

Internet Engineering Task Force (IETF)
Request for Comments: 7600
Category: Experimental
ISSN: 2070-1721

R. Despres
RD-IPtech
S. Jiang, Ed.
Huawei Technologies Co., Ltd
R. Penno
Cisco Systems, Inc.
Y. Lee
Comcast
G. Chen
China Mobile
M. Chen
BBIX, Inc.
July 2015

IPv4 Residual Deployment via IPv6 - A Stateless Solution (4rd)

Abstract

This document specifies a stateless solution for service providers to progressively deploy IPv6-only network domains while still offering IPv4 service to customers. The solution's distinctive properties are that TCP/UDP IPv4 packets are valid TCP/UDP IPv6 packets during domain traversal and that IPv4 fragmentation rules are fully preserved end to end. Each customer can be assigned one public IPv4 address, several public IPv4 addresses, or a shared address with a restricted port set.

Status of This Memo

This document is not an Internet Standards Track specification; it is published for examination, experimental implementation, and evaluation.

This document defines an Experimental Protocol for the Internet community. This document is a product of the Internet Engineering Task Force (IETF). It represents the consensus of the IETF community. It has received public review and has been approved for publication by the Internet Engineering Steering Group (IESG). Not all documents approved by the IESG are a candidate for any level of Internet Standard; see Section 2 of RFC 5741.

Information about the current status of this document, any errata, and how to provide feedback on it may be obtained at <http://www.rfc-editor.org/info/rfc7600>.

Copyright Notice

Copyright (c) 2015 IETF Trust and the persons identified as the document authors. All rights reserved.

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

Table of Contents

1. Introduction	4
2. Terminology	5
3. The 4rd Model	7
4. Protocol Specifications	9
4.1. NAT44 on CE	9
4.2. Mapping Rules and Other Domain Parameters	10
4.3. Reversible Packet Translations at Domain Entries and Exits	11
4.4. Address Mapping from CE IPv6 Prefixes to 4rd IPv4 Prefixes	17
4.5. Address Mapping from 4rd IPv4 Addresses to 4rd IPv6 Addresses	19
4.6. Fragmentation Processing	23
4.6.1. Fragmentation at Domain Entry	23
4.6.2. Ports of Fragments Addressed to Shared-Address CEs	24
4.6.3. Packet Identifications from Shared-Address CEs	26
4.7. TOS and Traffic Class Processing	26
4.8. Tunnel-Generated ICMPv6 Error Messages	27
4.9. Provisioning 4rd Parameters to CEs	27
5. Security Considerations	30
6. IANA Considerations	31
7. Relationship with Previous Works	31
8. References	33
8.1. Normative References	33
8.2. Informative References	34
Appendix A. Textual Representation of Mapping Rules	37
Appendix B. Configuring Multiple Mapping Rules	37
Appendix C. Adding Shared IPv4 Addresses to an IPv6 Network	39
C.1. With CEs within CPEs	39
C.2. With Some CEs behind Third-Party Router CPEs	41
Appendix D. Replacing Dual-Stack Routing with IPv6-Only Routing ...	42
Appendix E. Adding IPv6 and 4rd Service to a Net-10 Network	43
Acknowledgements	44
Authors' Addresses	44

1. Introduction

For service providers to progressively deploy IPv6-only network domains while still offering IPv4 service to customers, the need for a stateless solution, i.e., one where no per-customer state is needed in IPv4-IPv6 gateway nodes of the provider, has been discussed in [Solutions-4v6]. This document specifies one such solution, named "4rd" for IPv4 Residual Deployment. Its distinctive properties are that TCP/UDP IPv4 packets are valid TCP/UDP IPv6 packets during domain traversal and that IPv4 fragmentation rules are fully preserved end to end.

Using this solution, IPv4 packets are transparently tunneled across IPv6 networks (the reverse of IPv6 Rapid Deployment on IPv4 Infrastructures (6rd) [RFC5969], in which IPv6 packets are statelessly tunneled across IPv4 networks).

While IPv6 headers are too long to be mapped into IPv4 headers (which is why 6rd requires encapsulation of full IPv6 packets in IPv4 packets), IPv4 headers can be reversibly translated into IPv6 headers in such a way that, during IPv6 domain traversal, UDP packets having checksums and TCP packets are valid IPv6 packets. IPv6-only middleboxes that perform deep packet inspection can operate on them, in particular for port inspection and web caches.

In order to deal with the IPv4 address shortage, customers can be assigned shared public IPv4 addresses with statically assigned restricted port sets. As such, it is a particular application of the Address plus Port (A+P) approach [RFC6346].

Deploying 4rd in networks that have enough public IPv4 addresses, customer sites can also be assigned full public IPv4 addresses. 4rd also supports scenarios where a set of public IPv4 addresses are assigned to customer sites.

The design of 4rd builds on a number of previous proposals made for IPv4-via-IPv6 transition technologies (Section 7).

In some use cases, IPv4-only applications of 4rd-capable customer nodes can also work with stateful NAT64s [RFC6146], provided these are upgraded to support 4rd tunnels in addition to their IP/ICMP translation [RFC6145]. The advantage is then a more complete IPv4 transparency than with double translation.

