

Interdomain traffic engineering with BGP

B. Quoitin, S. Uhlig, C. Pelsser, L. Swinnen and O. Bonaventure

Abstract

Traffic engineering is performed by means of a set of techniques that can be used to better control the flow of packets inside an IP network. We discuss the utilization of these techniques across interdomain boundaries in the global Internet. We first analyze the characteristics of interdomain traffic on the basis of measurements from three different Internet Service Providers and show that a small number of sources are responsible for a large fraction of the traffic. Across interdomain boundaries, traffic engineering relies on a careful tuning of the route advertisements sent via the Border Gateway Protocol (BGP). We explain how this tuning can be used to control the flow of the incoming and of the outgoing traffic and identify their limitations.

I. INTRODUCTION

Initially developed as a network that connects a small number of research networks, the Internet has become a world-wide data network that is used for mission critical applications. Supporting such mission critical applications across the global Internet implies several important challenges. The first challenge is the size of the Internet. The Internet is a large decentralized network that connected about 160 million hosts in June 2002. Furthermore, these hosts are divided in about 13000 distinct domains managed by distinct companies. All these domains are interconnected to form the global Internet. The Border Gateway Protocol (BGP) is used to route the IP packets that are exchanged between domains. There are basically two types of domains. The *stub domains* contain hosts that produce or consume IP packets. These domains do not carry IP packets that are not produced by or destined to their hosts. The *transit domains* interconnect different domains together and carry IP packets that are produced by and/or destined to external domains. Additional details on the relationships between domains may be found in [SARK02].

The second challenge is that the research Internet, was designed with a best-effort service in mind where connectivity was the most important issue. Today, connectivity is considered to be granted and the best-effort service is used for mission critical applications with stringent Service Level Agreements (SLA). To meet these SLAs, several Internet Service Providers (ISP) rely on traffic engineering [ACE⁺02] to better control the flow of IP packets. Large ISPs often need to engineer the flow of packets inside their own domain to reduce congestion by better distributing the traffic on all their links. Several techniques have been developed during the last few years, some require the utilization of Multi-Protocol Label Switching (MPLS) to forward IP packets while others only require a tuning of the traditional IP routing protocols used inside the ISP network. Besides optimizing the flow of packets inside their network, most ISPs also need to better control the flow of their interdomain traffic, i.e. the IP packets that cross the boundaries between distinct ISPs. Today, MPLS is not used across interdomain boundaries and the only solution to engineer the flow of interdomain traffic is to tune the configuration of the BGP routing protocol. This tuning is often done on a trial-and-error basis and suffers from limitations as will be shown in the rest of this article.

In this article, we first introduce the operation of the BGP protocol in section II. We provide recent results about the characteristics of interdomain traffic in section III. Finally, we describe in details several interdomain traffic engineering techniques in section IV and show their limitations.

II. BGP ROUTING IN THE INTERNET

Internet routing is handled by two distinct protocols with different objectives. Inside a single domain, link-state intradomain routing protocols distribute the entire network topology to all routers and select the shortest path according to a metric chosen by the network administrator. Across interdomain boundaries, the interdomain routing protocol is used to distribute reachability information and to select the best route to each destination according to the policies specified by each domain administrator. For scalability reasons, the interdomain routing protocol is only aware of the interconnections between distinct domains, it does not know any information about the content of each domain.

The Border Gateway Protocol (BGP) [RL02]¹ is the current de facto standard interdomain routing protocol. In BGP terminology, a domain is called an Autonomous System (AS). BGP is a *path-vector protocol* that works by sending *route advertisements*. A route advertisement indicates the reachability of a network (i.e. a network address and a netmask representing a block of contiguous IP addresses - for instance, 192.168.0.0/24 represents a block of 256 addresses between 192.168.0.0 and 192.168.0.255) because this network belongs to the same AS as the advertising router or because a route advertisement for this network was received from another AS. Besides the reachable network and the IP address of the router that must be used to reach this network (known as the *next-hop*), a route advertisement also contains the *AS-path* which is the list of all the transit AS that must be used to reach the announced network. The length of the *AS-path* can be considered as the route metric. A route advertisement may also contain several other attributes such as *local-pref*, *med* and *communities*. An important point to note about BGP is that if a BGP router of AS_x sends a route announcement for network *N* to a neighbor BGP router of AS_y, this implies that AS_x accepts to forward the IP packets to destination *N* on behalf of AS_y.

The authors are with the Infonet group, University of Namur, Belgium, <http://www.infonet.fundp.ac.be>, Email : {bqu,suh,lsw,cpe,lsw,obo}@infonet.fundp.ac.be

¹A more readable presentation of BGP may be found in [Ste99].

