

Sampling Strategies for Epidemic-Style Information Dissemination

Milan Vojnović*, Varun Gupta†, Thomas Karagiannis*, and Christos Gkantsidis*
 * Microsoft Research † Carnegie Mellon University

Abstract—We consider epidemic-style information dissemination strategies that leverage the nonuniformity of host distribution over subnets (e.g., IP subnets) to optimize the information spread. Such epidemic-style strategies are based on random sampling of target hosts according to a sampling rule. In this paper, the objective is to optimize the information spread with respect to minimizing the total number of samplings to reach a target fraction of the host population. This is of general interest for the design of efficient information dissemination systems and more specifically, to identify requirements for the containment of worms that use subnet preference scanning strategies.

We first identify the optimum number of samplings to reach a target fraction of hosts, given global information about the host distribution over subnets. We show that the optimum can be achieved by either a dynamic strategy for which the per host sampling rate over subnets is allowed to vary over time or by a static strategy for which the sampling over subnets is fixed. These results appear to be novel and are informative about (a) what best possible performance is achievable and (b) what factors determine the performance gain over oblivious strategies such as uniform random scanning. We then consider several simple, online sampling strategies that require only local knowledge, where each host biases sampling based on its observed sampling outcomes and keeps only $O(1)$ state at any point in time. Using real datasets from several large-scale Internet measurements, we evaluate the significance of the factors revealed by our analytical results on the sampling efficiency.

I. INTRODUCTION

We consider the problem of reaching a target fraction of an initially unknown population of hosts by using epidemic-style information dissemination. We make the assumption that the nodes of the population take identities from a large space (such as the space of IPv4 addresses), and that two nodes communicate only when one of them discovers the other by random probing. Initially, the information to be disseminated exists in few nodes which use random probing to discover more nodes. Discovered nodes also initiate a random probing process to further propagate the information. In this work, we describe optimum static and dynamic strategies, as well as local sub-optimal strategies to efficiently disseminate the information to a given fraction of the population.

If this “information” corresponds to worm-like malicious software, the problem transforms to characterizing the performance of worm propagation. In this case, we are interested in characterizing the minimum number of scans that a worm needs to infect a target fraction of susceptible machines; such knowledge is of value to designers of worm countermeasures. In general, however, epidemic-style information dissemination is of interest in a plethora of areas, including web-service

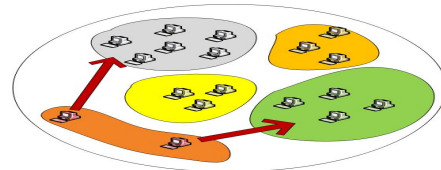


Fig. 1. The setting under consideration: strategies for efficient sampling of hosts that are distributed over different groups.

membership management [1], database maintenance [2], and streaming broadcasting [3]. In such systems, information dissemination is assisted by the underlying structure of the network (e.g., p2p overlays). However, it is of theoretical interest and practical value to understand whether the information in such systems can be efficiently disseminated without the support of an underlying network, when the population of the system and its distribution across the dissemination network are a-priori unknown. This is consistent for example with trying to reach hosts, a large fraction of which connect to the Internet through dynamic IPs [4], ergo resulting in high population variability over time and IP space.

Specifically, to design efficient information dissemination strategies we take into advantage the non-uniformity of host distribution over subnets. In particular, we assume that hosts are partitioned into groups or subnets such as IP address blocks or subnets defined by Autonomous Systems (ASes) (see for example Fig. 1). Our goal is then to minimize the samplings required to reach a predetermined fraction of hosts without any prior knowledge of their group partitioning, and only requiring the minimum amount of state possible.

We first examine the case where perfect global knowledge of the host distribution over subnets exists. We find that the optimum in this context may be achieved by either a static or a dynamic strategy, where surprisingly both strategies can be described by the same formula. Our analysis provides insights about the best possible performance that can be achieved. We show that while the selection of the right target set of subnets is crucial to achieve the optimum, further sampling optimization within this target set only bears minimal benefits.

While it may be argued that the knowledge required by the optimum strategy (i.e., the distribution of hosts over subnets) can be estimated through various sources such as, for example, intrusion detection systems, BGP tables or even reports of infected hosts from earlier worm attacks, we present evidence that such an inference might not be as straightforward. Using a set of diverse real Internet measurements, we highlight that different datasets provide distinct viewpoints of the Internet

topology, stressing the need for sampling strategies that require no prior knowledge of subnet partitioning, but instead they are able to deduce this information online. Motivated by this finding, we propose simple strategies that bias their sampling based on previous sampling attempts using just $O(1)$ state.

Our contributions can be summarized as follows:

- We describe the optimum static and dynamic strategies to minimize the number of samplings required to reach a target fraction of the host population (Sec. III-A and IV-A). Note that these appear to be novel results.
- We propose and evaluate simple local strategies (Sec. V) that require no prior knowledge of the host distribution over subnets, and yet significantly outperform uniform random scanning as well as local-subnet preference strategy.
- Using Internet measurement datasets, we evaluate and provide insights regarding the factors that determine the performance of the optimal and the sub-optimal strategies. We further highlight that different datasets provide distinct perspectives of the Internet topology both across the measurement traces, but also across time within the same trace (Sec. VI).

To the best of our knowledge, our work is the first to fully describe the optimum with respect to minimizing the number of samplings required to reach a fraction of hosts. While subnet preference sampling has been analyzed earlier [5] and the objective of minimizing samplings has been identified as important also in previous work [6], the authors there restricted themselves to studying sub-optimal strategies.

II. ASSUMPTIONS AND NOTATION

We first introduce the notation that we use in the remainder of the paper. We denote the size of the total address space with Ω (e.g., for IPv4, $\Omega = 2^{32}$). The address space is divided into J subnets with subnet j having an address space of size Ω_j . We will refer to the hosts that received the information item of our dissemination system as infected hosts and those that are interested in receiving the information items but have not yet received it as susceptible hosts. This nomenclature is commonplace in the literature. The number of hosts in a subnet j that are interested in the information item is denoted with N_j with N defined as the total number of hosts (i.e., $N = \sum_{j=1}^J N_j$). We denote with $I_j(t)$ the number of infected hosts in subnet j at time t . Similarly, we let $S_j(t) = N_j - I_j(t)$ be the number of susceptible hosts in subnet j at time t . The quantities n_j , $i_j(t)$ and $s_j(t)$ are the normalized versions defined as follows: $n_j = N_j/N$, $i_j(t) = I_j(t)/N$, and $s_j(t) = S_j(t)/N$. We denote with $i(t)$ and $s(t)$ the total fraction of infected and susceptible hosts, respectively, i.e. $i(t) = \sum_{j=1}^J i_j(t)$ and $s(t) = \sum_{j=1}^J s_j(t)$. Also, let $\omega_j = \Omega_j/\Omega$ be the fraction of the total address space occupied by the j th subnet. Without loss of generality, we assume that subnets are enumerated such that

$$\frac{S_1(0)}{\Omega_1} \geq \frac{S_2(0)}{\Omega_2} \geq \dots \geq \frac{S_J(0)}{\Omega_J} \geq 0 \quad (1)$$

and $S_1(0)/\Omega_1 > 0$.

