the period covered by the report, and the method used to generate that data.

We use four classes of unclean data for this paper:

1. **Bots:** An IP address identified as hosting some form of bot software or communicating with a botnet command and control host.
2. **Phishing:** An IP address identified as hosting a phishing site in order to fraudulently acquire private user information.
3. **Scanning:** An IP address identified as scanning using the methods developed by Gates *et al.* [7] and Jung *et al.* [11].
4. **Spamming:** An IP address identified as spamming using a behavioral spam detection technique³.

³This spam detection method is currently under review

These reports all describe phenomena associated with compromised hosts. Scanning and spamming are both common botnet uses, and phishing requires setting up a fraudulent web site.

We further divide reports as either *provided* or *observed*. Provided reports are collected from external parties, and can use different methodologies to observe the same effects. For example, a phishing list can acquire IP addresses by using spam traps [19] or by collecting user reports, (e.g., the submission form at the Castle-Cops PIRT service [1]). For the analyses within this paper, we use only one source per report and assume that the source’s collection methodology is consistent over the report period. In contrast to provided reports, observed reports are generated from network traffic logs reporting traffic covering a large edge network.

We use a simple notation to describe all reports; each report is differentiated by a *tag* which, for this paper, summarizes the period and source for the report. We express this using the notation \mathcal{R}_T . In this form, T is the tag (e.g., scan). A list of reports used in this paper is given in Table 1.

Because we expect uncleanliness to be a network property, we define a CIDR masking function $C_n(i)$. The CIDR masking function evaluates to the unique CIDR block with prefix length n that contains the IP address i (e.g., $C_{16}(127.1.135.14) = 127.1.0.0/16$). For convenience, when the CIDR masking function is applied on a report S , the result is set-valued and returns the set of all n -bit CIDR blocks in that set, that is:

$$C_n(S) \equiv \bigcup_{i \in S} C_n(i) \quad (1)$$

When determining whether or not an IP address resides within a set of CIDR blocks, we will use a CIDR inclusion relation, \sqsubset , to indicate that an IP address is resident in one of a set of CIDR blocks:

$$i \sqsubset S \rightarrow \exists n \text{ s.t. } C_n(i) \in S \quad (2)$$

With all sets and reports, we use bars to indicate cardinality, i.e., $|S|$ is the number of elements in the set S .

3.2 Reports

Table 1 is an inventory of all the reports used in this paper. Recall that provided reports have been given to us by other parties and that we generate observed reports using traffic logs from the observed network. Because we have greater control over observed reports, we can generate these reports over arbitrary periods. We have less control over when we receive provided reports.

The observed network is composed of over 20 million distinct IPv4 addresses and contains several servers that are heavily used by clients across the Internet. Given the size and activity of the observed network, we assume that IP addresses from the Internet crossing into it are a representative sample of the Internet as a whole.

In order to compensate for selection bias within observed reports, all reports have been filtered to only include addresses which are outside of the observed network and are not otherwise reserved (*e.g.*, all addresses specified in RFC 1918 have been removed from reports).

We classify four of the reports in this list as *unclean reports*, these are the reports we use as ground truth for identifying the four indicators discussed in §3.1: botnet membership, phishing sites, scanners and spammers. During the two week period of October 1st-14th, 2006, we have both provided and observed reports on all classes of unclean activity, consequently we use October 1st-14th to test temporal uncleanliness.

The next set of reports are used specifically to test the spatial and temporal uncleanliness hypotheses. The bot – test report describes a small botnet from 5 months before all the other activity analyzed in this paper, bot – test is used as an extreme case for prediction: if a five-month old report can accurately predict current unclean activity, then a more recent one should be more effective.

The control report consists of 47 million unique IP addresses observed during the week of September 25th, 2006. We compare the data from our other reports against randomly generated subsets of control in order to determine whether or not these reports exhibit spatial or temporal uncleanliness. We use the control report to more accurately reflect the structure of IPv4 space than we would using purely randomly chosen IP addresses. The report consists of IP addresses observed to engage in payload-bearing TCP activity, which reduces the risk of the address being spoofed. Furthermore, as noted in §3.1, the observed network includes a variety of servers used by hosts throughout the Internet, and by focusing exclusively on the IP addresses of the hosts without using any criteria apart from the unspoofed criterion, we expect the resulting report to approximate a random sample of active IP addresses on the Internet.

4 Spatial Uncleanliness

We define *spatial uncleanliness* as the propensity for occupied addresses (bots) to be clustered in unclean networks. In this section, we formulate and test the *spatial uncleanliness hypothesis*.

This section is divided as follows: §4.1 describes the methodology used to test for spatial uncleanliness. §4.2 describes the results of our tests and shows evidence for spatial uncleanliness.

4.1 Model and methodology

Recall our assumption that the likelihood of a host being compromised is a network property: if a network is unclean, then its administrators will not identify compromised machines or rapidly repair them. Consequently, we expect that multiple hosts within an unclean network will be compromised, and that compromised addresses will cluster within unclean networks. In order to test this hypothesis, we will compare the expected population of compromised hosts within equally sized CIDR blocks.