How the 4rd model fits in the Internet architecture is summarized in Section 3. The protocol specifications are detailed in Section 4. Sections 5 and 6 deal with security considerations and IANA considerations, respectively. Previous proposals that influenced

this specification are listed in Section 7. A few typical 4rd use cases are presented in Appendices A, B, C, D, and E.

2. 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 [RFC2119].

ISP: Internet Service Provider. In this document, the service it offers can be DSL, fiber-optics, cable, or mobile. The ISP can also be a private-network operator.

4rd (IPv4 Residual Deployment): An extension of the IPv4 service where public IPv4 addresses can be statically shared among several customer sites, each one being assigned an exclusive port set. This service is supported across IPv6-routing domains.

4rd domain (or Domain): An ISP-operated IPv6 network across which 4rd is supported according to the present specification.

Tunnel packet: An IPv6 packet that transparently conveys an IPv4 packet across a 4rd domain. Its header has enough information to reconstitute the IPv4 header at Domain exit. Its payload is the original IPv4 payload.

CE (Customer Edge): A customer-side tunnel endpoint. It can be in a node that is a host, a router, or both.

BR (Border Relay): An ISP-side tunnel endpoint. Because its operation is stateless (neither per CE nor per session state), it can be replicated in as many nodes as needed for scalability.

4rd IPv6 address: IPv6 address used as the destination of a Tunnel packet sent to a CE or a BR.

NAT64+: An ISP NAT64 [RFC6146] that is upgraded to support 4rd tunneling when IPv6 addresses it deals with are 4rd IPv6 addresses.

4rd IPv4 address: A public IPv4 address or, in the case of a shared public IPv4 address, a public transport address (public IPv4 address plus port number).

PSID (Port-Set Identifier): A flexible-length field that algorithmically identifies a port set.

4rd IPv4 prefix: A flexible-length prefix that may be a public IPv4 prefix, a public IPv4 address, or a public IPv4 address followed by a PSID.

Mapping rule: A set of parameters that are used by BRs and CEs to derive 4rd IPv6 addresses from 4rd IPv4 addresses. Mapping rules are also used by each CE to derive a 4rd IPv4 prefix from an IPv6 prefix that has been delegated to it.

EA bits (Embedded Address bits): Bits that are the same in a 4rd IPv4 address and in the 4rd IPv6 address derived from it.

BR Mapping rule: The Mapping rule that is applicable to off-domain IPv4 addresses (addresses reachable via BRs). It can also apply to some or all CE-assigned IPv4 addresses.

CE Mapping rule: A Mapping rule that is applicable only to CE-assigned IPv4 addresses (shared or not).

NAT64+ Mapping rule: The Mapping rule that is applicable to IPv4 addresses reachable via a NAT64+.

CNP (Checksum Neutrality Preserver): A field of 4rd IPv6 addresses that ensures that TCP-like checksums do not change when IPv4 addresses are replaced with 4rd IPv6 addresses.

4rd Tag: A 16-bit tag whose value allows 4rd CEs, BRs, and NAT64+s to distinguish 4rd IPv6 addresses from other IPv6 addresses.

3. The 4rd Model

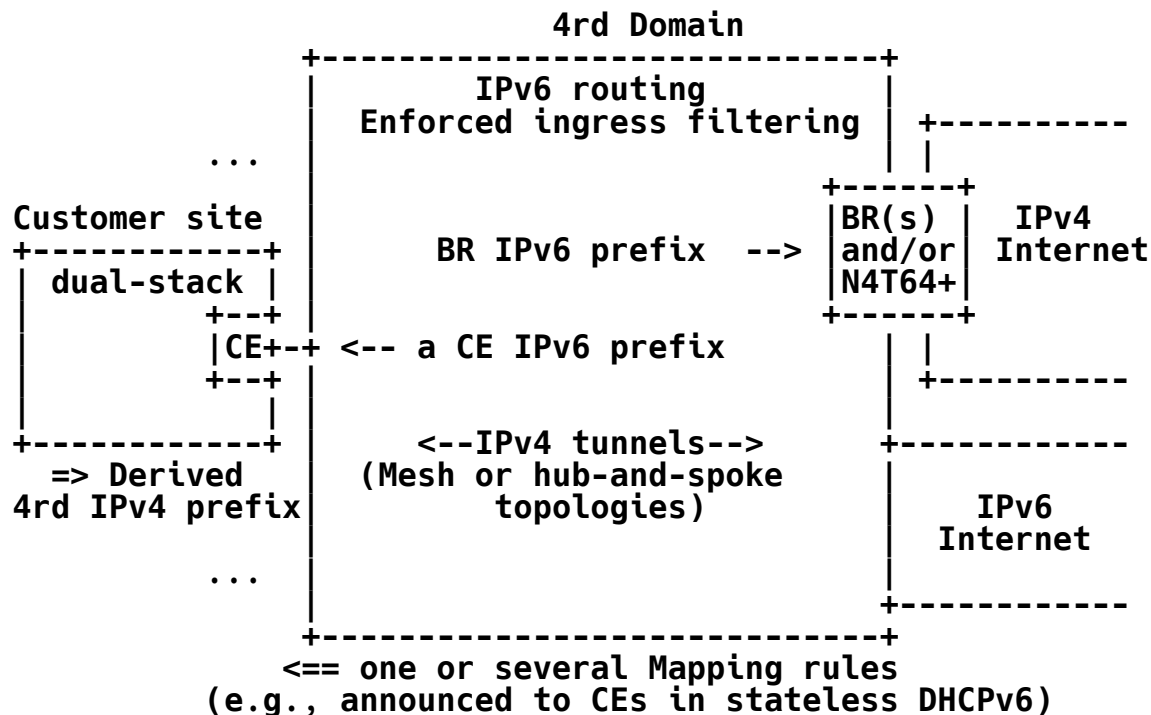


Figure 1: The 4rd Model in the Internet Architecture

How the 4rd model fits in the Internet architecture is represented in Figure 1.

