

PCP working group
Internet-Draft
Intended status: Standards Track
Expires: September 1, 2011

D. Wing, Ed.
Cisco
S. Cheshire
Apple
M. Boucadair
France Telecom
R. Penno
Juniper Networks
F. Dupont
Internet Systems Consortium
February 28, 2011

Port Control Protocol (PCP)
draft-ietf-pcp-base-06

Abstract

Port Control Protocol allows a host to control how incoming IPv6 or IPv4 packets are translated and forwarded by a network address translator (NAT) or simple firewall to an IPv6 or IPv4 host, and also allows a host to optimize its NAT keepalive messages.

Status of this Memo

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

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

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."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire on September 1, 2011.

Copyright Notice

Copyright (c) 2011 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](http://trustee.ietf.org/license-info) 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

1.	Introduction	4
2.	Scope	5
2.1.	Deployment Scenarios	5
2.2.	Supported Transport Protocols	5
2.3.	Single-homed Customer Premises Network	5
3.	Terminology	6
4.	Relationship of PCP Server and its NAT	8
5.	Common Request and Response Header Format	9
5.1.	Request Header	9
5.2.	Response Header	10
5.3.	Options	11
5.4.	Result Codes	13
6.	General PCP Operation	14
6.1.	General PCP Client: Generating a Request	14
6.2.	General PCP Server: Processing a Request	15
6.3.	General PCP Client: Processing a Response	16
6.4.	Multi-Interface Issues	16
6.5.	Epoch	17
6.6.	Version negotiation	18
6.7.	General PCP Options	19
6.7.1.	UNPROCESSED	19
7.	Introduction to MAP and PEER OpCodes	20
7.1.	For Operating a Server	20
7.2.	For Reducing NAT Keepalive Messages	21
7.3.	For Operating a Symmetric Client/Server	22
8.	MAP OpCodes	23
8.1.	OpCode Packet Formats	23
8.2.	OpCode-Specific Result Codes	26
8.3.	OpCode-Specific Client: Generating a Request	27
8.4.	OpCode-Specific Server: Processing a Request	27
8.5.	OpCode-Specific Client: Processing a Response	29
8.6.	Mapping Lifetime and Deletion	30

8.7.	Subscriber Renumbering	31
8.8.	PCP Options for MAP OpCodes	31
8.8.1.	REMOTE_PEER_FILTER	31
8.8.2.	PREFER_FAILURE	34
8.8.3.	THIRD_PARTY	35
8.9.	PCP Mapping State Maintenance	35
8.9.1.	Recreating Mappings	35
8.9.2.	Maintaining Mappings	36
9.	PEER OpCodes	36
9.1.	OpCode Packet Formats	37
9.2.	OpCode-Specific Result Codes	40
9.3.	OpCode-Specific Client: Generating a Request	40
9.4.	OpCode-Specific Server: Processing a Request	41
9.5.	OpCode-Specific Client: Processing a Response	41
9.6.	PCP Options for PEER OpCodes	42
9.6.1.	THIRD_PARTY	42
10.	THIRD_PARTY Option for MAP and PEER OpCodes	42
11.	Deployment Considerations	45
11.1.	Maintaining Same External IP Address	45
11.2.	Ingress Filtering	45
11.3.	Per-Subscriber Port Forwarding Quota	45
12.	Deployment Scenarios	46
12.1.	Dual Stack-Lite	46
12.1.1.	Overview	46
12.2.	NAT64	47
12.3.	NAT44 and NAT444	47
12.4.	IPv6 Simple Firewall	47
13.	Security Considerations	47
13.1.	Denial of Service	48
13.2.	Ingress Filtering	48
13.3.	Validating Target Address	48
14.	IANA Considerations	48
14.1.	Port Number	48
14.2.	OpCodes	48
14.3.	Result Codes	49
14.4.	Options	49
15.	Acknowledgments	49
16.	References	49
16.1.	Normative References	49
16.2.	Informative References	50
Appendix A.	Changes	52
A.1.	Changes from draft-ietf-pcp-base-05 to -06	52
A.2.	Changes from draft-ietf-pcp-base-04 to -05	53
A.3.	Changes from draft-ietf-pcp-base-03 to -04	54
A.4.	Changes from draft-ietf-pcp-base-02 to -03	54
A.5.	Changes from draft-ietf-pcp-base-01 to -02	55
A.6.	Changes from draft-ietf-pcp-base-00 to -01	55
Authors'	Addresses	56

1. Introduction

Port Control Protocol (PCP) provides a mechanism to control how incoming packets are forwarded by upstream devices such as NAT64, NAT44, and firewall devices, and a mechanism to reduce application keepalive traffic. PCP is primarily designed to be implemented in the context of both Carrier-Grade NATs (CGN) and small NATs (e.g., residential NATs). PCP allows hosts to operate server for a long time (e.g., a webcam) or a short time (e.g., while playing a game or on a phone call) when behind a NAT device, including when behind a CGN operated by their Internet service provider.

PCP allows applications to create mappings from an external IP address and port to an internal (target) IP address and port. These mappings are required for successful inbound communications destined to machines located behind a NAT or a firewall.

After creating a mapping for incoming connections, it is necessary to inform remote computers about the IP address and port for the incoming connection. This is usually done in an application-specific manner. For example, a computer game would use a rendezvous server specific to that game (or specific to that game developer), and a SIP phone would use a SIP proxy. PCP does not provide this rendezvous function. The rendezvous function will support IPv4, IPv6, or both. Depending on that support and the application's support of IPv4 or IPv6, the PCP client will need an IPv4 mapping, an IPv6 mapping, or both.

Many NAT-friendly applications send frequent application-level messages to ensure their session will not be timed out by a NAT. These are commonly called "NAT keepalive" messages, even though they are not sent to the NAT itself (rather, they are sent 'through' the NAT). These applications can reduce the frequency of those NAT keepalive messages by using PCP to learn (or control) the NAT mapping lifetime. This helps reduce bandwidth on the subscriber's access network, traffic to the server, and battery consumption on mobile devices.

Many NATs and firewalls have included application layer gateways (ALGs) to create mappings for applications that establish additional streams or accept incoming connections. ALGs incorporated into NATs additionally modify the application payload. Industry experience has shown that these ALGs are detrimental to protocol evolution. PCP allows an application create its own mappings in NATs and firewalls, removing the incentive to deploy ALGs in NATs and firewalls.

2. Scope

2.1. Deployment Scenarios

PCP can be used in various deployment scenarios, including:

- o Dual Stack-Lite [[I-D.ietf-softwire-dual-stack-lite](#)], and;
- o NAT64, both Stateful [[I-D.ietf-behave-v6v4-xlate-stateful](#)] and Stateless [[I-D.ietf-behave-v6v4-xlate](#)], and;
- o Carrier-Grade NAT [[I-D.ietf-behave-lsn-requirements](#)], and;
- o Basic NAT [[RFC3022](#)], and;
- o Network Address and Port Translation (NAPT) [[RFC3022](#)], such as commonly deployed in residential NAT devices, and;
- o Layer-2 aware NAT [[I-D.miles-behave-l2nat](#)] and Dual-Stack Extra Lite [[I-D.arkko-dual-stack-extra-lite](#)], and;
- o IPv6 firewall control [[RFC6092](#)].

2.2. Supported Transport Protocols

The PCP OpCodes defined in this document are designed to support transport protocols that use a 16-bit port number (e.g., TCP, UDP, SCTP, DCCP). Transport protocols that do not use a port number (e.g., IPsec ESP), and the ability to use PCP to forward all traffic to a single default host (often nicknamed "DMZ"), are beyond the scope of this document.

2.3. Single-homed Customer Premises Network

The PCP machinery assumes a single-homed host model. That is, for a given IP version, only one default route exists to reach the Internet. This is important because after a PCP mapping is created and an inbound packet (e.g., TCP SYN) arrives at the host the outbound response (e.g., TCP SYNACK) has to go through the same path so the proper address rewriting takes place on that outbound response packet. This restriction exists because otherwise there would need to be one PCP server for each egress, because the host could not reliably determine which egress path packets would take, so the client would need to be able to reliably make the same internal/external mapping in every NAT gateway, which in general is not possible.

3. Terminology

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

Internal Host, Target Host:

A host served by a NAT gateway, or protected by a firewall. This is the host that receives the incoming traffic created by a PCP MAP request, or the host that initiated an implicit dynamic mapping (e.g., TCP SYN) across a firewall or a NAT.

Remote Host:

A host with which an Internal Host is communicating.

Target Address, External Address, Remote Peer Address:

The address of an Internal Host served by a NAT gateway (typically a private address [[RFC1918](#)]) or an Internal Host protected by a firewall.

An External Address is the address of an Internal Host as seen by other Remote Hosts on the Internet with which the Internal Host is communicating, after translation by any NAT gateways on the path. An External Address is generally a public routable (i.e., non-private) address. In the case of an Internal Host protected by a pure firewall, with no address translation on the path, its External Address is the same as its Internal Address.

A Remote Peer Address is the address of a Remote Host, as seen by the Internal Host. A Remote Address is generally a public routable address. In the case of a Remote Host that is itself served by a NAT gateway, the Remote Address may in fact be the Remote Host's External Address, but since this remote translation is generally invisible to software running on the Internal Host, the distinction can safely be ignored for the purposes of this document.

Third Party:

In the common case, an Internal Host manages its own Mappings using PCP requests, and the Internal Address of those Mappings is the same as the source IP address of the PCP request packet.

In the case where one device is managing Mappings on behalf of some other device, the presence of the THIRD_PARTY option in the MAP request signifies that the specified address, not the source IP address of the PCP request packet, should be used as the Internal Address for the Mapping. This can occur when PCP is proxied (e.g., PCP to PCP proxy, UPnP IGD to PCP proxy) or if the

target host does not implement PCP.

Mapping:

A NAT mapping creates a relationship between an internal (target) IP transport address and an external IP transport address. More specifically, it creates a translation rule where packets destined to the external IP and port are translated to the target IP and port. In the case of a pure firewall, the "Mapping" is the identity function, translating an internal port number to the same external port number and vice versa, and this "Mapping" indicates to the firewall that traffic to and from this internal port number is permitted to pass. See also Port Forwarding.

Mapping Types:

There are three different ways to create mappings: implicit dynamic mappings, explicit dynamic mappings, and static mappings. Implicit dynamic mappings are created as a result of a TCP SYN or outgoing UDP packet. Explicit dynamic mappings are created as a result of PCP MAP requests. Both implicit and explicit dynamic mappings are dynamic in the sense that they are created on demand, as requested (implicitly or explicitly) by the Internal Host, and have a lifetime, after which they are automatically deleted unless the lifetime is extended by action by the Internal Host. Static mappings are created by manual configuration (e.g., command language interface or web page) and differ from dynamic mappings in that their lifetime is typically infinite (they exist until manually removed) but otherwise they behave exactly the same as dynamic mappings. For example, a PCP MAP request to create a mapping that already exists as a static mapping will return a successful result, confirming that the requested mapping exists.

Port Forwarding, Port Mapping:

Port forwarding (or port mapping) allows a host to receive traffic sent to a specific IP address and port.

In the context of a NAT or NAPT, the Internal Address and External Address are different. In the context of a pure firewall, the Internal Address and External Address are the same. In both contexts, if an internal host is listening to connections on a specific port (that is, operating as a server), the external IP address and port number need to be port forwarded to the internal IP address and port number, which may be the same, in the case of a pure firewall. In the context of a NAPT, it is possible that both the IP address and port are modified. For example with a NAPT, a webcam might be listening on port 80 on its internal address 192.168.1.1, while its publicly-accessible external address is 192.0.2.1 and port is 12345. The NAT does 'port forwarding' of one to the other.

PCP Client:

A PCP software instance responsible for issuing PCP requests to a PCP server. One or several PCP Clients can be embedded in the same host. Several PCP Clients can be located in the same local network of a given subscriber. A PCP Client can issue PCP request on behalf of a third party device of the same subscriber. An interworking function, from UPnP IGD to PCP, or from NAT-PMP [[I-D.cheshire-nat-pmp](#)] is another example of a PCP Client. A PCP server in a NAT gateway that is itself a client of another NAT gateway (nested NAT) may itself act as a PCP client to the upstream NAT.

PCP Server:

A network element which receives and processes PCP requests from a PCP client. See also [Section 4](#).

Interworking Function:

a functional element responsible for interworking another protocol with PCP. For example interworking between UPnP IGD [[IGD](#)] with PCP or NAT-PMP [[I-D.cheshire-nat-pmp](#)] and PCP.

subscriber:

an entity provided access to the network. In the case of a commercial ISP, this is typically a single home.

host:

a device which can have packets sent to it, as a result of PCP operations. A host is not necessarily a PCP client.

