

# The Impact of Internet Policy and Topology on Delayed Routing Convergence

Craig Labovitz, Ahba Ahuja  
Arbor Networks

{labovit, ahuja}@arbournetworks.com

Roger Wattenhofer, Srinivasan Venkatachary  
Microsoft Research

{rogerwa, cheenu}@microsoft.com

*Abstract—*

This paper examines the role inter-domain topology and routing policy play in the process of delayed Internet routing convergence. In recent work, we showed that the Internet lacks effective inter-domain path fail-over. Unlike circuit-switched networks which exhibit fail-over on the order of milliseconds, we found Internet backbone routers may take tens of minutes to reach a consistent view of the network topology after a fault. In this paper, we expand on our earlier work by exploring the impact of specific Internet provider policies and topologies on the speed of routing convergence. Based on data from the experimental injection and measurement of several hundred thousand inter-domain routing faults, we show that the time for end-to-end Internet convergence depends on the length of the longest possible backup autonomous system path between a source and destination node. We also demonstrate significant variation in the convergence behaviors of Internet service providers, with the larger providers exhibiting the fastest convergence latencies. Finally, we discuss possible modifications to BGP and provider routing policies which if deployed, would improve inter-domain routing convergence.

*Keywords—*Routing, Convergence, BGP, Network Measurement

## I. INTRODUCTION

The Internet's sustained exponential growth and the continued emergence of new and varied network applications provides testament to the scalability of the backbone infrastructure and protocols. The original TCP/IP decision to place network intelligence and state almost exclusively on end-nodes has enabled a diverse progeny of applications ranging from MP3 file exchange to collaborative learning. This scalability, however, comes at a price. Since its commercial inception in 1995, the Internet has lagged behind the public switched telephone network (PSTN) in availability, reliability and quality of service (QoS). This relative lack of reliability stems in part from the absence of intermediate backbone state and synchronization between routers. Despite the remarkable tolerance demonstrated by today's end-users for failures and delays in email and web services, the relative lack of Internet backbone reliability poses a significant challenge for emerging transaction-oriented and interactive applications like Internet telephony, online business and laboratories.

Although recent advances in the IETF's Differentiated Services working group promise to improve the performance of application-level services within some networks, across the wide-area Internet these QoS algorithms are usually predicated on the existence of a stable underlying forwarding infrastruc-

ture. In recent work, we showed that the Internet lacks effective inter-domain path fail-over [1]. Specifically, we found that multi-homed Internet sites may experience periods of degraded performance as well as complete loss of connectivity persisting fifteen minutes or more after a single fault.

We showed that most of the latency in Internet fail-over stems from *delayed convergence*, or the temporary routing table fluctuations generated during the operation of the path selection process on Internet backbone routers after a fault. Unlike switches in the public telephony network which exhibit failover on the order of milliseconds, our analysis found that inter-domain routers in the packet switched Internet may take several minutes to reach a consistent view of the network topology after a fault.

The current Internet inter-domain routing protocol, BGP, evolved from earlier distance vector routing algorithms. These protocols, including RIP [2], suffer from a number of well-documented problems, including slow convergence times [3]. Distance vector routing requires that each node maintain the distance from itself to each possible destination and the vector, or neighbor, to use to reach that destination. Whenever this connectivity information changes, the router transmits its new distance vector to each of its neighbors, allowing each to recalculate its routing table. The count-to-infinity problem [2] provides the canonical example used to illustrate the slow convergence in distance vector routing.

The adoption of the path vector in BGP is widely and incorrectly believed to have "solved" the routing table fluctuation problems exhibited by RIP. Instead, we showed in [1] that the adoption of the path vector exponentially exacerbates the number of potential routing table fluctuations. Specifically, we found that a default configuration (i.e. one without additional administratively added policies or filters) of  $n$  BGP autonomous systems connected in a complete graph may potentially explore  $n!$  routes, or all possible paths of all possible lengths between each AS after a fault. This upper theoretic bound on BGP convergence compares poorly with earlier routing protocols, such as RIP which have been shown to have  $O(n^3)$  computational complexity [4].

We based our earlier analysis on a simplified, abstract model of BGP interconnectivity. This model neglected the impact of routing policies, more realistic timing assumptions and inter-AS connectivity on the process of delayed convergence. Al-

though our initial model provides a useful theoretic upper bound on BGP distributed computation, it did not provide a practical framework in which to explore measured Internet convergence latencies. In this work, we expand on our earlier effort by exploring the measured convergence behaviors of “real” topologies, including more than 20 unique BGP route advertisements between more than 200 pairs of Internet service providers (ISPs). We also provide analysis of BGP behavior in general network topologies and under other more realistic assumptions. Our major results include:

- The time complexity for Internet fail-over convergence is upper bounded by  $30n$  seconds, where  $n$  is the length of the longest alternative ASPath between the source and any destination autonomous system for a route.
- On average, routes from customers of larger ISPs exhibit faster convergence than routes announced by customers of smaller Internet providers.
- Errant paths are frequently explored during delayed convergence. These “vagabond” paths likely stem from misconfiguration or software bugs.
- The majority of default-free Internet routes exhibit multiple alternative secondary paths. These paths often include several times the number of associated BGP autonomous systems in the ASPath as the steady state paths observed in routing table snapshots.

The remainder of this paper is organized as follows: In Section II, we provide some background and related work. Section III discusses our experimental data collection infrastructure. In Section IV, we present survey results on ISPs policy mechanisms and discuss the impact of these policies on the flow of routing information. In Section V, we present both empirical observations as well as quantitative analysis of the relationship between specific Internet topological configurations and the rate of convergence. We demonstrate a relationship between the convergence delay of a route announced between two providers and the longest ASPath allowed by the topology and policy between both providers. In Section VI, we present a proof of this relationship. Finally, we conclude in Section VII with a discussion of modifications to BGP which, if deployed, would significantly improve inter-domain routing convergence.

## II. BACKGROUND

In this Section, we provide a brief review of the more salient aspects of BGP inter-domain routing related to the discussion in this paper. We assume that the reader is familiar with Internet routing concepts and terminology discussed in [5], [6], [1].

As a path vector protocol, BGP updates include an *ASPath*, or a sequence of intermediate autonomous systems between source and destination routers that form the directed path for the route. BGP uses the *ASPath* for both loop detection and policy decisions. Upon receipt of a BGP update, each router evaluates the path vector and invalidates any route which includes the router’s own AS number in the path.

