

Network Working Group
Internet-Draft
Intended status: Informational
Expires: January 28, 2021

Y. Luo
L. Qu
China Telcom Co., Ltd.
X. Huang
Tencent
G. Mishra
Verizon Inc.
H. Chen
Futurewei
S. Zhuang
Z. Li
Huawei
July 27, 2020

Architecture for Use of BGP as Central Controller **draft-cth-rtgwg-bgp-control-05**

Abstract

BGP is a core part of a network including Software-Defined Networking (SDN) system. It has the traffic engineering information on the network topology and can compute optimal paths for a given traffic flow across the network.

This document describes some reference architectures for BGP as a central controller. A BGP-based central controller can simplify the operations on the network and use network resources efficiently for providing services with high quality.

Requirements Language

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [[RFC2119](#)].

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <https://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any

time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on January 28, 2021.

Copyright Notice

Copyright (c) 2020 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<https://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1.	Introduction	3
2.	Terminology	4
3.	Architectures	5
3.1.	Building Blocks	5
3.1.1.	TEDB	5
3.1.2.	SLDB	5
3.1.3.	TPDB	5
3.1.4.	CSPF	6
3.1.5.	TM	6
3.2.	One Controller	6
3.3.	Controller Cluster	8
3.4.	Hierarchical Controllers	10
4.	Application Scenarios	12
4.1.	Business-oriented Traffic Steering	12
4.1.1.	Preferential Users	12
4.1.2.	Preferential Services	13
4.2.	Traffic Congestion Mitigation	14
4.2.1.	Congestion Mitigation in Core	15
4.2.2.	Congestion Mitigation among ISPs	15
4.2.3.	Congestion Mitigation at International Edge	16
5.	Security Considerations	17
6.	IANA Considerations	17
7.	Acknowledgements	17
8.	Contributors	17
9.	References	17
9.1.	Normative References	18

9.2. Informative References	18
Authors' Addresses	19

1. Introduction

Border Gateway Protocol (BGP) [[RFC1771](#)] is an exterior gateway protocol (EGP). It is developed to exchange routing information among routers in different autonomous systems (ASes). Along its developments, BGP has been extended to provide numerous new functions. It collects the link states including traffic engineering (TE) information from other protocols such as IGP and distributes them among routers in different ASes [[RFC7752](#)]. It also controls the redirection of traffic flows [[RFC5575](#)]. Furthermore, it distributes MPLS labels [[RFC3107](#)]. For scalability, BGP is extended to have Route Reflector (RR) [[RFC4456](#)].

For segment routing (SR), BGP is extended to advertise SR policies with candidate paths to the policy headend routers, which are typically ingress routers [[I-D.ietf-idr-segment-routing-te-policy](#)]. The SR specific PCEP extensions are defined in [[I-D.ietf-pce-segment-routing](#)]. A stateful PCE can compute an SR traffic engineering (SR-TE) path satisfying a set of constraints, and initiate an SR-TE path on a headend router using the extensions.

An SDN controller (or controller for short) is the core of an SDN system or network. It is between network elements (NEs) such as routers or switches at one end and applications such as Operational Support System (OSS) or Network Management System (NMS) at the other end. The essential function of a controller is to steer traffic flows across the network for providing more services with higher quality. It manages network resources such as link bandwidth, computes expected paths for carrying traffic flows based on available network resources, programs the network elements for the creation of tunnels along the paths, and redirects traffic flows into corresponding tunnels.

Based on the current BGP, it is natural, beneficial and relatively simple to extend BGP to become a controller. Using BGP as a controller for a network will greatly simplify the operations on the network. It avoids deploying, operating and maintaining a new extra component or protocol such as PCE as a controller in the network.

This document describes some reference architectures for BGP as a central controller and introduces some scenarios to which the BGP controller can be applied.

2. Terminology

- o SR: Segment Routing
- o RR: Route Reflector
- o SID: Segment Identifier
- o SR-Path: Segment Routing Path
- o SR-Tunnel: Segment Routing Tunnel
- o TEDB: Traffic Engineering Database
- o LSDB: Link State Database
- o SLDB: SID/Label Database
- o TPDB: Tunnel and Path Database
- o CSPF: Constrained Shortest Path First
- o TM: Tunnel Manager
- o NMS: Network Management System
- o SRLB: SR Local Block
- o NE: Network Element
- o PCE: Path Computation Element
- o AS: Autonomous System
- o QoS: Quality of Service
- o ISP: Internet Service Provider
- o MAN: Metropolitan Area Network
- o OTT: Over the Top
- o OTTSP: Over the Top Service Provider, or Content Operator
- o AR: Access Router

3. Architectures

An architecture for the use of BGP as a central controller is based on the essential function of a controller. It is constructed from some building blocks or components. After introduction to building blocks, a few of reference architectures are described in this section.

3.1. Building Blocks

Some critical building blocks are briefed. They are Traffic Engineering Database (TEDB or TED for short), SID/Label Database (SLDB), Tunnel and Path Database (TPDB), Constrained Shortest Path First (CSPF), and Tunnel Manager (TM).

