

An Experimental Study of Internet Routing Convergence

Craig Labovitz
Microsoft Research
One Microsoft Way
Redmond, Washington 98052-6399
labovit@microsoft.com

Abha Ahuja, Abhijit Bose, Farnam Jahanian
University of Michigan
Department of Electrical Engineering and Computer Science
1301 Beal Ave.
Ann Arbor, Michigan 48109-2122
{ahuja, abose, farnam}@umich.edu

Technical Report
MSR-TR-2000-08

Microsoft Research
Microsoft Corporation
One Microsoft Way
Redmond, WA 98052

Abstract

This paper examines the latency in Internet path failure, failover and repair due to the convergence properties of inter-domain routing. Unlike switches in the public telephony network which exhibit failover on the order of milliseconds, we show 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. These delays stem from temporary routing table oscillations formed during operation of the BGP path selection process on Internet backbone routers. During these periods of *delayed convergence*, end-to-end Internet paths will experience intermittent loss of connectivity, as well as increased packet loss and latency. We present a two-year study of Internet routing convergence through the experimental instrumentation of key portions of the Internet infrastructure, including both passive data collection and fault-injection machines at major Internet exchange points. Based on data from the injection and measurement of several hundred thousand inter-domain routing faults, we describe several unexpected properties of convergence and show that the measured upper bound on Internet inter-domain routing convergence delay is an order of magnitude slower than previously thought. Our analysis also shows that the upper computational bound on the number of router states and control messages exchanged during the process of BGP convergence is exponential with respect to the number of autonomous systems in the Internet. Finally, we demonstrate that much of the observed convergence delay stems from both specific router vendor implementation decisions, as well as ambiguity in the BGP specification.

1 Introduction

In a brief number of years, the Internet has evolved from an experimental research and academic network to a commodity, mission-critical component of the public telecommunication infrastructure. During this period, we have witnessed an explosive growth in the size and topological complexity of the Internet and an increasing strain on its underlying infrastructure. As the national and economic infrastructure has become increasingly dependent on the global Internet, the end-to-end availability and reliability of data networks promises to have significant ramifications for an ever-expanding range of applications. For example, transient disruptions in backbone networks that previously impacted a handful of scientists may now cause enormous financial loss and disrupt hundreds of thousands of end users.

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). Factors contributing to these differences between the commercial Internet infrastructure and the PSTN have been discussed in various literature [21, 15]. 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 infrastructure.

The Internet backbone infrastructure is widely believed to support rapid restoration and rerouting in the event of individual link or router failures. At least one report places the latency of inter-domain Internet path failover on the order of 30 seconds or less based on qualitative end user experience [13]. These brief delays in inter-domain path failover are further believed to stem mainly from queuing and router CPU processing latencies [1, (message digests 11/98, 1/99)]. In this paper, we show that most of this conventional wisdom about Internet failover is incorrect. Specifically, we demonstrate that the Internet does not support effective inter-domain failover and that most of the delay in path restoration stems solely from delayed BGP convergence.

The slow convergence of distance vector (DV) routing is not a new problem [19]. DV 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.

DV routing can take a long time to converge after a topological change because routers do not have sufficient information to determine if their choice of next hop will cause routing loops to form. The count-to-infinity problem [19] is the canonical example used to illustrate the slow convergence in DV routing. Numerous solutions have been proposed to address this issue. For example, including the entire path to the destination, known as the *path vector* approach, is used in the Border Gateway Protocol (BGP), the inter-domain routing protocol in the Internet. Other attempts to solve the count-to-infinity problem or accelerate convergence in many common cases include techniques such as split horizon (with poison reverse), triggered updates, and the diffusing update algorithm [7].

Although the theoretical aspects of the delayed convergence problems associated with DV protocols are well known, this paper is the first to our knowledge to investigate and quantitatively measure the convergence behavior of BGP4 deployed in today's Internet. In [4], the authors showed that in the worst case, the original Bellman-Ford distance vector algorithm requires $O(N^3)$ iterations to find the shortest path lengths for a network with N nodes. However, we are not aware of any published result of a similar bound for path vector algorithms. The adoption of the path vector is widely and incorrectly believed to provide BGP with significantly improved convergence properties over traditional DV protocols, including RIP [12].

A number of recent studies, including Varadhan et al. [22] and Griffin and Wilfong [8] have explored BGP routing *divergence*. As we describe in the next Section, 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. Griffin et al. in [9] and Rexford et al. in [6] also describe modifications to BGP which guarantee that the protocol will not diverge. The authors of all these papers note that BGP divergence remains a theoretical finding and has not been observed in practice. Our work explores a complimentary facet of BGP routing – the convergence behavior of safe, or satisfiable routing policies. We show that even under constrained policies, the complexity of BGP convergence is exponential with respect to the number of autonomous systems.

Bhargavan et al. in [4], provide a stricter upper bound on the convergence of RIP. The authors account for implementation details of RIP, including poison reverse, triggered updates, and split-horizon, which provide for improved convergence behavior over previous analyses of Bellman-Ford algorithms. In [23], the authors simulated the convergence behaviors of several algorithms including a distributed Bellman-Ford and present metrics for comparing the convergence properties of these different protocols. In this work, we similarly focus on both measuring the convergence latencies of BGP and developing theoretical upper and lower bounds.

In [16], Labovitz et al. describe significant levels of measured Internet routing instability. The authors show that most Internet routing instability in 1997 was pathological and stemmed from software bugs and

artifacts of router vendor implementation decisions. In a later paper [17], Labovitz and his co-authors show that once ISPs deployed updated router software suggested by [16], the level of Internet routing instability dropped by several orders of magnitude. Finally, in [15] Labovitz et al. measured the rate of network failure, repair and availability. In this work, we present a complimentary study on both the impact and the rate at which inter-domain repair and failure information propagates through the Internet. We also measure the impact of Internet path changes on end-to-end network performance. Specifically, our major results include:

- Although the adoption of the path vector by BGP eliminates the DV count-to-infinity problem, the path vector exponentially exacerbates the number of possible routing table oscillations.
- The delay in Internet inter-domain path failovers now averages three minutes, and some percentage of failovers trigger routing table oscillations lasting up to fifteen minutes.
- The theoretical upper bound on the number of computational states explored during BGP convergence is $O(e(n-1)!)$, where n is the number of autonomous systems in the Internet. We note that this is a theoretical upper bound on BGP convergence and is unlikely to occur in practice.
- If we assume bounded delay on BGP message propagation, then the lower bound on BGP convergence is $\Omega((n-3) * 30)$ seconds and $O(n)$ states, where n is the number of autonomous systems in the Internet.
- The delay of inter-domain route convergence is due almost entirely to the unforeseen interaction of protocol timers with specific router vendor implementation decisions.
- Internet path failover has significant deleterious impact on end-to-end performance – measured packet loss grows by a factor of 30 and latency by a factor of four during path restoration.
- Minor changes to current vendor BGP implementations would, if deployed, reduce the lower bound on inter-domain convergence time complexity from $\Omega((n-3) * 30)$ to $\Omega(30)$ seconds, where n is the number of autonomous systems in the Internet.

The remainder of this paper is organized as follows: Section 2 provides additional background on BGP. Section 3 provides a description of our experimental measurement infrastructure. In Section 4, we present the results of our two year study of Internet routing convergence. We describe the measured convergence latencies of both individual ISPs and the Internet as a whole after several categories of injected routing faults. In Section 5, we provide analysis of our experimental results and a discussion of the theoretical convergence properties of BGP. We provide both a lower and upper bound on BGP convergence. Finally, we conclude with a discussion of specific modifications to vendor BGP implementations which, if deployed, would significantly improve Internet convergence latencies.

