# BGP routing policies in ISP networks

Matthew Caesar
UC Berkeley

Jennifer Rexford
Princeton University

## 1 Introduction

In the early days of the Internet, the problem of how to route packets to their final destination was much simpler than it is today. At the time, the requirements of the Internet's routing protocol were fairly simple, as the Internet was small by today's standards, operated by a single administrative entity (NSFNET), and shortest-path routing was typically used. Over time, as the Internet became more heavily commercialized and privatized, ISPs began to have vested interests in controlling the way traffic flowed for economic and political reasons. The Border Gateway Protocol (BGP) was borne out of the need for ISPs to control route selection (where to forward packets) and propagation (who to export routes to).

When BGP was first introduced, it was a fairly simple path-vector protocol. Over time, many incremental modifications to allow ISPs to control routing were proposed and added to BGP. The end result was a protocol weighted down with a huge number of mechanisms that can overlap and conflict in various unpredictable ways. These modifications can be highly mysterious since many of them, including the decision process used to select routes, are not part of the protocol specification [1]. Moreover, their complexity gives rise to several key problems, including unforeseen security vulnerabilities, widespread misconfiguration, and conflicts between policies at different ISPs.

Addressing BGP's problems is difficult, as changing certain aspects of BGP (for example changing the contents of updates or the way they are propagated) must be coordinated and simultaneously implemented in other ISPs to support the new design. Hence most modifications to the protocol have been made to the *decision process* BGP uses to choose routes. The result is a protocol where most of the complexity is in the decision process and the policies used to influence decisions, while the rest of the protocol remained fairly simple over time. Therefore, in order to understand BGP it is necessary to understand this decision process and the policies of ISPs that gave rise to its design. Understanding policies is also key to solving BGP's problems, understanding measurement data from BGP, or figuring out what to support when developing a new version of BGP.

The range of policies used by operators constitutes a huge space and hence it is impossible to list them all here. Instead, we try to list common goals of network operators and the knobs of BGP that can be used to express policies. In particular, we attempt to isolate certain *design patterns* commonly used by ISPs, the motivations behind them, and how they are implemented in an ISP's network using BGP's mechanisms. We taxonomize policies into three key categories: *business re-lationship* policy (Section 3) arising from economic or political relationships an ISP has with its neighbor, *traffic engineering* policy (Section 4) arising from the need to control traffic flow within an ISP and across peering links to avoid congestion and provide good service quality, and policies for *scalability* (Section 5) to reduce control traffic and avoid overloading routers. We also discuss several interesting avenues of research currently in progress that could improve upon the way BGP policies are handled today (Section 6). We start by giving an overview of BGP routing in the next section.

## 2 BGP routing in a single AS

The Internet consists of thousands of *Autonomous Systems* (ASes)—networks that are each owned and operated by a single institution. BGP is the routing protocol used to exchange reachability information across ASes. Usually each ISP operates one AS, though some ISPs may operate multiple ASes for business reasons (e.g. to provide more autonomy to administrators of an ISP's backbones in the United States and Europe) or historical reasons (a recent merger of two ISPs). Non-ISP businesses (enterprises) may also operate their own ASes so as to gain the additional routing flexibility that arises from participating in the BGP protocol.

Compared to enterprise networks, ISPs usually have more complex policies arising from the fact that they often have several downstream customers, connect to certain customers in multiple geographic locations, have complex traffic engineering goals, and run BGP at internal routers (rather than just border routers as enterprises often do). Although some of the observations we make apply to enterprise networks, our core focus in this paper is on ISP networks. In this section, we describe BGP from the standpoint of a single AS, describing first the protocol that transmits routes from one AS to another, then the decision process used to choose routes, and finally the mechanisms used at routers to implement policy.

