

PCE Working Group  
Internet-Draft  
Intended status: Informational  
Expires: February 28, 2014

D. Dhody  
Y. Lee  
Huawei Technologies  
N. Ciulli  
Nextworks  
L. Contreras  
O. Gonzalez de Dios  
Telefonica I+D  
August 27, 2013

**Cross Stratum Optimization enabled Path Computation  
draft-dhody-pce-cso-enabled-path-computation-04**

**Abstract**

Applications like cloud computing, video gaming, HD Video streaming, Live Concerts, Remote Medical Surgery, etc are offered by Data Centers. These data centers are geographically distributed and connected via a network. Many decisions are made in the Application space without any concern of the underlying network. Cross stratum application/network optimization focus on the challenges and opportunities presented by data center based applications and carriers networks together [[CSO-DATACNTR](#)].

Constraint-based path computation is a fundamental building block for traffic engineering systems such as Multiprotocol Label Switching (MPLS) and Generalized Multiprotocol Label Switching (GMPLS) networks. [[RFC4655](#)] explains the architecture for a Path Computation Element (PCE)-based model to address this problem space.

This document explains the architecture for CSO enabled Path Computation.

**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 <http://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 February 28, 2014.

## Copyright Notice

Copyright (c) 2013 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 (<http://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

<a href="#">1.</a>	Introduction . . . . .	<a href="#">2</a>
<a href="#">1.1.</a>	Requirements Language . . . . .	<a href="#">5</a>
<a href="#">2.</a>	Terminology . . . . .	<a href="#">5</a>
<a href="#">3.</a>	CS0 enabled PCE Architecture . . . . .	<a href="#">6</a>
<a href="#">4.</a>	Path Computation and Setup Procedure . . . . .	<a href="#">9</a>
<a href="#">4.1.</a>	Path Setup Using NMS . . . . .	<a href="#">11</a>
<a href="#">4.2.</a>	Path Setup Using a Network Control Plane . . . . .	<a href="#">11</a>
<a href="#">4.3.</a>	Path Setup using PCE . . . . .	<a href="#">12</a>
<a href="#">5.</a>	Other Consideration . . . . .	<a href="#">13</a>
<a href="#">5.1.</a>	Inter-domain . . . . .	<a href="#">13</a>
5.1.1.	One Application Domain with Multiple Network Domains	13
5.1.2.	Multiple Application Domains with Multiple Network Domains . . . . .	<a href="#">14</a>
<a href="#">5.1.2.1.</a>	ACG talks to multiple NCGs . . . . .	<a href="#">14</a>
5.1.2.2.	ACG talks to the primary NCG, which talks to the other NCG of different domains . . . . .	<a href="#">15</a>
<a href="#">5.2.</a>	Bottleneck . . . . .	<a href="#">16</a>
<a href="#">6.</a>	IANA Considerations . . . . .	<a href="#">16</a>
<a href="#">7.</a>	Security Considerations . . . . .	<a href="#">16</a>
<a href="#">8.</a>	Manageability Considerations . . . . .	<a href="#">16</a>
<a href="#">9.</a>	Acknowledgements . . . . .	<a href="#">16</a>
<a href="#">10.</a>	References . . . . .	<a href="#">16</a>
<a href="#">10.1.</a>	Normative References . . . . .	<a href="#">16</a>
<a href="#">10.2.</a>	Informative References . . . . .	<a href="#">16</a>

## **[1.](#) Introduction**

Many application services offered by Data Center to end-users make significant use of the underlying networks resources in the form of



bandwidth consumption used to carry the actual traffic between data centers and/or among data center and end-users. There is a need for cross optimization for both network and application resources.

[[CS0-PROBLEM](#)] describes the problem space for cross stratum optimization.

[NS-QUERY] describes the general problem of network stratum (NS) query in Data Center environments. Network Stratum (NS) query is an ability to query the network from application controller in Data Centers so that decision would be jointly performed based on both the application needs and the network status. Figure 1 shows typical data center architecture.