A 4rd domain is an IPv6 network that includes one or several 4rd BRs or NAT64+s at its border with the public IPv4 Internet and that can advertise its IPv4-IPv6 Mapping rule(s) to CEs according to Section 4.9.

BRs of a 4rd Domain are all identical as far as 4rd is concerned. In a 4rd CE, the IPv4 packets that need to reach a BR will be transformed (as detailed in Section 4.3) into IPv6 packets that have the same anycast IPv6 prefix, which is the 80-bit BR prefix, in their destination addresses. They are then routed to any of the BRs. The 80-bit BR IPv6 prefix is an arbitrarily chosen /64 prefix from the IPv6 address space of the network operator and appended with 0x0300 (16-bit 4rd Tag; see R-9 in Section 4.5).

Using the Mapping rule that applies, each CE derives its 4rd IPv4 prefix from its delegated IPv6 prefix, or one of them if it has several; see Section 4.4 for details. If the obtained IPv4 prefix has more than 32 bits, the assigned IPv4 address is shared among several CEs. Bits beyond the first 32 specify a set of ports whose use is reserved for the CE.

IPv4 traffic is automatically tunneled across the Domain, in either mesh topology or hub-and-spoke topology [RFC4925]. By default, IPv4 traffic between two CEs follows a direct IPv6 route between them (mesh topology). If the ISP configures the hub-and-spoke option, each IPv4 packet from one CE to another is routed via a BR.

During Domain traversal, each tunneled TCP/UDP IPv4 packet looks like a valid TCP/UDP IPv6 packet. Thus, TCP/UDP access control lists that apply to IPv6, and possibly some other functions using deep packet inspection, also apply to IPv4.

In order for IPv4 anti-spoofing protection in CEs and BRs to remain effective when combined with 4rd tunneling, ingress filtering [RFC3704] has to be in effect in IPv6 (see R-12 and Section 5).

If an ISP wishes to support dynamic IPv4 address sharing in addition to or in place of 4rd stateless address sharing, it can do so by means of a stateful NAT64. By upgrading this NAT to add support for 4rd tunnels, which makes it a NAT64+, CEs that are assigned no static IPv4 space can benefit from complete IPv4 transparency between the CE and the NAT64. (Without this NAT64 upgrade, IPv4 traffic is translated to IPv6 and back to IPv4, during which time the DF = MF = 1 combination for IPv4, as recommended for host fragmentation in Section 8 of [RFC4821], is lost.)

IPv4 packets are kept unchanged by Domain traversal, except that:

- o The IPv4 Time To Live (TTL), unless it is 1 or 255 at Domain entry, decreases during Domain traversal by the number of traversed routers. This is acceptable because it is undetectable end to end and also because TTL values that can be used with some protocols to test the adjacency of communicating routers are preserved [RFC4271] [RFC5082]. The effect on the traceroute utility, which uses TTL expiry to discover routers of end-to-end paths, is noted in Section 4.3.

- o IPv4 packets whose lengths are ≤ 68 octets always have their "Don't Fragment" (DF) flags set to 1 at Domain exit even if they had DF = 0 at Domain entry. This is acceptable because these packets are too short to be fragmented [RFC791] and so their DF bits have no meaning. Besides, both [RFC1191] and [RFC4821] recommend that sources always set DF to 1.
- o Unless the Tunnel Traffic Class option applies to a Domain (Section 4.2), IPv4 packets may have their Type of Service (TOS) fields modified after Domain traversal (Section 4.7).

4. Protocol Specifications

This section describes detailed 4rd protocol specifications. They are mainly organized by functions. As a brief summary:

- o A 4rd CE MUST follow R-1, R-2, R-3, R-4, R-6, R-7, R-8, R-9, R-10, R-11, R-12, R-13, R-14, R-16, R-17, R-18, R-19, R-20, R-21, R-22, R-23, R-24, R-25, R-26, and R-27.
- o A 4rd BR MUST follow R-2, R-3, R-4, R-5, R-6, R-9, R-12, R-13, R-14, R-15, R-19, R-20, R-21, R-22, and R-24.

4.1. NAT44 on CE

R-1: A CE node that is assigned a shared public IPv4 address MUST include a NAT44 [RFC3022]. This NAT44 MUST only use external ports that are in the CE-assigned port set.

NOTE: This specification only concerns IPv4 communication between IPv4-capable endpoints. For communication between IPv4-only endpoints and IPv6-only remote endpoints, the "Bump-in-the-Host" (BIH) specification [RFC6535] can be used. It can coexist in a node with the CE function, including scenarios where the IPv4-only function is a NAT44 [RFC3022].