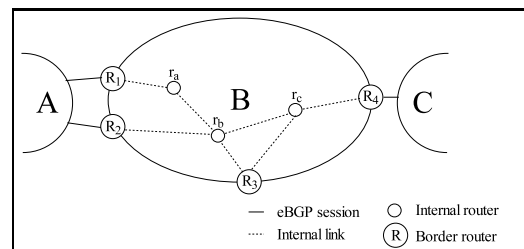### 2.1 Protocol for communicating routing state



Figure 1: Example topology with three ISPs A, B, and C.

Figure 1 shows a simple BGP network. BGP sessions are established between *border routers* that reside at the edges of an AS and border routers in neighboring ASes. These sessions are used to exchange routes between neighboring ASes. Border routers then distribute routes learned on these sessions to non-border (internal) routers as well as other border routers in the same AS using internal-BGP (iBGP). In addition, the routers in an AS usually run an Interior Gateway Protocol (IGP) to learn the internal network topology and compute paths from one router to another. Each router combines the BGP and IGP information to construct a forwarding table that maps each destination prefix to one or more outgoing links along shortest paths through the network to the chosen border router.

BGP is a relatively simple protocol with a few salient features. First, BGP is an *incremental* protocol, where after a complete routing table is exchanged between neighbors, only changes to that information are exchanged. These changes may be new route announcements, route withdrawals, or changes to route attributes. Second, BGP is a *path-vector* protocol where updates contain a list of ASes used to reach the destination. Third, routes are advertised at the *prefix* level, so an AS would advertise a separate update for each of its reachable prefixes. Fourth, BGP announcements contain several fields, including the prefix being updated, the next-hop used to reach the prefix, and a flag indicating whether the route is being advertised or withdrawn. Updates also contain several other route *attributes* that describe various characteristics of the route. An ISP implements its policies by modifying route attributes in updates and changing the way routers react to updates with certain route attributes, as discussed below.

## 2.2 Route selection

Table 1: Steps in the BGP decision process.

| Step | Attribute | Controlled by local or neighbor AS? |
|---|---|---|
| 1. | Highest LocalPref | local |
| 2. | Lowest AS path length | neighbor |
| 3. | Lowest origin type | neither |
| 4. | Lowest MED | neighbor |
| 5. | eBGP-learned over iBGP-learned | neither |
| 6. | Lowest IGP cost to border router | local |
| 7. | Lowest router ID (to break ties) | neither |

A BGP router in an ISP may have several alternate routes to reach a particular destination and must choose between them. A router makes its decision based on the values of attributes contained in the update messages. In the absence of policy, the router would choose the route with the minimum pathlength, with some arbitrary way to break ties between routes with the same pathlength. However, in order to give operators greater control over route selection, several additional attributes were added. The end result was the BGP *decision process*, consisting of an ordered list of attributes across which routes are compared, as shown in Table 1. The router goes down the list, comparing each attribute in the list across the two routes. If the routes have different values for the attribute, the router chooses the one that has the more desirable attribute, otherwise it moves on to compare the next attribute in the list. The route that is chosen is used by the router to forward packets. The ordering of attributes allows the operator to influence various stages of the decision process. For example, the Local Preference (LocalPref) is the first step in the decision process. By changing LocalPref, an operator can force a route with a longer AS path to be chosen over a shorter one. Alternatively, changing an attribute appearing later in the process like the Multi-Exit Discriminator (MED) can allow an operator to change preference between routes with the same AS path length, while ensuring routes with shorter AS paths are always chosen.

There are different ways a route attribute can be set. (a) *Local*: for example, LocalPref is an integer value set at and propagated throughout the local AS and filtered before sending to neighboring ISPs. (b) *Neighbor:* for example, MED is an integer used as a suggestion to a neighboring AS regarding which peering link should be used to reach the local AS. The MED attribute is typically used by two ASes connected by multiple links to indicate which peering link should be used to reach the AS advertising the MED attribute, and is not used to compare routes through two different next-hop ASes. (c) *Neither:* some attributes, for example whether the route was learned through an external BGP (eBGP) neighbor or from an internal router speaking BGP (iBGP), are set by the protocol and cannot be changed.

