

PCE Working Group
Internet-Draft
Intended status: Standards Track
Expires: 22 September 2022

A. Wang
China Telecom
B. Khasanov
Yandex LLC
S. Fang
R. Tan
Huawei Technologies, Co., Ltd
C. Zhu
ZTE Corporation
21 March 2022

PCEP Extension for Native IP Network
draft-ietf-pce-pcep-extension-native-ip-18

Abstract

This document defines the Path Computation Element Communication Protocol (PCEP) extension for Central Control Dynamic Routing (CCDR) based application in Native IP network. The scenario and framework of CCDR in native IP is described in [RFC8735] and [RFC8821]. This draft describes the key information that is transferred between Path Computation Element (PCE) and Path Computation Clients (PCC) to accomplish the End to End (E2E) traffic assurance in Native IP network under central control mode.

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 22 September 2022.

Copyright Notice

Copyright (c) 2022 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 Revised BSD License text as described in Section 4.e of the [Trust Legal Provisions](#) and are provided without warranty as described in the Revised BSD License.

Table of Contents

1.	Introduction	3
2.	Conventions used in this document	3
3.	Terminology	3
4.	Capability Advertisement	4
4.1.	Open message	4
5.	PCEP messages	4
5.1.	The PCInitiate message	5
5.2.	The PCRpt message	6
6.	PCECC Native IP TE Procedures	7
6.1.	BGP Session Establishment Procedures	7
6.2.	Explicit Route Establish Procedures	9
6.3.	BGP Prefix Advertisement Procedures	12
7.	New PCEP Objects	14
7.1.	CCI Object	14
7.2.	BGP Peer Info Object	15
7.3.	Explicit Peer Route Object	17
7.4.	Peer Prefix Advertisement Object	20
8.	End to End Path Protection	21
9.	Re-Delegation and Clean up	21
10.	BGP Considerations	22
11.	New Error-Types and Error-Values Defined	22
12.	Deployment Considerations	23
13.	Implementation Status	24
13.1.	Proof of Concept based on ODL	24
14.	Security Considerations	25
15.	IANA Considerations	25
15.1.	Path Setup Type Registry	25
15.2.	PCECC-CAPABILITY sub-TLV's Flag field	25
15.3.	PCEP Object Types	25
15.4.	PCEP-Error Object	26
16.	Contributor	27
17.	Acknowledgement	27

18. Normative References	27
Authors' Addresses	29

[1. Introduction](#)

Generally, Multiprotocol Label Switching Traffic Engineering (MPLS-TE) requires the corresponding network devices support Multiprotocol Label Switching (MPLS) or Resource ReSerVation Protocol (RSVP)/Label Distribution Protocol (LDP) technologies to assure the End-to-End (E2E) traffic performance. In Segment Routing either IGP extensions or BGP are used to steer a packet through an SR Policy instantiated as an ordered list of instructions called "segments". But in native IP network, there will be no such signaling protocol to synchronize the action among different network devices. It is necessary to use the central control mode that described in [\[RFC8283\]](#) to correlate the forwarding behavior among different network devices. [\[RFC8821\]](#) describes the architecture and solution philosophy for the E2E traffic assurance in Native IP network via Multi Border Gateway Protocol (BGP) solution. This draft describes the corresponding Path Computation Element Communication Protocol (PCEP) extensions to transfer the key information about BGP peer info, peer prefix advertisement and the explicit peer route on on-path routers.

[2. Conventions used in this document](#)

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "NOT RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [BCP 14](#) [\[RFC2119\]](#) [\[RFC8174\]](#) when, and only when, they appear in all capitals, as shown here.

[3. Terminology](#)

This document uses the following terms defined in [\[RFC5440\]](#): PCE, PCEP

The following terms are defined in this document:

- * CCDR: Central Control Dynamic Routing

- * E2E: End to End
- * BPI: BGP Peer Info
- * EPR: Explicit Peer Route
- * PPA: Peer Prefix Advertisement
- * QoS: Quality of Service

[4.](#) Capability Advertisemnt

[4.1.](#) Open message

During the PCEP Initialization Phase, PCEP Speakers (PCE or PCC) advertise their support of Native IP extensions.

This document defines a new Path Setup Type (PST) [[RFC8408](#)] for Native-IP, as follows:

- * PST = TBD1: Path is a Native IP path as per [[RFC8821](#)].

A PCEP speaker MUST indicate its support of the function described in this document by sending a PATH-SETUP-TYPE-CAPABILITY TLV in the OPEN object with this new PST included in the PST list.

[RFC9050] defined the PCECC-CAPABILITY sub-TLV to exchange information about their PCECC capability. A new flag is defined in PCECC-CAPABILITY sub-TLV for Native IP:

N (NATIVE-IP-TE-CAPABILITY - 1 bit - TBD2): If set to 1 by a PCEP speaker, it indicates that the PCEP speaker is capable for TE in Native IP network as specified in this document. The flag MUST be set by both the PCC and PCE in order to support this extension.

If a PCEP speaker receives the PATH-SETUP-TYPE-CAPABILITY TLV with the newly defined path setup type, but without the N bit set in PCECC-CAPABILITY sub-TLV, it MUST:

- * Send a PCErr message with Error-Type=10(Reception of an invalid object) and Error-Value TBD3(PCECC NATIVE-IP-TE-CAPABILITY bit is not set).
- * Terminate the PCEP session

5. PCEP messages

PCECC Native IP TE solution utilizing the existing PCE LSP Initiate Request message(PCInitiate) [[RFC8281](#)], and PCE Report message(PCRpt) [[RFC8281](#)] to accomplish the multi BGP sessions establishment, E2E TE path deployment, and route prefixes advertisement among different BGP sessions. A new PST for Native-IP is used to indicate the path setup based on TE in Native IP networks.

The extended PCInitiate message described in [[RFC9050](#)] is used to download or cleanup central controller's instructions (CCIs). [[RFC9050](#)] specifies an object called CCI for the encoding of central controller's instructions. This document specify a new CCI object-

type for Native IP. The PCEP messages are extended in this document to handle the PCECC operations for Native IP. Three new PCEP Objects (BGP Peer Info (BPI) Object, Explicit Peer Route (EPR) Object and Peer Prefix Advertisement (PPA) Object) are defined in this document. Refer to [Section 7](#) for detail object definitions.

5.1. The PCInitiate message

The PCInitiate Message defined in [[RFC8281](#)] and extended in [[RFC9050](#)] is further extended to support Native-IP CCI.

The format of the extended PCInitiate message is as follows:

```
<PCInitiate Message> ::= <Common Header>
                           <PCE-initiated-lsp-list>
```

Where:

```
<Common Header> is defined in [RFC5440]
```

```
<PCE-initiated-lsp-list> ::= <PCE-initiated-lsp-request>
                              [<PCE-initiated-lsp-list>]
```

```
<PCE-initiated-lsp-request> ::=
```

```
(<PCE-initiated-lsp-instantiation>|  
  <PCE-initiated-lsp-deletion>|  
  <PCE-initiated-lsp-central-control>)
```

```
<PCE-initiated-lsp-central-control> ::= <SRP>  
                                         <LSP>  
                                         (<cci-list>|  
                                          ((<BPI>|<EPR>|<PPA>)  
                                           <CCI>))
```

```
<cci-list> ::= <CCI>  
              [<cci-list>]
```

Where:

<cci-list> is as per
[I-D.ietf-pce-pcep-extension-for-pce-controller].
<PCE-initiated-lsp-instantiation> and
<PCE-initiated-lsp-deletion> are as per
[\[RFC8281\]](#).

The LSP and SRP objects are defined in [\[RFC8231\]](#).

When PCInitiate message is used create Native IP instructions, the SRP, LSP and CCI objects MUST be present. The error handling for missing SRP, LSP or CCI object is as per [\[RFC9050\]](#). Further only one of BPI, EPR, or PPA object MUST be present. The PLSP-ID within the

LSP object should be set by PCC uniquely according to the Symbolic Path Name TLV that included in the CCI object. The Symbolic Path Name is used by the PCE/PCC to identify uniquely the E2E native IP TE path.

