

# Routing Bottlenecks in the Internet: Causes, Exploits, and Countermeasures

Min Suk Kang  
ECE and CyLab  
Carnegie Mellon University  
Pittsburgh, PA USA  
minsukkang@cmu.edu

Virgil D. Gligor  
ECE and CyLab  
Carnegie Mellon University  
Pittsburgh, PA USA  
gligor@cmu.edu

## ABSTRACT

How pervasive is the vulnerability to link-flooding attacks that degrade connectivity of thousands of Internet hosts? Are some geographic regions more vulnerable than others? Do practical countermeasures exist? To answer these questions, we introduce the notion of the *routing bottlenecks* and show that it is a fundamental property of Internet design; i.e., it is a consequence of route-cost minimizations. We illustrate the pervasiveness of routing bottlenecks in an experiment comprising 15 countries and 15 cities distributed around the world, and measure their susceptibility to scalable link-flooding attacks. We present the key characteristics of routing bottlenecks, including size, link type, and distance from host destinations, and suggest specific structural and operational countermeasures to link-flooding attacks. These countermeasures can be deployed by network operators without needing major Internet redesign.

## Categories and Subject Descriptors

C.2.0 [Computer Communication Networks]: General—*Security and protection*; C.2.1 [Computer Communication Networks]: Network Architecture and Design—*Network topology*

## Keywords

DDoS attack; link-flooding attack; routing bottleneck; power law

## 1. INTRODUCTION

Recent experiments [26] and real-life attacks [6] have offered concrete evidence that link-flooding attacks can severely degrade, and even cut off, connectivity of large sets of adversary selected hosts in the Internet for uncomfortably long periods of time; e.g., hours. However, neither the root cause nor pervasiveness of this vulnerability has been analyzed to date. Furthermore, it is unknown whether certain network structures and geographic regions are more vulnerable

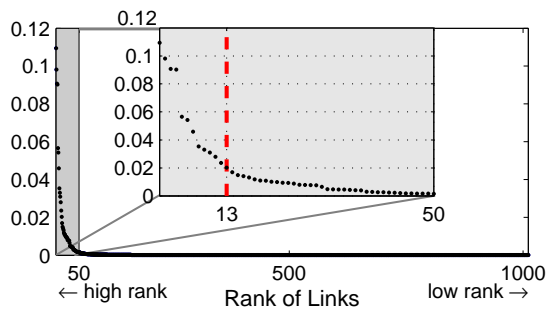
to these attacks than others. In this paper we address this gap in our knowledge about these attacks by (1) introducing the notion of the *routing bottlenecks* and its role in enabling link-flooding attacks at scale; (2) finding bottlenecks in 15 countries and 15 cities distributed around the world to illustrate their pervasiveness; and (3) measuring bottleneck parameters (e.g., size, link types, and distance to adversary-selected hosts) to understand the magnitude of attack vulnerability. We also present both structural and operational countermeasures that mitigate link-flooding attacks.

In principle, route diversity could enhance Internet resilience to link-flooding attacks against *large sets* of hosts (e.g., 1,000 hosts) since it could force an adversary to scale attack traffic to unattainable levels to flood all possible routes. In practice, however, the mere existence of many routes between traffic sources and selected sets of destination hosts cannot guarantee resilience whenever the vast majority of these routes are distributed across very few links, which could effectively become a routing bottleneck.

To define routing bottlenecks more precisely, let  $S$  denote a set of (source) IP addresses of hosts that originate traffic to a set of IP destination addresses, denoted by  $D$ .  $S$  represents any set of hosts distributed across the Internet. In contrast,  $D$  represents a set of hosts of a specified Internet region (e.g., a country or a city), which are chosen at random and independently of  $S$ . A *routing bottleneck* on the routes from  $S$  to  $D$  is a small set  $B$  of IP (layer-3) links such that  $B$ 's links are found in a majority of routes whereas the remaining links are found in very few routes.  $|B|$  is often over an order of magnitude smaller than both  $|S|$  and  $|D|$ . If all links are ranked by their frequency of occurrence in the routes between  $S$  and  $D$ , the bottleneck links,  $B$ , have a very high rank whereas the vast majority of the remaining links have very low rank. The sharper the skew in the frequency of link occurrence in these routes, the narrower the bottleneck. Routes to hosts  $D$  may have more than one bottleneck of size  $|B|$ .

**An Example.** To illustrate a real routing bottleneck, we represent route sources  $S$  by 250 PlanetLab nodes [42] distributed across 164 cities in 39 countries. For the route destinations,  $D$ , we select 1,000 web servers at random from a list of publicly-accessible servers obtained using the ‘computer search engine’ called Shodan (<http://www.shodanhq.com>) in *Country15* of the fifteen-country list  $\{Country1, \dots, Country15\}$ . This list is a permutation of the alphabetically ordered list of countries  $\{\text{Brazil, Egypt, France, Germany, India, Iran, Israel, Italy, Japan, Romania, Russia, South}$

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
CCS'14, November 3–7, 2014, Scottsdale, Arizona, USA.  
Copyright 2014 ACM 978-1-4503-2957-6/14/11 ...\$15.00.  
<http://dx.doi.org/10.1145/2660267.2660299>.



**Figure 1: Normalized link-occurrence distribution in routes from  $S = 250$  PlanetLab nodes to  $D = 1,000$  randomly selected servers in *Country15*.**

Korea, Taiwan, Turkey, and United Kingdom}.<sup>1</sup> We trace the routes between  $S$  and  $D$  for *Country15*, collect over  $1.9 \times 10^6$  link samples from those routes, and plot their link-occurrence distribution, as shown in Fig. 1. This figure clearly shows a very skewed link-occurrence distribution, which implies the existence of a narrow routing bottleneck; i.e.,  $|B| = 13$  links are found in over 72% of the routes; viz., Fig. 11a.

In this paper we argue that the pervasive occurrence of routing bottlenecks is a fundamental property of the Internet design. That is, power-law distributions that characterize the frequency of link occurrence in routes are a consequence of employing route-cost minimization, which is a *very desirable* feature of Internet routing; viz., Section 2. Fortunately, routing bottlenecks do not lead to traffic degradation during ordinary Internet use, because the bandwidth of bottleneck links is usually provisioned adequately for normal mode of operation. Hence, these bottlenecks should not be confused with the bandwidth bottlenecks in end-to-end paths [1, 23] since one does not always imply the other; viz., Section 6.3.

**Problem.** Unfortunately, however, bottleneck links provide a very attractive target to an adversary whose goal is to flood few links and severely degrade or cut off connectivity of targeted servers,  $D$ , in various cities or countries around the world. For example, an adversary could easily launch a traffic amplification attack using NTP monlists (400 Gbps) [38] and DNS recursors (120 Gbps) [6] to distribute an aggregate of 520 Gbps traffic across the 13-link bottleneck of Fig. 1. Such an attack would easily flood these links, even if each of them is provisioned with a maximum of 40 Gbps capacity, severely degrading the connectivity of the 1,000 servers of *Country15* from the Internet; viz., Fig. 11a. More insidious attacks, such as Crossfire [26], can flood bottleneck links persistently with attack traffic that is indistinguishable from legitimate traffic by routers and invisible to, and hence undetectable by, the targeted servers,  $D$ .

To counter link-flooding attacks that exploit routing bottlenecks, we first define the parameters that characterize these bottlenecks; e.g., size, link types, and average distance of bottleneck links from the targeted servers,  $D$ . Then we define a *connectivity-degradation metric* to provide a quantitative view of the risk exposure faced by these servers. The bottleneck parameters and metric are particularly impor-

<sup>1</sup>The permutation is a country ordering by link-occurrence skew. Finding it and de-anonymizing the country list would require repeating the measurements illustrated in Fig. 2a.

tant for applications in the targeted country or city where Internet-facing servers need stable connectivity; e.g., industrial control systems [9], financial [49], defense and other government services. For these applications, routing bottlenecks pose unexpected vulnerabilities, since diversity of IP-layer connections, which is often incorrectly believed to be sufficient for route diversity, only assures necessary conditions for route diversity but does not guarantee it. We illustrate the usefulness of our connectivity-degradation metric in assessing the vulnerabilities posed by real life routing bottlenecks found in fifteen countries and fifteen different cities around the world; viz., Section 2.

Analysis of routing bottleneck exploits explains why intuitive but naive countermeasures will not work in practice; e.g., reactive re-routing to disperse the traffic flooding bottleneck links across multiple local links; flow filtering at routers based on traffic intensity; reliance on backup links on exposed routes. More importantly, our analysis provides a precise *route-diversity metric*, which is based on autonomous-system (AS) path diversity, and illustrates the utility of this metric as a proxy for the bottleneck avoidance in the Internet. Finally, our analysis suggests operational countermeasures against link-flooding attacks, including inter- and intra-domain load balancing, and automatic intra-domain traffic engineering.

