

EMBER: A Global Perspective on Extreme Malicious Behavior*

Tamara Yu

Richard Lippmann

James Riordan

Stephen Boyer

MIT Lincoln Laboratory
244 Wood Street, Lexington, MA 02420
{ tamara, lippmann, james.riordan, boyer }@ll.mit.edu

ABSTRACT

Geographical displays are commonly used for visualizing wide-spread malicious behavior of Internet hosts. Placing dots on a world map or coloring regions by the magnitude of activity often results in cluttered maps that invariably emphasize population-dense metropolitan areas in developed countries where Internet connectivity is highest. To uncover atypical regions, it is necessary to normalize activity by the local computer population. This paper presents EMBER (Extreme Malicious Behavior viewer), an analysis and display of malicious activity at the city level. EMBER uses a metric called Standardized Incidence Rate (SIR) that is the number of hosts exhibiting malicious behavior per 100,000 available hosts. This metric relies on available data that (1) Maps IP addresses to geographic locations, (2) Provides current city populations, and (3) Provides computer usage penetration rates. Analysis of several months of suspicious source IP addresses from DShield identifies cities with extremely high and low malicious activity rates on a day-by-day basis. In general, cities in a few Eastern European countries have the highest SIRs whereas cities in Japan and South Korea have the lowest. Many of these results are consistent with news reports describing local cyber security policies. A simulation that models how malware spreads preferentially within cities to local IP addresses replicates the long-tailed distribution of city SIRs that was found in the data. This simulation result agrees with past analyses in suggesting that malware often preferentially spreads to local regions with already high levels of malicious activity.

Categories and Subject Descriptors

H.5 [Information Interfaces and Presentation]: User

*This work is sponsored by the United States Air Force under Air Force Contract FA8721-05-C-0002 and by the National Security Agency/Central Security Service (No. U/OO/NTC/0701). Opinions, interpretations, conclusions and recommendations are those of the authors and are not necessarily endorsed by the United States Government.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

VizSEC '10, September 14, 2010, Ottawa, Ontario, Canada.
Copyright 2010 ACM 978-1-4503-0013-1/10/09...\$10.00.

Interfaces; K.6.5 [Management of Computing and Information Systems]: Security and Protection—*Invasive software (e.g., viruses, worms, Trojan horses)*

General Terms

Measurement, Design, Security

Keywords

Visualization, malware, world map, city infection rate, normalized metric, long-tail distribution, situational awareness, computer network attacks

1. INTRODUCTION

World maps are commonly used for visualizing widespread malicious behavior of Internet hosts. However, their utility has often been dismissed by security and visualization experts. For example, some existing displays geo-locate malicious IP addresses and show them as separate dots on a world map. These displays may indicate that there is an overall threat, but they are overly cluttered with individual—often overlapping—dots, rendering it difficult to interpret regional differences in the threat. In addition, these displays often only indicate the locations of major population centers in developed countries because these are the locations of the greatest numbers of hosts in the world that can be compromised or used for malicious purposes. Even when these displays provide a heat map that indicates the local density of infections instead of individual dots, it is still difficult to determine if a region has more or fewer infections than expected for the local computer population.

Despite its limitations, a world map is an intuitive representation for global situational awareness. It is accessible to non-experts and well-suited for depicting the geopolitical context of cyber events. A geographical view can be used to identify threats that target specific regions, employ language or culture-specific social engineering, or exploit localization or pirated software. It can also be used to assess “network hygiene” of different regional ISPs. Some local organizations or ISPs explicitly engage in or harbor criminal activity, while others prevent, block, and rapidly detect and eliminate malicious activity. Furthermore, a geographical view can be used to explore the effectiveness of different policies on cyber security. For example, Japan’s Cyber Clean Center [11] is a government initiative that collaborates with ISPs to expunge botnets by alerting users when their computers become infected and providing instructions for

cleanup. China operates one of the most pervasive Internet filtering systems, dubbed “the Great Firewall”, to censor political dissent. Other countries such as Iran, Syria, Burma, and Saudi Arabia also practice heavy filtering on politically or socially sensitive materials [12]. South Korea has in recent years enacted the Cyber Defamation Law, which requires users to submit identification when posting on large websites, and the Anti-Piracy Law, which shuts down sites and suspends user accounts when copyright violations occur [13]. We created EMBER, the Extreme Malicious Behavior viewer, to provide high-level global situational awareness. Among other uses, EMBER makes it possible to study the effects of the above policies, targeted threats, and the spread of malware.

EMBER provides a geographical view of locations in the world with extremely high or low malicious activity, where the malicious activity level is normalized for the number of potential victim or malicious machines in each region. This paper presents the techniques behind the EMBER display, as well as its application to a large global dataset of firewall and intrusion logs provided by the SANS Internet Storm Center’s DShield project [1]. This work makes several contributions:

1. We demonstrate the application of a metric called the Standardized Incidence Rate (SIR) for identifying regions with extreme levels of malicious activity normalized by the local computer population. We demonstrate that the number of available computers in a region can be estimated from publicly available data sources.
2. We compare malicious activity across cities by SIR. Both highly infected cities and well-protected cities can be clearly identified on EMBER’s world map. The visualization reflects the varying infection rates of threats rather than the general distribution of available victims (computers) in the world.
3. Using the DShield dataset, we have identified several Eastern European cities as highly malicious. On the other hand, cities in Korea, Japan and Australia are consistently ranked as the best protected.
4. We observe that the distribution of SIRs for cities has a long tail similar to a power law. Simulations suggest that this is the result of malware preferentially spreading to regions with already high levels of malicious activity.

In the remainder of the paper, Sections 2 and 3 explain the EMBER approach and user interface. Section 4 presents a case study using the DShield dataset. The extreme cities, infection duration, and the distribution of city SIRs are reviewed in depth. In addition, a simulation is presented showing how long-tailed SIR distributions may be generated by malware preferentially spreading to local IP addresses. Section 5 surveys related work in characterizing and visualizing global malicious activity. Sections 6 and 7 provide a summary of the work and address future directions.

2. APPROACH

It is common to characterize malicious activity or malware distribution by network or Autonomous System. This approach provides insights on malware propagation strategies and network targets. EMBER takes an alternative perspective based on geopolitical divisions rather than routing boundaries. Our display can be used to show countries, cities, language or ethnic regions actively targeted or avoided by malware. It can also be used to correlate malware distribution with local cyber policies regarding piracy, indecency, antidefamation and censorship.

EMBER provides a geographical view of cities in the world with extremely high or low malicious activity. We choose to examine aggregated statistics by city because an analysis by country is often too coarse and because cities often delineate the range of organizational networks, ISPs, law enforcement domains, and traffic monitored by a single authority to detect and prevent malicious activity. To fairly compare cities of different sizes and levels of technological advancement, we use a city’s population and Internet penetration rate to estimate its network size (i.e., the overall host population), then normalize detected malicious hosts by overall host population. The result is a normalized measure of detected malicious activity or malice in the city that is as unbiased as possible with existing data sources.

2.1 IP Geo-Location

To analyze malicious activity by cities, detected malicious IP addresses are geo-located to determine their hosts’ physical locations. The free MaxMind GeoLite City [5] database provides country, state/region, city, latitude, longitude, and Autonomous System Number for IP addresses. Table 1 provides a breakdown of the post-processing results for IP addresses collected on a typical day, where 30,398 out of 603,546 IP addresses, or approximately 5%, could not be geo-located to a city and had to be discarded. Of the remaining 95%, a subset may be geo-located to the wrong cities. [20] presents the accuracy of GeoLite City on a city level, which averages about 75% but varies widely by country. For countries with low accuracy, mis-location of IP addresses can significantly distort their cities’ infection rates. To remedy this problem, one could use other (commercial) geo-location services that provide more complete and more accurate city coverage. Alternatively, one could adjust the infection rate to compensate for low geo-location accuracy. We did not apply either option to correct for mis-located IP addresses but intend to do so as part of our future work.