The collective results of the decision process across routers is to produce a set of *equally good* border routers for each prefix, where each router in the set is equivalent according to the first four steps of the decision process that compare BGP attributes. Each internal router then chooses the router in that set that is closest according to the Interior Gateway Protocol (IGP) path cost to reach that border router. For example in Figure 1, suppose prefix 6.0.1.0/24 is reachable to B via both A and C, but B's LocalPref is set higher for routes through A. The set of equally good border routers would then contain $R_1$ and $R_2$, and each router in B would select the route that was closest exit point (lowest IGP cost): $r_a$ and $R_1$ would choose the route through $R_1$, and all other routers would choose the route through $R_2$.

There are three steps a router applies to route announcements. First *import policy* is applied to determine which routes should be filtered and hence eliminated from consideration. Next, the router applies the *decision process* to select the most desirable route. Finally, an *export policy* is applied which determines which neighbors the chosen route will be exported to. An ISP may implement its policy by controlling any of these three steps, i.e., by modifying import policy to filter routes it doesn't want to use, modifying route attributes to prefer some routes over others, or by modifying export policy to avoid providing routes for certain neighbors to use. In addition, an ISP can modify route attributes of updates it advertises, which can influence how its neighbors perform route selection.

## 2.3 Implementing policy at a router

There are three classes of knobs that can be used to control import and export policies:

1. *Preference* influences which BGP route will be chosen for each destination prefix. Changing preference is done by adding/deleting/modifying route attributes in BGP updates. Table 1 shows which attributes can be modified during import to control preference locally, and which can be modified during export to change how much a neighbor prefers the route.

2. *Filtering* eliminates certain routes from consideration and also controls who they will be exported to. Filtering may be applied both before preference (inbound filtering) or after preference (outbound filtering). Filtering is done by instructing routers to ignore updates with attributes matching certain specified values or ranges.

3. *Tagging* allows an operator to associate additional state with a route, which can be used to coordinate decisions made by a group of routers in an AS, or to share context across AS boundaries. The key mechanism is the *community attribute* [2] [3], a variable length string used to tag routes. The community attribute is a highly expressive mechanism, lending itself to support a wide variety of complex policies that are difficult to express through other means. For example, one community value might affect how the receiving router sets LocalPref, while another might cause the route to be filtered at another router.

An ISP implements its policies by applying configuration commands at routers. These configurations typically consist of a set of lists of preference, filtering, and tagging rules, one list for each *session* the router has with a neighboring BGP-speaking router. Although the configuration language differs between vendors, a key primitive that is often provided is a *route-map*, a language construct used to modify route attributes and define conditions that determine which routes are exported to peers. It consists of two parts: a set of conditions indicating when the map is to be invoked (e.g. the prefix is a specified value, or the AS path matches a specified regular expression), and the action to be taken if the update matches the conditions (e.g. modify a specified attribute, or filter the route).

## 3 Business relationships

ISPs often wish to control next hop selection so as to reflect agreements or relationships they have with their neighbors. Three common relationships ISPs have are: *customer-provider*, where one ISP pays another to forward its traffic, *peer-peer*, where two ISPs agree that connecting directly to each other (typically without exchanging payment) would mutually benefit both, perhaps because roughly equal amounts of traffic flow between their networks, and *backup* relationships, where two ISPs set up a link between them that is to be used only in the event that the primary routes become unavailable due to failure.