3.1.1. TEDB

The Traffic Engineering Database (TEDB) stores the Traffic Engineering (TE) information about the network. It includes the unreserved bandwidth at each of eight priority levels for every link in the network.

TEDB can be an individual block, which is constructed from the link state information received. It may be embedded into the link state database (LSDB) in the BGP when the BGP creates/updates the LSDB from the link state information it receives.

3.1.2. SLDB

The SID/Label Database (SLDB) records and maintains the status of every Segment Identifier (SID) and label for every node, interface/link and/or prefix in the network, which the controller controls. The status of SID/label indicates whether the SID/Label is assigned. If it is assigned, then the object such as the node, link or prefix, to which it is assigned, is recorded.

SLDB can be an individual block, which is constructed from the link state information such as SR Local Block (SRLB) that the BGP receives. It may be embedded into the link state database (LSDB) in the BGP when the BGP creates the LSDB from the link state information it receives.

3.1.3. TPDB

The Tunnel and Path Database (TPDB) stores the information for every tunnel, which includes:

- o the parameters received for the tunnel from a user/application,

- o the path computed for the tunnel,
- o the resources such as link bandwidth reserved along the path for the tunnel,
- o the SID/labels assigned along the path for the tunnel, and
- o the status of the tunnel.

3.1.4. CSPF

The Constrained Shortest Path First (CSPF) computes a path for a tunnel such as SR tunnel or LSP tunnel that satisfies a set of given constraints using the information in TEDB.

3.1.5. TM

The Tunnel Manager (TM) receives a request for an operation on a tunnel from a user or an application such as Network Management System (NMS). The operation may be a creation of a new tunnel, a deletion of an existing tunnel, or a change to an existing tunnel.

When receiving a request for creating a new tunnel, the TM asks the CSPF to compute a path for the tunnel that satisfies the constraints given for the tunnel.

After obtaining the path for the tunnel from the CSPF, the TM requests the SLDB to assign SID/labels along the path for the tunnel and asks the TEDB to reserve the resources such as link bandwidth along the path for the tunnel.

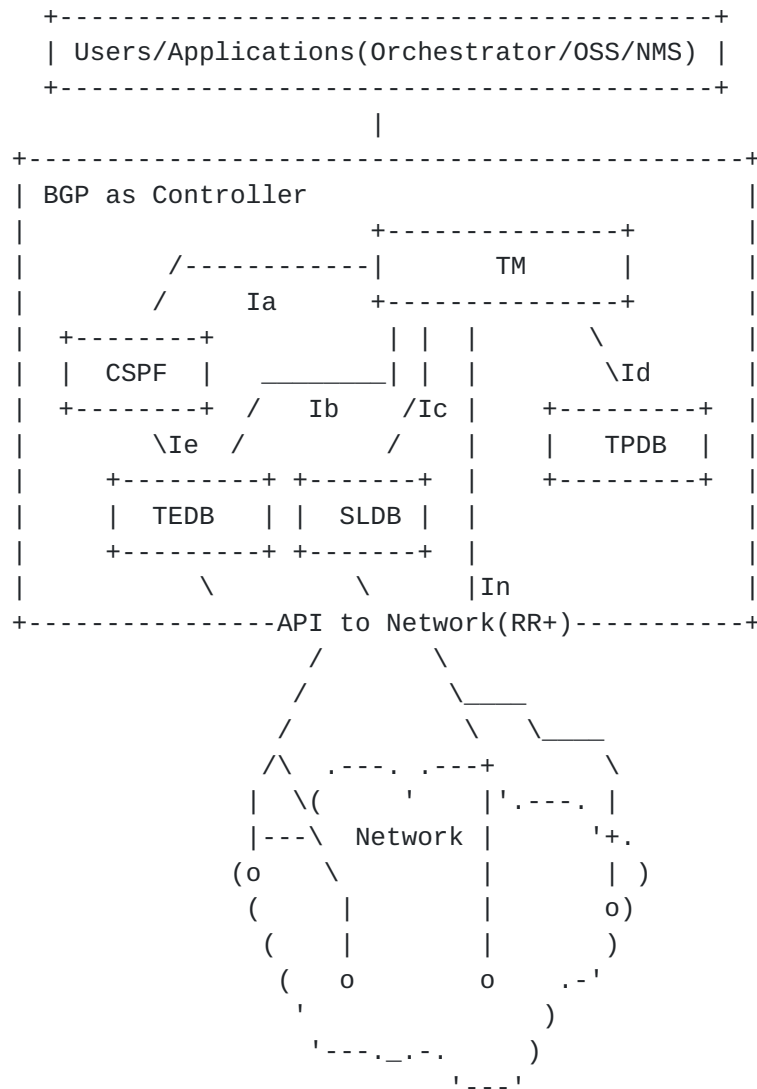
The TM in a central controller may set up the tunnel along the path in the network by programming each of the NEs along the path through the API to the network. In a SR network, the TM initiates a SR tunnel in the network by sending a sequence of SID/labels to the source NE of the tunnel.