2 Background

Autonomous systems (ASes) in the Internet today exchange inter-domain routing information through BGP. We assume that the reader is familiar with Internet architecture and the BGP routing concepts discussed in [20, 10]. We provide a brief review of the more salient attributes of BGP related to the discussion in this paper.

Unlike interior gateway protocols, which periodically flood an intra-domain network with all known topological information, BGP is an incremental protocol that sends update information only upon changes in network topology or routing policy. Routing information shared among BGP speaking peers has two forms - announcements and withdrawals. A route announcement indicates that a router has either learned of a new network attachment or has made a policy decision to prefer another route to a network destination. Route withdrawals are sent when a router makes a new local decision that a network is no longer reachable. Explicit withdrawals are those associated with a withdrawal message. Implicit withdrawals occur when an existing route is replaced by an announcement of a new, more preferred route without an intervening withdrawal message. We define route *failover* as the implicit withdrawal and replacement of a route with one having a different ASPath. For purposes of our discussion, we define a *steady-state network* as one where no BGP peers send updates for a given prefix for 30 minutes or more.

BGP limits the distribution of a router's reachability information to its peer, or neighbor routers. 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. The default BGP behavior 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 the 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 protocol. While BGP allows for path selection based on policy attributes, including local preference and multi-exit discriminator values; a review of BGP logs, discussions with Internet network operators, and a survey of policies registered in the Internet Routing Registry (IRR) indicates that the majority of ISP policies 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 the minimum route advertisement interval timer, abbreviated in this paper as MinRouteAdver, which specifies a minimum amount of time that must elapse between advertisement of routes to 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. In order to achieve a minimum of MinRouteAdver between announcements, this rate-limiter is applied as a jittered interval on a (prefix destination, peer) tuple basis. The BGP standard specifies that MinRouteAdver applies only to BGP announcements and not explicit withdrawals. This distinction stems from the goal of avoiding the long-lived "black holing" of traffic to unreachable destinations. Due to the delay introduced by MinRouteAdver on announcements throughout the Internet, BGP withdrawals are commonly believed to propagate and converge more quickly.

3 Methodology

We base our analysis on data collected from the experimental instrumentation of key portions of the Internet infrastructure. Over the course of two years, we injected over 150,000 routing faults into geographically and topologically diverse peering sessions with five major commercial Internet service providers. We then measured the impact of these faults through both end-to-end measurements and logging ISP backbone routing table changes.

Figure 1 shows a simplified diagram of our RouteViews measurement and fault injection infrastructure. We measured the impact of injected faults via both active and passive probe machines deployed at major US exchange points, as well as on the *(omitted for blind review)* campus. Our passive instrumentation included several “RouteViews” probe machines, which maintained default-free peering with over 25 Internet providers. These RouteViews machines time-stamped and logged all BGP updates received from peers to disk.

We injected faults consisting of BGP route transitions for both /19 and /24 prefix-length addresses. Although we injected faults from a number of diverse probe locations, we simplify the discussion in this paper by presenting data only from faults injected at the Mae-West exchange point and from the *(omitted for blind review)* campus. We note that data from other probe locations exhibited similar behaviors. 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. We generated faults over a two year 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 [15] found occur on the average of once a month.

Software from the *(omitted for blind review)* and *(omitted for blind review)* projects [2, 3] 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. This interval was chosen to avoid the effects of BGP route flap dampening [10] between fault injections. The faults simulated route failures, repairs and multi-homed failover. In the case of failover, we announced both a primary route for a given prefix with a short ASPath to one upstream BGP neighbor, and a longer ASPath route for the same prefix to a second provider. The announcement of two routes of different ASPath length represents a common method of customer multihoming to two Internet providers. In an effort to ensure that the downstream peers would always prefer the primary route if it existed, we prepended the long ASPath route announcement with three times the average number of AS numbers observed in steady-state path lengths. We then periodically failed the shorter ASPath route while maintaining the longer backup path.

While the RouteViews probes monitored the impact of BGP faults on core Internet routers, our active measurements monitored the impact on end-to-end performance. We configured these probe machines with a virtual interface addressed within the prefix blocks included in the injected BGP faults. These probe machines sent 512 byte MTU ICMP echo messages to 100 randomly selected web sites once a second. We

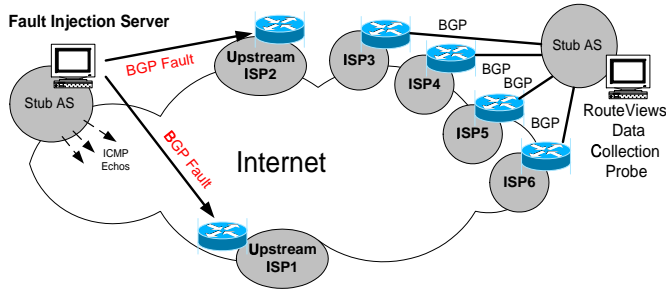


Figure 1: Diagram of the fault injection and measurement infrastructure.

randomly selected the web sites from a major Internet cache log of several hundred thousand entries.

We then correlated the data between our NTP synchronized fault injection probe machines and both our RouteViews and end-to-end measurement logs. These correlations provided data on the number of update messages generated for a particular route announcement and withdrawal, as well as the convergence delay for a particular ISP and all ISPs to reach steady state after a fault.

We also simulated routing convergence using software from the *(omitted for blind review)* project [3]. The *(omitted for blind review)* daemon supports the configuration of multiple BGP autonomous systems and associated routing tables within a single workstation process. As a complete routing protocol implementation, the software supports the generation of BGP update packets and the application of arbitrary BGP policies similar to those available on commercial routers. In simulation mode, the daemon exchanges packets internally and does not forward updates to the network. By programmatically introducing delay in message propagation and processing, we were able to simulate both the average and upper bound on BGP convergence for networks of varying degree and topology.

4 Experimental Results

In this section we present data collected with the experimental measurement infrastructure described in the previous section. We first provide a taxonomy for describing the four categories of routing events injected into the Internet during our study:

Tup A previously unavailable route is announced as available. This represents a route repair.

Tdown A previously available route is withdrawn. This represents a route failure.

Tshort An active route with a long ASPath is implicitly replaced with a new route possessing a shorter ASPath. This represents both a route repair and failover.

Tlong An active route with a short ASPath is implicitly replaced with a new route possessing a longer ASPath. This represents both a route failure and failover.

We define the *latency* of each injected event as the time between the injection of the fault and the routing tables of a given ISP, or all ISPs, we monitored to reach steady state for the injected prefix. In the following two subsections, we present data from our both our passive routing and active end-to-end measurements.

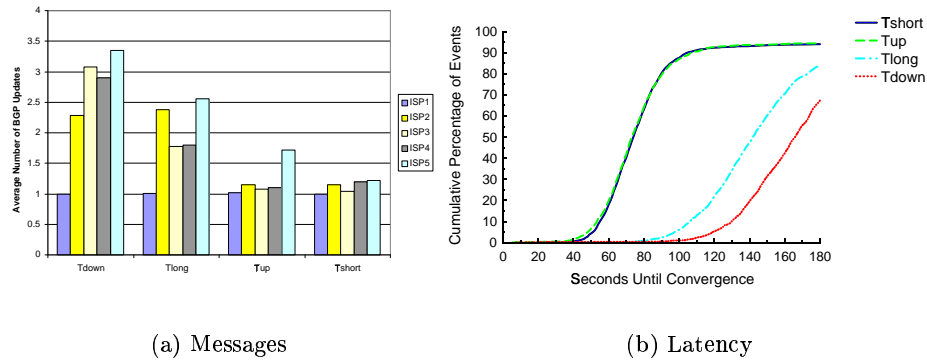


Figure 2: Average number of BGP updates from 5 ISPs triggered by Tdown, Tlong, Tup, Tshort events and convergence latency of cumulative percentage of Tup, Tshort, Tlong and Tdown events for all monitored ISPs over course of our two year study. Both data sets are for faults injected from (*omitted for blind review*).

