# Managing Routing Disruptions
# in Internet Service Provider Networks

Renata Teixeira
Laboratoire Lip6
Université Pierre et Marie Curie
Paris, France
renata.teixeira@lip6.fr

Jennifer Rexford
Dept. of Computer Science
Princeton University
Princeton, NJ, USA
jrex@cs.princeton.edu

*Abstract*— **Customers of Internet Service Providers (ISPs) are increasingly interested in running applications such as voice over IP, video games, and commercial transactions. This new range of applications cannot tolerate poor network performance (high delays or low available bandwidth) or network instability (periods of loss or variation in delay or available bandwidth). Unfortunately, routine events such as equipment failures or planned maintenance cause routing changes, which may lead to transient service disruptions or persistent performance problems. Operators of ISP networks are faced with the challenge of minimizing routing disruptions using current routing technology, which offers little control. In this paper, we discuss routing disruptions from an ISP perspective. First, we describe the causes and effects of routing changes. Then, we provide a set of network design guidelines and operational practices that network operators can use to reduce the impact of routing changes in their network.**

## I. INTRODUCTION

As the Internet becomes an ever more critical part of the world's communication infrastructure, Internet Service Providers (ISPs) are under increasing pressure to provide good, predictable performance to a wide range of applications. Unfortunately, routine events such as equipment failures and planned maintenance can lead to long-term changes in path properties (e.g., higher round-trip times and lower available bandwidth) and serious transient disruptions (e.g., high loss and delay during routing-protocol convergence). These service disruptions are a significant problem for applications such as voice over IP (VoIP), video games, and commercial transactions. For example, a recent study of VoIP performance found that most service disruptions occur during routing changes [1], not because the network lacks sufficient resources for carrying the traffic. A disruption lasting a few hundred milliseconds is long enough to interrupt a phone conversation or a video game, and other applications such as Web transactions are visibly affected by disruptions lasting a few seconds.

ISPs take great care in designing and operating their networks to prevent routing disruptions and limit their scope and impact. However, the IP routing protocols were not designed with these requirements in mind, leaving network operators in the difficult situation of "working around" the many limitations of today's technology. The operators must grapple with four main challenges:

1) **Indirect control over the flow of traffic:** The operators have only indirect control over how the routers select paths, making it difficult to satisfy complex objectives for balancing load and minimizing disruptions.
2) **Large reactions to small changes:** The routing protocols often over-react to small changes in the network topology and configurable parameters, at the expense of network robustness.
3) **Slow routing-protocol convergence:** During routing-protocol convergence, data packets are lost, delayed, and delivered out of order, causing a serious degradation in end-to-end performance.
4) **Poor support for planned events:** Although maintenance activities are planned in advance, the routing protocols cannot gracefully move the traffic to new paths beforehand.

The first two limitations make it difficult for operators to control the flow of traffic in steady state, whereas the last two limitations relate to transient disruptions that occur during routing-protocol convergence.

In this paper, we describe the state of the art for how the operators of large ISPs can reduce the frequency and impact of routing disruptions caused by internal network events. To set the stage, Section II presents a brief overview of routing in ISP networks. Then, Section III explains how the scope and impact of the routing disruptions depends on the underlying network events. Section IV describes how ISPs manage routing disruptions through the design and configuration of the network, as well as the operational practices for making changes to the network. Section V concludes the paper.

## II. ROUTING IN ISP NETWORKS

The Internet is an interconnection of thousands of Autonomous Systems (ASes), where each AS is a collection of routers and links managed by a single institution. Most of these ASes buy upstream connectivity from an Internet Service Provider (ISPs). ISPs face many unique routing challenges because of their role as "transit" providers. First, the large size of ISP backbones—often hundreds of routers and thousands of links—introduces scalability challenges in running routing protocols. Second, an ISP must connect to numerous customers and many other ISPs at diverse geographic locations; a pair of ISPs may connect in multiple geographic locations for richer connectivity and fault tolerance. Third, the ISP must know how to reach every destination prefix (block of IP

addresses) in the Internet. In today's Internet, an ISP typically stores routing information for more than 150,000 prefixes[1], which means that each forwarding table has more than 150,000 entries. To satisfy these requirements, ISPs typically run both intradomain and interdomain routing protocols on their many routers.