Although not specified in the BGP standard, most vendor implementations ultimately default to best path selection based on *ASPath* length. The number of ASes in the path is used in a manner similar to the metric count attribute in the RIP. While BGP allows for path selection based on policy attributes, including

local preference and multi-exit discriminator values, the majority of ISP policies ultimately default to the selection of the route with the shortest path. In the remainder of this paper, we base our analysis on such constrained shortest path first policies.

The BGP standard also includes a minimum route advertisement interval timer, abbreviated in this paper as *MinRouteAdver*, which specifies a minimum amount of time that must elapse between advertisements of routes for a particular destination from a given BGP peer. This timer provides both a rate-limiter on BGP updates as well as a window in which BGP updates with common attributes may be bundled into a single update for greater protocol efficiency. The standard recommends thirty seconds as the *MinRouteAdver* interval plus/minus some additional random jitter.

A number of recent studies, including Varadhan et al. [7] and Griffin and Wilfong [8] have explored BGP routing *divergence*. BGP allows the administrator of an autonomous system to specify arbitrarily complex policies. In BGP divergence, Griffin and Wilfong show that it is possible for autonomous systems to implement “unsafe,” or mutually unsatisfiable policies, which will result in persistent route oscillations. In [9], Gao et al. prove that adherence to specific common ISP policies, including provider and customer relationships, will guarantee convergence.

The authors of all the above papers note that BGP divergence remains a theoretical finding and has not been observed in practice<sup>1</sup>. Our work explores a complimentary facet of BGP routing – the convergence behavior of inter-domain routers under the default BGP path selection policies. In this paper and [1], we show that even under constrained policies, the BGP Internet routing exhibits an order of magnitude longer convergence latencies than previously believed.

An increasing number of Internet customers today choose to *multi-home*, or provision external connectivity through multiple ISPs. This provider redundancy is designed to secure against single link, router or even ISP failures. In [1], we showed that the convergence delay associated with route failure is equivalent to the delay of multi-homed failover. In the remainder of this paper, we focus our analysis on convergence following a route withdrawal, or  $T_{down}$  event, for clarity of presentation.

A number of studies, including [10], [11] have explored the inter-domain topology and diameter of the Internet. These studies typically build topological maps based on periodic snapshots of network routing tables or active traceroutes. Our work focuses on a complimentary aspect of Internet topology – the set of all *possible* paths between source and destination autonomous systems. Although in steady-state, Internet routers will normally select the shortest advertised *ASPath* to a given destination, we show in Section VI that BGP routers may explore all possible longer *ASPaths* to the destination following a failure.

## III. METHODOLOGY

Our study builds on the experimental infrastructure originally developed in [1]. Our measurement and fault injection apparatus consists of Unix-based probe machines maintaining geographically and topologically diverse BGP peering sessions with more

<sup>1</sup>Recent analysis of routing problems in a large ISP backbone may provide the first such real-world example of persistent routing oscillation. See NANOG 21 conference proceedings at <http://www.nanog.org>.

than 20 ISPs. While in [1] we observed the impact of faults injected into only two Internet providers, in this work we expand our instrumentation to inject BGP route transitions (i.e. announcements and withdraws) into more than 10 geographically and topologically diverse providers.

Software from the MRT and IPMA projects [12], [13] running on both FreeBSD PCs and Sun Microsystems workstations was used to generate BGP routing update messages at random intervals of roughly a two-hour periodicity. The faults simulated route failures and repairs. In [1], we showed that the convergence behavior of route failures is equivalent to multi-homed failover.

We generated faults over a six month period to provide statistical guarantees that our analysis was based on deliberately injected faults rather than normally occurring exogenous Internet failures, which the authors in [14] found occur on the average of once a month. Each cooperating provider agreed to both accept our fault-injection announcements and treat the address space as a customer address block with respect to policy and filtering. As we only injected routing information for addresses assigned to our research effort, these faults did not impact routing for commodity ISP traffic with the exception of the addition of some minimal level of extra routing control traffic.

While one set of probe machines actively injected faults, we observed the impact of these faults through passive instrumentation of an additional twenty ISP default-free routing tables. Again using software from the MRT and IPMA projects, we logged both periodic routing table snapshots and all BGP routing updates received by our “RouteTracker” probe machines from the 20 peers to disk. We then correlated the data between our NTP synchronized fault injection and measurement probe machines. These correlations provided data on the convergence delays between multiple source and destination ISP peers. We inferred the steady-state and convergence topologies between all probed ISP pairs using the ASPath information included in BGP update messages advertised to the passive Routeviews probe machine.

In addition to our experimental measurements, we surveyed a broad spectrum of Internet backbone providers about the details of their routing and peering policies. Responses from 15 backbone providers of varied network size and topologies provide the framework in which we discuss the impact of specific filtering and policy implementation mechanisms on the process of delayed convergence.

At the request of the providers participating in our study, we anonymize the AS numbers, IP addresses and names of all providers in our examples and accompanying discussions. We use the anonymized Internet provider names and AS numbers consistently throughout the paper.

#### IV. POLICY

In this Section, we explore how routing policies and policy implementation mechanisms can impact both the number and length of possible ASPaths associated with a given route.

The Internet retains a significant physical interconnection hierarchy with several “tiers” of service providers. In [9], Gao et al. describe the provider use of filters to implement several types of commercial peering relationships. In Figure 1, we present In-

<b>Transit relationships -- Inbound</b>	
Prefix filters to receive customer routes only	100%
<b>Peer relationships -- Outbound</b>	
Communities	73%
Prefix filters and ASPaths	13%
Prefix filters only	13%
<b>Peer relationships -- Inbound</b>	
Bogon Filtering Only	80%
No Filtering	20%

Fig. 1. Survey results of Internet provider filtering mechanisms.

ternet provider survey data on the specific mechanisms used to implement these route filtering policies. Each row lists the percentage of surveyed ISPs that implement the specific filtering mechanism on their border routers. We separate filtering mechanisms into three broad categories based on the type of inter-provider relationship. For each grouping, an ISP typically will implement only one mechanism in that set. We illustrate the impact of these policy mechanisms through several examples using the topology shown in Figure 2.

We first see in Figure 1 that all surveyed providers use prefix filters to limit *inbound* acceptance of customer announcements to only “legitimate” address space assigned to that customer. We provide an illustration of this policy in Figure 2 to provide the framework for later discussion in this Section. In this example we assume ISP D filters the peering session with ISP G to only accept ISP G’s backbone and customer routes. Following the hierarchy upwards. ISP A similarly filters the peering with D to only accept backbone and customer routes from ISP D. Since ISP G provides transit to ISP D, ISP A also accepts ISP G’s routes from ISP D.