The TM records the information for the tunnel in the Tunnel and Path Database (TPDB). The information includes the path computed for the tunnel, the resources such as bandwidth reserved along the path, the SID/labels assigned along the path for the tunnel, and the status of the tunnel.

3.2. One Controller

Figure below illustrates a reference architecture for using the BGP as a central controller, which controls a network. The BGP as a controller in the reference architecture controls a network through

The BGP controller comprises a number of modules, including a TM, a CSPF, a TEDB, a SLDB and a TPDB. The interfaces among these modules are listed as follows:



- o Interface Ia between the TM and the CSPF. Through this interface, the TM requests the CSPF to compute a path for a tunnel with a set of constraints, and the CSPF responses the TM with the path computed that satisfies the constraints.

- o Interface Ib between the TM and the TEDB. When a tunnel is to be created, through this interface, the TM reserves in the TEDB the TE resources such as link bandwidths on every link along the path computed for the tunnel. When a tunnel is deleted, the TM releases the TE resources such as link bandwidths on every link along the path for the tunnel.
- o Interface Ic between the TM and the SLDB. When a tunnel is to be created, through this interface, the TM reserves in the SLDB a SID/label for every link or some links along the path computed for the tunnel. When a tunnel is deleted, the TM releases the SID/label for every link or some links along the path for the tunnel.
- o Interface Id between the TM and the TPDB. the TM updates the information for every tunnel in the TPDB through this interface.
- o Interface Ie between the CSPF and the TEDB. Through this interface, the CSPF accesses the traffic engineering information such as link bandwidths when it computes a path for a tunnel.

There is an interface In between the BGP controller and the network. In fact, there is a control channel (or interface) between the BGP controller and every (edge) node in the network.

Initially, the TEDB obtains the original traffic engineering (TE) information such as link bandwidths from the network through the interface In (i.e., API to network) for every link in the network. The SLDB gets the original SID/label resources from the network through the interface for every node, link and prefix in the network.

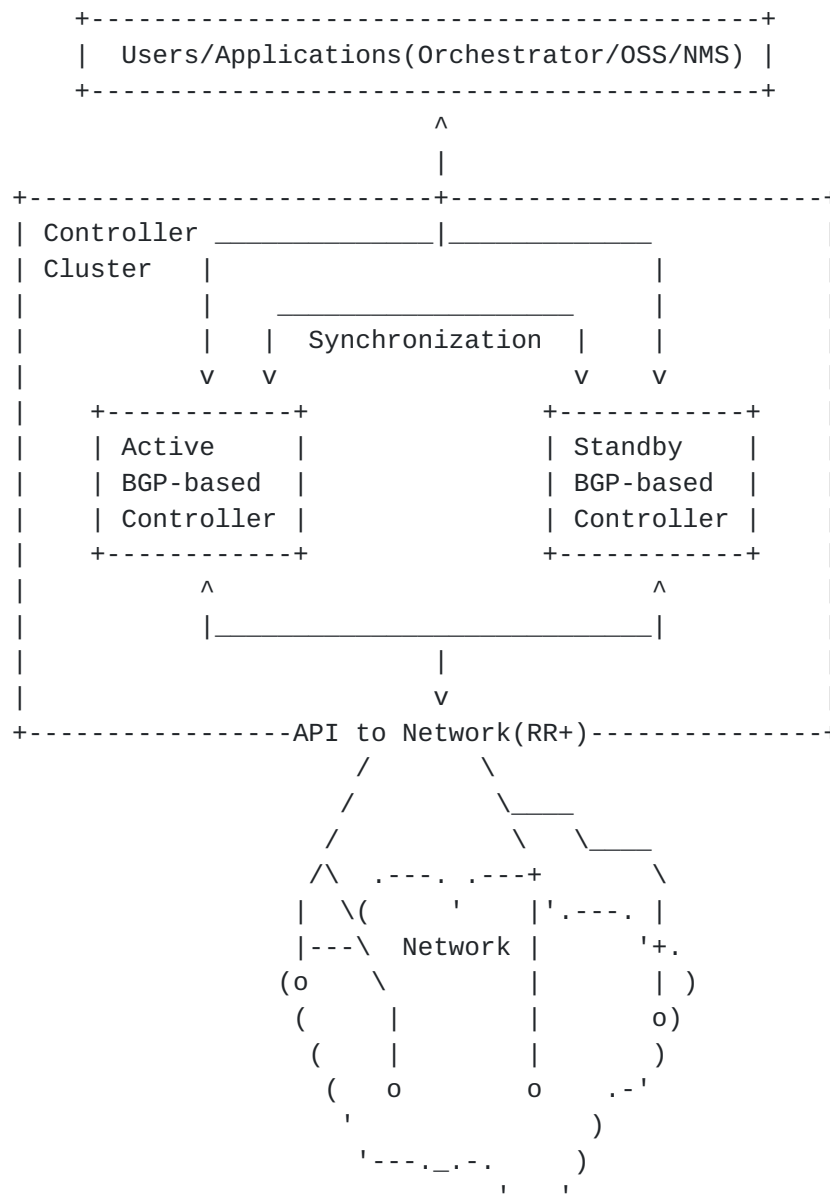
3.3. Controller Cluster

A critical issue in a network with a central controller is the failure of the controller, which is a single point of failure (SPOF). If the controller fails, the entire network may not work.

