

oBGP: an Overlay for a Scalable iBGP Control Plane

Iuniana Oprescu^{1,5}, Mickaël Meulle¹, Steve Uhlig²,
Cristel Pelsser³, Olaf Maennel⁴, and Philippe Owezarski⁵

¹ Orange Labs, 38-40, rue du Général Leclerc
92794 Issy-les-Moulineaux Cedex 9, France,

{mihaela.oprescu, michael.meulle}@orange-ftgroup.com,
² Deutsche Telekom Laboratories & Technische Universität Berlin
Ernst-Reuter-Platz 7, 10587 Berlin, Germany
steve@net.t-labs.tu-berlin.de

³ Internet Initiative Japan, Jinbo-cho Mitsui Bldg.,
1-105 Kanda Jinbo-cho, Chiyoda-ku, Tokyo 101-0051, Japan
cristel@iij.ad.jp

⁴ University of Loughborough, Department of Computer Science,
Haslegrave Bldg., Loughborough, LE3TU, United Kingdom
olaf@maennel.net

⁵ Université de Toulouse; UPS, INSA, INP, ISAE; LAAS;
7, Avenue du colonel Roche, F-31077 Toulouse, France
owe@laas.fr

Abstract. The Internet is organized as a collection of networks called Autonomous Systems (ASes). The Border Gateway Protocol (BGP) is the glue that connects these administrative domains. Communication is thus possible between users worldwide and each network is responsible of sharing reachability information to peers through BGP. Protocol extensions are periodically added because the intended use and design of BGP no longer fit the current demands. Scalability concerns make the required iBGP full mesh difficult to achieve in today's large networks and therefore network operators resort to confederations or Route Reflectors (RRs) to achieve full connectivity. These two options come with a set of flaws of their own such as persistent routing oscillations, deflections, forwarding loops etc.

In this paper we propose a new architecture for the redistribution of external routes inside an AS. Instead of relying on the usual statically configured set of iBGP sessions, we propose to use an overlay of routing instances that are collectively responsible for (i) the exchange of routes with other ASes, (ii) the storage of internal and external routes, (iii) the storage of the entire routing policy configuration of the AS and (iv) the computation and redistribution of the best routes towards Internet destinations to each router of the AS.

Keywords: routing, BGP, architecture, management

1 Introduction

The Border Gateway Protocol (BGP) is the glue that enables the computation of end-to-end paths in the Internet. BGP allows networks, called Autonomous Systems (ASes), to exchange their routing information and to implement independently customized routing policies. The Internet has reached a size of more than 35000 ASes and roughly 350000 blocks of IP addresses.

BGP and internal BGP in particular have been widely studied, and many extensions and improvements have been proposed to deal with matters like network convergence[2][3] or route diversity[5][6]. On the other hand, general design and architectural issues in iBGP have not been sufficiently confronted in our opinion.

In this paper, we propose a new solution for iBGP routing within an AS using a distributed overlay of routing software. The solution we elaborate is meant as a viable framework for replacing the current iBGP and we introduce a setting for scalable and flexible routing. By gathering the routing information in a platform, we aim to offer easier management of protocols and policies at the network level.

2 BGP Routing in a Nutshell

The Border Gateway Protocol[7] is in fact two protocols: internal BGP (iBGP) for handling messages inside an AS and external BGP (eBGP) for exchanging reachability information with other ASes in the Internet. This clear distinction makes it possible for an ISP to deploy a new iBGP in its network without any impact on neighboring ASes. We will further concentrate on describing some aspects of the iBGP mechanism.

Within a BGP router, the decision process takes into account the interactions with every neighbor. Roughly speaking, if n is the number of prefixes advertised in the Internet, an iBGP Routing Information Base (RIB) will contain about $n * m$ routes in the worst case, where m is the number of neighbors sending their full BGP table as seen in Fig. 1. The best path to a destination is selected, installed in the Forwarding Information Base (FIB) and actually used to perform packet forwarding. The router will also advertise its best path for a given prefix to the adjacent BGP peers.