We let η be the rate at which an infected host samples the address space for a susceptible host. Without loss of generality, we assume that $\eta = 1$. In the most general setting, at time t , an infected host in a subnet i decides to sample a node in a subnet j with probability $p_{ij}(t)$. Once it chooses subnet j , it samples

an address lying in subnet j 's address space *uniformly at random* and then initiates a contact to this address. Let β be the density of hosts, i.e., $\beta = \frac{N}{\Omega}$. Define $\beta_{ij}(t) = \beta \frac{1}{\omega_j} p_{ij}(t)$. We consider the many host limit where the total number of hosts N tends to be large while the following parameters are held fixed N/Ω , Ω_j/Ω , and $i_j(0)$, $j = 1, \dots, J$. Under assumption that each host initiates samplings at instances of a Poisson process with rate 1, the infected host population is described by a Markov process indexed with N . The host population frequencies converge with N uniformly on any compact time interval to the solution of the following system of ordinary differential equations (e.g., see Kurtz [7])

$$\frac{d}{dt} i_j(t) = \left(\sum_{i=1}^J \beta_{ij}(t) i_i(t) \right) s_j(t) \quad (2)$$

with $s_j(t) = n_j - i_j(t)$, for $j = 1, \dots, J$.

Given the density of hosts over subnets n_j and initial placement of infected hosts $i_j(0)$, seeking for an optimal sampling strategy is equivalent to finding the functions $\beta_{ij}(t)$, $t \geq 0$. The optimal strategy β^* would depend on the imposed constraints, such as, for example, for static subnet preference that β^* is time invariant. β_{ij} will depend only on j if one does not allow the nodes to scan their subnet faster or have privileged information about their subnet. In that case, all the subnets have the same subnet sampling bias.

In the sequel, we will denote with $u(t)$ the total number of samplings per host by time t , i.e.

$$u(t) = \int_0^t i(x) dx. \quad (3)$$

III. STATIC SUBNET PREFERENTIAL SAMPLING

We consider a class of sampling strategies for which the subnet preference probability distribution p is fixed in time. For this subset of sampling strategies, Eq. (2) boils down to:

$$\frac{d}{dt} i_j(t) = \beta_j i(t) s_j(t) \quad (4)$$

where $\beta_j := \beta p_j / \omega_j$, $j = 1, \dots, J$. From Eq. (4), it follows

$$i_j(u) = n_j - s_j(0) e^{-\beta_j u}, \quad j = 1, \dots, J \quad (5)$$

where $u(t)$, $t \geq 0$, is given by $u(0) = 0$ and

$$\frac{d}{dt} u(t) = 1 - \sum_{j=1}^J s_j(0) e^{-\beta_j u(t)}, \quad t \geq 0. \quad (6)$$

Note that for any static sampling strategy, the time dynamics of infected host population over subnets is entirely described by Eq. (5) and (6).

We briefly revisit the uniform random sampling, the well known S-I epidemics, which we will use recurrently as a reference as it is a commonplace sampling rule in practice. Note that with uniform random sampling $p_j = \omega_j$, i.e. a subnet j is sampled proportional to the address space size of the subnet j . Uniform random sampling is easy to analyze, for example, note that

$$s(u) = s(0) e^{-\beta u}, \quad u \geq 0. \quad (7)$$

We will see later that for real-life distributions of hosts over subnets, the uniform random sampling is grossly less efficient than an optimal strategy that minimises the samplings for a given target fraction of infected hosts. The following preliminary result provides a comparison of static sampling strategies and uniform random sampling; proof in Appendix [8].

Proposition 1 *The set of static sampling strategies is characterised as follows: (a) For any static sampling strategy specified by the subnet preference distribution p such that*

$$\frac{\sum_{j=1}^J p_j \frac{s_j(0)}{\Omega_j}}{\frac{s(0)}{\Omega}} < 1$$

the total fraction of infected hosts is smaller than under uniform random sampling, for any given fraction of samplings $u \geq 0$, i.e. $i^p(u) < i^\omega(u)$, for all $u > 0$, where $i^p(u)$ and $i^\omega(u)$ are the total fractions of infected hosts under static sampling p and uniform random sampling ω , respectively.

(b) For any subnet preference distribution p such that $p_j < \omega_j$ for some subnet j with $s_j(0) > 0$, $i^p(u) < i^\omega(u)$, for some $u > 0$.

Item (a) identifies a sufficient condition under which a static sampling strategy p is less efficient than uniform random sampling. These are sampling strategies that in average sense bias to sampling of rare subnets. Item (b) entails that any static sampling strategy other than uniform random sampling over a set of subnets with each containing susceptible hosts is worse than uniform random sampling over this set of subnets, for some total fraction of samplings u . In the next section, we identify optimal static sampling strategy that in the case of nonuniformly dense subnets requires smaller number of samplings than uniform random sampling, for a given target fraction of infected hosts.