5-tuple The 5 pieces of information that fully identify a flow: source IP address, destination IP address, protocol, source port number, destination port number.

[4. Relationship of PCP Server and its NAT](#)

The PCP server receives PCP requests. The PCP server might be integrated within the NAT or firewall device (as shown in Figure 1) which is expected to be a common deployment.

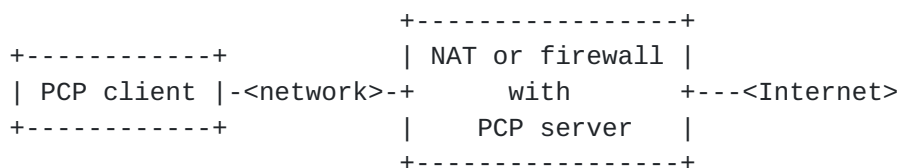


Figure 1: NAT or Firewall with Embedded PCP Server

It is also possible to operate the PCP server in a separate device from the NAT, so long as such operation is indistinguishable from the PCP client's perspective.

5. Common Request and Response Header Format

All PCP messages contain a request (or response) header containing an opcode, any relevant opcode-specific information, and zero or more options. The packet layout for the common header, and operation of the PCP client and PCP server are described in the following sections. The information in this section applies to all OpCodes. Behavior of the OpCodes defined in this document is described in [Section 8](#) and [Section 9](#).

5.1. Request Header

All requests have the following format:

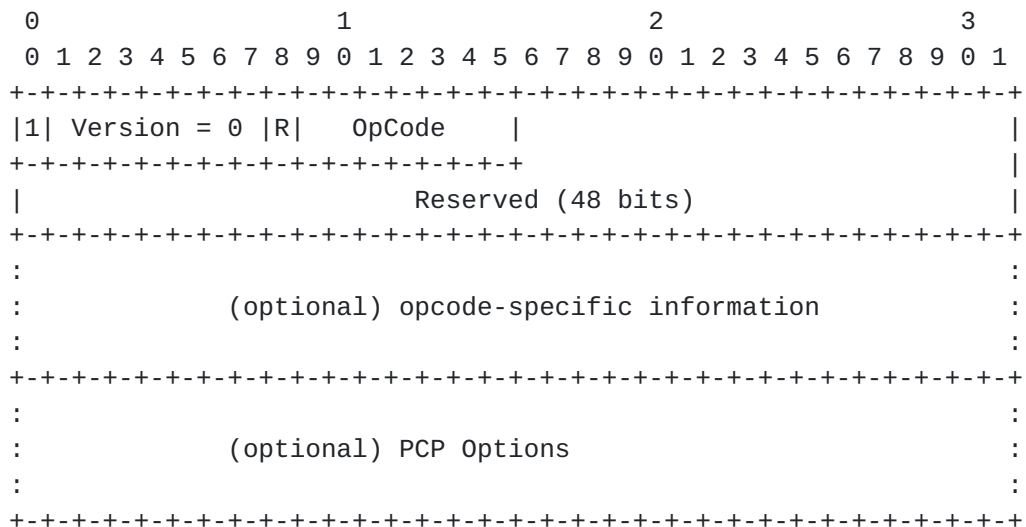


Figure 2: Common Request Packet Format

These fields are described below:

- 1 A single bit set to 1. This allows DTLS and non-DTLS to be multiplexed on same port, should a future update to this specification add DTLS support.

Version: This document specifies protocol version 0. Should later updates to this document specify different message formats with a version number greater than zero, the first two bytes of those new message formats will still contain the version number and opcode as shown here, so that a PCP server receiving a message format

newer or older than the version(s) it understands can still parse enough of the message to correctly identify the version number, and determine whether the problem is that this server is too old and needs to be updated to work with the PCP client, or whether the PCP client is too old and needs to be updated to work with this server.

R: Indicates Request (0) or Response (1). All Requests MUST use 0.

OpCode: Opcodes are defined in [Section 8](#) and [Section 9](#).

Reserved: 48 reserved bits, MUST be sent as 0 and MUST be ignored when received.

5.2. Response Header

All responses have the following format:

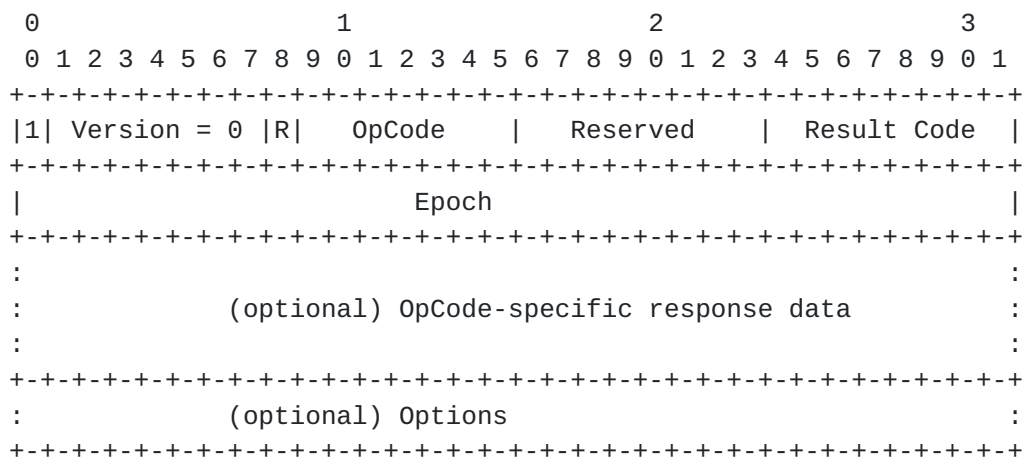


Figure 3: Common Response Packet Format

These fields are described below:

1 A single bit set to 1.

Ver: This document specifies protocol version 0. Should later updates to this document specify different message formats with a version number greater than zero, the first four bytes of those new message formats will still contain the version number, opcode, and result code, as shown here, so that a PCP client receiving a message format newer or older than the version(s) it understands can still parse enough of the message to correctly identify the version number, and determine whether the problem is that this client is too old and needs to be updated to work with the PCP server, or whether the PCP server is too old and needs to be

updated to work with this client.

R: Indicates Request (0) or Response (1). All Responses MUST use 1.

OpCode: The OpCode value, copied from the request.

Reserved: 8 reserved bits, MUST be sent as 0, MUST be ignored when received. This is set by the server.

Result Code: The result code for this response. See [Section 5.4](#) for values. This is set by the server.

Epoch: The server's Epoch value. See [Section 6.5](#) for discussion. This value is set both success and error responses.

5.3. Options

A PCP OpCode can be extended with an Option. Options can be used in requests and responses. It is anticipated that Options will include information which are associated with the normal function of an OpCode. For example, an Option could indicate DSCP [[RFC2474](#)] markings to apply to incoming or outgoing traffic associated with a PCP mapping, or an Option could include descriptive text (e.g., "for my webcam").

Options use the following Type-Length-Value format:

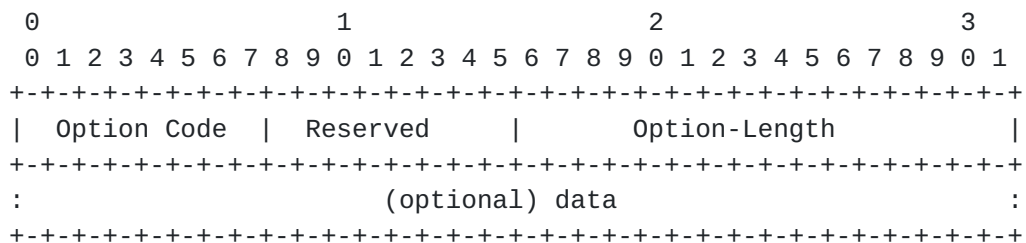


Figure 4: Options Header

The description of the fields is as follows:

Option Code: Option code, 8 bits. The first bit of the option code is the "0" (optional) bit. If clear, it indicates the option is mandatory to process (that is, non-optional). If set, it indicates the option is optional.

Reserved: MUST be set to 0 on transmission and MUST be ignored on reception.

Option-Length: Indicates in units of 4 octets the length of the enclosed data. Options with length of 0 are allowed.

data: Option data. The option data MUST end on a 32 bit boundary, padded with 0's when necessary.

A given Option MAY be included in a request or a response, as permitted by that Option. If a given Option was included in a request, and understood and processed by the PCP server, it MUST be included in the response. The handling of an Option by the PCP client and PCP server MUST be specified in an appropriate document and must include whether the PCP Option can appear (one or more times) in a request, and indicate the contents of the Option in the request and in the response. If several Options are included in a PCP request or response, they MUST be encoded in numeric order by the PCP client and are processed in the order received. The server MUST reject requests that have mis-ordered options with the MISORDERED_OPTIONS error, and this also includes checking optional-to-process options.

If, while processing an option, an error is encountered that causes a PCP error response to be generated, the PCP request causes no state change in the PCP server or the PCP-controlled device (i.e., it rolls back any changes it might have made while processing the request). The response MUST encode the Options in the same order, but may omit some PCP Options in the response, as is necessary to indicate the PCP server does not understand that Option or that Option is not permitted to be included in responses by the definition of the Option itself. Additional Options included in the response (if any) MUST be included at the end. A certain Option MAY appear more than once in a request or in a response, if permitted by the definition of the Option itself. If the Option's definition allows the Option to appear once but it appears more than once in a request, the PCP server MUST respond with the MALFORMED_OPTION result code; if this occurs in a response, the PCP client processes the first occurrence and ignores the other occurrences as if they were not present.

If the "0" bit in the OpCode is clear,

- o the PCP server MUST only generate a positive PCP response if it can successfully process the PCP request and this Option.
- o if the PCP server does not implement this Option, or cannot perform the function indicated by this Option (e.g., due to a parsing error with the option), it MUST generate a failure response with code UNSUPP_OPTION or MALFORMED_OPTION (as appropriate) and include the UNPROCESSED option in the response ([Section 6.7.1](#)).

If the "0" bit is set, the PCP server MAY process or ignore this Option, entirely at its discretion.

To enhance interoperability, newly defined Options SHOULD NOT be interdependent with each other. Option definitions MUST include the information below:

This Option:

name: <mnemonic>

number: <value>

purpose:

is valid for OpCodes: <list of OpCodes>

length: <rules for length>

may appear in: <requests/responses/both>

maximum occurrences: <count>

[5.4.](#) Result Codes

The following result codes may be returned as a result of any OpCode received by the PCP server. The only success result code is 0, other values indicate an error. If a PCP server has encountered multiple errors during processing of a request, it SHOULD use the most specific error message.

- 0 SUCCESS, success
- 1 MALFORMED_REQUEST, a general catch-all error.
- 2 UNSUPP_OPCODE, unsupported OpCode.
- 3 UNSUPP_OPTION, unsupported Option. This error only occurs if the Option is in the mandatory-to-process range.
- 4 MALFORMED_OPTION, malformed Option (e.g., exists too many times, invalid length).
- 5 UNSPECIFIED_ERROR, server encountered unspecified error.

- 6 UNSUPP_VERSION, unsupported version.
- 7 MISORDERED_OPTIONS, multiple options were in the request, but were not in the required lower..higher order.

Additional result codes, specific to the OpCodes and Options defined in this document, are listed in [Section 8.2](#), [Section 9.2](#), and [Section 10](#).

6. General PCP Operation

PCP messages MUST be sent over UDP, and the PCP server MUST listen for PCP requests on the PCP port number, 44323. Every PCP request generates a response, so PCP does not need to run over a reliable transport protocol.

PCP is idempotent, so if the PCP client sends the same request multiple times and the PCP server processes those requests, the same result occurs. The order of operation is that a PCP client generates and sends a request to the PCP server which processes the request and generates a response back to the PCP client.

6.1. General PCP Client: Generating a Request

This section details operation specific to a PCP client, for any OpCode. Procedures specific to the MAP OpCodes are described in [Section 8](#), and procedures specific to the PEER OpCodes are described in [Section 9](#).

Prior to sending its first PCP message, the PCP client determines which servers to use. The PCP client performs the following steps to determine its PCP server(s):

1. if a PCP server is configured (e.g., in a configuration file), the address(es) of the PCP server(s) used as the list of PCP server(s), else;
2. if DHCP indicates the PCP server(s), the address(es) of the indicated PCP server(s) are used as the the list of PCP server(s), else;
3. the address of the default router is used as the PCP server.