Our discussion draws on the example in Figure 1. This example illustrates the behavior of routing protocols running inside and among ASes using an internetwork with five ASes. The traffic from the web server to the user enters the Small ISP network at ingress router $C$, which has two choices of egress routers: $A$ and $B$. The interdomain routing protocol is responsible for selecting which egress router to use to forward traffic to the user; say that it selects egress $A$. Then, the intradomain routing protocol selects the path from $C$ to $A$, in this case $D$ is the next hop in the path to $A$. Router $C$ combines the information from the two routing protocols to construct a forwarding table that maps the user's prefix to next hop $D$ (via the two interfaces going from $C$ to $D$). We now summarize the path-selection process for each routing protocol, for more details refer to [2].
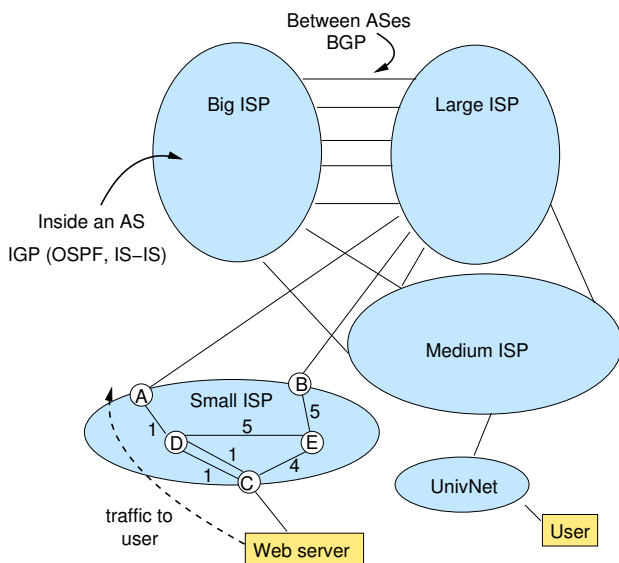


Fig. 1.   Example of an internetwork with five ASes and the details on the internal topology of the Small ISP.

### A. Intradomain Routing: Metric-Based

Inside each network the intradomain routing protocol, or Interior Gateway Protocol (IGP), is responsible for selecting the path between the ingress router and the egress router. Intradomain routing protocols compute shortest paths based on "metrics" or "link weights" statically configured by network operators. Network operators usually configure link weights to achieve some traffic-engineering goal such as balancing the traffic load or minimizing propagation delay in the network.

The most common IGPs in large ISPs today are OSPF and IS-IS, which are link-state routing protocols. Each router has

[1]For an up-to-date view of the number of routable prefixes, see the public data repository at http://www.route-views.org.

a complete view of the network topology. IGP messages are flooded periodically and in response to topology changes, such as link weight changes and equipment going up or down. Link weights only change when reconfigured by network operators, not dynamically to adapt to changes in network conditions such as congestion. Each router runs Dijkstra's algorithm to compute the shortest paths in terms of link weights to every other router and uses the results to build the forwarding table. In Figure 1, router $C$ selects the shortest path $C, D, A$ with distance 2 to reach router $A$.

### B. Interdomain Routing: Policy-Based

In the Internet, the Border Gateway Protocol (BGP) is the interdomain routing protocol responsible for exchanging reachability information about external destination prefixes, or prefixes that belong to other ASes. BGP is responsible for selecting the *AS path* (or which ASes have to be traversed) to reach a destination prefix. BGP is a path-vector protocol that allows each AS to apply local policies in selecting and propagating routes for each destination prefix. Network administrators use BGP policies to express business relationships between ISPs and customers, or to determine preference among choices of connections (e.g., to determine backup paths or to balance load) [3].