**Contributions.** In summary, we make the following contributions:

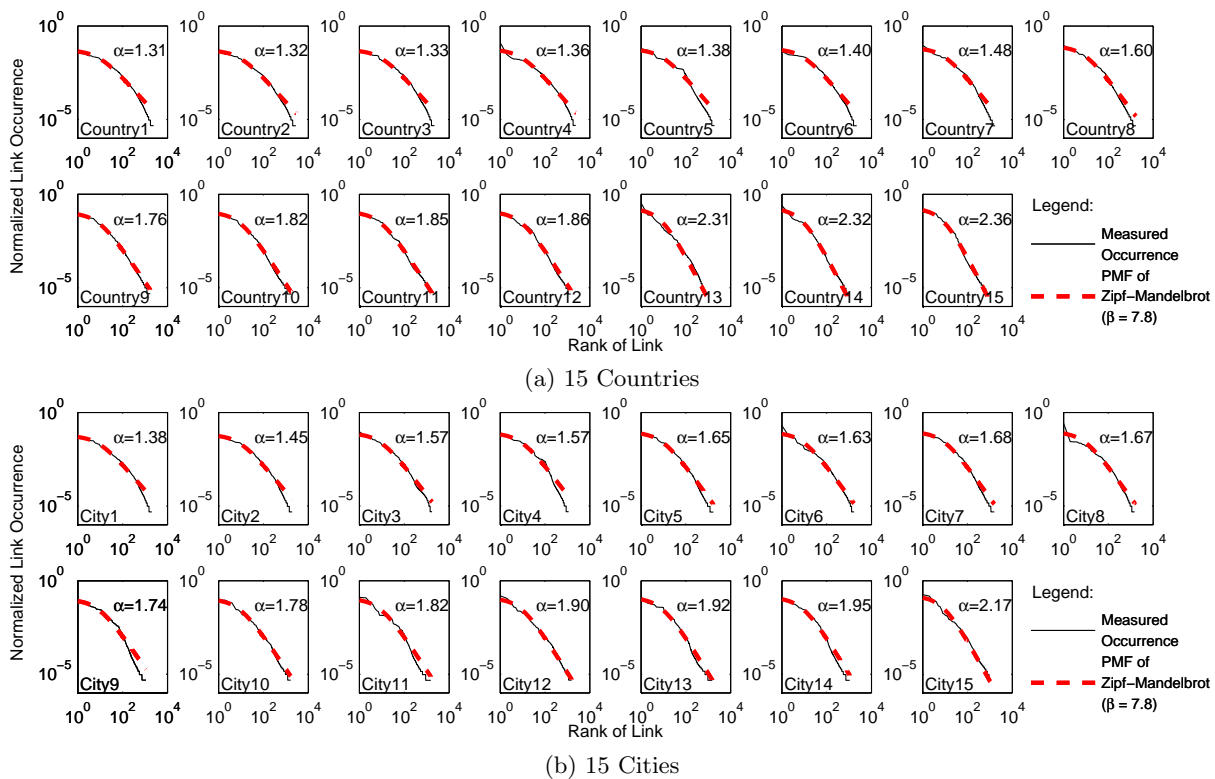
- We explain the root causes and characteristics of routing bottlenecks in the Internet, and illustrate their pervasiveness with examples found in 15 countries and 15 cities around the world.
- We present a precise quantitative measure of connectivity degradation to illustrate how routing bottlenecks enable an adversary to scale link-flooding attacks without much additional attack traffic.
- We present several classes of countermeasures against attacks that exploit routing bottlenecks, including both structural and operational countermeasures.

## 2. ROUTING BOTTLENECKS

### 2.1 Link-occurrence measurements

To determine the existence of routing bottlenecks, we measure the link-occurrence distribution in a large number of the routes towards a selected destination region. This requires that we perform *traceroutes* to obtain a series of *link samples* (i.e., IP addresses at either end of layer-3 links) on a particular route from a source host to a destination host in a selected Internet region. From the collected link samples on the routes, we construct the link-occurrence distribution by counting the number of samples for each link. Then we select the minimum set of links whose removal disconnects all routes to the destination region by removing redundant links.<sup>2</sup> Section 4.1 describes the selection algorithm used. In these measurements, we trace 250,000 routes by using *traceroute* from 250 source hosts (i.e., 250 PlanetLab nodes [42]) to 1,000 randomly selected web servers in each of 15 countries and 15 cities.

<sup>2</sup>For example, a link is redundant if the routes that traverse it, also traverse a link already selected for the minimum set.



**Figure 2: Normalized link occurrence/rank in traced routes to 1,000 randomly selected hosts in each of the 15 countries in (a) and 15 cities in (b).**

*Traceroute* is a common network monitoring tool whose use is often fraught with pitfalls [47]. Care was taken in analyzing the *traceroute* dataset so that our measurement results are not affected by the typical errors of *traceroute* use; e.g., alias resolution, load-balanced routes, accuracy of returned IP, hidden links in MPLS tunnels. For a detailed discussion, see Section 3.3. We perform multiple *traceroutes* for the same source-destination host pair to determine the *persistent* links; i.e., links that always show up in the multiple *traceroutes*. We collect only the samples of persistent links because non-persistent links do not lead to reliable exploitation of routing bottleneck. We have found extremely skewed link-occurrence distribution for the 1,000 randomly selected hosts in each of the 15 countries and 15 cities, which strongly indicates the existence of routing bottlenecks in all the countries and cities in which we performed our measurements.

## 2.2 Power-law in link occurrence distributions

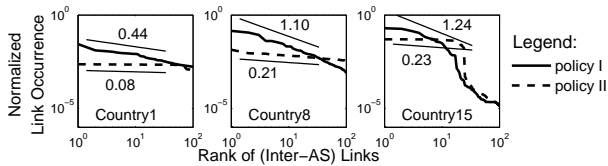
The analysis of link-occurrence distributions helps us understand both the cause of routing bottlenecks and their physical characteristics (e.g., size, type, distance from destination hosts) as well as countermeasures against flooding attacks that attempt to exploit them. To illustrate the skew of link-occurrence distributions, we present our measurements for 15 countries and cities around the world in Fig. 2a and Fig. 2b, respectively. The fifteen-city list  $\{City1, \dots, City15\}$  is a permutation of the following 15 cities, which are listed in alphabetical order: {Beijing, Berlin, Chicago, Guangzhou, Houston, London, Los Angeles, Moscow, New York, Paris, Philadelphia, Rome, Shanghai, Shenzhen, and Tianjin}. In

these figures, we illustrate the relation between the link occurrence normalized by the total number of measured routes and the rank of links on *log-log* scale, for 1,000 servers in each country and city. The normalized occurrence of a link is the portion of routes between  $S$  and  $D$  carried by the link; e.g., if a link carries 10% of routes between  $S$  and  $D$ , its normalized occurrence is 0.1.

We observe that the normalized link-occurrence distribution is accurately modeled by the *Zipf-Mandelbrot* distribution; namely,

$$f(k) \sim 1/(k + \beta)^\alpha,$$

where  $k$  is the rank of the link,  $\alpha$  is the exponent of the power-law distribution, and  $\beta$  is the fitting parameter. Exponent  $\alpha$  is a good measure of route concentration, or distribution skew, and hence of bottleneck size: the higher  $\alpha$ , the sharper the concentration of routes in a few links. Fitting parameter  $\beta$  captures the flatter portions of the distribution in the high-rank region; i.e., lower values on the x-axis. This region is not modeled as well by an ordinary Zipf distribution since its probability mass function would be a straight line in *log-log* scale on the entire range. The flatter portion of the distribution in high-rank region is due to the nature of link sampling via route measurement. That is, multiple links are sampled together when each route is measured and there exist no duplicate link samples in a route in general due to the loop-freeness property of Internet routes. Thus, the occurrences of extremely popular links are limited and the high-rank region is flattened. (Similarly flattened occurrence of high-ranked data samples was observed and explained in other measurements and modeling studies [19].)



**Figure 3: Normalized link occurrence/rank in simulated inter-AS links in three countries.**

To enable comparison of route concentration in a few links of different destination regions, we fix the fitting parameter  $\beta$  and find the values of exponent  $\alpha$  for the best fit across the fifteen countries; i.e.,  $\beta = 7.8$  causes the smallest fitting error. In Fig. 2a and Fig. 2b, the fifteen countries and cities are ordered by increasing value of  $\alpha$  in the range 1.31 – 2.36.

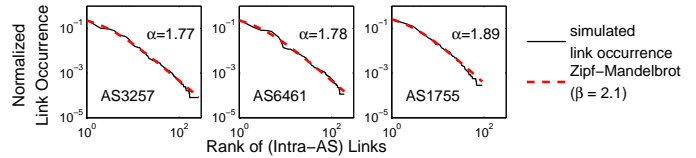
### 2.3 Causes

What causes routing bottlenecks, or high skew/power-law distribution of link occurrence? Often, power-law distributions (especially the Zipf-Mandelbrot distribution) arise from processes that involve some cost minimization. For example, research in linguistics shows that power laws defining the frequency of word occurrences in random English text arise from the minimization of human-communication cost [33, 55]. Thus, one would naturally expect that power-laws in link-occurrence distributions are caused by the *cost minimization* criteria for route selection and network design in the Internet; i.e., both intra- and inter-domain interconnections and routing. Extra cost minimization is provided by the “hot-potato” routing between domains. We note that a correlation between shortest-path routing and routing bottlenecks was already suggested in the Crossfire attack [26]. However, this suggestion does not show causality since shortest-path routing is not always minimum cost; e.g., BGP path selection criteria [16, 18].

#### 2.3.1 Cost minimization in inter-domain routing

**Inter-domain routing policy creates routing bottlenecks in inter-AS links:** BGP is the *de facto* routing protocol for inter-domain (i.e., AS-level) Internet connectivity. The *rule-of-thumb* BGP policy for choosing inter-AS paths is the minimization of the network’s operating cost. That is, whenever several AS paths to a destination are found, the minimum-cost path is selected; e.g., customer links are preferred over peer links and over provider links. If there exist multiple same-cost paths, the shortest path is selected. This policy is intended to minimize operating costs of routing in the Internet [16, 18].