2.2 City Host Population

Fair comparison of cities necessitates normalization of a city’s malicious hosts by its total available hosts, but estimating a city’s host population proves to be a difficult problem. One way to get an estimate is to actively probe networks in a city and count the hosts that respond. This is often not practical because many networks use firewalls to block suspicious or non-essential traffic, preventing others from discovering their internal structures.

Another way is to use Internet address registries provided by organizations such as the Internet Corporation for Assigned Names and Numbers (ICANN). This is simple but insufficient because registries could be out of date and reserved

	Justification	Mal. IPs	Subtotal
Discarded	No city	30,398	108,680
	No population	77,864	
	Low IPR	418	
Retained	Adjusted IPR	15,717	494,866
		479,149	
Total		603,546	

Table 1: Approximately 80% of the malicious IP addresses recorded by DShield on 15 February 2010 are included for analysis and display. We discard any IP address that cannot be ge-located to a city, that is in a city with no population data, or that is in a country with extremely low Internet Penetration Rate (IPR).

address blocks may be under or over-utilized.

It is common to infer host population size by counting unique IP addresses observed in traffic passively. This approach has several problems as well. Some networks use NAT devices that translate between internal and external IP addresses, often masking a large number of hosts with a small number of public-facing addresses. Over time, IP addresses of hosts may change due to DHCP churn, old hosts going offline, or new ones coming online. Estimation is further complicated by ISPs that route traffic through proxies. For example, AOL is known to route all North American traffic through a few gateways [20], and satellite Internet providers have to send their traffic through satellite terminals.

Another option is to infer host population by gathering Domain Name System (DNS) records and query statistics from DNS servers. This could be done at a coarse level by examining only the root DNS servers, or a fine-grain level if one has access to local DNS servers. A DNS “census” may allow us to estimate the number of hosts in a network. However, such detailed information is hard to obtain, and not all addressable hosts are in DNS records. It may also be possible to measure a network’s “NAT factor” by characterizing the DNS traffic it generates. While it is difficult to learn the host populations of all cities using only DNS data, this approach may complement others.

We decided to estimate the host population of a city using open source population and Internet usage data:

$$N_{city} = Population_{city} \times InternetPenetrationRate_{city}$$

This approach relies on population data, which is easy to obtain and should be reasonably accurate for medium to large cities. Given the percentage of the population that uses the Internet, we can approximate the number of computers connected to the Internet assuming a 1:1 ratio of Internet users to computers. Although there are complicating factors, such as users who own multiple computers and those who access shared computers (e.g., Internet cafes, schools, libraries), as well as servers, they do not appear to bias the results significantly for a particular country or city. EMBER currently uses population data from GeoNames [2], a site that provides regularly updated information on a comprehensive list of cities. Internet World Stats [3] provides Internet penetration rates by countries compiled from various online mea-

surement services and local communications providers and regulators. According to their website, some of its major data sources include Internet market research firms such as Nielsen and GfK Group, which use proprietary tools and techniques to monitor and measure global Internet activity. The site also uses data from the United Nations’ International Telecommunication Union, which is responsible for reporting and organizing efforts on worldwide connectivity.

2.3 Standardized Incidence Rate

The EMBER display uses a normalized metric called the Standardized Incidence Rate (SIR), expressed as the number of malicious machines (i.e., unique IP addresses) for every 100,000 actual machines that could be infected in a city:

$$SIR_{city} = \frac{IPs_{city}}{N_{city}} \times 100,000$$

This results in an integer ranging from zero if no hosts are malicious to 100,000 if all hosts are malicious. Typically, a small percentage of hosts are infected and the SIR is a number from 0 to 1,000. For example, if 1% of all hosts in a city were detected as being engaged in malicious activity, the SIR for that city would be 1,000 which is 1% of 100,000. If 0.1% of the hosts were engaged in malicious activity, the SIR would be 100. SIR has traditionally been used in cancer statistics [17]. The Microsoft Security Intelligence Report uses a similar metric called Computers Cleaned per Mil (CCM) that represents the number of computers cleaned per thousand executions of the Malicious Software Removal Tool (MSRT) [21]. These metrics all represent infection rates. In the context of EMBER, SIR alerts users to cities with disproportionately high or low levels of malicious activity that warrant further investigation.

2.4 Adjustments

To compensate for data flaws and statistical variability, we introduce a few adjustments to the computation of city SIR.

2.4.1 Internet Penetration Rate

Since the Internet penetration rates are compiled from a variety of sources and there is no ground truth, we examine how potential errors in the penetration rate may affect city scoring. We first assume the city penetration rate to be the same as the country penetration rate. For every rate reported, we calculate the SIR of a hypothetical city with 100,000 human population and 1,000 malicious hosts. Fig. 1 shows the resulting SIR scores as the upper red circles. We observe that SIR scores are highly sensitive for countries with low penetration rates; that is, cities in these countries could be substantially overrated or underrated if their penetration rates err slightly one way or another. Furthermore, undeveloped and developing countries often have the fastest growth online and the largest technological disparity between urban and rural areas. We expect city penetration rates to be higher than the national average for these countries. Finally, any penetration rate below 0.01 is presumed too small to be reliable. As of February 2010, There are 30 out of 241 countries and autonomous regions that have such low rates, including many small island nations, select African nations (e.g., the Democratic Republic of Congo, Ethiopia and Liberia), select Southeast Asian nations (e.g., Bangladesh, Myanmar and Cambodia), and North Korea.

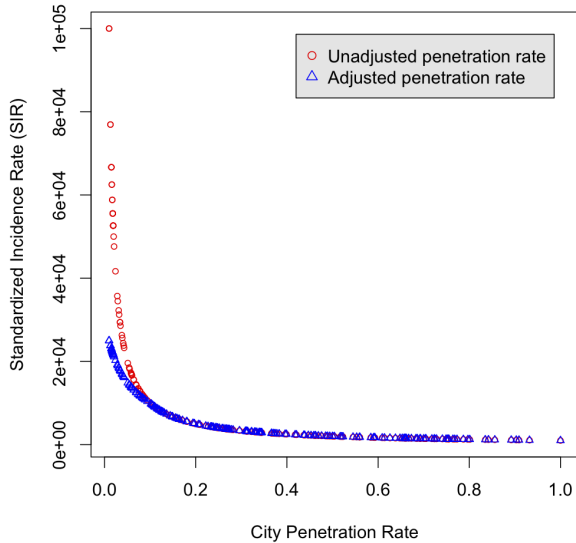


Figure 1: The Internet penetration rates of 241 countries in February 2010 and the resulting standardized incidence rates (SIR) of a hypothetical city.

We exclude cities in these countries from the display and analysis.

To reduce the effect of inaccurate city penetration rate estimates and account for variability caused by inaccurate rates, the following formula is used to derive city penetration from country penetration rates:

$$rate_{city} = \begin{cases} rate_{country}, & \text{if } rate_{country} \geq 0.1 \\ \frac{2}{3}rate_{country} + \frac{1}{30}, & \text{if } 0.1 > rate_{country} \geq 0.01 \\ 0, & \text{otherwise} \end{cases}$$

For most countries, the first case in which the city penetration rate is equivalent to the country penetration rate applies. The second case gives countries with relatively low penetration rates a multiplier, gradually diminishing from 4x to 1x. The third case discards countries with extremely low penetration rates below 0.01. When the same hypothetical city is scored using the adjusted rates, the resulting SIR curve has a much smaller range of variation as shown by the lower blue triangles in Fig. 1. The adjustments introduced in the formula affect a small percentage of hosts observed every day. Table 1 shows that on a typical day, less than 0.1% of the IP addresses (418 out of 603,546) are discarded because their countries' Internet penetration rates are too low. Only 2.6% (15,717 out of 603,546) originate from countries where the Internet penetration rates need correction.