With that list of PCP servers, the PCP client formulates its PCP request. The PCP request contains a PCP common header, PCP OpCode and payload, and (possibly) Options. It initializes a retransmission timer to 4 seconds. (As with all UDP or TCP clients on any operating

system, when several PCP clients are embedded in the same host, each uses a distinct source port number to disambiguate their requests and replies.) The PCP client sends a PCP message to each server in sequence, waiting for a response until its timer expires. Once a PCP client has successfully communicated with a PCP server, it continues communicating with that PCP server until that PCP server has not responded to 5 retransmissions, which causes the PCP client to attempt to re-iterate the procedure starting with the first PCP server on its list. If a hard ICMP error is received the PCP client SHOULD immediately abort trying to contact that PCP server (see [Section 2 of \[RFC5461\]](#) for discussion of ICMP and ICMPv6 hard errors). If no response is received from any of those servers, it doubles its retransmission timer and tries each server again. This is repeated 4 times (for a total of 5 transmissions to each server). If, after these transmissions, the PCP client has still not received a response, the PCP client SHOULD abort the procedure.

Upon receiving a response (success or error), the PCP client does not change to a different PCP server. That is, it does not "shop around" trying to find a PCP server to service its (same) request.

6.2. General PCP Server: Processing a Request

This section details operation specific to a PCP server.

Upon receiving a PCP request message, the PCP server parses and validates it. A valid request contains a valid PCP common header, one valid PCP Opcode, and zero or more Options (which the server might or might not comprehend). If an error is encountered during processing, the server generates an error response which is sent back to the PCP client. Processing an OpCode and the Options are specific to each OpCode.

If the received message is shorter than 4 octets, has the R bit set, or the first bit is clear, the request is simply dropped. If the version number is not supported, a response is generated containing the UNSUPP_VERSION response code and the protocol version which the server does understand (if the server understands a range of protocol versions then it returns the supported version closest to the version in the request).

If the OpCode is not supported, a response is generated with the UNSUPP_OPCODE response code. If the length of the request exceeds 1024 octets or is not a multiple of 4 octets, it is invalid. Invalid requests are handled by copying up to 1024 octets of the request into the response, setting the response code to MALFORMED_REQUEST, and zero-padding the response to a multiple of 4 octets if necessary.

Error responses have the same packet layout as success responses, with fields copied from the request copied into the response, and other fields assigned by the PCP server MUST be cleared to 0.

6.3. General PCP Client: Processing a Response

The PCP client receives the response and verifies the source IP address and port belong to the PCP server of an outstanding request. It validates the version number and OpCode matches an outstanding request. Responses shorter than 12 octets, longer than 1024 octets, or not a multiple of 4 octets are invalid and ignored, likely causing the request to be re-transmitted. The response is further matched by comparing fields in the response OpCode-specific data to fields in the request OpCode-specific data. After a successful match with an outstanding request, the PCP client checks the Epoch field to determine if it needs to restore its state to the PCP server (see [Section 6.5](#)).

If the response code is 0, the PCP client knows the request was successful.

If the response code is not 0, the request failed. If the response code is UNSUPP_VERSION, processing continues as described in [Section 6.6](#). If the response code is SERVER_OVERLOADED, clients SHOULD NOT send *any* further requests to that PCP server for the time indicated by that OpCode's response, if present (e.g., the lifetime field of a MAP response), or 30 seconds have elapsed. For other error response codes, The PCP client SHOULD NOT resend the same request for the time indicated by that OpCode's response, if present (e.g., the lifetime field of a MAP response), or 30 seconds have elapsed.

If the PCP client has discovered a new PCP server (e.g., connected to a new network), the PCP client MAY immediately begin communicating with this PCP server, without regard to hold times from communicating with a previous PCP server.

6.4. Multi-Interface Issues

Hosts which desire a PCP mapping might be multi-interfaced (i.e., own several logical/physical interfaces). Indeed, a host can be configured with several IPv4 addresses (e.g., WiFi and Ethernet) or dual-stacked. These IP addresses may have distinct reachability scopes (e.g., if IPv6 they might have global reachability scope as for GUA (Global Unicast Address) or limited scope such as ULA (Unique Local Address, [RFC4193](#))).

IPv6 addresses with global reachability scope SHOULD be used as the

source interface when generating a PCP request. IPv6 addresses with limited scope (e.g., ULA [RFC4193]), SHOULD NOT be used as the source interface when generating a PCP request. If IPv6 privacy addresses [RFC4941] are used for PCP mappings, a new PCP request will need to be issued whenever the IPv6 privacy address is changed. This PCP request SHOULD be sent from the IPv6 privacy address itself. It is RECOMMENDED that mappings to the previous privacy address be deleted.

Due to the ubiquity of IPv4 NAT, IPv4 addresses with limited scope (e.g., [RFC1918]) MAY be used as the source interface when generating a PCP request.

As mentioned in [Section 2.3](#), only single-homed CP routers are in scope. Therefore, there is no viable scenario where a host located behind a CP router is assigned with two GUA addresses belonging to the same global IPv6 prefix.

6.5. Epoch

Every PCP response sent by the PCP server includes an Epoch field. This field increments by 1 every second, and indicates to the PCP client if PCP state needs to be restored. If the PCP server resets or loses the state of its explicit dynamic Mappings (that is, those mappings created by PCP MAP requests), due to reboot, power failure, or any other reason, it MUST reset its Epoch time to 0. Similarly, if the public IP address(es) of the NAT (controlled by the PCP server) changes, the Epoch MUST be reset to 0. A PCP server MAY maintain one Epoch value for all PCP clients, or MAY maintain distinct Epoch values for each PCP client; this choice is implementation-dependent.

Whenever a client receives a PCP response, the client computes its own conservative estimate of the expected Epoch value by taking the Epoch value in the last packet it received from the gateway and adding 7/8 (87.5%) of the time elapsed since that packet was received. If the Epoch value in the newly received packet is less than the client's conservative estimate by more than one second, then the client concludes that the PCP server lost state, and the client MUST immediately renew all its active port mapping leases as described in [Section 8.9.1](#).

When the PCP server reduces its Epoch value, the PCP clients will send PCP requests to refresh their mappings. The PCP server needs to be scaled appropriately to accomodate this traffic. Because PCP lacks a mechanism to simultaneously inform all PCP clients of the Epoch value, the PCP clients will not flood the PCP server simultaneously when the PCP server reduces its Epoch value.

6.6. Version negotiation

A PCP client sends its requests using PCP version number 0. Should later updates to this document specify different message formats with a version number greater than zero it is expected that PCP servers will still support version 0 in addition to the newer version(s). However, in the event that a server returns a response with error code UNSUPP_VERSION, the client MAY log an error message to inform the user that it is too old to work with this server, and the client SHOULD set a timer to retry its request in 30 minutes (in case this was a temporary condition and the server configuration is changed to rectify the situation).

If future PCP versions greater than zero are specified, version negotiation is expected to proceed as follows:

1. If a client or server supports more than one version it SHOULD support a contiguous range of versions -- i.e., a lowest version and a highest version and all versions in between.
2. Client sends first request using highest (i.e., presumably 'best') version number it supports.
3. If server supports that version it responds normally.
4. If server does not support that version it replies giving a response containing the error code UNSUPP_VERSION, and the closest version number it does support (if the server supports a range of versions higher than the client's requested version, the server returns the lowest of that supported range; if the server supports a range of versions lower than the client's requested version, the server returns the highest of that supported range).
5. If the client receives an UNSUPP_VERSION response containing a version it does support, it records this fact and proceeds to use this message version for subsequent communication with this PCP server (until a possible future UNSUPP_VERSION response if the server is later updated, at which point the version negotiation process repeats).
6. If the client receives an UNSUPP_VERSION response containing a version it does not support then the client MAY log an error message to inform the user that it is too old to work with this server, and the client SHOULD set a timer to retry its request in 30 minutes.

6.7. General PCP Options

The following options can appear in certain PCP responses.

6.7.1. UNPROCESSED

If the PCP server cannot process a mandatory-to-process option, for whatever reason, it includes the UNPROCESSED Option in the response, shown in Figure 5. This helps with debugging interactions between the PCP client and PCP server. For simplicity, no more than 4 options can be encoded. This option MUST NOT appear more than once in a PCP response, no matter how many PCP options appeared in the request and were unprocessed by the PCP server. If only one Option code was unprocessed, that option code is placed in option-code-1 (and the other three fields are set to zero), if two Option codes were unprocessed, their option codes are placed in option-code-1 and option-code-2, and so on. If a certain Option appeared more than once in the PCP request, that Option value only appears once in the option-code fields. The order of the Options in the PCP request has no relationship with the order of the Option values in this UNPROCESSED Option. This Option MUST NOT appear in a response unless the associated request contained at least one mandatory-to-process Option. This Option MUST NOT appear more than once.

The UNPROCESSED option is formatted as follows:

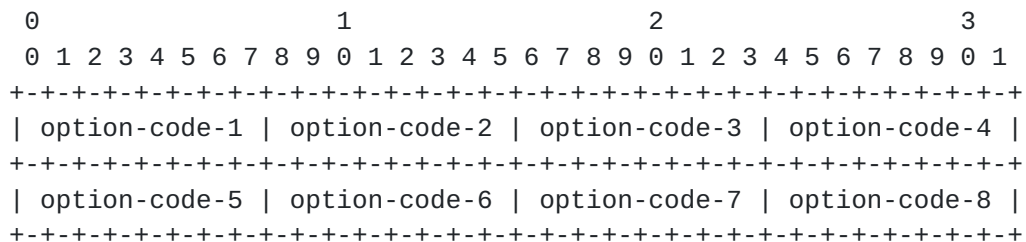


Figure 5: UNPROCESSED option

This Option:

name: UNPROCESSED

number: 1

purpose: indicates which PCP options in the request are not supported by the PCP server

is valid for OpCodes: all

length: 1

may appear in: responses, and only if the response code is non-zero.

maximum occurrences: 1

7. Introduction to MAP and PEER OpCodes

There are three uses for the MAP and PEER OpCodes defined in this document: a host operating a server (and wanting an incoming connection), a host operating a client (and wanting to optimize the application keepalive traffic), and a host operating a client and server on the same port. These are discussed in the following sections.

When operating a server ([Section 7.1](#) and [Section 7.3](#)) the PCP client knows if it wants an IPv4 listener, IPv6 listener, or both on the Internet. The PCP client also knows if it has an IPv4 interface on itself or an IPv6 interface on itself. It takes the union of this knowledge to decide to send a one or two MAP requests for each of its interfaces. Applications that embed IP addresses in payloads (e.g., FTP, SIP) will find it beneficial to avoid address family translation, if possible.

7.1. For Operating a Server

A host operating a server (e.g., a web server) listens for traffic on a port, but the server never initiates traffic from that port. For this to work across a NAT or a firewall, the application needs to (a) create a mapping from a public IP address and port to itself as described in [Section 8](#) and (b) publish that public IP address and port via some sort of rendezvous server (e.g., DNS, a SIP message, a proprietary protocol). Publishing the public IP address and port is out of scope of this specification. To accomplish (a), the application follows the procedures described in this section.

As normal, the application needs to begin listening to a port, and to ensure that it can get exclusive use of that port it needs to choose a port that is not in the operating system's ephemeral port range. Then, the application constructs a PCP message with the appropriate MAP OpCode depending on if it is listening on an IPv4 or IPv6 interface and if it wants a public IPv4 or IPv6 address.

The following pseudo-code shows how PCP can be reliably used to operate a server:

```
/* start listening on the local server port */
int s = socket(...);
internal_sockaddr = ...;
bind(s, &internal_sockaddr, ...);
listen(s, ...);
requested_external_sockaddr = 0;
pcp_send_map_request(internal_sockaddr,
    requested_external_sockaddr, &assigned_external_sockaddr,
    requested_lifetime, &assigned_lifetime);
update_rendezvous_server("Client 12345", assigned_external_sockaddr);
while (1) {
    int c = accept(s, ...);
    /* ... */
}
```

Figure 6: Pseudo-code for using PCP to operate a server

7.2. For Reducing NAT Keepalive Messages

A host operating a client (e.g., XMPP client, SIP client) sends from a port but never accepts incoming connections on this port. It wants to ensure the flow to its server is not terminated (due to inactivity) by an on-path NAT or firewall. To accomplish this, the applications uses the procedure described in this section.

Middleboxes such as NATs or firewalls need to see occasional traffic or will terminate their session state, causing application failures. To avoid this, many applications routinely generate keepalive traffic for the primary (or sole) purpose of maintaining state with such middleboxes. Applications can reduce such application keepalive traffic by using PCP.

Note: For reasons beyond NAT, an application may find it useful to perform application-level keepalives, such as to detect a broken path between the client and server, detect a crashed server, or detect a powered-down client. These keepalives are not related to maintaining middlebox state, and PCP cannot do anything useful to reduce those keepalives.

To use PCP for this function, the applications first connects to its server, as normal. Afterwards, it issues a PCP request with the PEER4 or PEER6 OpCode as described in [Section 9](#). The PEER4 OpCode is used if the host is using IPv4 for its communication to its peer; PEER6 if using IPv6. The same 5-tuple as used for the connection to the server is placed into the PEER4 or PEER6 payload.