To determine whether the rule-of-thumb routing policy (i.e., policy I) contributes to the creation of routing bottlenecks, we compare its effects with those of a *hypothetical* routing policy that distributes routes uniformly across possible inter-domain links (i.e., policy II). This hypothetical policy favors inter-domain links that serve fewer AS paths for a particular destination. To perform this comparison, we run AS-level simulations using the most recent (i.e., June 2014) CAIDA’s AS relationship, which is derived by Luckie *et al.* [30]. We simulate the hypothetical policy using Dijkstra’s shortest-path algorithm [10] with dynamically changing link weights, which are proportional to the number of BGP paths served.



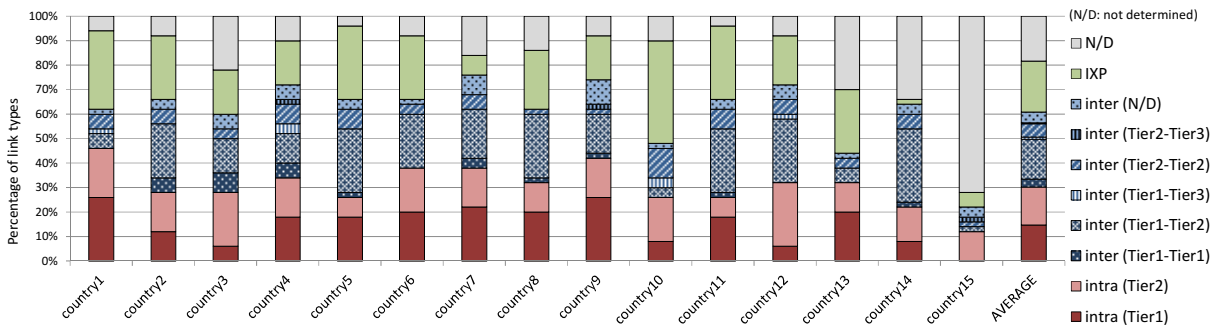
**Figure 4: Normalized link occurrence/rank in simulated AS-internal routes for three ISPs.**

Fig. 3 shows the normalized link occurrence/rank plots for inter-AS links when we create BGP paths from all stub ASes to the ASes in *Country1*, *Country8*, and *Country15* according to the two BGP policies. To clearly see the different skew of the link occurrence distribution of the two policies, we measure the slopes of link distributions in *log-log* scale in the high-rank region. Since the link-occurrence distribution of policy II is not modeled by Zipf-Mandelbrot distribution, we simply measure the slope in the high rank region to compare the skew. *Country1* has a barely observable skew in this region (i.e., the slope is less than 0.1) with policy II while it has a much higher skew (i.e., a slope of 0.44) with policy I. *Country8* and *Country15* have small skews (i.e., slopes of 0.21 – 0.23) with policy II and much higher skews of 1.10 – 1.24 with policy I. This suggests that, even though inter-domain Internet topology may have no physical bottlenecks (or very few, as in *Country8* or *Country15*), the BGP cost-minimization policy creates inter-domain routing bottlenecks.

#### 2.3.2 Cost minimization in intra-domain network topology and routing

**Internal AS router-level topology creates intra-domain routing bottlenecks:** Most ISPs build and manage hierarchical internal network structures for cost minimization reasons [46, 29] and these structures inherently create routing bottlenecks within ISPs. An ISP is composed of multiple points of presence (or PoPs) in different geographic locations and they are connected via few high-capacity *backbone* links. Within each PoP, many low-to-mid capacity *access* links connect the backbone routers to the border routers.

In general, ISPs aim to minimize the number of expensive long-distance high-capacity backbone links by multiplexing as much traffic as possible at the few backbone links; viz., HOT network model in [29]. As a result, backbone links naturally become routing bottlenecks. To show this, we carry out simulations using Tier-1 ISP topologies inferred by Rocketfuel [46]. We construct ingress-egress routes for all possible pairs of access routers using shortest-path routing [32]. Fig. 4 shows the simulated normalized link occurrence/rank for the three ASes belonging to different ISPs. In all three ASes, we find that the Zipf-Mandelbrot distribution fits accurately for high value of skew  $\alpha$  (i.e., 1.77 – 1.89) when  $\beta$  is deliberately fixed to 2.1 to yield a best fit and allow direct skew comparison. That is, a few AS internal links are extremely heavily used whereas most other internal links are very lightly used. Moreover, most of the heavily used links (i.e., 70%, 70%, and 90% of 10 most heavily used links in each of the three ISPs, respectively) are indeed backbone links that connect distant PoPs. We reconfirm the prevalence of intra-domain bottleneck links later in Section 2.4.1



**Figure 5: Percentage of link types of the 50 most occurred links for each of the 15 countries. Three link types (i.e., intra-AS links, inter-AS links, and IXP links) and three AS types (i.e., Tier-1, Tier-2, and Tier-3) are used for categorization.**

where we find that a large percentage (i.e., 30%) of links in routing bottlenecks are intra-AS links.

**Hot-potato routing policy in ISPs aggravates inter-domain routing bottlenecks:** The *hot-potato* routing policy is another example of a cost-minimization policy used by ISPs; i.e., this policy chooses the closest egress router among multiple egress routers to the next-hop AS [52]. As already reported [54], this policy causes a load imbalance at multiple inter-AS links connecting two ASes and thus aggravates the routing bottlenecks at the inter-AS links.

## 2.4 Characteristics of Bottleneck Links

In this subsection we investigate the characteristics of the links in the routing bottlenecks in terms of link types (e.g., intra-AS links, inter-AS links, or IXP links) and distance from the hosts in the target region (e.g., average router and AS hops) as a backdrop to the design of countermeasures against attacks that exploit bottleneck links. Our investigation suggests that the variety of link types found and their distribution make it *impractical* to design a single ‘one-size-fits-all’ countermeasure. Instead, in Section 5, we discuss several practical countermeasures that account for the specific bottleneck link types.

### 2.4.1 Link types

We consider three link types based on their roles in the Internet topology: intra-AS links, which connect two routers owned by the same AS; inter-AS links, which connect routers in two different ASes; and IXP links, which connect routers of different ASes through a switch fabric. Although the link types are clearly distinguished in the above definitions, the determination of link types via *traceroute* is known to be surprisingly difficult and error prone due to potential inference ambiguity [35]. For example, the *AS boundary ambiguity* [35] arises because routers at AS boundaries sometimes use IPs borrowed from their neighbor ASes for their interfaces. This is possible because the IPs at the both ends of the inter-AS links are in the same prefix. Borrowed IPs make it difficult to determine whether a link is an intra- or inter-AS link.

Our method of determining link type eliminates the AS boundary ambiguity by utilizing route diversity at the bottleneck links. Unlike previous measurements and analyses [35, 23], we measure a large number of disjoint incoming/outgoing routes to/from a bottleneck link. In other words, we gather all visible links 1-hop before/after the bot-

tleneck link, and this additional information helps us infer the link types at AS boundary without much ambiguity.<sup>3</sup>

Fig. 5 summarizes the percentage of the link types of the 50 most frequently found links for each of the 15 countries. The average percentage of all 15 countries is presented in the rightmost bar. Notice that the intra-AS and the inter-AS links are further categorized by the AS types; i.e., Tier-1, Tier-2, and Tier-3 ASes. The list of Tier-1 ASes is obtained from the 13 selected ASes in Renesys’ Baker’s Dozen.<sup>4</sup> ASes that have no customer but only providers or peers are Tier-3 ASes. The rest of the ASes are labeled as Tier-2 ASes.

Our investigation found two unexpected results. The first is that the intra-AS links are a major source of routing bottlenecks; see the rightmost bar in Fig. 5 where approximately 30% of routing bottleneck links are intra-AS links while the other 30% and 20% are inter-AS links and IXP links, respectively. (The balance of 20% is not determined due to lack of *traceroute* visibility). This high percentage of intra-AS bottleneck links contradicts the common belief that ISPs distribute routes over their internal links very well using complete knowledge of, and control over, their own networks. This result motivated us to investigate the practical challenges of route distribution within individual ISPs; viz., Section 5.3. The second unexpected result is that the majority of both intra-AS and inter-AS bottleneck links (i.e., 100% for intra-AS type and 81.2% for inter-AS type) is exclusively owned and managed by large ASes; e.g., Tier-1 or Tier-2 ASes. This implies that the Tier-1/Tier-2 ASes are the primary sources of bottleneck links.

### 2.4.2 Link distance

We also measure the *router-hop* and *AS-hop* distance of the bottleneck links from the hosts in the target regions. To measure a bottleneck link’s router-hop distance, we take the average of router-hop distances from the 1000 hosts in a region. A challenge in measuring the router-hop distance via *traceroute* is that some destinations used have firewalls in their local networks, which prevents discovery of the last few router hops from the destinations. When *traceroute* does not reach a destination we assume the presence of the destination immediately past the last hop found. Thus, the

<sup>3</sup>For IP to ASN mapping, we use the public IP-to-ASN mapping database by Cymru (<https://www.team-cymru.org/Services/ip-to-asn.html>).