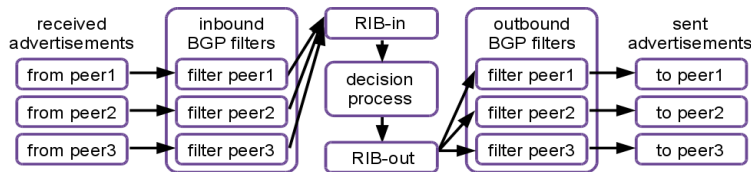


Fig. 1. The selection of the best route among the received routes

BGP requires that entries be kept for each reachable network: this constraint leads to large routing tables. There are many routes that cannot be aggregated

and even if the optimizations in vendor code reduce the size of memory needed, they still do not fix the problem. The natural growth of the Internet as well as its increasing connectivity and the tweaking of routing entries for traffic engineering purposes have inflated the size of BGP routing tables by a factor of more than 3 within the last decade [1]. The current trend of the routing table indicates continuous growth of the Internet and we expect future evolution to be similar, especially after the migration to the apparently inexhaustible IPv6 space.

Inside an AS, a single administrative entity manages all routers and distributes a consistent routing policy configuration. The goal of iBGP is to redistribute routing data inside the AS in accordance with the routing policy configured in each BGP router. iBGP routing originally required a full mesh between the routers within a single AS to guarantee that each router will be able to learn the best external route for forwarding IP packets.

The full mesh configuration can quickly turn out to be a scalability problem since the number of sessions grows with the square number of participants. There are two alternatives for avoiding the processing overhead induced by full mesh: confederations and route reflectors (RRs) illustrated in Fig. 2.

Confederations are sub-ASes meant to divide a large network into more manageable areas. A route reflector is a router that takes the role of a central point where a subset of the other routers peer. These designs are both prone to unpredictable effects such as persistent routing oscillations and forwarding loops affecting network convergence, sub-optimal routing due to network opacity or non-deterministic decisions influenced by the state of the network at the arrival time of the advertisements. We further detail these drawbacks in 2.1.

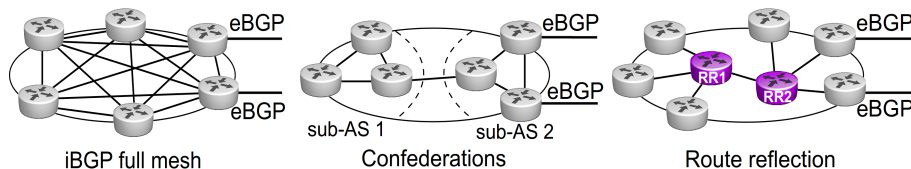


Fig. 2. ASes exemplifying a full mesh of iBGP sessions, confederations, route reflection

2.1 Plagues in Current iBGP

In specific architectures using route reflection, routing is victim of a series of aspects induced by the inherent design of iBGP.

Below is a brief display of the drawbacks in current iBGP:

- **scalability** is quantified by the number of protocol messages exchanged over time, the number of established sessions and especially by the size of the routing table. We estimate that the growth of the routing table can be handled by means of a partitioning model that we expose in this paper:

achieve scalability through division of the control plane by placing subsets of prefixes in different locations and then performing the computation of the BGP decision process in a distributed manner.

- **network opacity** occurs in architectures where route reflection schemes are used for propagating routes. Incomplete knowledge of the set of routes advertised by neighbor ASes leads to inconsistencies and issues such as routing oscillations and deflections that can cause forwarding loops. Extensive studies [2][3][4][10] give conditions and methods for defining correct iBGP configurations that avoid anomalies and achieve full mesh optimality.
- **poor route diversity** is a direct consequence of network opacity. The fundamental design of BGP route redistribution demands that each peer advertise only its best route. Diffusing a single route impacts the available choices and there is a noticeable loss of route diversity when comparing border routers to internal routers [5].

The graph in fig. 3 presents the diversity of neighbor ASes and BGP next-hops for the received prefixes on 5 random routers. The data reveals the fact that there is a large diversity in the received routes but this diversity is severely reduced by the BGP selection mechanism of the best route.