There are three key ways these relationships manifest themselves in policy:

**Influencing the decision process (by assigning LocalPrefs):** ISPs often prefer customer-learned routes over routes learned from peers and providers when both are available. This is often done because sending traffic through customers generates revenue for the ISP while sending traffic through providers costs the ISP money and sending to peers can reduce equality of the peering relationship and thereby give incentive to the party receiving more traffic to tear down the relationship or start charging the other party. Often an ISP will achieve this by assigning a non-overlapping range of LocalPref values to each type of peering relationship (for example LocalPref values in the range 90-99 might be used for customers, 80-89 for peers, 70-79 for providers, and 60-69 for backup links). LocalPref can then be varied within each range to do traffic engineering without violating the constraints associated with the business relationship, as described in Section 4). A similar approach can allow an ISP to set aside certain routes as backups. As another example, a large international ISP, due to pricing reasons arising from political boundaries, may wish to use one provider to service customers in one country, but to use a different provider for its customers in another country. This can be done by increasing LocalPref on routes that should be more highly desired.

**Controlling route export (by using the community attribute):** Routes learned from providers or peers are usually not exported to other providers or peers, because there is no economic incentive for an ISP to forward traffic it receives from one provider or peer to another. This can be done by tagging announcements with a community attribute signifying the business relationship of the session, and filtering routes with certain community attributes when exporting routes to peers. For example, suppose B wishes to not export routes learned from A to C as shown in Figure 1, perhaps because it does not get paid for transiting traffic from C to A. It can do this as follows. First, for every session routers $R_1$ and $R_2$ have with routers in A, B configures an import policy that appends the community attribute $X_{peer}$ to any route learned over these sessions, to indicate that the route was received from a peer—information which is ordinarily lost in BGP as the route propagates across the AS. After appending the community attribute, B exports the route onwards into its internal iBGP network. Second, B configures export policies at $R_3$ that match on this community attribute to determine which routes get exported to C. In particular, every session between $R_3$ and a router in C is configured with an export policy that filters any route with the community attribute $X_{peer}$.

**Defensive programming (by filtering routes and attributes):** An AS is vulnerable to the information in BGP announcements sent by neighboring domains. As such, misconfiguration, software bugs, and malicious attacks in other parts of the Internet can have a significant influence on routing in an AS [4]. Depending on the degree of trust an ISP has that its neighbors are properly administered, it may wish to protect itself in certain ways. The *import policy* on a BGP session is typically configured to filter invalid routes. For example, an

ISP may wish to prevent certain neighbors from influencing its choice of routes. It can do this by filtering out certain attributes (like MED and community attributes) from updates it receives from those neighbors. Also, ISPs can perform certain sanity checks on the AS path: for example a Tier-1 ISP should not accept any routes from its customers that contain another Tier-1 ISP in the AS path. In addition, an AS may configure its *export policies* to filter BGP announcements for private networks and other routes that should not be externally reachable (e.g., routes to router's administrative consoles).

# 4 Traffic engineering

While business relationships affect relative preferences for routes, there are often several routes available that are equally preferred. Moreover, ISPs often connect at multiple locations to reduce pathlengths and improve reliability, increasing the number of available routes. A secondary goal for many ISPs is to engineer their traffic by modifying preference within the same business class to meet or maximize certain performance criteria (e.g., achieve desired quality and availability). An ISP can do this by modifying the import policies applied by its routers, each of which can have a different configuration. Common traffic engineering goals include:

**Outbound traffic control (by changing LocalPref and IGP costs):** Operators can influence outbound traffic flow either by configuring import policies that affect which routes get in the set of equally-good border routers, or by modifying IGP link costs. One common goal is *early-exit routing* (also called hot-potato routing), where the ISP forwards the packet to its closest possible exit point, so as to reduce the number of links packets traverse and hence the resulting congestion in its internal network. Although early-exit routing is known to inflate end-to-end path lengths in the Internet, ISPs often exercise early-exit routing to reduce their costs and network congestion, and because BGP does not support alternatives like determining global shortest paths across multiple ISPs.