<sup>4</sup><http://www.renesys.com/2014/01/bakers-dozen-2013-edition/>

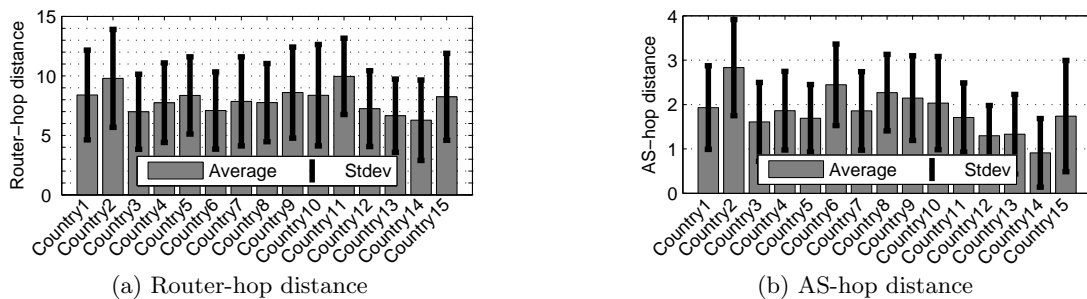


Figure 6: Router-hop (a) and AS-hop (b) distances of 50 bottleneck links for each of the 15 countries from the target regions.

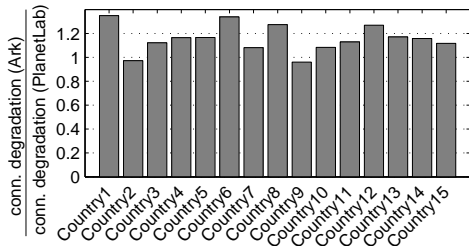


Figure 7: Connectivity degradation in the Ark dataset relative to the PlanetLab dataset for 50 flooding links selected from the routes measured by the PlanetLab nodes.

measured router-hop distance is a strict *lower-bound* of the average router-hop distance from destination hosts.

Fig. 6a shows the average and standard deviation of the router-hop distance of the 50 bottleneck links for each of the 15 countries. The average router-hop distance ranges from 6 to 10 router hops with average of 7.9 router hops and no significant differences were found across the 15 countries. Considering the average length of Internet routes is approximately 17 router hops [13], we conclude that the bottleneck links are located in the middle to the slightly closer to the target region on the routes to the target. The distance analysis is also consistent with the observation that the most bottleneck links are within or connecting Tier-1/Tier-2 ASes.

Fig. 6b shows the average and standard deviation of the AS-hop distance for the 15 countries. The average AS-hop distance from the target to the bottleneck links ranges from 1 to 3 AS hops with average of 1.84 AS hops. Again, the measured AS-hop distances are strict lower bounds of the average AS-hop distances due to limited *traceroute* visibility.

### 3. VALIDATION OF BOTTLENECK MEASUREMENTS

#### 3.1 Independence of Route Sources

One of the common pitfalls in Internet measurements is the dependency on vantage point; that is, the location where a measurement is performed can significantly affect the interpretation of the measurement [41]. Here we argue that our routing-bottleneck results are independent of the selection of route sources  $S$ . To show this, we validate our

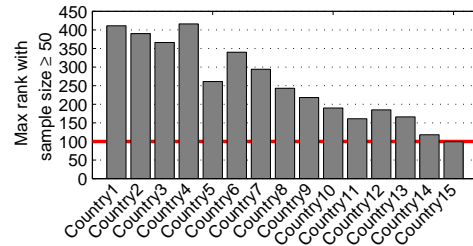


Figure 8: Maximum rank with sample size  $\geq 50$  for each of the 15 countries.

computation of routing-bottleneck results by comparing the connectivity degradation (defined in Section 4.2) calculated using the original source set  $S$  (i.e., 250 PlanetLab nodes) with that calculated using an *independent* source set  $S'$  (i.e., 86 Ark monitors),<sup>5</sup> as shown in Fig. 7. Notice that we select 50 bottleneck links for each country by analyzing the routes measured by PlanetLab nodes for *both*  $S$  and  $S'$ . The selection of these bottleneck links is discussed in Section 4.2. In most countries, the ratios of the two connectivity degradations are slightly higher than or very close to 1, which means that the bottlenecks of the PlanetLab dataset also become the bottlenecks of the independent Ark dataset. This also confirms the independence of the bottleneck-link flooding results from the choice of route-source distribution [26].

#### 3.2 Sufficiency of Link-Sample Size

Another common pitfall in Internet measurements aiming to discover statistical properties of datasets is the lack of a sufficiently large sample size; that is, it is possible that the sample size is insufficient to detect possible deviations from a discovered distribution. For reliable parameter estimates, the rule of thumb is that one needs to collect at least 50 samples for each element value [8]. Fig. 8 shows the maximum rank of the links (ordered by decreasing link occurrence) that are observed with at least 50 link samples for the 15 countries in our measurement. The figure shows that for all 15 countries, all the high ranked links (i.e., 0–100 rank) are observed with more than 50 link samples and thus

<sup>5</sup>The CAIDA’s Ark project uses 86 monitors distributed over 81 cities in 37 countries and performs *traceroute* to all routed /24’s. For consistent comparison, we use the Ark dataset that was measured on the same day when the PlanetLab dataset was obtained and select a subset of the measured traces in the Ark dataset that has the same AS destination,  $D$ , used in the PlanetLab dataset.

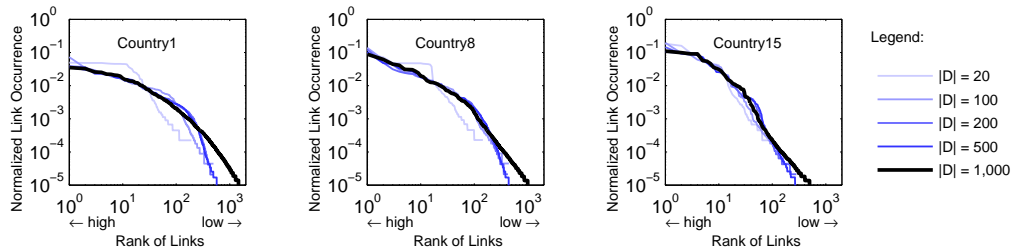


Figure 9: Normalized link occurrence/rank in traced routes to 1,000 randomly selected hosts in 3 countries.

the parameter estimates based on these links (i.e.,  $\alpha$  and  $\beta$  in Fig. 2a) are statistically sound.

Fig. 9 confirms that we have collected a sufficient number of link samples. In this figure, we illustrate the normalized link occurrence with the various sizes of disjoint  $D$  and observe how the link occurrence in the high rank region (i.e., rank 0–100) converges. We can conclude that  $|D| = 1,000$  is sufficient to discover the power-law distribution in the top-100 rank because it displays the same power-law distribution in the range as that observed with smaller size for  $D$ ; i.e.,  $|D| < 1,000$ . Thus, with a relatively small number of measurements one can learn the power-law distribution of the few but frequently observed high-rank links.

### 3.3 Traceroute Accuracy

In this subsection we review the common pitfalls in analyzing *traceroute* results and explain why they do not affect our measurement results.

**Inaccurate alias resolution:** As shown in many topology measurement studies, it is extremely important to accurately infer the group of interfaces located in the same router (or alias resolution) because its accuracy dramatically affects the resulting network topology [46, 37]. Highly accurate alias resolution still remains an open problem. Our measurements do *not* need alias resolution because we do not measure any router-level topology, but only layer-3 links (i.e., interfaces) and routes that use those links.

**Inaccurate representation of load-balanced routes:** Ordinary *traceroute* does not accurately capture load-balanced links and thus specially crafted *traceroute*-like tools (e.g., Paris *traceroute* [4]) are needed to discover these links. Our measurement does *not* need to discover load-balanced links because they cannot become the routing bottlenecks. Instead, we perform ordinary *traceroute* multiple times (e.g., 6 *traceroutes* in our measurement) for the same source-destination pair and ignore the links that do not always appear in multiple routes.

**Inconsistent returned IPs:** In response to *traceroute*, common router implementations return the address of the *incoming* interface where packets enter the router. However, very few router models return the *outgoing* interface used to forward ICMP messages back to the host launching *traceroute* [35, 36] and thus create measurement errors. However, our routing bottleneck measurement is not affected by this router behavior because (1) most of the identified router models that return outgoing interfaces are likely to be in small ASes since they are mostly Linux-based software routers or lower-end routers [35], and (2) we remove all load-balanced links that can be created by the routers which return outgoing interfaces [36].

**Hidden links in MPLS tunnels:** Some routers in MPLS tunnels might not respond to *traceroute* and this might cause serious measurement errors [59]. However, according to a recent measurement study in 2012 [11], in the current Internet, nearly all (i.e., more than 95%) links in MPLS tunnels are visible to *traceroute* since most current routers implement RFC4950 ICMP extension and/or *ttl-propagate* option to respond to *traceroute* [11].

## 4. ROUTING-BOTTLENECK EXPLOITS

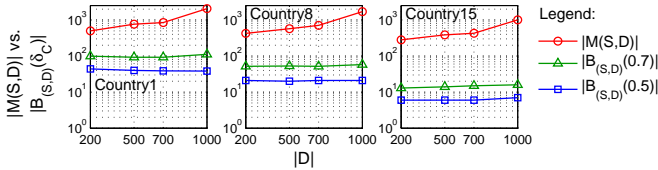
Bottleneck links provide a very attractive target for link-flooding attacks [6, 26]. By targeting these links, attacks become both *scalable* and *persistent*. Scalable because the number of targeted hosts can be increased substantially by flooding only few additional links, as shown later in this section. Persistent because adversaries can dynamically change the flooded link sets while maintaining the same targeted hosts, making the attacks undetectable by traditional anomaly detection methods. In this paper, we focus primarily on the scalability of link-flooding attacks; i.e., we discuss the selection of target bottleneck links, the expected degradation in connectivity to the targeted hosts  $D$ , and how to increase attack targets without much additional measurement effort and attack traffic. The persistency properties of routing-bottleneck exploits is discussed in detail in the Crossfire attack [26].