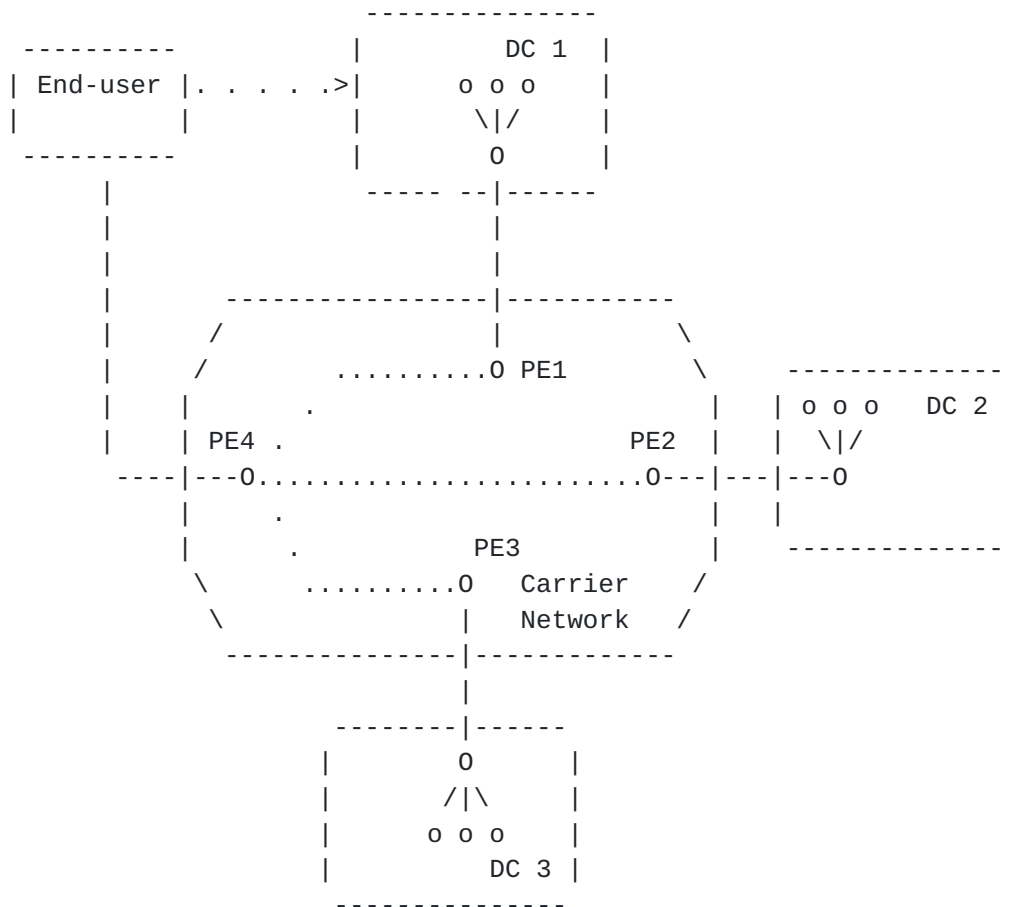


Figure 1: Data Center Architecture

Figure 2 shows the context of NS Query within the overarching data center architecture shown in Figure 1.



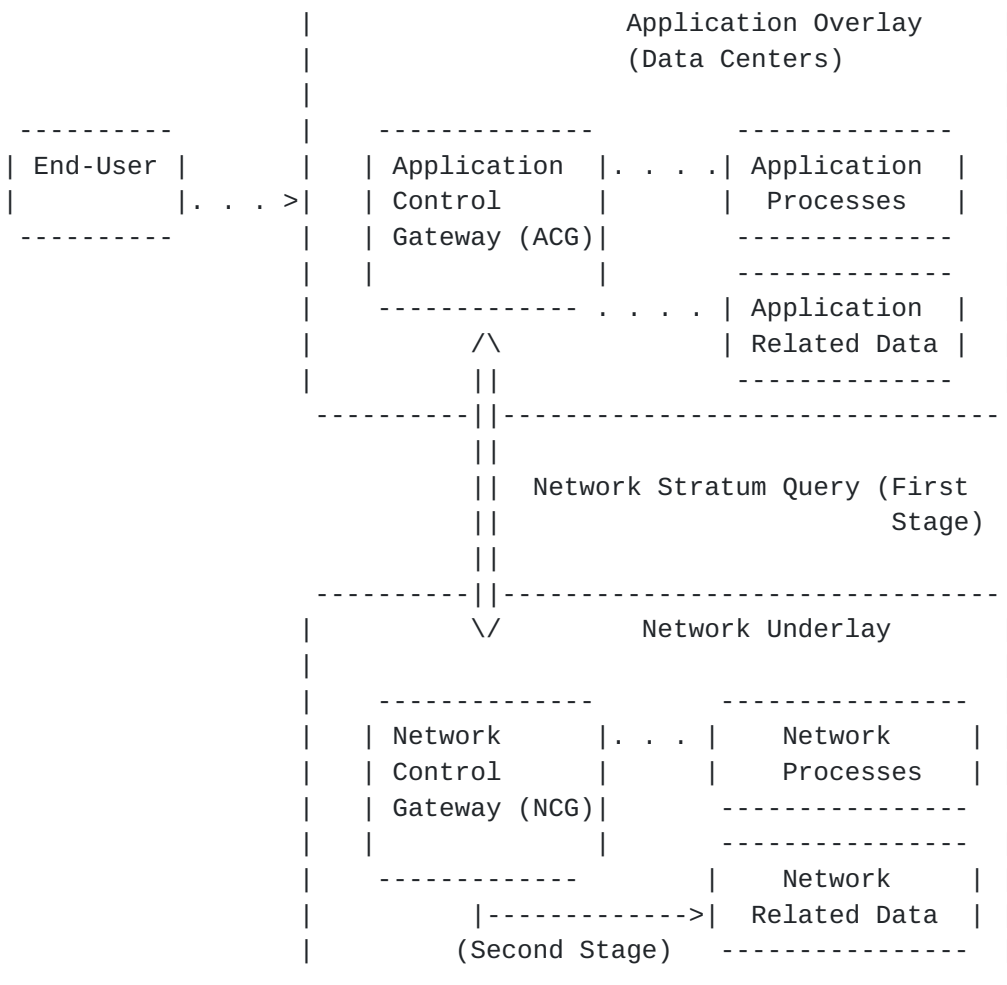


Figure 2: NS Query Architecture

NS Query is a two-stage query that consists of two stages:

- o A vertical query capability where an external point (i.e., the Application Control Gateway (ACG) in Data Center) will query the network (i.e., the Network Control Gateway (NCG)). The query can be initiated either by ACG to NCG or NCG to ACG depending on the mode of operation. ACG initiated query is an application-centric mode while NCG initiated query is a network-centric mode. It is anticipated that either ACG or NCG can be a final decision making point that chooses the end-to-end resources (i.e., both application IT resources and the network connectivity) depending on the mode of operation.
- o A horizontal query capability where the NCG gathers the collective information of a variety of horizontal schemes implemented in the network stratum.



As an example for vertical query (1st stage), [[ALTO-APPNET](#)] describes Application Layer Traffic Optimization (ALTO) information model and protocol extensions to support application and network resource information exchange for high bandwidth applications in partially controlled and controlled environments as part of the infrastructure to application information exposure (i2aex) initiative.

For the horizontal query (2nd stage), PCE can be an ideal choice, [[CS0-PCE-REQT](#)] describes the general requirement PCE should support in order to accommodate CS0 capability. This document is intended to fulfill the general PCE requirements discussed in the aforementioned reference.

This document describes how PCE Architecture as described in [[RFC4655](#)] can help in the second stage of NS query.

### **1.1. 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 [[RFC2119](#)].

## **2. Terminology**

The following terminology is used in this document.

ACG: Application Control Gateway.