4.1 Routing Measurements

We first look at the volume or number of BGP routing updates triggered by each injection of a routing event. We observe that the injection of a single routing event may trigger the generation of multiple route announcements and withdrawals from each ISP. In Figure 2(a), we show the average number of update messages generated by five ISPs for each category of routing event over the two year course of our study. Although we monitored the BGP routing tables of 25 ISPs, we graph only five ISPs in Figure 2(a) for clarity. We note that data from the other monitored providers exhibited similar behaviors.

The most salient observation we make from Figure 2(a) is that both Tdown and Tlong events on average triggered more than two times the number of update messages than both Tup and Tshort events. As we explain later in this Section, (Tlong, Tdown) and (Tup, Tshort) form equivalence classes with respect to both convergence latency and the number of update messages they trigger. We also note significant variation in the average number of updates generated by individual ISPs within each equivalence class. For example, we see that for ISP3, Tdown triggered twice the number of messages as Tlong. In contrast, Tlong events triggered more messages in ISP2 than Tdown. In all categories, ISP1 generated an average of only one BGP update. We provide probable explanations for these behaviors later in the next section.

In Figure 2(b), we explore the differences in latency among the four categories of routing events. Figure 2(b) shows the convergence latency for a cumulative percentage of Tdown, Tup, Tshort and Tlong events over all monitored ISPs. The horizontal axis represents the number of seconds from injection of the fault until all 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 180 seconds. All four events exhibited a heavy-tailed distribution of convergence latencies extending up to fifteen minutes for a small, but tangible percentage of events. Significantly, Figure 2(b) shows more than twenty percent of Tlong and forty percent of Tdown events oscillated for more than three minutes. We note that these observed latencies

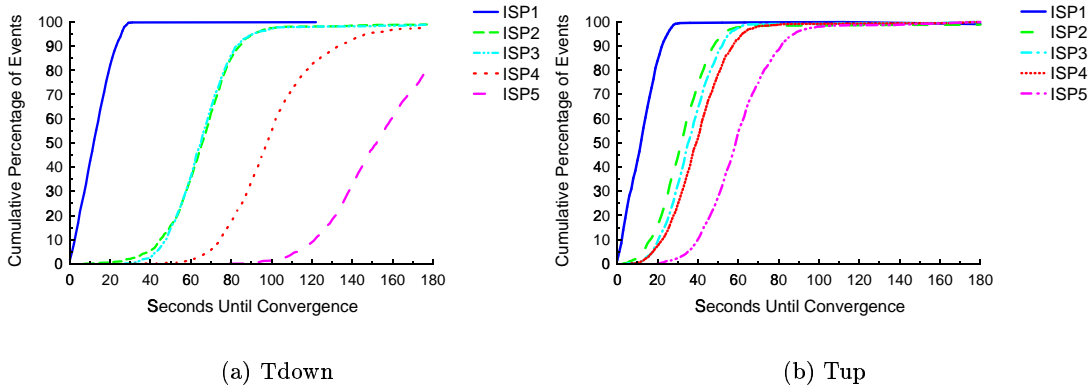


Figure 3: Convergence latency of a cumulative percentage of Tdown and Tup events injected at (*omitted for blind review*) for five major ISPs

are an order of magnitude longer than those reported in [1, 13].

As mentioned earlier, Figure 2 shows that (Tlong, Tdown) and (Tshort, Tup) form equivalence classes based both on the number of updates they trigger and their similar distribution of convergence latencies. Both Tdown and Tlong converged more slowly than Tup or Tshort: Tup and Tshort events converged within 90 seconds while only five percent of Tdown and Tlong events converged within 90 seconds, and twenty percent of Tdown/Tlong required longer than two minutes to converge. We note strong correlations between the relative number of update messages generated per equivalence class in Figure 2(a) and the convergence latencies of each category.

We now look at the latency for two categories of injected events on a per ISP basis. Figure 3 shows the convergence latency of a cumulative percentage of both Tdown (a) and Tup (b) events for five ISPs. The horizontal axis represents the delay in one second bins between the time of event injection and the BGP routing tables in each ISP reached steady state for that prefix. The vertical axis shows the cumulative percentage of all such events. As before, we present data from only five ISPs and limit the horizontal axis to 180 seconds for clarity of presentation.

We observe significant variation in the convergence latencies of the five ISPs in both graphs of Figure 3. The variations appear most pronounced in Figure 3(a) where a three-minute gap separates 80% of ISP1 events from ISP5. In our analysis, we looked for correlations between the convergence latencies of an ISP and both the geographic and network distance of that ISP. We define *network distance* as the number of traceroute hops from the point of fault injection to the peer border router interface of each ISP. In Figure 2 and Figure 3, ISP1 represents a special case – the only ISP into which we both injected events and monitored the convergence latencies. As such, the routing table of ISP1 did not exhibit EBGP route oscillations. As we explain in Section 5, at all times ISP1 either had the shortest ASPath route, or ignored updates from neighbor ISPs after detection of an ASPath loop.

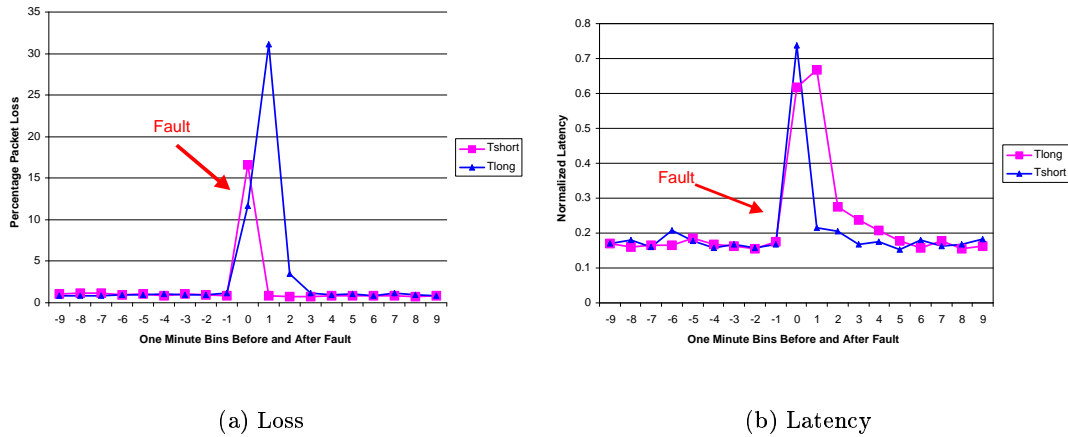


Figure 4: Average percentage end-to-end loss and normalized latency of 512 byte MTU ICMP echos sent to 100 web sites every second during the ten minutes immediately preceding and following the injection of Tshort and Tlong events at the Mae-West exchange point.

