Internet Engineering Task Force Internet-Draft Intended status: Informational Expires: December 5, 2012

# Considerations on the AS-Level Application-Layer Traffic Optimization draft-asai-cross-domain-overlay-04

## Abstract

Overlay routing technologies (a.k.a. peer-to-peer technologies) have been introduced to various distributed systems, such as content delivery networks and live media streaming systems. However, these systems cannot always utilize optimal network resources from the viewpoint of layer 3 network providers and operators of layer 3 network providers have difficulties in controlling and optimizing the traffic of these systems because these systems construct their own networks (i.e., overlay networks) over the layer 3 network without taking into account layer 3 network's routing policies. The ALTO Working Group has worked on application-layer traffic optimization to fill the gaps in routing policies between the layer 3 network and overlay networks by providing the underlay network topology and cost information to applications that build overlay networks. However, since ASes are operated under different policies and cost metrics for application-layer traffic optimization are assumed to be autonomously configured by each AS, there are considerations in applying application-layer traffic optimization techniques to cross-domain (inter-AS) traffic. This document summarizes general problems with cross-domain traffic of overlay networks and considerations on the AS-level application-layer traffic optimization from the viewpoint of inter-AS economics. This document also discusses the conceivable approaches to solve the problems and considerations.

#### Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

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 <a href="http://datatracker.ietf.org/drafts/current/">http://datatracker.ietf.org/drafts/current/</a>.

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

Asai, et al. Expires December 5, 2012

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 December 5, 2012.

Copyright Notice

Copyright (c) 2012 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 (<u>http://trustee.ietf.org/license-info</u>) 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

$\underline{1}$ . Introduction
<u>1.1</u> . Terminology
<u>1.1.1</u> . AS Relationships
<u>1.1.2</u> . Transit
<u>1.1.3</u> . Peering
<u>1.1.4</u> . Cross-Domain Links and Traffic
<u>1.1.5</u> . Overlay Network
<u>1.2</u> . Requirements Language
2. Cross-Domain Traffic Optimization Problems and
Considerations
2.1. General Problems with Cross-Domain Traffic of Overlay
Networks
2.2. Considerations on AS-Level Application-Layer Traffic
Optimization
<u>2.2.1</u> . Uneven policy configuration
<u>2.2.2</u> . Asymmetric economic policies <u>10</u>
<u>2.2.3</u> . Uncertain ingress routes
2.3. Summary of the Problems and Considerations 13
$\underline{3}$ . Conceivable Solution Approaches and Discussion $\underline{14}$
<u>3.1</u> . A1. ALTO servers without interconnection at local ASes <u>14</u>
3.2. A2. ALTO servers without interconnection at local and
remote ASes
3.3. A3. ALTO servers with mesh interconnection <u>16</u>
3.4. A4. ALTO servers with mesh interconnection and reverse
path probes $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\ldots$ $\frac{17}{2}$
3.5. A5. ALTO servers with hop-by-hop interconnection by
using a path vector algorithm
<u>3.6</u> . Summary of the Conceivable Solution Approaches <u>19</u>
<u>3.7</u> . Discussion on uneven policy configuration <u>19</u>
$\underline{4}$ . IANA Considerations
5. Security Considerations
<u>6</u> . Acknowledgements
$\underline{7}$ . Informative References
Authors' Addresses

#### **<u>1</u>**. Introduction

Many distributed systems, such as content delivery networks (CDNs) and live media streaming systems, have introduced application-layer overlay routing as represented by peer-to-peer technologies for their communication scheme to avoid excessive server load and to achieve effective and high-quality communication (e.g., high throughput, fault tolerance). Today, the traffic generated by these applications using application-layer overlay routing becomes a significant amount of the Internet traffic [RFC5693]. Since these applications construct their own network topology (a.k.a. overlay network) over the Internet, generally without taking into account the layer 3 network topology, these applications frequently utilize a larger amount of network resources than layer 3 network providers expect. Moreover, they may utilize a detoured path that is disallowed or unexpected by layer 3 network providers [H009].

The ALTO Working Group has worked on application-layer traffic optimization to fill the gaps between the layer 3 network and applications by providing the underlay network topology and cost information to these applications for their overlay network construction. However, there exist considerations on inter-AS policy conflicts when we deploy the AS-level application-layer traffic optimization because ASes are operated under different policies and cost metrics are assumed to be autonomously configured by each AS.