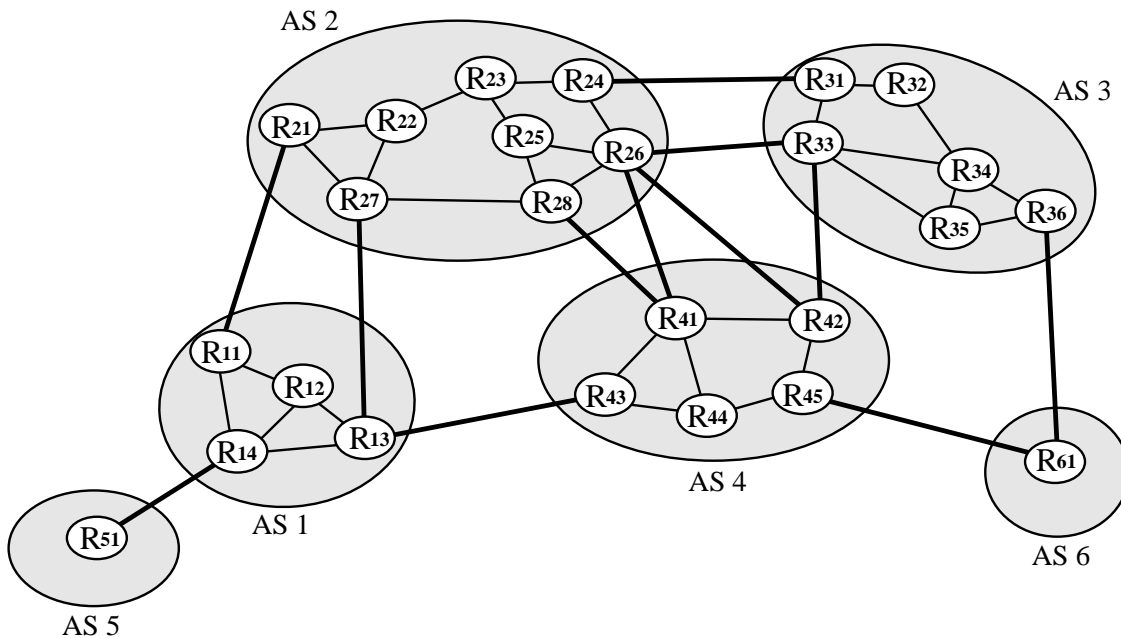


Fig. 1. A simple Internet

There are two variants of BGP. The eBGP variant is used to announce the reachable prefixes on a link between routers that are part of distinct AS (e.g. R_{51} and R_{14} in figure 1). The iBGP variant is used to distribute inside an AS the best routes learned from neighboring AS. For this, a BGP router inside an AS will establish an iBGP session with all the other BGP routers of the same AS. This will create a full-mesh of iBGP sessions inside the AS. For example, inside AS1 in figure 1, there will be a full mesh of iBGP sessions involving at least routers R_{11} , R_{13} and R_{14} . These iBGP sessions will be used for example by router R_{14} to announce to the other BGP routers of the AS the route advertisements received from AS5.

Inside a single domain, all routers are considered as “equal” and the intradomain routing protocol announces all known paths to all routers. In contrast, in the global Internet, all AS are not equal and an AS will rarely agree to provide a transit service for all its connected AS toward all destinations. Therefore, BGP allows a router to be selective in the route advertisements that it sends to neighbor eBGP routers. To better understand the operation of BGP, it is useful to consider a simplified view of a BGP router as shown in figure 2.

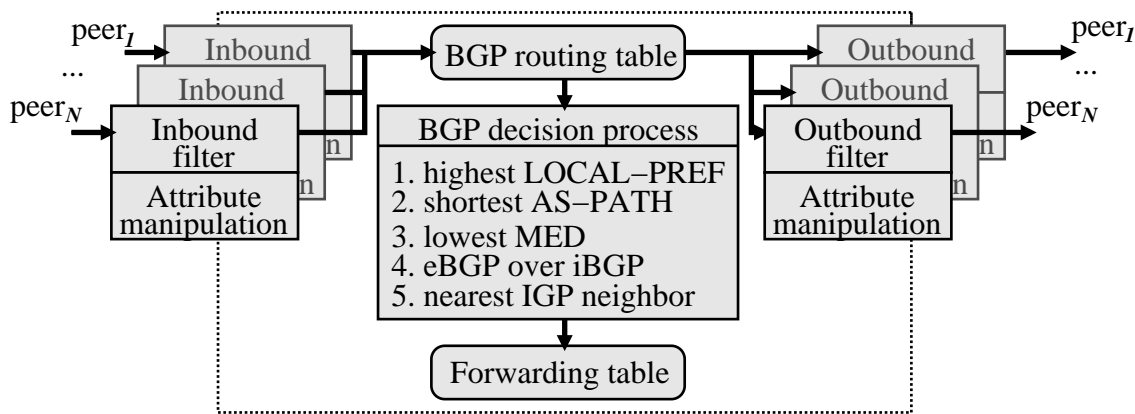


Fig. 2. Operation of a BGP router.