To measure the strength of a link-flooding attack, we first define an *ideal attack* that completely disconnects all routes from sources  $S$  to selected hosts of destinations  $D$ . Then, we define realistic attacks that can cause very substantial *connectivity degradation*.

### 4.1 Disconnection Attacks

Let  $S$  be the 250 PlanetLab nodes and  $D$  the 1,000 randomly selected hosts in the target region; e.g., a country or a city. For efficient disconnection attacks, the adversary needs to flood only *non-redundant* links; that is, flooding of a link should disconnect routes that have *not* been disconnected by the other already flooded links. Link redundancy can be avoided by flooding the *mincut* of the routes from  $S$  to  $D$ , namely the *minimum* set of links whose removal disconnects *all* the routes from  $S$  to  $D$ , which is denoted by  $M(S, D)$ . Note that our notion of the *mincut* differs from the *graph-theoretic* mincut. Our *mincut* is a set of links that cover all routes to chosen nodes whereas the graph-theoretic mincut is a set of (physical) link cuts for an arbitrary network partitioning. Thus, one cannot use well-known polynomial-time *mincut* algorithms of graph theory [53] for our purpose.

Finding  $M(S, D)$  can be formulated as the *set cover problem*: given a set of element  $U = \{1, 2, \dots, m\}$  (called the



**Figure 10: Measured sizes of mincuts,  $|M(S, D)|$ , and selected bottlenecks sizes,  $|B_{(S, D)}(\delta_C)|$ , for given degradation ratios  $\delta_C$  and varying  $|D|$  in 3 countries.**

universe) and a set  $\mathcal{K}$  of  $n$  sets whose union equals the universe, the problem is to identify the smallest subset of  $\mathcal{K}$  whose union equals the universe. Thus, our *mincut* problem can be formulated as follows: the set of all routes we want to disconnect is the universe,  $\mathcal{U}$ ; all IP-layer links are the sets in  $\mathcal{K}$ , each of which contains a subset of routes in  $\mathcal{U}$ , and their union equals  $\mathcal{U}$ ; the problem is to find the smallest set of links whose union equals  $\mathcal{U}$ . Since the set cover problem is NP-hard, we run a greedy algorithm [20] to calculate  $M(S, D)$ . The greedy algorithm, which is similar to the one used to find critical links in the Crossfire attack [26], iteratively selects (and virtually cuts) each link in  $M(S, D)$  until all the routes from  $S$  to  $D$  are disconnected.

Our experiments show that flooding an entire *mincut*,  $M(S, D)$ , in any of the fifteen countries and cities selected would be rather unrealistic. For example, approximately 83 Tbps would be required to flood a *mincut* of 2,066 links with 40 Gbps link capacity for a flooding attack against 1,000 servers in *Country1*.  $|M(S, D)|$  can be much larger than both  $|S|$  or  $|D|$  since in measuring the size of  $M(S, D)$  we exclude network links that are directly connected to hosts in  $S$  or  $D$ . Worse yet, Fig. 10 (top curves) shows that the *mincut* size,  $|M(S, D)|$ , grows as  $|D|$  grows. This implies that any practical link-flooding attack that disconnects *all* the hosts of a target region,  $D$ , must scale with an already large  $|M(S, D)|$ . However, as we show in the next section, an adversary does not need to flood an entire *mincut* to degrade connectivity of  $D$  hosts of a targeted region very substantially. Also, by taking advantage of the power-law distributions of bottleneck links, an adversary can scale attacks to very large target sets,  $D$ .

## 4.2 Connectivity Degradation Attacks

Feasible yet powerful connectivity-degradation attacks would flood much smaller sets of links to achieve substantial connectivity degradation to the routes from  $S$  to  $D$ . To measure the strength of such attacks we define a connectivity-degradation metric, which we call the *degradation ratio* [26], as follows:

$$\delta_{(S, D)}(B) = \frac{\text{number of routes that traverse } B}{\text{number of routes from } S \text{ to } D},$$

where  $B$  is the subset of the *mincut*  $M(S, D)$  links that are flooded by an attack.  $B$ 's size,  $|B|$ , is determined by an adversary's capability. Clearly, the maximum number of links that an adversary can flood is directly proportional to the maximum amount of traffic generated by attack sources controlled by the adversary; e.g., botnets or amplification servers. Here, we assume that the required bandwidth to flood a single link is 40 Gbps and thus the adversary should create  $40 \times n$  Gbps attack bandwidth to flood  $n$  links concurrently. Links with larger physical capacity (e.g., 100 Gbps)

have recently been introduced in the Internet backbone but the vast majority of backbone links still comprises links of 40 Gbps or lower capacity [24].

Fig. 11a and Fig. 11b show the expected degradation ratio calculated for each of the 15 countries and 15 cities for varying number of links to flood, or  $|B|$ , respectively. These countries and cities are ordered by increasing the averaged degradation ratio over the interval  $1 \leq |B| \leq 50$ . By definition, the degradation ratio for  $B$  (i.e.,  $\delta_{(S, D)}(B)$ ) is the sum of normalized occurrences of the links in  $B$ . Thus, degradation ratio can be accurately modeled by the cumulative distribution function (CDF) of the Zipf-Mandelbrot distribution since the normalized link occurrence follows this distribution. Parameters  $\alpha$  and  $\beta$  are listed in both figures. We note that the ordering of the degradation ratios in Fig. 11a and Fig. 11b is exactly the same as the ordering of the values of the distribution skew,  $\alpha$ , of the 15 countries and 15 cities in Fig. 2a and Fig. 2b, respectively. That is, countries and cities with low  $\alpha$  have a lower degradation ratio (i.e., are less vulnerable to flooding attacks) whereas countries and cities with high  $\alpha$  have high degradation ratio; i.e., are more vulnerable to flooding attacks. This confirms that the skew of the link-occurrence distribution,  $\alpha$ , of the Zipf-Mandelbrot distribution is a good indicator of vulnerability to link-flooding attacks.

Fig. 11a and Fig. 11b also show that the adversary can easily achieve significant degradation ratio (e.g., 40% - 82%) when flooding only few bottleneck links; e.g., 20 links. Given the proliferation of traffic amplification attacks achieving hundreds of Gbps or the extremely low costs of botnets, flooding several tens of bottleneck links of selected hosts in different countries around the world seems very practical.

### 4.2.1 Sizes of Bottlenecks

The size of a bottleneck selected for attack clearly depends on the *chosen degradation ratio*  $\delta_C$  sought by an adversary. This size is defined as:

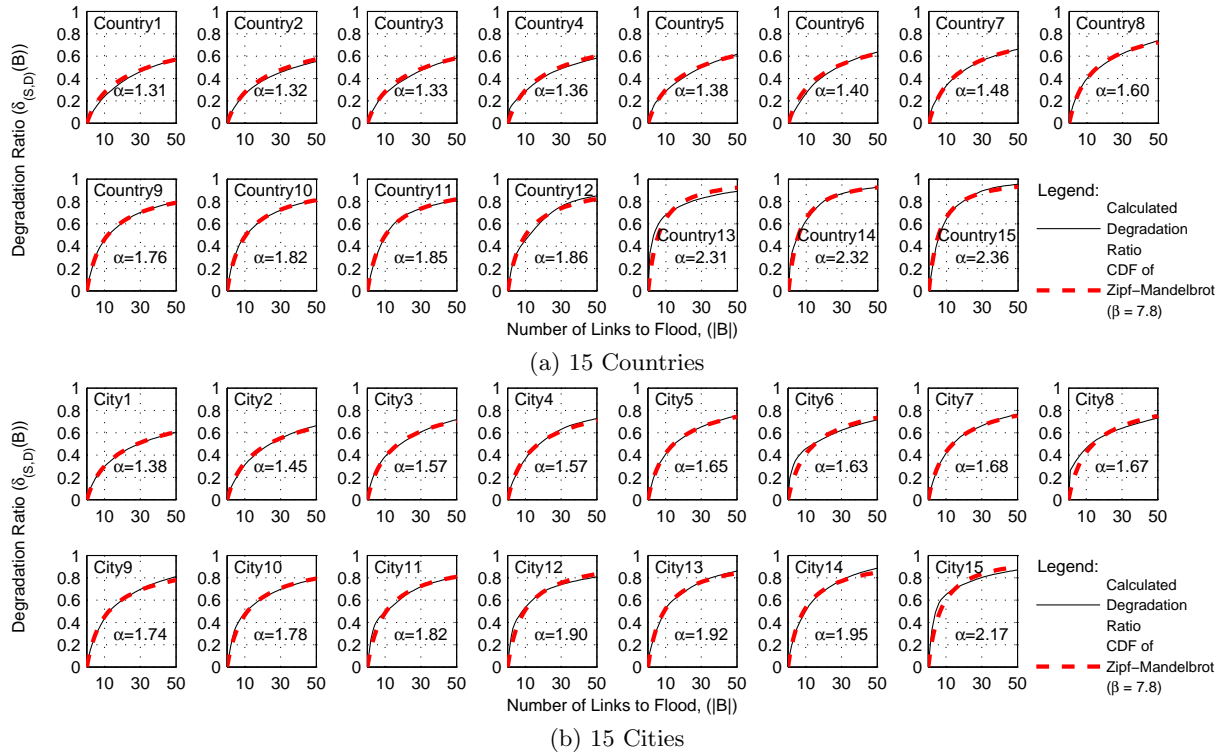
$$|B_{(S, D)}(\delta_C)| = \text{minimum } |B|, \text{ such that } \delta_{(S, D)}(B) \geq \delta_C.$$

The bottlenecks selected for attack,  $B_{(S, D)}(\delta_C)$ , are substantially smaller than their corresponding *mincuts*,  $M(S, D)$ . Fig. 10 shows the set sizes of the *mincuts* and the selected bottlenecks for chosen ratios  $\delta_C$  of 0.7 and 0.5 for varying sizes of  $D$ . The plots for the three countries show that  $|M(S, D)|$  is one to two *orders of magnitude* larger than  $|B_{(S, D)}(\delta_C)|$  in the entire range of measured  $|D|$  and  $\delta_C$ . In other words, the attack against the selected bottlenecks requires a much lower adversary's flooding capability than for a *mincut* while achieving substantial connectivity degradation; e.g., 70%.