To test for spatial uncleanliness, we begin with a measurement for comparative density. If we have two sets, S_1 and S_2 , and $|S_1| = |S_2|$, then we say that S_1 is *denser at n-bits* if the number of n-bit CIDR blocks that S_1 occupies is less than the number of n-bit CIDR blocks occupied by S_2 .

Throughout this paper, we use homogeneously sized CIDR blocks to model individual networks. While other categorization techniques are available we opt to use homogeneously sized CIDR blocks in order to control for population. Heterogeneous partitions, such as Krishnamurthy *et al.*'s network-aware clustering method [14]), result in networks that differ in size by several orders of magnitude.

In §1, we stated that spatial uncleanliness implies that if a host is compromised, there is a good chance another host on the same network will be compromised. Consequently, if we had a set of compromised host addresses, and a control set of randomly selected addresses with equal cardinality, we would expect that the compromised address set was *at least* as dense at all CIDR prefix lengths.

We therefore summarize the spatial uncleanliness hypothesis as follows: if we have a report which selects unclean traffic from the Internet, $\mathcal{R}_{\text{unclean}}$, then the IP addresses within that report will be more densely packed than a randomly selected set of IP addresses with equal cardinality.

To test the spatial uncleanliness hypothesis, we use the formulation given in Equation 3 below. Assuming that we have two reports, $\mathcal{R}_{\text{unclean}}$ which reports on unclean traffic, and $\mathcal{R}_{\text{control}}$ which is control data, and both reports are of equal cardinality, then:

$$\forall n \in [16, 32] |C_n(\mathcal{R}_{\text{unclean}})| \leq |C_n(\mathcal{R}_{\text{control}})| \quad (3)$$

Unclean reports					
Tag	Type	Class	Valid Dates	Size	Reporting method
bot	Provided	Bots	2006/10/01-2006/10/14	621,861	Bot addresses acquired through private reports from a third party
phish	Provided	Phishing	2006/05/01-2006/11/01	53,789	Addresses from a Phishing report list
scan	Observed	Scanning	2006/10/01-2006/10/14	151,908	IP addresses scanning the observed network
spam	Observed	Spam	2006/10/01-2006/10/14	397,306	IP addresses spamming the observed network
Reports for hypothesis testing					
bot – test	Provided	Bots	2006/05/10	186	Botnet addresses acquired through private communication
control	Observed	N/A	2006/09/25-2006/10/02	46,899,928	Control addresses acquired from the observed network

Table 1: Table of tags used in this report

Based on DDoS filtering work done by Collins and Reiter [3], we expect that 16 bit prefix lengths will be too imprecise for effective filtering and detection. Consequently, we limit our prefix lengths to between 16 and 32 bits.

4.2 Analysis

In order to test the spatial uncleanliness hypothesis, as formulated in Equation 3, we compare the population of addresses per n -bit CIDR blocks for an unclean report against the expected population for n -bit CIDR blocks across the Internet as a whole.

As discussed in §3.2 we model network populations by randomly selecting IP addresses from the $\mathcal{R}_{\text{control}}$ report. Kohler *et al.* [13] observe that IP addresses are not evenly distributed across IPv4 space; as a consequence, a purely random model will result in an artificially depressed density estimate. To compensate for this, we test two population estimates. The first, naive, estimate selects addresses evenly from across all /8's which are listed as populated by IANA⁴. The second, empirical, density estimate draws addresses from a pool of addresses observed to cross the network under observation from the week of September 25th–October 1st, 2006. In the empirical estimate, we create 1000 randomly generated subsets of $\mathcal{R}_{\text{control}}$ and group the resulting addresses.

Figure 2 plots the number of addresses per block

⁴<http://www.iana.org/assignments/ipv4-address-space>

for CIDR block prefix lengths of 16 to 32 bits. This plot compares the botnet density, \mathcal{R}_{bot} , against both the empirical and naive density estimates of equal size (621,861 addresses, as per Table 1). As this figure shows, the botnet population is more tightly packed than both empirical and naive estimates. In the case of the empirical estimate, botnet data results in nearly twice as many addresses per block for prefix lengths between 18 and 24 bits. The naive estimate is zero throughout these results. Based on the results from Figure 2, we use empirical estimation throughout the rest of this paper.

Figure 3 compares control data (empirically estimated populations) against each of our four datasets: spamming, scanning, botnet population and phishing. In comparison to the population plot in Figure 2, these plots represent the total number of n -bit blocks observed for that population; since each population is of equal size, the lowest line will have the highest density. For each plot in Figure 3, the control data consists 1000 random subsets of $\mathcal{R}_{\text{control}}$ and plotting the resulting distribution as a boxplot.

