# DefAT: Dependable Connection Setup for Network Capabilities

Soo Bum Lee, Virgil D. Gligor, Adrian Perrig

November 23, 2011

CMU-CyLab-11-018

CyLab Carnegie Mellon University Pittsburgh, PA 15213

# DefAT: Dependable Connection Setup for Network Capabilities

Soo Bum Lee Virgil D. Gligor Adrian Perrig CyLab, Carnegie Mellon University Email: {soobum, gligor, perrig}@cmu.edu

Abstract—Network-layer capabilities offer strong protection against link flooding by authorizing individual flows with unforgeable credentials (i.e., capabilities). However, the capabilitysetup channel is vulnerable to flooding attacks that prevent legitimate clients from acquiring capabilities; i.e., in Denial of Capability (DoC) attacks. Based on the observation that the distribution of attack sources in the current Internet is highly non-uniform, we provide a router-level scheme, named DefAT (Defense via Aggregating Traffic), that confines the effects of DoC attacks to specified locales or neighborhoods (e.g., one or more administrative domains of the Internet). DefAT provides precise access guarantees for capability schemes, even in the face of flooding attacks. The effectiveness of DefAT is shown in two ways. First, we illstrate the precise link-access guarantees provided by DefAT via ns2 simulations. Second, we show the effectiveness of DefAT in the current Internet via Interent-scale simulations using real Internet topologies and attack distribution.

## I. INTRODUCTION

Current service-flooding attacks rely on a large number of compromised machines that are organized as a "bot" network. Typical defense mechanisms that attempt to provide service-access guarantees despite such attacks assume absence of flooding in the underlying network links. Yet, a large-scale attack (e.g., a "botnet" with millions of "bots") can flood any chosen link in the Internet. In particular, defense mechanisms deployed at links near or at a network edge (e.g., Firewalls, IDSs) can be easily overwhelmed by such attacks. Worse yet, legitimate-looking attack packets can evade most of traditional techniques for handling address spoofing attacks at the network layer (e.g., IP tracebacks [1], [2], ingress filtering [3]).

Capability-based solutions, whereby distinct packet flows are separately authorized through capabilities obtained before flow initiation [4]–[6], provide congested routers with an effective way to prioritize legitimate flows and filter out unwanted traffic. Though promising, these solutions are still vulnerable to flooding attacks targeting the *capability-setup channel*, known as the Denial of Capability (DoC) attacks [7]. These attacks are possible because the initial capability-request packets are treated as best-effort packets, as opposed to the subsequent high-priority packets that carry capabilities. If DoC attacks cannot be countered, flow authorization via network-layer capabilities becomes impossible, and all access guarantees become meaningless at congested routers.

Previous solutions that attempt to protect capability requests

from flooding attacks (e.g., mechanisms based on aggregate request rates [6] or on proof of work [8]), though useful, are insufficient to provide dependable link-access guarantees for legitimate capability requests. For example, a fair-queueing mechanism, which fairly allocates buffer space to flow aggregates based on a router's confidence in precise identification of traffic origin [6], fails to provide *any* guarantee of link-access (viz., Section VII-A). Mechanisms based on proof of work (e.g., Portcullis [8]) provide only *weak* access guarantees during flooding attacks as they are (at best linearly) dependent on the number of global attack sources; e.g., a large number of bots could still flood a chosen link despite such guarantees. These previous schemes achieve relatively weak guarantees because they assume that attack sources are uniformly distributed in the network.

We observe, however, that malicious hosts, or bots are clustered: some domains include sufficiently strong security mechanisms that enable them to counter or deter contamination; others are easily contaminated by bots. Non-uniform distribution of attack sources is evident in a variety of worm propagation models [9], [10], evolutionary features of previous worms such as CodeRed I/II, Nimda and Slammer, and the disbribution of spam-bots [11] (viz., Section VIII-A). This non-uniformity actually enables us to achieve stronger guarantees. To be meaningful, these guarantees have to be independent of the number of attack sources (i.e., the size of a global botnet). In the worst case, they can only depend on attack sources in defined locales or neighborhoods (e.g., an administrative domain or a set of domains in the Internet). As a consequence, competing requests for a capability to a congested link that originate outside a contaminated locale should be unaffected, or only minimally affected, by a flooding attack, and should receive strong access guarantees. In contrast, initial capability requests originating from botcontaminated locales should receive weaker access guarantees, namely guarantees that depend only on the number of bots in the contaminated locale (but not on all bots of a multi-domain attack network). In short, our notion of dependable access to a flooded link provides differential guarantees for the capability setup channel. Differential access guarantees are desirable because they provide incentives for employing host security measures within administrative domains that prevent botnet (and other malware) contamination. In exchange, uncontaminated domains receive precise guarantees of link access for the capability setup channel, which support meaningful networklink and, ultimately, service-access guarantees.

Our scheme relies on three basic mechanisms. First, we define a new *path identification* mechanism that provides an unforgeable domain identifier to individual packets, and enables remote routers to identify a packet's domain of origin. Second, we define a *dynamic virtual queueing* mechanism that guarantees a minimum number of router buffer slots to domains originating flows through a router, which in effect, guarantees link access to those domains. Finally, we employ a *path aggregation* mechanism that optimizes router bandwidth allocation for legitimate capability requests based on domain contamination.

# II. BACKGROUND AND RELATED WORK

Lack of source address authenticity in the Internet Protocol (IP) enables attackers to forge the source addresses, and hence complicates/prevents address-based accounting during link flooding attacks. As a way to add authenticity to individual packets, capability solutions [4]–[6] have been proposed. Generally, a network-layer capability protocol requires a handshake between a client and a server, and during that phase, routers on the forwarding path collectively issue a connection capability; i.e., a series of router capabilities on the path. A router's capability, which is generated by hashing the source and destination IP address with the router's secret key, is cryptographically secure against forgeries since the router key is unavailable to an adversary.

However, the capability request protocol is still vulnerable to flooding (DoC) attacks [7]. That is, flooding with capability requests, which cannot be prioritized, successfully denies a legitimate access to a congested link. Portcullis [8] proposes a puzzle-based mechanism that provides a guaranteed link access during a flooding (DoC) attack. Though useful, the guarantee is linearly dependent on the number of bots, which can be substantial (e.g., the size of a botnet easily exceeds 1 million bots [12]). Alternatively, TVA's implementation of fair queueing on incoming traffic paths (i.e., hierarchical fair queueing) [6], which equally assigns queues to directly connected links and splits the queues recursively for distant links, places legitimate accesses of remote domains at a significant disadvantage since it provides fair service to the same level of queues (i.e., sub-queues split from a queue). More sophisticated application-layer solutions (e.g., CAPTCHA [13]) that attempt to distinguish between human- and machineinitiated traffic to prevent flooding attacks are impractical at the network-link level.

Attempts to block suspicious traffic upstream of a congested router by installing filters close to, or at, the domains originating attacks could protect legitimate flows that are independent of attacks. To be effective, cooperative filtering would require incentives that scale with the number of participating domains – a tall order since it depends on the attack itself. Furthermore, with only local information (the traffic rate of incoming links), a router cannot easily identify the links (or upstream links) that are responsible for the congestion; and even if such information is available, an adversary can launch a timed

attack where different groups of zombies/bots issue targeted requests by exploiting the time delay required for installing and releasing filters at upstream routers (e.g., on-off and rolling attacks).

# III. DESIGN OVERVIEW

In this section, we present an overview of our defense scheme by describing the basic mechanisms.

# A. Threat

The main threat we deal with in this work is a link flooding attack on the capability-setup channel, where attack sources collaboratively exhaust the link bandwidth allocated for connection establishment. We assume that both hosts and routers can be compromised and send/forward attack traffic. Compromised hosts are able to both flood a target link with capability request packets and disturb the path identification mechanism at a remote router by manipulating the header reserved for that purpose (viz., Section IV-A). Compromised routers can disturb path identification by either forwarding packets that contain false path-markings or adding invalid path-markings to the packets they forward.

# B. Path Identification

In this work, we consider routers that mark packets with path information. These path-markings create an unspoofable *origin* identifier because they cannot be controlled by endhosts.<sup>1</sup> In addition, path-markings enable remote routers to construct a traffic tree. The domain connectivity revealed in the traffic tree helps identify the distribution of attack sources in specified locales to which bandwidth allocation will be restricted (viz., Section VI).

The basic concept of route construction is similar to that of previous schemes [6], [14], yet we use a packet's AS (Autonomous System) path as a domain identifier for several reasons. First, a packet's AS-path, which is primarily determined by the number of AS hops (AS-path length) to the destination in the inter-domain routing protocol (e.g., BGP-4), is more stable than the routing path within an AS that may frequently change during flooding attacks due to link state changes (e.g., link failure). We use the AS-path of a packet as a persistent domain identifier. Second, a packet's AS-path can be constructed by the egress router of the source domain since the router contains the AS-path information of destination addresses in its routing table. This sourceconstructible domain identifier eliminates deployment issues that plagued previous path-marking schemes especially in the Internet core, and hence enables independent adoption of the marking scheme at the Internet border (e.g., provider/stub domains). We envision that prioritizing requests originating from path-marking domains would encourage early adoption of the marking scheme.

<sup>&</sup>lt;sup>1</sup>IP source routing may allow a client to select a path to a destination. However, strict and loose source routing are usually blocked at routers to avoid the associated processing overhead.

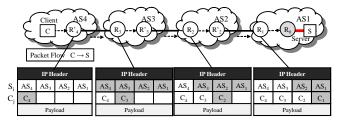


Fig. 1: Path Identifier.  $R_4'$  is the egress router of AS4 and  $R_3, R_2, R_1$  are the ingress routers of AS3, AS2, AS1 respectively.  $R_4'$  writes the path-identifier to the packet heading to server S in AS1, and ingress routers on the path validate the markings.  $C_j$  is the capability issued by  $R_j$ . Each ingress router can validate the shaded part of the markings.