Two routers establish a BGP session to exchange BGP update messages. A large backbone network has BGP sessions with multiple neighboring ASes and can also have multiple BGP sessions with each neighbor AS. As a result, a router may receive routes for a destination prefix from multiple neighbors. A BGP route has a number of attributes (such as local preference, next-hop, AS-path, origin type, and Multiple-Exit-Discriminator) that are conveyed in route advertisements and can be manipulated by local policies. The router invokes a *decision process* to select exactly one "best" route for each destination prefix among all the routes learned from its neighbors.

In the example in Figure 1, routers $A$ and $B$ learn routes to the user prefix from *external* BGP (eBGP) sessions. Each of $A$ and $B$ select the best of the eBGP-learned routes and propagate this choice to routers inside the AS using *internal* BGP (iBGP). Routers in the AS have the choice between these two routes. Since we assume that the network operator gives the same preference to both routes, then for routers inside the AS both routes look "equally good" at the BGP level (with the same local preference, number of AS hops, etc.). This leaves $C$, $D$, and $E$ with the dilemma of choosing between egress points $A$ and $B$. When there are multiple egress points for a prefix, routers direct traffic to the *closest* egress point—the router with the smallest intradomain distance (e.g., router $C$ selects egress point $A$ with distance 2 from $C$). This step in the BGP decision process is called *early-exit* or *hot-potato* routing.

### III. ROUTING DISRUPTIONS ARE INEVITABLE

Network operators configure both IGP link weights and BGP policies to control the flow of traffic in the network

to achieve good network performance and utilization. However, routine events such as link failures may trigger routing changes. This section first discusses the kinds of events that cause routing changes and the steps that each router has to take to adapt to these events. Then, we present a taxonomy of routing changes according to their scope and impact.

### A. Routing Reaction to Network Events

Events that disrupt routing can happen inside the network (e.g., equipment failures, routing software crash, planned maintenance, and routing configuration changes for traffic engineering) or outside (e.g., BGP updates from neighboring domains). In this paper, we focus on minimizing the impact of internal events, because these events are under the control of a single ISP. Network operators cannot fix the cause of external events directly. However, they often reconfigure IGP metrics or BGP policies to adjust to changes in the set of BGP-learned routes and fluctuations in the incoming traffic. This reconfiguration is an internal event and triggers routing changes. We call a *routing disruption* any transient or persistent perturbation of network performance that is caused by a routing change.

IP links represent logical connectivity, which can be implemented using a variety of link-level technologies (e.g., optical fiber, Ethernet, and FDDI). Optical links offer the possibility of automatic restoration after failures of link components (such as optical amplifiers or fiber segments), without triggering topology changes at the IP level. However, most ISPs rely on IP routing for restoration [4]. First, optical restoration is expensive in terms of spare capacity. Second, optical restoration cannot protect against interface and router failures, which represent a significant fraction of failures. Therefore, often failures of lower-layer components are visible as link failures at the IP layer. In fact, it is common that multiple IP links share a single optical amplifier or fiber segment, leading to multiple simultaneous failures at the IP level when an optical component fails.

Not all routing disruptions are caused by unexpected equipment failures. A large fraction of them are caused by routine maintenance. For instance, almost half of all intradomain events during a five-month period in the Sprint backbone happened during the maintenance window [4]. Maintenance activities happen more often than most people realize. Large ISP networks have hundreds of routers and thousand of links. In a network with 700 routers, if the operating system on the routers need to be upgraded (say) once a year, this means that on average two routers are under maintenance per day. In addition, construction activities require moving fibers, links may be added or removed, interfaces may crash, optical amplifiers may need repair or replacement, etc.