Application Stratum: The application stratum is the functional block which manages and controls application resources and provides application resources to a variety of clients/end-users. Application resources are non-network resources critical to achieving the application service functionality. Examples include: application specific servers, storage, content, large data sets, and computing power. Data Centers are regarded as tangible realization of the application stratum architecture.

ALTO: Application Layer Traffic Optimization.

CS0: Cross Stratum Optimization.

GMPLS: Generalized Multiprotocol Label Switching.

i2aex: Infrastructure to application information exposure.

LSR: Label Switch Router.

MPLS: Multiprotocol Label Switching.





NCG: Network Control Gateway.

Network Stratum: The network stratum is the functional block which manages and controls network resources and provides transport of data between clients/end-users to and among application resources. Network Resources are resources of any layer 3 or below (L1/L2/L3) such as bandwidth, links, paths, path processing (creation, deletion, and management), network databases, path computation, admission control, and resource reservation capability.

NMS: Network Management System

PCC: Path Computation Client: any client application requesting a path computation to be performed by a Path Computation Element.

PCE: Path Computation Element. An entity (component, application, or network node) that is capable of computing a network path or route based on a network graph and applying computational constraints.

PCEP: Path Computation Element Communication Protocol.

TE: Traffic Engineering.

TED: Traffic Engineering Database.

UNI: User Network Interface.

### **3. CS0 enabled PCE Architecture**

In the network stratum, the Network Control Gateway (NCG) serves as the proxy gateway to the network. The NCG receives the query request from the ACG, probes the network to test the capabilities for data flows to/from particular points in the network, and gathers the collective information of a variety of horizontal schemes implemented in the network stratum. This is a horizontal query (Stage 2 in Figure 2).

In this section we will describe how PCE fits in this horizontal scheme.

A Path Computation Element (PCE) is an entity that is capable of computing a network path or route based on a network graph, and of applying computational constraints during the computation.

(1) NCG and PCE are co-located.



In this composite solution, the same node implements functionality of both NCG and PCE. When a network stratum query is received from the ACG (stage 1), this query is broken into one or more Path computation requests and handled by the PCE functionality co-located with the NCG. There is no need for PCEP protocol here. In this case, an external PCE interface (e.g., CLI, SNMP, proprietary) needs to be supported. This is out of the scope of this document.

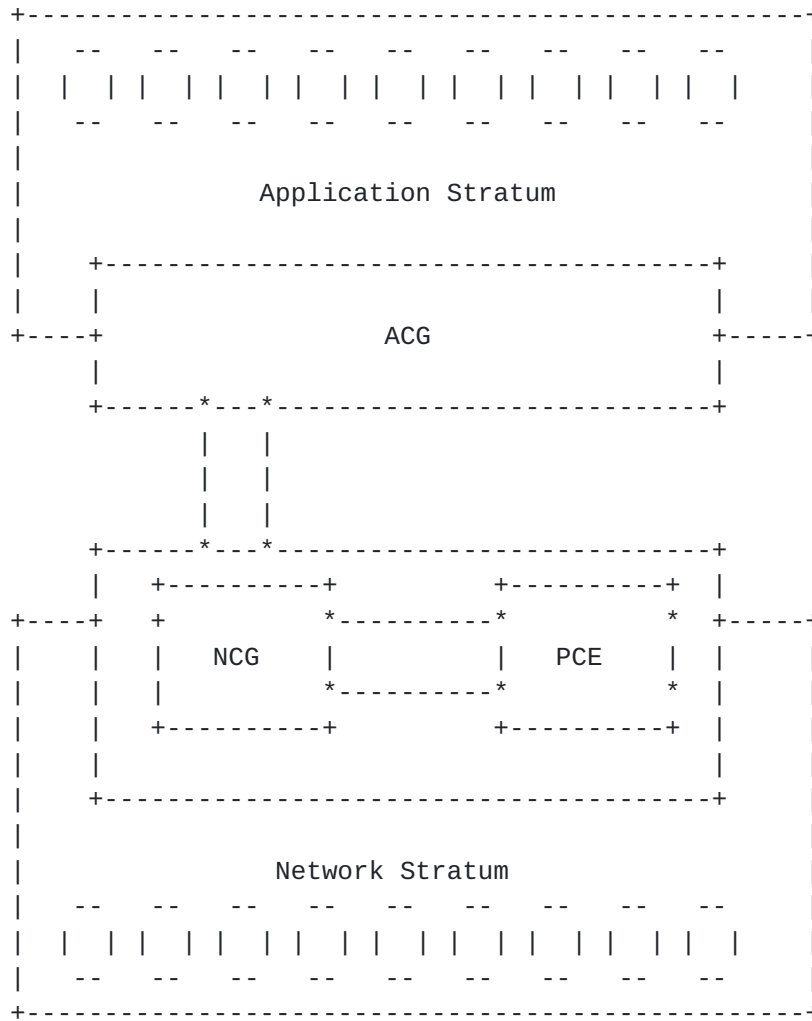


Figure 3: NCG and PCE Collocated

## (2) NCG and external PCE

In this solution, an external node implements PCE functionality. Network stratum query received from the ACG (stage 1) is converted into Path computation requests at the NCG and relayed to the external PCE using the PCEP [[RFC5440](#)]. In this case the NCG includes Path Computation Client (PCC) functionalities.