A controller cluster (i.e., a group of controllers) works as a single controller from user's point of view. A simple controller cluster consists of two controllers. One works as a active (or say primary) controller, and the other as a standby (or say secondary) controller. In normal operations, the active controller is responsible for the network it controls. It also synchronizes with the standby controller. When the active controller fails, the standby controller becomes a new active controller, which controls the network.

The Figure below illustrates a simple controller cluster containing two BGP-based controllers: Active BGP-based Controller and Standby BGP-based Controller. In normal operations, the active controller

interacts with users and/or applications. For example, it receives configurations for tunnels and the traffic flows to tunnels from users. The active controller instructs the network elements in the network to provide the services requested by users and/or applications. For example, after receiving the configurations for a tunnel and a traffic flow to the tunnel, the active controller computes a path for the tunnel, programs (or say instructs) the network elements along the path for creating the tunnel, and instructs the ingress of the tunnel to direct the traffic flow into the tunnel.

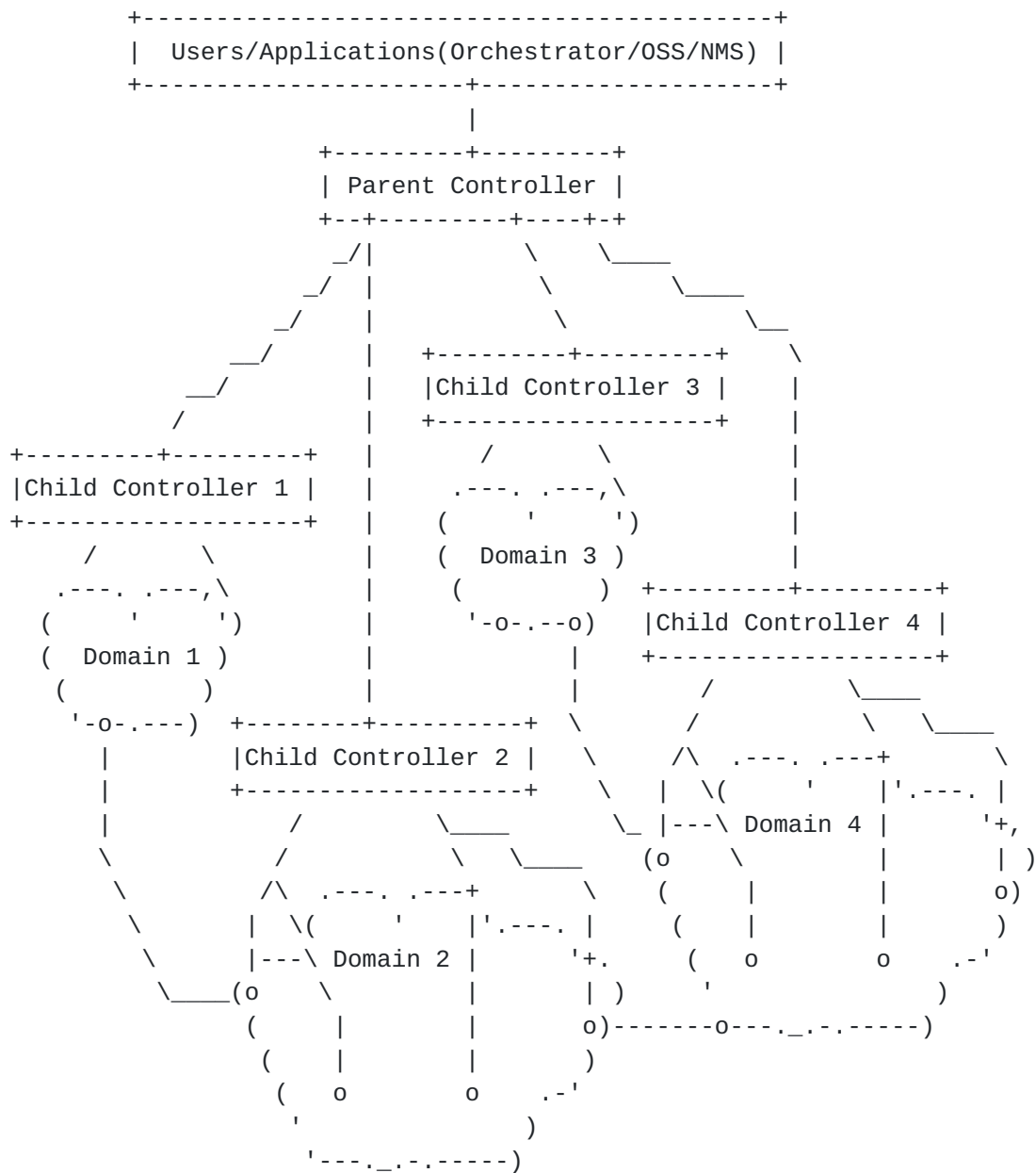


During this process, the status information about the network is updated in the active controller. The information includes: the