With the exception of ISP1, our data shows no correlation between convergence latency and geographic or network distance. For example ISP3, which is a national backbone in Japan, converged more quickly for both Tdown and Tup than a Canadian provider, ISP5. We show in Section 5 that convergence latencies are likely dependent on topological factors including the number of adjacent BGP peers and transit policies.

We also looked for temporal correlations between convergence delay and the time of day or week. In [16], Labovitz et al. describe a direct relationship between the hourly rate of routing instability and the diurnal bell curve exhibited by Internet bandwidth consumption and the corresponding load on backbone routers. Our analysis, however, found no such relationship between time of day/week and failover latency.

4.2 End-to-End Measurements

We now turn our attention from the convergence latencies of backbone routing tables to the impact of delayed convergence on end-to-end network paths. We show that even moderate levels of routing table oscillation will lead to increased packet loss, latency and out of order packets as routers drop packets for which they do not have a valid next hop, or queue packets while awaiting the completion of forwarding table cache updates. We expect end-to-end active measurements to provide a better measure of the application-level impact of routing convergence as not all routing table changes effect the forwarding path, and EBGP routing table measurements do not capture delays introduced by the convergence of smaller ISPs or IGP communication. In Figure 4(a), we show packet loss averaged over one minute intervals between our fault injection machine and 100 randomly selected web sites. The horizontal axis shows one-minute bins for the ten minutes both preceding and immediately following the injection of both a Tlong and Tshort failover event. Time 0 is the point of fault injection. The vertical axis represents the percentage loss for each one-minute bin

averaged both over all web sites and each corresponding bin in every ten-minute fault injection periods. We see in Figure 4(a) less than one percent packet loss throughout the ten-minute period before each fault. Immediately following the fault, the graphs for Tlong and Tshort events show a sharp rise to 17 and 32 percent loss, respectively, followed by sharply declining loss over the next three minutes. The wider curve of Tlong with respect to Tshort corresponds to the relative speeds of routing table convergence for both events shown in Figure 2. Specifically, Tlong exhibits a two minute period where loss exceeds twenty percent and Tshort a one minute period of greater than fifteen percent loss. These loss trends support the data in Figure 2, where eighty percent of Tlong and Tshort events converged within the same respective periods.

We also examine the impact of convergence on end-to-end path latency. Figure 4(b) shows the average normalized round-trip latency of ICMP echos in the ten-minute bins before and after a Tlong and Tshort event. Time 0 represents the instant of fault injection. We normalize the latency of echos on a per destination basis by dividing the latency of each echo by the average delay to that destination. As with the analysis of packet loss, we see that route failover has significant impact on end-to-end latencies. For both Tlong and Tshort, latencies more than tripled in the three minutes immediately following both categories of failover. Although Tshort exhibited an initially higher increase in latency, the curve for Tlong appears broader, extending for five minutes after the event. We note that the variation in end-to-end latency between Tup and Tdown corresponds with routing table convergence data presented in Figure 3.

Finally, we analyzed the end-to-end speed of repair, or Tup, by measuring the rate at which ICMP echos first began consistently returning from each web site after a repair. Although we omit the graph of Tup end-to-end behavior for brevity, we note that the majority (over 80%) of web sites began returning ICMP echoes within 30 seconds, and all web sites returned echos within one minute. These results also correspond with the routing convergence latencies reported in Figure 3.

We also note that our end-to-end and routing table measurements correspond to observations by other researchers. Delayed convergence provides a likely explanation for both the temporary routing table oscillations observed by Paxson in [18] as well as some of the instabilities observed by Labovitz et al. in [17].

5 Analysis

In this Section, we provide analysis and probable explanations for BGP convergence behaviors. We base our discussion on analysis of the BGP specification [20], simulation results, and the previously described experimental measurements. We begin by defining metrics to quantify the complexity and impact of BGP convergence. We choose our metrics to reflect router computational load, failover delay, and end-to-end path performance. We will refer to the following metrics throughout the remainder of our discussion:

Communication Complexity The overall number of update messages generated by all routers in the system. This includes both announcement and withdrawal messages.

State Complexity The number of distinct routing table states, or active routes explored by each node in

the system.

Time Complexity The delay of system in reaching convergence.

We simplify our analysis by choosing a full mesh, or complete graph of autonomous systems, as our model of the Internet. The choice of a complete graph over arbitrary topologies has significant impact on our analysis. We show in Section 5.2 that a complete graph provides the worst-case complexity of BGP convergence and, as such, overestimates the average case. We also note that the choice of a full mesh reflects current trends in the evolution of the Internet towards less hierarchy and a more meshed topology [14, 13]. We further simplify our analysis by excluding the delay and additional states introduced by IBGP communication and IGP convergence.

Since BGP does not place bounds on the delay of update propagation or processing, discussions of time complexity are only constructive if we assume bounded delays. We initially exclude the impact of Min-RouteAdver and associated timers on convergence. We will discuss time complexity and the impact of these timers in SubSection 5.3. Given the lack of bounds on message propagation, we assume messages may arrive in non-deterministic order subject only to the constraint that FIFO ordering is preserved between any pair of autonomous system peers.

We begin by observing that the adoption of the path vector in BGP provably provides a solution to the counting-to-infinity problems for constrained policies [9]. Intuitively, the adoption of the shortest path as a decision metric results in a monotonically increasing path selection algorithm, placing an upper bound on the number of computational steps. The length of the BGP ASPath chosen by a node is bounded by the length of a path that traverses all of its neighbors, or $(n - 1)$ in a complete graph of n autonomous systems. Further, the loop detection mechanism in BGP resolves the RIP routing table looping problem where a given node reuses information in new path that the node itself originally initiated. But, as we show below, the addition of the path vector exacerbates the bouncing problem [5].

5.1 BGP Convergence with Unbounded Propagation Delay

The most significant difference between the convergence behavior of traditional DV protocols and BGP is that DVs are strictly increasing, whereas BGP is monotonically increasing. Unlike RIP, BGP has $k!$ permutations of representing a path metric of length k . We show that in the worst case, long link, queuing or processing delays can result in an ordering of messages such that BGP will explore all possible paths of all possible lengths. We note that such an ordering represents the upper bound on BGP convergence and is unlikely to occur in practice.

In Figure 5, we show an example of BGP convergence involving a complete graph of a three node system, where all nodes are initially directly connected to route, R . The third column shows the routing table of each autonomous system at each computational stage. For each AS, we provide the matrix of current paths through each of it's neighbors. We denote the active route with an asterisk and a withdrawn, or invalid path with a dash. The fourth column provides the messages processed at each stage. The last column shows the

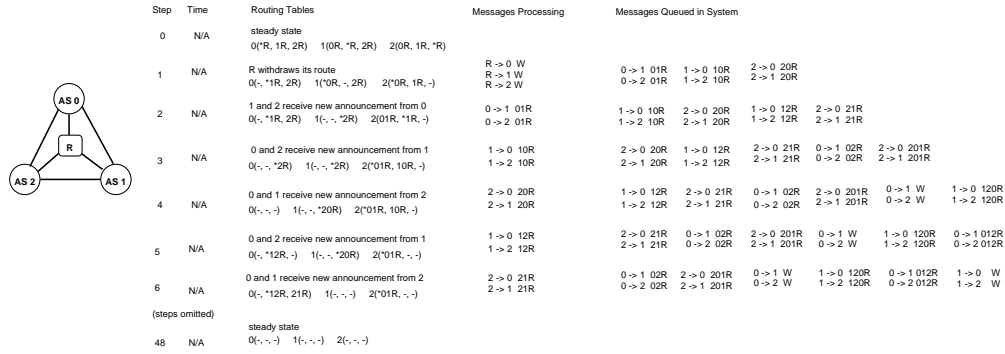


Figure 5: Example of BGP bouncing problem.

global queue of outstanding messages in the system. We process messages in serial fashion from this queue. In our example, we see at step 0 that *AS1* has one primary route (directly connected) and two backup paths (via *AS0* and *AS2*) to *R*.