A. Optimal Static Strategy

In this section, we identify the optimal static sampling strategy that minimises the total number of samplings to reach a given fraction of infected hosts i^0 . We show that the static strategy OPT-STATIC, specified by the following subnet preference probabilities is optimal:

OPT-STATIC
Each infected host scans a subnet j with probability:

$$p_j = \begin{cases} \alpha \omega_j \log \left(\frac{\frac{s_j^A(0)}{\omega_j^A}}{1 - \frac{i^0 - i(0)}{\sum_{k \in A} s_k(0)}} \right) & j \in A \\ 0 & j \notin A \end{cases} \quad (8)$$

where α is the normalization constant and A is the set of subnets $\{1, 2, \dots, J'\}$ with

$$J' = \max \left\{ j : \frac{s_j(0)}{\omega_j} > \frac{\sum_{k=1}^j s_k(0) - (i^0 - i(0))}{\sum_{k=1}^j \omega_k} \right\}.$$

Here we used the following notation:

$$s_j^A(0) := \frac{s_j(0)}{\sum_{k \in A} s_k(0)} \text{ and } \omega_j^A = \frac{\omega_j}{\sum_{k \in A} \omega_k}, j \in A.$$

The strategy specifies to sample a set A of initially densest subnets. The necessary condition for a set A to be optimal is that the initial density of susceptibles in every subnet in A must be larger than the final density of susceptibles in A . Indeed, if we target the set A then after target infection is reached, the final density of susceptibles in A is given by ¹

$$\frac{\sum_{j \in A} S_j(t^0)}{\sum_{j \in A} \Omega_j} = \frac{\sum_{j \in A} S_j(0) - (I^0 - I(0))}{\sum_{j \in A} \Omega_j}$$

where t^0 is the time when the fraction of infected hosts i^0 is reached. Furthermore, note that if a subnet j is in the target set A , all subnets whose initial density of susceptibles is larger than the initial density of susceptibles in j are also in A . Lastly, note that the strategy does not necessarily target the smallest densest set of subnets. One may need to target a larger set because even though that it may slow the initial phase, the density of susceptibles will still be sufficient as infection reaches target infection. Finally, note that in Eq. (8) the subnet preference probability for a subnet that is in the target set A is an expression containing a term logarithmic in the initial density of this subnet.

The next result establishes that OPT-STATIC is optimal over all static sampling strategies.

Theorem 2 *For any given target fraction of infected hosts, the strategy OPT-STATIC is optimal in minimising the total number of samplings over all static sampling strategies. The total number of required samplings for a target fraction of infected hosts i^0 is given by*

$$u_{\text{STA}}(i^0, A) = \frac{1}{\beta} \left(\sum_{j \in A} \omega_j \right) \left[\log \left(\frac{1}{1 - \frac{i^0 - i(0)}{\sum_{k \in A} s_k(0)}} \right) - D(\omega^A \| s^A(0)) \right]$$

where $D(\cdot \| \cdot)$ denotes Kullback-Liebler (KL) divergence.²

Proof: See Appendix [8]. ■

What does this tell us? We compare the required total number of samplings of the optimal static strategy to uniform random sampling. We first note the following easy result:

Proposition 3 *Under uniform random sampling of a subset of subnets A , we have that the per host total number of required samplings for a target fraction of infected hosts i^0 is given by*

$$u_{\text{UNI}}(i^0, A) = \frac{1}{\beta} \left(\sum_{j \in A} \omega_j \right) \log \left(\frac{1}{1 - \frac{i^0 - i(0)}{\sum_{j \in A} s_j(0)}} \right).$$

We note that the required total number of samplings for OPT-STATIC differs from that of the uniform random sampling only in the KL term $D(\cdot)$, presuming that both strategies target the set A specified by OPT-STATIC. The KL term measures the deviation between the initial susceptible host population sizes $s_j(0)$ and the subnet address space sizes ω_j over subnets in the target set. In particular, if the address sizes of subnets

¹In addition, at the end of the infection, all subnets in the target set will have equal density of susceptible nodes.

²Kullback-Liebler divergence between two probability measures p and q is defined by $D(p \| q) := \sum_i p(i) \log(p(i)/q(i))$.

are equal, then the KL measures the deviation of $s_j(0)$ from the uniform distribution on the target set A . Note that

$$D(\omega^A || s^A(0)) = \log \left(\frac{A_{\omega^A}(\rho(0))}{G_{\omega^A}(\rho(0))} \right)$$

where $\rho_j(0) := s_j(0)/\omega_j$ is the density of a subnet j , $A_{\omega^A}(\rho(0))$ and $G_{\omega^A}(\rho(0))$ are the arithmetic and geometric means, respectively.

$$A_{\omega^A}(\rho(0)) = \sum_{j \in A} \frac{\omega_j}{\sum_{k \in A} \omega_k} \rho_j(0)$$

$$G_{\omega^A}(\rho(0)) = \prod_{j \in A} \rho_j(0)^{\frac{\omega_j}{\sum_{k \in A} \omega_k}}$$

Indeed, we have that $D(\omega^A || s^A(0)) > 0$ unless all subnets are of same density, i.e., $s_i(0)/\omega_i = s_j(0)/\omega_j$, for all $i, j \in A$ (elementary property of means, Hardy, Littlewood, and Pólya [9, Section 2.5]).

The question of interest is how significant is the KL term in Theorem 2 relative to the logarithmic term therein. If the KL term is insignificant relative to the logarithmic term, then this suggests that provided that one identifies the target set A , then simple uniform random sampling over the set A would yield nearly-optimal performance. Furthermore, it is of interest to evaluate how critical is the selection of the target set on the resulting total number of samplings. We will address these questions in Section VI.

IV. DYNAMIC SAMPLING STRATEGIES

In this section we consider strategies for which subnet preference probabilities are allowed to vary over time.

A. Optimal dynamic strategy

We now consider what optimal performance can be achieved over the entire set of feasible sampling strategies. A priori, it may not be clear whether enlarging the set of strategies from static to dynamic will yield better performance. We show here that the answer is no. The following scheme is optimal over the entire set of dynamic sampling strategies:

OPT

At any point in time t with the subset of densest subnets $S(t)$, each infected host samples uniformly at random an address over the address space of subnets $S(t)$.

The optimality of this strategy would appear to be very intuitive. The strategy was claimed to be optimal in [6], however, we are unaware of a proof in the literature that shows that this is indeed an optimal strategy. We also characterize the host evolution for this strategy, which we later use to compare with sub-optimal strategies.

Theorem 4 OPT has the following properties:

- 1) The strategy is optimal in that at any point of its execution the fraction of infected hosts is maximised.
- 2) For any given target fraction of infected hosts i^0 , the total number of samplings is given by the relation in Theorem 2.