This document summarizes general problems with cross-domain (inter-AS) traffic generated by overlay networks and considerations on the AS-level application-layer traffic optimization from the viewpoint of inter-AS economics, which are not discussed in [RFC5693]. The main concerns on the AS-level application-layer traffic optimization are categorized to three groups; 1) uneven policy configuration between distinct layer 3 network domains (i.e., ASes), 2) asymmetric economic policies on transit links, and 3) uncertain ingress routes in BGP routing. The underlying problem inducing these concerns is that the detailed economic policies between interconnected ASes are non-disclosure information due to commercial contracts although these policies are important metrics for inter-AS traffic optimization. This document also discusses the conceivable approaches to solve the problems and considerations.

#### **<u>1.1</u>**. Terminology

We use the following terms in this document.

#### **<u>1.1.1</u>**. AS Relationships

AS relationships represent commercial relationships between interconnected ASes. AS relationships are categorized into two major types: transit and peering. There are typical inter-AS routing policies by each type of AS relationships [Wang03].

#### **<u>1.1.2</u>**. Transit

Transit is a type of AS relationships, in which a customer AS purchases Internet access from its transit provider(s) over transit link(s) by paying some amount of money according to the actual bandwidth usage. The Transit relationship is also called provider-customer relationship.

#### <u>1.1.3</u>. Peering

Peering is a type of AS relationships, in which two peering ASes are equal. Traffic exchanged over peering links is free of charge.

## **<u>1.1.4</u>**. Cross-Domain Links and Traffic

Domains correspond to ASes in this document. Cross-domain links and traffic denote inter-AS links and traffic, respectively.

#### **<u>1.1.5</u>**. Overlay Network

Overlay networks are constructed by application-layer nodes such as peer-to-peer application nodes over the layer 3 network (i.e., IP network) that is operated by network providers. The topology and routing of an overlay network are controlled by applications that build the overlay network but not by the layer 3 network providers.

## **<u>1.2</u>**. 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 <u>RFC 2119</u> [<u>RFC2119</u>].

## 2. Cross-Domain Traffic Optimization Problems and Considerations

This section discusses general problems with cross-domain traffic of overlay networks and considerations on the AS-level application-layer traffic optimization in terms of cross-domain traffic and its economics. First, we point out the general problems that the overlay networks do not take into account the layer 3 network economics and routing policies, and consequently, they cannot always utilize optimal network resources from the viewpoint of layer 3 network providers. Then, we present the considerations on the AS-level application-layer traffic optimization that are not discussed well in [RFC5693]. Since ASes are operated under different policies and cost metrics are assumed to be autonomously configured by each AS, the method for intra-domain optimization cannot be directly used and cross-domain concerns, such as uneven policy configuration between distinct ASes, asymmetric policies on transit links, and asymmetric underlay paths, must be solved to achieve the AS-level applicationlayer traffic optimization.

### 2.1. General Problems with Cross-Domain Traffic of Overlay Networks

The Internet consists of thousands of ASes operated by distinct network providers such as commercial ISPs, companies and universities. Each AS generally connects with multiple ASes, and there are distinct charging policies for each inter-AS link. These charging policies are roughly categorized into two major types of relationships; transit (with charge) and peering (without any charge). From the economic viewpoint, network providers want to reduce the traffic volume exchanged with transit providers as much as possible, and consequently, they manage BGP routing policies as explained in [Wang03]. Thus, the layer 3 paths are determined by network providers as intended.

However, overlay networks are not sometimes aware of these routing policies and generate more expensive cross-domain traffic. On the other hand, network providers cannot optimize the cross-domain traffic generated by applications on overlay networks. This is because the traffic is controlled by a set of application-specific algorithms that determines the overlay network topology and traffic delivery paths, such as peer, neighbor, or path selection algorithms.

+----+ provider | AS 1 |-----+ provider +----+ | transit | transit | v v v customer +----+ peering +----+ +----+ customer | AS 2 |<---->| AS 3 | | AS 4 | +----+ + +---++

AS 2 purchases Internet access from AS 1 via a transit link. On the contrary, the link between AS 2 and AS 3 is peering, which is a lower cost link from the viewpoint of AS 2 network operators.