As the full example includes over forty steps, we present only the first six steps and the last step in Figure 5 for clarity. The main goal of the example is to illustrate the exploration of ever increasing ASPath lengths and the generation of large numbers of update messages during convergence. At time 1, Route *R* is withdrawn following a fault. All three ASes then invalidate their directly connected paths of length 1, and choose secondary paths: *AS0* selects *1R*, *AS1* selects *0R* and *AS2* select *0R*. The three ASes also announce these new active routes to each of their neighbors. In the next steps, *AS0* detects a “looped” path from *AS1* and *AS2*, and invalidates both of these routes. Lacking a valid route to *R*, *AS0* then sends out withdrawal messages to both neighbors. Upon receipt of this withdraw, *AS1* and *AS2* failover to secondary routes (*AS1* via *20R*, and *AS2* via *10R*). In the final stages of the example, *AS1* and *AS2* detect the mutual dependency through each other via the exchange of looped BGP ASPaths. Finally, at stage 48 the system converges with all routes withdrawn.

5.2 Upper Bound on Convergence

In this Section, we provide a proof of the upper bound on convergence for a network of n BGP autonomous systems. As discussed in Section 1, we initially assume unbounded delay on message propagation. We begin with several observations:

Observation 1: For a complete graph of n nodes, there exist $O(e(n-1)!)$ distinct paths to reach a particular destination.

To show this, we note that there exists a total of $(n-1)$ paths of length 1 to reach a particular destination in a complete graph. Any other path of length greater than 1 must use these $(n-1)$ paths as the last hop in order to reach that destination. For example, there are exactly $(n-1) * (n-2)$ paths of length 2 in a complete graph. Therefore, the sum of all paths can be written as a series sum :

$$P[n] = (n-1) + (n-1)(n-2) + \dots + (n-1)!$$

The above expression can be rewritten as :

$$P[n] = (n-1)! [1 + 1/2! + 1/3! + \dots + 1/(n-2)!]$$

which is closely approximated by $p[n] = O(e*(n-1)!)$. This is an upper bound on the number of all possible paths to any destination in a complete graph of size n .

Observation 2: When a particular route is withdrawn, a path vector algorithm attempts to find an alternate path by iterating on the available paths of increasing length. We refer to this as a k -level iteration of the algorithm. At the k -th iteration, the algorithm looks at paths spanning at most k edges of the graph.

Observation 3: The conditions necessary for the worst case convergence are :

- i A complete graph, i.e. all nodes have a degree of $(n-1)$.
- ii All messages (both announcements and withdrawals) are put in a global queue, and only one message is allowed to be processed at a time. Such ‘serialization’ of the BGP algorithm may arise in practice if there are long link delays in a network.
- iii The messages generated in each k -level iteration are reordered in the queue. Those messages that invalidate the currently installed path at each node are favored and put in the queue ahead of the others.

We present an algorithm in the Appendix for reordering the messages as stipulated in condition (iii) while preserving all the essential elements of BGP. The following notations are used to represent messages: An announcement of a new path by node i sent to its neighboring node j is given as $i \rightarrow j[path]$ where path is the set of nodes starting with node i . Similarly a withdrawal message originated at node i is represented by $i \rightarrow j[\infty]$

With these definitions, it is straightforward to construct a sequence of messages between any two nodes i and j for each k -level iteration. Consider the routing table at node i of a network at time t : (*013, 103, ∞ , ∞)

The above shows all valid paths from i to a fixed destination with the symbol * denoting the active route. In this case, node i has two possible paths to the destination via its two neighboring nodes 0 and 1 respectively. Let us assume that node i receives a new announcement from its neighbor, node 1: $1 \rightarrow i[1i3]$. Since this newly announced path will create a routing loop, node i rejects it and also deletes path 103 from its routing table. The only effect of the announcement is the deletion of an alternative path from the routing table. No new update is generated at node i for its neighbors. We consider such announcements a necessity for rapid convergence of a network following the withdrawal of a route since the removal of path 103 prevents it from being propagated during the next k -level ($k=4$) iteration as a new path $i103$.

On the other hand, suppose that node i receives an announcement from a different neighbor, node 0 (instead of node 1): $0 \rightarrow i[0i3]$. This time, however, path 013 is withdrawn and a new path $i103$ is announced by node i . This leads to more iterations of the shortest path algorithm until every possible path containing $i103$ has been explored.

The above discussion points out an important characteristic of BGP. In the absence of a fixed timer such as MinRouteAdver, the order in which announcements are processed at a node influences the rate of convergence for a path-vector algorithm.

Observation 4: If condition (iii) is applied to every new announcement message generated at any k -level, the algorithm will continue until all possible paths have been explored. Once the set of all possible paths is exhausted, the algorithm will stop after processing the final withdrawal messages. This is the basis of our conjecture that the complexity for the worst case is $O(e(n-1)!)$.

Observation 5: The communication complexity, or the number of announcements and withdrawals, is much larger than the state complexity bound $O(e(n-1)!)$ as seen in Observation 1. Each announcement of a new path is forwarded to all $(n-1)$ neighbors of an AS, there by generating $(n-1)O(e(n-1)!)$ messages until convergence. The number of initial withdrawals is $(n-1)$ and in the worst case, the last k -level iteration (i.e. $k = n-1$) generates $(n-1)!$ messages, each of which ends in a withdrawal. Depending on the implementation details of BGP, this may result in $O(n-1)!$ withdrawals for the worst case. Therefore, for the worst-case BGP model, the number of messages (both withdrawals and announcements) grows faster than exponentially with n .

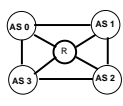
We present an algorithm that provides an ordering of messages as per condition (iii) while preserving the essential features of BGP in the Appendix. The algorithm forces the path-vector algorithm to explore all $k = 1, 2, \dots, (n-1) - length$ paths until convergence and results in the worst-case behavior of BGP. As pointed out in a later Section, the best case convergence for BGP can be achieved in $O(n)$ steps. Since the Internet is not a complete graph and the link delays vary widely, the convergence behavior in practice will be in between these two bounds. However, our study fills an important gap in the analysis of path-vector algorithms.

5.3 BGP Convergence with MinRouteAdver

In the previous Subsection, we described the behavior of BGP convergence under the assumption of unbounded message propagation delay. In practice, discussions with Internet providers and analysis of our data indicates that the vast majority of BGP messages propagate within 30 seconds. In this SubSection, we describe BGP convergence behaviors when all messages propagate within the bounds of the MinRouteAdver timer.