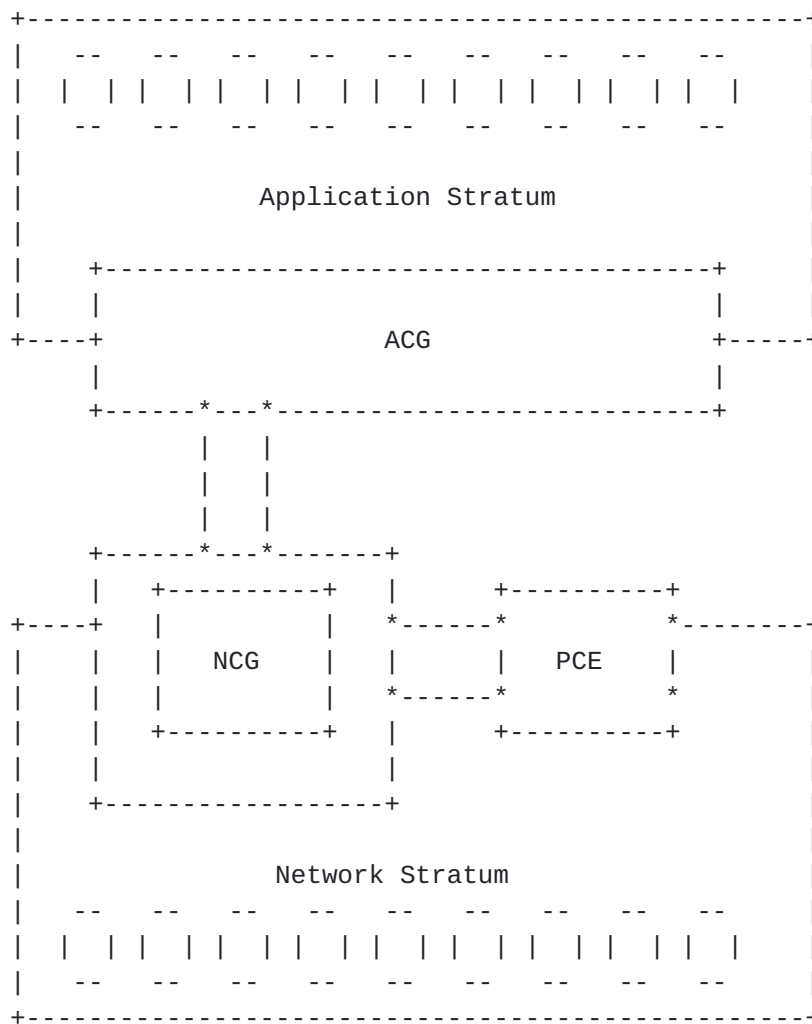


Figure 4: NCG and external PCE

PCE has the capability to compute constrained paths between a source and one or more destination(s), optionally providing the value of the metrics associated to the computed path(s). Thus it can fit very well in the horizontal query stage of CS0. A PCE MAY have further capability to do multi-layer and/or inter-domain path computation which can be further utilized. NCG which understands the vertical query and the presence of applications constraints can break the application request into suitable path computation request which PCE understands. In this scenario, the PCE MAY have no knowledge of applications and provide only network related metrics to the NCG: the NCG (or the ACG for an application-centric model) is in charge of correlating the network quotations with the application layer information to achieve the global CS0 objective.



With this architecture, NCG can request PCE different sets of computation mode that are not currently supported by PCE. For instance, NCG may request PCE a multi-destination and multi-source path computation request. This scenario arises when there are many possible Data Center choices for a given application request and there could be multiple sources for this request. Multi-destination with a single source (aka., anycast) is a default case for multi-destination and multi-source path computation.

In addition, with this architecture, NCG may have different sets of objectives and constraints than typical path computation requests. For instance, multi-criteria objective functions that combine the bandwidth requirement and latency may be very useful for some applications. [[PCE-SERVICE-AWARE](#)] describes the extension to PCEP to carry Latency, Latency-Variation and Loss as constraints for end to end path computation.

In a Stateful PCE (refer [[PCE-STATEFUL](#)]), there is a strict synchronization of the network states (in term of topology and resource information), and the set of computed paths and reserved resources in use in the network. In other words, the PCE utilizes information from the TED as well as information about existing paths (for example, TE LSPs) in the network when processing new requests. Stateful PCE will be very important tool to achieve the goals of cross stratum optimization as maintains the status of final path selected after cross (application and network) optimization.

As Stateful PCE would keep both LSP ID and the application ID associated with the LSP, it will make path computation more efficient in terms of resource usage and computation time. Moreover, Stateful PCE would have an accurate snapshot of network resource information and as such it can increase adaptability to the changes. This may be important for some application that requires a stringent performance objective.

In conclusion -

- o NCG can use the PCE to do path computation based on constraints from multiple sources and destinations.
- o Stateful PCE can help in maintaining the status of the final cross optimized path. It can also help NCG in maintaining the relationship of application request and setup path. In case of any change of the path, the Stateful PCE and NCG and cooperate and take suitable action.

#### **[4.](#) Path Computation and Setup Procedure**





Path computation flow is shown in Figure 5.

1. User for application would contact the application gateway ACG with its requirements.
2. ACG would further query the NCG to obtain the underlying network Status and quotations (offers) for the network connectivity services.
3. NCG would break the vertical request into suitable horizontal path computation request(s).
4. PCE would provide the result to NCG.
5. NCG would abstract the computation result and provide to ACG.
6. NCG and ACG would cooperate to finalize the path that needs to be setup.
7. Note that that the final decision can be made either in ACG or NCG depending on the mode of operation. With application centric mode, minimal data center/IT resource information would flow from ACG to NCG while ACG collects network abstracted information from NCG to choose the optimal application-network resources. With network centric mode, ACG would supply maximal data center/IT resource information to NCG so that NCG in conjunction with PCE would determine the optimal mixed set of application and network resources. In the latter case, the PCE COULD support application /IT- based constrained computation capability beyond network path computation. This requires further PCE capabilities to receive and process data center/IT resource information, possibly in conjunction with network information.