We define a packet's AS-path to its destination as the *pathidentifier* of the packet, and present it in the order of markings: from the origin to the destination. Thus, as illustrated in Fig. 1, the path-identifier seen at a congested router in  $AS_1$  is  $\{AS_4, AS_3, AS_2, AS_1\}$ . We implement this path-identifier in a shim header so that only upgraded routers interpret it. Throughout this paper, we denote the path-identifier whose markings start with  $AS_i$  by  $S_i$  and the BGP speaker of  $AS_i$  by  $R_i$ . In Section IV, we present a mechanism that protects pathidentifiers from potential attacks (e.g., spoofing and replay attacks).

### C. Link Access Guarantees

In defending against DoC attacks, our goal is to provide precise guarantees of link access to capability requests, where the guarantees are provided in a domain basis to confine the effects of attacks within the domains originating attack traffic. This goal is achieved by a new fair queueing mechanism that allocates separate buffer slots to individual domains. And, the guarantees provided by the queueing mechanism are optimized to favor the requests from uncontaminated domains by bots, using a path aggregation mechanism.

1) Fair Queueing Revisited: The use of a fair queueing scheme for link-access guarantees is intended to maximize service on the legitimate capability requests. Fair queueing schemes, if they can assign separate queues to individual path-identifiers, could provide fair bandwidth to the pathidentifiers without link under-utilization (which could occur whenever strict bandwidth reservation is made to individual path-identifiers). However, when the spatio-temporal dynamics of domains contributing to congestion (e.g., time-varying patterns of domain traffic) are considered, such queue assignment in a limited buffer is a challenging problem. For example, for a fixed buffer size, under-provisioning of the number of queues in a specific time period may fail to provide linkaccess guarantees to path-identifiers due to potential queue collisions among different path-identifiers. In contrast, overprovisioning of it would decrease the length of individual queues, hence weaken the guarantees (viz., Section V). Thus, we aim to design a fair queueing scheme that assigns a unique queue to each path-identifier and adjusts the individual queue lengths to fit the buffer size for link-access guarantees and their enhancement – a desired goal.

While a variety of traditional fair queueing schemes focus on the bandwidth fairness of flows in different queues that contain various sizes of packets, the Stochastic Fair Queueing (SFQ) scheme [15] offers queue length fairness via a buffer stealing mechanism, whereby a packet that finds a full buffer on its arrival would steal a buffer-slot from the longest queue. We note that the fixed size capability request packet would eliminate the intrinsic bandwidth unfairness of SFQ in the presence of different packet sizes [16]. Based on the bufferstealing idea, we improve SFQ in two respects. First, we avoid queue collisions among path-identifiers that are allowed but fairly distributed via stochastic queue assignment in SFQ. Second, we make queue management operations (e.g., queue assignment and buffer-slot preemption) scalable and efficient to easily adapt our scheme to diverse operating environments (e.g., link capacity, the number of required queues). Those improvements are made via a dynamic virtual queueing mechanism presented in Section V.

2) Path Aggregation: As more domains are contaminated by attack sources, link-access guarantees provided by our queueing scheme become weak as both available link bandwidth and buffer-slots to each path-identifier decrease. This undesirable dependency of guarantees on attack dispersion is unavoidable as long as all path-identifiers are equally treated. Protecting requests of uncontaminated domains essentially needs a differential treatment of path-identifiers based on the proportion of legitimate requests they deliver. Though the legitimacy of individual capability requests cannot be validated, the proportion of legitimate requests in a set of requests can be estimated using a couple of flow conformance tests. These tests consist of (1) a test on bandwidth conformance that represents the aggressiveness of requests and (2) a test on protocol conformance that indicates the legitimacy of authorized flows in various respects (viz., Section VI-A).

Conformance tests performed on each path-identifier enables differential assignment of bandwidth to path-identifiers that maximizes service to legitimate requests at the flooded link. Yet, in the presence of a large number of attack domains, such assignment cannot easily be made, nor can it tolerate imprecise measurement of domain contamination. Instead, we aggregate the path-identifiers of a highly contaminated locale and assign a new path-identifier to them. This, in effect, limits both available bandwidth and buffer space for those path-identifiers. We define this path aggregation problem as a constrained optimization problem and provide an efficient solution in Section VI-C.

# IV. PATH IDENTIFICATION

In this section, we first describe the basic path identification mechanism, and then enhance the mechanism with additional security features.

The basic path identification mechanism works as follows. When the egress router of a domain (i.e., the BGP speaker) forwards a packet that originates from its domain, it writes

the path-identifier (i.e., the AS-path to the destination) in the packet's header. AS ingress routers of the packet forwarding path validate the authenticity of a fraction of this path-identifier starting with the upstream AS that forwarded the packet and ending with the destination AS as shown in Fig. 1. Whenever AS ingress routers receive a non-marked packet, they write their own path-markings: the AS-path from their upstream AS to the destination AS.

As remote domains can validate only a part of pathmarkings, attack sources in non-path-marking domains may spoof path-identifiers unless the marking scheme (which includes the verification function) is sufficiently deployed. Even under wide deployment of the marking scheme, the authenticity of path-identifiers verified at a domain cannot be delegated to the downstream domains without a strong trust relationship established between those domains. This makes any manipulation of path-identifiers by compromised routers undetectable at remote routers. To protect path-identifiers from these attacks (i.e., spoofing and replay attacks), we present a secure path identification mechanism below.

# A. Unspoofable path-identifier

We first introduce potential attacks that disturb pathidentification at remote routers and present our defense mechanism against those attacks.

Let  $\{AS_n, \ldots, AS_2, AS_1\}$  be the path-identifier seen at the congested router, and let \* and # be any valid and forged sequence of markings respectively. If the domains up to  $AS_i$  are unprotected by our path identification (which includes both path marking and verification) scheme, both compromised sources in non-path-marking domains and compromised routers in  $AS_k$  can forge a path-identifier as  $\{\#, AS_i, *, AS_1\}$ .

In principle, a router can authenticate its path-markings to other routers by adding a digital signature to the path-markings. However, adding a different digital signature to every packet would impose significant computational overhead for both its generation and verification. Moreover, a per-packet signature, if employed, could be exploited by attackers to exhaust routers' computational resource (e.g., by flooding small-size packets). To reduce authentication overhead, we design an efficient path-identifier authentication mechanism, where each domain pre-distributes its domain-authenticator and uses it to authenticate its path-markings. One fundamental assumption for implementing this mechanism is that any protected AS has a public-private key pair certified by a trusted certificate authority (e.g., ICANN [17]).

1) Authenticator Distribution: When a BGP speaker advertises an address prefix that belongs to its domain, the BGP speaker adds an origin authentication number (OAN), which is unique in its domain and is digitally signed with the domain's private-key, to its route advertisement. All BGP routers that receive this route advertisement authenticate the OAN using the origin AS' public-key and hold the authenticated ASN (AS Number)-OAN pair for later path-identifier authentication.

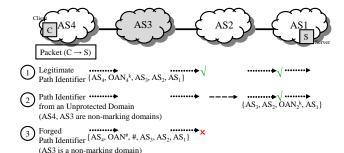


Fig. 2: Path-identifier Authentication. ① Path-identifier written at the packet's origin (AS4) can be validated at any domain (AS2, AS1) in the presence of a non-marking domain(s) (AS3) on the packet's forwarding path. ② If the origin AS does not participate in path-marking, the first participant (AS2) writes its markings and adds the incoming AS number (AS3) to distinguish the packets it forwards from the ones originating from it. ③ An invalid ASN-OAN pair (denoted by #) can be detected and filtered.

Since the number of ASN-OAN pairs is at most 65,535<sup>2</sup>, the space requirement for this validation is bounded, i.e., 262KB for 4-Byte OANs.

2) Origin Authentication: The BGP speaker of a packet's domain of origin writes its ASN-OAN pair followed by the AS-path to the destination in the path-identifier header. Fig. 2 illustrates the cases for origin authentication under different deployment scenarios of the marking scheme. Whenever no path-identifier is present in a packet, the ingress router of a marking AS constructs path-markings with its own ASN-OAN pair (viz., ② in Fig. 2). On receiving path-identifiers constructed as such, the AS ingress routers on the later path validate the origin's OAN and the partial AS-path as discussed above. In this way, the routers on the way to, or at the destination AS can identify any forged path-markings by adversaries even in the presence of consecutive non-pathmarking ASs on the path and filter packets carrying those forged path-markings.

While a compromised router in  $AS_i$  can still forge two valid types of path-identifiers such as  $\{AS_i, OAN_i^k, *\}$  and  $\{\#, AS_i, OAN_i^k, *\}$ , their effects can be limited to at most those of two path-identifiers by discarding the non-authenticated prefixes of path-identifiers.

# B. Preventing Replay Attacks

Under partial deployment of our path-marking scheme, attack sources in non-path-marking domains may forge path-identifiers ending with authenticated ASN-OAN pairs (since ASN-OAN pairs are not confidential to end-hosts) and use them in flooding a target link. Such replay attacks would

<sup>&</sup>lt;sup>2</sup>As of 2010, the number of advertised ASNs is about 35,000 out of 65,535 (16-bit) possible ASNs [18]

<sup>&</sup>lt;sup>3</sup>For path validation, routers need to keep AS-path information (from next hop to the destination AS) in their forwarding table (i.e., FIB). However, this would not require much space since the average number of ASs a packet traverses from its origin to destination is four.

significantly affect the requests from path-marking domains and hence prevent those domains from receiving incentive for early adoption of the path-marking scheme.