The following pseudo-code shows how PCP can be reliably used with a dynamic socket, for the purposes of reducing application keepalive messages:

```
int s = socket(...);
connect(s, &remote_peer, ...);
getsockname(s, &internal_address, ...);
external_address = 0;
pcp_send_peer_request(internal_address,
    requested_external_address, &assigned_external_address,
    remote_peer, requested_lifetime, &assigned_lifetime);
```

Figure 7: Pseudo-code using PCP with a dynamic socket

7.3. For Operating a Symmetric Client/Server

A host operating a client and server on the same port (e.g., Symmetric RTP [[RFC4961](#)] or SIP Symmetric Response Routing (rport) [[RFC3581](#)]) first establishes a local listener, (usually) sends the local and public IP addresses and ports to a rendezvous service (which is out of scope of this document), and (usually) initiates outbound connections from that same source address. To accomplish this, the application uses the procedure described in this section.

An application that is using the same port for outgoing connections as well as incoming connections MUST first signal its operation of a server using the PCP MAP OpCode, as described in [Section 8](#), and receive a positive PCP response before it sends any packets from that port.

Discussion: Although reversing those steps is tempting (to eliminate the PCP round trip before a packet can be sent from that port) and will work if the NAT has endpoint-independent mappings (EIM) behavior, reversing the steps will fail if the NAT does not have EIM behavior. With a non-EIM NAT, the implicit mapping created by an outgoing TCP SYN and the explicit mapping created using the MAP OpCode will cause different ports to be assigned (which is not desirable; after all, the application is using the same port for outgoing and incoming traffic on purpose) and they will generally also have different lifetimes. PCP does not attempt to change or dictate how a NAT creates its mappings (endpoint independent mapping, or otherwise) so there is no assurance that an implicit mapping will be EIM or non-EIM. Thus, it is necessary for applications to first signal its operation of a server using the PCP MAP OpCode.

The following pseudo-code shows how PCP can be used to operate a symmetric client and server:

```
/* start listening on the local server port */
int s = socket(...);
internal_sockaddr = ...;
bind(s, &internal_sockaddr, ...);
listen(s, ...);
requested_external_sockaddr = 0;
pcp_send_map_request(internal_sockaddr,
    requested_external_sockaddr, &assigned_external_sockaddr,
    requested_lifetime, &assigned_lifetime);
update_rendezvous_server("Client 12345", assigned_external_sockaddr);
send_packet(s, "Hello World");
while (1) {
    int c = accept(s, ...);
    /* ... */
}
```

Figure 8: Pseudo-code for using PCP to operate a symmetric client/server

8. MAP OpCodes

This section defines four OpCodes which control forwarding from a NAT (or firewall) to an internal target host. They are:

- MAP4=0: create a mapping between an internal target address and external IPv4 address (e.g., NAT44, NAT64, or firewall)
- MAP6=1: create a mapping between an internal target address and external IPv6 address (e.g., NAT46, NAT66, or firewall)

The internal target address is the source IP address of the PCP request message itself, unless the THIRD_PARTY option is used.

The operation of these OpCodes is described in this section.

8.1. OpCode Packet Formats

The two MAP OpCodes (MAP4, MAP6) share a similar packet layout for both requests and responses. Because of this similarity, they are shown together. For both of the MAP OpCodes, if the assigned external IP address and assigned external port both match the request's source IP address and MAP OpCode's internal IP address, the functionality is purely a firewall; otherwise it pertains to a network address translator which might also perform firewall

functions.

The following diagram shows the request packet format for MAP4 and MAP6. This packet format is aligned with the response packet format:

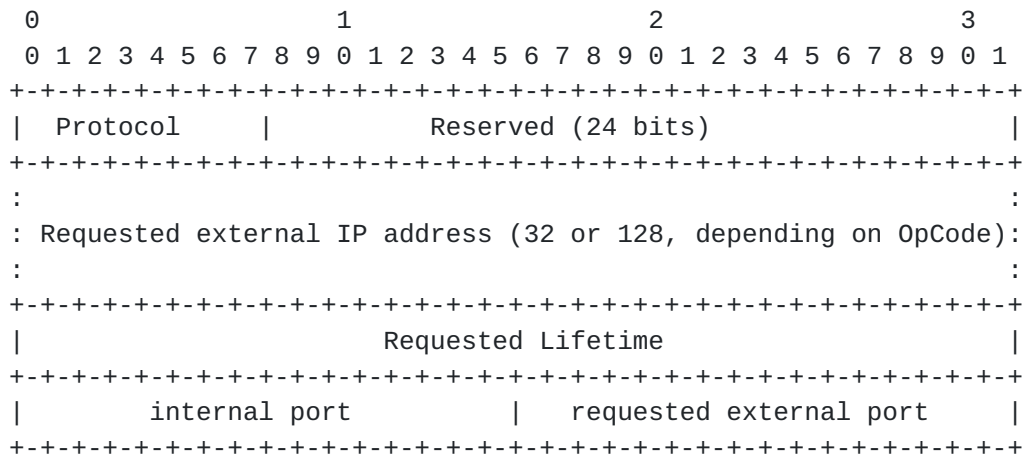


Figure 9: MAP OpCode Request Packet Format

These fields are described below:

Protocol: indicates protocol associated with this OpCode. Values are taken from the IANA protocol registry [[proto_numbers](#)]. For example, this field contains 6 (TCP) if the opcode is intended to create a TCP mapping. The value 0 means "all protocols" (supported by the PCP server), which is useful to create mappings for all protocols with a Basic NAT [[RFC3022](#)] or a firewall, and used to delete mappings for all protocols.

Reserved: 24 reserved bits, MUST be sent as 0 and MUST be ignored when received.

Requested External IP Address: Requested external IP address. This is useful for refreshing a mapping, especially after the PCP server loses state. If the PCP server can fulfill the request, it will do so. If the PCP client does not know the external address, or does not have a preference, it MUST use 0.

Requested lifetime: Requested lifetime of this mapping, in seconds.

Internal port: Internal port for the mapping. The value 0 MUST NOT be sent.

Requested external port: requested external port for the mapping.

This is useful for refreshing a mapping, especially after the PCP server loses state. If the PCP server can fulfill the request, it will do so. If the PCP client does not know the external port, or does not have a preference, it uses 0.

The following diagram shows the response packet format for MAP4 and MAP6 OpCodes:

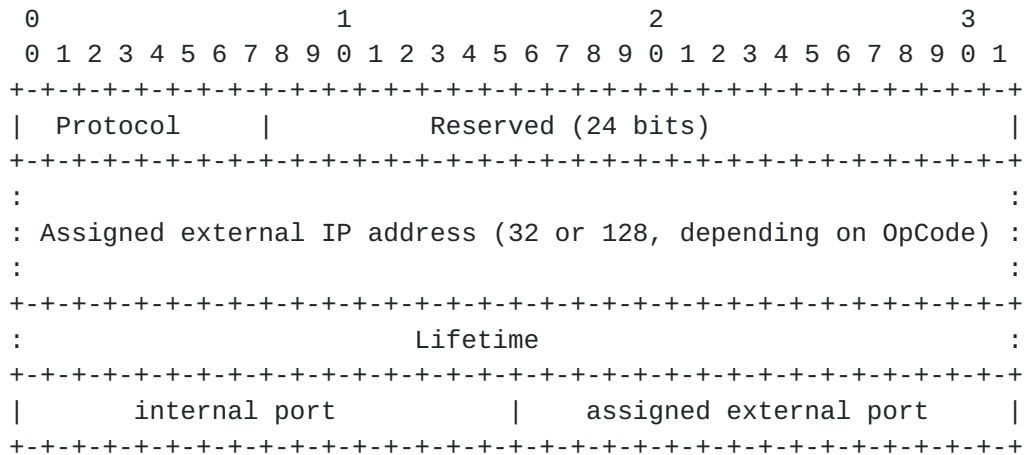


Figure 10: MAP OpCode Response Packet Format

These fields are described below:

Protocol: Copied from the request

Reserved: 24 reserved bits, MUST be sent as 0 and MUST be ignored when received.

Assigned external IP address: On success responses, this is the assigned external IPv4 or IPv6 address for the mapping; IPv4 or IPv6 address is indicated by the OpCode. On error responses, this MUST be 0.

Lifetime: On a success response, this indicates the lifetime for this mapping, in seconds. On an error response, this indicates how long clients should assume they'll get the same error response from the that PCP server if they repeat the same request.

Internal port: Internal port for the mapping, copied from request.

Assigned external port: On success responses, this is the assigned external port for the mapping. IPv4 or IPv6 address is indicated by the OpCode. If the NAT gateway can allocate the requested external port it SHOULD do so. This is beneficial for re-

establishing state lost when a NAT gateway fails or loses its state due to reboot. If the NAT gateway cannot allocate the requested external port but can allocate some other port, it MUST do so and return the allocated port in the response. Cases where a NAT gateway cannot allocate the requested external port include when the requested external port is prohibited by policy, already used by the NAT gateway for one of its own services (e.g., port 80 for the NAT gateway's own configuration pages) already allocated to another explicit mapping (by static manual allocation or by a prior PCP request by a different Internal Host) or the rare case where the requested external port was already allocated to an implicit mapping which cannot be 'promoted' to an explicit mapping for this Internal Host (a different Internal Host already made a prior outbound connection for which the NAT gateway happened to assign the external port requested in this explicit PCP request). On error responses, this MUST be 0.

8.2. OpCode-Specific Result Codes

In addition to the general PCP result codes ([Section 5.4](#)), the following additional result codes may be returned as a result of the four MAP OpCodes received by the PCP server. These errors are considered 'long lifetime' or 'short lifetime', which provides guidance to PCP server developers for the value of the Lifetime field for these errors. It is RECOMMENDED that short lifetime errors use 30 second lifetime and long lifetime errors use 30 minute lifetime.

- 19 SERVER_OVERLOADED, server is processing too many MAP requests from this client or from other clients, and requests this client delay sending other requests. This is a short lifetime error.
- 20 NETWORK_FAILURE, e.g., NAT device has not obtained a DHCP lease so cannot perform a MAP operation. This is a short lifetime error.
- 21 NO_RESOURCES, e.g., NAT device cannot create more mappings at this time. This is a system-wide error, and different from USER_EX_QUOTA. This is a short lifetime error.
- 22 UNSUPP_PROTOCOL, unsupported Protocol. This is a long lifetime error.
- 23 NOT_AUTHORIZED, e.g., PCP server supports mapping, but the feature is disabled for this PCP client, or the PCP client requested a mapping that cannot be fulfilled by the PCP server's security policy. This is a long lifetime error.

- 24 `USER_EX_QUOTA`, mapping would exceed user's port quota. This is a short lifetime error.
- 25 `CANNOT_PROVIDE_EXTERNAL_PORT`, indicates the port is already in use or otherwise unavailable (e.g., special port that cannot be allocated by the server's policy). This error is only returned if the request included the Option `PREFER_FAILURE`. This is a short lifetime error.
- 26 `UNABLE_TO_DELETE_ALL`, indicates the PCP server was not able to delete all mappings. This is a short lifetime error.

Additional result codes may be returned if the `THIRD_PARTY` option is used, see [Section 10](#).

8.3. OpCode-Specific Client: Generating a Request

This section describes the operation of a PCP client when sending requests with OpCodes `MAP4` and `MAP6`.

The request MAY contain values in the `requested-external-ip-address` and `requested-external-port` fields. This allows the PCP client to attempt to rebuild the PCP server's state, so that the PCP client could avoid having to change information maintained at the rendezvous server. Of course, due to other activity on the network (e.g., by other users or network renumbering), the PCP server may not be able to fulfill the request.

An existing mapping can have its lifetime extended by the PCP client. To do this, the PCP client sends a new `MAP` request indicating the internal IP address and port(s).

The PCP client SHOULD renew the mapping before its expiry time, otherwise it will be removed by the PCP server (see [Section 8.6](#)). In order to prevent excessive PCP chatter, it is RECOMMENDED to send a single renewal request packet when a mapping is halfway to expiration time, then, if no positive response is received, another single renewal request 3/4 of the way to expiration time, and then another at 7/8 of the way to expiration time, and so on, subject to the constraint that renewal requests MUST NOT be sent less than four seconds apart (a PCP client MUST NOT send an infinite number of ever-closer-together requests in the last few seconds before a mapping expires).

8.4. OpCode-Specific Server: Processing a Request

This section describes the operation of a PCP server when processing a request with the OpCodes `MAP4` or `MAP6`.

If the server is overloaded by requests (from a particular client or from all clients), it MAY simply discard requests, as the requests will be retried by PCP clients, or MAY generate the `SERVER_OVERLOADED` error response, or both.