A BGP router processes and generates route advertisements as follows. First, the administrator specifies, for each BGP peer, an input filter (figure 2, left) that is used to select the acceptable advertisements. For example, a BGP router could only select the advertisements with an `AS-Path` containing a set of trusted AS. Once a route advertisement has been accepted by the input filter, it is placed in the BGP routing table, possibly after having updated some of its attributes. The BGP routing table thus contains all the acceptable routes received from the BGP neighbors.

Second, on the basis of the BGP routing table, the BGP decision process (figure 2, center) will select the best route toward each known network. Based on the `next-hop` of this best route and on the intradomain routing table, the router will install a route toward this network inside its forwarding table. This table is then looked up for each received packet and indicates the outgoing interface that must be used to reach the packets' destination.

Third, the BGP router will use its output filters (figure 2, right) to select among the best routes in the BGP routing table the routes that will be advertised to each BGP peer. At most one route will be advertised for each reachable destination. The BGP router will assemble and send the corresponding route advertisement messages after a possible update of some of their attributes.

The input and output filters used in combination with the BGP decision process are the key mechanisms that allow a network administrator to support within BGP the business relationships between two AS. Many types of business relationships can be supported by BGP. Two of the most common relationships are the customer-to-provider and the peer-to-peer relationships [SARK02]. To understand how these two relationships are supported by BGP, consider figure 1. If AS5 is AS1's customer, then AS5 will configure its BGP router to announce its routes to AS1. AS1 will accept these routes and announce them to its peer (AS4) and upstream provider (AS2). AS1 will also announce to AS5 all the routes it receives from AS2 and AS4. If AS1 and AS4 have a peer-to-peer relationship on the link between R_{13} and R_{43} , then router R_{13} will only announce on this link the internal routes of AS1 and the routes received from AS1's customers (i.e. AS5). The routes received from AS2 will not be announced on the $R_{13} - R_{43}$ link by router R_{13} and thus AS1 will not carry traffic from AS4 toward AS2.

III. CHARACTERISTICS OF INTERDOMAIN TRAFFIC

An important element to consider when engineering the interdomain traffic of an AS are the characteristics of this traffic. Informal discussions with network operators on this topic indicate that often a small number of AS are responsible for a large fraction of the total traffic received or sent by a given ISP. However, while there are many studies on the topology of the Internet [SARK02] or the evolution of the BGP routing tables [BNC02] as well as many studies on the packet-level characteristics of the traffic [TMW97], few papers analyze together the traffic and its topological distribution [FP99], [FBR02], [UB02].

In the framework of a detailed analysis of interdomain traffic, we have collected several traces of all the traffic received or sent through the border routers of three stub ISPs. The first trace was collected during one week in December 2000 and covers all the traffic received by BELNET, the Belgian Internet provider for universities and research labs. During this period, BELNET received 2.1 terabytes of data from 4243 distinct AS. The second trace was collected during five consecutive days in April 2001 at the border routers of YUCOM, an ISP based in Belgium that provides dialup access to home users. During this period, it received IP packets corresponding to 1.1 terabytes of data from 7669 distinct AS. The last trace was collected during one day at the border routers of the Pittsburgh Super-computing Center (PSC) in March 2002. PSC provides access to Internet and Internet2 for organizations in western Pennsylvania. During the studied day, the border routers of PSC sent IP packets corresponding to 574 gigabytes of data to 11791 distinct AS. The difference in the number of AS for each studied ISP is mainly due to the number of AS in the Internet at the time of the measurement. In December 2000, BELNET had 6298 distinct AS in its routing table, while YUCOM knew 10560 AS in May 2001 and PSC knew almost 12000 AS in March 2002.

The analysis of those traces together with the routing tables of the studied ISPs provides useful information on the efficiency of interdomain traffic engineering. In figure 3, we show the cumulative distribution of the interdomain traffic received by BELNET and YUCOM and sent by PSC. We classified the traffic on the basis of its source and destination AS. A similar analysis could have been done at the prefix level (see [UB02] for BELNET). A first point to note about this figure is that the studied ISPs exchanged packets with most of the AS that compose the Internet during the studied periods. However, the studied ISPs do not exchange the same amount of traffic with each remote AS. The 10 (resp. 100) largest sources of traffic for YUCOM contribute to more than 30% (resp. 72%) of the traffic received by this ISP. Similarly, the 10 (resp. 100) largest sources of traffic for BELNET contribute to 22% (resp. 64 %) of the traffic it receives during one week. For PSC, the concentration of the traffic sinks is even more important as the 10 (resp. 100) largest destinations receive 38% (resp. 78%) of the total traffic sent by PSC. This concentration of the traffic within a relatively small number of sources means that interdomain traffic engineering does not need to influence all interdomain routes. [FBR02] mentions a similar distribution of the traffic for a large *tier-1* ISP.