Path-marking routers counter replay attacks via fast OAN renewals, which are efficiently implemented using a reverse hash chain [8]. Let  $OAN_i^0$  be the initial OAN of  $AS_i$ .  $AS_i$ constructs a hash chain of OANs by repeatedly hashing  $OAN_i^0$  with a cryptographic hash function (i.e.,  $OAN_i^k =$  $\operatorname{Hash}(\operatorname{OAN}_i^{k-1}||AS_i||k-1)$  for  $1 \leq k \leq M$ ), and distributes  $OAN_i^M$  when advertising a route. We engage  $AS_i$  and k-1in generating OAN to produce distinct OAN sequences for different ASs and initial OANs respectively. A BGP speaker uses  $OAN_i^k$  during a predefined interval; and changes it to  $OAN_i^{k-1}$  in the next interval. Hence, without breaking the hash function, an attacker cannot construct the valid sequence of  $OAN_i^k$ s to be used. A (ingress) router can authenticate  $\operatorname{OAN}_i^k$  by computing  $\operatorname{Hash}(\operatorname{OAN}_i^k||AS_i||k)$  and comparing it with  $\operatorname{OAN}_i^{k+1}$ . This OAN authentication is performed only once for every OAN renewal. Once  $OAN_i^k$  is used,  $OAN_i^{k+1}$ is invalidated. Note that if the OAN renewal period is less than the time required for replaying OANs, replay attacks will be effectively prevented. The length of a OAN hash chain (M) is determined in consideration of the OAN renewal period to avoid frequent OAN distribution. For example, if a 20-bit sequence number (M  $\approx$  1 million) and 500ms OAN renewal period are used, a domain needs to advertise its OAN once in every six days. We also note that routers in different domains need not be time-synchronized as an OAN carries its sequence number that is specific to the domain.

# V. DYNAMIC VIRTUAL QUEUEING

In this section, we describe a dynamic virtual queueing mechanism for link-access guarantees on path-identifiers. Our dynamic virtual queueing mechanism is designed to assign a separate queue to active path-identifiers and provide queue length fairness to the path-identifiers in a minmax manner. For these purposes, a router manages virtual queues rather than physically separate queues, that are distinguished by the path-identifier  $(S_i)$ , its count at time t  $(N_{S_i}(t))$  and packet location (memory address)  $(A_{S_i})$  in the buffer; i.e.,  $(S_i, N_{S_i}(t), A_{S_i})$ . Given these tuples and the buffer size  $L_Q$ , queue-length fairness on path-identifiers  $(\min \max_{S_i \in \mathcal{S}} N_{S_i}(t) \text{ for } \sum_{S_i \in \mathcal{S}} N_{S_i}(t) = L_Q) \text{ can be de-}$ scribed by the following buffer-slot preemption policy. If a packet finds the buffer full on its arrival, it preempts a bufferslot from the longest virtual queue. If the arrived packet belongs to the longest virtual queue, or its preemption produces another longest virtual queue, the packet would be dropped. This preemption policy ensures guaranteed buffer-slots to each path-identifier if the number of buffered path-identifiers is bounded. We assume that the number of buffered pathidentifiers can be statistically or deterministically bounded at a router (i.e., the minimum bandwidth to a legitimate pathidentifier can be determined).

# A. Implementing Buffer-slot Preemption

For efficient and scalable accounting of virtual queue lengths, we use a new Counting Bloom Filter (CBF) that holds the number of buffer-slots occupied by path-identifiers and provides lookup, add and remove operations in O(1) time (a modified version of CBF [19]). CBF consists of m counter arrays of size  $2^b$   $(a_1, a_2, \ldots, a_m)$  and m hash functions of b-bit output  $(H_1, H_2, \ldots, H_m)$ , where  $a_i$  is associated with  $H_i$ . For an input to CBF, each hash function maps its output to the corresponding array position; e.g.,  $a_i[H_i(S_1)]$  corresponds to the input  $S_1$  for  $1 \le i \le m$ .

Path-identifier accounting in CBF works as follows. All array values are initialized to zero. When a packet is added to the buffer, its path-identifier is fed into CBF. Then, CBF locates m array positions for the path-identifier, and increases the corresponding array values. The same applies to a packet removal from the buffer, yet the counter values are decreased. In this scheme, the limited hash output size (i.e.,  $2^b$ ) could cause hash-output collisions among path-identifiers. Such collisions would make corresponding array values increased by multiple path-identifiers, hence corrupted. However, unless all of the array values associated with  $S_i$  are corrupted, we can compute the count of buffered  $S_i$ 's by taking the minimum of the array values; i.e.,  $\min\{a_1[H_1(S_i)], a_2[H_2(S_i)], ..., a_m[H_m(S_i)]\}$ . Since the probability that all m array values of a path-identifier are corrupted is  $(1 - (1 - (1/2^b))^{|S|})^m$  for |S| buffered pathidentifiers [20], we can make the probability negligible by increasing the array size  $(2^b)$  or the number of arrays (m).

Path-identifiers that occupy more buffer slots than the guaranteed amount (i.e.,  $\lfloor \frac{L_Q}{|S|} \rfloor$ ) should be kept track of for possible preemption. To this end, a router maintains a table, named Path-Identifier Record (PIR), that holds *over-buffered* path-identifiers, their counts and corresponding packet locations. In PIR, a path-identifier is stored as the concatenation of its m hash outputs, defined as "path-signature." This enables fast buffer-slot preemption because the preempted packet's path-signature would directly locate array values that need to be decreased in CBF.

# B. Probabilistic Guarantees

If packet arrivals carrying path-identifier  $S_i$  are modeled as a Poisson process and k buffer-slots are allocated to  $S_i$ , the probabilistic lower bound of  $S_i$ 's link access (denoted by  $\mathcal{G}(|\mathcal{S}|,k,S_i)$ ) is provided as follows.

$$\mathcal{G}(|\mathcal{S}|, k, S_i) = \tag{V.1}$$

$$\begin{cases} \sum_{j=0}^{k-1} \frac{\left(k\rho_{S_i}\right)^j}{j!} e^{-k\rho_{S_i}} & \rho_{S_i} < 1\\ \frac{1}{\rho_{S_i}} (1 - \mathcal{G}_{\mathcal{L}}) (1 - \sum_{j=k}^{\infty} \frac{\left(k\rho_{S_i}\right)^j}{j!} e^{-k\rho_{S_i}} {j \choose k-1} \mathcal{G}_{\mathcal{L}}^{k-1}) & \rho_{S_i} \ge 1 \end{cases}$$

where  $\lambda_{S_i}$  is the request rate of  $S_i$ ,  $\rho_{S_i} = \frac{\lambda_{S_i}|\mathcal{S}|}{C_R}$  is the bandwidth utilization of  $S_i$ , and  $\mathcal{G}_{\mathcal{L}} = \sum_{j=0}^{k-1} \frac{(k\rho_{S_i})^j}{j!} e^{-k\rho_{S_i}}$ .

We justify the Poisson arrival model of capability requests with two reasons: (1) during the *short* interval that the guarantees are defined (i.e., the maximum queueing delay of a

router  $\Delta_Q$ ), the capability requests by different clients can be assumed independent; and (2) a single capability can be used for multiple *correlated* sessions that need to be established for most Web applications (that is, multiple correlated capability requests are unnecessary). Under this model, if  $\rho_{S_i} < 1$ , an arrival of  $S_i$  is guaranteed to be serviced if less than k arrivals of  $S_i$  has occurred in  $\Delta_Q$ . If  $\rho_{S_i} \geq 1$ , an arrival of  $S_i$  is guaranteed to be serviced only if its queue length is less than k. Thus, Eq. (V.1) can be easily proved. The probabilistic guarantee of  $S_i$ 's link-access is provided by setting  $|\mathcal{S}| = |\mathcal{S}|_{max}$  and  $k = \lfloor \frac{L_Q}{|\mathcal{S}|_{max}} \rfloor$ . We provide a full proof in Appendix A.

# C. Resource Requirements

1) Request Packet Buffer: A large buffer  $(L_Q)$  for capability request packets is preferable since it would not only improve the guarantees (viz., Eq. (V.1)) but also handle the requests from spontaneously created, short-lived paths. However, the size of the buffer should be bounded in consideration of the maximum allowed queueing delay to avoid unnecessary retries from flow sources. For example, if we assume 0.25 second maximum queueing delay and  $128B^4$  request packet size, for a 2.5 Gbps link  $^5$ , a router requires 4.0 MB buffer (when 5% of link bandwidth is allocated for capability requests [6]), and with which it can provide 8 guaranteed buffer slots up to 3.75K path-identifiers.

2) Path-Identifier Accounting: The memory requirement for CBF is determined by a target false-positive ratio. The false positive ratio of a CBF is determined by  $\left(1-\left(1-\frac{1}{2^b}\right)^{|\mathcal{S}|}\right)^m \approx \left(1-e^{-\frac{|\mathcal{S}|}{2^b}}\right)^m = \left(1-e^{-\frac{L_Q}{k\cdot 2^b}}\right)^m$  since  $L_Q = k\cdot |\mathcal{S}|$ . Hence, for a desired false positive ratio, the size of each counter array in CBF, which is same as the size of hash output  $(2^b)$ , is linear with the buffer size (i.e.,  $\Theta(L_Q)$ ). For example, a CBF with 8 hash functions of 14-bit outputs would require  $8\times 2^{14}$  (hash outputs)  $\times 2^8$  (counter) = 131KB memory space while producing a reasonably low false positive ratio  $(3.07\times 10^{-4}\%)$  in the presence of 3.75K path-identifiers.

PIR holds the path-identifiers whose count exceeds  $\lfloor \frac{L_Q}{|S|} \rfloor$  for possible preemption. Hence, the memory requirement is bounded by  $L_Q/(k+1) \times (16 \text{B (path-signature)} + 4 \text{B (address pointer)})$  (e.g., 60 KB for the above example), since the number of path-signatures in PIR has its maximum when all path-identifiers have k+1 packets in the buffer. Hence, the memory requirement for both CBF and PIR is  $\Theta(L_Q)$ .