If none of them are present, the receiving PCC MUST send a PCErr message with Error-type=6 (Mandatory Object missing) and Error-value=TBD4 (Native IP object missing). If there are more than one of BPI, EPR or PPA object are presented, the receiving PCC MUST send a PCErr message with Error-type=19(Invalid Operation) and Error-value=TBD5(Only one of the BPI, EPR or PPA object can be included in this message).

To cleanup the SRP object must set the R (remove) bit.

5.2. The PCRpt message

The PCRpt message is used to acknowledge the Native-IP instructions received from the central controller (PCE).

The format of the PCRpt message is as follows:

```
<PCRpt Message> ::= <Common Header>
                        <state-report-list>
```

Where:

```
<state-report-list> ::= <state-report>[<state-report-list>]
```

```
<state-report> ::= (<lsp-state-report>|
                    <central-control-report>)
```

```
<lsp-state-report> ::= [<SRP>]
                        <LSP>
                        <path>
```

```
<central-control-report> ::= [<SRP>]
                             <LSP>
                             (<cci-list>|
                             ((<BPI>|<EPR>|<PPA>)<CCI>))
```

Where:

<path> is as per [\[RFC8231\]](#) and the LSP and SRP object are also defined in [\[RFC8231\]](#).

The error handling for missing CCI object is as per [\[RFC9050\]](#). Further only one of BPI, EPR, or PPA object MUST be present.

If none of them are present, the receiving PCE MUST send a PCErr message with Error-type=6 (Mandatory Object missing) and Error-value=TBD4 (Native IP object missing). If there are more than one of BPI, EPR or PPA object are presented, the receiving PCE MUST send a PCErr message with Error-type=19(Invalid Operation) and Error-value=TBD5(Only one of the BPI, EPR or PPA object can be included in this message).

6. PCECC Native IP TE Procedures

The detail procedures for the TE in native IP environment are described in the following sections.

6.1. BGP Session Establishment Procedures

The PCInitiate message can be used to configure the parameters for a BGP peer session using the PCInitiate and PCRpt message pair. This pair of PCE messages is exchanged with a PCE function attached to each BGP peer which needs to be configured. After the BGP peer session has been configured via this pair of PCE messages the BGP session establishment process operates in a normal fashion. All BGP peers are configured for peer to peer communication whether the peers are E-BGP peers or I-BGP peers. One of the IBGP topologies requires that multiple I-BGPs peers operate in a route-reflector I-BGP peer topology. The example below shows two I-BGP route reflector clients interacting with one Route Reflector (RR), but Route Reflector topologies may have up to 100s of clients. Centralized configuration via PCE provides mechanisms to scale auto-configuration of small and large topologies.

The PCInitiate message should be sent to PCC which acts as BGP router and/or route reflector(RR).

The route reflector topology for a single AS is shown in Figure 1. The BGP routers R1, R3, and R7 are within a single AS. R1 and R7 are BGP router-reflector clients, and R3 is a Route Reflector. The PCInitiate message should be sent all of the BGP routers that need to be configured R1 (M3), R3 (M2 & M3), and R7 (M4).

PCInitiate message creates an auto-configuration function for these BGP peers providing the indicated Peer AS and the Local/Peer IP Address.

When PCC receives the BPI and CCI object (with the R bit set to 0 in SRP object) in PCInitiate message, the PCC should try to establish the BGP session with the indicated Peer AS and Local/Peer IP address.

When PCC creates successfully the BGP session that is indicated by

the associated information, it should report the result via the PCRpt messages, with BPI object and the corresponding SRP and CCI object included.

When PCC receives this message with the R bit set to 1 in SRP object in PCInitiate message, the PCC should clear the BGP session that indicated by the BPI object.

When PCC clears successfully the specified BGP session, it should report the result via the PCRpt message, with the BPI object included, and the corresponding SRP and CCI object.

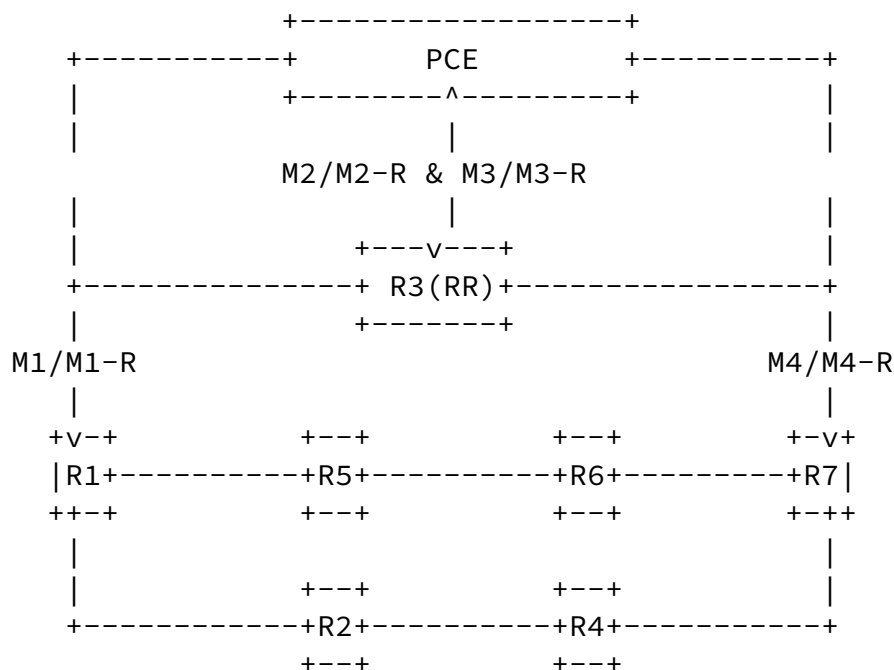


Figure 1: BGP Session Establishment Procedures(R3 act as RR)

The message number, message peers, message type and message key parameters in the above figures are shown in below table:

Table 1: Message Information

No.	Peers	Type	Message Key Parameters
M1	PCE/R1	PCInitiate	CC-ID=X1(Symbolic Path Name=Class A)
M1-R		PCRpt	BPI Object(Local_IP=R1_A,Peer_IP=R3_A)
M2	PCE/R3	PCInitiate	CC-ID=X2(Symbolic Path Name=Class A)
M2-R		PCRpt	BPI Object(Local_IP=R3_A,Peer_IP=R1_A)
M3	PCE/R3	PCInitiate	CC-ID=X3(Symbolic Path Name=Class A)
M3-R		PCRpt	BPI Object(Local_IP=R3_A,Peer_IP=R7_A)
M4	PCE/R7	PCInitiate	CC-ID=X4(Symbolic Path Name=Class A)
M4-R		PCRpt	BPI Object(Local_IP=R7_A,Peer_IP=R3_A)

If the PCC cannot establish the BGP session that required by this object, it should report the error values via PCErr message with the newly defined error type(Error-type=TBD6) and error value(Error-value=TBD7, Peer AS not match; or Error-Value=TBD8, Peer IP can't be reached), which is indicated in [Section 11](#)

If the Local IP Address or Peer IP Address within BPI object is used in other existing BGP sessions, the PCC should report such error situation via PCErr message with Err-type=TBD6 and error value(Error-value=TBD9, Local IP is in use; Error-value=TBD10, Remote IP is in use).

6.2. Explicit Route Establish Procedures

The explicit route establishment procedures can be used to install a route via PCE in the PCC/BGP Peer, using PCInitiate and PCRpt message pair. Although the BGP policy might redistribute the routes installed by explicit route, the PCE-BGP implementation needs to prohibit the redistribution of the explicit route. PCE explicit routes operate similar to static routes installed by network management protocols (netconf/restconf) but the routes are associated with the PCE routing module. Explicit route installations (like NM static routes) must carefully install and uninstall static routes in an specific order so that the pathways are established without loops.