traffic engineering information in their TEDBs, the SID/label information in their SLDBs, and the configurations, paths, resources and status for tunnels in their TPDBs. The active controller synchronizes this information with the standby controller. Thus these two controllers have the same status information about the network. When the active controller fails, the standby controller takes over the role of the active controller smoothly and becomes active controller.

3.4. Hierarchical Controllers

The Figure below illustrates a system with hierarchical controllers. There is one Parent Controller and four Child Controllers: Child Controller 1, Child Controller 2, Child Controller 3 and Child Controller 4.



The parent controller communicates with these four child controllers and controls them, each of which controls (or is responsible for) a domain. Child controller 1 controls domain 1, Child controller 2 controls domain 2, Child controller 3 controls domain 3, and Child controller 4 controls domain 4.

One level of hierarchy of controllers is illustrated in the figure above. There is one parent controller at top level, which is not a child controller. Under the parent controller, there are four child controllers, which are not parent controllers.

In a general case, at top level there is one parent controller that is not a child controller, there are some controllers that are both parent controllers and child controllers, and there are a number of child controllers that are not parent controllers. This is a system of multiple levels of hierarchies, in which one parent controller controls or communicates with a first number of child controllers, some of which are also parent controllers, each of which controls or communicates with a second number of child controllers, and so on.

The parent controller receives requests for creating end to end tunnels from users or applications. For each request, the parent controller is responsible for obtaining a path for the tunnel and creating the tunnel along the path through sending instructions to the corresponding child controllers.

4. Application Scenarios

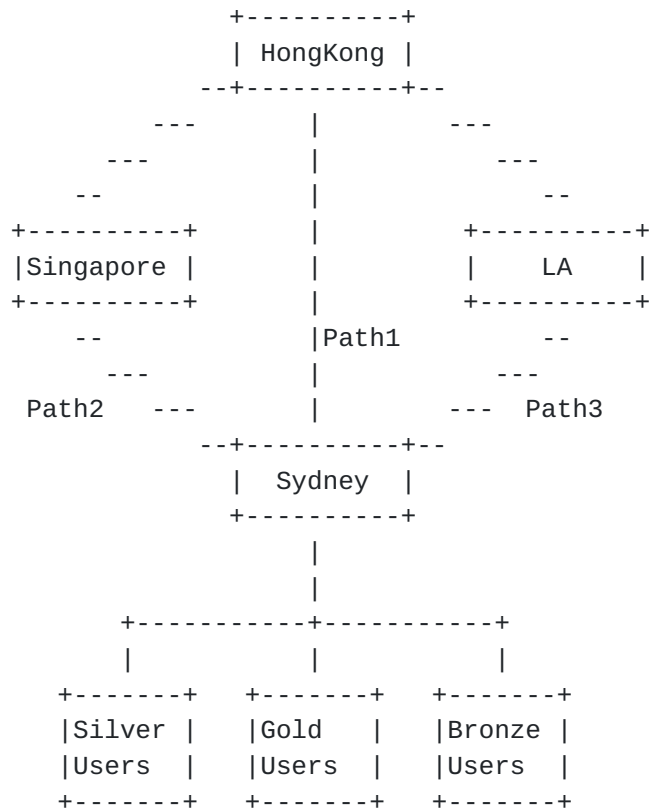
This section introduces a set of scenarios to which the controller can be applied.

4.1. Business-oriented Traffic Steering

It is reasonable in commercial sense to provide multiple paths to the same destination with differentiated experiences for preferential users/services. This is an efficient approach to maximize providers' network resource usage as well as their profit and offer more choices to network users.