Another important point to mention about the interdomain traffic exchanged by the studied ISPs is the distance (measured in AS hops) between the remote AS and each studied ISP. Figure 4 shows that the studied ISPs only exchange a small fraction of their traffic with their direct peers. Most of the packets are exchanged with AS that are only a few AS hops away. For the BELNET trace, most of the traffic is produced by sources located 3 and 4 AS hops away while YUCOM mainly receives traffic from sources that are 2 and 3 AS hops away. PSC on the other hand sends traffic to sources located at up to 4 AS hops away.

The results of this analysis show that interdomain traffic engineering could be achieved by influencing the routes to a few tens of AS that are located only a few AS hops away. A more detailed discussion about the implications of interdomain traffic characteristics on traffic engineering based on the study of a single AS may be found in [UB02].

IV. INTERDOMAIN TRAFFIC ENGINEERING

Interdomain traffic engineering requirements are diverse and often motivated by the need to balance the traffic on links with other AS and to reduce the cost of carrying traffic carried on these links. These requirements depend on the connectivity of an AS with others but also on the type of business handled by this AS.

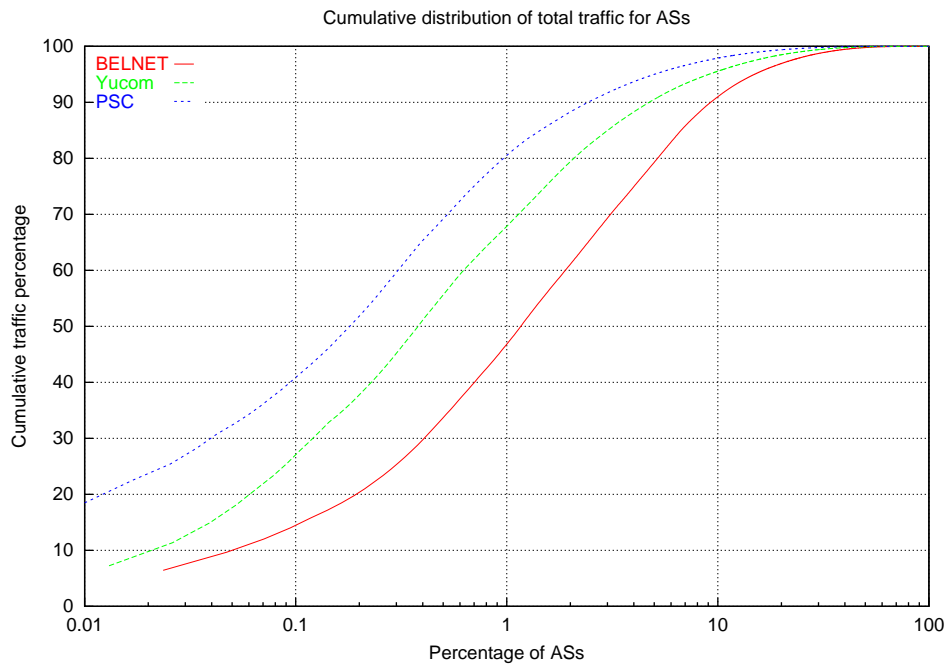


Fig. 3. Cumulative distribution of the traffic

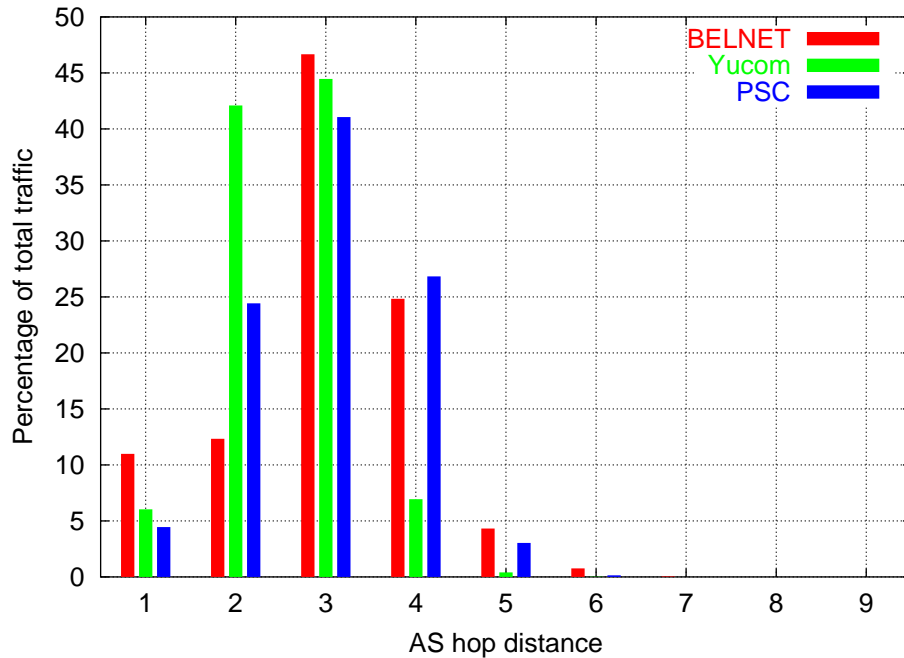


Fig. 4. Per-AS hop distribution of the traffic

The connectivity between AS is mainly composed of two types of relationships. The most frequent relationship between AS is the *customer-provider* relationship where a customer AS pays to use a link connected to its provider. This relationship is the origin of most of the interdomain cost of an AS. Each stub AS tries to maintain at least two of these links for performance and redundancy reasons [SARK02]. In addition, larger AS typically tries to obtain *peer-to-peer* relationships with other AS and then share the cost of the link with the other AS. Negotiating the establishment of those *peer-to-peer* relationships is often a complicated process since technical and economical factors, as exposed in [Bar00], need to be taken into account.

Moreover, an AS will want to optimize the way traffic enters or leaves its network, based on its business interests. Content-providers that host a lot of web or streaming servers and usually have several customer-provider relationships with larger AS will try to optimize the way traffic leaves their networks. On the contrary, access-providers that serve small and medium enterprises,

dialup or xDSL users typically wish to optimize how Internet traffic enters their networks. And finally, a transit AS will try to balance the traffic on the multiple links it has with its peers. Since content-providers, access-providers and transit AS have different traffic engineering requirements, the technique implemented may diminish the efficiency of the ones used by other AS.