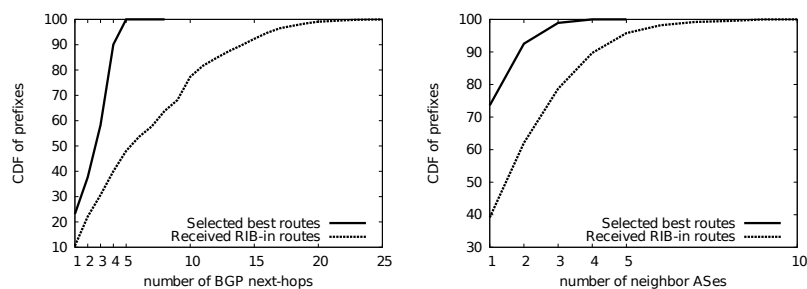


Fig. 3. Prefixes and routes on 5 random routers of a large ISP

Redundancy in case of failure is highly desirable and a secondary path could also be used for extra features such as load balancing or multipath routing. Protocol extensions introduce new capabilities for adding paths to BGP in [6], but there is no knowledge of the possible impact on current architectures.

- **management and troubleshooting** are often complex and challenging: inconsistency of the routing policies, path exploration meeting flap dampening [8] and difficulties in achieving network-wide traffic engineering are some of the issues encountered by network operators. A full view of the external routes and knowledge of the Interior Gateway Protocol (IGP) topology by one entity would make these processes easier. In oBGP, the interaction with the entire network is done through the overlay and the concentration of the BGP decision process on the nodes increases control over the network behavior.

We should, however, separate theory from practice. Some of the presented issues are commonly avoided with engineering tricks and configuration tweaking. Network operators adapt to inadvertencies by enforcing specific RR placement and building convenient topologies that behave correctly. Ideally, these aspects can be handled in an automated manner and this paper proposes an approach for better control over the network.

2.2 Previous Work

New routing paradigms like AIR[9], iBGPv2[10] or even PCE[11] propose different approaches for handling routing within an AS. LISP[12] tackles routing as a general problem and proposes a solution for the global Internet routing. In [13] Jing Fu exposes a centralized control scheme for IGP with faster routing convergence than link-state routing protocols. His results show it is possible to conceive a routing platform reaching performances comparable to native routing.

The need for separating routing from the routers is emphasised also by N. Feamster et al. in [15]. The presented work is a design overview of a Routing Control Platform (RCP) that aims to offer separate selection of routes on behalf of the routers while maintaining backward compatibility.

M. Caesar et al. later offer an implementation to the RCP concept. The prototype described in [16] has three modules: the IGP Viewer to collect topology information, the BGP engine that learns the BGP routes, performs the decision algorithm and then communicates the best paths to the routers and finally the Route Control Server that processes messages received from the other two modules and makes it possible to store one single copy of each BGP route, keep track of the routers to which each route has been assigned and maintain an order of preference of the egress point for each router [17]. We extend this work by going a step further in reaching scalability: in our approach, the prefix table is split, making possible parallel computation of routes while in the RCP solution, all the BGP information is concentrated in one point, even if there are multiple replicas of it.

Our hybrid solution integrates the division of the routing table within a centralized routing platform. Other projects advocate the idea of downsizing the routing table: ViAggre (Virtual Aggregation) is a configuration-only method for shrinking the size of the routing table in the Internet default-free zone. It proposes a “dirty slate” technique for distributing routing within an ISP network so that routers maintain only a part of the global routing table. One of the negative impacts of ViAggre[18] is a stretch imposed on traffic, diverting it from the native shortest path. Another inconvenient is the difficulty of the configuration. This same approach is advanced in [19] and X. Zhang et al. elaborate similar work in [20], but CRIO seems to bring more benefit to VPN routing.

The work of S. Uhlig et al. [21][22] emphasizes the fact that network operators need a smarter way to do route reflection. In [23][24] C. Pelsser et al. aim to build distributed route servers. We go beyond these proposals by providing scalability through the distribution of the control plane in iBGP routing.

3 oBGP: a Scalable Overlay for iBGP Routing

In today's IP networks, routing is highly distributed: each router in the AS makes its own decisions. We propose to separate the selection of paths (routing plane) from the actual forwarding of traffic (data plane) on distinct equipments. Offloading the control plane from the routers can be seen as a remedy to the explosion of the routing table size.