Another common goal is to reduce congestion on outbound links to neighbors. This can be done by *load balancing* traffic over several links when possible. Outbound traffic engineering can be done by changing LocalPref. For example, suppose B wishes to shift some traffic from its links to A to its link to C as shown in Figure 1, perhaps because the link to A is overutilized or because it is planning to take the link down for maintenance. B can reduce the traffic it sends to A and increase traffic it sends to C by decreasing LocalPref for routes traversing A or increasing LocalPref for routes traversing C.

**Inbound traffic control (by AS prepending and MED):** An ISP cannot rely on its neighbors to perform effective outbound traffic engineering, because its neighbors might not be aware of the ISP's traffic-engineering goals, internal topology, or load on internal links due to privacy reasons. Moreover, an ISP might not be willing to place such a high degree of trust in its neighbors. Hence, some mechanism to allow an ISP to control how much traffic it receives from each of its peering links is essential. Unfortunately, this is a highly challenging problem, as it requires the local ISP to influence route selection in remote ISPs, which in turn might wish to limit or completely ignore local ISP's goals. However, an ISP may convince its neighbor (perhaps through economic incentives) to allow the ISP to control how much traffic it receives on each link from the neighbor. This can be done by modifying the MED attribute. Shifting traffic between links to different neighbors is more challenging, as it requires controlling route selection in ASes multiple hops away, but can be done by prepending multiple copies of its AS number to the AS path in order to artificially inflate the AS-path length.

For example, suppose B wishes to shift some traffic from its link to A to its link to C. B can do this by prepending additional copies of its AS number onto the AS paths in BGP announcements it sends to A. This increases the AS-path length in these updates, which causes routes advertised by C to other ISPs to become more desirable in comparison. An alternate mechanism for controlling inbound traffic is the MED attribute, which can be used between a pair of ISPs connected via multiple peering links. For example, if B wanted to reduce the amount of traffic traversing router $R_1$, it could increase the value of the MED attribute $R_1$ advertises to $A$, causing the link to $R_2$ to become more preferred by A's routers and thereby decreasing $R_1$'s load.

**Remote control (by changing community attributes):** In certain cases, an ISP may need to remotely manage a router's configuration to implement a desired policy. For example, suppose B wishes to have all inbound traffic routed through A. If C has a LocalPref to prefer the direct route to B, no change in MED or AS prepending will force C to use alternate routes through A to B. B could request C to manually change its router configurations, but this can be time consuming for human operators if B changes its policy often (e.g. for traffic engineering purposes). Instead, C can allow B to control C's routing policy with respect to B's routes by configuring its routers to map certain community attributes to certain LocalPref values [2]. If desired, C can limit the degree of B's control to prevent certain policies of its own from being subverted. For example, C can configure its routers to map community value $X_1$ to a LocalPref of 60, and $X_2$ to a LocalPref of 75, allowing B to disable the route, but not allowing B to have it chosen over routes C wants to prefer more (by setting a higher LocalPref, like 85).

Achieving a specific level of load balance (e.g. balancing load to make spare capacity on both links equal) can be very difficult. The key challenge is to select the proper set of prefixes and change attributes for each appropriately: selecting too large a set will cause too much traffic to shift, overloading one of the links. It can also be tedious to express a long list of prefixes in a router configuration file. Some ISPs deal with this by changing preference for all prefixes whose AS path matches a regular expression, then tweaking the regular expression repeatedly to control how many prefixes match it. However, since this approach is done manually it is subject to misconfiguration, cannot be done in real time to adjust to changing load, and the outcome from a change can be difficult to predict. There are automated tools that an ISP can use to predict the effects of these actions [5].

In addition, there are a variety of "smart routing" tools [6] that small ASes at the edge of the Internet can use to balance load over multiple upstream providers. However, these tools generally are not appropriate for ISPs, as dynamically changing traffic can lead to BGP routing changes that are visible to other ASes, which can trigger flap damping (a mechanism that withdraws unstable routes) if the routes become too unstable. Moreover, these tools focus on load balancing over multiple outgoing links but do not consider the effect on traffic flow inside the AS [5].