Optimizing the way traffic enters or leaves a network means to favor one link over another to reach a given destination or to receive traffic from a given source. This type of interdomain traffic engineering can be performed by tweaking the BGP routers of the AS. In order to understand how BGP can be used to control the way traffic enters, leaves or crosses an AS, a better understanding of the BGP decision process is required. A BGP router receives one route toward each destination from each of its peers. To select the best route among this set of routes, a BGP router relies on a set of criteria called the decision process. Most BGP routers apply a decision process similar in principle to the one shown in figure 2. The set of routes with the same destination are analyzed by the criteria in the sequence indicated in figure 2. These criteria act as filters and the N^{th} criterion is only evaluated if more than one route has passed the $N - 1^{th}$ criterion. It should be noted that most BGP implementations allow the network administrator to optionally disable some of the criteria of the BGP decision process.

A. Control of the outgoing traffic

To control how the traffic leaves its network an AS must be able to choose which route will be used to reach a particular destination through its peers. To influence the selection of the best routes that BGP will retain in its routing table, an AS may rely on the setting of the `local-pref` attribute of routes which is the first to be compared by the decision process, as shown in figure 2. This attribute can be used to give different preferences to routes received by an AS. This preference is propagated to BGP speakers inside the AS (across iBGP sessions). In this way, preferences attached by different BGP routers to routes received from eBGP sessions can be compared during the BGP decision process inside the domain. The route with the highest `local-pref` gets preference.

For instance, an AS with multiple providers may receive different routes toward the same prefix. Each BGP router inside the AS can then carefully set the `local-pref` for each route in order to achieve some kind of load balancing. Indeed since the AS may measure the traffic leaving its network toward each prefix and try to balance the load over its upstream links with the help of local preferences.

Recently, a few commercial companies have devised new solutions under the name “smart routing” [Bor02]. These proprietary solutions address the requirements of content-providers and small AS that are multi-homed and want to enhance the performance of their outgoing traffic. On the basis of partial knowledge of the Internet topology and on measurement of the performance of upstream routes, they attach appropriate values of the `local-pref` attribute to the best routes.