- 3) For the total number of samplings $u \geq 0$, the fraction of susceptible hosts is given by, for $u_{n-1} \leq u < u_n$

$$s(u) = \sum_{j=n+1}^J s_j(0) + \left(\sum_{j=1}^n \omega_j \right) G_n e^{-\frac{\beta}{\sum_{j=1}^n \omega_j} u} \quad (9)$$

where

$$G_n = \prod_{j=1}^n \left(\frac{s_j(0)}{\omega_j} \right)^{\frac{\omega_j}{\sum_{k=1}^n \omega_k}}$$

and $u_0 = 0$ with

$$u_n = \frac{\sum_{j=1}^n \omega_j}{\beta} \log \left(\frac{G_n}{\frac{s_{j+1}(0)}{\omega_{j+1}}} \right) \quad (10)$$

for $n = 1, 2, \dots, J-1$, and $u_J = +\infty$.

The result entails the following corollary:

Corollary 5 For any given target fraction of infected hosts i^0 , there exists a static sampling strategy that achieves the smallest possible total number of samplings to infect the fraction of hosts i^0 over all dynamic sampling strategies. This static sampling strategy is OPT-STATIC.

We end this sub-section with a comparison with uniform random sampling.

Corollary 6 For uniform random sampling, we have

- 1) With the target fraction of infected hosts going to 1, the total number of samplings is asymptotically optimal.
- 2) In the prevailing limit, the fraction of susceptible hosts under uniform random sampling and optimal satisfy

$$\lim_{u \rightarrow +\infty} \frac{S_{UNI}(u)}{S_{OPT}(u)} = e^{D(\omega || s(0))} \quad (11)$$

where $D(\omega || s(0))$ is the Kullback-Lieber divergence between $(\omega_1, \dots, \omega_J)$ and $(s_1(0), \dots, s_J(0))$.

Item 1 is rather intuitive. Let $i^0 = 1 - 1/N$ be the fraction where all but one host are infected. Then, from Theorem 4, we have that the log term is logarithmic in N while the KL term is a constant, thus asymptotically negligible. Item 2 follows directly from Theorem 4, Proposition 3, and the fact Eq. (7).

B. Proportional Sampling: A Sub-Optimal Dynamic Strategy

The optimal strategies discussed so far would in many cases be difficult to implement as they require global knowledge about densities of susceptible hosts over subnets. In particular, the optimal static strategy requires knowing initial densities of susceptible hosts over subnets while the dynamic optimum strategy requires knowing the subset of densest subnets at any point in time during its execution. In this section, we consider a sampling strategy that is based on sampling a subnet proportional to the number of susceptible hosts in this subnet. It follows from the analysis below that, in general, proportional sampling is a sub-optimal strategy. It is not clear, though, how far is the proportional sampling from optimal for distributions of hosts over subnets in practice; we investigate this in Section VI. We next characterise a generalized version

of the proportional sampling (we call PROP) that is a mix of uniform random sampling and sampling proportional to the density of susceptibles per subnet. The strategy is specified by the parameter $0 \leq q \leq 1$ denoting the probability that a host samples a subnet proportional to the number of susceptibles in this subnet. ($q = 1$ corresponds to pure proportional sampling.)

PROP

An infected host with probability q samples a subnet proportional to the current number of susceptible hosts in this subnet, or else samples a subnet by uniform random sampling of the entire address space.

It turns out that under strategy PROP the fraction of susceptible hosts for any given total number of samplings $u \geq 0$ is given in a simple analytical form:

Theorem 7 *The fraction of susceptible hosts in a subnet j is given by*

$$s_j(u) = s_j(0) \frac{e^{-\beta(1-q)u}}{1 + \frac{s_j(0)}{\omega_j} \psi(u)} \quad (12)$$

where ψ is the implicit function

$$\sum_{k=1}^J \omega_k \log \left(1 + \frac{s_k(0)}{\omega_k} \psi(u) \right) = \beta qu. \quad (13)$$

Proof (Appendix [8]) shows that in fact the dynamics under PROP is entirely described by a scalar differential equation for ψ and that the evolution of the susceptible hosts in a subnet j is given by the function of ψ given in the theorem.

V. SAMPLING STRATEGIES THAT USE ONLY LOCAL KNOWLEDGE

We consider sampling strategies that are local in that each host biases its sampling over subnets based solely on the observed successes and failures of its own samplings. Moreover, we confine our attention to strategies that at any time keep state of only a fixed number of subnets with respect to the total number of subnets in the system. We consider several sampling strategies and describe their dynamics by differential systems. This enables our numerical evaluations in Section VI.

A. Local Subnet Preference

We first consider the well-known local subnet preference strategy (e.g., used by CodeRed-II worm), defined as follows:

LOC-PREF

A infected host in a subnet j , with probability q_j , samples an address uniformly over the address space of the subnet j , or else it samples an address uniformly over the entire address space.

This strategy can be seen as a dynamic sampling strategy specified by the following subnet preference probabilities

$$p_j(t) = \omega_j \left(1 - \sum_{k=1}^J q_k v_k(t) \right) + q_j v_j(t) \quad (14)$$

where $v_j(t) := i_j(t) / \sum_{k=1}^J i_k(t)$.

B. K-FAIL Strategy

We next consider another strategy that biases to subnets from which a host observes successful samplings. The strategy is described as follows ($K \geq 1$ is a configuration parameter):

K-FAIL

- 1) Initially, infected hosts use uniform random sampling.
- 2) When a host that performs uniform random sampling successfully samples a host in a subnet j , it continues to sample uniformly at random on the subnet j until K consecutive scan failures.
- 3) Upon K consecutive failures, the host switches to uniform random sampling, until it successfully samples a susceptible host when it goes to item 2.

We are unaware of a dissemination system in practice that uses this strategy. The strategy is perhaps closest related to the Zotob family of worms that used similar but different strategy. With some Zotob worms, an infected host starts with scanning its local subnet until K consecutive failures, and then switches to and remains indefinitely a uniform random scanner. Instead, K-FAIL strategy sticks to a subnet from which it successfully samples a susceptible host and may switch over between uniform random sampling mode and sticking to a subnet mode. We next describe dynamics of K-FAIL strategy. Each infected host is in one of K states: 0 denoting the state in which host performs uniform random sampling, or state k where $K - k$ denotes the number of successive failures that the host already incurred, $k = 1, \dots, K$. We denote with r_0 the fraction of infected hosts that are in state 0 and with $r_{j,k}$ the fraction of infected hosts in a subnet j that are in state k . The dynamics of the system is entirely described by the following differential equations:

$$\begin{aligned} \frac{d}{dt} s_j &= -\beta \left(\omega_j r_0 + \sum_{k=1}^K r_{j,k} \right) \frac{s_j}{\omega_j} \\ \frac{d}{dt} r_0 &= \sum_{i=1}^J r_{i,1} \left(1 - \beta \frac{s_i}{\omega_i} \right) - r_0 \beta \sum_{i=1}^J s_i \\ \frac{d}{dt} r_{j,k} &= -r_{k,j} + r_{k+1,j} \left(1 - \beta \frac{s_j}{\omega_j} \right), 1 \leq k < K \\ \frac{d}{dt} r_{j,K} &= -r_{j,K} + 2 \left(r_0 \beta s_j + \beta \frac{s_j}{\omega_j} \sum_{k=1}^K r_{j,k} \right). \end{aligned}$$

The equations capture the transitions of host states. We use this system in our numerical evaluations in Section VI.