The PCInitiate message should be sent to the on-path routers respectively. In the example, for explicit route from R1 to R7, the PCInitiate message should be sent to R1(M1), R2(M2) and R4(M3), as shown in Figure 2. For explicit route from R7 to R1, the PCInitiate

message should be sent to R7(M1), R4(M2) and R2(M3), as shown in Figure 3.

When PCC receives the EPR and the CCI object (with the R bit set to 0 in SRP object) in PCInitiate message, the PCC should install the explicit route to the the peer.

When PCC install successfully the explicit route to the peer, it should report the result via the PCRpt messages, with EPR object and the corresponding SRP and CCI object included.

When PCC receives the EPR and the CCI object with the R bit set to 1 in SRP object in PCInitiate message, the PCC should clear the explicit route to the peer that indicated by the EPR object.

When PCC clear successfully the explicit route that indicated by this object, it should report the result via the PCRpt message, with the EPR object included, and the corresponding SRP and CCI object.

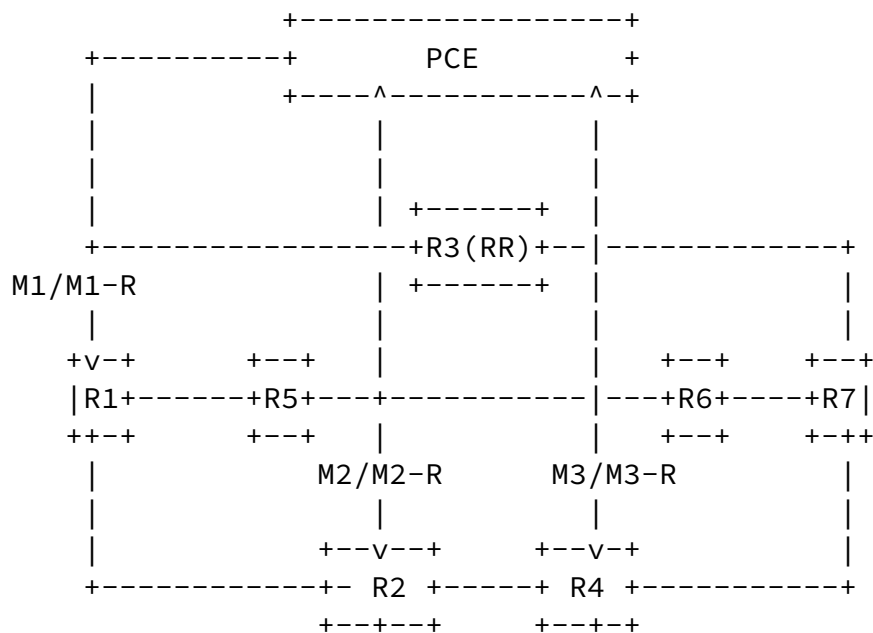


Figure 2: Explicit Route Establish Procedures(From R1 to R7)

The message number, message peers, message type and message key parameters in the above figures are shown in below table:

Table 2: Message Information

No.	Peers	Type	Message Key Parameters
M1	PCE/R1	PCInitiate	CC-ID=X1(Symbolic Path Name=Class A)
M1-R		PCRpt	EPR Object(Peer Address=R7_A,Next Hop=R2_A)
M2	PCE/R2	PCInitiate	CC-ID=X2(Symbolic Path Name=Class A)
M2-R		PCRpt	EPR Object(Peer Address=R7_A,Next Hop=R4_A)
M3	PCE/R4	PCInitiate	CC-ID=X3(Symbolic Path Name=Class A)
M3-R		PCRpt	EPR Object(Peer Address=R7_A,Next Hop=R7_A)

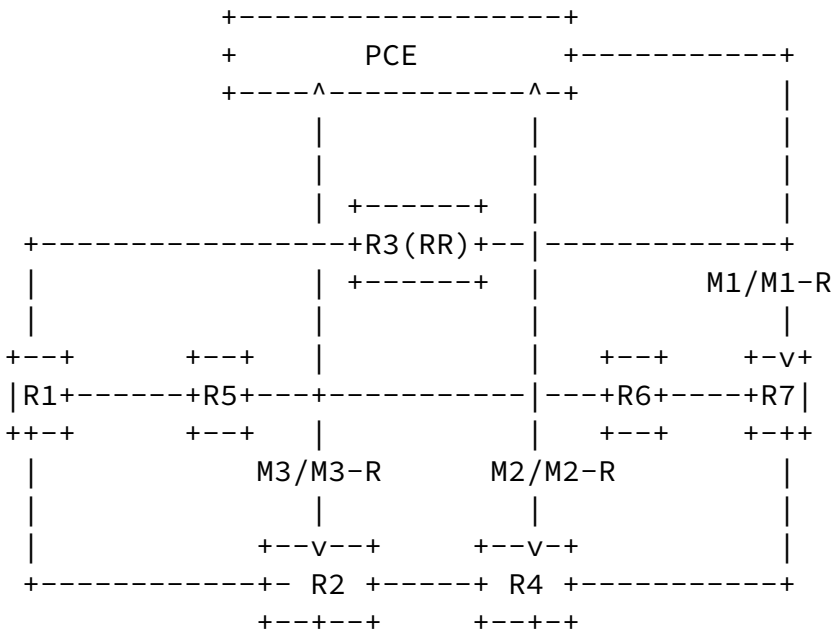


Figure 3: Explicit Route Establish Procedures(From R7 to R1)

The message number, message peers, message type and message key parameters in the above figures are shown in below table:

Wang, et al.

Expires 22 September 2022

[Page 11]

Internet-Draft

PCEP Extension for Native IP Network

March 2022

Table 3: Message Information

No.	Peers	Type	Message Key Parameters
M1	PCE/R7	PCInitiate	CC-ID=X1(Symbolic Path Name=Class A)
M1-R		PCRpt	EPR Object(Peer Address=R1_A,Next Hop=R4_A)
M2	PCE/R4	PCInitiate	CC-ID=X2(Symbolic Path Name=Class A)
M2-R		PCRpt	EPR Object(Peer Address=R1_A,Next Hop=R2_A)
M3	PCE/R2	PCInitiate	CC-ID=X3(Symbolic Path Name=Class A)
M3-R		PCRpt	EPR Object(Peer Address=R1_A,Next Hop=R1_A)

In order to avoid the transient loop during the deploy of explicit peer route, the EPR object should be sent to the PCCs in the reverse order of the E2E path. To remove the explicit peer route, the EPR object should be sent to the PCCs in the same order of E2E path.

To accomplish ECMP effects, the PCE can send multiple EPR objects to the same node, with the same route priority and peer address value but different next hop addresses.

The PCC should verify that the next hop address is reachable. Upon the error occurs, the PCC SHOULD send the corresponding error via

PCErr message, with an error information (Error-type=TBD6, Error-value=TBD12, Explicit Peer Route Error) that defined in [Section 11](#).

When the peer info is not the same as the peer info that indicated in BPI object in PCC for the same path that is identified by Symbolic Path Name TLV, an error (Error-type=TBD6, Error-value=17, EPR/BPI Peer Info mismatch) should be reported via the PCErr message.

[6.3](#). BGP Prefix Advertisement Procedures

The detail procedures for BGP prefix advertisement are shown below, using PCInitiate and PCRpt message pair.

The PCInitiate message should be sent to PCC that acts as BGP peer router only. In the example, it should be sent to R1(M1) or R7(M2) respectively.

When PCC receives the PPA and the CCI object (with the R bit set to 0 in SRP object) in PCInitiate message, the PCC should send the prefixes indicated in this object to the appointed BGP peer.

When PCC sends successfully the prefixes to the appointed BGP peer, it should report the result via the PCRpt messages, with PPA object and the corresponding SRP and CCI object included.

When PCC receives the PPA and the CCI object with the R bit set to 1 in SRP object in PCInitiate message, the PCC should withdraw the prefixes advertisement to the peer that indicated by this object.