When rethinking the current design, we place all the knowledge of routing data into a separate iBGP routing plane handled by an overlay of routing processes that do not forward traffic. We propose to implement BGP routing engines called *oBGP*. The oBGP nodes act as the border routers of the domain and connect to the external peers through multi-hop eBGP sessions. This approach allows the overlay to receive all the routes from the neighboring ASes and aggregate the announced routes to achieve a unified complete view.

oBGP routing software is intended to be executed by additional servers running on commodity hardware. The logical overlay is composed of routing processes (or nodes) that are jointly responsible of:

- collecting, splitting and storing the complete set of routes received from eBGP and the internally originated routes,
- storing the routing policies and configurations of all the routers in the AS,
- computing BGP best paths for each router,
- redistributing the computed paths to the client routers.

One of the main concerns of an iBGP architecture is its ability to scale: support the growing routing table and handle protocol messages over time. To achieve scalability, we design an oBGP solution where the routing information is divided in several sub-planes. In this approach, distinct subsets of overlay nodes each handle only a fraction of the entire set of prefixes in the routing table.

3.1 Overview

The next paragraph explains the passage of a route advertisement in the oBGP overlay from the arrival in the AS to the installation of the best path in the RIB. Fig. 4 shows the chronological steps of a route announced to the oBGP overlay.

The oBGP acts as a border router and the neighboring ASes connect to an oBGP node through multi-hop eBGP sessions. When a route towards a destination (e.g. the prefix 1.2.3.0/24) is advertised in the Internet, it reaches the first oBGP node that determines the corresponding sub-plane in charge of the prefix. The oBGP node then forwards the information to the nodes handling the correct sub-plane. After running the BGP decision process and applying the according configuration and IGP topology constraints, the nodes output a best path. The overlay distributes the best path to the client routers connected through sessions and they can immediately install it in their RIBs. Upon reception of the best route, the native routing mechanism takes course and installs the path to the prefix in the FIB for actual packet forwarding.

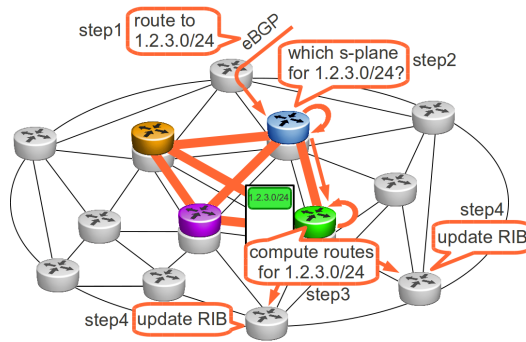


Fig. 4. The steps followed by an advertisement in the overlay network

The oBGP nodes need to be aware of the actual mapping of the reachable IP space within the overlay. To insure resiliency and avoid a single point of failure, a sub-plane is replicated on several oBGP nodes. Coordination between the copies of sub-planes is accomplished through an exchange of meta-data across the oBGP. The following paragraphs depict the sub-plane concept.

3.2 Distributed Storage

A router learns routes toward a given prefix from its neighbors, and in the general case routers of the same AS do not learn the same exact set of routes or the same quantity. The full visibility of BGP routes received from external ASes can be assimilated to a sum of queries on all border routers of an AS. Aggregating routes received on every border router is equivalent to the global view of the advertised Internet as seen by the domain.

oBGP manages to keep this external view intact by indexing it directly in the overlay according to a mapping mechanism. The oBGP nodes act as the collection of border routers of the AS and establish eBGP multi-hop sessions with neighbor ASes.

Storage of prefixes is distributed across the overlay and nodes divide between each other the computational load of the control plane. We define several chunks of the reachable address space that are allocated on distinct nodes. These large IP spaces are called routing sub-planes. The overlay is in charge of keeping a coherent state where no pair of sub-planes has overlapping prefixes and they are stored on different nodes. A structure similar to a distributed hash-table can be used for managing the sub-planes. The oBGP nodes guarantee the frontiers of the sub-plane, but another aspect to take into account is the replication of the information on the nodes covering the same sub-plane.