# VI. PATH AGGREGATION

In this section, we first describe a mechanism for estimating the proportion of legitimate requests of individual pathidentifiers, and then, a path-identifier aggregation mechanism that maximizes the *goodput ratio*, defined as the proportion of legitimate requests in all *serviced* requests, at a congested link. Aggregating path-identifiers produces an optimal traffic

tree to which applying our queueing mechanism maximizes goodput ratio at the congested link.

# A. Goodput Estimation

In absence of any other useful information regarding the origin of attack sources and the path-identifiers assigned to them, the request rate of path-identifier  $S_i$  ( $\lambda_{S_i}$ ) can be used as a unique measure for estimating the goodput ratio of  $S_i$ . We define the *bandwidth conformance* of path-identifier  $S_i$  as  $\min\{1, \frac{C_R}{\lambda_{S_i}|S|_{max}}\}$  to represent how the request rate of  $S_i$  conforms to the assigned bandwidth to it, and denote it by  $\mathcal{E}_{R_i}^{\mathcal{B}}$ , i.e.,  $\mathcal{E}_{R_i}^{\mathcal{B}} = \min\{1, \frac{C_R}{\lambda_{S_i}|S|_{max}}\}$  (recall that  $S_i$  is assigned to all packets originating from  $R_i$ ).

Additionally, we estimate domain contamination more accurately by identifying the following attack flows.

Unauthorized flows: A capability issued by a router during the connection establishment phase of a flow must be used at least once for actual data transmission unless it is denied afterward by application services, firewalls or IDSs. Thus, the proportion of unused capabilities could effectively measure domain contamination as it reflects the strong flow authorization results applied at the network ends.

High-rate flows: Flows that send high-rate traffic using valid capabilities would exhibit high packet-drop rates as indicated in [21]. Hence, if a router implements per-domain bandwidth control,<sup>6</sup> high-rate attack flows within a domain can be identified by capability drop rates [22].

High-fanout sources: If sources are allowed to establish an unlimited number of connections with other destinations through the congested link, they can deplete link's bandwidth with a large number of legitimate-looking flows [23]. This insidious attack will be prevented if a router limits the number of per-source capabilities as follows.

Let  $C_{f_{s,d}}$  be the capability for a flow  $f_{s,d}$  between a source s and a destination d.  $C_{f_{s,d}}$  consists of two parts, namely  $C_{f_{s,d}} = C_{f_{s,d}}^0 || C_{f_{s,d}}^1$ . Here,  $C_{f_{s,d}}^k$  is defined as:

$$\begin{array}{lcl} \boldsymbol{C}_{f_s,d}^0 & = & \operatorname{Hash}(\operatorname{IP}_s,\operatorname{IP}_d,K_R^1) \\ \boldsymbol{C}_{f_s,d}^1 & = & \operatorname{Hash}(\operatorname{IP}_s,\operatorname{f}(\operatorname{IP}_d),K_R^2) \end{array}$$

where  $IP_s$  and  $IP_d$  are the source and destination IP addresses,  $K_R^0$  and  $K_R^1$  are the router's secret keys, and  $f(\cdot)$  is a function whose output is randomly uniform on  $[0, n_{max}$ -1].

 $C_{f_{s,d}}^0$  provides identifier authenticity to flows [5], [6], and  $C_{f_{s,d}}^1$  restricts the number of per-source capabilities to  $n_{max}$  by taking  $f(IP_d)$  as a hash input. If  $C_{f_{s,d}}^1$  is used for estimating flow bandwidth, flows of high-fanout sources would be aggregated and turn into high-rate flows.

The above attack-flow identification measures help estimate the proportion of legitimate flows in flows carrying  $S_i$ , which we define as the *protocol conformance* of  $S_i$  and denote by  $\mathcal{E}_{R_i}^{\mathcal{P}}$ .

<sup>&</sup>lt;sup>4</sup>We reserve 88B shim header: 40B for path-identifiers (up to 10 AS markings), 8B for an origin authenticator and 40B for 5 capabilities.

<sup>&</sup>lt;sup>5</sup>2.5Gbps (OC–48) links are widely used for ISP's backbone links.

<sup>&</sup>lt;sup>6</sup>Flows in different domains could exhibit different drop rates due to different RTTs.

Based on the bandwidth and protocol conformances, the conformance estimate  $\mathcal{E}_{R_i}$  of  $S_i$ , representing the estimate of  $S_i$ 's goodput ratio, is defined as:

$$\mathcal{E}_{R_i} = e^{-\frac{\gamma \cdot \lambda_{S_i} |\mathcal{S}|_{max}}{C_R}} (\mathcal{E}_{R_i}^{\mathcal{B}} - \mathcal{E}_{R_i}^{\mathcal{P}}) + \mathcal{E}_{R_i}^{\mathcal{P}}$$

$$\mathcal{E}_{R_i}(t_i) = (1 - \alpha)\mathcal{E}_{R_i} + \alpha\mathcal{E}_{R_i}(t_{i-1})$$

where  $\gamma$  and  $\alpha$  are the weighting coefficients.

The conformance estimate of  $S_i$  is the weighted average of the bandwidth conformance and the protocol conformance, where the weighting factor exponentially favors the protocol conformance as sufficient requests have been made. In this way, we prevent a domain's conformance estimate from being highly biased by its (low) request-rate; e.g., unexpected packet drops of a low-rate path-identifier would produce a very low protocol conformance for the corresponding domain. We determine  $\mathcal{E}_{R_i}$  at time  $t_j$  by taking the moving average of  $\mathcal{E}_{R_i}$ s, and update it once in every aggregation period ( $\Delta_{aqq}$ ); i.e.,  $t_j - t_{j-1} = \Delta_{agg}$ .

# B. Aggregation Problem

For path aggregation, the congested router  $R_0$  builds the traffic tree  $\mathcal{T}_{R_0}$  using the path identifiers carried in the *active* flows and decomposes  $\mathcal{T}_{R_0}$  into a legitimate tree  $\mathcal{T}_{R_0}^{\mathcal{L}}$  and an attack tree  $\mathcal{T}_{R_0}^{\mathcal{A}}$ .  $\mathcal{T}_{R_0}^{\mathcal{L}}$  is constructed with legitimate path-identifiers that have a higher conformance estimate than a certain threshold  $(\mathcal{E}_{th})$ , and  $\mathcal{T}_{R_0}^{\mathcal{A}}$  is constructed with the other (non-legitimate) path-identifiers. Then, the router constructs a new traffic tree  $\mathcal{T}'_{R_0}$  by merging those two trees at the root (i.e., the disjoint union of  $\mathcal{T}^{\mathcal{L}}_{R_0}$  and  $\mathcal{T}^{\mathcal{A}}_{R_0}$ ). Path aggregation is performed on this new traffic tree  $T'_{R_0}$ , so that legitimate paths would never be aggregated with attack paths.

The congested router starts path aggregation from neighboring domains (i.e., domains with longest suffix-matching path-identifiers) to localize attack effects, and proceeds with aggregation until a desired number of path reductions are made (viz., Eq. (VI.1)). Aggregation is performed with respect to the conformance estimate of each path since link-access guarantees should not be biased by the path's request rate. Hence, if the number of access-guaranteed path-identifiers is  $|\mathcal{S}|_{max}$ , the path aggregation problem is to construct an optimal tree which has  $|S|_{max}$  distinct paths and to which providing link-access guarantees maximizes goodput ratio at the congested link. This can be defined as a constrained optimization problem below.

Let  $\mathcal{R}$  be the set of all nodes in  $\mathcal{T}'_{R_0}$ , and  $\mathcal{R}_i$  be the set of leaf nodes of a subtree rooted at  $R_i \in \mathcal{T}'_{R_0}$  (i.e.,  $\mathcal{T}_{R_i}$ ). Then, the optimization problem is defined as:

$$\max O(\mathcal{T}'_{R_0}) = \max \sum_{R_i \in \mathcal{R}} \frac{1}{|\mathcal{R}_i|} \sum_{R_j \in \mathcal{R}_i} \mathcal{E}_{R_j}$$
 (VI.1)

subject to 
$$\sum_{R_i \in \mathcal{R}} I_{\mathcal{R}_i} \leq |\mathcal{S}|_{max}$$
 and  $\bigsqcup_{R_i \in \mathcal{R}} \mathcal{R}_i = \mathcal{R}_0$ 

where  $I_{R_i}$  equals 1, if paths are aggregated at  $R_i$ , and 0, otherwise. For a non-aggregated path,  $I_{R_i}$  is 1 at the leaf node. Since  $\sum_{S_i \in S} I_{R_i}$  is the number of path identifiers seen at  $R_0$ , it should be bounded by  $|\mathcal{S}|_{max}$ .

In the above equation, aggregation at  $R_i$  decreases the total conformance estimate by  $\frac{|\mathcal{R}_i|-1}{|\mathcal{R}_i|}\sum_{R_j\in\mathcal{R}_i}\mathcal{E}_{R_j}$ . We define this value as the aggregation cost and denote it by  $C^{\mathcal{A}}(R_i)$ ; i.e.,  $C^{\mathcal{A}}(R_i) = \frac{|\mathcal{R}_i|-1}{|\mathcal{R}_i|}\sum_{R_j\in\mathcal{R}_i}\mathcal{E}_{R_j}$ . Hence, a set of nodes at which aggregating path-identifiers produces the minimum (total) aggregation cost, would be a solution to the above problem.

We note that, if the set of aggregating nodes (routers) are fixed, the optimization problem of Eq. (VI.1) is the same as the 0-1 knapsack problem<sup>7</sup> which is known to be NP-complete. In Eq. (VI.1), however, the set of aggregating nodes and the relative aggregation cost of a leaf node  $(\frac{|\mathcal{R}_i|-1}{|\mathcal{R}_i|}\mathcal{E}_{R_j},\,R_j\in\mathcal{R}_i)$  vary as aggregation proceeds to the root. This means the 0-1 knapsack problem should be solved repeatedly as the set of aggregating nodes is redefined. We present an efficient algorithm for this problem below.