Figure 1: An example of AS-level topology with AS relationships

We show an example of the problem with the unawareness of the layer 3 network economics and cross-domain traffic generated by overlay networks. An example of interconnections of ASes and their relationships is shown in Figure 1. Suppose remote nodes of a peer-to-peer content delivery network that provide a certain content file exist in both AS 3 and AS 4, and a local node that downloads the file in AS 2 is to retrieve the file from one of these remote nodes, a remote node in AS 3 should be selected to reduce transit charge for both ASes of the local node and the remote node, but today's peer-to-peer content delivery networks that are unaware of AS relationships often select other remote nodes.

Thus, overlay networks cannot always utilize optimal network resources but high-cost ones (i.e., transit links from/to transit providers) from the economic viewpoint of network providers. Moreover, especially on peer-to-peer overlay networks, the connectivity of most of end-point nodes (i.e., peers) is provided by residential ISPs, and most of residential ISPs are not transit providers but transit customers, and consequently, it is significantly important to control the transit traffic not to increase their charge to their providers though these kinds of application-layer traffic are hardly controlled by network providers. [<u>RFC5693</u>] also claims this problem with cross-domain traffic in terms of transit cost as well as congestion in intra-domain networks.

+----+ provider | AS 1 |-----+ provider +----+ | transit | transit V v customer +----+ peering +----+ peering +----+ customer | AS 2 |<---->| AS 3 |<---->| AS 4 | +---+ +---+ +---+ Node <====> Node <====> Node detoured path

According to the typical BGP routing policies, the path from AS 2 to AS 4 is to be AS 2->AS 1->AS 4. The path AS 2->AS 3->AS 4 is not usually allowed because AS 3 relays traffic from AS 2 to AS 4 without any charge if this path is allowed.

Figure 2: An example of AS-level detouring by overlay networks

Another problem with overlay networks is that overlay networks may create a detoured path that is disallowed or unexpected by layer 3 network providers [<u>Ho09</u>]. For example, in Figure 2, the traffic from AS 2 to AS 4 can pass through AS 3 if a node of an overlay network exists in AS 3 and relays the traffic, but this is usually disallowed by the routing policy of AS 3 to avoid free-riding.

In summary, overlay networks have following problems from the crossdomain economic viewpoint of network providers.

- o Overlay networks usually do not take into account the layer 3 network economics to construct their network topology and to exchange traffic between end-point nodes. Therefore, they cannot always utilize optimal cross-domain network resources from the economic viewpoint of network providers.
- o Overlay networks may enable AS-level traffic detouring that is disallowed by the layer 3 network routing policies. This problem possibly increases transit expenses or induces free-riding.

#### **2.2**. Considerations on AS-Level Application-Layer Traffic Optimization

The ALTO Working Group has worked on application-layer traffic optimization to fill the gaps in routing policies between the layer 3 network and overlay networks. It has worked on solving the problems stated in [RFC5693], but [RFC5693] misses some considerations on the AS-level (cross-domain) application-layer traffic optimization. We summarize the missing considerations as follows.

Internet-Draft

- o Uneven policy configuration between distinct administrative domains: ASes hardly cooperate with each other in fairly regulating uneven policies of distinct ASes because each AS operates its network under its own policy and inter-AS policies are complicated to achieve total optimization.
- Asymmetric economic policies on transit links: It is difficult to regulate the asymmetric economic policies on transit links because transit customers' policies run counter to transit providers; i.e., customers want to reduce the traffic exchanged with their providers to reduce their expense though providers want to increase the traffic exchanged with their customers to increase their income.
- o Uncertain ingress routes in BGP routing: In addition to the asymmetric economic policies on transit links, the policy-based routing of BGP often creates asymmetric paths. Local ASes hold their egress routes (i.e., next hop AS for a remote end-point node), but they do not know ingress routes (i.e., previous hop AS from a remote end-point node). Therefore, it is difficult to configure ingress cost for remote end-point nodes/networks because the layer 3 network providers do not have the ingress routes to determine which AS becomes the previous hop for the remote end-point nodes/networks.

The details of these considerations are described in the following sections.