### 4.2.2 Scaling the Number of Targets

Our experiments suggest that an adversary need not scale routing measurements and attack traffic much beyond those illustrated in this paper for much larger target-host sets (i.e.,  $|D| \gg 1,000$ ) in a chosen region to obtain connectivity-degradation ratios in the range illustrated in this paper. This is the case following two reasons. First, our measurements for multiple disjoint sets of selected hosts in a target region yield the *same* power-law distribution for different unrelated sizes of  $D$ ; viz., Fig. 9. Hence, increasing the number of routes from  $S$  to a much larger  $D$  will not increase the size of the bottlenecks appreciably. In fact, we have already





**Figure 11: Calculated degradation-ratio/number-of-links-to-flood for 1,000 servers in each of the 15 countries in (a) and 15 cities in (b).**

noted that, unlike the size of *mincuts*,  $|M(S, D)|$ , the size of the selected bottlenecks for a chosen degradation ratio  $\delta_C$ ,  $|B_{(S,D)}(\delta_C)|$ , does *not* change as  $|D|$  increases, as shown by the lower two curves of in Fig. 10. Second, we showed that routing-bottleneck discovery is independent of the choice of  $S$ , where  $|S| \gg |B|$ ; viz., Section 3.1. This implies that, to flood the few additional bottleneck links necessary for a much increased target set  $D$ , an adversary needs not increase the size of  $S$  and attack traffic appreciably.

## 5. COUNTERMEASURES

Defenses against attacks exploiting routing bottlenecks range from simple but naïve approaches to far-reaching structural countermeasures and operational countermeasures. We summarize these countermeasures, discuss their deployment challenges, and briefly evaluate their effectiveness. Naturally, defense mechanisms for *server-flooding* attacks (viz., [17]) are irrelevant to this discussion.

### 5.1 Naïve Approaches

The naïve approaches presented here are the most probable responses that the current networks would perform once the degradation attacks hit the routing bottleneck of any target region.

**Local rerouting:** Targeted networks can reactively change routes crossing flooded links so that the flooding traffic (including both legitimate and attack flows that are indistinguishable from legitimate flows) is distributed over multiple other local links. However, this might cause more collateral damage on the other local links after all.

**Traffic-intensity based flow filtering:** Typical mitigations for volumetric DDoS attacks detect and filter long-lived large flows *only* because otherwise they cannot run in real-time in large networks [28]. This countermeasure cannot detect nor filter attack flows in bottleneck links because these could be low-rate and thus indistinguishable from legitimate.

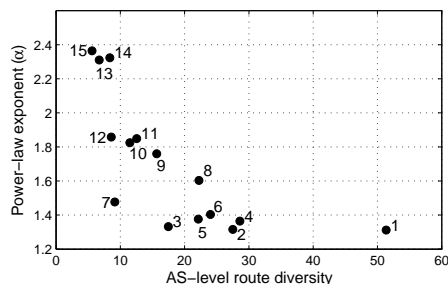
**Using backup links:** Typical backbone links are protected by the backup links, such that whenever links are cut, the backup links seamlessly continue to convey traffic. However, backup links cannot counter link-flooding attacks because they could be flooded too.

### 5.2 Structural Countermeasures

Structural countermeasures range from changes of physical Internet topology to those of inter-AS relationships. Although this type of countermeasures might require significant time (e.g., months) to implement, it could widen routing bottlenecks and decrease link-flooding vulnerability significantly. For example, if a country is connected to the rest of the world via only a handful of market-dominating ISPs, no matter how well routes are distributed, the country would inevitably experience routing bottlenecks. To remove these bottlenecks, the country would have to increase its route diversity through structural changes to its connectivity to the outside world.

Fig. 12 illustrates how AS-level structural changes could solve the routing-bottleneck problem. The x-axis is the metric called *AS-level route diversity* and it is calculated as

$$\frac{\{\text{number of intermediate ASes}\}}{\{\text{average AS hops from Tier-1 ASes to target region}\}}, \quad (1)$$



**Figure 12: Correlation between power-law exponent and AS-level route diversity in 15 countries.** (Legend: number  $i$  = Country $i$ , for  $i = 1, \dots, 15$ )

where the intermediate ASes are the ASes that connect the Tier-1 ASes with ASes located within each target region. The list of ASes within a target region are obtained from <http://www.nirsoft.net/countryip/>. Ratio (1) is a good proxy for measuring the AS-level route diversity because it represents the average number of possible ASes at every AS hop from the Tier-1 ASes to the target region. The y-axis is the power-law exponent  $\alpha$  obtained in Section 2.2. We see the clear correlation between the AS-level route diversity and the power-law exponent for the 15 countries, which supports our claim that the more AS-route diversity, the lower the power-law exponent.

We find that the Western European countries use significantly higher AS-level route diversity (i.e., 24.5 on average) than the rest of the countries (i.e., 16.5 on average), and thus are much less vulnerable to link-flooding attacks. This is undoubtedly due to long-standing policies (e.g., local-loop unbundling [5]) in European Union to stimulate ISP competition; e.g., to lower the cost of entry in ISP markets [15]. We believe that similar policies that promote ISP competition will increase route diversity and ultimately reduce the vulnerability to link-flooding attacks in other parts of the Internet.

## 5.3 Operational Countermeasures

Operational countermeasures could improve the management plane of various routing protocols (e.g., BGP or OSPF) to either decrease the skew of link-occurrence distribution or better react to the exploits. Although most of the countermeasures discussed here have been proposed in different contexts before (e.g., [54, 60, 22, 39]) their effectiveness in reducing routing bottlenecks is unknown. Hence, we briefly analyze these countermeasures here.

### 5.3.1 Dynamic inter-domain load balancing

As seen in Section 2.4.1, about 30% of the bottleneck links are inter-AS links. When an inter-AS link is flooded, at least one of the ASes should be able to quickly redirect the flooding traffic to relieve congestion. To do this, an AS would need to update its BGP announcements to its neighbours that use flooded links. However, inbound traffic redirection via updated BGP announcements [56] is not guaranteed since upstream ASes may have no positive incentives to re-route; i.e., upstream ASes would ignore these announcements whenever re-routing increases traffic cost. Even if neighbour ASes followed the updated BGP announcements, the long BGP convergence time (e.g., up to an hour [34])

would render them ineffective for timely response to link-flooding. Further delays would be incurred because outbound inter-AS level redirection requires human intervention in the current Internet. That is, inter-AS traffic redirection can only be manually configured since inter-AS links are selected by the coupling of BGP and IGP (e.g., OSPF) protocols [52]. Hence, added costs would become necessary to diffuse flooding traffic [54]. Therefore, timely and cost-effective reaction to inter-AS link flooding requires a dynamic mechanism that adaptively utilizes multiple parallel inter-AS links [54] and/or multiple AS-level route with different next-hop ASes [60]. The specific design of such mechanisms is beyond the scope of this paper.

### 5.3.2 Dynamic intra-domain load balancing

In principle, intra-domain load balancing can be an effective countermeasure because any balanced link *cannot* be the bottleneck. Recall that we remove any load-balanced links from our *traceroute* dataset for this reason; viz., Section 2.1. Many of today’s networks, especially those of large ISPs, deploy *intra-domain* load-balancing mechanisms based on the Equal-Cost Multi-Path (ECMP) algorithm [22]; e.g., approximately 40% of Internet routes [4] are load balanced by it. However, ECMP is insufficient to prevent routing bottlenecks. For example, ECMP requires that all intra-domain routes subject to load balancing have equal-cost paths – a condition which cannot always be satisfied for the alternate routes that happen to be available during flooding attacks. Moreover, given that only about 30% of the bottleneck links are intra-domain (viz., Section 2.4.1), ECMP cannot solve the overall problem. To prevent the degradation attacks from targeting their internal links, large ISPs should identify commonly used but not load-balanced links and dynamically reconfigure their networks (e.g., by updating link weights) so that the identified links are load-balanced with other links.

### 5.3.3 Intra-domain traffic engineering

One of the most widely used traffic engineering mechanisms is MPLS. As of 2013, at least 30% of Internet routes travel through MPLS tunnels [11] and they are mostly deployed in the large ISPs. Unlike the local rerouting solution discussed in Section 5.1, MPLS reconfiguration can perform fine-grained traffic steering to avoid collateral damage on the other links.