If the request contains `internal-port=0`, the server MUST generate a `MALFORMED_REQUEST` error.

If the requested lifetime is 0, it indicates a request to delete the mapping immediately. On a deletion request, the requested external port field is ignored by the server. PCP MAP requests only control mappings created by MAP requests. So, if the PCP client attempts to delete a static mapping (i.e., a mapping created outside of PCP itself), the PCP server deletes all of the PCP-created mappings but MUST respond with `UNABLE_TO_DELETE_ALL` result code, with the other fields encoded as described above. If the PCP client attempts to delete a mapping that does not exist, the success response code is returned. If the PCP client is not authorized to delete this mapping, `NOT_AUTHORIZED` is returned. If the deletion request was properly formatted, a positive response is generated with lifetime of 0 and the server copies the protocol and internal port number from the request into the response; this positive response is generated even if there is no mapping (because the mapping could have been already deleted by a previous PCP transaction).

If the requested lifetime is not zero, it indicates a request to create a mapping or extend the lifetime of an existing mapping.

Processing of the lifetime is described in [Section 8.6](#).

If the PCP-controlled device is stateless (that is, it does not establish any per-flow state, and simply rewrites the address and/or port in a purely algorithmic fashion), the PCP server simply returns an answer indicating the external IP address and port yielded by this stateless algorithmic translation. This allows the PCP client to learn its external IP address and port as seen by remote peers. Examples of stateless translators include stateless NAT64 and 1:1 NAT44, both of which modify addresses but not port numbers.

If an Option with value greater than 128 exists but that option does not make sense (e.g., the `PREFER_FAILURE` option is included in a request with `lifetime=0`), the request is invalid and generates a `MALFORMED_OPTION` error.

By default, a PCP-controlled device MUST NOT create mappings for a protocol not indicated in the request. For example, if the request was for a TCP mapping, a UDP mapping MUST NOT be created.

If the `THIRD_PARTY` option is not present in the request, the source IP address of the PCP packet is used when creating the mapping. If the `THIRD_PARTY` option is present, the PCP server validates the indicated target IP address belongs to the same subscriber. This validation depends on the PCP deployment scenario; see [Section 13.3](#) for the validation procedure. If the internal IP address in the PCP request does not belong to the subscriber, an error response MUST be generated with result code `NOT_AUTHORIZED`.

Mappings typically consume state on the PCP-controlled device, and it is RECOMMENDED that a per-subscriber or per-host limit be enforced by the PCP server to prevent exhausting the mapping state. If this limit is exceeded, the response code `USER_EX_QUOTA` is returned.

If all of the proceeding operations were successful (did not generate an error response), then the requested mappings are created as described in the request and a positive response is built. This positive result contains the same `OpCode` as the request, but with the "R" bit set.

As a side-effect of creating a mapping, ICMP messages associated with the mapping MUST be forwarded (and also translated, if appropriate) for the duration of the mapping's lifetime. This is done to ensure that ICMP messages can still be used by hosts, without application programmers or PCP client implementations needing to signal PCP separately to create ICMP mappings for those flows.

8.5. [OpCode-Specific Client: Processing a Response](#)

This section describes the operation of the PCP client when it receives a PCP response for the `OpCodes` `MAP4` or `MAP6`.

A response is matched with a request by comparing the protocol, internal IP address, and internal port. Other fields are not compared, because the PCP server sets those fields.

If a successful response, the PCP client can use the external IP address and port(s) as desired. Typically the PCP client will communicate the external IP address and port(s) to another host on the Internet using an application-specific rendezvous mechanism such as DNS SRV records.

If the response code is `IMPLICIT_MAPPING_EXISTS`, it indicates the PCP client is attempting to use `MAP` when an implicit dynamic connection already exists for the same internal host and internal port. This can occur with certain types of NATs. When this is received, if the PCP client still wants to establish a mapping, the PCP client MUST choose a different internal port and send a new PCP request

specifying that port.

On an error response, clients SHOULD NOT repeat the same request to the same PCP server within the lifetime returned in the response.

8.6. Mapping Lifetime and Deletion

The PCP client requests a certain lifetime, and the PCP server responds with the assigned lifetime. The PCP server MAY grant a lifetime smaller or larger than the requested lifetime. The PCP server SHOULD be configurable for permitted minimum and maximum lifetime, and the RECOMMENDED values are 120 seconds for the minimum value and 24 hours for the maximum. It is NOT RECOMMENDED that the server allow lifetimes exceeding 24 hours, because they will consume ports even if the internal host is no longer interested in receiving the traffic or no longer connected to the network.

Once a PCP server has responded positively to a mapping request for a certain lifetime, the port forwarding is active for the duration of the lifetime unless the lifetime is reduced by the PCP client (to a shorter lifetime or to zero) or until the PCP server loses its state (e.g., crashes). However, if the PCP lifetime has reached zero yet there is still active inside-to-outside traffic, the PCP server MAY, if it desires, keep the mapping active until the inside-to-outside traffic has stopped.

An application that forgets its PCP-assigned mappings (e.g., the application or OS crashes) will request new PCP mappings. This will consume port mappings. The application will also likely initiate new implicit dynamic mappings (e.g., TCP connections) without using PCP, which will also consume port mappings. If there is a port mapping quota for the internal host, frequent restarts such as this may exhaust the quota. PCP provides no explicit protection against such port consumption. In such environments, it is RECOMMENDED that applications use shorter PCP lifetimes to reduce the impact of consuming the user's port quota. An operating system or framework that issues a mapping request to "delete all" (protocol=0, port=0, lifetime=0) on reboot protects itself against this resource exhaustion by voluntarily relinquishing all of its old mappings before beginning to request new ones. The PCP server MAY choose to allocate the same (recently relinquished) mappings when mappings are re-requested by the booting OS. Some port mapping APIs (such as the "DNSServiceNATPortMappingCreate" API provided by Apple's Bonjour on Mac OS X, iOS, Windows, Linux, etc.) automatically monitor for process exit (including application crashes) and automatically send port mapping deletion requests if the process that requested them goes away without explicitly relinquishing them.

In order to reduce unwanted traffic and data corruption, a port that was mapped using the MAP OpCode SHOULD NOT be assigned to another internal target, or another subscriber, for 120 seconds (MSL, [RFC0793]). However, the PCP server MUST allow the same internal target to re-acquire the same port during that same interval.

When a PCP client first acquires a new IP address, it may want to remove mappings that may have been instantiated for a previous host. To do this, the PCP client sends a MAP request with protocol, external port, internal port, and lifetime set to 0.

8.7. Subscriber Renumbering

The customer premises router might obtain a new IPv4 address or new IPv6 prefix. This can occur because of a variety of reasons including a reboot, power outage, DHCP lease expiry, or other action by the ISP. If this occurs, traffic forwarded to the subscriber might be delivered to another customer who now has that address. This affects both implicit dynamic mappings and explicit dynamic mappings. However, this same problem occurs today when a subscriber's IP address is re-assigned, without PCP and without an ISP-operated CGN. The solution is the same as today: the problems associated with subscriber renumbering are caused by subscriber renumbering and are eliminated if subscriber renumbering is avoided. PCP defined in this document does not provide machinery to reduce the subscriber renumbering problem.

When a new Internal Address is assigned to a host embedding a PCP client, the NAT (or firewall) controlled by the PCP server will continue to send traffic to the old IP address. Assuming the PCP client wants to continue receiving traffic, it needs to install new mappings for its new IP address. The requested external port field will not be fulfilled by the PCP server, in all likelihood, because it is still being forwarded to the old IP address. Thus, a mapping is likely to be assigned a new external port number and/or public IP address. Note that this scenario is not expected to happen routinely on a regular basis for most hosts, since most hosts renew their DHCP leases before they expire (or re-request the same address after reboot) and most DHCP servers honor such requests and grant the host the same address it was previously using before the reboot.

8.8. PCP Options for MAP OpCodes

8.8.1. REMOTE_PEER_FILTER

This Option indicates packet filtering is desired. The remote peer port and remote peer IP Address indicate the permitted remote peer's source IP address and port for packets from the Internet. The remote

length: 2 if used with MAP4, 5 if used with MAP6

may appear in: requests

maximum occurrences: as many as fit within maximum PCP message size

Because of interactions with dynamic ports this Option MUST only be used by a client that is operating a server (that is, using the MAP OpCode), as this ensures that no other application will be assigned the same ephemeral port for its outgoing connection. Other use by a PCP client is NOT RECOMMENDED and will cause some UNSAF NAT traversal mechanisms [[RFC3424](#)] to fail where they would have otherwise succeeded, breaking other applications running on this same host.

The prefix-length indicates how many bits of the IPv6 address or IPv4 address are used for the filter. For MAP4, a prefix-length of 32 indicates the entire IPv4 address is used. For MAP6, a prefix-length of 128 indicates the entire IPv6 address is used. For MAP4 the minimum prefix-length value is 0 and the maximum value is 32. For MAP6 the minimum prefix-length value is 0 and the maximum value is 128. Values outside those range cause an MALFORMED_OPTION response code.

If multiple occurrences of REMOTE_PEER_FILTER exist in the same MAP request, they are processed in the same order received, and they MUST all be successfully processed or return an error (e.g., MALFORMED_OPTION if one of the options was malformed). As with other PCP errors, returning an error causes no state to be changed in the PCP server or in the PCP-controlled device. If an existing mapping exists (with or without a filter) and the server receives a MAP request with REMOTE_PEER_FILTER, the filters indicated in the new request are added to any existing filters. If a MAP request has a lifetime of 0 and contains the REMOTE_PEER_FILTER option, the error MALFORMED_OPTION is returned.

To remove all existing filters, the prefix-length 0 is used. There is no mechanism to remove a specific filter.

To change an existing filter, the PCP client sends a MAP request containing two REMOTE_PEER_FILTER options, the first option containing a prefix-length of 0 (to delete all existing filters) and the second containing the new remote peer's IP address and port. Other REMOTE_PEER_FILTER options in that PCP request, if any, add more allowed remote hosts.

The PCP server or the PCP-controlled device is expected to have a limit on the number of remote peers it can support. This limit might be as small as one. If a MAP request would exceed this limit, the entire MAP request is rejected with the result code

EXCESSIVE_REMOTE_PEERS, and the state on the PCP server is unchanged.

If this option appears in a request, the following addition result code could be returned:

- 27 EXCESSIVE_REMOTE_PEERS, indicates the PCP server was not able to create the filters in this request. This result code MUST only be returned if the MAP request contained the REMOTE_PEER_FILTER Option. This is a long lifetime error.

8.8.2. PREFER_FAILURE

This option indicates that if the PCP server is unable to allocate the requested port, then instead of returning an available port that it *can* allocate, the PCP server should instead allocate no port and return result code CANNOT_PROVIDE_EXTERNAL_PORT.

This option is intended solely for use by UPnP IGD interworking [[I-D.bpw-pcp-upnp-igd-interworking](#)], where the semantics of IGD version 1 do not provide any way to indicate to an IGD client that any port is available other than the one it requested. A PCP server MAY support this option, if its designers wish to support downstream devices that perform IGD interworking. PCP servers MAY choose to rate-limit their handling of PREFER_FAILURE requests, to protect themselves from a rapid flurry of 65535 consecutive PREFER_FAILURE requests from clients probing to discover which external ports are available. PCP servers that are not intended to support downstream devices that perform IGD interworking are not required to support this option. PCP clients other than IGD interworking clients SHOULD NOT use this option because it results in inefficient operation, and they cannot safely assume that all PCP servers will implement it. The option is provided only because the semantics of IGD version 1 offer no viable alternative way to implement an IGD interworking function. It is anticipated that this option will be deprecated in the future as more clients adopt PCP natively and the need for IGD interworking declines.

This Option:

name: PREFER_FAILURE

number: 3

is valid for OpCodes: MAP4, MAP6

is included in responses: MUST

length: 0

may appear in: requests

maximum occurrences: no

8.8.3. THIRD_PARTY

The THIRD_PARTY option is used by both the MAP OpCode and the PEER OpCode, and defined in [Section 10](#).

8.9. PCP Mapping State Maintenance

If an event occurs that causes the PCP server to lose state (such as a crash or power outage), the mappings created by PCP are lost. Such loss of state is rare in a service provider environment (due to redundant power, disk drives for storage, etc.). But such loss of state is more common in a residential NAT device which does not write information to its non-volatile memory.

The Epoch allows a client to deduce when a PCP server may have lost its state. If this occurs, the PCP client can attempt to recreate the mappings following the procedures described in this section.

8.9.1. Recreating Mappings

The PCP server SHOULD store mappings in persistent storage so when it is powered off or rebooted, it remembers the port mapping state of the network. Due to the physical architecture of some PCP servers, this is not always achievable (e.g., some non-volatile memory can withstand only a certain number of writes, so writing PCP mappings to such memory is generally avoided).