Routers react to internal network events in stages:

- **Detection**: If the event corresponds to an IGP or BGP reconfiguration, then the detection is immediate. However, in case of equipment failures, it may take some time before the router detects it either by receiving an explicit alarm from the underlying hardware (in case of Packet over SDH/SONET links) or by detecting consecutive losses of IGP keepalive messages. Until the failure is detected, the router continues forwarding data packets into the failed link, leading to packet loss.
- **Propagation**: Upon detecting the link failure, the incident routers generate link-state advertisement messages to inform the other routers about the change. To avoid overloading the network with multiple messages, routers employ a timer to limit the message-generation rate.
- **Path recomputation**: Upon receiving a message reporting a change, routers need to recompute their best paths to all other routers. Since the path recomputation consumes many CPU cycles, routers impose a minimum time interval between two consecutive computations. Messages that arrive during this interval are grouped. If BGP routes are affected, then the router also needs to re-run the BGP decision process for any affected destination prefixes. Some vendors' BGP implementations do not react immediately to changes in the IGP distances. Instead, these routers have a scan process that runs periodically to sequence through the BGP routing table and revisit the BGP routing decision for each prefix. This process runs only once a minute in some router implementations [5].
- **Forwarding-table update**: Finally, routers update their forwarding tables. In high-speed routers, each line card has its own forwarding table for rapid forwarding of data packets. Hence, a router typically needs to update multiple forwarding tables with new entries.

The period of time between when an event happens and the last router updates its routing information is called *routing convergence*, whereas the time for the last router to update its forwarding table is called *forwarding-plane convergence*. Although routing and forwarding-plane convergence are related, forwarding-plane convergence might take less time (because some of the routing protocol messages do not lead to changes in the forwarding path) or more time (because of extra delay to update the forwarding table after routers select best paths, especially when many paths change at the same time). During forwarding-plane convergence, packets may be caught in forwarding loops, which may cause them to be lost, delayed, or delivered out of order. The transient period of disruption associated with convergence varies depending on the scope of the routing change.

### B. Taxonomy of Routing Changes Triggered by Internal Events

The convergence delay and the volume of traffic affected depends on the scope of the routing change:

**Local IGP routing changes**: In Figure 1, when one of the links between $C$ and $D$ fails (e.g., because of a failure of the interface associated with the link), packets from the web server to the user may get lost for a short period of time before router $C$ detects the failure and starts forwarding all packets using the remaining link. When a router has multiple outgoing links on shortest paths, forwarding-plane convergence usually takes less than a second [4]. Even after $C$ updates its forwarding table, a performance disruption may continue due to congestion on the remaining link from $C$ to $D$.

**IGP routing changes affecting multiple routers**: If both links between $C$ and $D$ fail simultaneously (e.g., because they

share the same fiber or optical amplifier), then multiple routers need to update their forwarding tables. IGP convergence in a network with hundreds of routers takes several seconds. After IGP convergence, there may be persistent performance problems either because the path $C, E, D, A$ is congested or because it leads to longer propagation delay.

**Egress-point changes**: In addition to an IGP routing change, the failure of both links between $C$ and $D$ causes some routers to change their selection of egress point for some destination prefixes. The routing change causes the distance from $C$ to $A$ to increase from 2 to 10. Although the BGP route through $A$ is still available, the IGP distance change would cause $C$ to select the route through egress point $B$ because of hot-potato routing. Even if the BGP-level route does not change, the convergence of BGP routes after an IGP change can take a couple of minutes [5]. In a large ISP network, a single IGP routing change can affect tens of thousands of destination prefixes at the same time. The traffic shift may cause an abrupt increase in traffic along the new path, much to the surprise of downstream neighbors, as well as changes in the end-to-end path characteristics (e.g., delay and available bandwidth) seen by the user.