2.4.2 Population Thresholds

EMBER allows a user to select cities with the highest or lowest SIR values for display and inclusion on lists that represent cities with extreme SIR values. It would be misleading to simply order all cities by SIR and display them with their rank. Such a simple approach does not account for the differing effects of changes in the underlying number of infected hosts in each city on the SIR. The SIR scores for cities with low host populations below 100,000 change much more dramatically when a few hosts are infected or patched than

SIR scores for cities with populations well above 100,000 because the SIR is a ratio with the city population in the denominator and the number of infected hosts in the numerator. For example, the SIR for a city with 10,000 potential infected hosts is 100 when 10 hosts are infected. This SIR doubles from 100 to 200 when ten more hosts are infected. The SIR for a city with 1 million hosts is 100 when 1000 hosts are infected. The SIR increases by only 1% when ten more hosts are infected.

The problem of SIR variability is addressed by only including cities with a computer population above a threshold on EMBER high and low SIR lists and displays. There is a tradeoff in selecting the population threshold. If the value is too high, many cities will be excluded and this could potentially allow knowledgeable adversaries to avoid detection by targeting cities with low computer populations. If the threshold is too low, it could lead to a misleading display including cities that had high or low SIRs by chance. In the current version of EMBER, we only display and list cities with populations greater than 100,000 for both extreme SIR value displays and require cities to have at least 20 infections before they are included on the lowest-SIR list. The city computer population limit of 100,000 affects primarily the list of high-SIR cities. First, it reduces the sensitivity of high-SIR scores to changes caused by adding or removing a few infected hosts per city. Most of these cities on the high-SIR list have SIR scores above 300-500. Since the populations of these cities is restricted to be above 100,000, the number of infected hosts in these cities will be above 300-500 and adding or removing a few infected machines (e.g. < 10) will vary the final SIR score by less than ± 10 which is only 2% to 3% for SIR scores ranging from 300 to 500. The population limit of 100,000 also focuses the analysis on larger cities where human population estimates, computer penetrations estimates, and IP address mapping are all more accurate. It also keeps the resulting map less cluttered and easier to interpret when many cities are displayed.

The requirement that there needs to be 20 malicious IP addresses in a city before that city is included on an extreme SIR list affects primarily the list of low-SIR cities. If there are 20 malicious IP addresses in a city, then adding or removing one infected host varies the SIR by less than $\pm 5\%$. When cities have low SIR values ranging from 1 to 5, this limit indirectly restricts the population of cities on the low-SIR list to be from 400,000 to 2,000,000. It is thus more restrictive in terms of the smallest city population allowed on the low-SIR list. Although these choices appear to provide relatively consistent results across days, further statistical analyses of SIR day-to-day variability is planned to validate and possibly improve these thresholds.

2.4.3 SIR Ranking

When presenting the lists of cities with extreme SIR values, it would be misleading to simply order cities by SIR. This approach suggests to users that the absolute ranks on the high and low-SIR lists are too important even though many other cities on the list have SIRs that may be statistically equivalent. We address this problem by grouping cities with comparable SIRs and adjusting the display and the list to reflect ties in the ranking.

There are several ways to implement grouping. We can cluster cities with similar SIR values and order the clusters by cluster means to produce the ranking. This approach requires a priori knowledge of the number of clusters or ranks. We can also use a statistical approach to test whether two cities' SIRs are equivalent by comparing confidence intervals. Confidence intervals are determined using per-city interdecile ranges of SIR variability measured over the past 10 days. For the high-SIR list, the ranking procedure begins by identifying the city with the highest SIR as the prototype for rank 1 and assigns the same rank to any city with $SIR \geq SIR_{max} - \frac{R}{2}$, where R is the median interdecile range. Then, from the remaining unranked cities, the city with the highest SIR is selected as the prototype for the next rank. Comparable cities are then found for that rank. This process is repeated until all cities are ranked. The ranking procedure for the low-SIR list is analogous. This approach is simple to implement and interpret.

3. EMBER

EMBER aggregates malicious activity by city and displays a circle of variable size and color (e.g. glowing ember) for each city that indicates the normalized malicious activity in that city. The interface (Fig. 2) consists of a calendar on the right to select a day for analysis, a panel on the upper center to select metrics and parameters, a central world map to display extreme cities, a table on the left for details of the cities, and a histogram of displayed scores on the bottom.

Four different metrics are available for selecting cities to display in the world map. Users may set a threshold for cities shown using the upper middle control panel in the application. If thresholded by rows, only cities with the N highest (or for SIR also lowest) values are displayed. If thresholded by ranks, all cities with the N highest ranks are displayed. Cities are displayed as dots in the map. Dot colors are proportional to the metric value, ranging from orange for the highest to blue for the lowest. Dot sizes are proportional to the rank. The four supported metrics are described below.

Alerts (total number of alerts per city). This metric measures the total number of alerts triggered by IP addresses from each city. It shows which cities are responsible for the greatest number of alerts across the world. It does not distinguish between a city with a single malicious host that scans many other hosts across the Internet and a city with many malicious hosts that contact only one host each. It also is not normalized for the city population. This metric is useful for identifying high-volume scanners.

IPs (total number of unique malicious IP addresses detected per city). This metric measures the number of unique IP addresses that have been detected performing malicious activity in each city. It shows the cities that contain largest numbers of compromised hosts. Because this display is not normalized for the size of each city, the display tends to focus on cities with the largest populations where the potential numbers of infected computers are also largest. This metric is useful for assessing the spread of threats.

High Standardized Incidence Rate. This display shows the cities with the highest SIR scores. As described above, the SIR is expressed as the number of malicious IP addresses

for every 100,000 actual machines that could be infected in a city. This metric is useful for identifying cities with higher-than-expected levels of infection. Cities with high SIRs may be allowing malicious activity, or they may be more heavily targeted than other cities. High SIRs could also be attributed to poor network hygiene or unusual sensor coverage.

Low Standardized Incidence Rate. This display shows the cities with the lowest SIR scores. This metric is useful for identifying cities with lower-than-expected levels of infection; that is, they prevent malicious activity well, are less targeted than other cities, or are not covered by malicious activity sensors as well as other cities.

4. CASE STUDY

EMBER is a general-purpose tool that can be applied to any dataset containing lists of IPv4 addresses. This section describes a case study of the EMBER display and analysis using a large dataset collected over several months from the SANS Internet Storm Center. We will describe the data source, our observations and interpretations.

4.1 Data Source

The dataset used in the case study comes from the SANS Internet Storm Center's DShield project. Each day, DShield collects millions of intrusion detection logs from firewalls and intrusion detection systems all around the world, and publishes on their website a daily feed of source IP addresses observed to exhibit malicious behavior. This may include hosts engaged in phishing, serving exploits or malware, sending SPAM emails, hosting scam pages, serving as repositories for pirated software or pornography, serving as command and control centers for botnets, or serving as botnet zombies. According to a 2003 study [31] using the same source, scans were the dominant class of malicious activity detected, with port 80 worms such as Code Red and Nimda accounting for 20-60% of the intrusion attempts each day. The Internet landscape has changed since the study. However, because the publicly available part of the dataset only contains the source of malicious activity and not the type of activity, we have no further insight on the threat types represented in this dataset without access to the full dataset.

There may be additional limitations to the dataset. The data could contain many false alarms. Even though the data is from sensors around the world, it may not provide complete coverage of all regions due to insufficient log contribution, higher level filtering, or blocking of IP addresses or ranges. It would be useful to verify the accuracy of the data used to detect malicious activity either by correlating across different data sources or by active verification. In addition, we cannot easily account for spoofed IP addresses or gateway devices that serve many hosts using only one IP address. Finally, we do not have any data from the IPv6 address space. Despite its limitations, the dataset serves as a good test case, and has led to some interesting findings.