# 5  Scalability

ISPs wish to protect themselves from instability and excessive routing table growth due to misconfigurations or faults occurring in other ISPs. A properly configured set of BGP policies can provide some defense from these events. Common goals include:

**Limiting routing table size (by filtering)**: ISPs often want to limit routing table size because overflow can cause the router to reboot [7]. This can be a particularly important issue for smaller ISPs which may have less expensive routers with less memory capacity.

1. *Protection from other ISPs:* ISPs can protect themselves from excessive advertisements from neighbors are: (a) filtering long prefixes (e.g., longer than /24) to encourage use of aggregation [8]. (b) As a safety check, routers often maintain a fixed per-session prefix limit that limits the number of prefixes a neighbor can advertise. (c) Default routing: an ISP with a small number of routes may not need the entire routing table, and may instead configure a default route through which most destinations can be reached.

2. *Protecting other ISPs:* An ISP can reduce the number of prefixes it advertises by using *route aggregation*, where instead of advertising two adjacent prefixes (e.g., 4.1.2.0/24 and 4.1.3.0/24) to a neighbor, they can be filtered in the export policy and a less specific prefix (e.g. 4.1.2.0/23) advertised [9]. Doing this effectively may require knowledge of the neighbor's connectivity as illustrated in the following example.
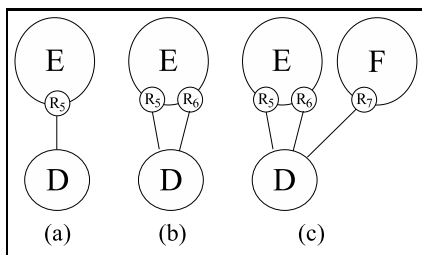


Figure 2: Example topology where adding new customer D triggers E to generate (a) no new advertisements (b) internal advertisement (c) internal and external advertisements.

Suppose E (Figure 2) owns prefix 6.0.0.0/8. E has allocated the subnet 6.1.0.0/16 to router $R_5$, and has allocated smaller subnets to its customers connected to $R_5$, including a new customer D which is allocated subnet 6.1.1.0/24. When adding D as a new customer, E may need to make changes to its routers' configuration, and the configuration it chooses impacts whether new advertisements are generated. There are three cases:

1. *No new advertisements:* Suppose D's sole provider is E, and D connects to just one router $R_5$ in E. In this case, $R_5$ is already announcing 6.1.0.0/16, obviating the need for $R_5$ to announce more specific subnets like 6.1.1.0/24. Hence, E just adds a statically configured route at $R_5$ to forward all traffic in 6.1.1.0/24 to D, and so no advertisements will be sent from E to its neighbors, nor will any new advertisements be sent internally within E.

2. *Internal advertisement:* Suppose instead D connects to two routers $R_5$ and $R_6$ in E. In this case, both $R_5$ and $R_6$ need to advertise the prefix 6.1.1.0/24 within E, so all routers within E know they can reach D via either $R_5$ or $R_6$. However, E can aggregate the advertisement into its address space and hence E will not send BGP updates to its neighbors. This is done by configuring $R_5$ and $R_6$ to tag a community attribute onto advertisements of prefix 6.1.1.0/24, and configuring all border routers to filter routes with that community attribute.

3. *External and internal advertisement:* Suppose D connects to both E and F. In this case E should not aggregate the prefix into its own address space: if it did, then F would then be advertising a longer prefix route to reach D, and since the longest prefix match is always chosen in BGP, all routers in the Internet will prefer F's route over E's route. If D wishes traffic to flow over both links, it must request that E not perform aggregation on its prefix. E can avoid aggregating the prefix by configuring its routers peering with D to append a certain community attribute, and configure its border routers to export routes containing that community attribute.