4.2. Mapping Rules and Other Domain Parameters

R-2: CEs and BRs MUST be configured with the following Domain parameters:

- A. One or several Mapping rules, each one comprising the following:
 - 1. Rule IPv4 prefix
 - 2. EA-bits length
 - 3. Rule IPv6 prefix
 - 4. Well-Known Ports (WKPs) authorized (OPTIONAL)
- B. Domain Path MTU (PMTU)
- C. Hub-and-spoke topology (Yes or No)
- D. Tunnel Traffic Class (OPTIONAL)

"Rule IPv4 prefix" is used to find, by a longest match, which Mapping rule applies to a 4rd IPv4 address (Section 4.5). A Mapping rule whose Rule IPv4 prefix is longer than /0 is a CE Mapping rule. BR and NAT64+ Mapping rules, which must apply to all off-domain IPv4 addresses, have /0 as their Rule IPv4 prefixes.

"EA-bits length" is the number of bits that are common to 4rd IPv4 addresses and 4rd IPv6 addresses derived from them. In a CE Mapping rule, it is also the number of bits that are common to a CE-delegated IPv6 prefix and the 4rd IPv4 prefix derived from it. BR and NAT64+ Mapping rules have EA-bits lengths equal to 32.

"Rule IPv6 prefix" is the prefix that is used as a substitute for the Rule IPv4 prefix when a 4rd IPv6 address is derived from a 4rd IPv4 address (Section 4.5). In a BR Mapping rule or a NAT64+ Mapping rule, it MUST be a /80 prefix whose bits 64-79 are the 4rd Tag.

"WKPs authorized" may be set for Mapping rules that assign shared IPv4 addresses to CEs. (These rules are those whose length of the Rule IPv4 prefix plus the EA-bits length exceeds 32.) If set, well-known ports may be assigned to some CEs having particular IPv6 prefixes. If not set, fairness is privileged: all IPv6 prefixes concerned with the Mapping rule have port sets having identical values (no port set includes any of the well-known ports).

"Domain PMTU" is the IPv6 Path MTU that the ISP can guarantee for all of its IPv6 paths between CEs and between BRs and CEs. It MUST be at least 1280 octets [RFC2460].

"Hub-and-spoke topology", if set to Yes, requires CEs to tunnel all IPv4 packets via BRs. If set to No, CE-to-CE packets take the same routes as native IPv6 packets between the same CEs (mesh topology).

"Tunnel Traffic Class", if provided, is the IPv6 traffic class that BRs and CEs MUST set in Tunnel packets. In this case, evolutions of the IPv6 traffic class that may occur during Domain traversal are not reflected in TOS fields of IPv4 packets at Domain exit (Section 4.7).

4.3. Reversible Packet Translations at Domain Entries and Exits

R-3: Domain-entry nodes that receive IPv4 packets with IPv4 options MUST discard these packets and return ICMPv4 error messages to signal IPv4-option incompatibility (Type = 12, Code = 0, Pointer = 20) [RFC792]. This limitation is acceptable because there are a lot of firewalls in the current IPv4 Internet that also filter IPv4 packets with IPv4 options.

R-4: Domain-entry nodes that receive IPv4 packets without IPv4 options MUST convert them to Tunnel packets, with or without IPv6 fragment headers, depending on what is needed to ensure IPv4 transparency (Figure 2). Domain-exit nodes MUST convert them back to IPv4 packets.

An IPv6 fragmentation header MUST be included at tunnel entry (Figure 2) if and only if one or several of the following conditions hold:

- * The Tunnel Traffic Class option applies to the Domain.
- * TTL = 1 OR TTL = 255.
- * The IPv4 packet is already fragmented, or may be fragmented later on, i.e., if MF = 1 OR offset > 0 OR (total length > 68 AND DF = 0).

In order to optimize cases where fragmentation headers are unnecessary, the NAT44 of a CE that has one SHOULD send packets with TTL = 254.

- R-5: In Domains whose chosen topology is hub-and-spoke, BRs that receive 4rd IPv6 packets whose embedded destination IPv4 addresses match a CE Mapping rule **MUST** do the equivalent of reversibly translating their headers to IPv4 and then reversibly translate them back to IPv6 as though packets would be entering the Domain.