B. Control of the incoming traffic

The first method that can be used to control the traffic that enters an AS is to rely on selective advertisements and announce different route advertisements on different links. For example in figure 1, if AS1 wanted to balance the traffic coming from AS2 over the links $R_{11} - R_{21}$ and $R_{13} - R_{27}$, then it could announce only its internal routes on the $R_{11} - R_{21}$ link and only the routes learned from AS5 on the $R_{13} - R_{27}$ link. Since AS2 would only learn about AS5 through router R_{27} , it would be forced to send the packets whose destination belongs to AS5 via router R_{27} . However, a drawback of this solution is that if the link $R_{13} - R_{27}$ fails, then AS2 would not be able to reach AS5 through AS1. This is not desirable and it should be possible to utilize link $R_{11} - R_{21}$ for the packets toward AS5 at that time without being forced to change the routes that are advertised on this link.

A variant of the selective advertisements is the advertisement of more specific prefixes. This advertisement relies on the fact that an IP router will always select in its forwarding table the most specific route for each packet (i.e. the matching route with the longest prefix). For example, if a forwarding table contains both a route toward $16.0.0.0/8$ and a route toward $16.1.2.0/24$, then a packet whose destination is $16.1.2.200$ would be forwarded along the second route. This fact can also be used to control the incoming traffic. In the following example, we assume that prefix $16.0.0.0/8$ belongs to AS3 and that several important servers are part of the $16.1.2.0/24$ subnet. If AS3 prefers to receive the packets toward its servers on the $R_{24}-R_{31}$ link, then it would advertise both $16.0.0.0/8$ and $16.1.2.0/24$ on this link and only $16.0.0.0/8$ on its other external links. An advantage of this solution is that if link $R_{24}-R_{31}$ fails, then subnet $16.1.2.0/24$ would still be reachable through the other links.

Another method would be to allow an AS to indicate a ranking among the various route advertisements that it sends. Based on the utilization of the length of the `AS-Path` as the third criteria in the BGP decision process, a possible way to influence the selection of routes by a distant AS is to artificially increase the length of the `AS-Path` attribute. Coming back to figure 1, assume that AS3's primary interdomain link uses link $R_{61} - R_{45}$ while link $R_{61} - R_{36}$ is only used as backup primary link. In this case, AS6 would announce its routes normally on the primary link (i.e. with an `AS-Path` of AS6) but would attach add its own AS number several times instead of once in the `AS-Path` attribute (e.g. AS6 AS6 AS6) on the $R_{61} - R_{36}$ link. The route advertised on the primary link will be considered as the best route by all routers that do not rely on manually configured settings for the weight and `local-pref` attributes. This technique can be combined with selective advertisements. For example, an AS could divide its address space in two prefixes $p1$ and $p2$ and advertise prefix $p1$ without prepending and prefix $p2$ with prepending on its first link and the opposite of its second link.

The last method to allow an AS to control its incoming traffic is to rely on the `multi-exit-discriminator` attribute also known as the MED. This attribute can only be used by an AS multi-connected to another AS to influence the link that should be used

by the other AS to send packets toward a specific destination. It should however be noted that the utilization of the MED attribute is usually subject to a negotiation between the two peering AS and some AS do not accept to take the MED attribute into account in their decision process.

C. Community-based Traffic Engineering

In addition to these techniques, several ISPs have been using the `communities` attribute to give their customers a finer control on the redistribution of their routes. The `communities` attribute is an optional attribute that can be attached to routes. This attribute can contain several 32 bits wide community values. Community values are often used to attach optional information to routes such as a code representing the city where the route was received or a code indicating whether the route was received from a peer or a customer. The community values can also be used for traffic engineering purposes. In this case, predefined community values can be attached to routes in order to request actions such as not announcing the route to a specified set of peers, prepending the `as-path` when announcing the route to a specified set of peers or setting the `local-pref`. However, this technique relies on an ad hoc definition of community values and on manual configurations of BGP filters which makes it difficult to use and subject to errors.

The IETF is currently considering the definition of a new standard type of extended communities that are called “redistribution communities” [BCH⁺02] to solve the drawbacks of the utilization of classical `communities` to do traffic engineering. These redistribution communities can be attached to routes to influence the redistribution of those routes by the upstream AS. The redistribution communities attached to a route contain both the traffic engineering action to be performed and the BGP peers that are affected by this action. One of the supported actions allows an AS to indicate to its upstream peer that it should not announce the attached route to some of its BGP peers.

Another type of action allows an AS to request its upstream to perform `AS-Path` prepending when redistributing a route to a specified peer. To understand the usefulness of such redistribution communities, let us consider again figure 1, and assume that AS6 receives a lot of traffic from AS1 and AS2 and that it would like to receive the packets from AS1 (resp. AS2) on the R_{45} - R_{61} (resp. R_{36} - R_{61}) link. AS6 cannot achieve such a traffic distribution by performing `AS-Path` prepending itself. However, this becomes possible with the redistribution communities by requesting AS4 to perform the prepending when announcing the AS6 routes to external peers. AS6 could thus advertise to AS4 its routes with a redistribution community that indicates that this route should be prepended two times when announced to AS2. With this redistribution communities, AS4 would advertise path `AS4:AS4:AS6` to AS2 and path `AS4:AS6` to AS1. AS2 would thus receive two routes toward AS6: `AS4:AS4:AS6` and `AS3:AS6` and would select the route via AS3. AS1 on the other hand would select the `AS4:AS6` route which is shorter than the `AS2:AS3:AS6` route.