Figure 3(i) is a plot of the comparative volume for \mathcal{R}_{bot} . As this plot shows, the population of \mathcal{R}_{bot} is more densely packed than the expected population drawn from $\mathcal{R}_{\text{control}}$. Figure 3(ii) plots the volume of $\mathcal{R}_{\text{phish}}$ reported from May to October, 2006. We use a five month sample due to the smaller size of the phishing reports in comparison to the other reports. As noted in Table 1, the 6 month phishing report is approximately an order of magnitude smaller than the other unclean reports. As with Figure 3(i), addresses in the phishing

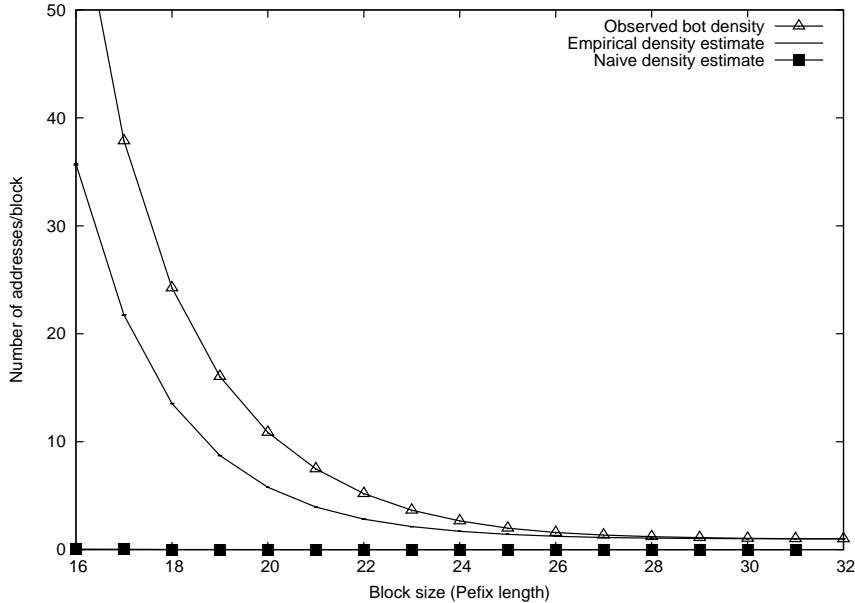


Figure 2: Density of botnets per netblock, compared against empirical and naive control sets

report are more tightly packed than addresses selected from the control report.

Figure 3(iii) plots the volume of $\mathcal{R}_{\text{spam}}$ from October 1st to 14th, 2006. Figure 3(iv) plots the volume of $\mathcal{R}_{\text{scan}}$ for the same period. Each of these reports is more tightly packed than the comparative control reports.

As Figures 2 and 3 show, unclean reports have an n -bit density greater than or equal to or greater than the n -bit density of the control reports for all values of n . Consequently, this data supports the spatial uncleanliness hypothesis: compromised hosts are disproportionately concentrated in certain networks.

5 Temporal Uncleanliness

We now address temporal uncleanliness: the propensity for networks to remain unclean for extended periods of time. In order to test for temporal uncleanliness we compare the ability of a report of unclean addresses to predict future compromised addresses; in particular, whether or not a report of bot addresses can predict future bots, spamming, scanning and phishing.

This section is divided as follows: §5.1 describes our

method for measuring the presence of temporal uncleanliness, and §5.2 shows the results.

5.1 Model and methodology

To observe temporal uncleanliness, we examine the *predictive* capacity of reports of unclean data. Consider three reports: $\mathcal{R}_{\text{test}}$, $\mathcal{R}_{\text{control}}$ and $\mathcal{R}_{\text{result}}$. If $\mathcal{R}_{\text{test}}$ and $\mathcal{R}_{\text{control}}$ are of equal cardinality, then $\mathcal{R}_{\text{test}}$ is a better predictor of the report $\mathcal{R}_{\text{result}}$ at prefix length n if:

$$\frac{|C_n(\mathcal{R}_{\text{test}}) \cap C_n(\mathcal{R}_{\text{result}})|}{|C_n(\mathcal{R}_{\text{control}}) \cap C_n(\mathcal{R}_{\text{result}})|} > 1 \quad (4)$$

If temporal uncleanliness exists, then we expect that unclean reports will consistently be better predictors of future unclean reports than a control report. However, we note that due to spatial uncleanliness, an unclean report will have fewer n -bit CIDR blocks than an equivalent control report. As a consequence, as block size increases, the control report will have a larger number of imprecise successes. Therefore, there will be some prefix length below which a control report will always be a better predictor than the test report.