Index of Virtual Prefixes: The mapping of the sub-planes on the oBGP nodes takes into account the split factor $n = 4$ and attempts to evenly allocate each chunk of $total/n$ prefixes to a sub-plane. This strategy turns out to be very

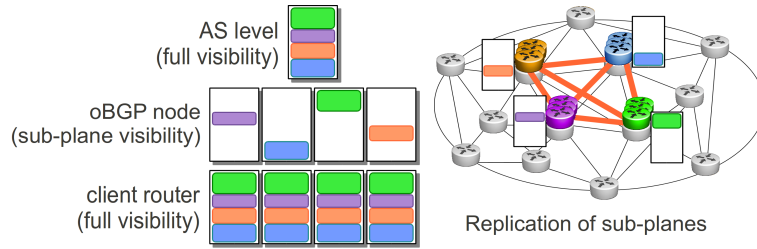


Fig. 5. The routing table is split between the $n = 4$ nodes of the overlay

coarse grained and thus we introduce smaller containers for the IP space called Virtual Prefixes as in [20].

Table 1 shows an example of a possible configuration of the sub-planes: the reachable IP space is divided in $n = 4$ sub-planes and each sub-plane covers the equivalent of a $/2$ prefix (corresponding roughly to 2^{30} possible hosts). To better control the load incurred by the oBGP nodes handling the sub-planes, the network operator may choose to define several virtual prefixes as is the case for sub-plane 1 that contains 2 virtual prefixes. The virtual prefixes may be swapped between the oBGP nodes in order to achieve a balanced load on the sub-planes. Data⁶ in columns 3 and 4 shows that the density of prefixes advertised in the Internet can be almost uniformly distributed across the previously defined sub-plane space. If the distribution varies in time, we deem necessary to use a dynamic algorithm.

Table 1. Sub-planes containing virtual prefixes

sub-plane ID	virtual prefixes	# of prefixes	% of total
sub-plane 1	64.0.0.0/4	53250	17.85%
	32.0.0.0/3	21408	7.17%
sub-plane 2	160.0.0.0/3	38425	12.88%
	192.0.0.0/4	37667	12.62%
sub-plane 3	80.0.0.0/4	34552	11.58%
	96.0.0.0/3	35679	11.96%
sub-plane 4	208.0.0.0/4	40207	16.82%
	128.0.0.0/3	17719	5.93%
	0.0.0.0/3	9411	3.15%

We envision as future work to develop an on-line procedure that allocates smaller virtual prefixes to the oBGP nodes to obtain a fine grain arrangement.

⁶ Dataset of November 2010, based on a total of 354682 prefixes

It is also possible to enforce a rule allowing for popular prefixes to be cached based on a statistical computation of the frequency of occurrence (i.e. cache the popular prefixes that are more stable as opposed to swapping more often the less popular prefixes).

3.3 Selection and Propagation of BGP Routes

The main purpose for offloading the control plane into an overlay is to achieve scalability of the routing table, but the separation of the decision process from the actual forwarding of routes has several other benefits such as complete visibility of the routes advertised to the AS.

The oBGP nodes gather information through eBGP and at the same time they are part of the IGP topology which allows them to be aware of the metrics toward the next-hop. This feature is important because the customized computation of the best BGP route for a given prefix for a particular router will take into account the full view of the BGP routes and the interior cost for reaching the next-hop. The optimal routes are what the client router would choose if it had full view. Complete knowledge of both topologies allows the routing engines to make a correct selection and avoid situations like routing loops.

Having a federating entity makes policy management easier: a global policy can be configured on the oBGP nodes and then applied to all the routes entering or exiting the AS. The overlay ensures AS-wide coordination while still allowing for specific policies to be safely implemented on individual routers. We can state that the BGP decision process is neighbor-specific and the algorithm will provide the best output for every individual router connected to the overlay.

Propagation of routes from the overlay to the client routers relies on the classic iBGP sessions. Once the oBGP nodes compute the best routes for the allotted prefixes, the result is pushed to the connected routers very much in the same way that a client router would receive advertisements from RRs. Fig. 6 shows a router connected using an iBGP session to each oBGP node responsible of a specific sub-plane.