4.2 Extreme Cities

The DShield data indicates that the SIR is much higher in Eastern Europe countries including Moldova, Romania, Macedonia, and Bulgaria than in any other part of the world.

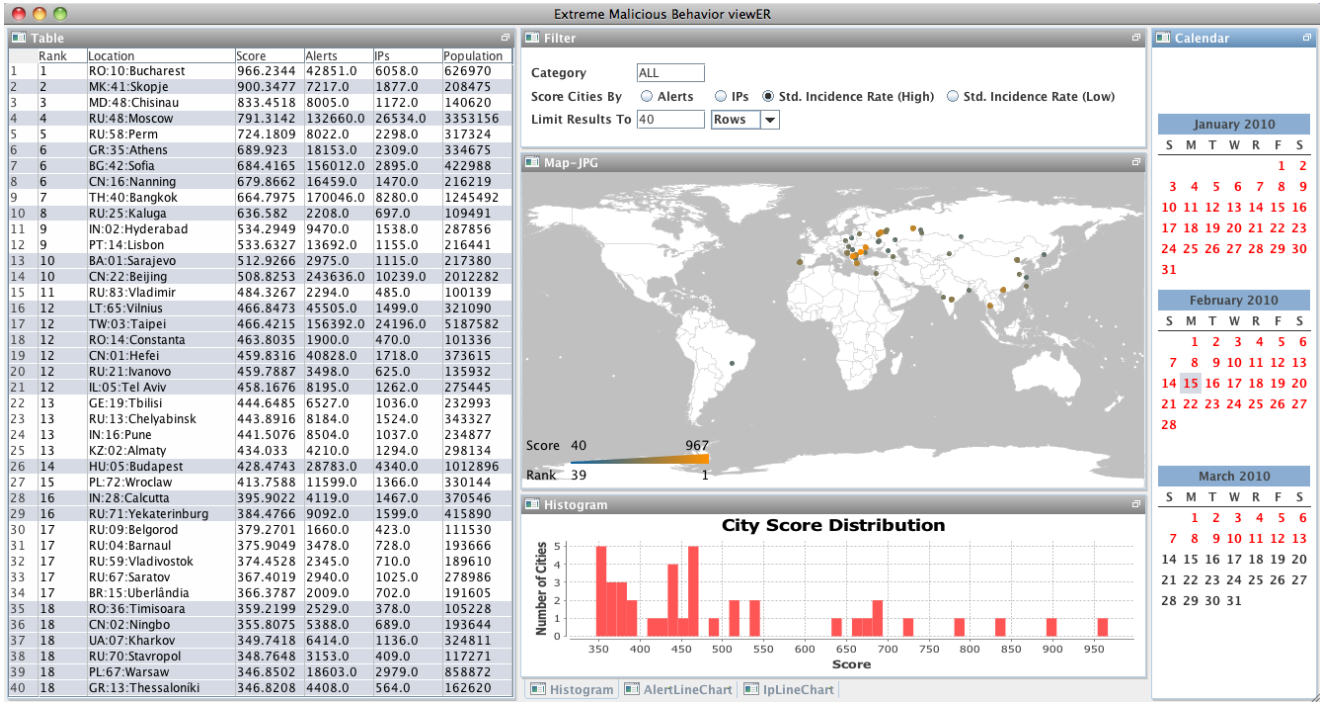


Figure 2: The EMBER Interface showing 40 cities with the highest SIR scores on 15 February 2010.

Fig. 3(a) offers a snapshot of the highest-SIR cities on a typical day. Fig. 4, which shows the range of SIR scores for each city that has ranked in the top 10 for at least ten days between 1 January and 13 March 2010, further demonstrates that the list of high-SIR cities is quite consistent over time. In this figure, boxes represent interquartile ranges of city SIR scores, middle dark bands represent median values, whisker ends represent upper and lower bounds, and circles represent outliers. This plot illustrates that most malicious cities have relatively narrow SIR ranges except Athens and Bangkok, which tend to flare up occasionally. Bucharest and Chisinau are clearly the worst offenders with SIRs consistently near 1000, surpassing lower-ranked cities by almost a factor of 2. The same plot also shows the number of days each city has been in the top 10 (see the grey dots connected by lines). Not surprisingly, Bucharest and Chisinau are among the six cities that have made the top-10 list almost every day. News reports suggest that malicious activity could be related to organized criminals in these countries, where cyber law enforcement is lax and economic incentive for the “over-educated and under-employed specialists” is strong [18]. Other cities with high levels of malicious activity include Moscow, and Nanning, China.

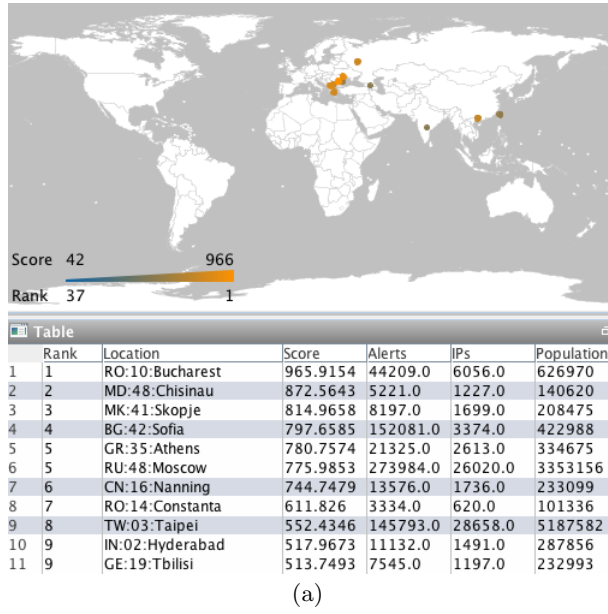
Top cities on the high-SIR list are generally well differentiated. Their color gradation is clearly visible on the map in Fig. 3(a), and the table contains very few ties. On the other hand, for low-SIR the top cities are not as clear-cut. Fig. 3(b) shows the typical one-day result, where nearly a dozen cities are tied for the lowest level of malicious activity. Most notably, cities in South Korea consistently have the lowest SIRs. This finding contradicts the Microsoft Security Intelligence Report [21], which ranks South Korea 5th in locations with the highest infection rates. However, because the CCM

metric used by the report actually measures the number of computers cleaned per 1,000 MSRT executions, this suggests that South Korea in fact cleans up computer infections very quickly, likely contributing to its low SIR scores. The low level of malicious activity in South Korea may also be attributed to its strong cyber antidefamation laws, Internet content controls, and personal identification requirements on large web sites [13].

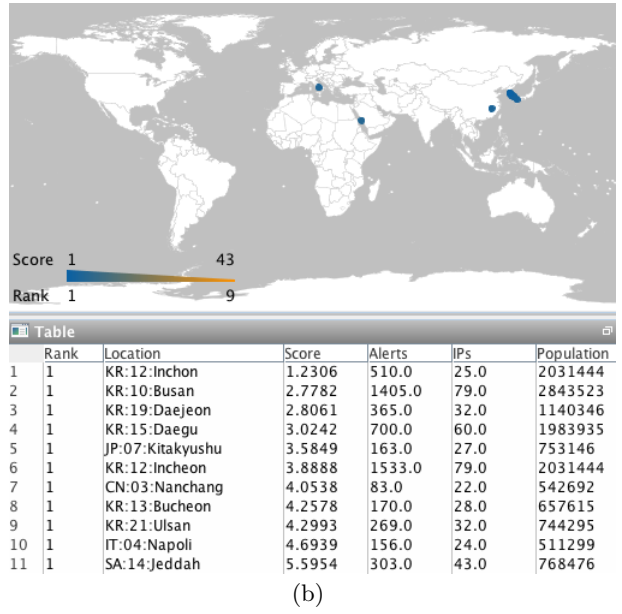
4.3 Infection Duration

Evaluating cities on a day-to-day basis (without differentiating IP addresses by infection duration) enables detection of sharp changes, such as a botnet launching a coordinated attack. It is also useful to analyze the longevity of infected IP addresses to better understand the long-term dynamics of infection and cleanup in these cities. For each IP address observed to behave maliciously on a given day, we examine its history to determine the date of initial infection (d). Since hosts could be turned off and malware could be dormant, we cannot expect infected hosts to be active every day. An IP address is considered a new infection on date d only if it is absent from the previous t days. The longevity of an infected IP address, or infection duration, is measured as the number of days since d .