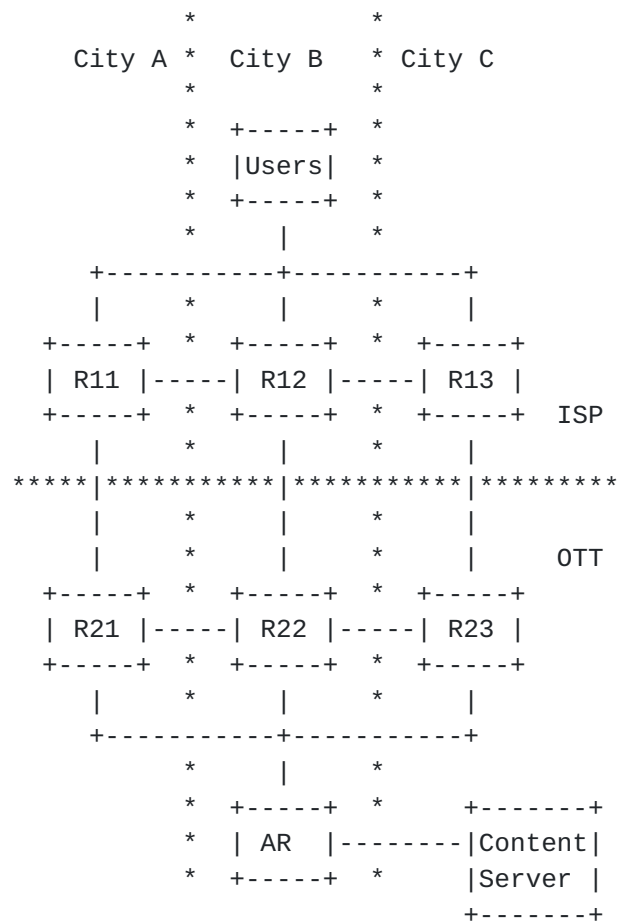
4.1.1. Preferential Users

In the Figure below for an ISP network, there are three kinds of users in Sydney, saying Gold, Silver and Bronze, and they wish to visit website located in HongKong. The ISP provides three different paths with different experiences according to users' priority. The Gold Users may use Path1 with less latency and loss. The Silver Users may use the Path2 through Singapore with less latency but maybe some congestion there. The Bronze Users may use Path3 through LA with some latency and loss.



[4.1.2.](#) Preferential Services

As depicted in the Figure below, the OTTSP has 3 exits with one ISP, which are located in City A, City B and City C. The content is obtained from Content Server and send to the exits through AR. An OTTSP may make its steering strategy based on different services. For example, the OTTSP in the Figure may choose exit R21 for video service and exit R22 for web service, which REQUIRES a mechanism/ system exists to identify different services from traffic flow.



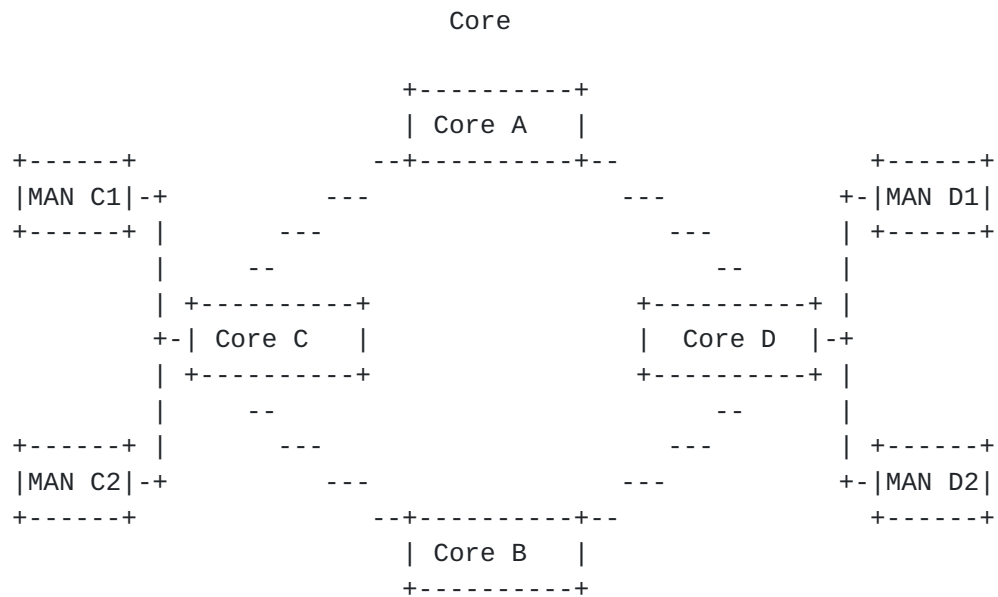
4.2. Traffic Congestion Mitigation

It is a persistent goal for providers to increase the utilization ratio of their current network resources, and to mitigate the traffic congestion. Traffic congestion is possible to happen anywhere in the ISP network(MAN, IDC, core and the links between them), because internet traffic is hard to predict. For example, there might be some local online events that the network operators didn't know beforehand, or some sudden attack just happened. Even for the big events that can be predicted, such as annual online discount of e-commerce company, or IOS update of Apple Inc, we could not guarantee there is no congestion. Since the network capacity expansion is usually an annual operation, there could be delay on any links of the engineering. As a result, the temporary traffic steering is always needed. The same thing happens to the OTT networks as well.

It should be noted that, the traffic steering is absolutely not a global behavior. It just acts on part of the network, and it's temporary.