C. Subnet Preference Strategy

We consider a subnet preference sampling strategy where each host maintains a candidate set of a fixed number $K \geq 1$ of subnets. Each host splits its effort to sampling subnets from its candidate set and uniform random sampling.

Again, we are unaware of a system that uses this strategy. It appears natural to consider strategies for which each host keeps a (small) list of subnets and uses sorting strategies of

the candidate subnets in the list, such as for example, move-to-front sort heuristic, in order to bias its sampling to subnets at the head of the list. We instead consider the below introduced K-CANDSET strategy that keeps no order of the items in the list for ease of analysis. Indeed, for the candidate set of size at most 1, i.e. $K = 1$, the two schemes are equivalent. In this case, it can be considered as a generalization of local subnet preference sampling by letting the preference subnet not to be fixed to local subnet but to change to subnets from which the host observes successful samplings. The strategy is introduced in the following, where $0 \leq q < 1$ is a configuration parameter.

K-CANDSET

- 1) Init: infected hosts set their candidate sets according to a policy (default: empty set).
- 2) With probability q , a host samples a subnet by picking uniformly at random from its candidate set.
- 3) Otherwise, the host samples by uniform random sampling of the entire address space. If the sampling to a subnet k is successful and subnet k is not in the candidate set of this host, the following happens:
 - a) **If** the candidate set of this host is smaller than K , then subnet k is added to the candidate set,
 - b) **else**, the subnet k replaces a subnet from the candidate set, evicted uniformly at random.
- 4) A host infected by an instigator host, inherits the candidate set of the instigator host (the candidate set of the instigator updated after successful sampling of this host).

The dynamics of the above described sampling strategy is entirely specified by the following system of ordinary differential equations. In order to keep the notation simple, we display equations only for the special case $K = 1$, but note that similar equations are easily derived for the general case. Let r_k denote the fraction of infected hosts of type k , i.e., those with the candidate set $\{k\}$ and let r_0 denote the fraction of hosts of type 0, i.e., with empty candidate set.

$$\begin{aligned} \frac{d}{dt} s_j &= -\beta \left((1-q)\omega_j \sum_{k=0}^J r_k + q(\omega_j r_0 + r_j) \right) \frac{s_j}{\omega_j} \\ \frac{d}{dt} r_0 &= -r_0 \sum_{k=1}^J \beta s_k \\ \frac{d}{dt} r_j &= 2(1-q) \left(\sum_{k=0}^J r_k \right) \beta s_j + \\ &+ r_j \left(q\beta \frac{s_j}{\omega_j} - (1-q) \sum_{k=1}^J \beta s_k \right). \end{aligned}$$

Again, we will use these equations to derive our numerical results in Section VI.

VI. EXPERIMENTAL RESULTS

In this section we perform an extensive evaluation of the strategies described throughout sections III- V. We first outline

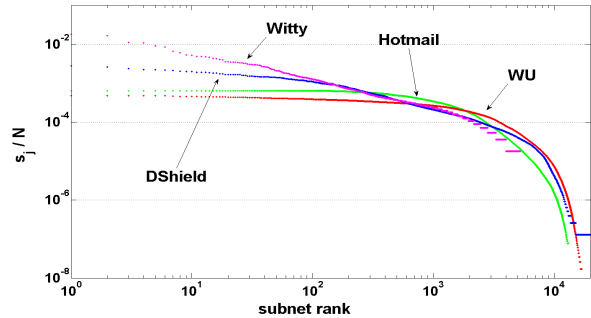


Fig. 2. Subnet density for all datasets.

the set of analysed datasets which cover diverse Internet measurements reflecting the distributions of IP addresses over the IP space (Section VI-A). We then examine the factors that determine the dynamics of the optimum strategy (Section VI-B) and finally we evaluate the performance of the proposed strategies that require no side knowledge about the distribution of hosts over subnets (Section VI-C).

A. Datasets

Our datasets consist of measurement traces of IP addresses. We aggregate IPs into groups or subnets of various sizes such as for example /8 and /16 subnets, or into groups based on Autonomous Systems (AS). Without loss of generality we will use /16 subnet groups for the remainder of the paper unless otherwise specified. Throughout our evaluation, we make use of the following datasets:³

- *WU* : The dataset refers to IIS logs collected at the Windows Update system [10]. In our experiments we will use the 117 million IP addresses observed during the first day of the measurement. (Populated /16s: 17503.)
- *Hotmail* : The dataset consists of approximately 103 million IP addresses which were observed over a period of three months from logs of user-logins at the Hotmail service [4]. (Populated /16s: 13135.)
- *DShield* : The dataset consists of roughly 7.6 million IP addresses that were collected by a set of firewall and intrusion detection systems provided by DShield [11] and were used in [12] and [13]. This dataset may contain spoofed source IP addresses. (Populated /16s: 22861.)
- *Witty* : A list of IPs (roughly 55 thousands) corresponding to hosts spreading the Witty worm provided by CAIDA [14]. (Populated /16s: 5271.)

Fig. 2 presents the density of hosts in /16 subnets as seen in each of our datasets. Density here refers to the fraction N_j/N , while the x -axis presents to the rank of each subnet with respect to its density (i.e., $x = 1$ refers to a densest subnet). Note that while overall the shapes of the curves appear similar, densities of distinct subnets appear quite different across the datasets. We will extensively examine these differences in the following section, as they significantly impact the performance of the sampling strategies, e.g. the selection of the target set A described in Section III-A.

³Due to space limitations results will be presented for only some datasets interchangeably. Our findings however apply to all datasets.