Although the connectivity of customers clearly influences the way policy is configured, the existence of alternate links is not discovered or signaled by the BGP protocol and hence must be manually detected and accounted for by human operators.

**Limit routing changes (by suppressing routes that flap)**: Routing instability is undesirable, as it can introduce jitter and packet loss in applications like Voice over IP, interfere with TCP's round-trip-time calculations, and increase load on routers thereby potentially reducing their reaction time to routing events. The key mechanism used to improve routing stability is *flap damping*. Flap damping is a mechanism that limits propagation of unstable routes. It works by maintaining a penalty value associated with the route that is incremented whenever an update is received. When the penalty value surpasses a configurable threshold, the route is *suppressed* for some time, i.e., it is made unavailable to the decision process and hence will not be selected. An ISP can lower the

penalty threshold to improve route stability at the cost of worsening availability. ISPs sometimes more aggressively dampen longer prefixes than shorter prefixes, with the motivation that damping a shorter prefix can have a large effect on reachability, and sometimes disable damping for certain important networks (e.g., BGP routes to the root Domain Name System servers) [10].

# 6 Looking forward

BGP's rich feature set of tunable knobs and complex cross-protocol interactions make it highly subject to a variety of problems, including misconfiguration, oscillations, and protocol divergence. The challenge of supporting many different complex policies in BGP without significantly complicating the protocol or degrading its performance has led to much research activity. Four key areas of research related to BGP policy are:

**Configuration checking:** Due to the complexity of Internet routing it can be difficult to predict the way policies interact and configuration mistakes can become prevalent. Interdependence of policy across ISPs and within a single ISP can trigger problems like persistent route oscillations. Configuration checking tools can avoid misconfigurations by verifying certain consistency criteria hold [11], and modeling tools can predict side-effects of configuration changes on routers within an ISP [5]. Across ISPs, uncoordinated routing policy can worsen route convergence and stability. The Routing Arbiter [12] project introduced a distributed architecture for publishing and coordinating routing policies so as to avoid these problems. Other work has attempted to coordinate route policy selection across ISPs without revealing private details of policies [13].

**Detecting policy violations** An ISP may wish to verify that its neighbors are not trying to subvert its routing decisions for their own gain. For example, ISPs may check that their neighbors send consistent route announcements across all peering points [14]. Otherwise, a neighbor selectively advertising a prefix at certain peering points could force the ISP to carry the traffic a longer distance across its own network. (Consistent export is explicitly required in some peering contracts [15, 16] to prevent such selfish activity. However, detecting violations is challenging as it requires observing route announcements at multiple locations in the ISP network.) ISPs may also detect policy subversion by monitoring inbound traffic rates. For example, a customer with two providers might wish to balance load equally over both providers, but one provider may begin advertising a more attractive route (e.g., with a shortest AS path or more-specific prefix) to attract more traffic (and, hence, more revenue) to its connection.

**Language design:** Routing Policy Specification Language (RPSL) [17] is a vendor-neutral language proposed to describe an ISP's policy. It was envisioned these descriptions could be bound together in a database and checked for consistency [12]. RPSL, though mature, is somewhat low-level and mechanism oriented. It may be possible to substantially improve upon RPSL by designing router configuration languages with higher level constructs that allow diverse policies while precluding certain misconfigurations, enforcing certain consistency properties to hold, simplifying configuration of certain common design patterns [18], however the design of such a language remains an open problem.

**New architectures:** There are several routing architectures aimed at fixing problems in and extending functionality of BGP. HLP [19] is a proposed replacement for eBGP. The design philosophy of HLP is to expose common policies that can typically be inferred in BGP today and optimize the routing protocol based on the resulting structure, with the aim to improve scalability and convergence of interdomain routes. Routing Control Platform (RCP) [20] is a logically centralized system that computes and distributes routes to routers inside an ISP. The centralization allows policies to be applied at the AS level, while the RCP takes care of decomposing the policy into router-level constructs. This simplifies the configuration and application of policies and avoids misconfiguration.