Overall, MinRouteAdver introduces a dampening effect and synchronization to the converging system. In a complete graph after a failure of a given route, each AS will select a new active route, announce that route, and then suppress all further announcements for a minimum of 30 seconds. We show that MinRouteAdver effectively creates multiple 30 second *rounds* for the propagation of new path information. For the purpose of our analysis, we make the following assumptions:

1. MinRouteAdver is only applied to announcements as suggested in the BGP specification [20].
2. Only the BGP receiver performs loop-detection. The BGP specification calls for ASPath loop detection, but does not specify if detection should be performed by the sender or the receiver. Most commercial router implementations perform only receiver-side loop detection.
3. Path selection breaks ties between routes of equal length ASPath by a deterministic mechanism, such as lowest neighbor IP address. In the following analysis, we use lowest AS number in lieu of neighbor



Step	Time	Routing Tables	Messages Processing	Messages Queued in System
0	N/A	steady state 0(*R, 1R, 2R, 3R) 1(*R, *R, 2R, 3R) 2(*R, 1R, *R, 3R) 3(*R, 1R, 2R, *R)		steady state
1	N/A	R withdraws its route 0(*, *1R, 2R, 3R) 1(*0R, -, 2R, 3R) 2(*0R, 1R, -, 3R) 3(*0R, 1R, 2R, -)		0->1 01R 1->0 10R 2->0 20R 3->0 30R 0->2 01R 1->2 10R 2->1 20R 3->1 30R 0->3 01R 1->3 10R 2->3 20R 3->2 30R
2	N/A	announcement from 0 0(*, *1R, 2R, 3R) 1(*, -, *2R, 3R) 2(*0R, *1R, -, 3R) 3(*0R, *1R, 2R, -)	0->1 01R 1->0 10R 2->0 20R 3->0 30R 0->2 01R 1->2 10R 2->1 20R 3->1 30R 0->3 01R 1->3 10R 2->3 20R 3->2 30R	
3	N/A	announcement from 1 0(*, -, *2R, 3R) 1(*, -, *2R, 3R) 2(*01R, 10R, -, 3R) 3(*01R, 10R, 2R, -)	1->0 10R 1->2 10R 1->3 10R 2->0 20R 2->1 20R 2->3 20R 3->0 30R 3->1 30R 3->2 30R	
4	N/A	announcement from 2 0(*, -, -, *3R) 1(*, -, -, *3R) 2(*01R, 10R, -, *3R) 3(*01R, 10R, 20R, -)	2->0 20R 2->1 20R 2->3 20R 3->0 30R 3->1 30R 3->2 30R	
5	30	Min Route Timer expires announcement from 3 0(*, -, -, -) 1(*, -, -, *20R, 30R) 2(*01R, 10R, -, 30R) 3(*01R, 10R, 20R, -)	3->0 30R 3->1 30R 3->2 30R	0->1 W 1->0 120R 2->0 201R 3->0 301R 0->2 W 1->2 120R 2->1 201R 3->1 301R 0->3 W 1->3 120R 2->3 201R 3->2 301R
6	N/A	withdrawal from 0 0(*, -, -, -) 1(*, -, -, *20R, 30R) 2(*, -, -, *30R) 3(*, -, -, *20R, -)	0->1 W 0->2 W 0->3 W	1->0 120R 1->2 120R 1->3 120R 2->0 201R 2->1 201R 2->3 201R 3->0 301R 3->1 301R 3->2 301R
7	N/A	announcement from 1 0(*, -, -, -) 1(*, -, -, *20R, 30R) 2(*, -, -, *30R) 3(*, -, -, *20R, -)	1->0 120R 1->2 120R 1->3 120R	2->0 201R 2->1 201R 2->3 201R 3->0 301R 3->1 301R 3->2 301R
8	N/A	announcement from 2 0(*, -, -, -) 1(*, -, -, *30R) 2(*, -, -, *30R) 3(*, -, -, *20R, -)	2->0 201R 2->1 201R 2->3 201R	3->0 301R 3->1 301R 3->2 301R
9	60	Min Route Timer expires announcement from 3 0(*, -, -, -) 1(*, -, -, -) 2(*, -, -, *301R) 3(*, -, -, *20R, 201R, -)	3->0 301R 3->1 301R 3->2 301R	1->0 W 1->2 W 1->3 W 2->0 2301R 2->1 2301R 2->3 2301R 3->0 3120R 3->1 3120R 3->2 3120R
10	N/A	withdrawal from 1 0(*, -, -, -) 1(*, -, -, -) 2(*, -, -, *301R) 3(*, -, -, *201R, -)	1->0 W 1->2 W 1->3 W	2->0 2301R 2->1 2301R 2->3 2301R 3->0 3120R 3->1 3120R 3->2 3120R
11	N/A	announcement from 2 0(*, -, -, -) 1(*, -, -, -) 2(*, -, -, *301R) 3(*, -, -, -)	2->0 2301R 2->1 2301R 2->3 2301R	3->0 3120R 3->1 3120R 3->2 3120R
12	90	Min Route Timer expires announcement from 3 0(*, -, -, -) 1(*, -, -, -) 2(*, -, -, -) 3(*, -, -, -)	3->0 3120R 3->1 3120R 3->2 3120R	2->0 W 2->1 W 2->3 W 3->0 W 3->1 W 3->2 W
13	N/A	process withdrawals 0(*, -, -, -) 1(*, -, -, -) 2(*, -, -, -) 3(*, -, -, -)	2->0 W 2->1 W 2->3 W 3->0 W 3->1 W 3->2 W	

Figure 6: Example of BGP bouncing problem with MinRouteAdver.

IP address.

In Figure 6, we provide an example of BGP convergence with bounded delay and MinRouteAdver. As before in Figure 5, we show five columns reflecting the current computational step, time, the routing state of each node, the messages processed at that step, and the messages in the global queue. At step 1, nodes 0 through 3 detect the failure of route R and withdraw their primary route. We then proceed as in the prior unbounded delay convergence example, but unlike Figure 5, we note that MinRouteAdver suppresses the generation of new update messages from each node until the timer expires after step 4.

Figure 7 provides data illustrating the reduction in state and communication complexity due to MinRouteAdver. We will discuss the modified MinRouteAdver simulation results provided in Figure 7(c) in the Conclusion. The first two tables in Figure 7 show the results of (*omitted for blind review*) simulations for both unbounded (a) and bounded BGP (b) message propagation for complete graphs of size 4-7. The “states column” corresponds to state complexity and “messages” represents communication complexity. We observe that a seven node system generates over 23,000 BGP messages in the worst case, as compared to the 140 messages generated for the same number of nodes under the constraints of MinRouteAdver.

Proof of Lower Bound on BGP Convergence using MinRouteAdver

We now show that with the adoption of MinRouteAdver, the lower bound on convergence for BGP requires at least $(n - 3)$ rounds of the MinRouteAdver timer, where n is the number of autonomous systems. In the following discussion, we refer to the complete graph of five nodes shown in Figure 6.

Observation 1: The best case algorithm for MinRouteAdver, when applied to a complete graph of size n results in complete withdrawal of at most one node at the end of the first round.

Nodes	Time	States	Messages	Nodes	Time	States	Messages	Nodes	Time	States	Messages
4	N/A	12	41	4	30	11	26	4	30	11	26
5	N/A	60	306	5	60	26	54	5	30	23	54
6	N/A	320	2571	6	90	50	92	6	30	39	92
7	N/A	1955	23823	7	120	85	140	7	30	59	140

(a) Unbounded (b) MinRouteAdver (c) Modified

Figure 7: Simulation results for convergence with unbounded delay, MinRouteAdver, and modified MinRouteAdver.

We examine the consequences of a withdrawal of route R which is initially directly connected to every node in the graph. The routing table at each node is represented in step 0 of Figure 6.