Although all surveyed providers share the same inbound policy mechanisms, the choice of outbound route filtering mechanisms differed markedly. The outbound peer relationship category in Figure 1 shows that 73 percent of surveyed ISPs advertise routes to peers based on community attributes. We also see that 26 percent of ISPs control their outbound advertisements to peers through some combination of prefix and ASPath filters. Thirteen percent implement both ASPath and prefix filter mechanisms and another 13 percent only prefix filter. For example, in Figure 2 we assume ISP A implements both prefix and ASPath filters on its outbound route announcements to peers. For routes learned from customer D, ISP A’s filters ensure that both A’s advertisements to peers match D’s address space and that all advertised routes come directly from ISP D (i.e. ISP D is the first AS in the ASPath). Based on these filters, we observe that ISP A may advertise routes with the paths “D G” and “D”, but not “C D G”.

The combination of ASPath and prefix filters prevents the unintentional creation of back-up transit paths. If ISP A only implements prefix filters, then after a failure between ISP A and ISP D, ISP A might learn ISP D’s routes from ISP C with an ASPath of “C D” and “C D G”. Lacking ASPath filters, ISP A will advertise these “C D G” routes to peer ISP B and, thus, provide transit to ISP B for ISP C.

As another example, we consider a tier-2 provider, ISP D, multi-homed to two upstream providers, ISP A and ISP C. ISP D also maintains peer relationships with tier-2 providers, ISP E

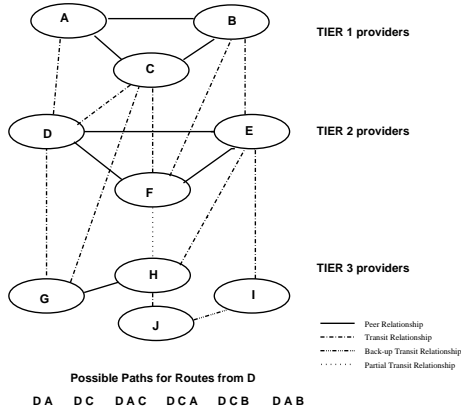


Fig. 2. Example illustrating impact of provider policies on the propagation of routing information.

and ISP F, and provides transit to ISP G. We list the set of possible paths for routes originating in ISP D below the example. We observe that both ISP C and A do not reannounce any of the routes learned from peers (non-transit/customer routes) to other peers. Thus, the paths “D A C B” and “D C A B” are not valid. However, both the “D A C B” and “D C A B” would be valid after a failure if ISP A and ISP C provide back-up transit for each other. In the following sections, we will see how these relationships control the number of possible paths explored during convergence.

In general, our survey data shows that smaller tier providers tend to possess a higher degree of peer and transit interconnectivity than larger providers. For example, by definition tier-1 providers do not purchase transit or generally maintain backup transit relationships with other providers. We will show in the next Section that these difference between tier-1 and tier-2 policies often lead to more numerous and longer alternative ASPaths for customers of tier-2 ISPs.

Finally, Figure 1 illustrates an important aspect of current filtering practices – the relative lack of filters applied to routes received from peers. Eighty percent of providers filter only “bogon” peer advertisements, or prefixes which represent private address space, default, unallocated address space, etc. Due to the technical and contractual difficulties of maintaining filter lists for the large number of routes advertised by peers, most providers resort to trusting their peers to send only valid information. We note that in a few well-publicized incidents, this trust has proven catastrophic for Internet routing [15].

## V. EXPERIMENTAL RESULTS

In this section, we present both empirical observations and analysis of the data collected by our fault injection experiments. We first provide several examples illustrating the process of delayed convergence. In the second subsection, we explore the impact of specific topological factors on convergence latencies.

### A. Effect of Topology on Convergence

During the six months of our study, we analyzed the routing topologies between more than 200 pairs of Internet providers. We graph only three representative topologies in Figure 3 for clarity. We note that all of the other monitored topologies in our

study exhibited related behaviors.

In the below examples we focus on the interaction of MinRouteAdver timers between neighboring BGP speaking routers. We explore the role these timers and the ASPath length of alternative paths play in the time required for convergence.

In general, we expect initial routing information to propagate more slowly via longer paths than shorter as each BGP speaker along the path adds between 0 and 30 seconds of MinRouteAdver delay. As most router vendors have implemented the MinRouteAdver timer on a per peer (instead of per prefix basis), the exact setting of MinRouteAdver may depend on other route instabilities. Specifically, we note that the per peer MinRouteAdver timer value initially follows a uniform probability distribution between 0 and 30 seconds. After propagation of the initial update, however, all subsequent MinRouteAdver timer values will delay updates for at least 30 seconds. As we discuss later in this Section, path selection depends both on individual MinRouteAdver settings as well as on the interaction, or *MinRouteAdver interference*, of multiple paths in a topology. Finally, in Section VI we show that that length of the longest possible ASPath between two nodes provides the upper bound on convergence delay.

The three subfigures in Figure 3 show a subset of the primary and alternative paths announced to our Routeviews machine by a single Japanese Internet provider, ISP4, after we withdrew routes  $R_1$ ,  $R_2$ , and  $R_3$  from our Mae-West exchange point BGP peering sessions with providers IS1, ISP2 and ISP3. For convenience, we will refer to these three ISPs at Mae-West as our *immediate providers*. The arrows represent the flow of routing information as inferred from the BGP ASPath update information announced by ISP4.

In each diagram, we label the *steady-state path*, or the path normally selected by ISP4 in the absence of a fault. The steady-state paths include IS1-ISP4 in Figure 3(a), ISP2-ISP4 in (b) and ISP3-ISP4 in (c). Similarly, we label backup paths chosen by ISP4 in each diagram with the letter P followed by integers denoting the frequency with which we observed that backup path (i.e. P1 to P6 in Figure 3(c)). For clarity, we graph only the most common backup paths observed during our study. In addition to the paths illustrated, ISP4 announced an additional 11 unique paths for  $R_2$  and 7 additional paths for  $R_3$  after 23 and 27 percent of the faults, respectively. We note that ISP4 only announced a single backup path throughout the course of our study for the topology in Figure 3(a).

#### A.1 Routes from ISP1

In steady state, we first observe from Figure 3 that ISP4 maintains a direct BGP peering session with all three of our immediate providers. Active ICMP ping and traceroute measurements show these three steady-state paths exhibit similar loss and latency characteristics. Although the steady-state paths are similar, we observe significant variation in both convergence latencies and the topologies explored by ISP4 for each of the three routes following the injection of a fault.