# 7 Conclusion

Although BGP policies can be highly complex, there are a number of common design patterns that are typically used by ISPs. In this article we discussed several common patterns and how they can be realized using BGP policy mechanisms. We believe that by recognizing these patterns exist we can more efficiently develop tools that directly support them, such as languages that preclude errors, or analysis tools that check correctness, or architectures that are designed for common cases.

# References

[1] Y. Rekhter, T. Li, "A border gateway protocol 4," IETF RFC 1771, March 1995.

[2] E. Chen, T. Bates, "An application of the BGP community attribute," IETF RFC 1998, August 1996.

[3] B. Quoitin, O. Bonaventure, "A survey of the utilization of the BGP community attribute," expired Internet Draft *draft-quoitin-bgp-comm-survey-00.txt*, February 2002.

[4] O. Nordstrom, C. Dovrolis, "Beware of BGP attacks," in *ACM SIGCOMM CCR*, April 2004.

[5] N. Feamster, J. Winick, J. Rexford, "A model of BGP routing for network engineering," in ACM SIGMETRICS, New York, NY, June 2004.

[6] S. Hares, J. Johnson, A. Britt, R. Bays, M. Lloyd, D. Golding, B. Ross, "Smart routing technologies," in *NANOG 25*, June 2002, http://www.nanog.org/mtg-0206/smart.html.

[7] L. Wang, X. Zhao, D. Pei, R. Bush, D. Massey, A. Mankin, S. Wu, L. Zhang, "Observation and analysis of BGP behavior under stress," in IMW, November 2002.

[8] S. Bellovin, R. Bush, T. Griffin, J. Rexford, "Slowing routing table growth by filtering based on address allocation policies," unpublished, March 2005, http://www.cs.princeton.edu/~jrex/papers/filter.pdf

[9] E. Chen, J. Stewart, "A framework for inter-domain route aggregation," IETF RFC 2519, February 1999.

[10] "RIPE routing-WG recommendations for coordinated flap damping parameters," October 2001, http://www.ripe.net/ripe/docs/routeflap-damping.html

[11] N. Feamster, H. Balakrishnan, "Detecting BGP configuration faults with static analysis," in *NSDI*, May 2005.

[12] R. Govindan, C. Alaettinoglu, G. Eddy, D. Kessens, S. Kumar, W. Lee, "An Architecture for Stable, Analyzable Internet Routing," IEEE Network Magazine, Jan-Feb 1999.

[13] R. Mahajan, D. Wetherall, T. Anderson, "Negotiation based routing between neighboring domains," in *NSDI*, May 2005.

[14] N. Feamster, Z. Mao, J. Rexford, "BorderGuard: detecting cold potatoes from peers," in *IMC 2004*, October 2004.

[15] America Online, "Settlement-free interconnection policy," web site `http://www.atdn.net/settlement_free_int.shtml`

[16] UUNET, "MCI policy for settlement-free interconnection," web site `http://global.mci.com/uunet/peering/`

[17] A. Alaettinoglu, C. Villamizar, E. Gerich, D. Kessens, D. Meyer, T. Bates, D. Karrenberg, M. Terpstra, "Routing policy specification language (RPSL)," IETF RFC 2622, June 1999.

[18] T. Griffin, A. Jaggard, V. Ramachandran, "Design principles of policy languages for path vector protocols," in ACM SIGCOMM, Karlsruhe, Germany, August 2003.

[19] L. Subramanian, M. Caesar, C. Ee, M. Handley, Z. Mao, S. Shenker, I. Stoica, "Towards a Next Generation Inter-domain Routing Protocol," in *HotNets-III*, November 2004

[20] M. Caesar, D. Caldwell, N. Feamster, J. Rexford, A. Shaikh, J. van der Merwe, "Design and implementation of a routing control platform," in *NSDI*, May 2005.