## **<u>2.2.1</u>**. Uneven policy configuration

The ALTO Working Group has proposed a protocol to distribute end-toend network cost between peers to

applications [I-D.ietf-alto-protocol]. This protocol does not intend to define the cost computation algorithm, but it assumes that the cost is computed by network providers. Two oracle-based cost computation algorithms, [<u>Aggarwal07</u>] and [<u>Xie08</u>], have been proposed and evaluated in the research area. [Aggarwal07] computes the ASlevel cost according to AS hop count between two end-point nodes. So, it ignores the information on AS relationships (i.e., transit cost). [Xie08] computes the AS-level cost according to the configured parameters (e.g., `local preference' in BGP) in routers. This takes into account AS relationships. However, there is a problem with this algorithm when it is applied to the real Internet. Charging policies for exchanged inter-AS traffic volume are so complicated that different ASes hardly cooperate with each other in computing and fairly balancing cost. Even in the BGP routing managed by network providers, the hot potato problem stated in [RFC4277] shows the difficulty in regulating policies of distinct ASes.

	++ provider					
	AS 1	+				
provider	++ 3	transit				
	5   transit (1\$/Mbps)	(2\$/Mbps)				
	30 v	v 10				
customer	++	++ customer				
	AS 2	AS 3				
	++	++				

Each number represents egress cost. Transit charge on each transit link is noted in a set of parentheses.

Figure 3: An example of uneven cost configuration

For example, suppose egress cost of each inter-AS link is configured autonomously (i.e., each AS sets cost according to its own policies) as shown in Figure 3, then the cost of the path from AS 2 to AS 1 becomes larger than that of the path from AS 3 to AS 1 though the path from AS 2 to AS 1 seems to be a cheaper link than the other. Thus, oracle-based approaches are exposed to a fairness or evenness issue among multiple autonomous domains.

In summary, inter-AS policies are so complex that ASes cooperate with each other in fairly regulating policies of distinct ASes in terms of cross-domain cost configuration.

### **<u>2.2.2</u>**. Asymmetric economic policies

There is a difficulty in regulating the asymmetric economic policies on inter-AS links, especially between transit customers and providers. One of the causes of this difficulty is same as that of the issue on the uneven policy configuration; i.e., because each AS configures its own desired policy. Another cause of this difficulty is that the policies on the transit links are not symmetric. So, one party's policy does not match the other's. The same asymmetric nature is also found in the BGP routing, but the policy regulation on transit links becomes more complex in the overlay routing than the BGP routing. This is because overlay network nodes that have the same functionality or contents possibly exist in multiple ASes while the functionality (i.e., end-to-end connectivity to the destination) of the BGP routing is mapped to a single AS. The BGP path-vector routing is performed under the control of the layer 3 network providers by route import and export policies, and consequently, the computed paths in the BGP routing are based on the benefit principle. However, the overlay routing might destroy the control and the principle.

+---+ | AS 1 | Remote node provider +----+ Λ 5 | transit | Good for AS 1 30 V | Not good for AS 2 customer +----+ V | AS 2 | Local node provider +----+ Λ 5 | transit | Good for AS 2 30 V | Not good for AS 3 customer +----+ V | AS 3 | Remote node +---+

Each number in this figure represents cost. Note that cost for each type of AS relationships is already simplified and regulated here; 5 for provider to customer and 30 for customer to provider. This asymmetric cost regulation follows the typical import policy in the BGP routing (i.e., local preference) although it is quite difficult to achieve this regulation between distinct ASes.

Figure 4: An example of asymmetric cost configuration

For example, suppose the cost of each inter-AS link configured as shown in Figure 4 is egress cost, then the end-to-end cost from AS 1 to AS 2 becomes smaller than that from AS 3 to AS 2. On the other hand, suppose the cost of each inter-AS link configured as shown in Figure 4 is ingress cost, then the end-to-end cost from AS 1 to AS 2 becomes larger than that from AS 3 to AS 2. This means that the path from AS 3 to AS 2 is preferred than the other from the viewpoint of AS 2 but the path from AS 1 to AS 2 is preferred than the other from the viewpoint of AS 1 and AS 3. Unlike the BGP routing, overlay networks may have the same functionality or contents at their nodes both in AS 1 and AS 3, and consequently, it is required to consider the conflicts of asymmetric economic policies on transit links between multiple ASes.

In summary, the existence of identical functionality at multiple ASes in overlay networks raises the problem with asymmetric economic policies on transit links.