When PCC withdraws successfully the prefixes that indicated by this object, it should report the result via the PCRpt message, with the PPA object included, and the corresponding SRP and CCI object.

The allowed AFI/SAFI for the IPv4 BGP session should be 1/1(IPv4 prefix) and the allowed AFI/SAFI for the IPv6 BGP session should be 2/1(IPv6 prefix). If mismatch occur, an error(Error-type=TBD6, Error-value=TBD18, BPI/PPR address family mismatch) should be reported via PCErr message.

When the peer info is not the same as the peer info that indicated in BPI object in PCC for the same path that is identified by Symbolic Path Name TLV, an error (Error-type=TBD6, Error-value=TBD19, PPA/BPI peer info mismatch) should be reported via the PCErr message.

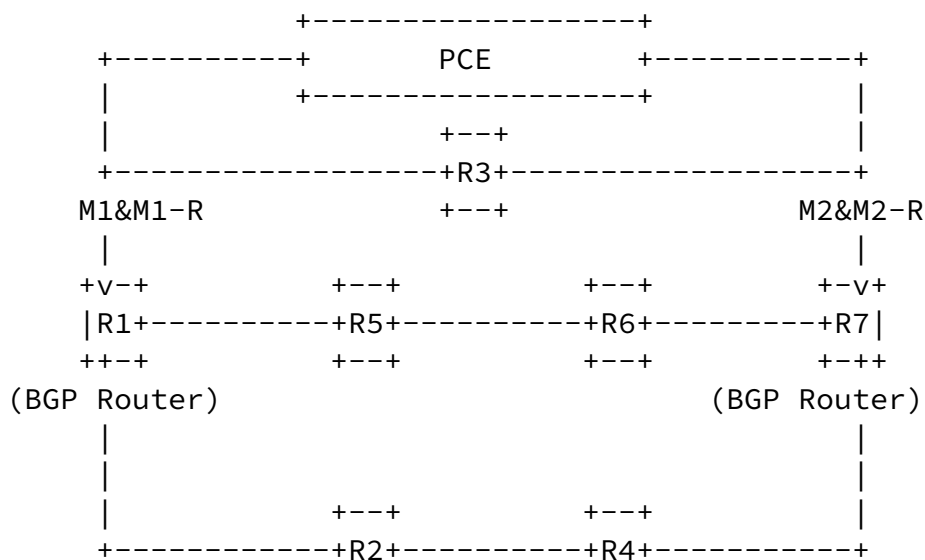


Figure 4: BGP Prefix Advertisement Procedures

Table 4: Message Information

No.	Peers	Type	Message Key Parameters
M1	PCE/R1	PCInitiate	CC-ID=X1(Symbolic Path Name=Class A)
M1-R		PCRpt	PPA Object(Peer IP=R7_A,Prefix=1_A)
M2	PCE/R7	PCInitiate	CC-ID=X2(Symbolic Path Name=Class A)
M2-R		PCRpt	PPA Object(Peer IP=R1_A,Prefix=7_A)

7. New PCEP Objects

One new CCI Object and three new PCEP objects are defined in this draft. All new PCEP objects are as per [\[RFC5440\]](#)

7.1. CCI Object

The Central Control Instructions (CCI) Object is used by the PCE to specify the forwarding instructions is defined in [\[RFC9050\]](#). This document defines another object-type for Native-IP.

CCI Object-Type is TBD13 for Native-IP as below

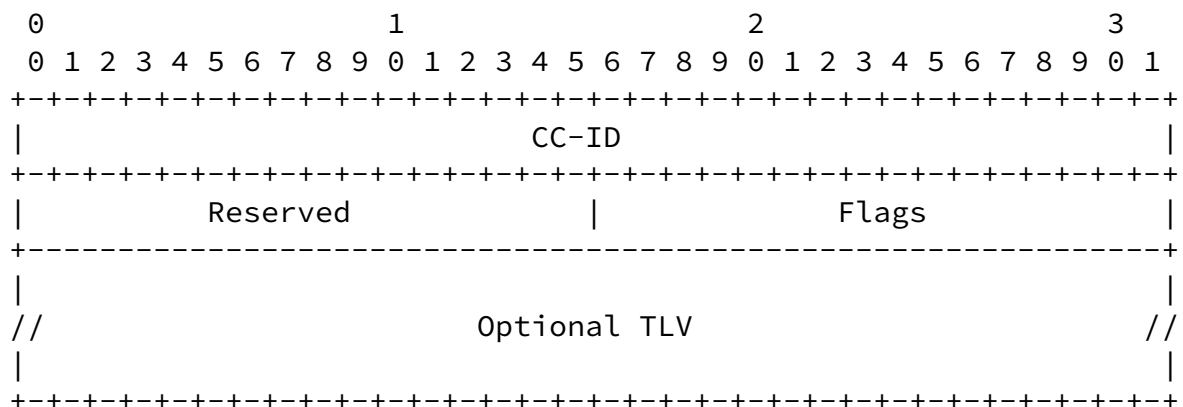


Figure 5: CCI Object for Native IP

Figure 1

The field CC-ID is as described in [\[RFC9050\]](#). Following fields are defined for CCI Object-Type TBD13

Reserved: is set to zero while sending, ignored on receipt.

Flags: is used to carry any additional information pertaining to the CCI. Currently no flag bits are defined.

The Symbolic Path Name TLV [\[RFC8231\]](#) MUST be included in the CCI Object-Type TBD13 to identify the E2E TE path in Native IP environment and MUST be unique.

7.2. BGP Peer Info Object

The BGP Peer Info object is used to specify the information about the peer that the PCC should establish the BGP relationship with. This object should only be included and sent to the head and end router of the E2E path in case there is no Route Reflection (RR) involved. If the RR is used between the head and end routers, then such information should be sent to head router, RR and end router respectively.

By default, there MUST be no prefix be distributed via such BGP session that established by this object.

By default, the Local/Peer IP address SHOULD be dedicated to the usage of native IP TE solution, and SHOULD NOT be used by other BGP sessions that established by manual or non PCE initiated configuration.

BGP Peer Info Object-Class is TBD14

BGP Peer Info Object-Type is 1 for IPv4 and 2 for IPv6

The format of the BGP Peer Info object body for IPv4(Object-Type=1) is as follows:

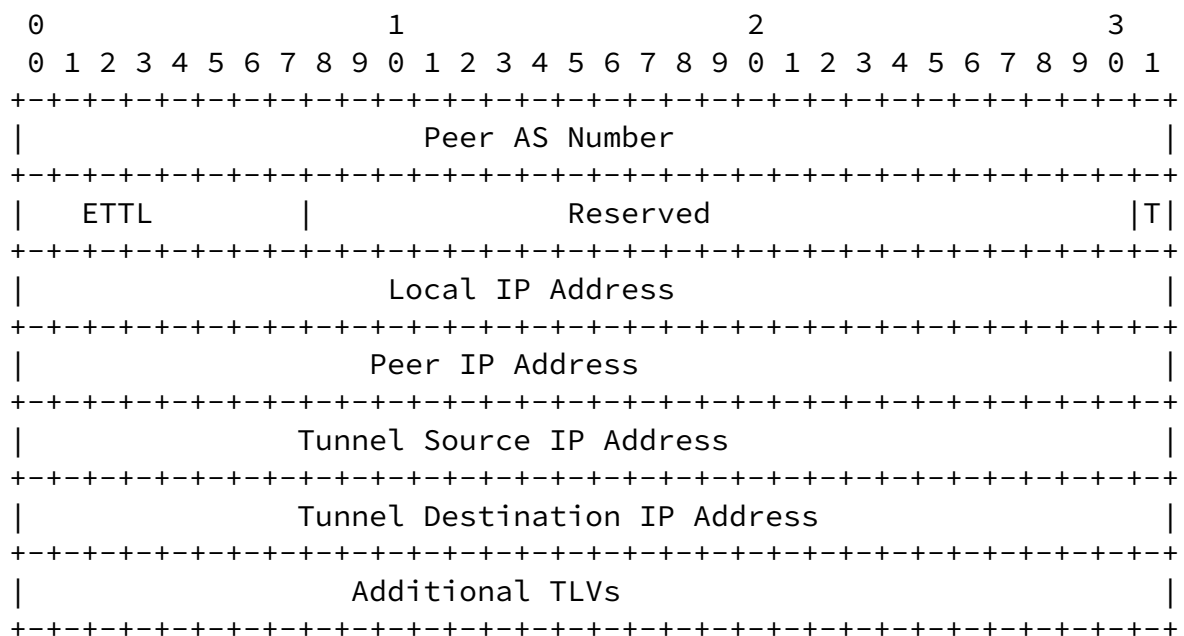


Figure 6: BGP Peer Info Object Body Format for IPv4

The format of the BGP Peer Info object body for IPv6(Object-Type=2) is as follows:

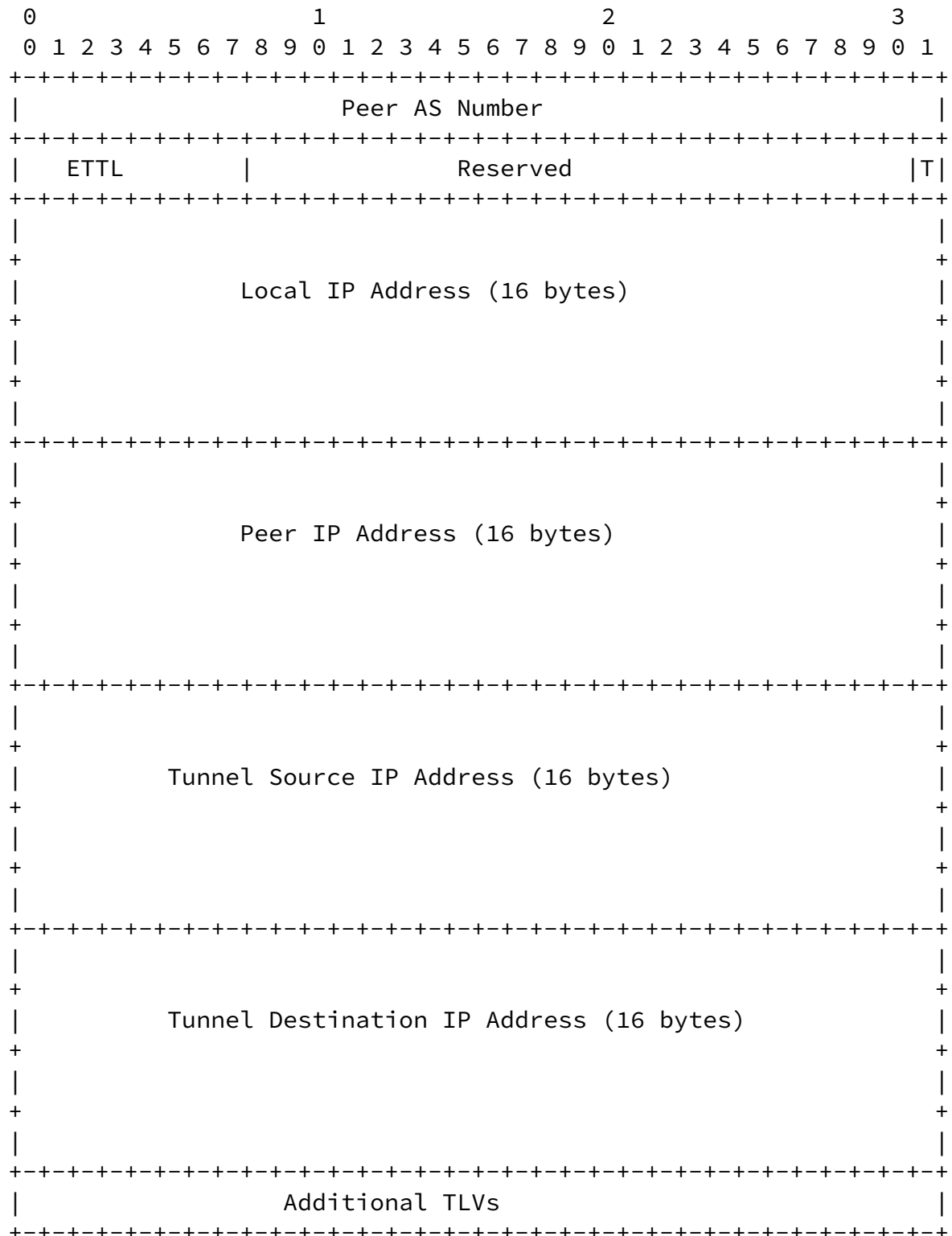


Figure 7: BGP Peer Info Object Body Format for IPv6

Peer AS Number: 4 Bytes, to indicate the AS number of Remote Peer.

ETTL: 1 Byte, to indicate the multihop count for EBGp session. It should be 0 and ignored when Local AS and Peer AS is same.

Reserved: is set to zero while sending, ignored on receipt.

T bit: Indicates whether the traffic that associated with the prefixes advertised via this BGP session is transported via IPinIP tunnel (when T bit is set) or not (when T bit is clear).

Local IP Address(4/16 Bytes): IP address of the local router, used to peer with other end router. When Object-Type is 1, length is 4 bytes; when Object-Type is 2, length is 16 bytes.

Peer IP Address(4/16 Bytes): IP address of the peer router, used to peer with the local router. When Object-Type is 1, length is 4 bytes; when Object-Type is 2, length is 16 bytes;

Tunnel Source IP Address(4/16 Bytes): IP address of the tunnel source, should be owned by the local router. When Object-Type is 1, length is 4 bytes; when Object-Type is 2, length is 16 bytes.

Tunnel Destination IP Address(4/16 Bytes): IP address of the tunnel destination, should be owned by the peer router. When Object-Type is 1, length is 4 bytes; when Object-Type is 2, length is 16 bytes. Should be different from the Peer IP Address.

Additional TLVs: TLVs that associated with this object, can be used to convey other necessary information for dynamic BGP session establishment. Their definition are out of the current document.

When PCC receives BPI object, with Object-Type=1, it should try to establish BGP session with the peer in AFI/SAFI=1/1; when PCC receives BPI object with Object-Type=2, it should try to establish the BGP session with the peer in AFI/SAFI=2/1. Other BGP capabilities, for example, Graceful Restart (GR) that enhance the BGP performance should also be negotiated and used by default.

[7.3.](#) Explicit Peer Route Object

The Explicit Peer Route object is defined to specify the explicit

peer route to the corresponding peer address on each device that is on the E2E assurance path. This Object should be sent to all the devices that locates on the E2E assurance path that calculated by PCE.

The path established by this object should have higher priority than other path calculated by dynamic IGP protocol, but should be lower priority than the static route configured by manual or NETCONF or by other means.

Explicit Peer Route Object-Class is TBD15.