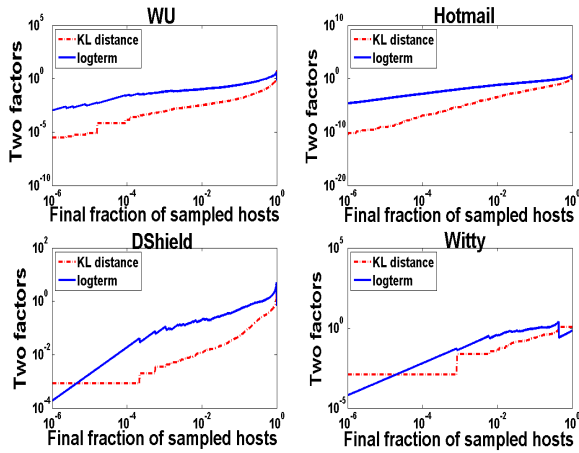


Fig. 3. Logterm vs. KL term in the four datasets. The KL term does not appear to be significant.

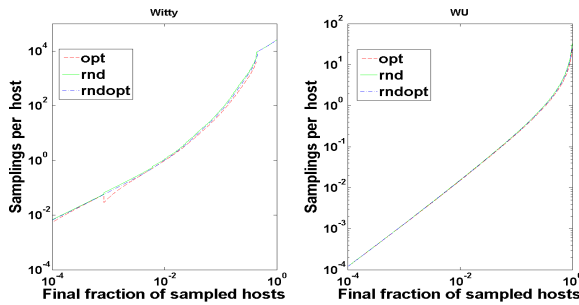


Fig. 4. Optimal strategy (*opt*) vs. uniform random scanning of the target set (*rnd*) and the optimal set of densest subnets (*rndopt*). All curves fall on top of one another.

B. Optimal strategy evaluation

We first examine the factors that affect the total number of samplings for the optimal sampling strategy. Note that in Theorem 2 we found that the total number of samplings per susceptible host depends on the logarithmic and the KL term therein. We now examine the significance of these two individual factors. This is of interest as if it turns out that the KL term can be neglected relative to the logarithmic term, then the implication is that simple uniform random sampling of the target set A would already be near-optimal and fine-tuning of the subnet preference probabilities may not be needed. Fig. 3 specifically examines these two factors in the four datasets as the fraction of infected hosts i^0 grows. Indeed, in all cases we observe that the KL term is several orders of magnitude smaller than the logarithmic term, especially in the larger datasets for smaller i^0 . Furthermore, while for the *DShield* and *Witty* datasets the KL term appears larger for a range of values for small i^0 , this is only a side-effect of the smaller number of IP addresses in these datasets and does not reflect true operational regions (e.g., for $i^0 = 10^{-6}$ in the *witty* dataset there are no hosts to be infected since the total number of IP addresses is roughly 55 thousand).

The above observation suggests that *optimization of scanning rates over subnets of the optimal target set yields insignificant or moderate gain compared to simply using uniform random scanning over the optimal target set*. This is further

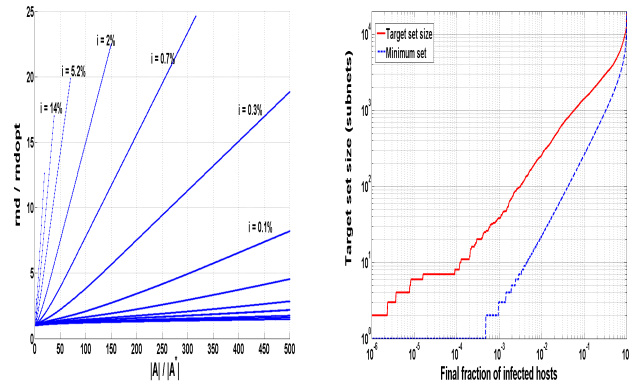


Fig. 5. LEFT: Augmenting the target set size and its effect on the number of samplings. RIGHT: The optimal target set size and the minimum number of sets to cover a fraction of the population.

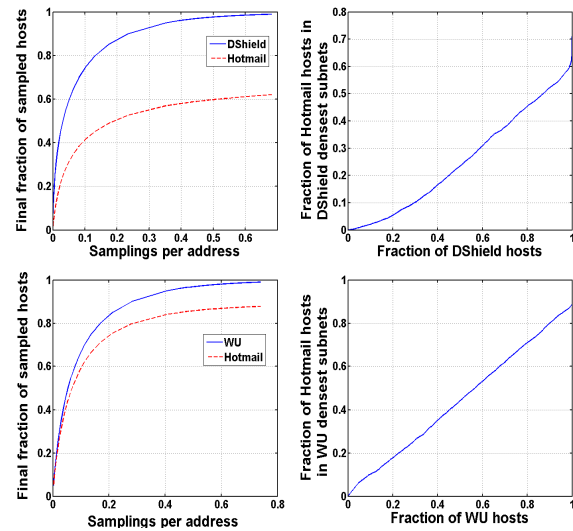


Fig. 6. Examples of employing a non-optimal target set with a prior distribution. TOP: Prior distribution: *DShield*. True: *Hotmail*. BOTTOM: Prior distribution: *WU*. True: *Hotmail*.

validated by Fig. 4 which compares the optimal strategy vs. uniform random sampling within the target set and uniform random sampling on the smallest set of densest subnets that covers the target fraction of infected hosts (*rnd* and *rndopt* curves in the figure). In Fig. 4, all curves fall on top of one another, showing that these strategies perform similarly if the optimal target set has been well identified.

However, *the choice of the target set is critical* with respect to the total number of samplings, since a poor choice of the target set A may have a dramatic impact on the required number of samplings to reach a target fraction of the host population. To evaluate the discrepancy from the optimal strategy with a poor choice of A , we perform the following experiment. For a given target fraction of infected hosts i^0 , we augment the optimal target set A^* by adding subnets in decreasing order of their densities, and we then examine the incurred penalty. For example, if the optimum target set is $A^* = \{1, 2, \dots, j\}$, we define sets A as $A^* \cup \{j + 1\}$, $A^* \cup \{j + 1, j + 2\}$, etc. In our evaluations, we enlarge the target set by a factor k , i.e. $k = \frac{|A|}{|A^*|}$. This effect is shown