#### **2.2.3**. Uncertain ingress routes

In case that cost between end-point nodes or networks is configured by each AS, underlay paths must also be considered to compute the cost. Since local ASes hold their egress routes (i.e., next hop AS for a remote end-point node/network), egress cost can be configured by looking at the cost of the link to the next hop AS. However,

ingress cost is difficult to be computed because local ASes do not have ingress routes to determine which AS becomes the previous hop for a remote end-point node/network. Note that the ingress routes are difficult to be expected from the eqress routes because the policy-based routing of BGP often creates asymmetric paths. BGP operators usually work on traffic engineering with reactive methods for the ingress traffic.

provider +----+ provider +-----| AS 1 |-----+ transit | +----+ | transit V customer +----+ peering +----+ | AS 2 |<---->| AS 3 | +----+ provider | Routing table @AS 2 transit | +----+ V | Destination | Next | +----+ customer | +----+ | AS 4 | | AS 1 | AS 1 | +----+ provider | | AS 3, 4, 5 | AS 3 | transit | +----+ | +-----+ V V Routing table @AS 5 +----+ customer +----+ | AS 5 | | Destination | Next | +----+ +----+ | AS 1, 2 | AS 1 | | AS 3, 4 | AS 4 | +----+

The underlay path between AS 2 and AS 5 is asymmetric. From AS 2 to AS 5, the path goes through AS 3 and AS 4 as the routes from peering are generally preferred to those from transit providers. On the other hand, from AS 5 to AS 2, the path goes through AS 1 as the shortest paths are generally selected as the best when both exits are transit. Routing tables at AS 2 and AS 5 are also represented in this figure to point out the uncertainness of ingress routes.

Figure 5: An example of uncertain ingress routes with asymmetric paths in the BGP routing

For example, the BGP routing often creates asymmetric paths as shown in Figure 5. Suppose that AS 2 is the local AS and AS 5 is one of the remote ASes and the network provider of AS 2 tries to configure ingress cost from AS 5, AS 2 holds an egress route to AS 5 that goes to AS 3 over a peering link in its routing table but AS 2 does not have any ingress routes, and consequently, AS 2 cannot determine

which AS (AS 1 or AS 3) is the previous hop from AS 5 because the reverse path is determined by routing tables of other ASes. In this figure, the path between AS 2 and AS 5 is asymmetric, and consequently, the next hop AS to AS 5 (AS 3) becomes different from the previous hop AS to AS 5 (AS 1).

### **<u>2.3</u>**. Summary of the Problems and Considerations

This section summarizes the general problems with cross-domain traffic of overlay networks pointed out in <u>Section 2.1</u> and considerations on the AS-level application-layer traffic optimization described in <u>Section 2.2</u>. These problems and considerations of overlay networks towards the AS-level application-layer traffic optimization are summarized into the following five categories.

- o C1. AS relationships unawareness of overlay networks
- o C2. AS-level detouring by overlay networks
- o C3. Uneven policy configuration of distinct ASes
- o C4. Asymmetric economic policy on transit links
- o C5. Uncertain ingress routes in determining ingress cost

# 3. Conceivable Solution Approaches and Discussion

This section discusses the conceivable approaches to solve the problems and considerations, C1 to C5, summarized in <u>Section 2.3</u>. This document does not intend to define specifications and protocols, and consequently, this section shows the overview of each approach to open up further discussion. We assume that the cost or policy information is provided to applications via ALTO servers defined in [<u>RFC5693</u>]. The conceivable solution approaches are listed as follows.

- o A1. ALTO servers without interconnection at local ASes
- o A2. ALTO servers without interconnection at local and remote ASes
- o A3. ALTO servers with mesh interconnection
- o A4. ALTO servers with mesh interconnection and reverse path probes
- o A5. ALTO servers with hop-by-hop interconnection by using a path vector algorithm

The detail and points to be solved of each approach are given in the following sections.

## 3.1. A1. ALTO servers without interconnection at local ASes

The simplest way to take into account the AS relationships in overlay networks is that network providers configure cost information for end-point nodes in other ASes as well as end-point nodes in a local AS using ALTO servers. In this approach, a local end-point node discovers and uses ALTO servers in the local AS. . A local AS can configure egress cost for end-point nodes in other ASes to its ALTO servers by looking at its routing table and the inter-AS policy to the next hop AS, but cannot configure ingress cost because of the same reason as C5. The ALTO servers in the local AS do not have ingress and egress cost of remote ASes.