Determining the threshold t , or the “grace period,” requires some experimentation. On the one hand, because IP addresses are often dynamically assigned, the longer the grace period is, the less likely that two occurrences of an IP address separated by the grace period actually address the same host. On the other hand, a grace period too short would cause current infections with low frequency of activity be mislabeled as new. We tested t values ranging from 3 to 9 days and found that the new infection count would only



(a)



(b)

Figure 3: EMBER displays cities with extremely high (a) or low (b) malicious activity daily with SIR scores, ranks, total alerts, total malicious IP addresses and host populations for 2 March 2010.

get marginally smaller as the grace period increases. This suggests that the majority of the malicious IP addresses detected each day are short-lived and that it is reasonable to set t to 5 without significantly over-counting new infections.

Fig. 5 shows the complementary cumulative distribution function of infection duration for all IP addresses detected on 13 March 2010. A high churn rate in IP addresses is the presumed cause of the short durations of most malicious IP addresses. Approximately 75% of the addresses on this day are new infections, compared to 10% that have existed for over a week. This distribution is consistent from day to day. Furthermore, our observation agrees with the findings of [27], which concludes through the analysis of a botnet takeover that DHCP churn can lead to gross over-estimation of botnet drone population using unique IP address count (by as much as 36.5% in a day). There may still be other factors that also contribute to the high turnover of malicious IP addresses. [28] suggests that botnet servers and phishing sites are often cleaned up or abandoned within a few days. Some malware employ tactics such as “fast-flux” and spoofing to change their IP addresses constantly to evade detection on watchful networks.

4.4 Persistent Infections

As mentioned before, a small percentage of IP addresses observed each day are persistent (e.g., over a week). They may be individual infected hosts that are truly persistent or gateway devices that mask large local area networks. The former implies that persistent IP addresses can be used to find rogue cities that neglect malicious activity whereas the latter implies false positives. We explore this idea by computing the SIR on persistent IP addresses exclusively. Fig. 6(a) displays cities with the highest *persistent-SIRs*. Moscow and Athens are notably absent from the group even though they have high SIRs overall (Fig. 3(a)). Tbilisi of Georgia and

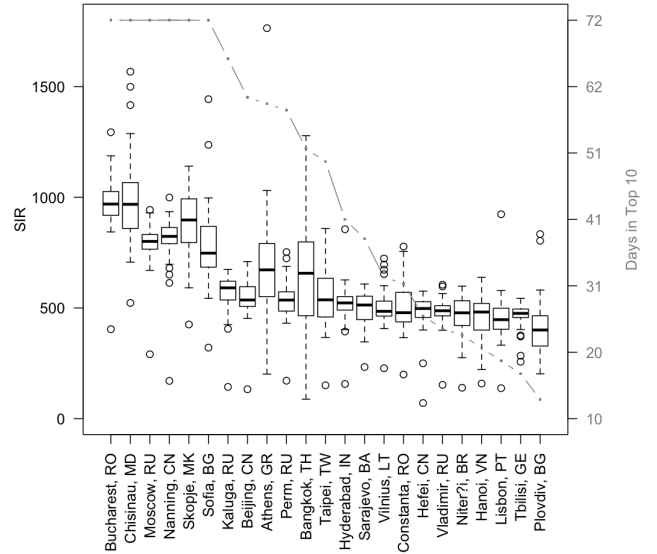
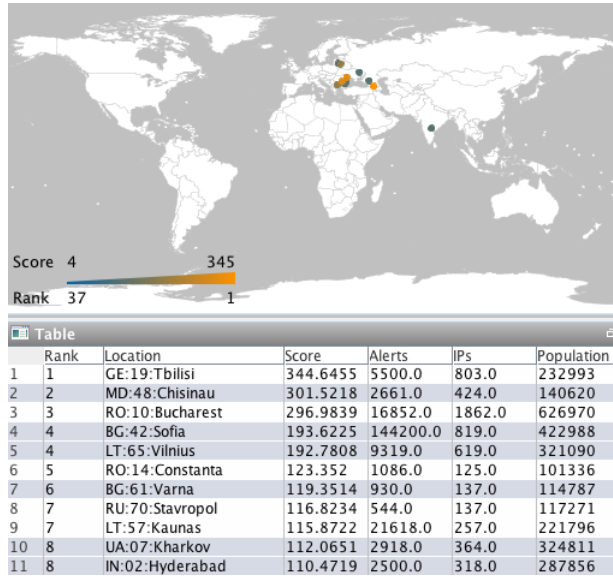
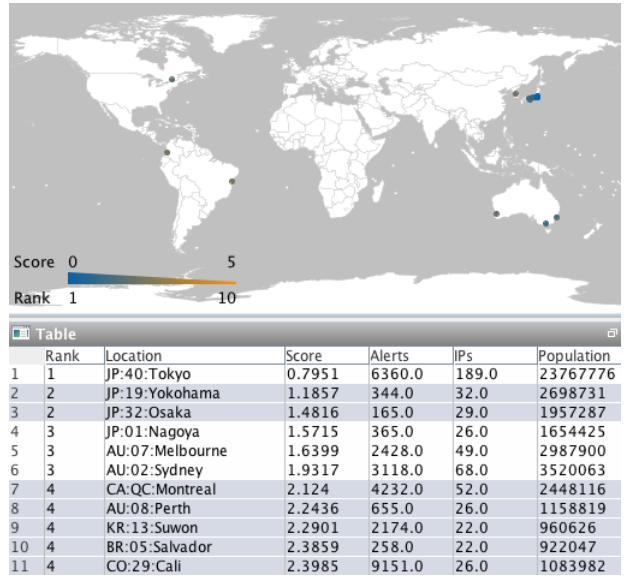


Figure 4: A boxplot of SIRs for extreme malicious cities (by the High-SIR metric) from 1 January to 13 March 2010. Cities shown have ranked in the top 10 for ten or more days during the observed period. The grey line indicates the number of days each city has ranked in the top 10.



(a)



(b)

Figure 6: EMBER displays cities with extremely high (a) or low (b) *sustained* malicious activity, that is, activity caused by persistent malicious IP addresses with uptime over 7 days. The date selected is 1 March 2010.

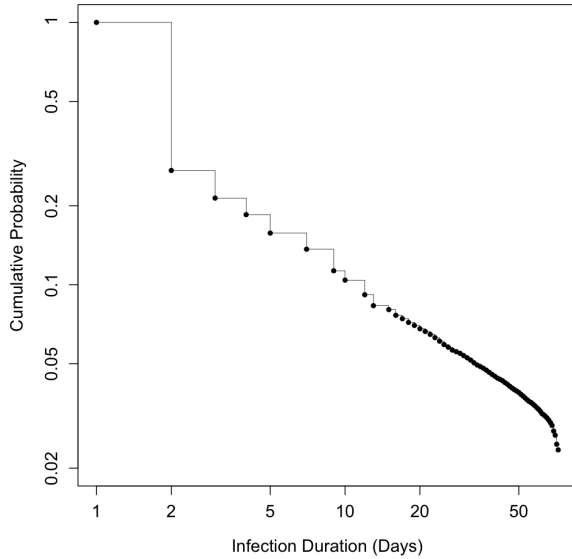


Figure 5: Complementary cumulative distribution function for infection durations of IP addresses on 13 March 2010. The history of address appearances goes back to 1 January 2010 (72 days).