Figure 4(a) provides a numerical breakdown of the alternate path data graphically represented in Figure 3. The first column shows the distribution of paths announced by ISP4 after  $R_1$  is withdrawn from ISP1. The next two columns show similar re-

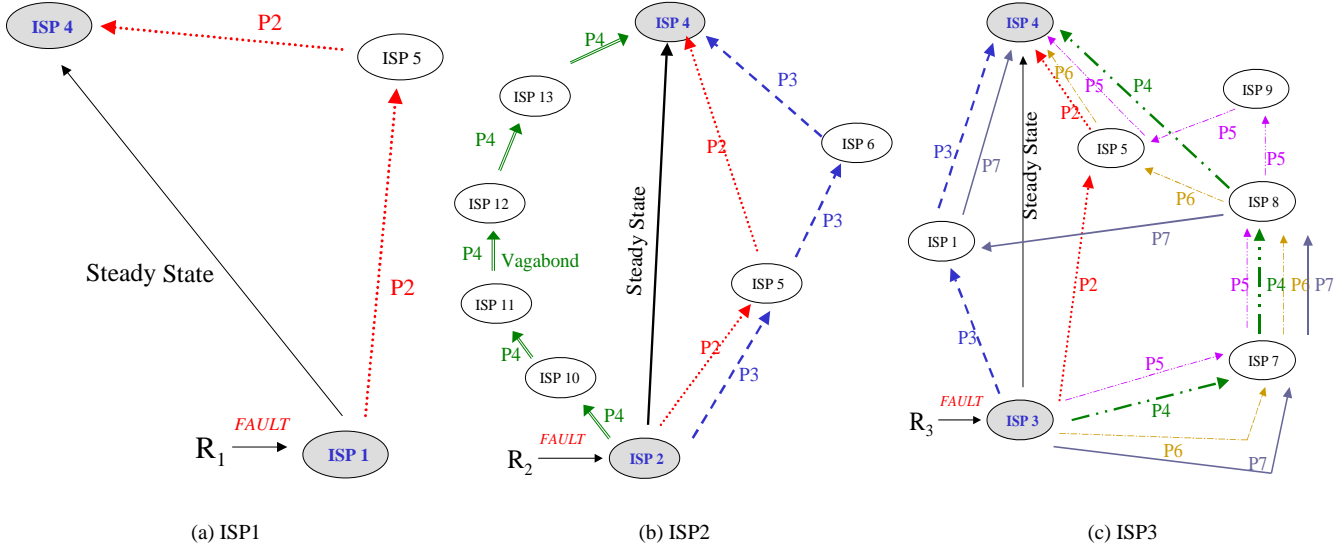


Fig. 3. Subset of secondary paths explored during the process of delayed convergence for routes from a Japanese provider to three different ISP routers at the California Mae-West exchange point.

Delayed Convergence of R1	Delayed Convergence of R2	Delayed Convergence of R3
96% Average: 92 (min/max 63/340) seconds	63% Average: 79 (min/max 44/208) seconds	36% Average: 110 (min/max 78/135) seconds
Announce AS4 AS5 AS1 (44 seconds)	AS4 AS5 AS2 (35 seconds)	Announce AS4 AS5 AS3 (52 seconds)
Withdraw (92 seconds)	Withdraw (79 seconds)	Withdraw (110 seconds)
4% Average: 32 (min/max 27/38) seconds	7% Average: 88 (min/max 80/94) seconds	35% Average: 107 (min/max 91/133) seconds
Withdraw (32 seconds)	Announce AS4 AS5 AS2 (33 seconds)	Announce AS4 AS1 AS3 (39 seconds)
	Announce AS4 AS6 AS5 AS2 (61 seconds)	Announce AS4 AS5 AS3 (68 seconds)
	Withdraw (88 seconds)	Withdraw (107 seconds)
	7% Average: 54 (min/max 29/9) seconds	2% Average: 140.00 (min/max 120/142)
	Withdraw (54 seconds)	Announce AS3 AS8 AS7 AS3 (27 seconds)
	23% <b>Other</b>	Announce AS3 AS5 AS8 AS7 AS3 (86 seconds)
		Withdraw (140 seconds)
		27% <b>Other</b>

Fig. 4. Ordered list of the most common ASPath sets announced by ISP4 during the process of delayed convergence following the withdrawal of routes  $R_1$ ,  $R_2$  and  $R_3$  from three Internet providers at the Mae-West exchange point.

results for routes  $R_2$  and  $R_3$  withdrawn from ISP2 and ISP3, respectively. Each entry provides a partial list of ASPath sets ordered by their measured frequency, and an associated average, minimum and maximum convergence delay. For example, the first entry for ISP1 indicates that after 96 percent of  $R_1$  failures, ISP4 announces a backup path “AS4 ASA5 AS1” (in an average of 44 seconds) followed by a withdrawal an average of 92 seconds after the fault. This backup path announcement corresponds to the primary backup path, P2, labeled in Figure 3(a). The second entry for  $R_1$  in Figure 4 denotes that after four percent of failures, ISP4 withdrew  $R_1$  without an intervening announcement of any backup path. We will provide probable explanations for these behaviors later in this Section.

## A.2 Routes from ISP2

In Figure 3(b) and the second column of Figure 4, we provide a slightly more complex example of the process of delayed convergence. After 63 percent of  $R_2$  failures, we again see that ISP4 first fails-over to the primary backup path, P2, followed by a withdrawal. After less than 7 percent of faults, however, we observe that ISP4 fails-over to backup path P2 followed by P3 before a final withdrawal. For another 7 percent of faults, ISP4 immediately withdraws the route. Finally, after the remaining 23 percent of faults, we observe more than 45 sequences of “other” ASPath set announcements, each with a frequency of well less than one percent. These other ASPath sequences include an additional 11 unique paths as well as inter-mixed withdrawals from ISP4. Analysis of this “other” category and discussions with In-

ternet providers suggests that the majority of these rare backup routes represent transient paths due to router misconfiguration errors.