This approach does not require new specifications in comparison with intra-domain application-layer traffic optimization. This approach does not completely solve any of C1-C5, but this is ``better than nothing'' for C1. The inter-AS traffic from the local AS to other ASes can be optimized by looking at the egress cost configured in the local AS' ALTO servers.

+----+ +----+ Remote AS | | Local AS | +----+ \_\_\_ | ALTO server | / \ / \ +----+ | R |===| R | \* egress cost ∖\_\_\_/ | +----+ | +----+ | | | Local node |>==== Optimized for Local AS ===>| Remote node | | | +-----+ | +-----+ | +----+ +---------+

Figure 6: ALTO servers without interconnection at local ASes

The system overview of this approach is shown in Figure 6.

## **3.2.** A2. ALTO servers without interconnection at local and remote ASes

To extend the optimization to remote ASes from A1, a local end-point node can look up ALTO servers in remote ASes in addition to those in the local AS. In this approach, a local end-point node discovers and uses ALTO servers both in the local and remote ASes. Both local and remote ASes can configure their egress cost by the same way as A1.

This approach requires a new specification, in comparison with intradomain application-layer traffic optimization, to discover ALTO servers in remote ASes. To simplify the specification and the discovery procedure, this discovery may be proxied by remote endpoint nodes. This approach does not completely solve any of C1-C5, but this is also ``better than nothing'' and better than A1 for C1. The inter-AS traffic from a remote AS to other ASes as well as the inter-AS traffic from the local AS to other ASes can be optimized by looking at the egress cost configured in both ASes' ALTO servers.

+----+ +----+ | | Local AS Remote AS | +----+ \_\_\_ +----+ +->| ALTO server | / \ / \ +->| ALTO server | | +----+ | R |===| R | | +----+ | \* egress cost \\_\_\_/ \\_\_\_/ | \* egress cost +----+ 1 1 1 1 | +-----+>==== Optimized for Local AS ===>+----++ | | | Local node | | | Remote node | | | +----+<===== Optimized for Remote AS ==<+----+ | +----+ +-------+

Figure 7: ALTO servers without interconnection at local and remote ASes

The system overview of this approach is shown in Figure 7.

#### **3.3.** A3. ALTO servers with mesh interconnection

A1 and A2 do not completely solve any of C1-C5. Moreover, A2 requires a specification to discover ALTO servers of remote ASes. This approach adds a specification to exchange egress cost between ALTO servers of local and remote ASes to provide eqress cost of both ASes to overlay networks, instead of the discovery of ALTO servers of remote ASes. The inter-AS traffic from the local AS to others and the inter-AS traffic from a remote AS to other ASes are optimized by this approach.

One problem with this approach is scalability that the ALTO servers of the local AS must establish the interconnection with ALTO servers of remote ASes to exchange egress cost configuration.

+		+	+		+
Local AS	+ ex	change e	gress cost	+	Remote AS
1					1
	V		I	V	
+	+			+	+
+->  ALTC	server	/ \	/ \	ALTO s	erver
+	+	R  =:	==  R	+	+
* eg	ress cost	\/	\/	* egre	ss cost
+	+>==== 0p	timized <sup>.</sup>	for Local	AS ===>+	+
Local node					Remote node
+	+<==== 0p	timized <sup>.</sup>	for Remote	e AS ==<+	+
+		+	+		+

Internet-Draft

Figure 8: ALTO servers with mesh interconnection

The system overview of this approach is shown in Figure 8.

#### **3.4.** A4. ALTO servers with mesh interconnection and reverse path probes

Any of A1, A2, and A3 cannot achieve the optimization of ingress inter-AS traffic (i.e., cannot solve C5) because these approaches cannot determine the ingress routes. In this approach, ALTO servers resolve ingress routes with reverse path probes between these interconnected ALTO servers. This approach requires a specification to resolve the ingress routes in addition to A3's specification, which is to exchange egress cost between ALTO servers of local and remote ASes. Thus, this approach provides four types of cost to applications; 1) the ingress cost of a local AS, 2) the egress cost of a local AS, 3) the ingress cost of a remote AS, and 4) the egress cost of a remote AS. Applications can optimize the inter-AS traffic for both local and remote ASes using these four types of cost according to benefit principle; for example, attaching weight to the ingress and egress cost of a local AS and detaching weight on the ingress and egress cost of a remote AS if a end-point node in the remote AS is the beneficiary, vice versa. In summary, this approach solves C1, C4, and C5. Note that this approach has the same scalability problem as A3.