D. Discussion

The sections above have described several techniques that can be used by ISPs to engineer their interdomain traffic. However, there are some limitations to be considered when deploying those techniques.

A first point to note is that the control of the outgoing traffic with BGP is based on the selection of the best route among the available routes. This selection can be performed on the basis of various parameters, but it is limited by the diversity of routes received from upstream providers which depends on the connectivity and the policy of these AS.

The control of the incoming traffic is based on a careful tuning of the advertisements sent by an AS. This tuning can cause several problems. First, an AS that advertises more specific prefixes or has divided its address space in distinct prefixes to announce them selectively will advertise a larger number of prefixes than required. All these prefixes will be propagated throughout the global Internet and will increase the size of the BGP routing tables of potentially all AS in the Internet. [BNC02] reports that more specific routes constitute more than half of the entries in a BGP table. Faced with this increase of their BGP routing tables, several large ISPs have started to install filters to ignore the BGP advertisements corresponding to more specific prefixes. The deployment of those filters implies that the more specific prefixes will not be announced by those large ISPs and thus the technique will become much less effective.

When considering the manipulation of the `AS-Path` attribute, we have mentioned that it can be used on backup links. It is sometimes also used to better balance the traffic ([BNC02] reports that `AS-Path` prepending affected 6.5 % of the BGP routes in November 2001). However, figure 4 has shown that the AS that produce most of the traffic are located at only a few AS hops away. This implies that only a few values of prepending are possible and thus the granularity of this technique is limited.

The redistribution communities can provide a finer granularity than `AS-Path` prepending or selective announcements. In practice, it can be expected that those communities will be used to influence the redistribution of routes toward large transit ISPs with a large number of customers. For example, consider as an example YUCOM discussed in section III. This ISP has two major upstream providers that allow it to reach the entire Internet. These two providers are then each connected to several tier-1 ISPs that provide most of their connectivity. Figure 5 provides a closer look at the Internet topology as seen by YUCOM on the basis of the BGP advertisements that it received from its two providers. In this figure, we consider the three largest tier-1 ISPs that were connected to YUCOM’s providers. Based on the BGP advertisements, it appears that one of these ISPs was connected to both providers while the two others were only connected to a single provider of YUCOM. In addition to this topological information, figure 5 also reports the number of distinct AS announced by each tier-1 ISP on their links with the two providers of YUCOM.

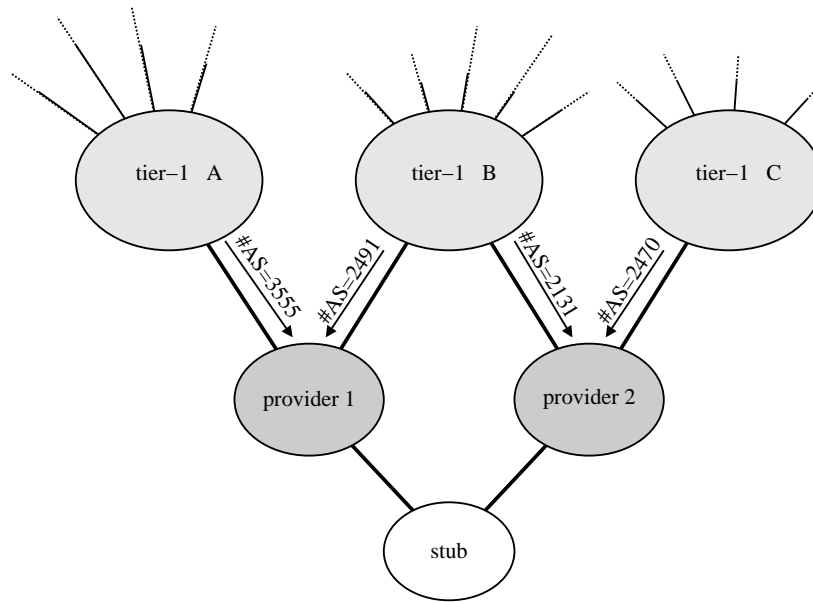


Fig. 5. Subset of the interdomain topology seen from the studied ISP and number of different AS advertised by each tier-1 ISP

Figure 5 reveals two interesting informations. First, each *tier-1* ISP provides connectivity and thus announces routes toward a large number of AS. In total, the three largest *tier-1* providers announce more than 8500 AS. Second, the studied ISP learns routes toward more than 2000 different AS reachable via *tier-1 B* via its two upstream providers. By using redistribution communities targeted at those large *tier-1* ISPs, our ISP could influence the redistribution of its routes to a large number of AS with only a few communities. For example, the studied ISP could utilize a single redistribution community to request its first upstream provider to announce its local routes with *AS-Path* prepending only toward *tier-1 B*. The result of this modified advertisement by the first provider will be that the traffic coming from AS attached to *tier-1 B* would be received through the other provider.