Throughout the six months of our study, we observed frequent examples of these misconfigured, or *vagabond* paths between the majority of the 200 pairs of Internet providers we monitored. We define a vagabond path as a backup route which persists for a brief, fixed period of time (usually only several days) and does not conform to any intended or published policies. We based our classification of vagabond paths on automated tool inferences and lengthy discussions with Internet providers. For example, over the course of two weeks we observed the backup path P4 in Figure 3(b) after less than five percent of faults for *R2*. In the course of its travels from Mae-West to Japan, P4 traverses an additional four ISPs and transits one small Mediterranean country. ISP2 responded to our inquiries about this unusual path and pointed the blame at a single border router access-list configuration error which was subsequently resolved. We were able to partially automate detection of this and other vagabond paths as many of these erroneous routes transited the same misconfigured router, or Internet provider. This example of a vagabond path emphasizes an important aspect of Internet routing – the disproportionate impact a single ISP misconfiguration error may have on global Internet routing. As we described in Section IV, most large ISPs only filter customers and do not check the validity of announcements from peers.

### A.3 Routes from ISP3

Finally, we observe a yet more complex example of delayed convergence in Figure 3(c). In this diagram, we illustrate ISP4’s exploration of six backup paths ranging between lengths of 3 and 5 after a fault for route *R3* from ISP3. After the majority of faults, ISP4 fails-over to path P2 followed by a withdraw. At a slightly lower frequency, ISP4 fails over first to P3, then P2, and finally withdraws the route. During convergence after the remaining 29 percent of faults, we measured ISP4’s exploration of an additional 145 ASPath set combinations which include an additional 14 unique ASPaths. Analysis of these remaining paths and discussions with providers indicate that approximately 65 percent represent vagabond paths, while 35 percent are “legitimate” backup paths. We next provide probable explanations for these different fail-over behaviors.

Analysis of our data shows that both the probability of selection and order of backup paths chosen during delayed convergence depends primarily on the interaction of the *MinRouteAdver* timer on routers along each path. As described in [1], the most widely deployed commercial router software today implements *MinRouteAdver* on a per peer basis. As a result, the initial *MinRouteAdver* timer value applied to the propagation of a new route failure is dependent on earlier routing instability propagated across each peering session. For example, for the majority of failures in Figure 3(a), the withdrawal for *R1* propagates along the steady-state path *before* ISP4 learns the backup P2 is also invalid. After four percent of faults, however, ISP1’s initial *MinRouteAdver* associated with ISP4 delays the withdrawal to ISP4 longer than the time required for a withdrawal to propagate along path P2. In this latter case, ISP4 will first invalidate the backup path P2. As the primary steady-state path still ex-

ists, ISP4 will not send out any update until it finally receives the ISP1 withdrawal and invalidates the primary path after an average of 32 seconds. This latter sequence of events is much less probable than an initial fail-over to P2 as it requires ISP1’s *MinRouteAdver* timer to ISP4 to be longer than the combined *MinRouteAdver* timers of ISP1 to ISP5, and ISP5 to ISP4.

As noted earlier in this Section, backup path selection also may depend on the interaction, or *MinRouteAdver interference*, of multiple paths in a topology. For example, Figure 3(c) shows that ISP5 maintains a steady-state active path through ISP3 and secondary paths through both its neighbors, ISP8 and ISP9. After a fault for *R3*, ISP5 on average will first receive a withdrawal from ISP3. In this scenario, ISP5 will then fail-over deterministically to one of its backup paths, announce the new active path and re-initialize its *MinRouteAdver* timer to each peer. We observe that this initial fail-over and reset of ISP5’s *MinRouteAdver* timers will delay the propagation of subsequent fail-over announcements. Specifically, after receipt of a withdrawal from ISP9 or ISP8, ISP5’s *MinRouteAdver* timer will suppress the announcement of fail-over to either P4 or P5 for an additional 30 seconds.

In most instances, the interaction of *MinRouteAdver* timers provides significant benefit in maintaining scalability of the distributed operation of the BGP protocol. As we observed in [1], *MinRouteAdver* timers add a level of synchronization to the system and thereby limits the number of BGP update messages and computational states. In earlier simulations, we found that the synchronization due to *MinRouteAdver* reduces computational complexity of BGP convergence from  $O(n!)$  to  $O(n)$ . Although *MinRouteAdver* reduces the computation complexity, the timers also introduce significant additional latency during delayed convergence. In the next Subsection, we explore the dominant relationships between convergence latency and network topology.

### B. Topology Impact on Convergence

In the previous Subsection, we showed that both the order and selection of backup paths explored during delayed convergence depends on a number of factors, including the initial setting of the *MinRouteAdver* timers between peers and, to a lesser degree, link and router processing delay. In Figure 5, we present a cumulative distribution graph of the convergence latencies for routes in the three topologies shown in Figure 3. The horizontal axis represents the number of seconds from injection of the fault until each ISPs’ BGP routing tables reach steady state for that prefix; the vertical axis shows the cumulative percentage of all such events. For clarity we limit the horizontal axis to 190 seconds. All ISPs exhibit a long-tailed distribution of convergence latencies extending up to fifteen minutes for a small, but tangible percentage of events. Analysis of these long-tailed events finds that most correspond to the exploration of vagabond paths.

In Figure 5, we first observe that the distribution of convergence latencies of ISP2 and ISP1 appear similar, while ISP3 exhibits significantly slower convergence times. We note that 80 percent of ISP1 and ISP2 failures converged in 100 seconds, while only 20 percent converged from ISP3 in the same time period. If we neglect the impact of vagabond paths, at a high level these convergence latencies correspond with the relative complexity of the secondary topologies explored in Figure 3.

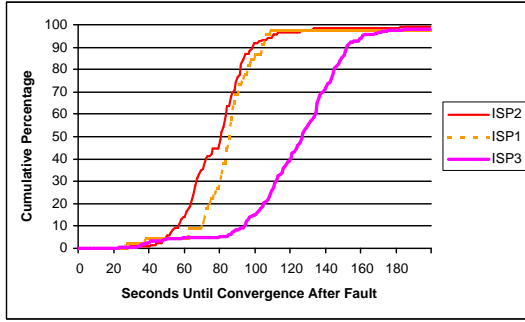


Fig. 5. Convergence latency of route from ISP1, ISP2 and ISP3 after a cumulative percentage of faults.

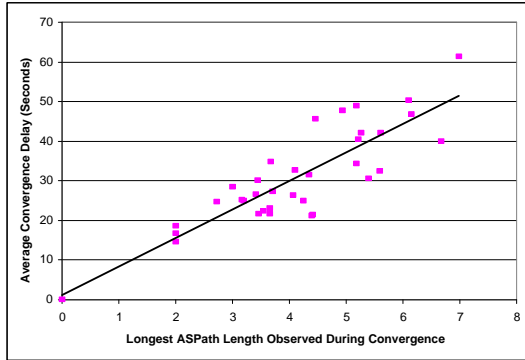


Fig. 6. Scatter plot of the average convergence latency and longest ASPath length explored during the process of delayed convergence between more than 200 pairs of Internet providers.