# C. Aggregation Algorithm

Whenever aggregation is necessary (i.e.,  $|\mathcal{S}| > |\mathcal{S}|_{max}$ ), aggregation is performed as summarized in Algorithm 1. Let  $\mathcal{O}$  be the solution set and  $\mathcal{C}$  be the candidate set. Initially,  $\mathcal{O}$  is empty and  $\mathcal{C}$  has all intermediate (i.e., non-leaf) nodes in  $\mathcal{T}'_{R_0}$ as its elements. Then, the algorithm works in a *greedy* fashion: for each iteration, the node that causes the lowest cost-decrease to  $\mathcal{C}$  is added to  $\mathcal{O}$ , and this continues until the constraint on the number of path identifiers in Eq. (VI.1) is satisfied. Though Algorithm 1 is a greedy approximation algorithm, it ensures that the total cost of the candidate set decreases minimally at each iteration. As a consequence, its approximation error from the optimal aggregation cost is bounded by the number of incoming links of the last added node to  $\mathcal{O}$ . We provide the proof on the error bound in Appendix B.

# Algorithm 1 Aggregation

- 1: Set  $\mathcal{O} = \emptyset$  and  $\mathcal{C} = \{R_i | R_i \in \mathcal{T}'_{R_0} \mathcal{R}_0\}$ . 2: Move the lowest aggregation cost node in  $\mathcal{C}$  to  $\mathcal{O}$ .
- 3:  $R_i \in \mathcal{C}$  replaces the current solution set if it satisfies the following replacement conditions:
  - $C^{\mathcal{A}}(R_i) < \sum_{R_i \in \mathcal{O}} C^{\mathcal{A}}(R_j)$
  - $C^{\mathcal{A}}(R_i) > \max_{R_i \in \mathcal{O}} C^{\mathcal{A}}(R_j)$
- 4: Repeat steps 2 and 3 until the constraint on the number of pathidentifiers (in Eq. (VI.1)) is satisfied.

# VII. SIMULATION RESULTS

In this section, we present our ns2 simulation results for various attack scenarios to evaluate our design. Network topologies for simulations are configured to capture the worst case effect of different attacks and to ascertain how well our

 $^{7} rac{C^{\mathcal{A}}(R_{i})}{|\mathcal{R}_{i}|}$  can be considered as the unit value of an element,  $|\mathcal{R}_{i}|$  as the size of an element, and  $|\mathcal{S}| - |\mathcal{S}|_{max}$  as the knapsack size in the 0-1 knapsack

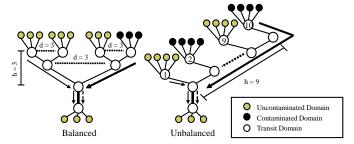


Fig. 3: Topology used in simulation. Legend: "d" is the number of sibling nodes and "h" is the tree height.

design goals are satisfied. The balanced tree shown in Fig. 3 is used for simulations that evaluate the access guarantees and the effectiveness of aggregation. The unbalanced tree is used to show that our scheme effectively provides access guarantees to domains independently of their location on a routing path. We assign 5% of link capacity to the capability request channel as in [6]. In most simulations, the total request rate of legitimate sources is set close to the link capacity of request channel (i.e.,  $\rho_{S_i} \approx 1$  for legitimate domains) to accurately capture the effects of attacks. Requests are randomly placed during the specified simulation interval to approximate Poisson arrivals.

We compare our simulation results with those of TVA [6], which protects capability requests using a hierarchical fair-queueing mechanism.

# A. Link-Access Guarantees

To evaluate the local effect of flooding attacks in our scheme, we use a 27-path balanced tree, where 30 legitimate sources are attached to each leaf node, and attack sources are increased at a leaf node. In this simulation, we set the number of access-guaranteed paths ( $|S|_{max}$ ) to 27 and the buffer size to that of 108 packets so that 4 buffer-slots are guaranteed to each path. Each source randomly starts 100 different sessions (which is equivalent to 100 times more sources) between 0 and 10 seconds. This source configuration is used for entire simulations. We also run simulations with a TVA [6] router configured to have 1000 queues of length 4 (as TVA requires distinct queues for individual sources in the current implementation) for comparative evaluation.

As Fig. 4 shows, the request drop ratios of legitimate paths are stable over the wide range of attack sizes with both our scheme and TVA. That is, both schemes effectively localize flooding attacks when compared with the *no defense* case. Note that a per-client defense would have the same result as that of no defense when bots are used to flood the link. Yet, our scheme outperforms TVA with a much smaller buffer (108 vs. 4000 buffer-slots). This is because our scheme dynamically adjusts virtual-queue lengths in a *min-max* manner, which in effect allows more than the guaranteed buffer-slots to path-identifiers unless their bursts are synchronized (in which case, only the guaranteed buffer-slots hold).

To illustrate the robustness of the guarantees that our scheme

provides, we configure an extreme adversarial scenario where 60 paths of a 64-path balanced tree (i.e., h=3 and d=4 in Fig. 3) send a large number of requests, and observe the service ratio of the remaining 4 paths. Fig. 5 shows the probabilistic guarantee  $(\mathcal{G}(|\mathcal{S}|,k,S_i),$  viz., Eq. (V.1)), the stationary service probability  $(P(|\mathcal{S}|,k,S_i))^8$ , and the simulation result  $(Pr(|\mathcal{S}|,k,\mathcal{S}^{\mathcal{L}}))$  for the set of legitimate path-identifiers  $\mathcal{S}^{\mathcal{L}}$ , under specified bandwidth utilizations – the ratio of request rate to an allocated bandwidth. Even under this extreme attack scenario, the service ratio of legitimate paths is close to the theoretical stationary packet service probability, which is much higher than the probabilistic guarantees, as illustrated in the figure.

Next, we show that link-access guarantees provided by our scheme are independent of attack location. For this simulation, we use a 40-path unbalanced tree shown in Fig. 3. We attach 30 legitimate sources to each leaf node, and 200 attack sources to each of eight attack nodes; four of these nodes are placed at different locations for each simulation and the remaining four nodes are placed at the farthest location from the flooded link. In this scenario, we simulate the queue implementation for  $\mathcal{G}(34, 8, S_i)$ ,  $\mathcal{G}(64, 4, S_i)$  and  $\mathcal{G}(64, 8, S_i)$ , and those for the corresponding 4 and 8-slot queues in a TVA router (i.e., 4000 and 8000 total buffer-slots respectively). Fig. 6 shows the request drop ratios of legitimate paths, where the horizontal axis represents the index of attack location (viz., unbalanced tree in Fig. 3). With our scheme, the request drop ratios are uniform over different attack locations. This means our scheme provides almost same protection against flooding attacks regardless of the attackers' location. In contrast, TVA's performance is highly dependent upon attackers' location since TVA assigns more buffer space to nearby domains (viz., Section II).

### B. Differential Guarantees

Path-identifier aggregation, which optimizes domain bandwidth allocation when attack sources are widely dispersed across domains, occurs whenever the number of active paths  $(|\mathcal{S}|)$  becomes greater than the number of access-guaranteed paths  $(|\mathcal{S}|_{max})$ . In Fig. 6, the result of the queue implementation for  $\mathcal{G}(34,8,S_i)$  illustrates the effectiveness of aggregation. As aggregation increases bandwidth allocation to legitimate paths by a factor of  $\frac{|\mathcal{S}|-|\mathcal{S}|_{max}}{|\mathcal{S}|_{max}}$  (i.e.,  $6/34\approx 17.6\%$  in that simulation), the request drop ratio of those paths decreases 76.8% (from 6.43% to 1.49%) when compared with that of the queue implementation for  $\mathcal{G}(64,4,S_i)$  (under which no path aggregation occurs). This is  $far\ below$  the stationary drop probability of legitimate paths (i.e.,  $1-P(|\mathcal{S}|,8,S_i)\approx 5.32\%$ ) which would result when physically separate queues are assigned to those paths.

We also evaluate the effectiveness of the protocol conformance measure in aggregating attack paths. For this, we

<sup>8</sup>For k guaranteed buffer-slots, the stationary packet service probability of  $S_i$  is determined by  $P(|\mathcal{S}|,k,S_i)=1-\frac{\rho_{S_i}^k(1-\rho_{S_i})}{1-\rho_{S_i}^{k+1}}$ . This is derived from the blocking probability of a M/M/1/k queueing system.

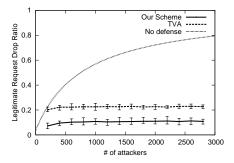


Fig. 4: Request drop ratio of legitimate paths. Error bars represent 95% confidence intervals.

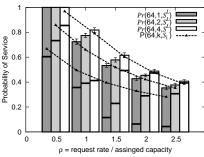


Fig. 5: Request service probability of legitimate paths with respect to bandwidth utilization  $(\rho)$ . The solid horizontal lines inside bars represent the probabilistic guarantees  $(\mathcal{G}(|\mathcal{S}|, k, S_i))$ .

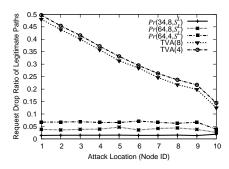


Fig. 6: Request drop ratio of legitimate paths with respect to attack location in the unbalanced tree. TVA(k) represents the result of TVA with queue-length k.

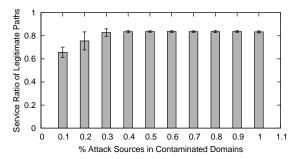


Fig. 7: Aggregation by protocol conformance: The request service ratio of legitimate paths increases as the fraction of bots becomes higher.