**BGP path change**: Take again the example in Figure 1 and imagine the scenario in which both links between $C$ and $D$ fail, and router $A$ is using a route through Big ISP. When router $C$ shifts from the route via $A$ to the one via $B$, the AS path of the new route is different; hence $C$ must send a BGP update to its customer to report the change. Besides the impact of the egress change, this event may impact the upstream path. Sending BGP updates to neighbors may have unpredictable effects, because it may impact whether or where the neighbor chooses to direct its traffic. Previous studies have shown BGP may take several *minutes* to converge after a topology or policy change [6]. BGP's long convergence time is mainly due to the exploration of alternate routes before selecting a stable route. The details of path exploration depend on timing details at routers throughout the Internet. Interarrival times of around 30 seconds are quite common for external routing changes, since many routers use a 30-second minimum-route-advertisement timer for eBGP sessions [6].

Network operators would prefer that internal events be contained inside the network as much as possible (only causing IGP changes), which cause shorter transient disruptions and have no impact in where the traffic enters or leaves the network. Even though network operators often know the time and the location of some of these events in advance, they are largely unable to prevent disruptions because of the way routing protocols react to internal events. For example, if an operator knows that a change will last only a few minutes (the time to reboot a router or exchange an interface card, for instance), then there is no need to trigger an egress-point or BGP path change. Or, if an operator reconfigures IGP weights for redistributing load inside the network, it would be unfortunate if this change caused BGP updates to neighboring domains (which may change where traffic enters the network) or egress-point changes (which changes where traffic leaves the network). In the next section, we discuss ways of minimizing the impact of internal events by reducing
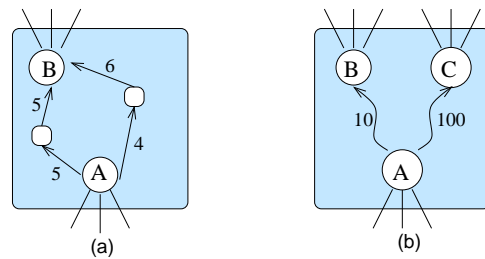


Fig. 2. Preventing BGP egress changes with redundant paths or by having a preferred egress router.

both the duration of routing convergence and the scope of the routing change.

## IV. "Working Around" Routing Disruptions

Preventing routing disruptions helps avoid shifts in traffic, extra delays in forwarding-plane convergence, and externally-visible BGP updates. However, current routing technologies do not offer network operators the control necessary to contain routing disruptions completely inside the domain and eliminate transient perturbation during routing convergence. In this section, we describe the best common practices for backbone design and configuration, and operational activities to reduce the negative effects of routing changes.

### A. Backbone Design and Configuration

To reduce the likelihood of link congestion and high propagation delay, the ISP can *design a backbone with enough bandwidth and path diversity* to accommodate the offered traffic after an equipment failure or during planned maintenance. *Overprovisioning* also reduces the need for the network operators to modify the routing-protocol configuration to move traffic to different paths when the offered traffic changes (say, due to natural time-of-day fluctuations). Reducing the frequency of configuration changes reduces the number of routing disruptions in the network.

In addition to accommodating extra traffic, path diversity can help reduce the effects of routing-protocol convergence. If the backbone has multiple shortest paths between a pair of routers, the failure of a link on one of these paths has only local consequences, as discussed earlier in Section III-B. For example, after the failure of one link between routers $C$ and $D$ in Figure 1, router $C$ immediately switches to forwarding all traffic over the other link, well before the rest of the routers have learned about the link failure. As such, the ISP should *favor backbone designs with multiple shortest paths*. The Sprint backbone follows this approach and, as such, typically does not experience transient forwarding loops during routing-protocol convergence [7]. Having multiple shortest paths also reduces the likelihood of egress-point changes due to hot-potato routing. As shown in Figure 2(a), router $A$ has two shortest paths (with an IGP distance of 10) to egress point $B$. This decreases the likelihood that a single internal event would change the IGP distance to reach $B$, which would tend to prevent egress-point changes [5].