Specifically, we observe that in Figure 3(a) ISP4 explored a single backup path, P1, of length two. In Figure 3(b), ISP4 explored backup paths P1 and P2 of lengths 2 and 3, respectively. In contrast, ISP4 explored significantly longer and more complex topologies, including paths of length 5, during the delayed convergence of  $R3$ . This relatively higher degree of topological complexity corresponds to the longer convergence latencies shown for ISP3 in Figure 5.

Analysis of the convergence latencies for more than 200 (source, destination) inter-domain paths shows that the average convergence latency for a route failure corresponds to the length of the longest possible backup path allowed by policy and topology between two providers. We quantitatively illustrate this relationship between path length and convergence delay with a scatter plot in Figure 6. The vertical axis provides the average convergence delay for each (source, destination) pair observed during our study. The horizontal axis provides the longest, non-vagabond ASPath length we observed announced during the process of delayed convergence. We include a trend line in the graph to better illustrate data inter-relationships.

Although the data in Figure 6 contains significant variability, we observe a linear relationship between the longest ASPath length for a route between two ISPs and the average failure convergence delay for that route. We present a proof of this relationship in the next Section. A probable explanation for the variability in Figure 6 and differences between convergence latencies for ISP2 and ISP1 in Figure 5 involve differences in initial the MinRouteAdver settings due to previous routing instability,

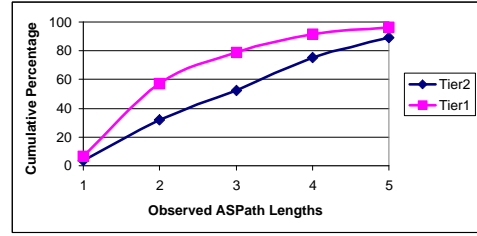


Fig. 7. Cumulative percentage distribution plot of ASPath lengths seen by customers of tier-1 and tier-2 ISPs during convergence.

intra-domain routing protocol convergence, and the processing and link delays on intermediate backbone routers.

We now explore the relationship between backup path length and Internet provider policies. Analysis of the more than 200 paths observed during our study shows a relationship between the tier size of an ISP and both the convergence delay and the length of ASPaths exhibited by routes from customers of that ISP. For example, in Figure 3 ISP1 represents one of the largest tier-1 Internet backbone providers; ISP2 represents a moderate sized US-based tier-2 provider; and ISP3 represents a regional, tier-3 network.

We illustrate this relationship between the tier of a provider and convergence behaviors in Figure 7. This graph shows the cumulative percentage of backup ASPath lengths observed during the six months of our study for routes originating from customers of two categories of providers: tier-1 and tier-2. We group providers into tier-1 and tier-2 categories based on published network topologies, published peer policies and provider self-designation. In Figure 7, we observe upstream routers on average explore shorter ASPath for routes originating in tier-1 than tier two providers. While the longest non-vagabond path explored during convergence for tier-1 customer routes is 9 ASes, tier-2 customer routes grew as long as 12 ASes. As we discussed earlier, a probable explanation for these differences in the backup paths from different tiers of providers is that smaller ISPs typically purchase transit from multiple upstream providers. These upstream transit providers may also, in turn, purchase connectivity from additional tier-1 and tier-2 providers. Smaller ISPs also implement policies unnecessary in larger providers, such as backup transit.

## VI. PROOF

In this Section, we present a formal model for the Internet, for which we prove the relationship between the longest backup inter-domain path for a route between two Internet providers and the convergence latency of that route. We first present a model of network topology and BGP communication.

### A. Model of BGP Convergence

We model the Internet as a directed graph  $G$ , with a set  $V$  of nodes and a set  $E$  of links  $(u, v)$ , with  $u, v \in V$ . The nodes  $V$  represent the set of autonomous systems in the Internet. We fix a destination  $X$  throughout this Section. The links  $E$  are connections; a link  $e = (u, v)$  exists if and only if node  $u$  will inform node  $v$  (i.e. send an update) about its best route to destination  $X$  (but not vice versa). We define the *out-neighbors*

$N^-(u) = \{v \in V | \exists e = (u, v) \in E\}$ , and similarly the *in-neighbors*  $N^+(v) = \{u \in V | \exists e = (u, v) \in E\}$ .

As discussed in Section II, we assume the BGP protocol decision process selects the best (active) route to a destination  $X$  based on ASPath length (minimum number of hops) only. For simplicity we do not consider other decision or tie-breaking criteria. A node  $u$  informs all out-neighbors about its active route to destination  $X$ . In addition to its active route to destination  $X$ , a node  $u$  stores also alternative routes; at most one per in-neighbor.

If the active route of a node  $u$  to destination  $X$  changes, node  $u$  will announce its new active route to all out-neighbors (it will send an update message with complete route information); if node  $u$  loses a connection to destination node  $X$ ,  $u$  will send a  $\text{withdraw}(X)$  message.

We are given a client node  $X$  that is not connected to the rest of the network, except a single link to  $X$ 's AS node  $A$ . We let this single link go down at time 0, thus terminating all possible routes from any node in the Internet to  $X$ . We are interested in the time it takes until the network is stable again (every node knows that there is no route to  $X$ ). We call this the  $T_{\text{down convergence time}}$ . In addition, we are also interested in the time it takes from establishment of a connection from  $X$  to  $A$  at time 0 until the BGP routing tables at all nodes are filled accordingly. We call this the  $T_{\text{up convergence time}}$ .

*Definition VI.1:* A simple path  $p$  is an ordered sequence of  $k$  nodes  $(u_1, \dots, u_k)$  such that for any pair  $i, j$  with  $i \neq j$  we have  $u_i \neq u_j$ , and there is a link  $e$  such that  $e = (u_i, u_{i+1}) \in E$ , for  $i = 1, \dots, k-1$ . The length of a path is  $|p| = k$ , the number of nodes of the path  $p$ . Note that each possible route at node  $u$  is an inverse path to node  $u$ .

The following definition will give us a comfortable handle to argue about update messages that are causal. It captures a formal definition for ‘‘the time to forward a message on a path’’.