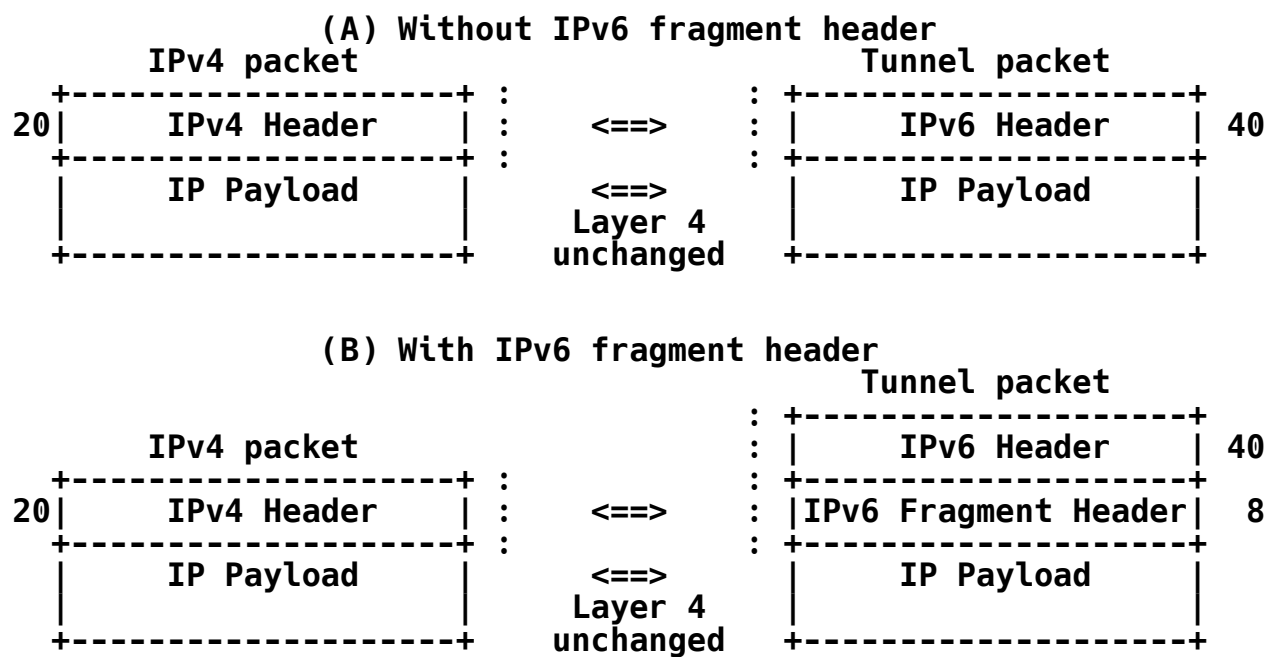


Figure 2: Reversible Packet Translation

IPv6 Field	Value (fields from IPv4 header)
Version	6
Traffic Class	TOS OR Tunnel Traffic Class (Section 4.7)
Addr_Prot_Cksm	Sum of addresses and Protocol (Note 1)
Payload length	Total length - 12
Next header	44 (fragment header)
Hop limit	IF Time to Live = 1 or 255 THEN 254 ELSE Time to Live (Note 2)
Source address	See Section 4.5
Dest. address	See Section 4.5
2nd next header	Protocol
Fragment offset	IPv4 fragment offset
M	More Fragments flag (MF)
IPv4_DF	Don't Fragment flag (DF)
TTL_1	IF Time to Live = 1 THEN 1 ELSE 0 (Note 2)
TTL_255	IF Time to Live = 255 THEN 1 ELSE 0 (Note 2)
IPv4_TOS	Type of Service (TOS)
IPv4_ID	Identification

Table 2: IPv4-to-IPv6 Reversible Header Translation
(with Fragment Header)

IPv4 Field	Value (fields from IPv6 header)
Version	4
Header length	5
TOS	Traffic Class
Total length	Payload length + 20
Identification	0
DF	1
MF	0
Fragment offset	0
Time to Live	Hop count
Protocol	Next header
Header checksum	Computed as per [RFC791] (Note 3)
Source address	Bits 80-111 of source address
Dest. address	Bits 80-111 of destination address

Table 3: IPv6-to-IPv4 Reversible Header Translation
(without Fragment Header)

IPv4 Field	Value (fields from IPv6 header)
Version	4
Header length	5
TOS	Traffic Class OR IPv4_TOS (Section 4.7)
Total length	Payload length + 12
Identification	IPv4_ID
DF	IPv4_DF
MF	M
Fragment offset	Fragment offset
Time to Live (Note 2)	IF TTL_255 = 1 THEN 255 ELSEIF TTL_1 = 1 THEN 1 ELSE hop count
Protocol	2nd next header
Header checksum	Computed as per [RFC791] (Note 3)
Source address	Bits 80-111 of source address
Destination address	Bits 80-111 of destination address

Table 4: IPv6-to-IPv4 Reversible Header Translation
(with Fragment Header)

NOTE 1: The need to save in the IPv6 header a checksum of both IPv4 addresses and the IPv4 protocol field results from the following facts: (1) header checksums, present in IPv4 but not in IPv6, protect addresses or protocol integrity; (2) in IPv4, ICMP messages and null-checksum UDP datagrams depend on this protection because, unlike other datagrams, they have no other address-and-protocol integrity protection. The sum MUST be performed in ordinary two's complement arithmetic.

IP-layer Packet length is another field covered by the IPv4 header checksum. It is not included in the saved checksum because (1) doing so would have conflicted with [RFC6437] (flow labels must be the same in all packets of each flow); (2) ICMPv4 messages have good enough protection with their own checksums; (3) the UDP length field provides to null-checksum UDP datagrams the same level of protection after Domain traversal as without Domain traversal (consistency between IP-layer and UDP-layer lengths can be checked).

NOTE 2: TTL treatment has been chosen to permit adjacency tests between two IPv4 nodes situated at both ends of a 4rd tunnel. TTL values to be preserved for this are TTL = 255 and TTL = 1. For other values, TTL decreases between two IPv4 nodes as though the traversed IPv6 routers were IPv4 routers.

The effect of this TTL treatment on IPv4 traceroute is specific: (1) the number of routers of the end-to-end path includes traversed IPv6 routers; (2) IPv6 routers of a Domain are listed after IPv4 routers of Domain entry and exit; (3) the IPv4 address shown for an IPv6 router is the IPv6-only dummy IPv4 address (Section 4.8); (4) the response time indicated for an IPv6 router is that of the next router.