Explicit Peer Route Object-Type is 1 for IPv4 and 2 for IPv6

The format of Explicit Peer Route object body for IPv4(Object-Type=1) is as follows:

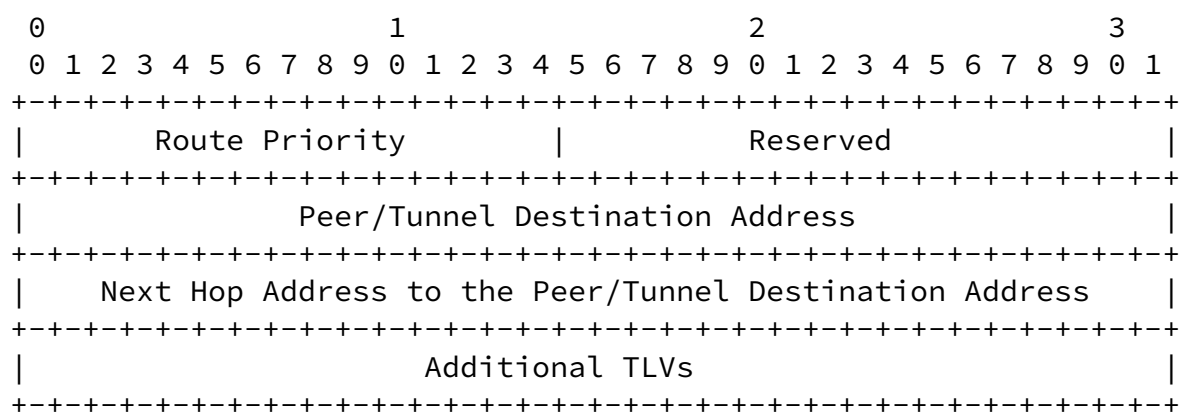


Figure 8: Explicit Peer Route Object Body Format for IPv4

The format of Explicit Peer Route object body for IPv6(Object-Type=2) is as follows:

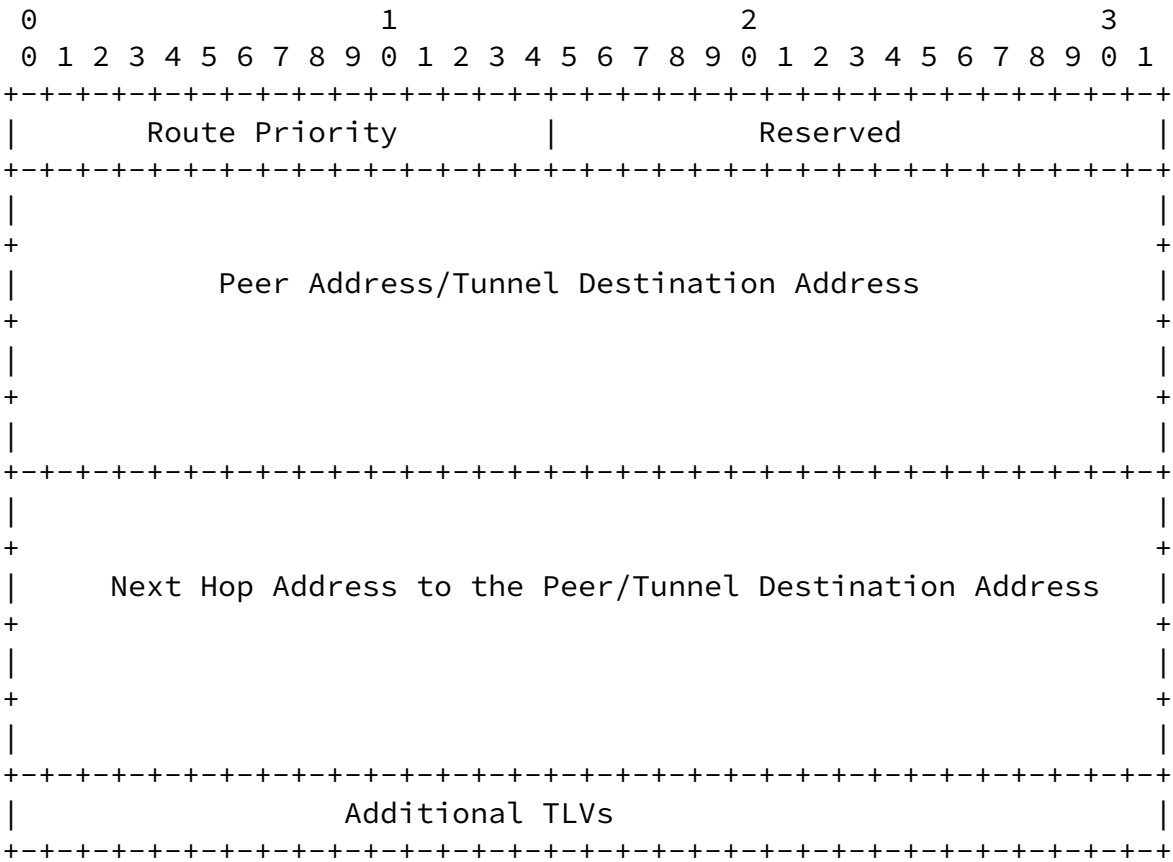


Figure 9: Explicit Peer Route Object Body Format for IPv6

Route Priority: 2 Bytes, The priority of this explicit route. The higher priority should be preferred by the device. This field is used to indicate the backup path at each hop.

Reserved: is set to zero while sending, ignored on receipt.

Peer/Tunnel Destination Address: To indicate the peer address(4/16 Bytes). When T bit is set in the associated BPI object, use the tunnel destination address in BPI object; when T bit is clear, use the peer address in BPI object.

Next Hop Address to the Peer/Tunnel Destination Address: To indicate the next hop address(4/16 Bytes) to the corresponding peer/tunnel destination address.

Additional TLVs: TLVs that associated with this object, can be used to convey other necessary information for explicit peer path establishment. Their definitions are out of the current document.

[7.4.](#) Peer Prefix Advertisement Object

The Peer Prefix Advertisement object is defined to specify the IP prefixes that should be advertised to the corresponding peer. This object should only be included and sent to the head/end router of the end2end path.

The prefixes information included in this object MUST only be advertised to the indicated peer, MUST NOT be advertised to other BGP peers.

Peer Prefix Advertisement Object-Class is TBD16

Peer Prefix Advertisement Object-Type is 1 for IPv4 and 2 for IPv6

The format of the Peer Prefix Advertisement object body is as follows:

0

1

2

3

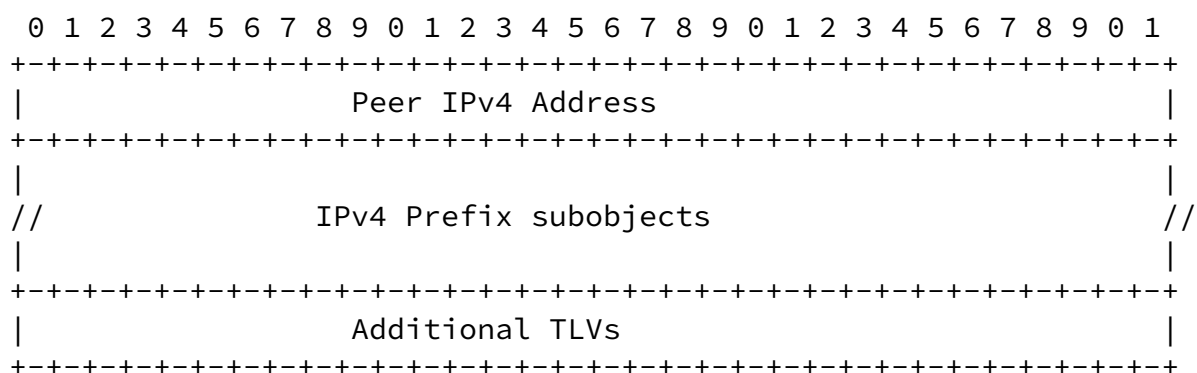


Figure 10: Peer Prefix Advertisement Object Body Format for IPv4

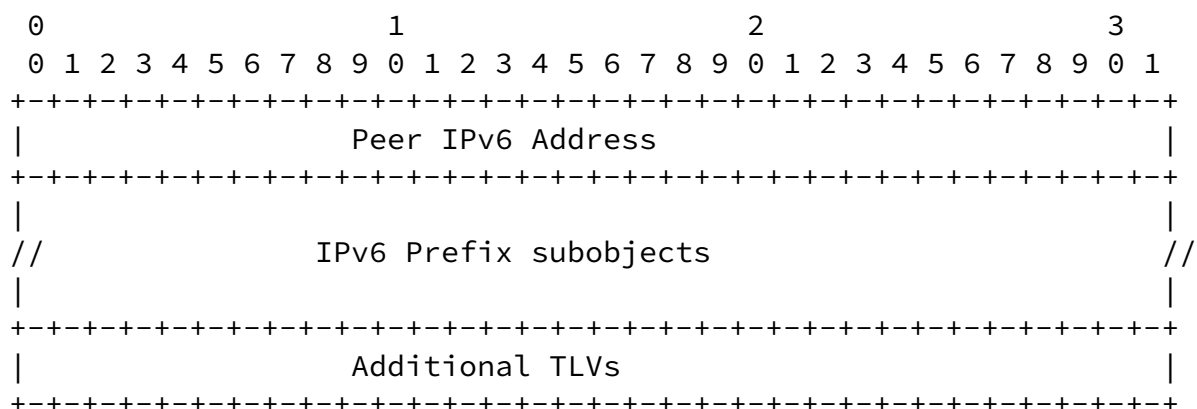


Figure 11: Peer Prefix Advertisement Object Body Format for IPv6

Peer IPv4 Address: 4 Bytes. Identifies the peer IPv4 address that the associated prefixes will be sent to.