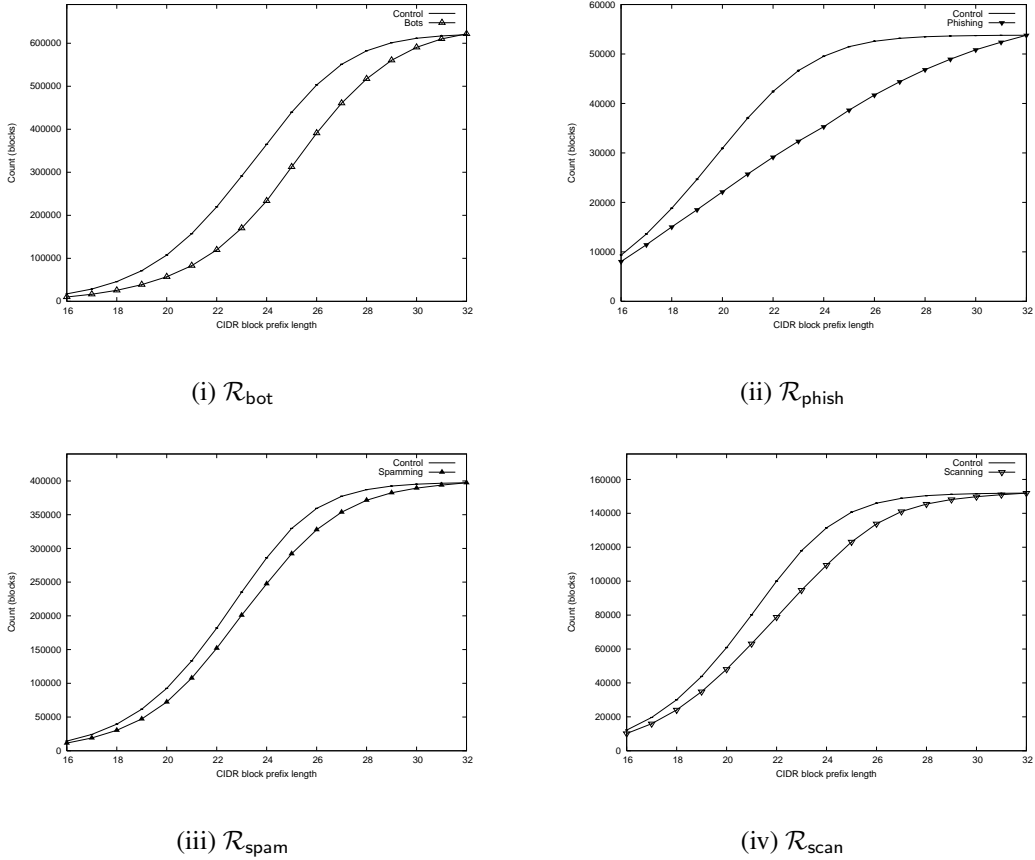


Figure 3: Comparative density of Unclean netblocks against $\mathcal{R}_{\text{control}}$

For testing, we use the form of the temporal uncleanliness hypothesis given in the equation below. Given that $\mathcal{R}_{\text{unclean}}$ and $\mathcal{R}_{\text{control}}$ have equal cardinality, then

$\exists n \in [16, 32]$ s.t.

$$|C_n(\mathcal{R}_{\text{unclean}}) \cap C_n(\mathcal{R}_{\text{result}})| > |C_n(\mathcal{R}_{\text{control}}) \cap C_n(\mathcal{R}_{\text{result}})|$$

That is, there exists a prefix length where a previously generated report of unclean activity is more predictive of present unclean activity than a control report of equal cardinality. As with spatial uncleanliness, we limit our analyses to CIDR blocks of at least 16 bits.

5.2 Analysis

We now test the temporal uncleanliness hypothesis formulated in Equation 5. To do so, we compare the effective

predictiveness of $\mathcal{R}_{\text{bot-test}}$ on the unclean reports during the period of October 1st-14th, 2006.

Figure 4 shows the relative predictive capacity of $\mathcal{R}_{\text{bot-test}}$ against future unclean reports; for these figures, $\mathcal{R}_{\text{phish}}$ is a sub report of $\mathcal{R}_{\text{phish}}$ from Table 1. This report is considerably smaller than the other reports, with 2302 addresses. This results in a smaller degree of intersection with the randomly generated reports from the control report.

As in §4.2, we generate the reference line by plotting a boxplot showing the variance of 1000 randomly selected test reports. In contrast with Figure 3, the small cardinality of $\mathcal{R}_{\text{bot-test}}$ ensures that the variations observed by the boxplot are visible. We consider the $\mathcal{R}_{\text{bot-test}}$ to be a better predictor than $\mathcal{R}_{\text{control}}$ if the cardinality of its intersection with the corresponding unclean report is higher than the intersection with randomly selected addresses in 95% of the observed cases.

As Figure 4 shows, $\mathcal{R}_{\text{bot-test}}$ is a better predictor than

$\mathcal{R}_{\text{control}}$ for botnets, spamming and scanning at various prefix lengths. Also of note is the impact of spatial uncleanliness: in these three figures, $\mathcal{R}_{\text{bot-test}}$ is a better predictor for prefix lengths of approximately 19-20 bits and longer. At shorter prefix lengths, randomly selected addresses become better predictors. Using the 95% threshold, $\mathcal{R}_{\text{bot-test}}$ is a stronger predictor of future botnet activity between 20 and 25 bits, spamming between 19 and 32 bits, and scanning between 20 and 24 bits. For prefix lengths longer than these values, the two reports are equally predictive due to the low probability of seeing CIDR blocks from either report intersect.

Figure 4(ii) plots the predictive capacity of $\mathcal{R}_{\text{bot-test}}$ against $\mathcal{R}_{\text{phish}}$. In contrast to the other plots in Figure 4, this plot indicates that $\mathcal{R}_{\text{bot-test}}$ is not a good predictor of future phishing activity in comparison to randomly selected control sets.

We have two hypotheses as to why phishing this is so: Ramachandran *et al.* [22] describe how botnet owners place a higher premium on addresses that have not yet been identified as bots. Because phishing sites need to be publicized, a phishing IP address becomes public knowledge, marked on blacklists and consequently highly unattractive for the owner of a botnet.

An alternative explanation is that, in contrast to botnets, phishing sites are generally hosted on web servers, and a phisher may prefer to host phishing sites in a actual datacenter to ensure robustness during a flash crowd. At the minimum, a phishing site must be publicly accessible, while a bot can exist behind a NAT or a firewall and still be useful. Therefore, phishers may prefer sites that are already hosting web servers and have the resources to handle a high traffic load.