However, the widely used offline MPLS reconfiguration cannot be very effective since it can reconfigure tunnels only on a time scale ranging from tens of minutes to hours and days [14, 57]. Worse yet, the online MPLS reconfiguration, called the auto-bandwidth mechanism [39], which automatically allocates required bandwidth to each tunnel and change routes, is susceptible to sustained congestion. This is because it cannot detect congestion *directly* but only via reduced traffic rates caused by congestion. Thus, even an auto-bandwidth mechanism would require human intervention to detect link-flooding attacks [48] thereby slowing reaction time considerably. Therefore, large ISPs need an automated control system that monitors link congestion and quickly reconfigures the related MPLS tunnels to be steered through other underutilized links. We note that the recently proposed real-time traffic engineering techniques that leverage software-defined networking (SDN) have limited scope; i.e., in datacenter [43] and private wide-area networks [25, 21]. It is currently unknown whether SDN can implement

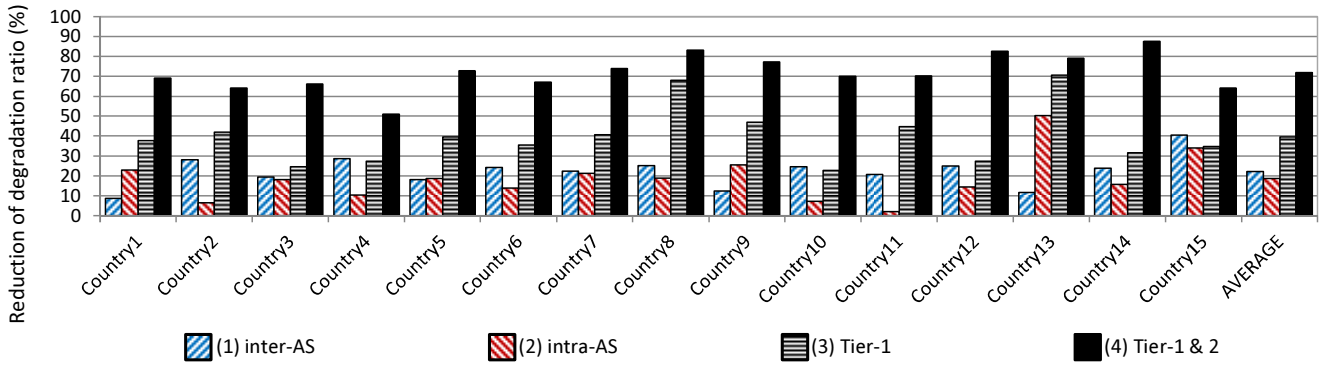


Figure 13: Reduction of degradation ratios due to four defense strategies when 20 bottleneck links are flooded for each of the 15 countries.

fast and robust traffic engineering in large public networks; i.e., Tier-1/Tier-2 networks.

#### 5.3.4 Effectiveness of operational countermeasures

We evaluate the reduction of degradation ratios due to the following four defense strategies using the operational countermeasures: (1) inter-domain load balancing at all inter-AS links, (2) intra-domain load-balancing and traffic engineering at all intra-AS links, (3) all operational countermeasures at all Tier-1 ASes, and (4) all operational countermeasures at all Tier-1 and Tier-2 ASes. In this evaluation, the types of all flooded links are known. Fig. 13 shows the reduction of degradation ratios in percentage for 15 countries. It shows that the defense strategies that protect a specific type of links (i.e., strategy (1) and (2)) are not very effective in general (approximately 20% reduction on average) because adversaries can still find bottleneck links from the other types of links. However, the defense strategies deployed by Tier-1 and/or Tier-2 ASes (i.e., strategy (3) and (4)) show much higher effectiveness: when all Tier-1 ASes implement all the operational countermeasures, the degradation ratio is reduced by 40%; and when all Tier-2 ASes also join the defense, 72% of reduction is achieved on the average. This confirms our previous observation that the large Tier-1 and Tier-2 ASes are primarily responsible for routing bottlenecks.

#### 5.3.5 Cost of operational countermeasures

In general, ISPs would incur negligible cost for the proposed operational countermeasures. It is well known that ISPs’ internal networks have substantial route diversity [51], and thus intra-domain countermeasures would not incur significant costs. Inter-domain load balancing would also incur low costs in most cases because large Tier-1 or Tier-2 ISPs could easily find multiple inter-domain links to balance traffic while maintaining the same next-hop AS.

### 5.4 Application Server Distribution

One might distribute application servers in different geographic locations, possibly using content distribution (e.g., Akamai [40]) and overlay networks (e.g., RON [3], SOS [27]) to distribute routes. The application servers have to be distributed in such a way that inherent route diversity is fully utilized; i.e., analysis must show that no routing bottleneck arises. However, this might not be practical for some domains such as industrial process systems, financial services,

or defense services where constrained geography may restrict application distribution.

## 6. RELATED WORK

### 6.1 Internet Topology Studies

A large body of research investigates the topology of Internet. Two long-term projects have measured the router-level Internet topology via *traceroute*-like tools: CAIDA’s Archipelago project [7] and DIMES project [45]. Rocketfuel [46] is another project that use approximately 800 vantage points for *traceroute* to infer major ISP’s internal topology. Together these studies provide important insights into the layer-3 topology of the Internet.

Our routing bottleneck measurement differs from the topology studies in two important ways. First, we do not measure or even infer the router-level topology but simply observe *how the routes are distributed* on the underlying router-level topology. Second, we do not need nor attempt to observe all the routes covering the entire address space but focus on the route-destination regions of potential adversary interest.

### 6.2 Topological Connectivity Attacks

Faloutsos *et al.* analyzed *traceroute* data and concluded that the node degree of the routers and ASes have power-law distribution [12]. Albert *et al.* confirmed the power-law behavior of the node-degree distribution and concluded that the Internet suffers from an ‘Achilles’ heel’ problem; i.e., targeted removal attacks against the small number of hub nodes with high node degree will break up the entire Internet into small isolated pieces [2].

The Achilles’ heel argument has triggered several counter-arguments. Some find that node-removal attacks are unrealistic because the number of required nodes to be removed is impractically high [31, 58]. Li *et al.* argue that the power-law behavior in node-degree distribution does not necessarily imply the existence of hub nodes in the Internet by showing that power-law node-degree distribution can be generated without hub nodes [29].

Our routing-bottleneck study discovers a new power-law distribution in the Internet. However, this power-law is *completely different* from that of the above-mentioned work for two reasons. First, we measure a power law for the link usage in Internet routes whereas the above-mentioned work finds power laws in the node-degree distribution. Second,

the scope of our power-law analysis is different; i.e., it is focused on, and limited to, a *chosen* route-destination region whereas the above-mentioned work analyzes the power-law characteristics of the entire Internet.

### 6.3 Bandwidth Bottleneck Studies

In networking research, the term ‘bottleneck’ has been traditionally used to represent the link with the smallest available bandwidth on a route; i.e., the link that determines the end-to-end route throughput. To distinguish it from a routing bottleneck, we call this link the *bandwidth* bottleneck. Several attempts have been made to measure bandwidth bottlenecks in the Internet; viz., BFind [1] and PathNek [23]. However, routing and bandwidth bottlenecks are fundamentally different as they do not necessarily imply each other. That is, routing bottlenecks are unrelated to the available bandwidth or provisioned link capacity, but closely related to the number of routes served by each link. Conversely, bandwidth bottlenecks can occur in the absence of any routing bottlenecks.

### 6.4 Control-Plane and Link-Flooding Attacks

Attacks that cause instability of the control plane in Internet routing [44] and link-flooding attacks [50, 26] have been recently proposed and launched in real life already [6]. Our past work on the Crossfire attack [26] presents a specific strategy to identify and flood a few targeted links for eight selected target areas in the US. In contrast to Crossfire, where we focus on the feasibility of flooding a small set of critical links, here we explore a *fundamental vulnerability* of today’s Internet, namely, pervasive routing bottlenecks that can be exploited by any flooding attack. We show the ubiquity of routing bottlenecks in various countries and cities around the world via extensive measurements, and identify their root cause. We also explore the characteristics of bottleneck links; e.g., link type and distance to targets. Last, we provide several practical countermeasures.

## 7. CONCLUSIONS

We introduce the notion of the routing bottlenecks, define it using power-law distributions of link occurrence in routes to chosen destinations, and show that it arises from route-cost minimization in the Internet. We identify the key characteristics of these bottlenecks in terms of size, link type, and distance from host destinations, and measure the degradation of host connectivity caused by attacks that flood bottleneck links. We show that the routing bottlenecks are pervasive and certain geographic regions (i.e., countries and cities) around the world are more susceptible than others to scalable link-flooding attacks. We present structural and operational countermeasures against these attacks and discuss their effectiveness. We argue that deployment of the proposed countermeasures does not require major Internet redesign and their cost is insignificant.

## 8. ACKNOWLEDGMENTS

We are grateful to Vyas Sekar, Soo Bum Lee, and the reviewers for their insightful comments and suggestions. This work was supported in part by the National Science Foundation under grant CCF-0424422 and a grant from PNC Bank. The views and conclusions contained in this document are solely those of the authors and should not be interpreted

as representing the official policies, either expressed or implied, of any sponsoring institution, the U.S. government or any other entity.