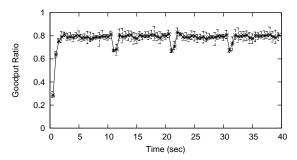


Fig. 8: Time variation of goodput ratio at the congested link. Legend: Error bars represent the minimum and maximum of goodput ratio.

configure a 64-path balanced tree such that the same number of nodes are attached to leaf nodes to make the request rates of all paths identical. Then, we set  $|\mathcal{S}|_{max}$  to 34 (which limits the number of attack path-identifiers by at most two) and increase the fraction of attack sources whose capability requests are denied at the destination host, from 10 to 100% in half of the leaf nodes. Note that the bandwidth conformance measure alone cannot distinguish attack paths from legitimate ones when the same request rates occur in all paths.

As Fig. 7 shows, aggregation is more precisely performed on attack paths (which leads to higher service ratios of legitimate paths) as the fraction of attack sources in contaminated domains grows. When domains are lightly contaminated (i.e., the fraction of attack sources is less than 40% in this simulation), legitimate paths can be aggregated. This is because aggregating attack paths near the attack target (i.e., multi-level aggregation of those attack paths) produces a higher aggregation cost than aggregating legitimate paths near their origins. Relatively high cost of multi-level aggregation also causes high service-ratio variation to legitimate paths, as a result of imprecise distinction between legitimate and attack paths.

# C. Rolling Attacks

Another simulation we performed is that of the "rolling attacks", whereby attack sources change their location to

exploit delays in the response time of any defense mechanism. For this simulation, we attach 16 attack nodes at 4 different locations in the unbalanced tree (i.e., at node 1,2,9 and 10) of Fig. 3 and place 200 attack sources in each attack node. We configure a rolling attack such that attack sources attached to node 1 and 10 flood the target for 10 seconds and the other attack sources for the next 10 seconds with a 20-second period.

In Fig. 8, we illustrate the time variation of goodput ratio (viz., Section VI) at the congested link averaged over 10 runs. The goodput ratio is very low at the beginning of the simulation, since attack requests go through the target link before being preempted by legitimate ones. However, as buffer-preemption occurs (as soon as the buffer is filled) and aggregation starts (around t=2), the goodput ratio rises sharply. Changing attack location significantly decreases the goodput ratio as the number of attack path-identifiers at the congested router increases four times (i.e., from 2 aggregated path-identifiers to 8 path-identifiers). However, these effects disappear whenever a new aggregation decision is made on the switched attack paths in  $\Delta_{agg}$  (which is set to  $20 \cdot \text{RTT} \approx 2 \text{ seconds}$  in this simulation).

# VIII. INTERNET-SCALE SIMULATIONS

In this section, we present large-scale simulation results to evaluate and compare the effectiveness of different defense

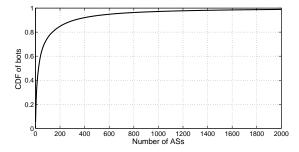


Fig. 9: Bot distribution vs. the number of ASs.

mechanisms (i.e., DefAT, Portcullis and TVA) against DoC attacks. For this purpose, we construct network topologies using real packet-routes and bot distribution in the Internet. Then, we compare the link-access times of legitimate capability-requests under different defense mechanisms.

### A. Datasets

For Internet-scale simulations, we use two real datasets: CAIDA Skitter-Map [24] and Composite Blocking List (CBL) [11]. A Skitter-Map contains the full routing-paths measured from a root-DNS server to a large set of randomly-chosen hosts ( $300\sim400$  thousand) in the Internet. A Skitter-Map is used as a reference topology from which simulation topologies are generated for given attack sizes. Among several distinct maps (that are constructed at different locations), we use widely different topologies for simulations, in order to observe the dependence of defense mechanisms on network topologies (more specifically the locales of legitimate and attack sources) .

A CBL contains a list of the IP addresses of active spambots. We first cluster the IP addresses in the CBL by their AS using GeoLite ASN [25], and obtain a reference distribution of bots (clustered by AS) as illustrated in Fig. 9. The figure shows that 300 ASs are responsible for about 90 % of bots and 600 ASs are for over 95 %. When the number of active ASs is considered (which is over 35,000), only a small fraction of ASs host most bots. This evidences the highly non-uniform distribution of bots in the Internet.

Based on this bot-distribution, if the size of simulated attacks is determined, we harvest attack sources from the same subnet as the one appears in the reference topology (i.e., Skitter-Map) and place these attack sources in the topology such that they have the same bot-distribution as the reference distribution. Then, we randomly choose legitimate sources from the Skitter-Map and add them to the simulation topology. Thus, the distribution of legitimate sources would be similar with that of AS sizes (in terms of the allocated IP address space).

# B. Scenarios

Our simulator runs in a discrete-time fashion, where individual packets advance a single router-hop in a time tick. If we assume 5ms delay for a single link (i.e., a clock-tick is 5ms),

the end-to-end delay for a source located 30-hops from the destination would be 150ms. Routers keep the packets arrived during a tick and handle the packets according to its admission policy. In simulations, the bottleneck-link capacity is set to 1000 requests per tick, which corresponds to 2.8Gbps (i.e., slightly higher capacity than OC-48) if 5ms clock-tick and 5 % bandwidth reservation for the capability-request channel are assumed. In all simulations, we configure attack sources to send 10 times more capability-request packets than legitimate sources, hence the relative strength of attack at the target-link would be 10 times the ratio of attack sources to legitimate ones. Attack sources start sending packets from the beginning of simulation and keep flooding the target during the entire simulation interval. Meanwhile, legitimate sources start their transmission after the target-link is fully congested to avoid the case that packets from closely located sources to the target get through the link (i.e., are serviced) before congestion occurs and they finish their transmission. Since packet arrivals from legitimate sources are delayed proportional to their distance to the target-link, these packet arrivals would have the same distribution as that of path length (viz., Section VIII-C). Before presenting simulation results, we first briefly describe individual defense mechanisms (i.e., Portcullis, TVA, and DefAT) used for comparison.

- Portcullis: A Portcullis client, once identifying its capability request being rejected (due to link congestion), starts solving a computational puzzle that requires to spend a certain amount of time (i.e., proves its computational effort) and increases the puzzle level until it receives a valid capability. In order to solve a higher level puzzle, the client needs to spend twice the time spent for the current level puzzle. Portcullis routers prioritize the packets that carry higher-level puzzle solutions. Hence, once a legitimate source solves a higher-level puzzle than attack sources, its request is guaranteed to be serviced at Portcullis routers. Portcullis provides the best persource link-access guarantee if attack sources cannot be distinguished from legitimate sources and are uniformly distributed over the Internet.
- TVA: TVA implements fair queueing on the incoming domains (i.e., ASs) from which packets arrive. The original scheme is improved later to handle remotely originated packets, which, whether being legitimate or not, become aggregated with others as they proceed to the target. For this purpose, TVA adopts a hierarchical fair queueing mechanism, where the TVA router allocates a fair amount of queue space to immediate upstream domains and these queues are split recursively for their upstream domains to provide fairness. Hence, the queue size for a source domain is determined by its distance (in terms of AS hops) to the target domain and the number of paths that are aggregated on its path to the destination. In the simulator, fair queueing is implemented for all outstanding requests arrived during a time tick via keeping all requests during the interval and randomly choosing

excess requests that need to be dropped. We assume that all capability-requests carry valid path-identifiers though path-identifier authenticity is not considered in TVA. Hence, our implementation approximates the best fairness that TVA can achieve. In comparative simulations, we use this advanced version of TVA.

• DefAT: A DefAT router provides link-access guarantees to source ASs as explained throughout this paper. However, for fair comparative simulations, path-aggregation is disabled at DefAT routers because it can significantly favors the results of DefAT depending on how we set the number of access-guaranteed paths. Though the number of access-guaranteed paths can be optimized to maximize goodput, we leave it as a configurable parameter as discussed before. We assume that path-identifiers cannot be spoofed or replayed using the security protection mechanisms provided in Section IV.

# C. Topology

We choose three Skitter-Maps constructed at different locations (i.e., f-root, h-root, and apan-jp) and generate simulation topologies with two parameters: the number of legitimate sources and the number of attack sources. In topology generation, we set the number of legitimate sources to 10K and change the attack size from 50K to 300K. Fig. 10 shows the topology statistics (i.e., AS-path length from source to destination and average degree of ASs located at the same distance from the target) for 300K attack sources. The length of AS-path, which is the number of ASs on a path including the source AS, spans from 1 to 10, yet is mostly concentrated between 3 and 5 in all three topologies. The average ASdegree (i.e., the number of immediate upstream ASs that send/forward traffic) is widely different in the topologies, hence it would better characterize topologies. Note that the number of paths (left vertical-axis) and AS-degree (right vertical-axis) are shown in a normal scale and a log scale respectively. For example, f-root is constructed at an AS that has a very low degree yet whose provider-ASs have a very high degree. On the other hand, h-root is constructed at a highdegree AS whose 1-hop and 2-hop neighboring ASs have high AS-degrees as well. Finally, apan-jp topology has mid-degrees both at the target (where the topology is constructed) and its provider. Simulations using these topologies would produce different results when a router's defense scheme prioritize traffic based on its confidence on traffic source (e.g., more buffer allocation to closer domains in TVA). We note that for different attack strengths, topologies do not change significantly mainly because attack sources are highly clustered by their locale. We summarize the more statistical data on the above topologies in Table I.

Legend:  $l_{avg}$ : average path length,  $l_{max}$ : longest path length,  $l_{avg}^{AS}$ : average AS-path length,  $l_{max}^{AS}$ : longest AS-path length,  $N_r$ : number of routers.

TABLE I: Topology statistics.