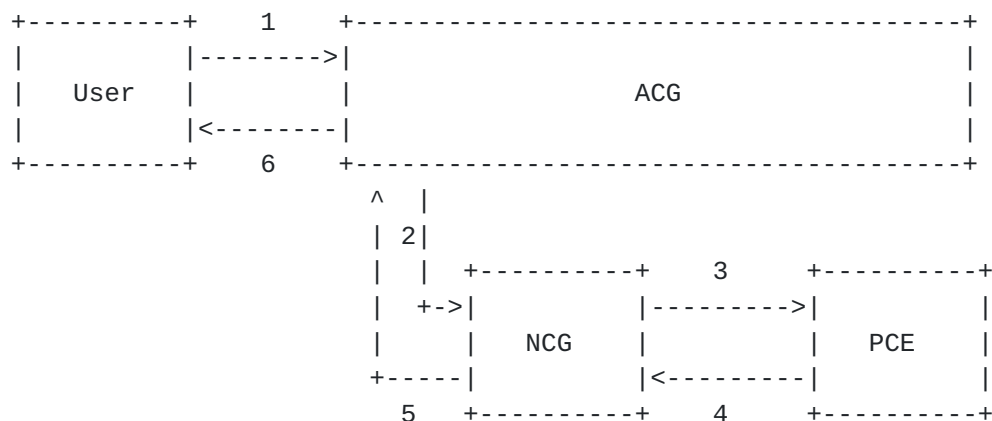


Figure 5: Path Computation Flow



In this section we would analyze the mechanisms to finally setup the cross stratum optimized path.

#### 4.1. Path Setup Using NMS

After ACG and NCG have decided the path that needs to be set, NCG can send a request to NMS asking it relay the message to the head end LSR (also a PCC) to setup the pre computed path. Once the path signaling is completed and the LSP is setup, PCC should relay the status of the LSP to the Stateful PCE.

In this mechanism we can reuse the existing NMS to establish the path. Any updates or deletion of such path would be made via the NMS.

Head end LSR (PCC) 'H' is always the owner of the path.

See Figure 6 for this scenario.

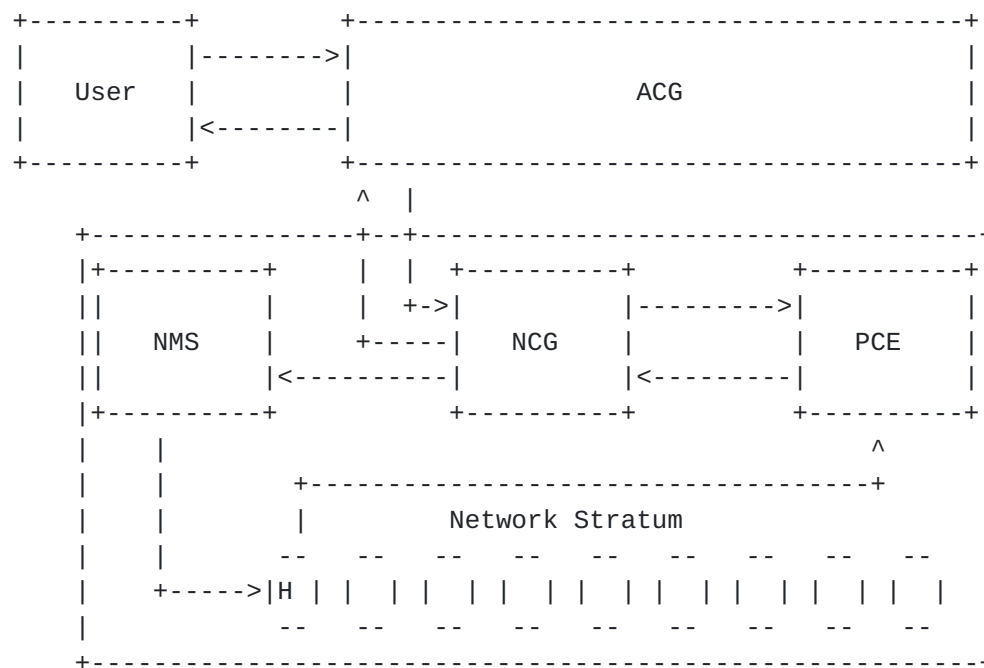


Figure 6: Path Setup Using NMS

#### 4.2. Path Setup Using a Network Control Plane

A network control plane (e.g. GMPLS) MAY be used to automatically establish the cross optimized path between the selected end points. This control plane MAY be triggered via -



- o NCG to Control Plane: GMPLS UNI or other protocols
- o Control Plane to Head end Router: GMPLS Control Channel Interface (CCI). Suitable protocol extensions are needed to achieve this.

See Figure 7 for this scenario.

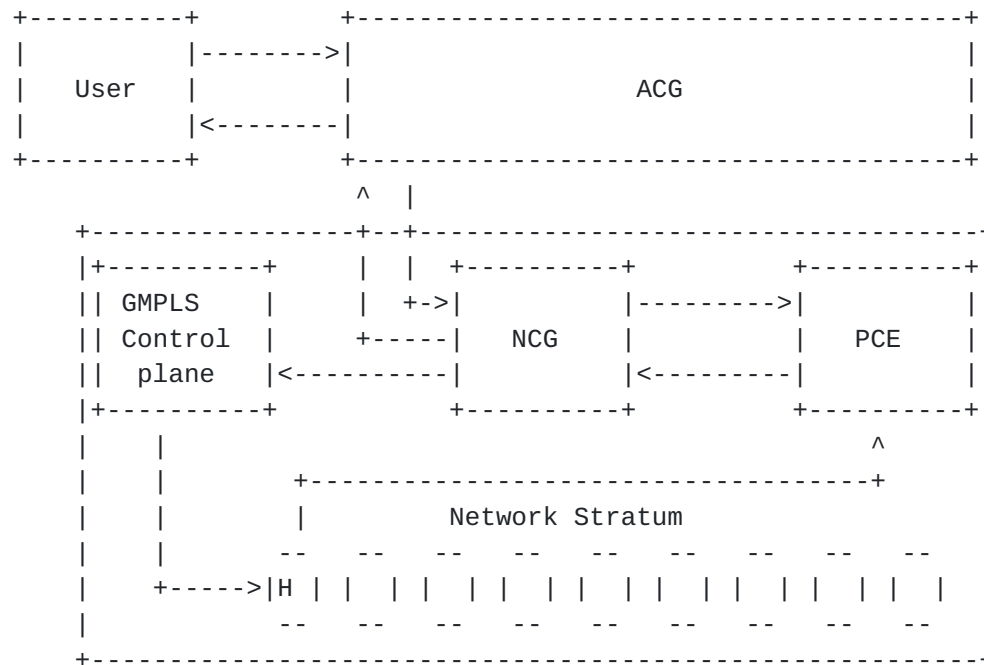


Figure 7: Path Setup Using Centralized Control Plane

After cross optimization, ACG and NCG will select the suitable end points, (the path is already calculated by PCE), this path is conveyed to the head end LSR which signals the path and notify the status to the Stateful PCE. Later NCG can send suitable message to tear down the path.