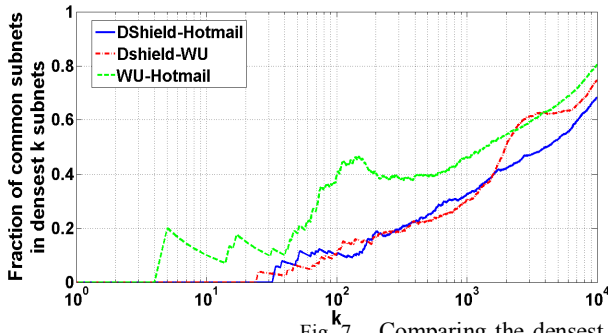
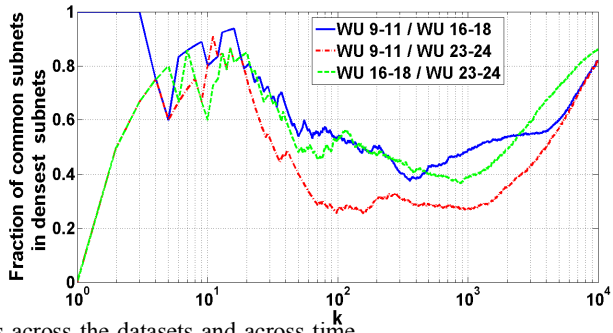


Fig. 7. Comparing the densest subnets across the datasets and across time.



in Fig. 5(left) where for various i^0 of the *WU* dataset, we show the ratio of the number of samplings by uniform random sampling of the set A to that of uniform random sampling of the optimum target set A^* vs. the target fraction of infected hosts. For reference, Fig. 5(right) presents the size of the optimal target set A^* and the minimum possible set of subnets to cover a given fraction of the target host population i^0 . We observe that as i^0 increases, the penalty factor becomes quite significant especially when i is larger than 1% (e.g., to reach roughly 14% of the population, a tenfold increase of the target set, will produce 5 times more samplings). When i^0 is small augmenting the target set size does not incur a high penalty, since i^0 is already reachable by the densest sets already in A^* . Fig. 5 also shows that the penalty increases roughly linearly with k for a given i^0 . (In Fig. 5-left, $i \equiv i^0$.)

The criticality of the target set is also highlighted when the optimal static sampling strategy is configured using a prior distribution of hosts over subnets, and then applied over a host population following a different distribution. We compare the final fraction of sampled hosts when the above strategy is applied to a) the population distributed according to the prior distribution, and b) a population distributed according to another "true" reference distribution. This could be seen as having imperfect knowledge about the true distribution, which can result in using a non-optimal target set A . Fig. 6 presents the results of two such experiments: when using *DShield* as a prior and *Hotmail* as the true distribution, and using *WU* as a prior and *Hotmail* as the true distribution. For reference, we also plot the cumulative fraction of hosts for the true distribution residing in the subnets corresponding to the densest subnets in the prior distribution with a given total fraction of population. Fig. 6 highlights that the discrepancy with respect to the optimal is significant.

How sub-optimal might the target set A be having partial knowledge? There are two important factors which determine how useful a prior distribution might be:

- 1. The fraction of the population in the true distribution for which the corresponding subnets in the prior are empty.** The importance of this fact is highlighted when comparing *Hotmail* vs. *Dshield*. Approximately 30% of *Hotmail* hosts reside in subnets which are empty in the *DShield* data set ($x = 1, y = 0.7$ top right plot in Fig. 6). This implies that when using *DShield* as a prior for *Hotmail*, we can never reach more than 70% of *WU*'s population.

- 2. The difference in distribution of hosts among the**

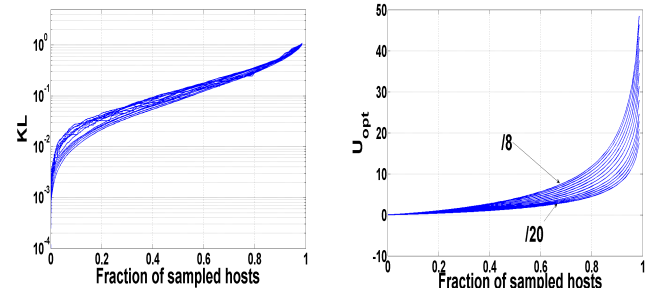


Fig. 8. The impact of the group size on the KL term and the samplings under the optimal strategy.

populated subnets. The importance of this fact is highlighted when using *WU* vs. *Hotmail*. In both cases, approximately 10% of hosts in one data set reside in subnets which are empty in the other data set. However, the curves comparing the distributions look drastically different. The plot of *Hotmail*'s distribution in the densest *WU* subnets as a function of cumulative fraction of hosts in *WU* is almost a straight line and very close to $y = x$ (Fig. 6, bottom, right). Thus, using *WU* as a prior for *Hotmail* performs better than using *DShield*; yet, the other way round (i.e., *Hotmail* as a prior for *WU*) does not (plot omitted due to lack of space; similar to Fig. 6, top).

To further examine the difference between the distributions according the various datasets, we compare the analysed traces by examining the fraction of common densest subnets within the set of top- k densest subnets for each distribution. Fig. 7 presents the extent to which the densest subnets vary from one distribution to another by pairwise comparing the datasets. Similarly, Fig. 7(right) presents the same result but for different time viewpoints for the *WU* dataset where timing information is available. Surprisingly, the discrepancy appears to be substantial. For example, for $k = 10$ the sets of the 10 densest subnets of the *WU* and the *Hotmail* distributions have only one subnet in common, while no common subnets exists when considering the *DShield* distribution! While the densest subnets across time for the same distribution appear to be more stable the discrepancy may still be nontrivial. *These observations clearly raise the problem of accurately estimating the true subnet distribution in the Internet.*

Finally we consider the impact of host partitioning in subnets on performance of the optimal strategy. To this extent, Fig. 8 presents the KL term and the number of samplings under the optimal strategy for various subnet definitions of the *WU* trace. Specifically, IP addresses are grouped to $/n$ subnets, with n ranging from 8 to 20. We find that while grouping