*Definition VI.2:* We denote the *time of path*  $p$  with  $t(p)$ . If  $|p| = 1$  then  $t(p) = 0$ . Else let  $p = (u_1, \dots, u_{k-1}, u_k)$ . Let  $m$  be the first update message that was received by node  $u_k$  which was sent over link  $e = (u_{k-1}, u_k)$  and sent after time  $t(p')$  where  $p' = (u_1, \dots, u_{k-1})$ . Then  $t(p)$  is the time at which update message  $m$  is received by node  $u_k$ .

Due to the operation of the MinRouteAdver timer, a node  $u$  does not send two update messages over a link  $e$  within time MinRouteAdver.

Since the MinRouteAdver timer depends on update messages for *any* destination, it is difficult to accurately estimate how long node  $u$  will wait until MinRouteAdver allows node  $u$  to send another update message over link  $e$ . However, our measurements yield that the actually experienced waiting time is on average half the MinRouteAdver time (30 seconds). An estimate is therefore

$$t(p) \approx \text{MinRouteAdver}/2 \cdot |p| \approx 15|p| \text{ seconds} \quad (1)$$

On the other hand, the MinRouteAdver timer helps provide an upper bound on  $t(p)$ . Given that the number of in-neighbors of a node is at most the number of ASs, and given that each in-neighbor sends at most one update every 30 seconds, we can safely assume that there is never congestion of the update messages, and each update message is processed almost immedi-

ately after reception. Moreover, the time for transmission and the processing time of an update message are negligible terms compared with the MinRouteAdver time. Therefore

$$t(p) \leq (\text{MinRouteAdver} + \epsilon) \cdot |p| \approx 30|p| \text{ seconds} \quad (2)$$

## B. $T_{\text{up}}$ Convergence

We show that  $T_{\text{up}}$  convergence latency is equal to the time to forward a message on the shortest path.

We connect client node  $X$  to its AS node  $A$  at time 0.

*Theorem VI.3:* A node  $u$  learns its active route to  $X$  at time  $t(p)$ , where  $p$  is the shortest path from  $A$  to  $u$ .

*Proof:* Let  $p = (A = u_0, \dots, u_{k-1}, u_k = u)$  be the shortest path from  $A$  to  $u$ . The proof goes by induction, the induction hypothesis is that node  $u_i$  on path  $p$  knows its active route at time  $t(p_i)$ , where  $p_i$  is the prefix of the path  $p$  up to node  $u_i$ . Node  $A$  learns the shortest route to  $X$  at time 0, which corresponds to  $t(p_0) = 0$ , thus the base case of the induction is true. From the induction hypothesis we know that node  $u_{k-1}$  knows its active route to  $X$  at time  $t(p_{k-1})$ . With the Definition VI.2 we know that there is an update message  $m$ , sent from  $u_{k-1}$  to  $u_k$ , which announces this active route, sent after time  $t(p_{k-1})$ . The update message  $m$  is received by node  $u_k$  at exactly time  $t(p_k) = t(p)$ . ■

In order to optimize the  $T_{\text{up}}$  convergence time we need to both minimize MinRouteAdver delay as well as reduce the diameter (all shortest paths) of a network.

## C. $T_{\text{down}}$ Convergence Time

We disconnect  $X$  and  $A$  at time 0. In the following we present a sufficient and a necessary condition that will give us an upper and a lower bound on the time  $T_{\text{down}}$ .

*Theorem VI.4:* We are given a destination  $X$  with AS  $A$ , and a node  $u$ . We have  $T_{\text{down}} \leq \max_{p \in P} t(p)$ , where  $P$  is the set of all paths from AS node  $A$  to node  $u$ .

*Proof:* Let  $p_i = (A = u_0, u_1, \dots, u_i)$ ,  $i \geq 0$  be a simple path. Let  $u_k = u$ , thus  $p = p_k \in P$ . The proof is by induction on index  $i$ . The induction hypothesis is that after time  $t(p_i)$ , node  $u_i$  does not have a route inverse to the simple path  $p_i$ . For node  $u_0 = A$  this is true by definition, at time 0 node  $A$  knows that the connection with destination  $X$  is down. From the induction hypothesis, node  $u_i$  does not have a route inverse to the simple path  $p_i$  after time  $t(p_i)$ . Directly from definition VI.2 we know that node  $u_i$  is sending an update message  $m$  to node  $u_{i+1}$ , which is received at time  $t(p_{i+1})$ . Message  $m$  announces either a route that is different from  $p_i$  or withdraws all routes. Therefore, after time  $t(p_{i+1})$  node  $u_{i+1}$  fulfills the induction hypothesis. And thus at time  $t(p)$ , node  $u$  does not have a route inverse to path  $p$ . From the definition of BGP we know that the only routes node  $u$  can possibly consider are the simple paths from  $A$  to  $u$ , the set  $P$ . The Theorem follows directly. ■

With Equation (2) we get  $T_{\text{down}} \leq (\text{MinRouteAdver} + \epsilon)|p|$  seconds, where  $\epsilon \ll \text{MinRouteAdver}$  is a small constant. Since there is at most one update message on each link every 30 seconds, we can also upper bound the number of update messages on each link by  $(\text{MinRouteAdver} + \epsilon)|p|/\text{MinRouteAdver} \approx |p|$ . We get immediately:

*Corollary VI.5:* The number of update messages caused by a single failure is upper bounded by  $|p| \cdot |E|$ , where  $|p|$  is the size of the longest path, and  $|E|$  is the number of links. Both  $|p|$  and  $|E|$  can be upper bounded, resulting that the number of update messages is less than  $|V|^3$ , where  $|V|$  is the number of nodes (ASs).

*Definition VI.6:* Let  $p_i = (A = u_0, u_1, \dots, u_i)$  be a simple path, with  $i = 0, \dots, k$ . The simple path  $p_k$  is called *vital* if and only if at each node  $u_i$  ( $i = 0, \dots, k$ ), node  $u_i$ 's active route was the inverse of path  $p_i$  right before time  $t(p_i)$ .

*Theorem VI.7:* We are given a destination  $X$  with AS  $A$ , and a node  $u$ . We have  $T_{down} \geq t(p)$ , where  $p$  is any *vital* path from  $A$  to  $u$ .

*Proof:* Let  $p$  be a *vital* path from node  $A$  to node  $u$ . With Definition VI.6 we know that node  $u$  has the inverse of path  $p$  as its active route right before time  $t(p)$ . Therefore, right before time  $t(p)$  node  $u$  still has a route to the destination  $X$ , which lets  $T_{down}$  to be at least  $t(p)$ . ■

In general the upper bound given in Theorem VI.4 and the lower bound given in Theorem VI.7 can differ substantially. The lower (upper) bound is in the order of the time to forward a message on a shortest (longest) path.