NOTE 3: Provided the sum of obtained IPv4 addresses and protocol matches Addr_Prot_Cksm. If not, the packet MUST be silently discarded.

4.4. Address Mapping from CE IPv6 Prefixes to 4rd IPv4 Prefixes

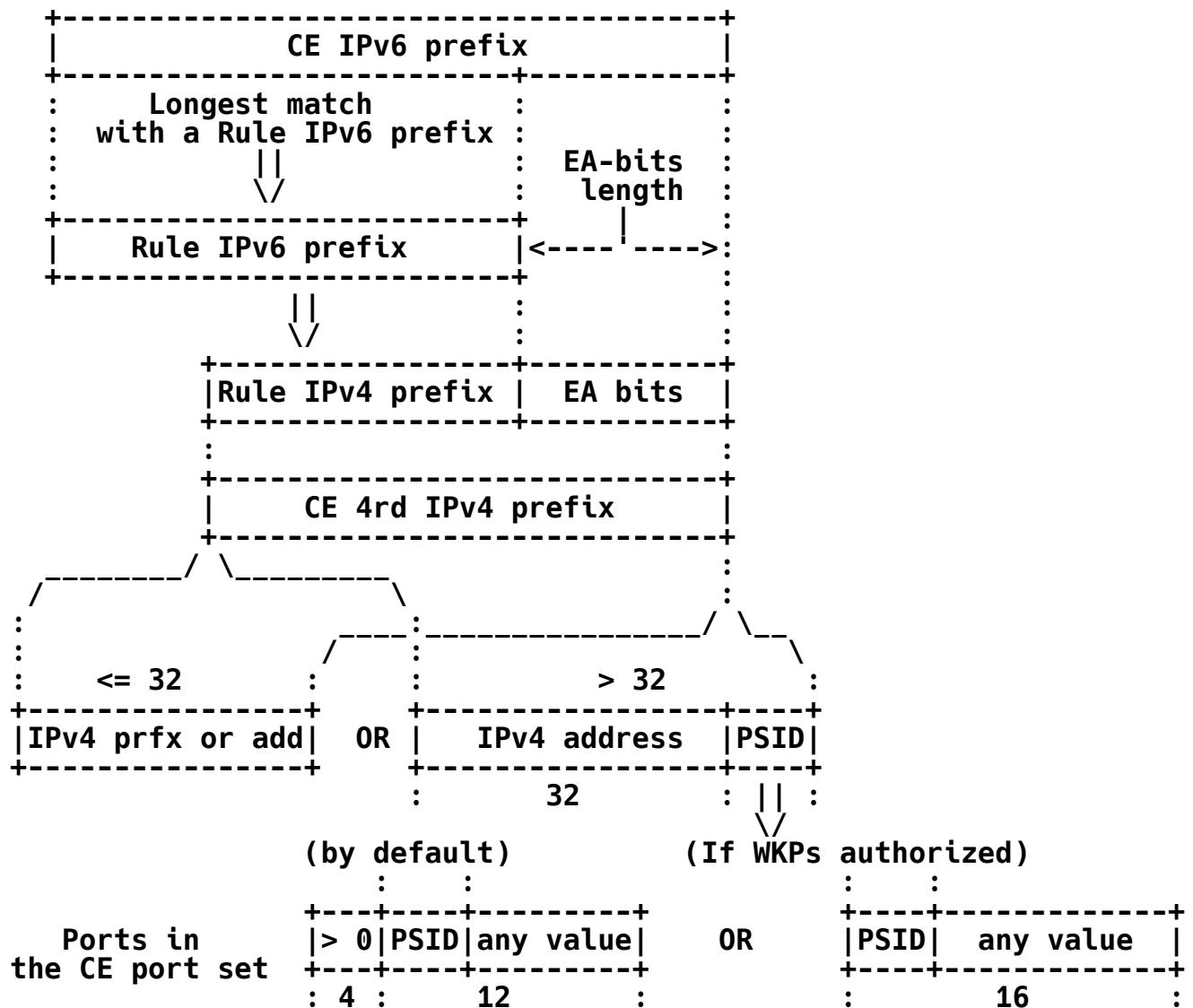


Figure 4: From CE IPv6 Prefix to 4rd IPv4 Address and Port Set

R-7: A CE whose delegated IPv6 prefix matches the Rule IPv6 prefix of one or several Mapping rules MUST select the CE Mapping rule for which the match is the longest. It then derives its 4rd IPv4 prefix as shown in Figure 4: (1) The CE replaces the Rule IPv6 prefix with the Rule IPv4 prefix. The result is the CE 4rd IPv4 prefix. (2) If this CE 4rd IPv4 prefix has less than 32 bits, the CE takes it as its assigned IPv4 prefix. If it has exactly 32 bits, the CE takes it as its IPv4 address. If

it has more than 32 bits, the CE MUST take the first 32 bits as its shared public IPv4 address and bits beyond the first 32 as its Port-Set identifier (PSID). Ports of its restricted port set are by default those that have any non-zero value in their first 4 bits (the PSID offset), followed by the PSID, and followed by any values in remaining bits. If the WKP authorized option applies to the Mapping rule, there is no 4-bit offset before the PSID so that all ports can be assigned.