IPv4 Prefix subobjects: List of IPv4 Prefix subobjects that defined in [\[RFC3209\]](#), identify the prefixes that will be sent to the peer that identified by Peer IPv4 Address List.

Peer IPv6 Address: 16 Bytes. Identifies the peer IPv6 address that the associated prefixes will be sent to.

IPv6 Prefix subobjects: List of IPv6 Prefix subobjects that defined in [\[RFC3209\]](#), identify the prefixes that will be sent to the peer that identified by Peer IPv6 Address List.

Additional TLVs: TLVs that associated with this object, can be used

to convey other necessary information for prefixes advertisement. Their definitions are out of the current document.

8. End to End Path Protection

[RFC8697] defines the path associations procedures between sets of Label Switched Path (LSP). Such procedures can also be used for the E2E path protection. To accomplish this, the PCE should attach the ASSOCIATION object with the EPR object in the PCInitiate message, with the association type set to 1 (Path Protection Association). The Extended Association ID that included within the Extended Association ID TLV, which is included in the ASSOCIATION object, should be set to the Symbolic Path Name of different E2E path. This PCInitiate should be sent to the head-end of the E2E path.

The head-end of the path can use the existing path detection mechanism(for example, Bidirectional Forwarding Detection [[RFC5880](#)]), to monitor the status of the active path. Once it detects the failure, it can switch the backup protection path immediately.

9. Re-Delegation and Clean up

In case of a PCE failure, a new PCE can gain control over the central controller instructions. As per the PCEP procedures in [[RFC8281](#)], the State Timeout Interval timer is used to ensure that a PCE failure does not result in automatic and immediate disruption for the services. Similarly, as per [[RFC9050](#)], the central controller instructions are not removed immediately upon PCE failure. Instead, they could be re-delegated to the new PCE before the expiration of this timer, or be cleaned up on the expiration of this timer. The allows for network clean up without manual intervention. The PCC MUST support the removal of CCI as one of the behaviors applied on expiration of the State Timeout Interval timer.

10. BGP Considerations

This draft defines the procedures and objects to create the BGP sessions and advertises the associated prefixes dynamically. Only the key information, for example peer IP addresses, peer AS number

are exchanged via the PCEP protocol. Other parameters that are needed for the BGP session setup should be derived from their default values, as described in [Section 7.2](#). Upon receives such key information, the BGP module on the PCC should try to accomplish the task that appointed by the PCEP protocol and report the status to the PCEP modules.

There is no influence to current implementation of BGP Finite State Machine(FSM). The PCEP cares only the success and failure status of BGP session, and act upon such information accordingly.

The error handling procedures related to incorrect BGP parameters are specified in [Section 6.1](#), [Section 6.2](#), and [Section 6.3](#). The handling of the dynamic BGP sessions and associated prefixes on PCE failure is described in [Section 9](#).

[11](#). New Error-Types and Error-Values Defined

A PCEP-ERROR object is used to report a PCEP error and is characterized by an Error-Type that specifies that type of error and an Error-value that provides additional information about the error. An additional Error-Type and several Error-values are defined to represent some the errors related to the newly defined objects, which are related to Native IP TE procedures.

Error-Type	Meaning	Error-value
TBD6	Native IP TE failure	
		0: Unassigned
		TBD7: Peer AS not match
		TBD8:Peer IP can't be reached
		TBD9:Local IP is in use
		TBD10:Remote IP is in use
		TBD11:Exist BGP session broken
		TBD12:Explicit Peer Route Error
		TBD17:EPR/BPI Peer Info mismatch
		TBD18:BPI/PPA Address Family mismatch
		TBD19:PPA/BPI Peer Info mismatch

Figure 12: Newly defined Error-Type and Error-Value

12. Deployment Considerations

The information transferred in this draft is mainly used for the light weight BGP session setup, explicit route deployment and the prefix distribution. The planning, allocation and distribution of the peer addresses within IGP should be accomplished in advanced and they are out of the scope of this draft.

[RFC8232] describes the state synchronization procedure between stateful PCE and PCC. The communication of PCE and PCC described in this draft should also follow this procedures, treat the three newly defined objects that associated with the same symbolic path name as the attribute of the same path in the LSP-DB.

When PCE detects one or some of the PCCs are out of control, it should recompute and redeploy the traffic engineering path for native IP on the active PCCs. When PCC detects that it is out of control of the PCE, it should clear the information that initiated by the PCE. The PCE should assure the avoidance of possible transient loop in such node failure when it deploy the explicit peer route on the PCCs.

If the established BGP session is broken after some time, the PCC should also report such error via PCErr message with Err-type=TBD6 and error value(Error-value=TBD11, Existing BGP session is broken). Upon receiving such PCErr message, the PCE should clear the prefixes advertisement on the previous BGP session, clear the explicit peer route to the previous peer address; select other Local_IP/Peer_IP pair to establish the new BGP session, deploy the explicit peer route to the new peer address, and advertises the prefixes on the new BGP session.

13. Implementation Status

[Note to the RFC Editor - remove this section before publication, as well as remove the reference to [RFC 7942](#).]

This section records the status of known implementations of the protocol defined by this specification at the time of posting of this Internet-Draft, and is based on a proposal described in [[RFC7942](#)]. The description of implementations in this section is intended to assist the IETF in its decision processes in progressing drafts to RFCs. Please note that the listing of any individual implementation here does not imply endorsement by the IETF. Furthermore, no effort has been spent to verify the information presented here that was supplied by IETF contributors. This is not intended as, and must not be construed to be, a catalog of available implementations or their features. Readers are advised to note that other implementations may exist.

According to [[RFC7942](#)], "this will allow reviewers and working groups to assign due consideration to documents that have the benefit of running code, which may serve as evidence of valuable experimentation and feedback that have made the implemented protocols more mature. It is up to the individual working groups to use this information as they see fit".

[13.1.](#) Proof of Concept based on ODL

.At the time of posting the -18 version of this document, there are no known implementations of this mechanism. A proof of concept for the overall design has been verified using another SBI protocol on the Open DayLight (ODL) controller.

Wang, et al.

Expires 22 September 2022

[Page 24]

Internet-Draft

PCEP Extension for Native IP Network

March 2022

[14.](#) Security Considerations

The setup of BGP sessions, prefix advertisement, and explicit peer route establishment are all controlled by the PCE. See [\[RFC4271\]](#) and [\[RFC4272\]](#) for BGP security considerations. Security consideration part in [\[RFC5440\]](#) and [\[RFC8231\]](#) should be considered. To prevent a bogus PCE sending harmful messages to the network nodes, the network devices should authenticate the validity of the PCE and ensure a secure communication channel between them. Mechanisms described in [\[RFC8253\]](#) should be used.

[15.](#) IANA Considerations

[15.1.](#) Path Setup Type Registry

[RFC8408] created a sub-registry within the "Path Computation Element Protocol (PCEP) Numbers" registry called "PCEP Path Setup Types". IANA is requested to allocate a new code point within this registry, as follows:

Value	Description	Reference
TBD1	Native IP TE Path	This document

[15.2.](#) PCECC-CAPABILITY sub-TLV's Flag field

[RFC9050] created a sub-registry within the "Path Computation Element Protocol (PCEP) Numbers" registry to manage the value of the PCECC-CAPABILITY sub-TLV's 32-bits Flag field. IANA is requested to allocate a new bit position within this registry, as follows:

Value	Description	Reference
TBD2(N)	NATIVE-IP-TE-CAPABILITY	This document

[15.3.](#) PCEP Object Types

IANA is requested to allocate new registry for the PCEP Object Type:

Wang, et al.