Careful design and configuration of the network can also reduce the likelihood that internal events trigger a change in egress point. For example, in Figure 2(b), router $A$ has a small IGP distance of 10 to reach egress point $B$ and a much larger IGP distance of 100 to reach $C$. This significantly reduces the likelihood that small variations in IGP distances would trigger $A$ to switch to the BGP route learned from $C$. In addition to selecting appropriate IGP weights, network designers can also avoid egress-point changes by selecting the location of connections to neighboring domains. When adding new peering connections to neighboring domain, network designers can use a model of the network and metrics of network sensitivity to events [8] to *prioritize adding connections at locations that are sensitive to disruptions*, with the goal of alleviating the problem. This technique puts more routers in a situation like router $A$ in Figure 2(b). When this is not possible, the ISP can *avoid placing services that are sensitive to transient routing disruptions at locations that have high chance of experiencing egress-point changes*. For instance, an ISP can place a gaming server or key VoIP clients at a location that is more robust to internal changes.

The network designer can also reduce routing disruptions through *careful selection of the underlying technologies in the backbone*. For example, IGPs like OSPF and IS-IS depend on lost keepalive timers to detect link failures, leading to potentially long delays (with large timer values) or high overhead (with low timer values). Instead, the designer could build the backbone out of Packet Over SDH/SONET (POS) links that have explicit alarms for reporting failures. Today, POS links are commonly used in large ISP backbones, meaning that most link failures are detected relatively quickly at the link layer rather than through lost keepalive messages [4]. [2]

### B. Operational Practices

Given a network design, the network operators *adjust the configuration of the routing protocols to match the offered load to the available network resources*. Unfortunately, operators do not have direct control over the flow of traffic. Instead, they configure the IGP link weights and BGP routing policies to achieve their network-wide objectives, such as balancing link load, limiting propagation delays, and creating multiple shortest paths between pairs of routers. The operators *collect measurements of the offered traffic* and *use traffic-engineering tools that predict the effects of changes to the routing-protocol configuration* [9]. These tools can run optimization routines to select link weights and routing policies that satisfy the network objectives. However, the optimization problems are NP hard [9], even for the simplest objective functions, forcing the use of local-search techniques. Finding a good setting of the IGP metrics is especially difficult when routing must be

[2]As another example, the network designer could use a tunneling technology, such as Multi-Protocol Label Switching (MPLS) to direct traffic across the backbone, rather than running BGP on each internal router. These tunnels would avoid transient forwarding disruptions during the BGP convergence process when multiple routers switch to new egress points. MPLS also includes support for fast reroute through the construction of backup paths, which can reduce the effects of IGP convergence.
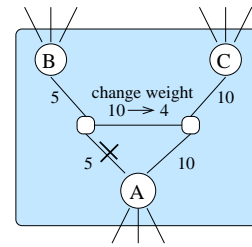


Fig. 3. $A$ still picks egress $B$ during maintenance

robust to equipment failures, since the optimization routine must explore the effects of many different failure scenarios [9].

Although overprovisioning reduces the frequency of configuration changes, some critical failures or maintenance activities require the operators to adapt the configuration. *For planned maintenance activities, the operators can try to schedule the change to occur during periods of lower traffic load*, when the network has enough capacity to accommodate the resulting shifts in traffic on to new paths. Still, the IGP topology change might lead to a shift in traffic from one egress point to another, leading to longer convergence delays and BGP routing changes in other domains. *The operators can use traffic-engineering tools to determine if a planned maintenance activity would trigger changes in egress points, and tune the IGP link weights accordingly*. For example, in Figure 3 the router $A$ selects egress point $B$ with an IGP distance of 10 over egress $C$ with a distance 20. However, if the left link from $A$ needs to be disabled for upgrading, the distance to $B$ would increase to 25, making $C$ the closer egress point. This egress-point change can be avoided by changing the weight of the middle link from 10 to 4 before the maintenance activity; this ensures that the alternate path to $B$ has distance 19—smaller than the distance to $C$.