Using centralized control plane can make the NCG responsible for the LSP. Head end LSR signals and maintains the status but the establishment and tear-down are initiated by the control plane. This would have an obvious advantage in managing the setup paths. The Stateful PCE will maintain the TED as well as the status of setup LSP. NCG through centralized control plane can further setup/teardown/modify/re-optimize those paths.

#### 4.3. Path Setup using PCE

A Stateful PCE extension MAY be developed to communicate the cross optimized path to the head end LSR. Current PCEP protocol requires



PCC to trigger Path request and PCE to provide reply. Even in Stateful PCE, PCC must delegate the LSP to a PCE, a PCE never initiate path setup. An extension to PCEP protocol MAY let PCE notify to PCC (Head end LSR) to establish the path.

NCG via PCE and PCEP protocol can establish and tear-down LSP as shown in Figure 8. [[PCE INITIATED](#)] is one such attempt to extend PCEP.

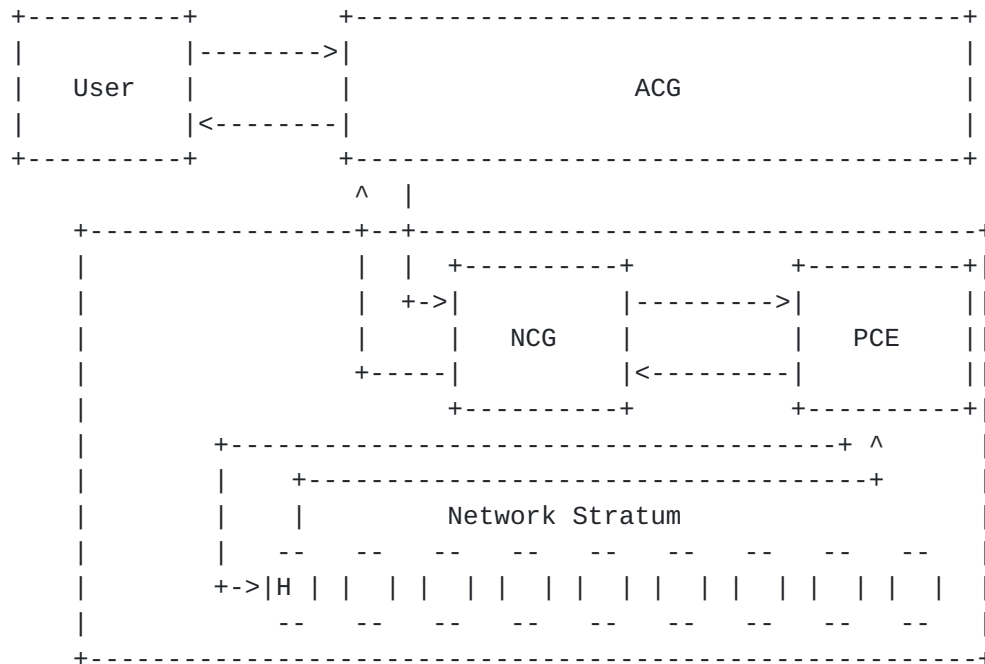


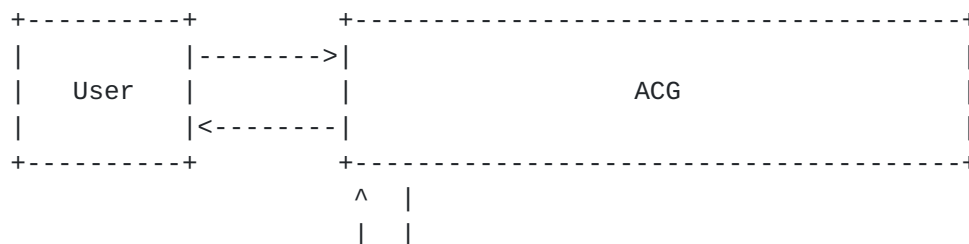
Figure 8: Path Setup using PCE

## 5. Other Consideration

### 5.1. Inter-domain

#### 5.1.1. One Application Domain with Multiple Network Domains

Underlying network connecting the datacenters MAYBE made up of multiple domains (AS and Area). In this case an inter-domain path computation is required.







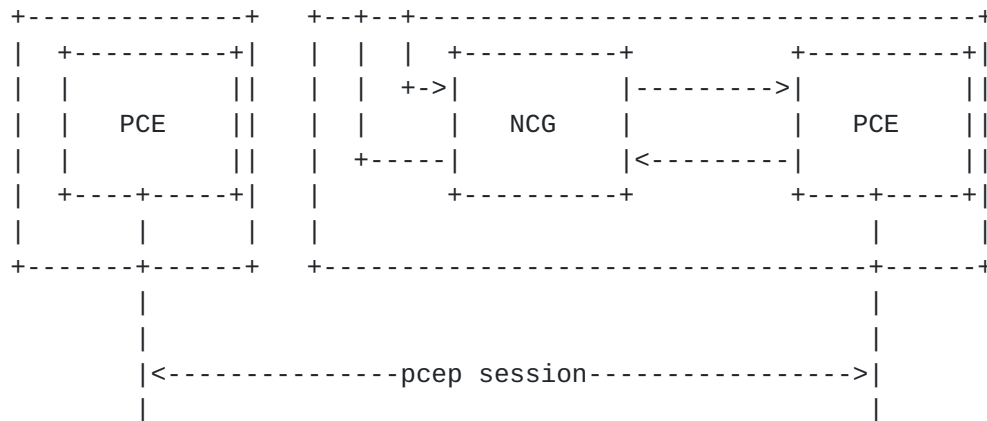


Figure 9: Multi-domain Scenario

[RFC5441] describes an inter-domain path computation with cooperating PCEs which can be enhanced and utilized in CS0 enabled path computation.