NOTE: The choice of the default PSID position in port fields has been guided by the following objectives: (1) for fairness, avoid having any of the well-known ports 0-1023 in the port set specified by any PSID value; (2) for compatibility with RTP/RTCP [RFC4961], include in each port set pairs of consecutive ports; (3) in order to facilitate operation and training, have the PSID at a fixed position in port fields; (4) in order to facilitate documentation in hexadecimal notation, and to facilitate maintenance, have this position nibble-aligned. Ports that are excluded from assignment to CEs are 0-4095, instead of just 0-1023, in a trade-off to favor nibble alignment of PSIDs and overall simplicity.

- R-8: A CE whose delegated IPv6 prefix has its longest match with the Rule IPv6 prefix of the BR Mapping rule MUST take as its IPv4 address the 32 bits that, in the delegated IPv6 prefix, follow this Rule IPv6 prefix. If this is the case while the hub-and-spoke option applies to the Domain, or if the Rule IPv6 prefix is not a /80, there is a configuration error in the Domain. An implementation-dependent administrative action MAY be taken.

A CE whose delegated IPv6 prefix does not match the Rule IPv6 prefix of either any CE Mapping rule or the BR Mapping rule, and is in a Domain that has a NAT64+ Mapping rule, MUST be noted as having the unspecified IPv4 address.

4.5. Address Mapping from 4rd IPv4 Addresses to 4rd IPv6 Addresses

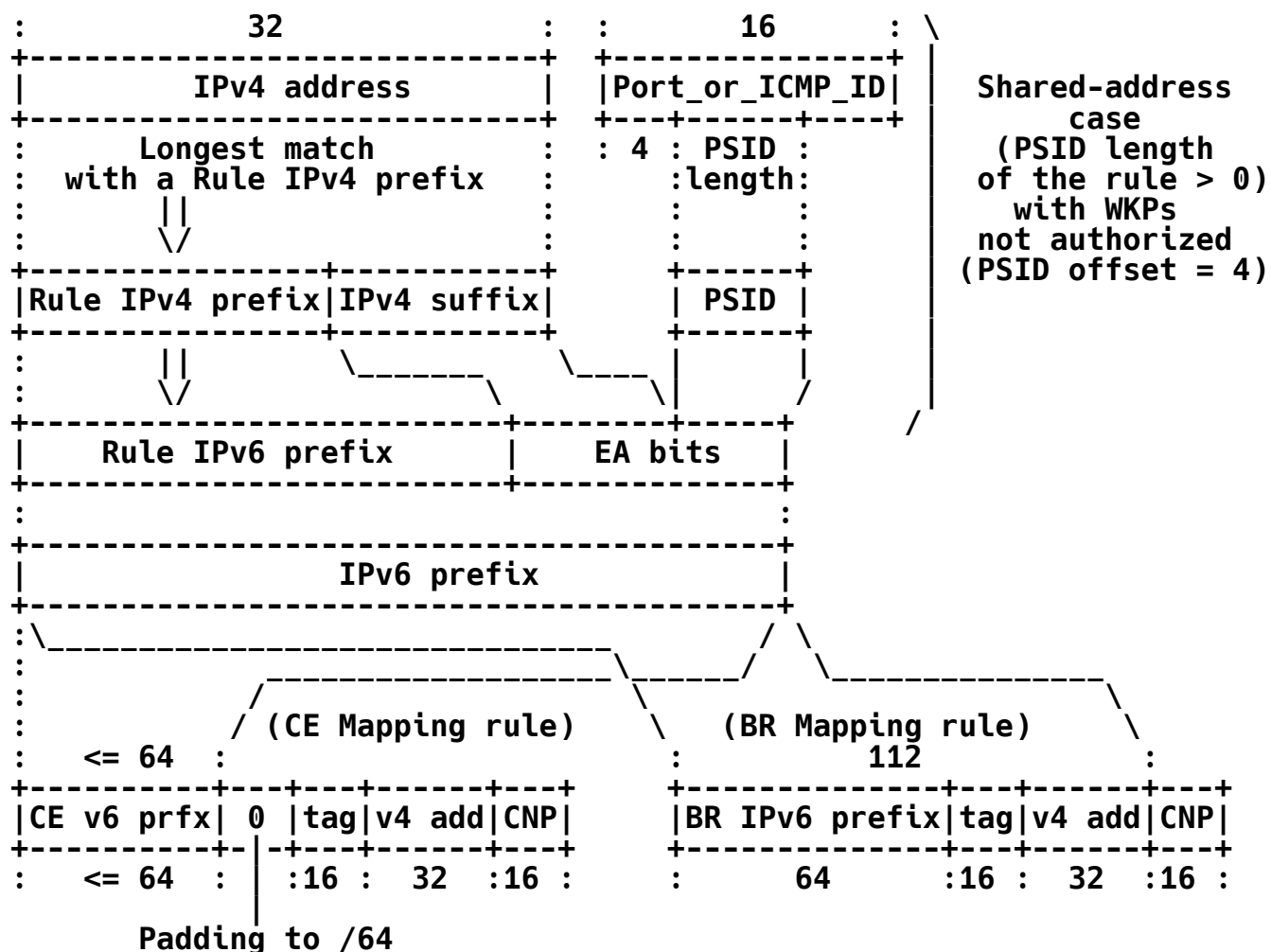


Figure 5: From 4rd IPv4 Address to 4rd IPv6 Address

- R-9: BRs, and CEs that are assigned public IPv4 addresses, shared or not, MUST derive 4rd IPv6 addresses from 4rd IPv4 addresses via the steps below or their functional equivalent (Figure 5 details the shared public IPv4 address case):

NOTE: The rules for forming 4rd-specific Interface Identifiers (IIDs) are to obey [RFC7136]:

"Specifications of other forms of 64-bit IIDs MUST specify how all 64 bits are set."

and

"the whole IID value MUST be viewed as an opaque bit string by third parties, except possibly in the local context."