Expires 22 September 2022

[Page 25]

Internet-Draft

PCEP Extension for Native IP Network

March 2022

Object-Class Value	Name	Reference
44	CCI Object Object-Type TBD13: Native IP	This document
TBD14	BGP Peer Info Object-Type 1: IPv4 address 2: IPv6 address	This document
TBD15	Explicit Peer Route Object-Type 1: IPv4 address 2: IPv6 address	This document
TBD16	Peer Prefix Advertisement Object-Type 1: IPv4 address 2: IPv6 address	This document

[15.4.](#) PCEP-Error Object

IANA is requested to allocate new error types and error values within the "PCEP-ERROR Object Error Types and Values" sub-registry of the PCEP Numbers registry for the following errors::

Error-Type	Meaning	Error-value
6	Mandatory Object missing	

TBD4:Native IP object missing

[10](#) Reception of an invalid object

TBD3:PCECC NATIVE-IP-TE-CAPABILITY bit is

[19](#) Invalid Operation

TBD5:Only one of the BPI,EPR or PPA objec

TBD6 Native IP TE failure

TBD7:Peer AS not match

TBD8:Peer IP can't be reached

TBD9:Local IP is in use

TBD10:Remote IP is in use

TBD11:Exist BGP session broken

TBD12:Explicit Peer Route Error

TBD17:EPR/BPI Peer Info mismatch

TBD18:BPI/PPA Address Family mismatch

TBD19:PPA/BPI Peer Info mismatch

Wang, et al.

Expires 22 September 2022

[Page 26]

Internet-Draft PCEP Extension for Native IP Network

March 2022

[16.](#) Contributor

Dhruv Dhody has contributed the contents of this draft.

[17.](#) Acknowledgement

Thanks Mike Koldychev, Susan Hares, Siva Sivabalan, Adam Simpson for his valuable suggestions and comments.

[18.](#) Normative References

[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>>.

[RFC3209] Awduche, D., Berger, L., Gan, D., Li, T., Srinivasan, V., and G. Swallow, "RSVP-TE: Extensions to RSVP for LSP Tunnels", [RFC 3209](#), DOI 10.17487/RFC3209, December 2001, <<https://www.rfc-editor.org/info/rfc3209>>.

- [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>>.
- [RFC4272] Murphy, S., "BGP Security Vulnerabilities Analysis", [RFC 4272](#), DOI 10.17487/RFC4272, January 2006, <<https://www.rfc-editor.org/info/rfc4272>>.
- [RFC5440] Vasseur, JP., Ed. and JL. Le Roux, Ed., "Path Computation Element (PCE) Communication Protocol (PCEP)", [RFC 5440](#), DOI 10.17487/RFC5440, March 2009, <<https://www.rfc-editor.org/info/rfc5440>>.
- [RFC5880] Katz, D. and D. Ward, "Bidirectional Forwarding Detection (BFD)", [RFC 5880](#), DOI 10.17487/RFC5880, June 2010, <<https://www.rfc-editor.org/info/rfc5880>>.
- [RFC8174] Leiba, B., "Ambiguity of Uppercase vs Lowercase in [RFC 2119](#) Key Words", [BCP 14](#), [RFC 8174](#), DOI 10.17487/RFC8174, May 2017, <<https://www.rfc-editor.org/info/rfc8174>>.
- [RFC8231] Crabbe, E., Minei, I., Medved, J., and R. Varga, "Path Computation Element Communication Protocol (PCEP) Extensions for Stateful PCE", [RFC 8231](#), DOI 10.17487/RFC8231, September 2017, <<https://www.rfc-editor.org/info/rfc8231>>.

- [RFC8232] Crabbe, E., Minei, I., Medved, J., Varga, R., Zhang, X., and D. Dhody, "Optimizations of Label Switched Path State Synchronization Procedures for a Stateful PCE", [RFC 8232](#), DOI 10.17487/RFC8232, September 2017, <<https://www.rfc-editor.org/info/rfc8232>>.
- [RFC8253] Lopez, D., Gonzalez de Dios, O., Wu, Q., and D. Dhody, "PCEPS: Usage of TLS to Provide a Secure Transport for the Path Computation Element Communication Protocol (PCEP)", [RFC 8253](#), DOI 10.17487/RFC8253, October 2017, <<https://www.rfc-editor.org/info/rfc8253>>.
- [RFC8281] Crabbe, E., Minei, I., Sivabalan, S., and R. Varga, "Path Computation Element Communication Protocol (PCEP)

Extensions for PCE-Initiated LSP Setup in a Stateful PCE Model", [RFC 8281](#), DOI 10.17487/RFC8281, December 2017, <<https://www.rfc-editor.org/info/rfc8281>>.

- [RFC8283] Farrel, A., Ed., Zhao, Q., Ed., Li, Z., and C. Zhou, "An Architecture for Use of PCE and the PCE Communication Protocol (PCEP) in a Network with Central Control", [RFC 8283](#), DOI 10.17487/RFC8283, December 2017, <<https://www.rfc-editor.org/info/rfc8283>>.
- [RFC8408] Sivabalan, S., Tantsura, J., Minei, I., Varga, R., and J. Hardwick, "Conveying Path Setup Type in PCE Communication Protocol (PCEP) Messages", [RFC 8408](#), DOI 10.17487/RFC8408, July 2018, <<https://www.rfc-editor.org/info/rfc8408>>.
- [RFC8697] Minei, I., Crabbe, E., Sivabalan, S., Ananthakrishnan, H., Dhody, D., and Y. Tanaka, "Path Computation Element Communication Protocol (PCEP) Extensions for Establishing Relationships between Sets of Label Switched Paths (LSPs)", [RFC 8697](#), DOI 10.17487/RFC8697, January 2020, <<https://www.rfc-editor.org/info/rfc8697>>.
- [RFC8735] Wang, A., Huang, X., Kou, C., Li, Z., and P. Mi, "Scenarios and Simulation Results of PCE in a Native IP Network", [RFC 8735](#), DOI 10.17487/RFC8735, February 2020, <<https://www.rfc-editor.org/info/rfc8735>>.
- [RFC8821] Wang, A., Khasanov, B., Zhao, Q., and H. Chen, "PCE-Based Traffic Engineering (TE) in Native IP Networks", [RFC 8821](#), DOI 10.17487/RFC8821, April 2021, <<https://www.rfc-editor.org/info/rfc8821>>.

- [RFC9050] Li, Z., Peng, S., Negi, M., Zhao, Q., and C. Zhou, "Path Computation Element Communication Protocol (PCEP) Procedures and Extensions for Using the PCE as a Central Controller (PCECC) of LSPs", [RFC 9050](#), DOI 10.17487/RFC9050, July 2021, <<https://www.rfc-editor.org/info/rfc9050>>.

Authors' Addresses

Aijun Wang
China Telecom
Beiqijia Town, Changping District
Beijing
Beijing, 102209
China
Email: wangaj3@chinatelecom.cn

Boris Khasanov
Yandex LLC
Ulitsa Lva Tolstogo 16
Moscow
Email: bhassanov@yahoo.com

Sheng Fang
Huawei Technologies,Co.,Ltd
Huawei Bld., No.156 Beiqing Rd.
Beijing
China
Email: fsheng@huawei.com

Ren Tan
Huawei Technologies,Co.,Ltd
Huawei Bld., No.156 Beiqing Rd.
Beijing
China
Email: tanren@huawei.com

Chun Zhu
ZTE Corporation
50 Software Avenue, Yuhua District
Nanjing
Jiangsu, 210012
China
Email: zhu.chun1@zte.com.cn