In order to determine whether the temporal uncleanliness hypothesis does hold for phishing, we now consider a test that uses phishing data exclusively. Figure 5 plots the intersection of $\mathcal{R}_{\text{phish-test}}$ against the same phishing set as in Figure 4(ii). In this case, $|\mathcal{R}_{\text{phish-test}}| = 1386$. We note that this figure shows strong evidence for temporal uncleanliness in phishing.

Since these results show that five month old reports can be used to more effectively predict the population of future reports than randomly selected IP addresses from a week before, we conclude that the temporal uncleanliness hypothesis is supported by this data. Furthermore, in Equation 5, we chose a range of IP blocks arbitrarily, we can now establish a lower limit for the prefix length of 20 bits, an an upper limit in excess of 24 bits.

We have also shown that phishing activity and botnet activity are not related in the way that bots, scanning and spamming are. As noted elsewhere [21, 15], scanning

and spamming are commonly implemented with botnets, so we would expect that \mathcal{R}_{bot} , $\mathcal{R}_{\text{scan}}$ and $\mathcal{R}_{\text{spam}}$ are related. However, the inability of $\mathcal{R}_{\text{bot-test}}$ to predict future phishing activity suggests that a measurement for uncleanliness will have to be multidimensional: phishing sites are still taken over, but it may be that phishers have different criteria for the machines they occupy than botnet owners.

6 Blocking Tests

The spatial and temporal uncleanliness hypotheses together provide a method for identifying compromised hosts. Spatial uncleanliness implies that if an address within a network is occupied, then we can expect other networks within the same netblock to be occupied. Temporal uncleanliness indicates that if we have seen an address in the past used for an attack, then we can expect addresses from the same network to do so in the future.

We now address the issue of whether unclean networks can be *effectively* blocked; that is, whether or not blocking a set of unclean networks will adversely affect traffic into an active network. To do so, we will examine the impact of blocking a set of unclean networks from two weeks of network traffic. In this section, we compare the expected false positive and false negative values over a range of *operating characteristic* values. For this work, the operating characteristic is n , the CIDR block prefix length.

We begin by collecting traffic logs of all traffic that crosses the observed network from all IP addresses $i \in C_{24}(\mathcal{R}_{\text{bot-test}})$ for the test period of October 1st-14th 2006. This report, $\mathcal{R}_{\text{candidate}}$ consists of all IP addresses crossing the observed network which share a /24 in common with any of the IP addresses in $\mathcal{R}_{\text{bot-test}}$. This allows us to test the effectiveness of filtering from the /24 to the /32 range; we pick this range because, as seen in Figure 3, 24 bits is the minimum block size at which $\mathcal{R}_{\text{bot-test}}$ is an unambiguously better predictor of future uncleanliness than control data. We further constrain $\mathcal{R}_{\text{candidate}}$ to those addresses which generate at least one TCP record during this time period.

The source data used for this analysis is CISCO NetFlow⁵ traffic records, which are a compact summarization of traffic information, but do not contain payload. As a consequence, our analysis will have some degree of uncertainty as we cannot directly validate the payload. We will therefore differentiate addresses in two ways: by membership in one of the unclean reports and