Kharkov and Lviv of Ukraine do not have the highest SIRs overall, but their presence on the rogue list suggests that a significant portion of their malicious IP addresses are persistent. Finally, Chisinau, Moldova and Bucharest, Romania are the worst cities either way. A small, randomly selected subset of persistent IP addresses have been manually cross-validated with the email reputation system SenderBase [16], Malware Domain List [4], which tracks domains used for distributing and controlling malware, and Google Safe Browsing [24], which identifies malicious websites. Many persistent IP addresses in cities with the highest persistent-SIRs are either identified as spammers in SenderBase or collocated with spammers. Otherwise, persistent IP addresses are rarely flagged by Malware Domain List or Google Safe Browsing. The cities from EMBER with the five highest SIRs were also included in the map generated by the FIRE system on their associated web site (see [28]) on 5 April 2010.

Through identifying cities with low persistent-SIRs, we hope to find cities that are well-protected. Fig. 6(b) displays cities that have the lowest persistent-SIRs on a given day. This picture is markedly different from Fig. 3(b): The financial and technological centers of the world—Tokyo, Sydney and New York—now rank as best protected. South Korean cities fall out of the list because their malicious IP address counts are below the threshold of 20 (see Section 2.4.2). Unlike cities with high persistent-SIRs, these results are more difficult to interpret. While some of the persistent IP addresses are confirmed malicious by the above data sources, many others appear to be gateway devices for large corporations such as Akamai based on registry records (i.e., *whois*). These gateways can artificially reduce SIRs because many malicious IP addresses can be mapped to a single IP address. Further analysis would be required to determine if these cities are best-protected because they have effec-

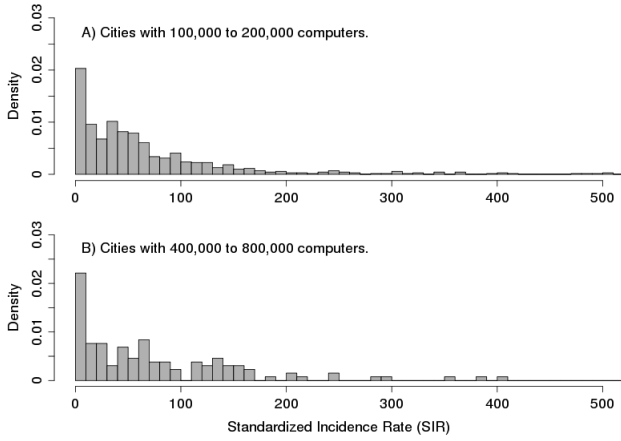


Figure 7: Histograms of SIRs on one day for A) Cities with computer populations between 100,000 and 200,000 and B) Cities with computer populations between 400,000 and 800,000.

tive cyber policies, the technological know-how, or dense concentration of gateway devices for large corporations and governmental networks.

4.5 Standardized Incidence Rate Distribution

The standardized incidence rate (SIR) has been used extensively in the past to analyze health statistics and determine if any city, state or country has higher or lower incidence rates than expected [9]. When analyzing cancer rates, it is assumed that the probability that any individual has cancer (normalized for age) is the same across locations. When comparing SIR rates of cities, the overall SIR is computed across the overall population and deviations from this overall distribution are computed assuming that the number of cancer cases in a city has a binomial distribution with a mean equal to the overall SIR multiplied by the city or population divided by 100,000 (e.g. [23]). Before analyzing the DShield data we assumed that similar assumptions could be used to model computer infection rates. This assumption proved to be incorrect.

Fig. 7 shows histograms of SIRs for cities with computer populations ranging from 100,000 to 800,000 hosts. If all hosts had the same probability of infection, then these distributions would be approximately Gaussian-shaped with means near 80 and a standard deviation of roughly 7.5 in the upper histogram and 3.7 in the lower histogram. These distributions are clearly not Gaussian, but have long tails characteristic of data that is often modeled using a power law such as the human population of cities and the number of hyperlinks to web sites (e.g. [10]). The x axis in Fig. 7 was terminated to better show the distribution shapes. The long tails actually extend to above 1,000. The shapes of the distributions in Fig. 7 do not vary substantially as the city population sizes increase by a factor of four from 100,000-200,000 in the upper to 400,000-800,000 in the lower histogram. Although not shown, the distribution shape is similar for cities with computer populations above roughly 10,000.

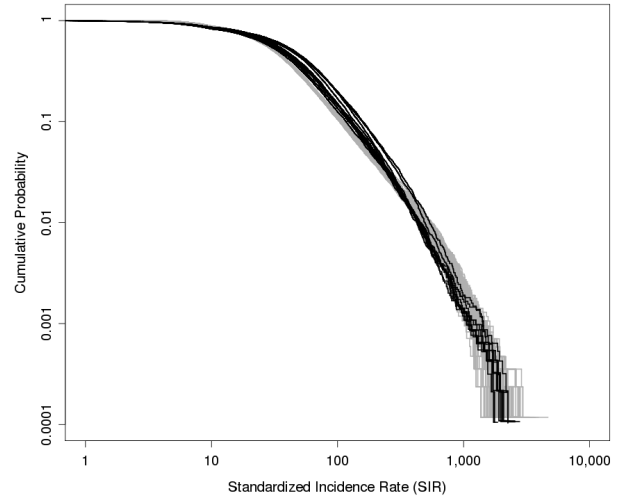


Figure 8: Complementary cumulative density function for SIRs of cities with computer populations greater than 10,000 for ten Wednesdays from 13 January to 17 March 2010. A separate black line is shown for each day. Gray lines are from simulations where new malware preferentially spreads to cities with most existing infections.

Fig. 8 provides another summary of the SIR distribution. It shows the SIR complementary cumulative density function (CCDF) for cities with populations greater than 10,000 using data from eleven Wednesdays from 13 Jan to 17 March 2010. Note that both axes in this plot are logarithmic. The CCDFs shown in black are representative of distributions for other weekdays. They show that SIR values range over more than three orders of magnitude and that roughly 40% of the cities have SIR values contained in a long tail extending above a value of 60 where the distribution tail appears to begin. Maximum likelihood estimates described in [10] averaged over these curves indicate that above a lower limit of roughly 60 the distribution falls off with a power law exponent of -2.52. Although a Kolmogorov-Smirnov statistical test described in [10] indicates that these SIR distributions do not exactly follow a power-law, in practice they have a long tail and similar characteristics to power law distributions as noted above.

4.6 Simulation to Explain SIR Distribution

The observation of SIR values that have an approximately power-law behavior suggests that infected computers on the Internet detected by DShield could be generated by preferential attachment models of how computer infections spread for recent worms and botnets. Preferential attachment models have been used, for example, to simulate the distribution of the number of Internet hyperlink connections on web sites (e.g. [22]). We explored a simple model that assumes that newly infected hosts in a city are more likely to be caused by locally infected hosts in the same city than by more remote hosts in more distant cities. This is consistent with the behavior of many recent worms and botnets which preferentially attempt to spread to computers that are geographically close because they either preferentially scan IP addresses in local class A and B networks or they take advan-

tage of IP address, email, chat, and other social networking information discovered on already-infected computers (e.g. [26, 7, 25]).

Our simulation is much simpler than many others (e.g. [26, 7, 25]), but is the only simulation we are aware of that has been used to model the number of hosts that are infected for worms and botnets on a per-city basis. The simulation models the number of hosts that are infected in 9,157 cities with computer populations above 10,000 that are the same as those used to produce the SIR distribution shown above. These cities contain a total of roughly 692 million hosts and the host populations used in simulations are the same as those used to compute SIR values.