In the event of a withdrawal message from node R , every node in the system, except node 0 will choose the path “ $0R$ ” as the active route; node 0 will announce path “ $1R$ ”. Under the MinRouteAdver timer, node 0 will receive $(n - 2)$ announcements from its neighbors and will try to replace its alternate paths (i.e. paths “ $1R$ ”, “ $2R$ ”, “ $3R$ ” etc.) with the newly received information. However, each of these new updates results in a loop and therefore, node 0 removes these paths. Node 0 then sends a withdrawal message to all its neighbors, as it no longer has a valid path to R .

The above description applies to any complete graph of size n since the path though that node “ $0.(n - 1)$ ” is the best route to reach node $(n - 1)$ from any other node at iteration $k = 2$. This explains why node 0 will always be withdrawn irrespective of the size of the graph.

Observation 2: The primary effect of a MinRouteAdver timer is to impose a monotonically increasing path metric for successive k -level iterations.

This is the most important contribution of the MinRouteAdver timer and also helps to intuitively explain rapid convergence of general graphs in the event of a route failure. At the end of the k -th MinRouteAdver round, only the next higher level (i.e. $k+1$) paths are announced. This means that there should be no pending path announcements of length k in the global queue when a $(k + 1)$ -length path has been announced by any node. Under a MinRouteAdver timer, a node must process **all** $(n - 1)$ announcements from its neighbors before it can send out a new announcement. The order in which it processes each announcement does not matter since it receives only one message from each of its neighbors. This newly received path may either result in a loop or replace the existing path to that neighbor. If it replaces an existing path, we need to show that the path being replaced is a shorter path than the path replacing it. If this is true for all nodes, each of the nodes will send out a longer path in the next MinRouteAdver timer. To see this, let us consider the 5-node example again. The initial routing tables at steady state prior to route withdrawal are shown in step 0 of Figure 6.

Upon receiving the withdrawals from node R , twelve messages are generated as shown in step 1 of Figure 6. Let us consider the messages waiting to be processed at node 1. Its routing table currently consists of

paths of length two: $1(*0R, \infty, 2R, 3R)$. However, each of the arriving messages at node 1 replaces the corresponding 2-length path to a 3-length path. As a result, once all $(n - 1)$ messages have been processed at node 1 under the MinRouteAdver timer, its routing table becomes as represented in step 5 of Figure 6; and a new longer ($k = 4$) path “120R” is announced to its neighbors in the next iteration. Let us contrast this situation with the case when no MinRouteAdver timer is allowed. In this case, node 1 will process **only one** message before it announces a new path. If the particular message $0 \rightarrow 1 [01R]$ was processed (without the MinRouteAdver timer), the routing table at node 1 would become: $1(\infty, \infty, *2R, 3R)$ resulting in another 3-length path “12R” to be announced to its neighbors. The same is true for all the other nodes in the example graph. One can also confirm this observation for $k = 4$ in the above example, and for longer paths in graphs of size $n > 5$.

Therefore, the effect of the MinRouteAdver timer is to impose a global state synchronization which results in deletion of all k -length paths before a new longer path is announced by any node.

Observation 3: Since $k_{max} = n - 1$ and each MinRouteAdver timer deletes paths of length k at a particular k -level iteration, there will be at least $(n - 1)$ MinRouteAdver rounds for the best-case algorithm. (This follows readily from **Observation 2**.)

Observation 4: The above estimate for the number of MinRouteAdver rounds can be further reduced to $(n - 3)$ for a complete graph of size n greater than 3. This result follows from the observation that for 2-node and 3-node complete graphs, BGP converges on a route withdrawal within a single MinRouteAdver period. We re-emphasize that the above observations are valid when the best-case algorithm with the MinRouteAdver timer is applied to a complete graph. The degree to which MinRouteAdver preserves the monotonicity of each k -level iteration in incomplete graphs is a topic of our current research.

5.4 MinRouteAdver and Experimental Results

Armed with a model of BGP convergence, we now return to the results presented in Section 4. Why do Tup/Tshort converge more quickly than Tdown/Tlong? The explanation lies in the observation that, like the comparison between DV algorithms and BGP, Tup/Tshort are strictly increasing while Tdown/TLong are monotonically increasing. Intuitively, once a node receives an update during Tup and selects an active path, the node will never choose a route with a longer path. In contrast, since the Tdown implicit metric of infinity is longer than all possible ASPaths, each node will failover to secondary paths until all paths have been eliminated. If we assume bounded delays, then Tup has a computational complexity of $O(1)$ and Tdown of $O(N)$.

Unlike Tup/Tshort, Figure 2(b) shows a slight variation between the relative latencies of Tlong and Tdown. Due to the effects of MinRouteAdver, we might expect Tlong to converge at the same rate or slower than Tdown. Analysis of the data, however, shows that if the prepended ASPath associated with a Tlong is not sufficiently long, then this route might be preferred over shorter paths at some point during convergence.

In effect, these Tlongs would resemble both Tshort and Tdown events and represent the average of the two. In our experiments, we observed a small number of paths with lengths four times the steady-state average following Tdown and Tlong events. As described in Section 3, we only associated a path of only three times the steady-state average with the injected Tlongs.

Although we did not associate a sufficiently long ASPath with Tlong to render Tshort completely indistinguishable from Tup, or Tdown indistinguishable from Tlong; Tshort/Tup enjoy the property that routing information associated with the shortest ASPath will usually propagate faster than routing information associated with longer paths. This speed advantage arises because longer ASPaths, by definition, are formed by routing information traveling through more BGP autonomous systems, each of which adds some additional latency. Although convergence following Tshort theoretically may have introduced added oscillations over Tup as the system explored ASPaths longer than Tlong, such oscillations are unlikely in practice.

In Figure 3, we described significant variations between the convergence latencies of five ISPs. We noted that these differences were independent of both geographic and network distance. As we showed in Section 5.3, if the Internet were truly a complete mesh we would expect all ASes to exhibit the same convergence behaviors. Instead, analysis of the data shows that these variations directly relate to a number of topological factors, including the length and number of possible paths between a AS and a given destination. The number of available paths is a factor of peering relationships, transit policies/agreements and the implementation of filters by both the AS and downstream ASes.

Analysis of Figure 3(a) also shows that the Tdown convergence times of between 0 and 180 seconds directly relate to the number of MinRouteAdver rounds. Our data shows a strong correlation between the average ASPath length during Tdown events and convergence latency. Specifically, as the point of injection ISP1 always announced routes of length one; ISP3 averaged 2.6, and ISP5 averaged ASPaths of length 6. These results corresponds with our $30(n - 3)$ lower bound on MinRouteAdver convergence times.

Finally, we examine the 0 to 30 second convergence latencies exhibited in Figure 3(b). As described earlier, Tup events are strictly increasing and do not typically generate multiple announcements. Figure 2 shows that most ISPs average one update message following a Tup event. Since MinRouteAdver does not impact the first announcement of a route, we might expect Tup latencies to be significantly less than 30 seconds, as they would reflect only the network latency and router processing delays along a single path. Discussion with a major router vendor, however, indicates that at least one widely deployed router implements MinRouteAdver on a per peer basis instead of the (destination prefix, peer) tuple. We emphasize that this implementation choice is in accordance with the BGP specification [20] and improves router memory utilization. A per peer timer, however, introduces some portion of the MinRouteAdver delay to Tup/Tshort updates. If a router has previously sent any update to a given peer within the last 30 seconds, then a new Tup announcement destined for the same peer will also be delayed until the expiration of the per-peer MinRouteAdver timer.