However, maintenance still causes transient disruptions as the routers converge to the new IGP paths in a distributed fashion. Rather than simply disabling the equipment, the operators can *"prepare" the network for the impending topology change by increasing the link weight to a very high value*. This has two main advantages. First, the incident router learns the IGP change immediately, instead of waiting to detect a link failure. Second, the link can continue to carry packets in flight while the routers converge to new paths. Operators can safely disable the link after the convergence process completes. The inverse procedure can be used to reactivate the link, so that operators can test whether the link works properly before routers start using it to forward data traffic.

### C. Summary

In summary, an ISP can address the four challenges enumerated in Section I by:

1) gaining control over the flow of traffic by using traffic-engineering tools to predict the effects of changes to the IGP topology and routing configuration,
2) avoiding large reactions to small changes by selecting peering locations and IGP weights to minimize the likelihood of hot-potato routing changes,

3) reducing convergence delay by selecting designs that reduce the scope of routing changes and selecting technologies such as tunneling and POS, and

4) preparing for maintenance by carefully "costing out" the equipment beforehand (to reduce convergence delay) and adjusting the routing-protocol configuration (to prevent congestion on the new paths).

## V. Conclusions

This paper discusses the disruptive effects of routing changes. Although routing protocols were designed to adapt quickly to topology and configuration changes, current applications demand even smaller periods of disruption. This requirement is even harder to satisfy after the immense growth of the Internet infrastructure. The guidelines presented in this paper are useful for network designers and operators to reduce routing disruptions using current routing technology. Many of these techniques are simply clever "hacks" to work around a system that was not designed to be managed.

A complete solution requires changes to the routing-protocol implementations (by the router vendors) and enhancements to the protocols themselves (by standards bodies such as the IETF). These enhancements should reduce the need for network operators to follow some of the guidelines discussed in Section IV to deal with routing disruptions. Ultimately, we believe that instead of proposing incremental enhancements to the protocol that fix one aspect of the problem at a time, the IETF and the research community could investigate alternative approaches that give operators more direct and effective control over the selection of paths to meet the demanding performance requirements of Internet applications.

## Acknowledgments

## References

[1] C. Boutremans, G. Iannacconne, and C. Diot, "Impact of Link Failures on VoIP Performance," in *Proc. of NOSSDAV workshop*, ACM Press, May 2002.

[2] S. Halabi and D. McPherson, *Internet Routing Architectures*. Cisco Press, second ed., 2001.

[3] M. Caesar and J. Rexford, "BGP routing policies in ISP networks," *IEEE Network Magazine*, pp. 5–11, November/December 2005.

[4] G. Iannaccone, C.-N. Chuah, S. Bhattacharyya, and C. Diot, "Feasibility of IP Restoration in a Tier-1 Backbone," *IEEE Network Magazine*, March 2004.

[5] R. Teixeira, A. Shaikh, T. Griffin, and J. Rexford, "Dynamics of hot-potato routing in IP networks," in *Proc. ACM SIGMETRICS*, June 2004.

[6] C. Labovitz, A. Ahuja, A. Bose, and F. Jahanian, "Delayed Internet routing convergence," *IEEE/ACM Trans. on Networking*, June 2001.

[7] A. Sridharan, S. B. Moon, and C. Diot, "On the correlation between route dynamics and routing loops," in *Proc. Internet Measurement Conference*, October 2003.

[8] R. Teixeira, T. Griffin, A. Shaikh, and G. Voelker, "Network sensitivity to hot-potato disruptions," in *Proc. ACM SIGCOMM*, August 2004.

[9] J. Rexford, "Route optimization in IP networks," in *Handbook of Optimization in Telecommunications* (P. Pardalos and M. Resende, eds.), Kluwer Academic Publishers, 2005. To appear.