The simulation begins by randomly and uniformly selecting across all hosts a certain number of infections. It is difficult to determine the correct number of initial infected hosts to use in this simulation because we are modeling all the different worms and botnets seen in DShield data, “hit-lists” and other approaches are often used to create multiple initial infections (e.g. [30]), and worms can be propagated physically via laptop and infected media as well as by scanning. In our simulations we started by infecting 1% of the total number of infected hosts seen in the DShield data. Uniformly spreading these over the entire host population preferentially places infected hosts in the larger cities.

After initialization, infections are added one at a time until the number of infected hosts estimated to be in the DShield data for one day (497,250) is reached. This number was set to the number of IP addresses seen in DShield data reduced by 37% to account for DHCP churn as suggested by the daily DHCP churn analysis in [28]. An infection is added with probability λ preferentially to cities with more infections by setting the probability of infection for hosts in each city proportional to the number of infections in that city. An infection is added with probability $1 - \lambda$ to all computers by setting the probability of infection to be uniform across all computers. This models malware that preferentially spreads to IP address that are geographically local. Values of λ ranging from 0.7 to 0.8 provide roughly similar results suggesting that 70% to 80% of worm spread in local. Over time, the effectiveness of scanning for new hosts in each city is reduced to account for scans that reach already-infected hosts.

The gray lines in Fig. 8 that range around the actual CCDF curves represent CCDF curves generated by 100 runs of this simulation with different random starting seeds. These CCDFs provide a good visual fit to the actual CCDFs and suggest that the simple generative model of worms and botnets preferentially spreading to nearby hosts in the same city may help explain the power-law like distribution of SIRS. The good fit to the actual SIR distribution requires values of λ above roughly 0.7, a correct target goal for total number of infected hosts, and cities with host populations that also have a long tail as in this data.

5. RELATED WORK

There have been several efforts to identify centers of extreme malicious activity. The FIRE system [28] automatically identifies rogue networks by monitoring the longevity of botnet, drive-by-download and phishing servers and scor-

ing Autonomous Systems (AS) based on the number of sustained rogue servers. FIRE generates a *mal score* by normalizing the total number of rogue IP addresses in an AS by a factor based on the number of /24 prefixes announced by that AS. EMBER differs from FIRE in that (1) the analysis is performed at the city-level instead of the AS-level, (2) it produces the top protected cities as well as the top rogue cities, (3) it does not differentiate types of malicious activities, and (4) it has been applied to all malicious IP addresses and the subset of sustained malicious IP addresses.

HostExploit has compiled lists of the best and worst hosting service providers using data drawn from eleven different sources [15]. It calculates per AS an *HE Index*, which considers the number of bad instances in the AS normalized by the size of the AS. Similar to EMBER, adjustments have to be made to reduce the HE Index’s sensitivity for small ASs. EMBER discards cities with less than 100,000 computers or less than 20 infections, whereas HostExploit uses a Bayesian ratio that moves the score of an AS toward the average as the size gets smaller. The report also provides top-10 lists by threat type (e.g., spamming, botnet command and control and malware distribution). Unlike EMBER, HostExploit does not consider statistical ties when ranking ASs by their HE Index scores.

The semi-annual Microsoft Security Intelligence Report [21] provides a global assessment of cyber threats using data collected from the Microsoft Windows Malicious Software Removal Tool (MSRT) running on computers worldwide. The report uses a metric called Computers Cleaned per Mil (CCM) for comparing infection rates across regions. CCM measures the number of computers cleaned per 1,000 executions of MSRT. This approach misses a class of high-risk systems that do not run MSRT, either because the feature has been disabled by malware or the software is unpatched or pirated. It also does not distinguish a host performing a thousand executions from a thousand hosts performing one execution each. These factors introduce a significant bias in the final assessment of a location’s level of malice or protection. The approach is also not extensible to non-Windows systems.

Several studies have used the DShield dataset to characterize cyber attack trends globally. Most recently, [32] presents an analysis of attack sources based on intrusion records collected between 2004 and 2006. It examines source IP address distribution across the IPv4 address space and concludes that a small set of address ranges contain the majority of the malicious IP addresses. However, the analysis does not account for the uneven allocation of IP address space. As a result, the top malicious domains found were major ISPs at that time.

A world map display is widely used to provide global cyber situational awareness. The DShield [1] website uses a map to show high-level attack statistics and port trends by continent. Shadowserver [6] tracks botnet activities and regularly posts maps to show geographic distributions of command-and-control servers and drones, where each dot represents the relative number of infected hosts at a location. A display from the FIRE system showing persistent malware sites with dots all of the same size is available on a website noted

in [28]. Team Cymru [29] uses a global heatmap to show density of Conficker infected IP addresses across regions. While these approaches are good for conveying the extent of malicious activity worldwide, it is difficult to determine any regional differences because dots or heatmap regions tend to concentrate around population centers, where there are also the greatest numbers of hosts to compromise.

Other visualization techniques have also been adapted for presenting and analyzing cyber activities. [19] monitors network traffic using a treemap, where traffic is sliced and diced by continent, country, AS, and IP address prefix, and block size is proportional to the traffic quantity within. A novel feature of the visualization is that it maintains the relative geographic positions of continents and countries in the treemap to facilitate detection of regional patterns. The visualization does not account for the extent different address blocks are utilized and does not normalize network activity by computer population. STARMINE [14] integrates geographical, temporal and logical visualizations for monitoring cyber attacks. In a 3-D environment, it displays an IP Matrix, a world map, and a volume-over-time chart on three orthogonal planes. Lines connect attacks' originating IP addresses on the IP Matrix to their locations on the world map. While this approach allows users to visualize attacks from different perspectives simultaneously, the display quickly becomes too cluttered to read as data volume increases.

Many simulations have been developed to explore worm and botnet propagation including recent studies that explore the spread of worms that preferentially spread to local IP addresses (e.g. [26, 7, 25]). To the best of our knowledge, none of these simulations have modeled infection rates worldwide across many cities and considered the resulting SIR distributions.

6. FUTURE WORK

The analysis and display of EMBER could benefit from a number of enhancements. Our approach could be improved by obtaining more accurate geo-location data to map IP addresses to cities. We could also analyze how inaccuracy affects different regions and correct for it accordingly.

We would like to explore other ways of improving the estimation of city host population. One of the challenges is to identify gateway devices that mask large networks behind limited public-facing IP addresses and to more accurately account for their true size. Another challenge is the limited data coverage on IPv6 address space even though it has been adopted by many countries and its usage continues to grow. Finally, as computing platforms expand beyond desktop machines and servers to smart phones and gaming devices, attacks also increase in non-traditional networks such as public WiFi, cellular, and satellite networks. We need to expand our research on infection statistics and modeling to those areas as well.

We plan to apply EMBER to additional datasets and perform more systematic cross-validation against other data sources. It would be interesting to test on threat-specific datasets such as botnets and phishing sites. Moreover, we would like to continue to improve the interface to better

illustrate changes over time.

Our simple simulation does not address details of malware dynamics and it does not rule out other factors that may contribute to the long-tailed SIR distribution. These include rates of patching and disinfection, interactions between multiple malware families, underlying vulnerability distributions, use of pirated software, size and makeup of social networks, and enforcement and strength of cybercrime laws. It is likely that many of these factors causally contribute to creating cities with extreme SIR values and the long-tailed SIR distributions. Further analyses and simulations are necessary to understand temporal mechanisms leading to the long-tailed SIR distributions and to model the dynamics of multiple interacting malware families simultaneously in cities worldwide.