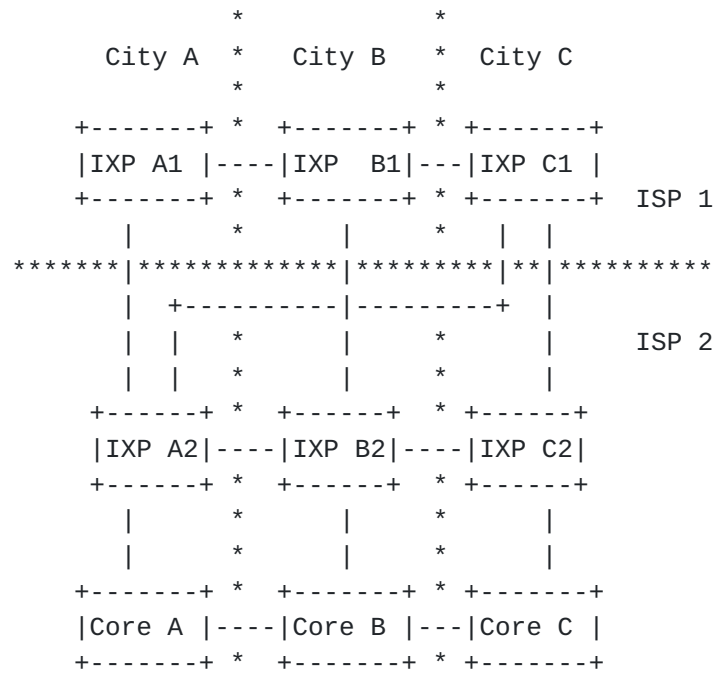
4.2.1. Congestion Mitigation in Core

As depicted in the Figure below, traffic from MAN C1 to MAN D2 follows the path Core C->Core B->Core D as the primary path, but somehow the load ratio becomes too much. It is reasonable to transfer some traffic load to less utilized path Core C->Core A->Core D when the primary path has congestion.



4.2.2. Congestion Mitigation among ISPs

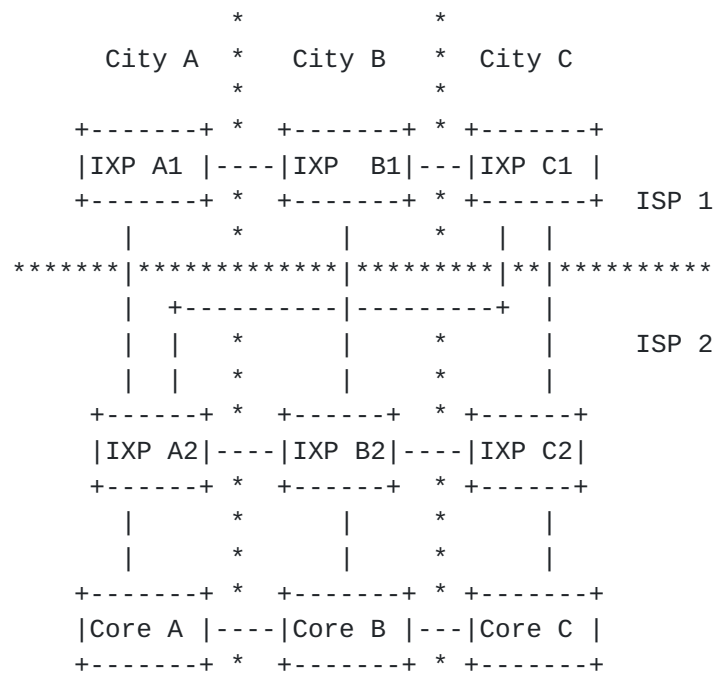
As depicted in the Figure below, ISP1 and ISP2 are interconnect by 3 exits which are located in 3 cities respectively. The links between ISP1 and ISP2 in the same city are called local links, and the rest are long distance links. Traffic from IXP C1 to Core A in ISP 2 usually passes through link IXP C1->IXP A2->Core A. This is a long distant route, directly connecting city C and city A. Part of traffic could be transferred to link IXP.



4.2.3. Congestion Mitigation at International Edge

An ISP usually interconnects with more than 2 transit networks at the international edge, so it is quite common that multiple paths may exist for the same foreign destination. Usually those paths with better QoS properties such as latency, loss, jitter and etc are often preferred. Since these properties keep changing from time to time, the decision of path selection has to be made dynamically.

As depicted in the Figure below, the traffic to the foreign destination H from IP core network (AS C1) has two choices on transit network, saying Transit A and Transit B. Under normal conditions, Transit B is the primary choice, but Transit A will be preferred when the QoS of Transit B gets worse. As a result, the same traffic will go through Transit A instead.



5. Security Considerations

The interactions with a BGP-based controller are similar to those with any other SDN controller. The security implications of SDN controller have not been fully discussed or described. Therefore, protocol and applicability for solutions around this architecture must take proper account of these concerns.

6. IANA Considerations

This document does not require any IANA actions.

7. Acknowledgements

The authors would like to thank Chris Bowers, Jeff Tantsura for their valuable suggestions and comments on this draft.