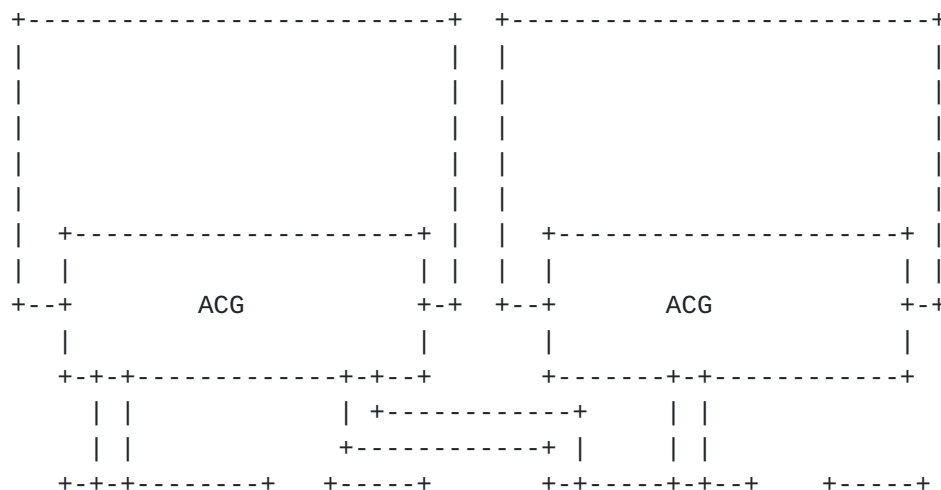
#### 5.1.2. Multiple Application Domains with Multiple Network Domains

Underlying network connecting the datacenters MAY be made up of multiple domains (AS and Area) as well as applications domains and ACG MAY be distributed. In such case multiple ACG and NCG will be involved in cross optimizing. This needs to be analyzed further.

##### 5.1.2.1. ACG talks to multiple NCGs

As shown in Figure 10, ACG where the request originates may communicate with multiple NCG to get the network information from multiple domains to be cross optimized.

Application stratum





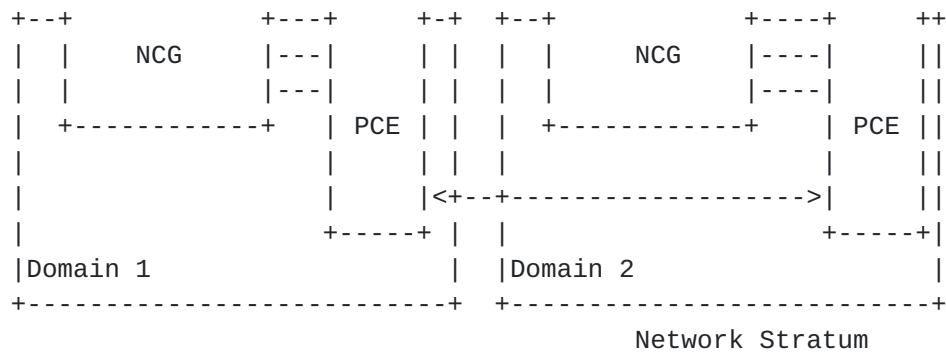


Figure 10: ACG talks to multiple NCG

#### 5.1.2.2. ACG talks to the primary NCG, which talks to the other NCG of different domains

As shown in Figure 11, ACG communicated only to the primary NCG, which may gather network information from multiple NCG and then communicate consolidated information to ACG.