A last point to note concerning the techniques that require changes to the attributes of BGP advertisements is that any (small) change to an attribute will force the route advertisement to be redistributed to potentially the entire Internet. Although it would be possible to define techniques relying on measurements to dynamically change the BGP advertisements of an AS for traffic engineering purposes, a widespread deployment of such techniques would increase the number of BGP messages exchanged and could lead to BGP instabilities. Any dynamic interdomain traffic engineering technique that involves frequent changes to the values of BGP attributes should be studied carefully before being deployed.

V. CONCLUSION

In this paper, we have described several techniques that are used today to control the flow of packets in the global Internet. We have first described the current organization of the Internet and the key role played by BGP.

We have also discussed the characteristics of interdomain traffic based on long traces collected by three distinct ISPs. Two common characteristics appeared in those three ISPs. First, although the Internet is composed of about 13000 AS today, a small percentage of those AS contribute to a large fraction of the traffic received or sent by those ISPs. Second, those highly active sources of destinations are located only a few AS hops away, although the adjacent AS are only responsible for a small fraction of the total traffic.

We have explained how BGP is tuned today for interdomain traffic engineering purposes. We have shown that an AS has more control on its outgoing than on its incoming traffic. Several techniques can be used to control the incoming traffic, but they have limitations. The selective advertisements and the more specific prefixes have the drawback of increasing the size of the BGP routing tables. With *AS-Path* prepending, it can be difficult to select the appropriate value of prepending to achieve a given goal. Finally, we have shown how the redistribution communities could allow an AS to flexibly influence the redistribution of its routes toward non-directly connected *tier-1* ISPs.

ACKNOWLEDGMENTS

This work was supported by the European Commission within the IST ATRIUM project. We would like to thank C. Rapier (PSC), B. Piret (YUCOM) and M. Roger (BELNET) for their traffic traces. We also thank S. De Cnodder, C. Filsfils, A. Danthine and the anonymous reviewers for their useful comments.

REFERENCES

- [ACE⁺02] D. Awduche, A. Chiu, A. Elwalid, I. Widjaja, and X. Xiao. Overview and principles of internet traffic engineering. Internet Engineering Task Force, RFC3272, May 2002.
- [Bar00] S. Bartholomew. The art of peering. *BT Technology Journal*, 18(3), July 2000.
- [BCH⁺02] O. Bonaventure, S. De Cnodder, J. Haas, B. Quoitin, and R. White. Controlling the redistribution of bgp routes. Internet draft, draft-ietf-ptomaine-bgp-redistribution-01.txt, work in progress, August 2002.
- [BNC02] A. Broido, E. Nemeth, and K. Claffy. Internet expansion, refinement and churn. *European Transactions on Telecommunications*, January 2002.
- [Bor02] S. Borthick. Will route control change the internet ? *Business Communications Review*, September 2002.
- [FBR02] N. Feamster, J. Borkenhagen, and J. Rexford. Controlling the impact of BGP policy changes on IP traffic. Technical report, Toronto, June 2002.
- [FP99] W. Fang and L. Peterson. Inter-as traffic patterns and their implications. In *IEEE Global Internet Symposium*, December 1999.
- [RL02] Y. Rekhter and T. Li. A border gateway protocol 4 (bgp-4). Internet draft, draft-ietf-idr-bgp4-17.txt, work in progress, May 2002.
- [SARK02] L. Subramanian, S. Agarwal, J. Rexford, and R. Katz. Characterizing the internet hierarchy from multiple vantage points. In *INFOCOM 2002*, June 2002.
- [Ste99] J. Stewart. *BGP4 : interdomain routing in the Internet*. Addison Wesley, 1999.
- [TMW97] K. Thompson, G. Miller, and R. Wilder. Wide-area internet traffic patterns and characteristics. *IEEE Network Magazine*, 11(6), November/December 1997.
- [UB02] S. Uhlig and O. Bonaventure. Implications of interdomain traffic characteristics on traffic engineering. *European Transactions on Telecommunications*, January 2002.