However, maintaining this state is not essential for correct operation. When the PCP server loses state and begins processing new PCP messages, its Epoch is reset to zero (per the procedure of [Section 6.5](#)).

A mapping renewal packet is formatted identically to an original mapping request; from the point of view of the client it is a renewal of an existing mapping, but from the point of view of the PCP server it appears as a new mapping request.

As the result of receiving a packet where the Epoch field indicates that a reboot or similar loss of state has occurred, the client renews its port mappings.

The discussion in this section focuses on recreating inbound port

mappings after loss of PCP server state, because that is the more serious problem. Losing port mappings for outgoing connections destroys those currently active connections, but does not prevent clients from establishing new outgoing connections. In contrast, losing inbound port mappings not only destroys all existing inbound connections, but also prevents the reception of any new inbound connections until the port mapping is recreated. Accordingly, we consider recovery of inbound port mappings the more important priority. However, clients that want outgoing connections to survive a NAT gateway reboot can also achieve that using PCP. After initiating an outbound TCP connection (which will cause the NAT gateway to establish an implicit port mapping) the client should send the NAT gateway a port mapping request for the source port of its TCP connection, which will cause the NAT gateway to send a response giving the external port it allocated for that mapping. The client can then store this information, and use it later to recreate the mapping if it determines that the NAT gateway has lost its mapping state.

8.9.2. Maintaining Mappings

A PCP client can refresh a mapping by sending a new PCP request containing information from the earlier PCP response. The PCP server will respond indicating the new lifetime. It is possible, due to failure of the PCP server, that the public IP address and/or public port, or the PCP server itself, has changed (due to a new route to a different PCP server). To detect such events more quickly, the PCP client may find it beneficial to use shorter lifetimes (so that it communicates with the PCP server more often). If the PCP client has several mappings, the Epoch value only needs to be retrieved for one of them to verify the PCP server has not lost port forwarding state.

If the client wishes to check the PCP server's Epoch, it sends a PCP request for any one of the client's mappings. This will return the current Epoch value. In that request the PCP client could extend the mapping lifetime (by asking for more time) or maintain the current lifetime (by asking for the same number of seconds that it knows are remaining of the lifetime).

9. PEER OpCodes

This section defines two OpCodes for controlling dynamic connections. They are:

PEER4=2: Set or query lifetime for flow from IPv4 address to a remote peer's IPv4 address.

PEER6=3: Set or query lifetime for flow from IPv6 address to a remote peer's IPv6 address.

The operation of these OpCodes is described in this section.

9.1. OpCode Packet Formats

The PEER OpCodes provide a single function: the ability for the PCP client to query and (possibly) extend the lifetime of an existing mapping.

The two PEER OpCodes (PEER4 and PEER6) share a similar packet layout for both requests and responses. Because of this similarity, they are shown together.

The following diagram shows the request packet format for PEER4 and PEER6. This packet format is aligned with the response packet format:

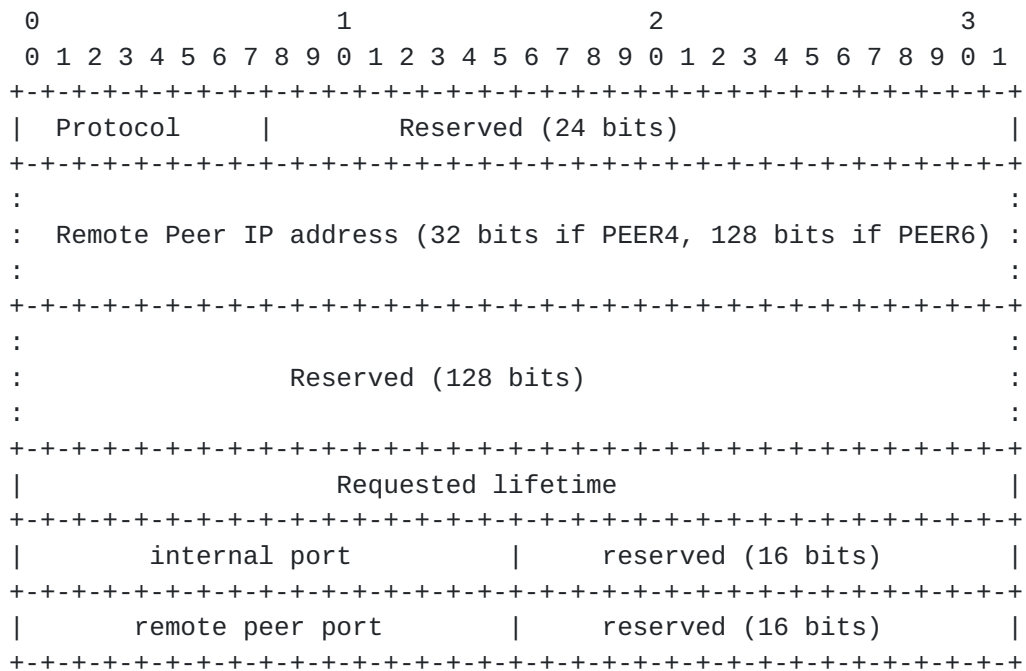


Figure 12: PEER OpCode Request Packet Format

These fields are described below:

Protocol: indicates protocol associated with this OpCode. Values are taken from the IANA protocol registry [[proto_numbers](#)]. For example, this field contains 6 (TCP) if the OpCode is describing a TCP peer.

Reserved: 24 reserved bits, MUST be 0 on transmission and MUST be ignored on reception.

Remote Peer IP Address: Remote peer's IP address, from the perspective of the PCP client.

Reserved: 128 reserved bits, MUST be 0 on transmission and MUST be ignored on reception.

Requested lifetime: Requested lifetime of this mapping, in seconds. Unlike the MAP OpCode, there is no special meaning of 0.

internal port: Internal port for the of the 5-tuple.

Reserved: 16 reserved bits, MUST be 0 on transmission and MUST be ignored on reception.

Remote Peer Port: Remote peer's port of the 5-tuple.

Reserved: 16 reserved bits, MUST be 0 on transmission and MUST be ignored on reception.

The following diagram shows the response packet format for PEER4 and PEER6:

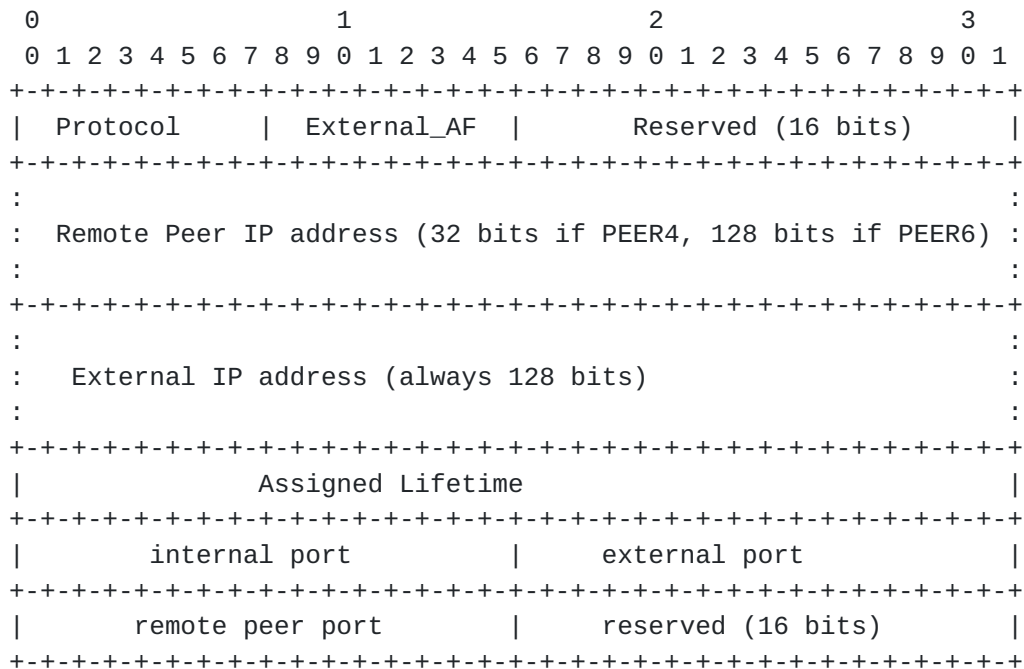


Figure 13: PEER OpCode Response Packet Format

Protocol: Copied from the request.

External_AF For success responses, this contains the address family of the external IP address associated with this peer connection. Values are from IANA's address family numbers (IPv4 is 1, IPv6 is 2). For error responses, the value MUST be 0.

Reserved: 16 reserved bits, MUST be 0 on transmission, MUST be ignored on reception.

remote Peer IP address Copied from the request.

External IP Address For success responses, this contains the external IP address, assigned by the NAT (or firewall) to this mapping. If firewall, this will match the internal IP address. This field is always 128 bits long. If External_AF indicates IPv4, the IPv4 address is encoded in the first 32 bits of the External IP Address field and the remaining 96 bits are zero. For error responses, this MUST be 0.

Assigned Lifetime: For success responses, this is the assigned lifetime, in seconds. For error responses, this is 0.s

internal port: copied from request.

external port: For success responses, this is the external port number, assigned by the NAT (or firewall) to this mapping. If firewall or 1:1 NAT, this will match the internal port. For error responses, this MUST be 0.

remote peer port: Copied from request.

Reserved: 16 reserved bits, MUST be 0 on transmission, MUST be ignored on reception.

9.2. OpCode-Specific Result Codes

In addition to the general PCP result codes ([Section 5.4](#)) the following additional result codes may be returned as a result of the two PEER OpCodes received by the PCP server.

50 NONEXIST_PEER, the connection to that peer does not exist in the mapping table.

Additional result codes may be returned if the THIRD_PARTY option is used, see [Section 10](#).

9.3. OpCode-Specific Client: Generating a Request

This section describes the operation of a client when generating the OpCodes PEER4 or PEER6.

The PEER4 or PEER6 OpCodes MUST NOT be sent until establishing bi-directional communication with the remote peer. For TCP, this means completing the TCP 3-way handshake. This is because the PCP-controlled device may not be able to extend the lifetime of a mapping until after bi-directional communications has been established.

The PEER4 and PEER6 OpCodes contain a description of the socket, from the perspective of the PCP client. This is important when the PCP-controlled device is performing address family translation (NAT46 or NAT64), because the destination address from the perspective of the PCP client is different from the destination address on the other side of the address family translation device.

9.4. OpCode-Specific Server: Processing a Request

This section describes the operation of a server when receiving a request with the OpCodes PEER4 or PEER6.

On receiving the PEER4 or PEER6 OpCode, the PCP server examines the mapping table. If a mapping does not exist, the NONEXIST_PEER error is returned.

If the PCP-controlled device can extend the lifetime of a mapping, the PCP server uses the smaller of its configured maximum lifetime value and the requested lifetime from the PEER request, and sets the lifetime to that value.

The PEER4 or PEER6 OpCodes MUST NOT reduce the lifetime of an existing mapping. If the mapping is terminated by the TCP client or server (e.g., TCP FIN or TCP RST), the mapping will eventually be destroyed normally; the earlier use of PEER does not extend the lifetime in that case.

If all of the proceeding operations were successful (did not generate an error response), then a SUCCESS response is generated, with the assigned-lifetime containing the lifetime of the mapping.

Note: it is implementation-specific if the PCP-controlled device destroys the mapping when the lifetime expires, or if inside->outside traffic keeps the mapping alive.

9.5. OpCode-Specific Client: Processing a Response

This section describes the operation of a client when processing a response with the OpCodes PEER4 or PEER6.

A response is matched with a request by comparing the protocol, external AF, internal IP address, internal port, remote peer address and remote peer port. Other fields are not compared, because the PCP server changes those fields to provide information about the mapping created by the OpCode.

If the error response NONEXIST_PEER, this could have occurred if the PCP client sent its PEER request before the PCP-controlled device had installed the mapping, or because the mapping has been destroyed (e.g., due to a TCP FIN). If the PCP client believes the mapping should exist, the PCP client SHOULD retry the request after a brief delay (e.g., 5 seconds).

Other error responses SHOULD NOT be retried.

If a successful response, the PCP client uses the assigned lifetime value to reduce its frequency of application keepalives for that particular NAT mapping. Of course, there may be other reasons, specific to the application, to use more frequent application keepalives. For example, the PCP assigned-lifetime could be one hour but the application may want to ensure the server is still accessible (e.g., has not crashed) more frequently than once an hour.