6 Conclusion

As the national and economic infrastructure become increasingly dependent on the global Internet, the availability and scalability of IP-based networks will emerge as among the most significant problems facing the continued evolution of the Internet. This paper has argued that the lack of inter-domain failover due to delayed BGP routing convergence will potentially become one of the key factors contributing to the gap between the needs and expectations of today's data networks. In this paper, we demonstrated that multi-homed failover now averages three minutes, and may trigger oscillations lasting as long as fifteen minutes. Further, we showed that these delays will grow linearly with the addition of new autonomous systems to the Internet in the best case, and exponentially in the worst. These results suggest a strong need to reevaluate applications and protocols, including emerging QoS and VoIP standards [11], which assume a stable underlying inter-domain forwarding infrastructure and fast IP path restoration.

While MinRouteAdver significantly reduces the computational and communication complexity of BGP convergence, the timer also artificially creates multiple thirty-second rounds, which will delay end-to-end failover in most cases. As we showed in Section 5.3, these rounds form due to the delay in the exchange of path vectors containing mutually dependent routes. Although the BGP specification describes ASPath loop detection, [20] does not specify where the detection should occur. If loop detection is performed on both the sender and receiver side, in the best case all mutual dependencies will be discovered and eliminated within a single round. Figure 7(c) provides simulation results of MinRouteAdver modified to perform sender-side loop detection. We note that for all node sizes, modified MinRouteAdver converges within a single thirty second round. We also observe that although the communication complexity remained the same, modified MinRouteAdver exhibits improved state complexity over unmodified MinRouteAdver.

We discussed this proposed modification to MinRouteAdver with a number of router vendors, and at least one indicated that all future versions of a widely deployed router will include both sender and receiver-side ASPath loop detection. The elimination of rounds, however, requires that the router does not apply MinRouteAdver to withdrawals as specified in [20]. At least one major router vendor has made an implementation decision to apply MinRouteAdver to both announcements and withdrawals. A discussion of the motivation and engineering tradeoffs for applying MinRouteAdver to withdrawals is outside the scope of this paper.

But even with the above suggested changes to ASPath loop detection, BGP path changes will still trigger temporary oscillations and require many seconds longer than the current PSTN restoration times. We can certainly improve BGP convergence through the addition of synchronization, diffusing updates [7] and additional state information [5], but all of these changes to BGP come at the expense of a more complex protocol and increased router overhead. The extraordinary growth and success of the Internet is arguably due to the scalability and simplicity of the underlying protocols.

References

- [1] North American Network Operators Group (NANOG). <http://www.nanog.org>.
- [2] Omitted for double-blind review.
- [3] Omitted for double-blind review.
- [4] K. Bhargavan, D. Obradovic, and C. Gunter. Formal Verification of Distance Vector Routing Protocols. *To appear in Proc. of IEEE INFOCOM*, March 2000.
- [5] C. Cheng, R. Riley, S. Kumar, and J.J Garcia Aceves. A Loop-Free Extended Bellman-Ford Routing Protocol Without Bouncing Effect. *Proc. of the ACM SIGCOMM*, pages 224–236, August 1989.
- [6] Lixin Gao and J. Rexford. Stable Internet Routing Without Global Coordination. *To appear in the Proc. of ACM SIGMETRICS*, June 2000.
- [7] J. Garcia-Luna-Aceves. Loop-free Routing Using Diffusing Computations. *IEEE/ACM Transactions on Networking*, February 1993.
- [8] T. Griffin, F. Shepard, and G. Wilfong. An Analysis of BGP Convergence Properties. *Proceedings of the ACM SIGCOMM*, August 1999.
- [9] T. Griffin, F. Shepard, and G. Wilfong. Policy Disputes in Path-Vector Protocols. *To appear in Proc. of IEEE Infocom*, March 2000.
- [10] B. Halabi. *Internet Routing Architecture*. Cisco Press, 1997.
- [11] D. Hampton, H. Salama, and D. Shah. The IP Telephony Border Gateway Protocol (TBGP), June 1999. draft-ietf-iptel-glp-tbgp-01.txt.
- [12] John Hawkinson. Cisco Routing FAQ. <http://www.faqs.org/faqs/cisco-networking-faq>.
- [13] M. Kaat. Overview of 1999 IAB Network Layer Workshop, November 1999. <http://www.ietf.org/internet-drafts/draft-ietf-iab-ntwlyrws-over-01.txt>.
- [14] C. Labovitz. *Scalability of the Internet Backbone Routing Infrastructure*. PhD thesis, University of Michigan, August 1999.
- [15] 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.
- [16] C. Labovitz, G.R. Malan, and F. Jahanian. Internet Routing Instability. *IEEE/ACM Transactions on Networking*, August 1997.
- [17] C. Labovitz, G.R. Malan, and F. Jahanian. Origins of Pathological Internet Routing Instability. *Proc. of the IEEE INFOCOM*, March 1999.
- [18] V. Paxson. End-to-End Internet Packet Dynamics. *Proc. of the ACM SIGCOMM*, 1997.
- [19] R. Perlman. *Interconnections, Second Edition*. Addison-Wesley, Reading Massachusetts, 1999.
- [20] Y. Rekhter and T. Li. A Border Gateway Protocol 4 (BGP4), September 1999. draft-ietf-idr-bgp4-09.txt.
- [21] F. Schneider, S. Bellovin, and A. Inouye. Building Trustworthy Systems: Lessons from the PTN and Internet. *Internet Computing*, pages 64–72, November 1999.

- [22] Kannan Varadhan, Ramesh Govindan, and Deborah Estrin. Persistent Route Oscillations in Inter-Domain Routing. Technical Report USC CS TR 96-631, Department of Computer Science, University of Southern California, February 1996.
- [23] William Zaumen and J. J Garcia-Luna Aceves. Dynamics of Distributed Shortest-Path Routing Algorithms. *Proc. of the ACM SIGCOMM*, August 1991.

A Algorithm for Upper Bound on BGP Convergence

The reordering algorithm requires that all n nodes of a network are labeled from 0 to $n - 1$ and the node directly connected to the destination is the $(n - 1)$ th node. The steps of the algorithm are as follows (all messages are placed in a global queue and processed one at a time):

- (a) reorder the initial withdrawal messages from $(n - 1)$ th node in increasing order, i.e.,

$$(n - 1) \rightarrow 0[\infty][r], (n - 1) \rightarrow 1[\infty][r], (n - 1) \rightarrow 2[\infty][r], \dots, (n - 1) \rightarrow (n - 2)[2][r]$$

Once a set of messages are reordered (denoted by $[r]$ next to the path information), these must be processed before any new messages can be reordered.

- (b) do $k = 1, 2, \dots, (n - 1)$ until convergence, the following:

- (i) process each reordered message, place any resulting withdrawal or announcement message at the end of the queue. Repeat (i) until all reordered messages have been processed.
- (ii) perform a 2-pass radix sort on the remaining messages ($i \rightarrow j[path]$) in the queue. e.g. the following set of messages

$$0 \rightarrow [013], 0 \rightarrow 2[013], 1 \rightarrow 0[103], 1 \rightarrow 2[103], 2 \rightarrow 0[203], 2 \rightarrow 1[203]$$

will be sorted as

$$1 \rightarrow 0[103][r], 2 \rightarrow 0[203][r], 0 \rightarrow 1[013][r], 2 \rightarrow 1[203][r], 0 \rightarrow 2[013][r], 1 \rightarrow 2[103][r]$$

- (iii) reorder the resulting messages by interleaving them, i.e. picking a message from each bucket in turn until all buckets are empty.

In steps (ii) and (iii), if there are messages with same values of i and j , the sorting and reordering must preserve the order in which these messages appeared in the queue at the end of step (i).