Theorem VI.7 enables us to raise the lower bound by observing that paths are *vital* in general. The routing table of node  $u$  (at any time) consists of at most one entry per in-neighbor node. The entry with the shortest route is the currently active. In the proof of Theorem VI.4 we have seen that invalidating a route needs forwarding a message along the path of the route. Since forwarding a message on a shorter path usually takes less time than forwarding a message on a longer path, the update message invalidating the active route will often arrive first. In the following Corollary we take this observation to the extreme.

*Corollary VI.8:* If  $t(p') < t(p'')$  for  $|p'| < |p''|$ , and ties are broken with advantage of the currently *vital* route, then  $T_{down} = t(p)$ , where  $p$  is the longest path from node  $A$  to node  $u$ .

Alternatively to the strong assumption of Corollary VI.8 we could also presume that the graph  $G$  is acyclic; an assumption support by recent studies [16].

*Corollary VI.9:* Let the graph  $G$  be acyclic, and let the time to transfer a message (including MinRouteAdver) be at least one time unit. Then  $T_{down} \geq |p|$  time units, where  $p$  is the longest path from node  $A$  to node  $u$ .

*Proof:* Let  $p_0 = p$  be the longest path from node  $A$  to node  $u$ . Let  $u_i$  be the first node where the message of path  $p_i$  was received and did not invalidate the active route  $r$  of node  $u_i$ . Let  $p_{i+1}$  be the path whose prefix is the path associated with the active route  $r$  of node  $u_i$ , and whose postfix is the postfix of path  $p_i$ . Let finally  $p_k$  be the path that is received by node  $u$  without being interrupted by another path. Then  $p_k$  is *vital*, and  $t(p_k) \geq t(p) \geq |p|$ . ■

With Equation (1), a precondition of either corollary is sufficient to show that  $T_{down} \geq 15|p|$  seconds.

Computing the longest simple path for a given graph is known to be NP-complete. Even worse, there is provably no good heuristic to approximate the longest simple path unless  $P = NP$  (the longest path problem is in APX [17]).

In order to optimize  $T_{down}$  convergence time, we need to minimize the longest paths in a network. However, we encounter a

trade-off between the length of possible paths and the degree of connectivity in a network. Specifically, in a highly connected network, the longest path is linear in the number of nodes of the network. We will elaborate on this trade-off in future work.

## VII. CONCLUSIONS AND FUTURE WORK

In this paper and previous work, we have shown that the Internet currently lacks the level of reliability and fault-tolerance required for the successful deployment of many emerging mission-critical network services. We argue that current Internet convergence latencies of up to fifteen minutes after a single multi-homed fault will prove untenable for new interactive and transaction-oriented network applications like Internet telephony.

This paper demonstrated that the time complexity for multi-homed Internet path fail-over scales linearly with the length of the longest possible backup path for that route. We showed that the length of inter-domain backup paths depends on a number of inter-provider contractual and policy implementation details. Our results show that customers sensitive to fail-over latency should multi-home to larger providers, and that smaller providers should limit their number of transit and backup transit interconnections. Finally, the large number of erroneous vagabond paths we observed during our study suggests a significant need for better route validation and authentication mechanisms.

In ongoing work, we are exploring possible solutions to the problems of delayed convergence and path authentication. Solutions under exploration include the use of adaptive MinRouteAdver timers and the association of additional information with BGP withdrawal messages. Our hope is to identify mechanisms which ameliorate the delayed convergence problem while preserving the scalability and flexibility of the Internet routing protocols.

## REFERENCES

- [1] C. Labovitz, A. Ahuja, A. Bose, and F. Jahanian, "Delayed Internet Routing Convergence," *Proc. of the ACM SIGCOMM*, Aug. 2000. A version also to appear in *IEEE/ACM Transactions on Networking*.
- [2] R. Perlman, *Interconnections, Second Edition*, Addison-Wesley, Reading Massachusetts, 1999.
- [3] J. Garcia-Luna-Aceves, "Loop-free Routing Using Diffusing Computations," *IEEE/ACM Transactions on Networking*, Feb. 1993.
- [4] K. Bhargavan, D. Obradovic, and C. Gunter, "Formal Verification of Distance Vector Routing Protocols," *International Conference on Theorem Proving and Higher Order Logics*, Aug. 2000.
- [5] Y. Rekhter and T. Li, "A Border Gateway Protocol 4 (BGP4)," Sept. 1999, draft-ietf-idr-bgp4-09.txt.
- [6] B. Halabi, *Internet Routing Architecture*, Cisco Press, 1997.
- [7] Kannan Varadhan, Ramesh Govindan, and Deborah Estrin, "Persistent Route Oscillations in Inter-Domain Routing," Tech. Rep. USC CS TR 96-631, Department of Computer Science, University of Southern California, Feb. 1996.
- [8] T. Griffin and G. Wilfong, "An Analysis of BGP Convergence Properties," *Proceedings of the ACM SIGCOMM*, Aug. 1999.
- [9] Lixin Gao and J. Rexford, "Stable Internet Routing Without Global Coordination," *Proc. of ACM SIGMETRICS*, June 2000.
- [10] Bill Cheswick and Hal Burch, "Internet Mapping Project," <http://www.cs.bell-labs.com/who/ches/map/>.
- [11] Ramesh Govindan and Hongsuda Tangmunarunkit, "Heuristics for Internet Map Discovery," *Proc. of IEEE INFOCOM*, 2000.
- [12] "Internet Performance Measurement and Analysis Project (IPMA)," <http://www.merit.edu/ipma>.
- [13] "Multithreaded Routing Toolkit (MRT) Project," <http://www.mrtd.net>.

- [14] C. Labovitz, A. Ahuja, and F. Jahanian, "Experimental Study of Internet Stability and Wide-area Network Failures," *Proc. International Symposium on Fault-Tolerant Computing*, June 1999.
- [15] R. Barrett, S. Haar, and R. Whitestone, "Routing snafu causes internet outage," *Interactive Week*, 1997.
- [16] Lixin Gao, "On Inferring Autonomous System Relationships in the Internet," *IEEE Global Internet*, Nov. 2000.
- [17] David R. Karger, Rajeev Motwani, and G. D. S. Ramkumar, "On approximating the longest path in a graph," *Algorithmica*, vol. 18, no. 1, pp. 82–98, May 1997.