	$l_{avg}$	$l_{max}$	$l_{avg}^{AS}$	$l_{max}^{AS}$	$N_r$
f-root	14.99	29	5.44	10	48,624
h-root	13.55	31	4.84	10	42,679
apan-jp	17.15	33	4.80	9	36,621

# D. Comparative Simulations

We compare the link-access times of legitimate capabilityrequests with different defense mechanisms employed at the target link. In f-root topology, DefAT provides earlier link access to over 90 % of legitimate requests than other mechanisms in the presence of 100K attack sources, and 80 % of those requests are almost unaffected by the attack when compared with the reference access time curve; i.e., that of no attack (viz., Fig. 11(a)). With Portcullis, all legitimate requests get a link access when legitimate sources start solving a higher level puzzle than attack sources. The figure shows that about a half of requests are serviced at around 150 tick, yet the remaining requests are serviced at 300 tick as they had to spend twice the time (i.e., 300 ticks) to solve the next level puzzle. As a consequence, the link access time curve of Portcullis has two sharp increases like a step function; i.e., the link access-times show bimodal distribution. TVA favors only a small fraction of legitimate requests (about 45 % of legitimate requests in froot topology) because requests from close ASs (to the target) have higher buffer allocation than those of remote ASs.

In h-root topology, slightly faster link-access times are observed with DefAT. This is because the path-diversity in this topology is higher than f-root as can be seen in Fig. 10. Higher path-diversity of legitimate requests enables those requests to get more buffer allocation in a DefAT router as DefAT provides guarantees to individual source ASs (this would further reduce the buffer allocation to highly clustered attack requests). Portcullis provides almost identical linkaccess times for legitimate requests despite topology changes since its request admission is primarily determined by clients' computational effort (the level of puzzle that clients solved) rather than simulation topologies. With TVA, a significantly different result is observed: only 22 % of legitimate requests have an earlier link access than no defense, yet service to remaining 78 % of them are more delayed than in f-root topology. This explains that the effectiveness of TVA is highly dependent on the network topology. These results are consistent under different attack sizes (i.e., topologies with 100K, 200K, and 300K attack sources) as Fig. 11, 12 and 13

Now, we observe how the link-access times of legitimate requests are affected by the attack size by increasing the attack size from 50K to 300K. With DefAT, about 80 % of legitimate requests are unaffected regardless of the attack size, yet the remaining 20 % of requests (which originate from attack domains) take longer time to get through the congested link as the attack size grows. This is an expected result since we do not attempt to distinguish legitimate requests from

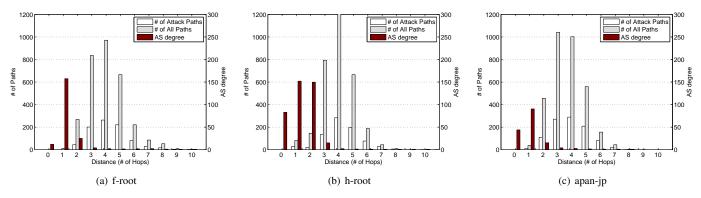


Fig. 10: Simulation topology with 300K attack sources.

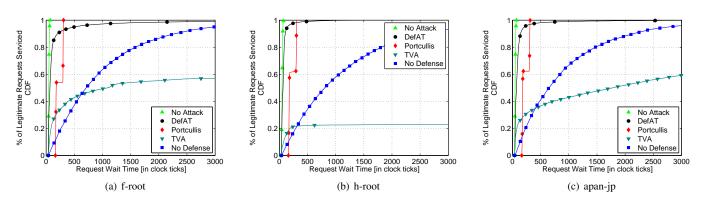


Fig. 11: Link-access time for legitimate capability-requests under 100K attack sources.

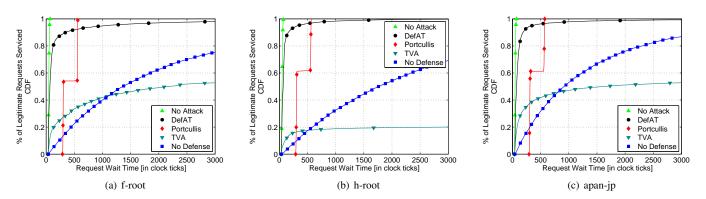


Fig. 12: Link-access time for legitimate capability-requests under 200K attack sources.

attack requests if they originate from the same domain. In all three topologies, the link-access times show a consistent result: the link-access times of 20 % of legitimate requests become longer as the attack size grows though they are slightly differ in different topologies. With Portcullis, the link-access times are doubled if the attack size reaches at a certain threshold. However, the threshold is not proportional to the number of attack sources because attack sources, in order to congest a link with a higher level of puzzles (attack sources solve a computational puzzle as well), should double their size. Otherwise, they cannot fully congest the link. Thus, with Portcullis, link-access times are highly dependent on the

attack size even though they are not on network topologies. In contrast, TVA's performance is highly dependent on network topologies as illustrated in Fig. 14(c), 15(c) and 16(c). This phenomenon can be explained as: despite the high AS-degree of the congested domain (viz., Fig. 10(b)), most of legitimate requests are aggregated with attack requests if they originate remotely (note that in a limited size buffer, queue cannot split indefinitely). TVA works well only if the high-degree AS is directly connected with or closely located at the target ASs. This is why TVA shows relatively better performance in apanjp topology. However, TVA's advantage to legitimate sources is marginal as TVA allocates more buffer-space to closely

located ASs regardless whether they originate legitimate or attack requests. In Fig. 15(c) and 16(c), several leaps in the CDF (e.g., between 1000 and 1500 ticks in Fig. 15(c)) indicate more queue-space becomes available to legitimate requests in a short time interval. This happens when some legitimate paths disappear after finishing transmission and eventually enable a TVA queue to be split into multiple separate queues for other remaining paths. Note that a TVA queue splits recursively (towards the source ASs in the traffic tree) unless the number of distinct incoming-paths (i.e., immediate children) to the queue (i.e., intermediate node in the traffic tree) exceeds the available queue size.

# IX. CONCLUSIONS

In this paper, we present a defense scheme against link flooding attacks targeting connection setups in capability systems. Our design of a new authenticated path-identification mechanism provides individual packets with unforgeable domain identifiers to which link-access guarantees are provided at remote routers. Guarantees of link access, defined as the probabilistic lower bounds of link access, are provided in a domain basis and they are provided differentially based on domain contaminations. We show the effectiveness of our design in two ways. First, NS2 simulations support our analytical results: (1) link-access guarantees that are independent of global attack sources and their location, and (2) resilience against attack dispersion via differential guarantees. Second, Internet-scale simulations, using the real network topologies and bot distributions, provide strong evidences on the nonuniform distribution of bots and how DefAT localizes their effects on legitimate capability-requests. More specifically, the simulation results show that over 80 % of legitimate requests are unaffected or minimally affected by large-scale attacks, which could not be achieve with previous per-source or per-aggregate defense mechanisms. We note that differential link-access guarantees would provide positive incentives to administrative domains that employ strong security measures against malware contamination.

# ACKNOWLEDGMENT

This research was supported in part by US Army Research Laboratory and the UK Ministry of Defence under Agreement Number W911NF-06-3-0001 and by the US Army Research Office under Contract W911NF-07-1-0287 at the University of Maryland. The work on Internet simulations was supported by Northrop Grumman Cyber Research Consortium under Contract NGIT2009100109 at CyLab, Carnegie Mellon University. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the US Army Research Laboratory, US Army Research Office, the U.S. Government, the UK Ministry of Defense, the UK Government, or Northrop Grumman.

### REFERENCES

[1] D. X. Song and A. Perrig, "Advanced and Authenticated Marking Schemes for IP Traceback," in *INFOCOM*, 2001.

- [2] S. Savage, D. Wetherall, A. R. Karlin, and T. Anderson, "Practical network support for IP traceback," in SIGCOMM, 2000.
- [3] P. Ferguson, "Network Ingress Filtering:Defeating Denial of Service Attacks which employ IP Source Address Spoofing," RFC 2827, 2000.
- [4] T. Anderson, T. Roscoe, and D. Wetherall, "Preventing Internet denial-of-service with capabilities," in *Hotnets-II*, 2003.
- [5] A. Yaar, A. Perrig, and D. Song, "SIFF: A Stateless Internet Flow Filter to Mitigate DDoS Flooding Attacks," in *Proceedings of the IEEE Security and Privacy Symposium*, 2004.
- [6] X. Yang, D. Wetherall, and T. Anderson, "A DoS-limiting network architecture," in *IEEE/ACM TRANSACTIONS ON NETWORKING*, 2008.
- [7] K. Argyraki and D. R. Cheriton, "Network Capabilities: The Good, the Bad and the Ugly," *HotNets IV*, 2005.
- [8] B. Parno, D. Wendlandt, E. Shi, A. Perrig, B. Maggs, and Y.-C. Hu, "Portcullis: Protecting Connection Setup from Denial-of-Capability Attacks," in SIGCOMM, 2007.
- [9] S. Staniford, V. Paxson, and N. Weaver, "How to Own the Internet in Your Spare Time," in *Proceedings of the 11th USENIX Security* Symposium, 2002.
- [10] D. Dagon, C. Zou, and W. Lee, "Modeling Botnet Propagation Using Time Zone," Network and Distributed System Security Symposium, 2006.
- [11] "http://cbl.abuseat.org/."
- [12] "http://www.computerworld.com/s/article/9076278/."
- [13] L. von Ahn, M. Blum, N. Hopper, and J. Langford, "CAPTCHA: Using hard AI problems for security," in *Proceedings of Eurocrypt*, 2003.
- [14] A. Yaar, A. Perrig, and D. Song, "Pi: A Path Identification Mechanism to Defend against DDoS Attacks," in *In IEEE Symposium on Security and Privacy*, 2003.
- [15] P. E. McKenney, "Stochastic fairness queueing," in INFOCOM, 1990.
- [16] M. Shreedhar and G. Varghese, "Efficient fair queueing using deficit round robin," in SIGCOMM '95, 1995, pp. 231–242.
- [17] "http://www.icann.org."
- [18] "http://www.potaroo.net/tools/asn32/."
- [19] L. Fan, P. Cao, J. Almeida, and A. Z. Broder, "Summary cache: A scalable wide-area web cache sharing protocol," in *IEEE/ACM Transactions* on *Networking*, 1998.
- [20] W. C. Feng, D. D. Kandlur, D. Saha, and K. G. Shin, "Stochastic Fair Blue: A Queue Management Algorithm for Enforcing Fairness," in *INFOCOM*, 2001, pp. 1520–1529.
- [21] R. Mahajan, S. Floyd, and D. Wetherall, "Controlling high-bandwidth flows at the congested router," in ICNP '01, 2001.
- [22] S. B. Lee and V. D. Gligor, "Floc: Dependable link access for legitimate traffic in flooding attacks," in *International Conference on Distributed Computing Systems*, 2010, pp. 327–338.
- [23] A. Studer and A. Perrig, "The coremelt attack," in ESORICS, Saint Malo, France, September 2009.
- [24] "http://www.caida.org/."
- [25] "http://www.maxmind.com/app/asnum."