+	+	+	+
Local AS +	exchange	cost+	Remote AS
	I		
V <	<path discov<="" td=""><td>ery probe&gt; v</td><td></td></path>	ery probe> v	
+	+	+	+
+->  ALTO server	/ \	/ \   ALTO se	rver
+	-+   R  ==	=  R   +	+
* egress cos	st \/	\/ * egres	s cost
* ingress co	ost	* ingre	ss cost
++>======	== Optimized	for both ====>+-	+
Local node	I		Remote node
++<======	== Optimized	for both =====<+-	+
+	+	+	+

Figure 9: ALTO servers with mesh interconnection and reverse path probes

The system overview of this approach is shown in Figure 9.

# 3.5. A5. ALTO servers with hop-by-hop interconnection by using a path vector algorithm

A1-A4 cannot solve C2 because these simple cost-based approaches have limitations to reflect inter-AS policies in the BGP routing while BGP is a path-vector routing protocol that is one of policy-based routing protocols. A solution approach to solve C2 is to introduce a pathvector algorithm that propagates policies hop-by-hop between interconnected ALTO servers like the BGP routing, and to provide cost for detoured paths. This approach enables flexible policy propagation, but requires many new specifications, protocols, and algorithms to propagate policies and to compute a cost map that is provided to applications.



| AS 3 via AS 2 | (AS 2) | 30 | 30 | +------------+

Figure 10: ALTO servers with hop-by-hop interconnection by using a path vector algorithm

An example of topology and a cost map adopting this approach is shown in Figure 10. ALTO servers of each AS establish interconnection, and then they exchange and propagate policies with each other by a path-

vector algorithm. The ALTO server of the local AS then computes a cost map for sources/destinations and detoured paths from all the received policies. Thus, this approach can provide the cost for detoured paths as well as the layer 3 network paths. In this example, the cost of the detoured path from/to AS 3 via AS 2 becomes lower than that of the end-to-end path from/to AS 3 via AS 1. This means that the AS 2 does not disallow the detoured path. If AS 2 does not prefer detouring, the cost of the detoured path becomes higher than the others according to the propagated policies. Note that the sources/destination and detoured paths are aggregated by per-AS in this example, but they may be per-prefix.

### **3.6**. Summary of the Conceivable Solution Approaches

+	+	+	•	+	-+		+	-+
Approach	C1		C2	C3		C4	C5	Ι
+	+	+		+	-+		+	-+
A1	x				I			I
1	I	Ι					1	
A2	x	i		*	i		i	İ
A3	x			*	T			Ι
							1	
A4	X	Ì		*	Í	х	x	Ì
A5	x		х	*	T	х	x	Ι
+	+	+		+	-+		+	-+

x: Solved or better than nothing; \*: Solved or better than nothing with additional methods discussed in Section 3.7

Table 1: Summary of the conceivable solution approaches

## 3.7. Discussion on uneven policy configuration

Any of A1 A5 do not provide a method to solve C3. This section discusses possible methods to solve C3. We first note that A1 never achieves to solve C3 because it only provides the egress cost from the local AS. The possible methods to solve C3 are listed and summarized as follows.

o Global regulation: This method is to define global regulation for cost configuration, and to provide a cost computation function; for example, cost X for X\$/Mbps, or more simply cost 2, 1, and 0 for transit (from/to provider), peering, and transit (from/to customer), respectively. However, it is not realistic because the inter-AS policy is too complex to define the global regulation. Moreover, the economic policies between interconnected ASes cannot

be disclosed by each network provider due to commercial contracts.

- o Hop-by-hop regulation: In A3 A5, ALTO servers of each AS establish interconnection. In these cases, the interconnected ALTO servers can introduce a regulation algorithm for the exchanged cost. The cost will be regulated among interconnected ALTO servers and this might be ``better than nothing'', but C3 still exists among ALTO servers that are not interconnected.
- o Inference-based regulation: AS relationships inference can be one of the methods to regulate the complex economic policies. This is an alternative to the global regulation method to solve the problem with non-disclosure AS relationships. AS relationships inference algorithms have been proposed in the research field, such as [Asai10], [Dimitropoulos07], and [Ga001]. Global regulation is achieved by using these inference algorithms. The inferred AS relationships for some links may not be accurate, but this might also be ``better than nothing''.