## 9. REFERENCES

- [1] A. Akella, S. Seshan, and A. Shaikh. An empirical evaluation of wide-area Internet bottlenecks. In *Proc. IMC*, 2003.
- [2] R. Albert, H. Jeong, and A.-L. Barabási. Error and attack tolerance of complex networks. *Nature*, 406(6794), 2000.
- [3] D. Andersen, H. Balakrishnan, F. Kaashoek, and R. Morris. Resilient Overlay Networks. In *Proc. ACM SOSP*, 2001.
- [4] B. Augustin, T. Friedman, and R. Teixeira. Measuring load-balanced paths in the Internet. In *Proc. IMC*, 2007.
- [5] M. Bourreau and P. Dogan. Unbundling the local loop. *European Economic Review*, 49(1), 2005.
- [6] P. Bright. Can a DDoS break the Internet? Sure... just not all of it. In *Ars Technica*, April 2, 2013.
- [7] CAIDA Monitors. The Archipelago Measurement Infrastructure, 2002.
- [8] A. Clauset, C. R. Shalizi, and M. E. Newman. Power-law distributions in empirical data. *SIAM Review*, 51(4), 2009.
- [9] DHS. Project Shine. *ICS-CERT Monitor, Quarterly Newsletter, Oct-Dec.*, 2012.
- [10] E. W. Dijkstra. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1), 1959.
- [11] B. Donnet, M. Luckie, P. Mérindol, and J.-J. Pansiot. Revealing MPLS tunnels obscured from traceroute. *ACM SIGCOMM CCR*, 42(2):87–93, 2012.
- [12] M. Faloutsos, P. Faloutsos, and C. Faloutsos. On power-law relationships of the Internet topology. *ACM SIGCOMM CCR*, 29(4), 1999.
- [13] A. Fei, G. Pei, R. Liu, and L. Zhang. Measurements on delay and hop-count of the Internet. In *IEEE GLOBECOM’98-Internet Mini-Conference*, 1998.
- [14] A. Feldmann, A. Greenberg, C. Lund, N. Reingold, J. Rexford, and F. True. Deriving traffic demands for operational IP networks: methodology and experience. *ACM SIGCOMM Computer Communication Review*, 30(4), 2000.
- [15] B. Fung. What Europe can teach us about keeping the Internet open and free. In *The Washington Post*, September 20, 2013.
- [16] L. Gao. On inferring autonomous system relationships in the Internet. *IEEE/ACM Transactions on Networking (TON)*, 9(6), 2001.
- [17] M. Geva, A. Herzberg, and Y. Gev. Bandwidth Distributed Denial of Service: Attacks and Defenses. *IEEE Security & Privacy*, 12(1), January 2014.
- [18] P. Gill, M. Schapira, and S. Goldberg. A survey of interdomain routing policies. *ACM SIGCOMM Computer Communication Review*, 44(1), 2014.
- [19] K. P. Gummadi, R. J. Dunn, S. Saroiu, S. D. Gribble, H. M. Levy, and J. Zahorjan. Measurement, modeling, and analysis of a peer-to-peer file-sharing workload. *ACM SIGOPS Operating Systems Review*, 37(5), 2003.
- [20] D. S. Hochbaum. Approximating covering and packing problems: set cover, vertex cover, independent set, and related problems. In *Approximation algorithms for NP-hard problems*. PWS Publishing Co., 1996.
- [21] C.-Y. Hong, S. Kandula, R. Mahajan, M. Zhang, V. Gill, M. Nanduri, and R. Wattenhofer. Achieving high utilization with software-driven WAN. In *Proc. of ACM SIGCOMM*, 2013.
- [22] C. E. Hopps. Analysis of an equal-cost multi-path algorithm. *RFC 2992*, 2000.
- [23] N. Hu, L. E. Li, Z. M. Mao, P. Steenkiste, and J. Wang. Locating Internet bottlenecks: algorithms, measurements, and implications. In *Proc. SIGCOMM*, 2004.

- [24] Internet2 Network – Layer3 / IP Connectors Map, 2013.
- [25] S. Jain, A. Kumar, S. Mandal, J. Ong, L. Poutievski, A. Singh, S. Venkata, J. Wanderer, J. Zhou, M. Zhu, et al. B4: Experience with a globally-deployed software defined WAN. In *Proc. of ACM SIGCOMM*, 2013.
- [26] M. S. Kang, S. B. Lee, and V. D. Gligor. The Crossfire Attack. In *IEEE Symposium on Security and Privacy*, 2013.
- [27] A. D. Keromytis, V. Misra, and D. Rubenstein. SOS: secure overlay services. In *Proc. of ACM SIGCOMM*, 2002.
- [28] R. Krishnan, M. Durrani, and P. Phaal. Real-time SDN Analytics for DDoS mitigation. *Open Networking Summit*, 2014.
- [29] L. Li, D. Alderson, W. Willinger, and J. Doyle. A first-principles approach to understanding the Internet’s router-level topology. *ACM SIGCOMM Computer Communication Review*, 34(4), 2004.
- [30] M. Luckie, B. Huffaker, A. Dhamdhare, V. Giotsas, et al. AS relationships, customer cones, and validation. In *Proc. IMC*, 2013.
- [31] D. Magoni. Tearing down the Internet. *IEEE Journal on Selected Areas in Communications*, 21(6), 2003.
- [32] R. Mahajan, N. Spring, D. Wetherall, and T. Anderson. Inferring link weights using end-to-end measurements. In *Proc. ACM SIGCOMM Workshop on Internet measurement*, 2002.
- [33] B. Mandelbrot. Information theory and psycholinguistics. *BB Wolman and E*, 1965.
- [34] Z. M. Mao, R. Govindan, G. Varghese, and R. H. Katz. Route flap damping exacerbates Internet routing convergence. *ACM SIGCOMM Computer Communication Review*, 32(4), 2002.
- [35] Z. M. Mao, J. Rexford, J. Wang, and R. H. Katz. Towards an accurate AS-level traceroute tool. In *Proc. of ACM SIGCOMM*, 2003.
- [36] P. Marchetta, V. Persico, E. Katz-Bassett, and A. Pescapé. Don’t trust traceroute (completely). In *Proc. CoNEXT Student workshop*, 2013.
- [37] P. Marchetta, V. Persico, and A. Pescapé. Pythia: yet another active probing technique for alias resolution. In *Proc. CoNEXT*, 2013.
- [38] M. Mimoso. 400 Gbps NTP Amplification attack alarmingly simple. In *Threatpost*, Feb. 13, 2014.
- [39] MPLS Traffic Engineering (TE)–Automatic Bandwidth Adjustment for TE Tunnels. [http://www.cisco.com/c/en/us/td/docs/ios/12\\_0s/feature/guide/fsteaut.html#wp1015347](http://www.cisco.com/c/en/us/td/docs/ios/12_0s/feature/guide/fsteaut.html#wp1015347).
- [40] E. Nygren, R. K. Sitaraman, and J. Sun. The Akamai network: a platform for high-performance Internet applications. *ACM SIGOPS Operating Systems Review*, 44(3), 2010.
- [41] V. Paxson. Strategies for sound Internet measurement. In *Proc. IMC*, 2004.
- [42] PlanetLab. <http://www.planet-lab.org/>.
- [43] J. Rasley, B. Stephens, C. Dixon, E. Rozner, W. Felter, K. Agarwal, J. Carter, and R. Fonseca. Planck: Millisecond-scale Monitoring and Control for Commodity Networks. In *Proc. of ACM SIGCOMM*, 2014.
- [44] M. Schuchard, A. Mohaisen, D. Foo Kune, N. Hopper, Y. Kim, and E. Y. Vasserman. Losing control of the Internet: using the data plane to attack the control plane. In *Proc. NDSS*, 2010.
- [45] Y. Shavitt et al. The DIMES Project, 2008.
- [46] N. Spring, R. Mahajan, and D. Wetherall. Measuring ISP topologies with Rocketfuel. *ACM SIGCOMM CCR*, 32(4), 2002.
- [47] R. Steenbergen. A practical guide to (correctly) troubleshooting with traceroute. *North American Network Operators Group*, pages 1–49, 2009.
- [48] R. Steenbergen. MPLS RSVP-TE Auto-Bandwidth - Lessons Learned. *NANOG 58*, 2013.
- [49] C. Strohm and E. Eric. Cyber Attacks on U.S. Banks Expose Computer Vulnerability. In *Bloomberg*, Sept. 27, 2012.
- [50] A. Studer and A. Perrig. The Coremelt attack. In *Proc. ESORICS*, 2009.
- [51] R. Teixeira, K. Marzullo, S. Savage, and G. M. Voelker. In search of path diversity in ISP networks. In *Proc. IMC*, 2003.
- [52] R. Teixeira, A. Shaikh, T. Griffin, and J. Rexford. Dynamics of hot-potato routing in IP networks. *ACM SIGMETRICS Performance Evaluation Review*, 32(1), 2004.
- [53] H. Thomas et al. *Introduction to algorithms*. MIT press, 2009.
- [54] P. Verkaik, D. Pei, T. Scholl, A. Shaikh, A. C. Snoeren, and J. E. Van Der Merwe. Wresting Control from BGP: Scalable Fine-Grained Route Control. In *USENIX ATC*, 2007.
- [55] P. Vogt. Minimum cost and the emergence of the Zipf-Mandelbrot law. In *Proc. Artificial Life IX*, 2004.
- [56] J. H. Wang, D. M. Chiu, J. C. Lui, and R. K. Chang. Inter-AS inbound traffic engineering via ASPP. *IEEE Transactions on Network and Service Management*, 4(1), 2007.
- [57] N. Wang, K. Ho, G. Pavlou, and M. Howarth. An overview of routing optimization for Internet traffic engineering. *IEEE Communications Surveys Tutorials*, 10(1), 2008.
- [58] Y. Wang, S. Xiao, G. Xiao, X. Fu, and T. H. Cheng. Robustness of complex communication networks under link attacks. In *Proc. ICAIT*, 2008.
- [59] W. Willinger, D. Alderson, and J. C. Doyle. Mathematics and the Internet: A Source of Enormous Confusion and Great Potential. *Notices of the AMS*, 56(5), 2009.
- [60] W. Xu and J. Rexford. MIRO: Multi-path Interdomain ROuting. In *Proc. of ACM SIGCOMM*, 2006.