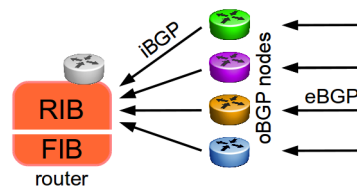


Fig. 6. The routers connect to the overlay through iBGP sessions

Once the best route to a given prefix reaches the router, the RIB is updated and the next-hop is programmed in the FIB. Client routers receive directly

the announcements of the best routes that the oBGP has chosen, therefore the number of received routes is smaller and the size of the RIB is reduced.

4 Deployment of oBGP

An oBGP overlay network can be progressively deployed on top of an existing iBGP architecture, using a step by step approach. The overlay is a logical topology and can be optimized according to the underlying physical graph and the location of the oBGP nodes in the network.

As a first step, the network operator has to decide of an initial number of sub-planes and a mapping of each available oBGP node to one sub-plane while taking into account the expected redundancy of the overlay. An initial partitioning of the IP space specifies the granularity of the virtual prefixes and their correspondence to the defined sub-planes. The overlay nodes can be configured with specific policies to apply and finely tune the selection of routes for the individual client routers in the AS.

The second step is twofold: setting up the topology between the nodes of the same sub-plane and interconnecting the different sub-planes. Note that the oBGP approach does not specify a topology and several arrangements of the nodes can be studied in order to optimize the performances. The oBGP nodes participate to the IGP topology to retrieve knowledge of routing costs between any given pair of routers within the AS. The overlay implements redistribution of routes between the sub-planes and replication of the routes in a given sub-plane using reliable flooding mechanisms.

The final step is more delicate and consists of safely migrating the eBGP sessions to the overlay and removing iBGP sessions between routers. Integrity and coherence of routes must be guaranteed during the transition, until the oBGP manages to take over the route redistribution.

For each router of the AS, establish one iBGP session with at least one oBGP node of each sub-plane (the closest possible). At this point, routers can receive the routes from the oBGP overlay. To enable the router to redistribute the eBGP-received routes to its iBGP neighbors, we temporary turn each iBGP session of the router into a route reflector to client iBGP session. As soon as a pair of routers is migrated, the session between them can be removed.

When all border routers have been migrated, internal routers having only iBGP sessions are migrated the same way we migrate eBGP connected routers.

5 Gain through Design

This section summarizes some of the solution strengths that derive from the design and presents some feasible extensions.

The de-correlation of route selection from propagation allows routers to gain a broader visibility: the decision process takes into account both BGP and IGP information leading to an optimal selection at the AS-level. The overlay is aware

of all external routes received through eBGP and of the IGP topology, which implies that the decision process is based on complete knowledge of the routes.

Another advantage is the increased diversity of routes that can be received by the client routers. Indeed, we first federate knowledge of the routes at the overlay level and then distribute the computational load according to sub-planes so that in the end, client routers connect to each of the sub-planes and receive the optimal routes. Controlling only a subset of the IP space means less memory needed to store a sub-plane instead of the entire RIB and less BGP reachability information to process on each oBGP node. The speedup in the selection of the best routes comes from the distribution of prefixes on several nodes, allowing for parallel computation.

The construction of the overlay leaves room for future additional features like flexible load sharing of routing data. An algorithm for dynamically partitioning and mapping the reachable space enables a re-organization of the various prefixes in the sub-planes.

6 Conclusions and Future Work

In this paper we present a new framework for scalable iBGP routing. The oBGP concept is illustrated: an overlay responsible for performing the BGP decision process on behalf of the client routers within the AS. We expose some of the major drawbacks in current iBGP and how the oBGP routing platform solves these issues. We provide the design principles and advantages of oBGP then reveal a possible scenario for deployment. Through the construction of this approach, oBGP provides ground for implementations of extra features and proposes a new direction in the study of the iBGP control plane scalability.

As previously mentioned, the advantage of splitting the routing table can be overshadowed by the computational overhead induced in the overlay. An important point of future work consists of determining the optimal threshold for which it is appealing to compute paths with oBGP. Research perspectives include refining the split algorithm and improving it to gracefully handle the dynamic re-organization of the virtual prefixes on the oBGP nodes. We will also evaluate the architecture, study the relevance of parameters and quantify scalability, convergence time, correctness and compliance to the routing policy.