7. CONCLUSIONS

We have presented EMBER, a geographical display for extreme malicious behavior worldwide. EMBER scores cities by the Standardized Incidence Rate (SIR), which is the number of infections normalized by the local host population. This metric identifies cities with significantly higher or lower than expected level of malicious activity. This is useful for finding generally risky or well-protected regions, as well as regions that are targeted or avoided by specific threats. We have applied EMBER to a test dataset from DShield, which consists of millions of firewall and intrusion logs collected daily from sensors distributed worldwide. The dataset shows that cities in several Eastern European countries have the most malicious activity whereas cities in Korea, Japan and Australia appear to be best protected.

EMBER displays extreme variation in SIRs across the world and make many viewers curious about whether there are policies, behaviors, or other principles that are causally related to low and high SIRs. If so, perhaps these could be followed or made part of public policy to lower SIRs. Our simple simulation demonstrates that statistical variability caused by malware spreading preferentially locally can cause the extreme range of SIRs seen in the data. These simulations, however, only model the long-term distribution of malware. They don't model dynamics of malware spread or consider rates of patching and disinfection, underlying vulnerability distributions, the effect of pirated software, the size and makeup of social networks, enforcement and strength of cybercrime laws, and other factors. It is likely that many of these factors causally contribute to extreme SIR values. For example we already noted that many South Korean cities have low SIRs and that South Korea has strong cyber antidefamation laws, Internet content controls, and personal identification requirements when adding content to large web sites [13]. The percentage of pirated software used in a country also appears to be correlated with SIRs. Data from [8] indicate that countries for the 10 cities with the highest persistent SIRs from Fig. 6(a) have high percentages of pirated software (median 73%) while countries for the 10 cities with the lowest persistent SIRs from Fig. 6(b) have much lower percentages of pirated software (median 27%). Future analyses and research should explore the relationship between the these and other factors and the worldwide spread of malware. Computation and visualization of the SIR metric is a beginning of this search.

8. ACKNOWLEDGMENTS

The authors wish to thank DShield.org for the test data and Kevin Carter, Carolyn Buractaon and Rob Cunningham for their feedback.

9. REFERENCES

- [1] DShield, 2010. <http://www.dshield.org>.
- [2] GeoNames, 2010. <http://www.geonames.org>.
- [3] Internet World Stats, 2010. <http://www.internetworldstats.com>.
- [4] Malware Domain List, 2010. <http://www.malwaredomainlist.com>.
- [5] MaxMind GeoLite City, 2010. <http://www.maxmind.com/app/geolitecity>.
- [6] Shadowserver, 2010. <http://www.shadowserver.org>.
- [7] A. Bose and K. G. Shin. On Capturing Malware Dynamics in Mobile Power-Law Networks. In *SecureComm '08: Proceedings of the 4th international conference on Security and Privacy in Communication Networks*, pages 1–10, New York, NY, USA, 2008. ACM.
- [8] Business Software Alliance. Sixth Annual BSA-IDC Global Software 08 Piracy Study, May 2009. <http://global.bsa.org/globalpiracy2008>.
- [9] Centers for Disease Control and Prevention. U.S. Cancer Statistics: An Interactive Atlas, March 2010. http://apps.ncccd.cdc.gov/DCPC_INCA.
- [10] A. Clauset, C. Shalizi, and M. Newman. Power-Law Distributions in Empirical Data. *ArXiv*, 706, 2007.
- [11] Cyber Clean Center. FY2008 Cyber Clean Center (CCC) Activity Report, October 2009. https://www.ccc.go.jp/en_index.html.
- [12] R. Deibert, J. Palfrey, R. Rohozinski, and J. Zittrain, editors. *Access Denied: The Practice and Policy of Global Internet Filtering*. MIT Press, 2008.
- [13] Digital Nation Team. Free Speech in South Korea – Is the Internet a Poison or Cure? PBS FRONTLINE Blog, April 2009. <http://www.pbs.org>.
- [14] Y. Hideshima and H. Koike. STARMINE: a Visualization System for Cyber Attacks. In *APVis '06: Proceedings of the 2006 Asia-Pacific Symposium on Information Visualisation*, pages 131–138, Darlinghurst, Australia, Australia, 2006. Australian Computer Society, Inc.
- [15] HostExploit. Top 50 Bad Hosts and Networks, December 2009. <http://hostexploit.com/downloads.html>.
- [16] IronPort Systems. Reputation-based Mail Flow Control, 2002. <http://www.senderbase.org>.
- [17] A. Jemal, R. Siegel, E. Ward, Y. Hao, J. Xu, and M. J. Thun. Cancer Statistics, 2009. *CA: A Cancer Journal for Clinicians*, 59(4):225–249, 2009.
- [18] N. Kshetri. Positive Externality, Increasing Returns, and the Rise in Cybercrimes. *Commun. ACM*, 52(12):141–144, 2009.
- [19] F. Mansmann, D. A. Keim, S. C. North, B. Rexroad, and D. Sheleheda. Visual Analysis of Network Traffic for Resource Planning, Interactive Monitoring, and Interpretation of Security Threats. *IEEE Transactions on Visualization and Computer Graphics*, 13(6):1105–1112, 2007.
- [20] MaxMind. MaxMind GeoLite City Accuracy for Selected Countries, November 2008. http://www.maxmind.com/app/geolite_city_accuracy.
- [21] Microsoft. Microsoft Security Intelligence Report Volume 7, November 2009. <http://go.microsoft.com/?linkid=9693456>.
- [22] M. Mitzenmacher. A Brief History of Generative Models for Power Law and Lognormal Distributions. *Internet mathematics*, 1(2):226–251, 2004.
- [23] H. Ng, G. Filardo, and G. Zheng. Confidence Interval Estimating Procedures for Standardized Incidence Rates. *Computational statistics & data analysis*, 52(7):3501–3516, 2008.
- [24] N. Provos, P. Mavrommatis, M. A. Rajab, and F. Monrose. All your iFRAMEs Point to Us. In *SS'08: Proceedings of the 17th conference on Security symposium*, pages 1–15, Berkeley, CA, USA, 2008. USENIX Association.
- [25] A. Rajab, J. Zarfoss, F. Monrose, and A. Terzis. A Multifaceted Approach to Understanding the Botnet Phenomenon. *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*, page 52, 2006.
- [26] S. H. Sellke, N. B. Shroff, and S. Bagchi. Modeling and Automated Containment of Worms. *IEEE Trans. Dependable Secur. Comput.*, 5(2):71–86, 2008.
- [27] B. Stone-Gross, M. Cova, L. Cavallaro, B. Gilbert, M. Szydlowski, R. Kemmerer, C. Kruegel, and G. Vigna. Your Botnet is My Botnet: Analysis of a Botnet Takeover. In *CCS '09: Proceedings of the 16th ACM conference on Computer and communications security*, pages 635–647, New York, NY, USA, 2009. ACM.
- [28] B. Stone-Gross, C. Kruegel, K. Almeroth, A. Moser, and E. Kirda. FIRE: FInding Rogue nEtworks. In *2009 Annual Computer Security Applications Conference*, pages 231–240. IEEE, December 2009.
- [29] Team Cymru. Conficker Worm Visualizations, 2009. <http://www.team-cymru.org/Monitoring/Malevolence/conficker.html>.
- [30] N. Weaver, V. Paxson, S. Staniford, and R. Cunningham. A Taxonomy of Computer Worms. *Proceedings of the 2003 ACM Workshop on Rapid Malcode (WORM)*, pages 11–18, 2003.
- [31] V. Yegneswaran, P. Barford, and J. Ullrich. Internet Intrusions: Global Characteristics and Prevalence. *ACM SIGMETRICS Performance Evaluation Review*, 31(1), 2003.
- [32] C. Zesheng, J. Chuanyi, and P. Barford. Spatial-Temporal Characteristics of Internet Malicious Sources. In *2008 IEEE INFOCOM - The 27th Conference on Computer Communications*, pages 2306–2314. IEEE, 2008.