If the PCP client wishes to keep this mapping alive beyond the indicated lifetime, it **SHOULD** issue a new PCP request prior to the expiration. That is, inside->outside traffic is not sufficient to ensure the mapping will continue to exist. It is **RECOMMENDED** to send a single renewal request packet when a mapping is halfway to expiration time, then, if no positive response is received, another single renewal request 3/4 of the way to expiration time, and then another at 7/8 of the way to expiration time, and so on, subject to the constraint that renewal requests **MUST NOT** be sent less than four seconds apart (a PCP client **MUST NOT** send an infinite number of ever-closer-together requests in the last few seconds before a mapping expires).

9.6. PCP Options for PEER OpCodes

9.6.1. THIRD_PARTY

The **THIRD_PARTY** option is used by both the MAP OpCode and the PEER OpCode, and defined in [Section 10](#).

10. THIRD_PARTY Option for MAP and PEER OpCodes

This Option is used when a PCP client wants to control a mapping to an internal target host other than itself. This is used with both MAP and PEER OpCodes.

A PCP server will only process this option if sent by an authorized PCP client, otherwise will return an error. Determining which PCP clients are authorized to use the **THIRD_PARTY** option depends on the deployment scenario. For Dual-Stack Lite deployments, the PCP server only supports this option if the source IPv4 address is the B4's source IP address. For other scenarios, the subscriber has only one IPv4 address and this Option serves no purpose (and will only generate error messages from the server). If a subscriber has more than one IPv4 address, the ISP **MUST** determine its own policy for how to identify the trusted device within the subscriber's home. This might be, for example, the lowest- or highest-numbered host address for that user's IPv4 prefix.


```

      0               1               2               3
    0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+
:
: Target Internal IP address (32 bits of 128 bits, depending      :
:                               on Option length)                  :
+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+--+

```

The fields are described below:

Mapping Internal IP Address: Internal IP address of the mapping. If the length of this Option is 1, this is a 32-bit IPv4 address. If the length of this Option is 4, this is a 128-bit IPv6 address. This can contain the special value "0" (all zeros), which indicates "all addresses associated with this subscriber" which is used to delete all pre-existing mappings with the MAP Opcode.

This Option:

name: THIRD_PARTY

number: 4

purpose: Indicate the MAP or PEER request is for a host other than the host sending the PCP option.

is valid for OpCodes: MAP4, MAP6, PEER4, PEER6

length: 1 if OpCode is MAP4 or PEER4, 4 if OpCode is MAP6 or PEER6

may appear in: request. May appear in response only if it appeared in the associated request.

maximum occurrences: 1

The following additional result codes may be returned as a result of using this Option.

51 UNAUTH_TARGET_ADDRESS, indicting the target IP address specified is not permitted (e.g., the address does not belong to this subscriber, or is otherwise prohibited.). If this is a MAP request, this is a long-term error.

52 UNAUTH_SOURCE_ADDRESS, indicates the source address of this PCP message is not authorized by the PCP server to use the THIRD_PARTY option.

A PCP server is configured to permit or to restrict the use of the

THIRD_PARTY option. If this option is permitted, any host on a subscriber's network can create, modify, or destroy mappings for another host on the network, which is generally undesirable. If third party mappings are restricted, only a certain host on the subscriber's network can perform these operations. If a PCP server is configured to restrict third party mappings, and receives a PCP MCP request with a Target Address that does not match the source IP address of that request, it MUST generate a UNAUTH_TARGET_ADDRESS response.

It is RECOMMENDED that PCP servers embedded into customer premise equipment be configured to restrict third party mappings. With this configuration, if a user wants to create a third party mapping, the user needs to interact out-of-band with their customer premise router (e.g., using its HTTP administrative interface).

It is RECOMMENDED that PCP servers embedded into service provider NAT and firewall devices be configured to permit the THIRD_PARTY option. With this configuration, if a user wants to create an explicit dynamic mapping or query an implicit dynamic mapping for another host within their network, the user needs to interact out-of-band with their customer premise router (e.g., using its HTTP administrative interface). This is because the service provider's PCP server only allows a PEER or MAP request containing the THIRD_PARTY option if it has the IP address of the subscriber's customer premise router. To do this, the PCP server needs certain knowledge about the network's subscribers. It needs to determine the IP address of the subscriber's customer premise router and to determine the IP subnet assigned to the subscriber. This knowledge might be dynamic (e.g., database query into the service provider's user database for every incoming PCP request), might be a table (e.g., subscribers with a certain IPv4 network prefix all have an IPv4 /24, other IPv4 prefixes have an IPv4 /32, certain IPv6 prefixes have an IPv6 /32, and so on), or might be very static (e.g., all subscribers have one IPv4 address). In many common deployments, there is only one IPv4 address assigned to a subscriber, and thus the Target Address will always match the source address of the PCP message. If there are multiple IPv4 or multiple IPv6 addresses assigned to a subscriber, the PCP server allows the highest-numbered address to use the THIRD_PARTY option. Thus, on a network supporting PCP with multiple addresses assigned to a subscriber, the highest-numbered host SHOULD be the subscriber's customer premise router. Upon receiving a MAP or PEER request where the Target Address does not match the source IP address of the request, the PCP server determines if the source IP address of the request is the subscriber's highest numbered address, following the procedure above. If not, the PCP server MUST generate an UNAUTH_SOURCE_ADDRESS error. Then the PCP server determines if the Target Address belongs to the same subscriber as the source IP

address of the PCP packet, using the procedure described above. If not, the PCP server MUST generate an UNAUTH_TARGET_ADDRESS error.

If authorized to do so, a PCP client can delete all the PCP-created dynamic explicit mappings (i.e., those created by PCP) for all hosts belonging to the same subscriber. This is done by sending a PCP MAP request including the THIRD_PARTY option with its Target Address field set to 0.

11. Deployment Considerations

11.1. Maintaining Same External IP Address

It is REQUIRED that the PCP-controlled device assign the same external IP address PCP-created explicit dynamic mappings and to implicit dynamic mappings. It is RECOMMENDED that static mappings (e.g., those created by a command language interface on the PCP server or PCP-controlled device) also be assigned to the same IP address.

Once all internal hosts belonging to a given subscriber have no implicit dynamic mappings and have no explicit dynamic mappings in the PCP-controlled device, a subsequent PCP request for that internal host MAY be assigned to a different external IP address. Generally, this re-assignment would occur when a CGN device is load balancing newly-seen hosts to its public IPv4 address pool.

11.2. Ingress Filtering

To prevent spoofing of PCP requests, ingress filtering [[RFC2827](#)] MUST be performed by devices between the PCP clients and PCP server. For example, with a PCP server integrated into a customer premise router, the Ethernet switch needs to perform ingress filtering. As another example, with a PCP server deployed by a service provider, the service provider's aggregation router (the first device connecting to subscribers) needs to do ingress filtering.

11.3. Per-Subscriber Port Forwarding Quota

On PCP-controlled devices that create state when a mapping is created (e.g., NAT), the PCP server SHOULD maintain a per-subscriber mapping quota for PCP-created mappings. It is implementation-specific if the PCP server has a separate or combined quota for both implicit dynamic mappings (e.g., TCP SYNs) and explicit dynamic mappings (PCP).

12. Deployment Scenarios

12.1. Dual Stack-Lite

The interesting components in a Dual-Stack Lite deployment are the B4 element (which is the customer premises router) and the AFTR (Address Family Transition Router) element. The AFTR element terminates the IPv6-over-IPv4 tunnel and also implements the Carrier-Grade NAT44 function. The B4 element does not need to perform a NAT function (and usually does not perform a NAT function), but it does operate its own DHCP server and is the local network's default router.

12.1.1. Overview

Various PCP deployment scenarios can be considered to control the PCP server embedded in the AFTR element:

1. UPnP IGD and NAT-PMP [[I-D.cheshire-nat-pmp](#)] are used in the LAN: an interworking function is required to be embedded in the B4 element to ensure interworking between the protocol used in the LAN and PCP. UPnP IGD-PCP Interworking Function is described in [[I-D.bpw-pcp-upnp-igd-interworking](#)].
2. Hosts behind the B4 element will either include a PCP client or UPnP IGD client, or both.
 - A. if a UPnP IGD client, the B4 element will need to include an interworking function from UPnP IGD to PCP.
 - B. if a PCP client, the PCP client will communicate directly with the PCP server.
3. The B4 element includes a PCP client which is invoked by an HTTP-based configuration (as is common today). The internal IP address field in the PCP payload would be the internal host used in the port forwarding configuration.

In Dual Stack-Lite, the B4 element encapsulates its PCP messages into the IPv6 tunnel towards the AFTR element. It is expected the B4 element will also perform as a proxy from PCP to PCP [[I-D.bpw-pcp-proxy](#)], and may also proxy from other protocols to PCP (e.g., [[I-D.bpw-pcp-upnp-igd-interworking](#)]). When proxying for other hosts, the B4 element will necessarily use the THIRD_PARTY option with the MAP and PEER OpCodes.

12.2. NAT64

Hosts behind a NAT64 device can make use of PCP in order to perform port reservation (to get a publicly routable IPv4 port).

12.3. NAT44 and NAT444

Residential subscribers in NAT44 (and NAT444) deployments are usually given one IPv4 address, but may also be given several IPv4 addresses. These addresses are not routable on the IPv4 Internet, but are routable between the subscriber's home and the ISP's CGN. To accommodate multiple hosts within a home, especially when provided insufficient IPv4 addresses for the number of devices in the home, subscribers operate a NAPT device. When this occurs in conjunction with an upstream NAT44, this is nicknamed "NAT444".

12.4. IPv6 Simple Firewall

Many IPv6 deployments will include a simple firewall [[RFC6092](#)], which permits outgoing packets to initiate bi-directional communication but blocks unsolicited incoming packets, which is similar to PCP's security model that allows a host to create a mapping to itself. In many situations, especially residential networks that lack an IT staff, the security provided by an IPv6 simple firewall and the security provided by PCP are compatible. In such situations, the IPv6 simple firewall and the IPv6 host can use the MAP6 OpCode to allow unsolicited incoming packets, so the host can operate a server.

13. Security Considerations

The PCP client's source port SHOULD be randomly generated [[RFC6056](#)]. The PCP server MUST only listen for requests from its internal interfaces, and MUST NOT listen for requests on its Internet-facing interfaces.

This document defines Port Control Protocol and two types of OpCodes, PEER and MAP. The PEER OpCode allows querying and extending (if permitted) the lifetime of an existing implicit dynamic mapping, so a host can reduce its keepalive messages. The MAP OpCode allows creating a mapping so a host can receive incoming unsolicited connections from the Internet in order to run a server.

The PEER OpCode does not introduce any new security considerations.

On today's Internet, ISPs do not typically filter incoming traffic for their subscribers. However, when an ISP introduces stateful address sharing with a NAPT device, such filtering will occur as a

side effect. Filtering will also occur with IPv6 CPE [[I-D.ietf-v6ops-cpe-simple-security](#)]. The MAP OpCode allows a PCP client to create a mapping so that a host can receive inbound traffic and operate a server. Security considerations for the MAP OpCode are described in the following sections.

[13.1.](#) Denial of Service

Because the state created in a NAT or firewall, a per-subscriber quota will likely exist for both implicit dynamic mappings (e.g., outgoing TCP connections) and explicit dynamic mappings (PCP). A subscriber might make an excessive number of implicit or explicit dynamic mappings, consuming an inordinate number of ports, causing a denial of service to other subscribers. Thus, [Section 11.3](#) recommends that subscribers be limited to a reasonable number of explicit dynamic mappings.

[13.2.](#) Ingress Filtering

It is important to prevent a subscriber from creating a mapping for another subscriber, because this allows incoming packets from the Internet and consumes the other user's mapping quota. Both implicit dynamic mappings (e.g., outgoing TCP connections) and explicit dynamic mappings (PCP) need ingress filtering. Thus, PCP does not create a new requirement for ingress filtering.

[13.3.](#) Validating Target Address

The THIRD_PARTY Option contains a Target Address field, which allows a PCP client to create an explicit dynamic mapping for another host. Hosts within a subscriber's network cannot create, modify, or delete mappings of other hosts, except by using the administrative interface of the customer premise router (e.g., HTTP interface), as described in [Section 10](#).

[14.](#) IANA Considerations

IANA is requested to perform the following actions:

[14.1.](#) Port Number

IANA has assigned UDP port 44323 for PCP.

[14.2.](#) OpCodes

IANA shall create a new protocol registry for PCP OpCodes, initially populated with the values in [Section 8](#) and [Section 9](#).

New OpCodes in the range 1-95 can be created via Standards Action [[RFC5226](#)], and the range 96-128 is for Private Use [[RFC5226](#)].

[14.3.](#) Result Codes

IANA shall create a new registry for PCP result codes, numbered 0-255, initially populated with the result codes from [Section 5.4](#), [Section 8.2](#), [Section 8.8.1](#), [Section 9.2](#), and [Section 10](#).

Additional Result Codes can be defined via Specification Required [[RFC5226](#)].