⁵<http://www.cisco.com/go/netflow>

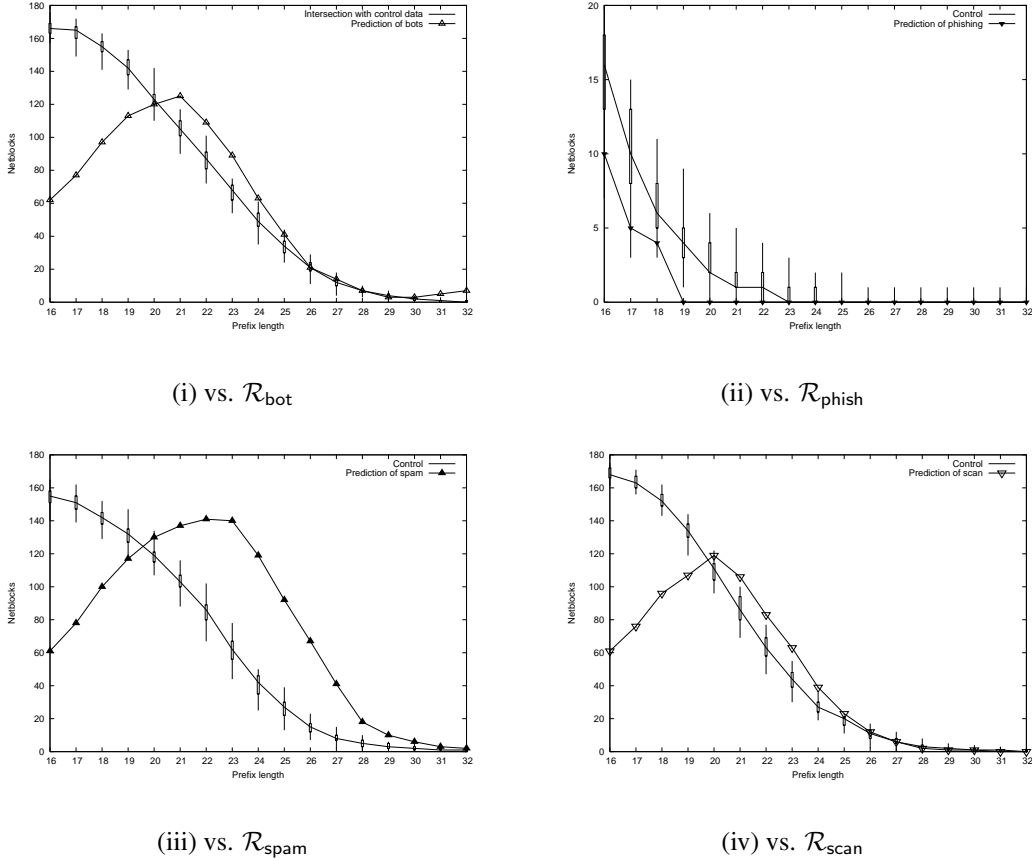


Figure 4: Comparative predictive capacity of $\mathcal{R}_{\text{bot-test}}$ against control data

by behavior.

We partition the $\mathcal{R}_{\text{candidate}}$ into three reports: $\mathcal{R}_{\text{unknown}}$, $\mathcal{R}_{\text{hostile}}$ and $\mathcal{R}_{\text{innocent}}$. A full inventory of the reports used in this analysis is given in Table 2.

$\mathcal{R}_{\text{hostile}}$ consists of any IP address in $\mathcal{R}_{\text{candidate}}$ that is also present in the unclean reports (i.e., scanning, spamming, phishing or botnet membership). The hostile set is identified purely by intersecting these reports, and once an IP address is identified as hostile it cannot be present in the remaining two reports. $\mathcal{R}_{\text{unknown}}$ is comprised of the addresses in $\mathcal{R}_{\text{candidate}}$ address which are *not* present in one of the unclean reports, but have no *payload bearing* flows. We define a flow as payload-bearing if it is a TCP flow with at least 36 bytes of payload and at least one ACK flag. Due to TCP options, a 3-packet SYN scan will often have 36 bytes of payload, even though this data is still part of the TCP handshake. Hand-examination of the flow logs found multiple examples

of 36-byte SYN-only scans to apparently random ports on distributed targets.

The IP addresses in $\mathcal{R}_{\text{unknown}}$ are not proven to be hostile but are highly suspicious. Due to the lack of payload in flow data, we cannot definitively categorize members of this report into either of the other two reports and consequently we remove them from the false positive and false negative calculations.

The population of $\mathcal{R}_{\text{innocent}}$ consequently consists of any IP address which does conduct payload-bearing TCP activity and is not present in any of the unclean reports.

Our prediction scenario assumes that an organization received $\mathcal{R}_{\text{bot-test}}$ and is blocking $C_n \mathcal{R}_{\text{bot-test}}$ for some value of $n \in [24, 32]$. The success of this defensive mechanism is based on how many hostile and innocent addresses are blocked by the attack mechanism (as noted above, while the unknown population is calculated and

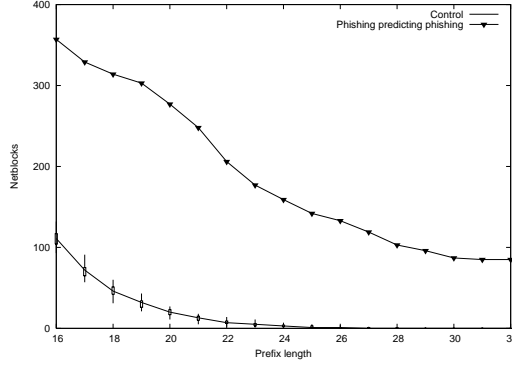


Figure 5: Comparative predictive capacity of phishing reports

Reports used for prediction testing					
Tag	Type	Class	Valid Dates	Size	Reporting method
unclean	Provided	Special	2006/10/01-2006/10/14	1,158,103	The union of the four unclean reports, note that there is overlap
candidate	Observed	N/A	2006/10/01-2006/10/14	1030	IP Addresses crossing the network border and which are in the same /24's as $\mathcal{R}_{\text{unclean}}$
hostile	Observed	N/A	2006/10/01-2006/10/14	287	Members of $\mathcal{R}_{\text{candidate}}$ also present in $\mathcal{R}_{\text{unclean}}$
unknown	Observed	N/A	2006/10/01-2006/10/14	708	Members of $\mathcal{R}_{\text{candidate}}$ not in $\mathcal{R}_{\text{unclean}}$, but engaged in suspicious activity
innocent	Observed	N/A	2006/10/01-2006/10/14	35	Members of $\mathcal{R}_{\text{candidate}}$ not present in $\mathcal{R}_{\text{hostile}}$ or $\mathcal{R}_{\text{unknown}}$

Table 2: Table of reports used for prediction test

analyzed in this exercise, it is not scored). The score for the defensive mechanism is the relative success, measured in true and false positives of the filter as a function of n . We define a false positive as a member of $\mathcal{R}_{\text{innocent}}$ blocked by the filter, and true positive as a member of $\mathcal{R}_{\text{hostile}}$ blocked by the filter.