8. Contributors

Nan Wu
Huawei
Email: eric.wu@huawei.com

9. References

9.1. Normative References

- [RFC1771] Rekhter, Y. and T. Li, "A Border Gateway Protocol 4 (BGP-4)", [RFC 1771](#), DOI 10.17487/RFC1771, March 1995, <<https://www.rfc-editor.org/info/rfc1771>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), DOI 10.17487/RFC2119, March 1997, <<https://www.rfc-editor.org/info/rfc2119>>.
- [RFC3107] Rekhter, Y. and E. Rosen, "Carrying Label Information in BGP-4", [RFC 3107](#), DOI 10.17487/RFC3107, May 2001, <<https://www.rfc-editor.org/info/rfc3107>>.
- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", [RFC 4271](#), DOI 10.17487/RFC4271, January 2006, <<https://www.rfc-editor.org/info/rfc4271>>.
- [RFC4456] Bates, T., Chen, E., and R. Chandra, "BGP Route Reflection: An Alternative to Full Mesh Internal BGP (IBGP)", [RFC 4456](#), DOI 10.17487/RFC4456, April 2006, <<https://www.rfc-editor.org/info/rfc4456>>.
- [RFC5575] Marques, P., Sheth, N., Raszuk, R., Greene, B., Mauch, J., and D. McPherson, "Dissemination of Flow Specification Rules", [RFC 5575](#), DOI 10.17487/RFC5575, August 2009, <<https://www.rfc-editor.org/info/rfc5575>>.
- [RFC7752] Gredler, H., Ed., Medved, J., Previdi, S., Farrel, A., and S. Ray, "North-Bound Distribution of Link-State and Traffic Engineering (TE) Information Using BGP", [RFC 7752](#), DOI 10.17487/RFC7752, March 2016, <<https://www.rfc-editor.org/info/rfc7752>>.

9.2. Informative References

- [I-D.ietf-idr-bgppls-segment-routing-epe]
Previdi, S., Talaulikar, K., Filsfils, C., Patel, K., Ray, S., and J. Dong, "BGP-LS extensions for Segment Routing BGP Egress Peer Engineering", [draft-ietf-idr-bgppls-segment-routing-epe-19](#) (work in progress), May 2019.
- [I-D.ietf-idr-flowspec-path-redirect]
Velde, G., Patel, K., and Z. Li, "Flowspec Indirection-id Redirect", [draft-ietf-idr-flowspec-path-redirect-11](#) (work in progress), May 2020.

[I-D.ietf-idr-segment-routing-te-policy]

Previdi, S., Filsfils, C., Talaulikar, K., Mattes, P., Rosen, E., Jain, D., and S. Lin, "Advertising Segment Routing Policies in BGP", [draft-ietf-idr-segment-routing-te-policy-09](#) (work in progress), May 2020.

[I-D.ietf-isis-segment-routing-extensions]

Previdi, S., Ginsberg, L., Filsfils, C., Bashandy, A., Gredler, H., and B. Decraene, "IS-IS Extensions for Segment Routing", [draft-ietf-isis-segment-routing-extensions-25](#) (work in progress), May 2019.

[I-D.ietf-pce-segment-routing]

Sivabalan, S., Filsfils, C., Tantsura, J., Henderickx, W., and J. Hardwick, "PCEP Extensions for Segment Routing", [draft-ietf-pce-segment-routing-16](#) (work in progress), March 2019.

[I-D.ietf-rtgwg-bgp-routing-large-dc]

Lapukhov, P., Premji, A., and J. Mitchell, "Use of BGP for routing in large-scale data centers", [draft-ietf-rtgwg-bgp-routing-large-dc-11](#) (work in progress), June 2016.

[I-D.ietf-spring-segment-routing]

Filsfils, C., Previdi, S., Ginsberg, L., Decraene, B., Litkowski, S., and R. Shakir, "Segment Routing Architecture", [draft-ietf-spring-segment-routing-15](#) (work in progress), January 2018.

Authors' Addresses

Yujia
China Telcom Co., Ltd.
109 West Zhongshan Ave, Tianhe District
Guangzhou 510630
China

Email: luoyuj@sdu.edu.cn

Liang
China Telcom Co., Ltd.
109 West Zhongshan Ave, Tianhe District
Guangzhou 510630
China

Email: ouliang@chinatelecom.cn

Xiang
Tencent

Email: terranhuang@tencent.com

Gyan S. Mishra
Verizon Inc.
13101 Columbia Pike
Silver Spring MD 20904
USA

Phone: 301 502-1347
Email: gyan.s.mishra@verizon.com

Huaimo Chen
Futurewei
Boston, MA
USA

Email: Huaimo.chen@futurewei.com

Shunwan Zhuang
Huawei
Huawei Bld., No.156 Beiqing Rd.
Beijing 100095
China

Email: zhuangshunwan@huawei.com

Zhenbin Li
Huawei
Huawei Bld., No.156 Beiqing Rd.
Beijing 100095
China

Email: lizhenbin@huawei.com