# APPENDIX

# A. Proof of Probabilistic Guarantees

If  $\rho_{S_i} \leq 1$ , a packet carrying  $S_i$  is guaranteed to be serviced if less than k arrivals of  $S_i$  have occurred in  $\Delta_Q$  before its arrival. Let  $N_{S_i}(\Delta_Q)$  be the # of requests in  $\Delta_Q$ . Then, the probability of service guarantee on  $S_i$  is given as follows.

$$\Pr(N_{S_i}(\Delta_Q) < k) = \sum_{j=0}^{k-1} \frac{(\lambda_{S_i} \Delta_Q)^j}{j!} e^{-\lambda_{S_i} \Delta_Q}$$

$$\geq \sum_{j=0}^{k-1} \frac{(k \cdot \rho_{S_i})^j}{j!} e^{-k \cdot \rho_{S_i}} \quad (A-1)$$

In contrast, for  $\rho_{S_i} > 1$ , a per-packet guarantee cannot be provided, since at least  $\frac{\rho_{S_i}-1}{\rho_{S_i}}$  of requests must be dropped regardless of the buffer size. In this case, only a fraction of its requests can be guaranteed to be serviced (i.e.,  $\frac{1}{\rho_{S_i}}$ ), hence the probabilistic lower bound of link access is defined

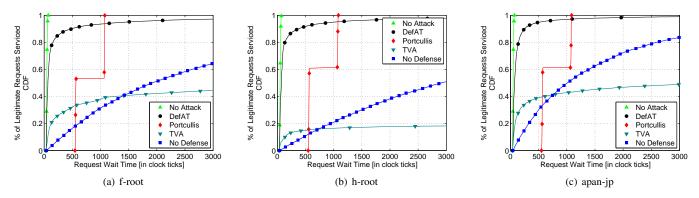


Fig. 13: Link-access time for legitimate capability-requests under 300K attack sources.

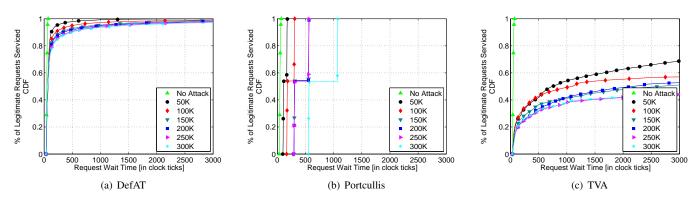


Fig. 14: Link-access time for legitimate capability-requests under the f-root topology and different botnet sizes.

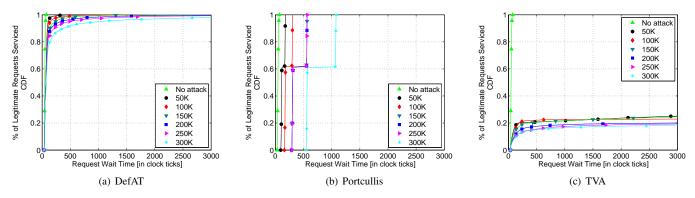


Fig. 15: Link-access time for legitimate capability-requests under the h-root topology and different botnet sizes.

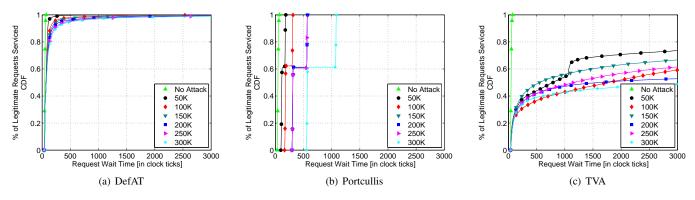


Fig. 16: Link-access time for legitimate capability-requests under the apan-jp topology and different botnet sizes.

as the product of  $\frac{1}{\rho_{S_i}}$  and the probability that the allocated bandwidth is fully utilized. Let  $P_f(S_i)$  be the probability that a packet arrival of  $S_i$  finds k buffered  $S_i$ s. The probability of full bandwidth utilization is greater than  $P_f(S_i)$ . Let  $\mathcal{G}_{\mathcal{L}} = \sum_{j=0}^{k-1} \frac{(k \cdot \rho_{S_i})^j}{j!} e^{-k \cdot \rho_{S_i}}$ . Then,

$$P_{f}(S_{i}) = 1 - \Pr(\# \text{ of } S_{i}\text{'s in the buffer } < k)$$

$$\geq 1 - \left(\mathcal{G}_{\mathcal{L}} + \sum_{j=k}^{\infty} \frac{(k \cdot \rho_{S_{i}})^{j}}{j!} e^{-k \cdot \rho_{S_{i}}} \binom{j}{j-k+1} (1 - \mathcal{G}_{\mathcal{L}})^{j-k+1} \mathcal{G}_{\mathcal{L}}^{k-1}\right)$$

$$\geq 1 - \left(\mathcal{G}_{\mathcal{L}} + \sum_{j=k}^{\infty} \frac{(k \cdot \rho_{S_{i}})^{j}}{j!} e^{-k \cdot \rho_{S_{i}}} \binom{j}{j-k+1} (1 - \mathcal{G}_{\mathcal{L}}) \mathcal{G}_{\mathcal{L}}^{k-1}\right)$$

$$= (1 - \mathcal{G}_{\mathcal{L}}) \left(1 - \sum_{j=k}^{\infty} \frac{(k \cdot \rho_{S_{i}})^{j}}{j!} e^{-k \cdot \rho_{S_{i}}} \binom{j}{k-1} \mathcal{G}_{\mathcal{L}}^{k-1}\right).$$
 (A-2)

By Eqs. (A-1) and (A-2), the Eq. (V.1) follows.

# B. Proof of Error Bound

We first define two types of aggregating node. In  $\mathcal{T}'_{R_0}$ , the node whose all children nodes are leaf nodes is defined as the "leaf aggregator" and the any other non-leaf node is defined as "intermediate aggregator." The last added node to the solution set can be either a leaf aggregator or an intermediate aggregator.

If the last added node  $R_i$  to the optimal set  $(\mathcal{O})$  is a leaf aggregator, the error from the optimal solution is bounded by  $\sum_{R_j \in \mathcal{R}_i} \mathcal{E}_{R_j} \leq |\mathcal{R}_i| \cdot \mathcal{E}_{th}$ , where  $|\mathcal{R}_i|$  is the number of incoming links of  $R_i$ .

If the last added node to  $\mathcal{O}$  is an intermediate aggregator, we can consider two different cases. Let  $R_i$  be an intermediate aggregator, and  $R_{i1}, \ldots, R_{in}$  be the one-hop children of  $R_i$ . By the definition of aggregation cost, the following inequality can be shown.

$$C^{\mathcal{A}}(R_i) = \frac{|\mathcal{R}_i| - 1}{|\mathcal{R}_i|} \sum_{R_j \in \mathcal{R}_i} \mathcal{E}_{R_j} \ge \sum_{j=1}^n C^{\mathcal{A}}(R_{ij}) \qquad (B-1)$$

The above inequality means that the last node added to  $\mathcal{O}$  is either (a) the node whose all immediate children aggregators are already aggregated, or (b) the node whose aggregation cost is less than the total aggregation cost of the current solution set.

case (a):

$$C^{\mathcal{A}}(R_i) - \sum_{j=1}^{n} C^{\mathcal{A}}(R_{ij})$$

$$\leq \frac{1}{|\mathcal{R}_{i1}|} \sum_{R_j \in \mathcal{R}_{i1}} \mathcal{E}_{R_j} + \dots + \frac{1}{|\mathcal{R}_{in}|} \sum_{R_j \in \mathcal{R}_{in}} \mathcal{E}_{R_j}$$

$$\leq n \cdot \mathcal{E}_{th}$$

Like the leaf aggregator, if aggregation is performed at an intermediate aggregator  $R_i$ , the sum of aggregation costs of  $R_i$ 's children are deducted from the total cost. Therefore,

the maximum increase of aggregation cost at an intermediate aggregator is bounded by n.

case (b):

By (B-1), aggregation can occur at a node if either all its children nodes are in the solution set, or  $C^{\mathcal{A}}(R_i) < \sum_{R_{oi} \in \mathcal{O}} C^{\mathcal{A}}(R_{oi})$ , where  $\mathcal{O} = \{R_{o1}, R_{o2}, \dots R_{on}\}$  is the optimal solution set and  $R_{on}$  is the last node added to the current solution set.

$$C^{\mathcal{A}}(R_i) < \sum_{R_{oi} \in \mathcal{O}} C^{\mathcal{A}}(R_{oi})$$

$$\Leftrightarrow C^{\mathcal{A}}(R_i) < \sum_{j=1}^{n-1} C^{\mathcal{A}}(R_{oj}) + C^{\mathcal{A}}(R_{on})$$

$$\Leftrightarrow C^{\mathcal{A}}(R_i) - \sum_{j=1}^{n-1} C^{\mathcal{A}}(R_{oj}) < C^{\mathcal{A}}(R_{on})$$

Hence, the increase of aggregation cost cannot be greater than the product of  $\mathcal{E}_{th}$  and the incoming-link degree of the last added node to the solution set.