In summary, within the interconnected ALTO servers, the hop-by-hop regulation is the best. The inference-based regulation can be used to assist in regulating uneven policy configuration in the other cases or finding anomalous cost configurations.

# **<u>4</u>**. IANA Considerations

No need to describe any request regarding number assignment.

# **<u>5</u>**. Security Considerations

This document is neither a requirements document nor a protocol specification. However, since the solution approaches exchange the inter-AS economic policies with ALTO servers operated by other ASes (i.e., external network domains), two security considerations are discussed as follows.

- o The ALTO servers operated by other ASes may falsify the received cost map or policies. The protocol specifications of the solution approaches must include anti-falsification and verification mechanisms (e.g., digital signing) for the exchanged cost map or policies.
- o The exchanged cost map or policies may contain the non-disclosure inter-AS information. The protocol specifications of the solution approaches should consider the schemes to aggregate and filter the exchanged cost map or policies in order not to reveal the nondisclosure information.

# <u>6</u>. Acknowledgements

Moritz Steiner (Bell-Labs), Piotr Wydrych (AGH University of Science and Technology), Russ White (Cisco Systems), Stefano Previdi (Cisco Systems), Volker Hilt (Alcatel-Lucent Bell-Labs), and many others provided informative discussions and valuable comments.

## 7. Informative References

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, March 1997.
- [RFC4277] McPherson, D. and K. Patel, "Experience with the BGP-4 Protocol", <u>RFC 4277</u>, January 2006.
- [RFC5693] Seedorf, J. and E. Burger, "Application-Layer Traffic Optimization (ALTO) Problem Statement", <u>RFC 5693</u>, October 2009.
- [I-D.ietf-alto-protocol]

Alimi, R., Penno, R., and Y. Yang, "ALTO Protocol", <u>draft-ietf-alto-protocol-11</u> (work in progress), March 2012.

## [Aggarwal07]

- Aggarwal, V., Feldmann, A., and C. Scheideler, "Can ISPs and P2P users cooperate for improved performance?", SIGCOMM Comput. Commun. Rev., vol. 37, no. 3, pp. 29-40, 2007.
- [Asai10] Asai, H. and H. Esaki, "Estimating AS Relationships for Application-Layer Traffic Optimization", 3rd Workshop on Economic Traffic Management, LNCS Vol. 6236, pp. 51-63, 2010.
- [Dimitropoulos07]

Dimitropoulos, X., Krioukov, D., Fomenkov, M., Huffaker, B., Hyun, Y., claffy, k., and G. Riley, "AS Relationships: Inference and Validation", ACM SIGCOMM Comput. Commun. Rev., Vol. 37, No. 1, pp. 29-40, 2001.

- [Gao01] Gao, L., "On inferring autonomous system relationships in the Internet", IEEE/ACM Transactions on Networking, Vol. 9, No. 6, pp. 733-745, 2001.
- [Ho09] Ho, Haddow, T., Ledlie, J., Draief, D., and P. Pietzuch, "Deconstructing internet paths: an approach for AS-level detour route discovery", Proceedings of the 8th international conference on Peer-to-peer systems, p. 6, 2009.
- [Wang03] Wang, F. and L. Gao, "On Inferring and Characterizing Internet Routing Policies", IMC '03: Proceedings of the 3rd ACM SIGCOMM conference on Internet measurement, pp. 15-26, 2003.

Internet-Draft Considerations on AS-Level ALTO June 2012

[Xie08] Xie, H., Yang, Krishnamurthy, A., Liu, and A. Silberschatz, "P4P: provider portal for applications", SIGCOMM '08: Proceedings of the ACM SIGCOMM 2008 conference on Data communication, pp. 351-362, 2008.

Authors' Addresses

Hirochika Asai The University of Tokyo 7-3-1 Hongo Bunkyo-ku, Tokyo 113-8656 JP

Phone: +81 3 5841 6748 Email: panda@hongo.wide.ad.jp

Hiroshi Esaki The University of Tokyo 7-3-1 Hongo Bunkyo-ku, Tokyo 113-8656 JP

Phone: +81 3 5841 6748 Email: hiroshi@wide.ad.jp

Tsuyoshi Momose Cisco Systems G.K. 2-1-1 Nishi-Shinjuku Shinjuku-ku, Tokyo 163-0409 JP

Phone: +81 3 5324 4154 Email: tmomose@cisco.com