References

1. Geoff Huston: <http://www.potaroo.net/>
2. A. Rawat, M. Shayman: Preventing persistent oscillations and loops in iBGP configuration with route reflection. *Comput.Netw.*, vol.50, no.18, pp. 3642–3665 (2006)
3. T. Griffin, G. Wilfong: On the correctness of iBGP configuration. *SIGCOMM Computer Comm. Review*, vol. 32, no.4, pp. 17–29 (2002)
4. T. Griffin, G. Wilfong: Analysis of the MED Oscillation Problem in BGP. *Proceedings of the 10th IEEE International Conference on Network Protocols* (2002)
5. S. Uhlig and S. Tandel: Quantifying the BGP routes diversity inside a tier-1 network. *Networking* (2006)

6. D. Walton, A. Retana, E. Chen, J. Scudder: Advertisement of Multiple Paths in BGP. Internet draft, draft-ietf-idr-add-paths-04 (2010)
7. Y. Rekhter, T. Li and S. Hares: A Border Gateway Protocol 4 (BGP-4). RFC 4271, IETF (2006)
8. C. Villamizar, R. Chandra, R. Govindan: BGP Route Flap Damping. RFC 2439, IETF (1998)
9. J.J. Garcia-Luna-Aceves and D. Sampath: Scalable integrated routing using prefix labels and distributed hash tables for MANETs. IEEE 6th International Conference on Mobile Adhoc and Sensor Systems 188–198 (2009)
10. M.-O. Buob, S. Uhlig and M. Meulle: Designing Optimal iBGP Route-Reflection Topologies. IFIP Networking 542–553 (2008)
11. A. Farrel, J.P. Vasseur, J. Ash: A Path Computation Element PCE-Based Architecture. RFC 4364, IETF (2006)
12. R. Hinden: New Scheme for Internet Routing and Addressing (ENCAPS) for IPNG. RFC 1955, IETF (1996)
13. J. Fu, P. Sjödin and G. Karlsson: Intra-domain routing convergence with centralized control. Comput. Netw. vol.53, no.18, 2985–2996 (2009)
14. IETF ForCES Working Group: <http://tools.ietf.org/wg/forces/>
15. N. Feamster, H. Balakrishnan, J. Rexford, A. Shaikh, and J. van der Merwe: The case for separating routing from routers. Proceedings of the ACM SIGCOMM workshop on Future directions in network architecture (2004)
16. M. Caesar, D. Caldwell, N. Feamster, J. Rexford, A. Shaikh and J. van der Merwe: Design and implementation of a routing control platform. Proc. of the 2nd Symposium on Networked Systems Design and Implementation (2005)
17. Zebra Route Server: <http://www.zebra.org/zebra/Route-Server.html>
18. H. Ballani, P. Francis, T. Cao and J. Wang: ViAggre: Making Routers Last Longer! Proc. of workshop on Hot Topics in Networks (2008)
19. P. Francis, X. Xu, H. Ballani, D. Jen, R. Raszuk and L. Zhang: FIB Suppression with Virtual Aggregation. Internet draft, draft-ietf-grow-va-03 (2010)
20. X. Zhang, P. Francis, J. Wang and K. Yoshida: Scaling IP Routing with the Core Router-Integrated Overlay. Proc. of the IEEE International Conference on Network Protocols 147–156 (2006)
21. S. Uhlig, C. Pelsser, B. Quoitin and O. Bonaventure: Vers des réflecteurs de route plus intelligents. Colloque Francophone sur l'Ingénierie des protocoles (2005)
22. O. Bonaventure, S. Uhlig, B. Quoitin: The case for more versatile BGP Route-Reflectors. Internet draft, draft-bonaventure-bgp-route-reflectors-00 (2004)
23. C. Pelsser, A. Masuda and K. Shiomoto: Scalable Support of Interdomain Routes in a Single AS. Proc. of IEEE Globecom (2009)
24. C. Pelsser, A. Masuda and K. Shiomoto: A novel internal BGP route distribution architecture. Proc. of the IEICE General Conference (2009)