To calculate the true and false positive rates, we define a membership function, m :

$$m(i, S) = \begin{cases} 1 & C_{32}(i) \sqsubset C_{32}(S) \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

For any prefix length n , we calculate the population as a function of n by summing the unique IP addresses that appear within the $\mathcal{R}_{\text{bot-test}}$

$$\text{pop}(n) = \sum_{i \sqsubset C_n(\mathcal{R}_{\text{bot-test}})} m(i, \mathcal{R}_{\text{candidate}} \cap (\mathcal{R}_{\text{innocent}} \cup \mathcal{R}_{\text{hostile}})) \quad (7)$$

As noted above, this calculation explicitly avoids the use of $\mathcal{R}_{\text{unknown}}$. We calculate the true positive and true negative values by calculating a similar value over the various reports:

$$\text{TP}(n) = \sum_{i \sqsubset C_n(\mathcal{R}_{\text{bot-test}})} m(i, \mathcal{R}_{\text{candidate}} \cap \mathcal{R}_{\text{hostile}}) \quad (8)$$

$$FP(n) = \sum_{i \in C_n(\mathcal{R}_{\text{bot-test}})} m(i, \mathcal{R}_{\text{candidate}} \cap \mathcal{R}_{\text{innocent}}) \quad (9)$$

Table 3 summarizes the effectiveness of this prediction method. As this table shows, all three populations increase as the bit length increases. At $n = 24$, 90% of the incoming addresses are correctly identified as hostile. If we assume that unknown addresses are hostile, true positive rate is 97%. Furthermore, the false positive rate remains relatively low until $n = 26$.

n	$TP(n)$	$FP(n)$	pop(n)	$\mathcal{R}_{\text{unknown}}$
24	287	35	322	708
25	172	22	194	344
26	81	1	82	200
27	38	1	39	105
28	18	0	18	60
29	7	0	7	29
30	1	0	1	14
31	1	0	1	7
32	1	0	1	0

Table 3: Observed true and false positive counts

Of note with this dataset are the volume of uncertain addresses (i.e., the population of $\mathcal{R}_{\text{unknown}}$). At a 24 bit prefix length, $|C_{24}(\mathcal{R}_{\text{bot-test}}) \cap C_{24}(\mathcal{R}_{\text{unknown}})|$ yields approximately 700 addresses. We first note that unknown addresses have engaged in TCP communications, but have not exchanged payload - consequently, blocking these addresses does not impact traffic.

Of more concern is that all of the addresses in $\mathcal{R}_{\text{unknown}}$ engage in some form of suspicious behavior (that is, suspicious apart from trying to connect with the network and not exchanging payload). Hand examination found many addresses trying to open communications from ephemeral ports to ephemeral ports (notably port TCP/51736) and slow scanning (the scan detection mechanism is calibrated to identify scans that take place over an hour, scans observed in this dataset would often contact less than 30 addresses a day over the course of the test period).

The strength of this blocking method is predicated on the relatively sparse amount of traffic issuing from these netblocks. As Table 3 shows, 1030 IP addresses were blocked when n was set to 24 bits. $|C_{24}(\mathcal{R}_{\text{bot-test}})| = 173$, which yields a potential set of 44,288 addresses that can be blocked. Consequently, less than 2% of the total IP addresses available in those /24s communicated with

the observed network during this time.

Some of the effectiveness of this method may be attributed to the demographics of the botnet and the network $\mathcal{R}_{\text{bot-test}}$ consists primarily of addresses outside the English-speaking world, with 70% of the addresses coming from Turkey. In addition, the network under observation is primarily an edge network; that is, all traffic at its border is either originating from an address within that border or going to an IP address within that border. Therefore, while we have shown that a five-month old botnet can still be used to effectively predict and halt hostile traffic, issues of demographics and a network’s target audience must also be evaluated.

7 Conclusion

In this paper, we have demonstrated that it is possible to effectively predict future hostile activity from past network activity. To do so, we have defined a network-based quality of uncleanness, which is an indicator of how likely a network is to contain compromised hosts.

As an initial work in this field, we have focused on testing basic hypotheses about uncleanness, which we have defined with the spatial and temporal uncleanness hypotheses. Using reports of network activity and traffic logs of a large network we have shown evidence of spatial and temporal uncleanness. We have also shown that an uncleanness measure may involve multiple dimensions, such as botnets and phishing.

Finally, we have demonstrated that spatial and temporal uncleanness, coupled with the limited audience of an edge network, can be effectively used to block hostile traffic in the future. Given the demographics issues noted in §6, uncleanness may best be used as a risk indicator - by showing that a network is demonstrating in unclean behavior, security personnel can evaluate whether the risk of hostile activity from the network is worth the benefit of receiving commerce and communication from that network under normal circumstances.

Our immediate goal following this work is to develop a more rigorous metric for uncleanness, in particular a multidimensional uncleanness metric to measure the aggregate probability that an address is occupied. The elements of this metric involve the components discussed in this work as well as other predictive indicators of vulnerability (communication with botnet C&C nodes).

We also believe that spatial uncleanness, in particular, has useful implications for network log analysis. If we know that a host from one network is attacking, scanning or otherwise interfering with the traffic on an

observed network, it is reasonable to examine other traffic from that network to see if there is coordinated hostile activity.