[14.4.](#) Options

IANA shall create a new registry for PCP Options, numbered 0-255 with an associated mnemonic. The values 0-128 are mandatory-to-process, and 128-255 are optional-to-process. The initial registry contains the options described in [Section 8.8](#) and [Section 10](#), and the option values 0 and 255 are reserved.

New PCP option codes in the range 0-63 and 128-192 can be created via Standards Action [[RFC5226](#)], and the range 64-127 and 192-255 is for Private Use [[RFC5226](#)].

[15.](#) Acknowledgments

Thanks to Alain Durand, Christian Jacquenet, Jacni Qin, and Simon Perreault for their comments and review. Thanks to Simon Perreault for highlighting the interaction of dynamic connections with PCP-created mappings.

[16.](#) References

[16.1.](#) Normative References

[I-D.ietf-behave-v6v4-xlate]

Li, X., Bao, C., and F. Baker, "IP/ICMP Translation Algorithm", [draft-ietf-behave-v6v4-xlate-23](#) (work in progress), September 2010.

[I-D.ietf-behave-v6v4-xlate-stateful]

Bagnulo, M., Matthews, P., and I. Beijnum, "Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers", [draft-ietf-behave-v6v4-xlate-stateful-12](#) (work in progress), July 2010.

[I-D.ietf-softwire-dual-stack-lite]

Durand, A., Droms, R., Woodyatt, J., and Y. Lee, "Dual-Stack Lite Broadband Deployments Following IPv4 Exhaustion", [draft-ietf-softwire-dual-stack-lite-06](#) (work in progress), August 2010.

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

[RFC2827] Ferguson, P. and D. Senie, "Network Ingress Filtering: Defeating Denial of Service Attacks which employ IP Source Address Spoofing", [BCP 38](#), [RFC 2827](#), May 2000.

[RFC4193] Hinden, R. and B. Haberman, "Unique Local IPv6 Unicast Addresses", [RFC 4193](#), October 2005.

[RFC5226] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", [BCP 26](#), [RFC 5226](#), May 2008.

[RFC6056] Larsen, M. and F. Gont, "Recommendations for Transport-Protocol Port Randomization", [BCP 156](#), [RFC 6056](#), January 2011.

[proto_numbers]

IANA, "Protocol Numbers", 2010, <<http://www.iana.org/assignments/protocol-numbers/protocol-numbers.xml>>.

16.2. Informative References

[I-D.arkko-dual-stack-extra-lite]

Arkko, J., Eggert, L., and M. Townsley, "Scalable Operation of Address Translators with Per-Interface Bindings", [draft-arkko-dual-stack-extra-lite-05](#) (work in progress), February 2011.

[I-D.bpw-pcp-proxy]

Boucadair, M., Penno, R., Wing, D., and F. Dupont, "Port Control Protocol (PCP) Proxy Function", [draft-bpw-pcp-proxy-00](#) (work in progress), February 2011.

[I-D.bpw-pcp-upnp-igd-interworking]

Boucadair, M., Penno, R., Wing, D., and F. Dupont, "Universal Plug and Play (UPnP) Internet Gateway Device (IGD)-Port Control Protocol (PCP) Interworking Function", [draft-bpw-pcp-upnp-igd-interworking-02](#) (work in progress), February 2011.

- [I-D.cheshire-nat-pmp]
Cheshire, S., "NAT Port Mapping Protocol (NAT-PMP)",
[draft-cheshire-nat-pmp-03](#) (work in progress), April 2008.
- [I-D.ietf-behave-lsn-requirements]
Yamagata, I., Miyakawa, S., Nakagawa, A., and H. Ashida,
"Common requirements for IP address sharing schemes",
[draft-ietf-behave-lsn-requirements-00](#) (work in progress),
October 2010.
- [I-D.ietf-v6ops-cpe-simple-security]
Woodyatt, J., "Recommended Simple Security Capabilities in
Customer Premises Equipment for Providing Residential IPv6
Internet Service", [draft-ietf-v6ops-cpe-simple-security-16](#)
(work in progress), October 2010.
- [I-D.miles-behave-l2nat]
Miles, D. and M. Townsley, "Layer2-Aware NAT",
[draft-miles-behave-l2nat-00](#) (work in progress),
March 2009.
- [IGD] UPNP Gateway Committee, "WANIPConnection:1",
November 2001, <[http://upnp.org/specs/gw/
UPnP-gw-WANIPConnection-v1-Service.pdf](http://upnp.org/specs/gw/UPnP-gw-WANIPConnection-v1-Service.pdf)>.
- [RFC0793] Postel, J., "Transmission Control Protocol", STD 7,
[RFC 793](#), September 1981.
- [RFC1918] Rekhter, Y., Moskowitz, R., Karrenberg, D., Groot, G., and
E. Lear, "Address Allocation for Private Internets",
[BCP 5](#), [RFC 1918](#), February 1996.
- [RFC2474] Nichols, K., Blake, S., Baker, F., and D. Black,
"Definition of the Differentiated Services Field (DS
Field) in the IPv4 and IPv6 Headers", [RFC 2474](#),
December 1998.
- [RFC3022] Srisuresh, P. and K. Egevang, "Traditional IP Network
Address Translator (Traditional NAT)", [RFC 3022](#),
January 2001.
- [RFC3424] Daigle, L. and IAB, "IAB Considerations for UNilateral
Self-Address Fixing (UNSAF) Across Network Address
Translation", [RFC 3424](#), November 2002.
- [RFC3581] Rosenberg, J. and H. Schulzrinne, "An Extension to the
Session Initiation Protocol (SIP) for Symmetric Response
Routing", [RFC 3581](#), August 2003.

- [RFC4941] Narten, T., Draves, R., and S. Krishnan, "Privacy Extensions for Stateless Address Autoconfiguration in IPv6", [RFC 4941](#), September 2007.
- [RFC4961] Wing, D., "Symmetric RTP / RTP Control Protocol (RTCP)", [BCP 131](#), [RFC 4961](#), July 2007.
- [RFC5461] Gont, F., "TCP's Reaction to Soft Errors", [RFC 5461](#), February 2009.
- [RFC6092] Woodyatt, J., "Recommended Simple Security Capabilities in Customer Premises Equipment (CPE) for Providing Residential IPv6 Internet Service", [RFC 6092](#), January 2011.

[Appendix A](#). Changes

[A.1](#). Changes from [draft-ietf-pcp-base-05](#) to -06

- o DS-Lite: consensus was encapsulation mode. Included a suggestion that the B4 will need to proxy PCP-to-PCP and UPnP-to-PCP.
- o defined THIRD_PARTY option to work with the PEER OpCode, too. This meant moving it to its own section, and having both MAP and PEER OpCodes reference that common section.
- o used "target" instead of "internal", in the hopes that clarifies internal address used by PCP itself (for sending its packets) versus the address for Mappings.
- o Options are now required to be ordered in requests, and ordering has to be validated by the server. Intent is to ease server processing of mandatory-to-implement options.
- o Swapped Option values for the mandatory- and optional-to-process Options, so we can have a simple lowest..highest ordering.
- o added MISORDERED_OPTIONS error.
- o re-ordered some error messages to cause MALFORMED_REQUEST (which is PCP's most general error response) to be error 1, instead of buried in the middle of the error numbers.
- o clarified that, after successfully using a PCP server, that PCP server is declared to be non-responsive after 5 failed retransmissions.

- o tightened up text (which was inaccurate) about how long general PCP processing is to delay when receiving an error and if it should honor OpCode-specific error lifetime. Useful for MAP errors which have an error lifetime. (This all feels awkward to have only some errors with a lifetime.)
- o Added better discussion of multiple interfaces, including highlighting WiFi+Ethernet. Added discussion of using IPv6 Privacy Addresses and [RFC1918](#) as source addresses for PCP requests. This should finish the section on multi-interface issues.
- o added some text about why server might send SERVER_OVERLOADED, or might simply discard packets.
- o Dis-allow internal-port=0, which means we dis-allow using PCP as a DMZ-like function. Instead, ports have to be mapped individually.
- o Text describing server's processing of PEER is tightened up.
- o Server's processing of PEER now says it is implementation-specific if a PCP server continues to allow the mapping to exist after a PEER message. Client's processing of PEER says that if client wants mapping to continue to exist, client has to continue to send recurring PEER messages.

A.2. Changes from [draft-ietf-pcp-base-04](#) to -05

- o tweaked PCP common header packet layout.
- o Re-added port=0 (all ports).
- o minimum size is 12 octets (missed that change in -04).
- o removed Lifetime from PCP common header.
- o for MAP error responses, the lifetime indicates how long the server wants the client to avoid retrying the request.
- o More clearly indicated which fields are filled by the server on success responses and error responses.
- o Removed UPnP interworking section from this document. It will appear in [[I-D.bpw-pcp-upnp-igd-interworking](#)].

A.3. Changes from [draft-ietf-pcp-base-03](#) to -04

- o "Pinhole" and "PIN" changed to "mapping" and "MAP".
- o Reduced from four MAP OpCodes to two. This was done by implicitly using the address family of the PCP message itself.
- o New option THIRD_PARTY, to more carefully split out the case where a mapping is created to a different host within the home.
- o Integrated a lot of editorial changes from Stuart and Francis.
- o Removed nested NAT text into another document, including the IANA-registered IP addresses for the PCP server.
- o Removed suggestion (MAY) that PCP server reserve UDP when it maps TCP. Nobody seems to need that.
- o Clearly added NAT and NAPT, such as in residential NATs, as within scope for PCP.
- o HONOR_EXTERNAL_PORT renamed to PREFER_FAILURE
- o Added 'Lifetime' field to the common PCP header, which replaces the functions of the 'temporary' and 'permanent' error types of the previous version.
- o Allow arbitrary Options to be included in PCP response, so that PCP server can indicate un-supported PCP Options. Satisfies PCP Issue #19
- o Reduced scope to only deal with mapping protocols that have port numbers.
- o Reduced scope to not support DMZ-style forwarding.
- o Clarified version negotiation.

A.4. Changes from [draft-ietf-pcp-base-02](#) to -03

- o Adjusted abstract and introduction to make it clear PCP is intended to forward ports and intended to reduce application keepalives.
- o First bit in PCP common header is set. This allows DTLS and non-DTLS to be multiplexed on same port, should a future update to this specification add DTLS support.

- o Moved subscriber identity from common PCP section to MAP* section.
- o made clearer that PCP client can reduce mapping lifetime if it wishes.
- o Added discussion of host running a server, client, or symmetric client+server.
- o Introduced PEER4 and PEER6 OpCodes.
- o Removed REMOTE_PEER Option, as its function has been replaced by the new PEER OpCodes.
- o IANA assigned port 44323 to PCP.
- o Removed AMBIGUOUS error code, which is no longer needed.

A.5. Changes from [draft-ietf-pcp-base-01](#) to -02

- o more error codes
- o PCP client source port number should be random
- o PCP message minimum 8 octets, maximum 1024 octets.
- o tweaked a lot of text in [section 7.4](#), "Opcode-Specific Server Operation".
- o opening a mapping also allows ICMP messages associated with that mapping.
- o PREFER_FAILURE value changed to the mandatory-to-process range.
- o added text recommending applications that are crashing obtain short lifetimes, to avoid consuming subscriber's port quota.

A.6. Changes from [draft-ietf-pcp-base-00](#) to -01

- o Significant document reorganization, primarily to split base PCP operation from OpCode operation.
- o packet format changed to move 'protocol' outside of PCP common header and into the MAP* opcodes
- o Renamed Informational Elements (IE) to Options.
- o Added REMOTE_PEER (for disambiguation with dynamic ports), REMOTE_PEER_FILTER (for simple packet filtering), and

PREFER_FAILURE (to optimize UPnP IGD interworking) options.

- o Is NAT or router behind B4 in scope?
- o PCP option MAY be included in a request, in which case it MUST appear in a response. It MUST NOT appear in a response if it was not in the request.
- o Response code most significant bit now indicates permanent/temporary error
- o PCP Options are split into mandatory-to-process ("P" bit), and into Specification Required and Private Use.
- o Epoch discussion simplified.

Authors' Addresses

Dan Wing (editor)
Cisco Systems, Inc.
170 West Tasman Drive
San Jose, California 95134
USA

Email: dwing@cisco.com

Stuart Cheshire
Apple, Inc.
1 Infinite Loop
Cupertino, California 95014
USA

Phone: +1 408 974 3207
Email: cheshire@apple.com

Mohamed Boucadair
France Telecom
Rennes, 35000
France

Email: mohamed.boucadair@orange-ftgroup.com

Reinaldo Penno
Juniper Networks
1194 N Mathilda Avenue
Sunnyvale, California 94089
USA

Email: rpenno@juniper.net

Francis Dupont
Internet Systems Consortium

Email: fdupont@isc.org