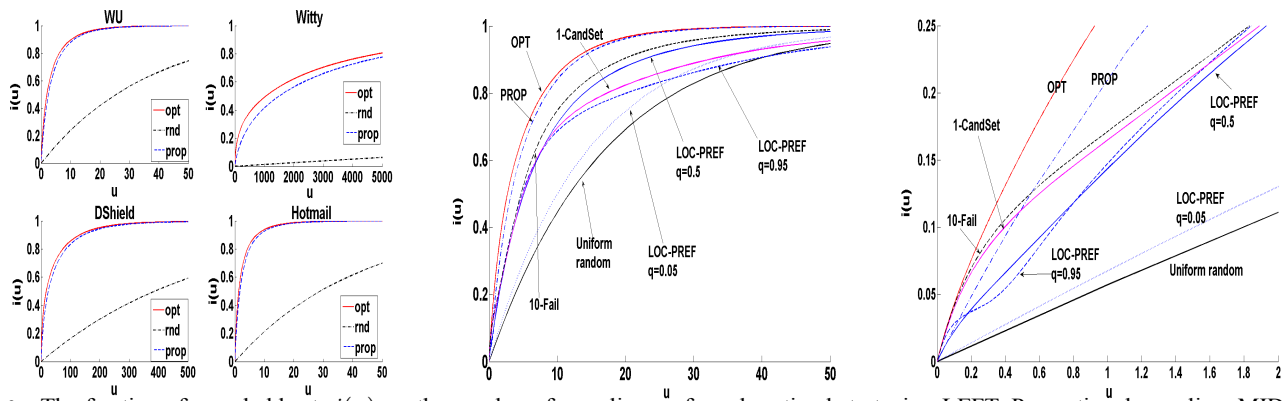


Fig. 9. The fraction of sampled hosts $i(u)$ vs. the number of samplings u for sub-optimal strategies. LEFT: Proportional sampling. MIDDLE: Local-knowledge strategies. RIGHT: Zooming in middle figure for small $i(u)$.

does not appear to have a significant effect on the KL term, the number of sampling increase as the groups become larger (i.e., n smaller). (This is expected from our analytical results as splitting subnets into sub-subnets indeed enlarges the set over which the optimisation is done.) This effect highlights that smaller groupings appear more attractive since effort will be concentrated in small populated subnets, in contrast to performing redundant samplings in large subnets which would be a lot sparser in general, but this may add to the complexity of identifying the optimum target set of subnets.

C. Sub-optimal strategies

Here, we examine the performance of our proposed strategies with respect to the optimal and the uniform random sampling strategies. We start by studying proportional sampling (PROP Section IV-B) in Fig 9(left), where we plot PROP vs. the optimal (OPT) and uniform random sampling (RND) for all datasets. In all cases, PROP follows closely the optimal. This is an interesting property as it suggest that already proportional sampling brings us very close to the optimal and may inform design of online sampling strategies.

Similarly, Fig 9(middle) and Fig 9(right) show the performance of sampling strategies that only take advantage of local knowledge (Section V) for the *WU* dataset. For reference, we also present the optimal, uniform random and PROP strategies. Fig 9(right) zooms in a particular range of Fig 9(middle) for small i . From these figures we can make the following observations for this specific dataset:

- Local subnet preference strategies perform close to (RND) for small q (local sampling probability). For larger q this strategy appears to suffer from persistently sampling exhausted subnets especially for larger i , thus performing close to (RND) or worse, while showing better performance for small $i(u)$.
- K-FAIL appears to consistently outperform all other local strategies and asymptotically follows the optimal strategy.
- CANDSET appears to suffer from similar issues with the local preference strategy, by persisting to scan subnets that have been exhausted.
- For smaller i (Fig 9, right), we see that both CANDSET and K-FAIL perform very close to the optimal (up to roughly $i = 0.1$). Note however, that $i = 0.1$ for the *WU* dataset

corresponds to a substantial number of hosts (over 10 million).

Overall, all our proposed strategies perform significantly better than uniform random sampling and local subnet preference strategy in the majority of the cases.

VII. CONCLUDING REMARKS

This paper studies the problem of epidemic-style information dissemination using random sampling. We identify optimal static and dynamic strategies to reach a target fraction of the host population in minimum number of samplings. We also propose and evaluate simple strategies that use no prior information and constant state, and provide significant gain. Future work may further investigate the space of simple strategies that perform near optimal.

REFERENCES

- [1] W. Vogels and C. Re, "WS-Membership Failure Management in a Web-Services World," in *Proc of the Twelfth International World Wide Web Conference*, 2003.
- [2] A. Demers, D. Greene, C. Hauser, W. Irish, J. Larson, S. Shenker, H. Sturgis, D. Swinehart, and D. Terry, "Epidemic algorithms for replicated database maintenance," in *Proceedings of the sixth annual ACM Symposium on Principles of distributed computing*, 1987.
- [3] L. Massoulie, A. Twigg, C. Gkantsidis, and P. Rodriguez, "Decentralized broadcasting algorithms," in *Proc. of the IEEE INFOCOM*, 2007.
- [4] Y. Xie, F. Yu, K. Achan, E. Gillum, M. Goldszmidt, and T. Wobber, "How Dynamic are IP Addresses?" in *Proc. of the ACM SIGCOMM*, 2007.
- [5] C. C. Zou, D. Towsley, and W. Gong, "On the performance of internet worm scanning strategies," *Performance Evaluation*, vol. 63, 2006.
- [6] Z. Chen and C. Ji, "A Self-Learning Worm Using Importance Scanning," in *ACM WORM '05*, Fairfax, Virginia, USA, 2005.
- [7] T. Kurtz, *Approximation of Population Processes*. CBMS-NSF Regional Conference Series in Applied Mathematics, 1981, vol. 36.
- [8] M. Vojnovic, V. Gupta, T. Karagiannis, and C. Gkantsidis, "Sampling Strategies for Epidemic-Style Information Dissemination," Microsoft Research, MSR-TR-2007-82, Tech. Rep., 2007.
- [9] G. Hardy and J. E. Littlewood and G. Pólya, *Inequalities*, 2nd ed. Cambridge Mathematical Library, 1952.
- [10] C. Gkantsidis, T. Karagiannis, P. Rodriguez, and M. Vojnovic, "Planet Scale Software Updates," in *Proc. of the ACM SIGCOMM*, 2006.
- [11] "DSshield: Cooperative Network Security Community," <http://www.dshield.org/index.d.htm>.
- [12] P. Barford, R. Nowak, R. Willett, and V. Yegneswaran, "Toward a Model for Source Addresses of Internet Background Radiation," in *Proc. of the Passive and Active Measurement Conference*, 2006.
- [13] Z. Chen and C. Ji, "Measuring network-aware worm spreading ability," in *Proc. of the IEEE INFOCOM*, 2007.
- [14] "Cooperative Association for Internet Data Analysis," <http://www.caida.org>.