References

- [1] CastleCops. Castlecops phishing incident reporting & termination (PIRT) squad. Accessible at <http://www.castlecops.com/pirt>, fetched on January 29th, 2007.
- [2] M. Collins, C. Gates, and G. Kataria. A model for opportunistic network exploits: The case of P2P worms. In *Proceedings of the 2006 Workshop on Economics and Information Security*, 2006.
- [3] M. Collins and M. Reiter. An empirical analysis of target-resident dos filters. In *Proceedings of the 2004 IEEE Symposium on Security and Privacy*, 2004. May 9 – 12, 2004.
- [4] D. Cook, J. Hartnett, K. Manderson, and J. Scanlan. Catching spam before it arrives: domain specific dynamic blacklists. In *ACSW Frontiers '06: Proceedings of the 2006 Australasian workshops on Grid computing and e-research*, Darlinghurst, Australia, Australia, 2006.
- [5] F. Freiling, T. Holz, and G. Wicherski. Botnet tracking: Exploring a root-cause methodology to prevent distributed denial-of-service attacks. In *Proceedings of the 2005 European Symposium on Research in Computer Security*, 2005.
- [6] C. Gates, J. McNutt, J. Kadane, and M. Kellner. Detecting scans at the isp level. Technical Report CMU/SEI-2006-TR-005, Software Engineering Institute, 2006.
- [7] C. Gates, J. McNutt, J. Kadane, and M. Kellner. Scan detection on very large networks using logistic regression modeling. In *ISCC '06: Proceedings of the 11th IEEE Symposium on Computers and Communications*, Washington, DC, USA, 2006.
- [8] T. Holz. Learning more about attack patterns with honeypots. In *Sicherheit 2006: Sicherheit - Schutz und Zuverlässigkeit, Beiträge der 3. Jahrestagung des Fachbereichs Sicherheit der Gesellschaft für Informatik e.v. (GI), 20.-22. Februar 2006 in Magdeburg*, 2006.
- [9] T. Holz, S. Marechal, and F. Raynal. New threats and attacks on the world wide web. *IEEE Security & Privacy*, 4(2), 2006.
- [10] J. Jung, B. Krishnamurthy, and M. Rabinovich. Flash crowds and denial of service attacks: Characterization and implications for CDNs and web sites. In *Proceedings of the International World Wide Web Conference*, May 2002.
- [11] J. Jung, V. Paxson, A. Berger, and H. Balakrishnan. Fast Portscan Detection Using Sequential Hypothesis Testing. In *IEEE Symposium on Security and Privacy 2004*, Oakland, CA, May 2004.
- [12] J. Jung and E. Sit. An empirical study of spam traffic and the use of DNS black lists. In *IMC '04: Proceedings of the 4th ACM SIGCOMM conference on Internet measurement*, New York, NY, USA, 2004.
- [13] E. Kohler, J. Li, V. Paxson, and S. Shenker. Observed structure of addresses in ip traffic. *IEEE/ACM Transactions on Networking*, 14(6), 2006.
- [14] B. Krishnamurthy and J. Wang. On network-aware clustering of web clients. In *Proceedings of the 2000 ACM Special Interest Group in Communications SIGCOMM Conference*, 2000.
- [15] B. Laurie and R. Clayton. Proof-of-work proves not to work. In *Proceedings of the 2004 Workshop on Economics and Information Security*, 2004.
- [16] E. Levy. The making of a spam zombie army: Dissecting the sobig worms. *IEEE Security and Privacy*, 1(4), 2003.
- [17] John McHugh and Carrie Gates. Locality: A new paradigm for thinking about normal behavior and outsider threat. In *Proceedings of the 2003 New Security Paradigms Workshop*, Ascona, Switzerland, 2003. August 18 – 21, 2003.
- [18] J. Mirkovic, G. Prier, and P. Reiher. Attacking ddos at the source. In *ICNP '02: Proceedings of the 10th IEEE International Conference on Network Protocols*, Washington, DC, USA, 2002.
- [19] K. Plöbl, H. Federrath, and T. Nowey. Protection mechanisms against phishing attacks. In *Proceedings of the second annual conference on Trust, Privacy and Security in Digital Business*, volume 3592 of *Lecture Notes in Computer Science*, August 2005.
- [20] The Spamhaus Project. Zen blocklist. Available at <http://www.spamhaus.org/zen>, Fetched on January 29th, 2007.

- [21] M. Rajand, J. Zarfoss, F. Monrose, and A. Terzis. A multifaceted approach to understanding the botnet phenomenon. In *Proceedings of the 2006 ACM Internet Measurement Conference*, 2006.
- [22] A. Ramachandran, N. Feamster, and D. Dagon. Revealing botnet membership using DNSBL counter-intelligence. In *Proceedings of the 2006 USENIX workshop on steps for reducing unwanted traffic on the internet (SRUTI)*, 2006.
- [23] Bleeding Edge Threats. Bleeding snort ruleset. Available at <http://www.bleedingsnort.com/index.php/about-bleeding-edge-threats/all-bleeding-edge-threats-signatures/>, Fetched on January 29th, 2007.
- [24] P. Walt. Agencies feel botnets' light footprint. *Government Computer News*, January 2007.