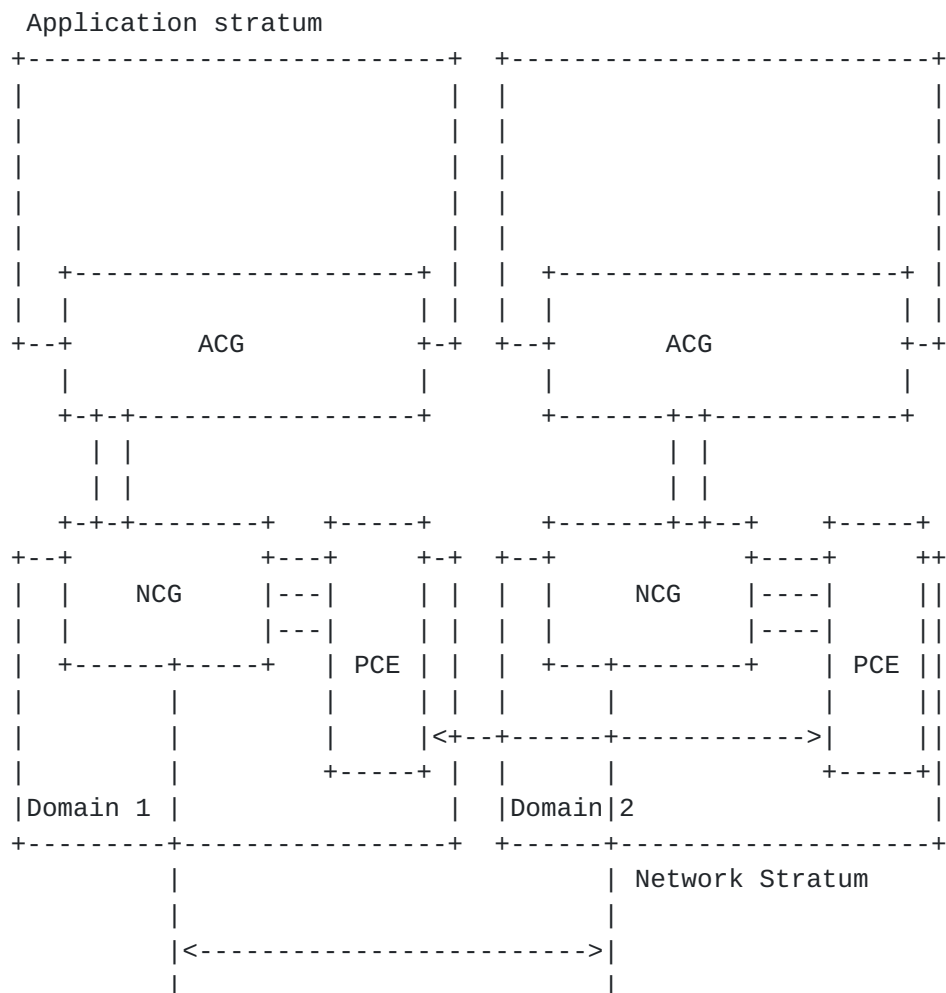




Figure 11: Primary NCG talks to other NCG

## **5.2. Bottleneck**

In optical networks all PCE messages are sent over control channel, in Stateful PCE cases its observed that in case of a major link or node failure lot of PCEP messages are sent from all PCC to PCE. This use lot of bandwidth of the control channel.

PCE MAY become a common point of failure and bottleneck. PCE/NCG/ACG failure as well as the link-failure disrupting connectivity could be highly disruptive to the system.

The solution should focus on reducing such bottleneck.

## **6. IANA Considerations**

TBD

## **7. Security Considerations**

TBD

## **8. Manageability Considerations**

TBD

## **9. Acknowledgements**

The research work of N. Ciulli and L.M. Contreras has been partially supported by the GEYSERS project ([www.geysers.eu](http://www.geysers.eu)), funded by the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement n. 248657.

## **10. References**

### **10.1. Normative References**

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

### **10.2. Informative References**

[RFC4655] Farrel, A., Vasseur, J., and J. Ash, "A Path Computation Element (PCE)-Based Architecture", [RFC 4655](#), August 2006.



- [RFC5440] Vasseur, JP. and JL. Le Roux, "Path Computation Element (PCE) Communication Protocol (PCEP)", [RFC 5440](#), March 2009.
- [RFC5441] Vasseur, JP., Zhang, R., Bitar, N., and JL. Le Roux, "A Backward-Recursive PCE-Based Computation (BRPC) Procedure to Compute Shortest Constrained Inter-Domain Traffic Engineering Label Switched Paths", [RFC 5441](#), April 2009.
- [CS0-DATACNTR]  
Lee, Y., Bernstein, G., So, N., Kim, T., Shiimoto, K., and O. Gonzalez-de-Dios, "Research Proposal for Cross Stratum Optimization (CS0) between Data Centers and Networks. ([draft-lee-cross-stratum-optimization-datacenter-00](#))", March 2011.
- [CS0-PROBLEM]  
Lee, Y., Bernstein, G., So, N., Hares, S., Xia, F., Shiimoto, K., and O. Gonzalez-de-Dios, "Problem Statement for Cross-Layer Optimization. ([draft-lee-cross-layer-optimization-problem-02](#))", January 2011.
- [NS-QUERY]  
Lee, Y., Bernstein, G., So, N., McDysan, D., Kim, T., Shiimoto, K., and O. Gonzalez-de-Dios, "Problem Statement for Network Stratum Query. ([draft-lee-network-stratum-query-problem-02](#)) ", April 2011.
- [CS0-PCE-REQT]  
Tovar, A., Contreras, L., Landi, G., and N. Ciulli, "Path Computation Requirements for Cross-Stratum-Optimization. ([draft-tovar-cso-path-computation-requirements-00](#))", October 2011.
- [PCE-SERVICE-AWARE]  
Dhody, D., Manral, V., Ali, Z., Swallow, G., and K. Kumaki, "Extensions to the Path Computation Element Communication Protocol (PCEP) to compute service aware Label Switched Path (LSP). ([draft-ietf-pce-pcep-service-aware-01](#))", July 2013.
- [PCE-STATEFUL]  
Crabbe, E., Medved, J., Varga, R., and I. Minei, "PCEP Extensions for Stateful PCE. ([draft-ietf-pce-stateful-pce-06](#))", August 2013.
- [ALTO-APPNET]





Lee, Y., Bernstein, G., Varga, T., Madhavan, S., and D. Dhody, "ALTO Extensions to Support Application and Network Resource Information Exchange for High Bandwidth Applications. ([draft-lee-alto-app-net-info-exchange-02](#))", July 2013.

[PCE\_INITIATED]

Crabbe, E., Minei, I., Sivabalan, S., and R. Varga, "PCEP Extensions for PCE-initiated LSP Setup in a Stateful PCE Model. ([draft-crabbe-pce-pce-initiated-lsp-02](#))", July 2013.

Authors' Addresses

Dhruv Dhody  
Huawei Technologies  
Leela Palace  
Bangalore, Karnataka 560008  
INDIA

EMail: [dhruv.dhody@huawei.com](mailto:dhruv.dhody@huawei.com)

Young Lee  
Huawei Technologies  
1700 Alma Drive, Suite 500  
Plano, TX 75075  
USA

EMail: [leeyoung@huawei.com](mailto:leeyoung@huawei.com)

Nicola Ciulli  
Nextworks

EMail: [n.ciulli@nextworks.it](mailto:n.ciulli@nextworks.it)

Luis M. Contreras  
Telefonica I+D

EMail: [lmcm@tid.es](mailto:lmcm@tid.es)



Oscar Gonzalez de Dios  
Telefonica I+D  
Don Ramon de la Cruz  
Madrid 28006  
Spain

EMail: ogondio@tid.es