- (1) If hub-and-spoke topology does not apply to the Domain, or if it applies but the IPv6 address to be derived is a source address from a CE or a destination address from a BR, find the CE Mapping rule whose Rule IPv4 prefix has the longest match with the IPv4 address.

If no Mapping rule is thus obtained, take the BR Mapping rule.

If the obtained Mapping rule assigns IPv4 prefixes to CEs, i.e., if the length of the Rule IPv4 prefix plus EA-bits length is $32 - k$, with $k \geq 0$, delete the last k bits of the IPv4 address.

Otherwise, if the length of the Rule IPv4 prefix plus the EA-bits length is $32 + k$, with $k > 0$, take k as the PSID length and append to the IPv4 address the PSID copied from bits p to $p+3$ of the Port_or_ICMP_ID field where (1) p , the PSID offset, is 4 by default and 0 if the WKPs authorized option applies to the rule; (2) the Port_or_ICMP_ID field is in bits of the IP payload that depend on whether the address is source or destination, on whether the packet is ICMP or not, and, if it is ICMP, whether it is an error message or an Echo message. This field is:

- a. If the packet Protocol is not ICMP, the port field associated with the address (bits 0-15 for a source address and bits 16-31 for a destination address).
- b. If the packet is an ICMPv4 Echo or Echo reply message, the ICMPv4 Identification field (bits 32-47).

- c. If the packet is an ICMPv4 error message, the port field associated with the address in the returned packet header (bits 240-255 for a source address and bits 224-239 for a destination address).

NOTE 1: Using Identification fields of ICMP messages as port fields permits the exchange of Echo requests and Echo replies between shared-address CEs and IPv4 hosts having exclusive IPv4 addresses. Echo exchanges between two shared-address CEs remain impossible, but this is a limitation inherent in address sharing (one reason among many to use IPv6).

NOTE 2: When the PSID is taken in the port fields of the IPv4 payload, implementation is kept independent from any particular Layer 4 protocol having such port fields by not checking that the protocol is indeed one that has such port fields. A packet may consequently go, in the case of a source mistake, from a BR to a shared-address CE with a protocol that is not supported by this CE. In this case, the CE NAT44 returns an ICMPv4 "protocol unreachable" error message. The IPv4 source is thus appropriately informed of its mistake.

- (2) In the result, replace the Rule IPv4 prefix with the Rule IPv6 prefix.
- (3) If the result is shorter than a /64, append to the result a null padding up to 64 bits, followed by the 4rd Tag (0x0300), and followed by the IPv4 address.

NOTE: The 4rd Tag is a 4rd-specific mark. Its function is to ensure that 4rd IPv6 addresses are recognizable by CEs without any interference with the choice of subnet prefixes in CE sites. (These choices may have been done before 4rd is enabled.)

For this, the 4rd Tag has its "u" and "g" bits [RFC4291] both set to 1, so that they maximally differ from these existing IPv6 address schemas. So far, u = g = 1 has not been used in any IPv6 addressing architecture.

With the 4rd Tag, IPv6 packets can be routed to the 4rd function within a CE node based on a /80 prefix that no native IPv6 address can contain.

- (4) Add to the result a Checksum Neutrality Preserver (CNP). Its value, in one's complement arithmetic, is the opposite of the sum of 16-bit fields of the IPv6 address other than the IPv4 address and the CNP themselves (i.e., five consecutive fields in address bits 0-79).

NOTE: The CNP guarantees that Tunnel packets are valid IPv6 packets for all Layer 4 protocols that use the same checksum algorithm as TCP. This guarantee does not depend on where the checksum fields of these protocols are placed in IP payloads. (Today, such protocols are UDP [RFC768], TCP [RFC793], UDP-Lite [RFC3828], and the Datagram Congestion Control Protocol (DCCP) [RFC5595]. Should new ones be specified, BRs will support them without needing an update.)

- R-10: A 4rd-capable CE SHOULD, and a 4rd-enabled CE MUST, always prohibit all addresses that use its advertised prefix and have an IID starting with 0x0300 (4rd Tag), by using Duplicate Address Detection [RFC4862].
- R-11: A CE that is assigned the unspecified IPv4 address (see Section 4.4) MUST use, for packets tunneled between itself and the Domain NAT64+, addresses as detailed in Figure 6: part (a) for its IPv6 source, and part (b) as IPv6 destinations that depend on IPv4 destinations. A NAT64+, being NAT64 conforming [RFC6146], MUST accept IPv6 packets whose destination conforms to Figure 6(b) (4rd Tag instead of "u" and 0x00 octets). In its Binding Information Base, it MUST remember whether a mapping was created with a "u" or 4rd-tag destination. In the IPv4-to-IPv6 direction, it MUST use 4rd tunneling, with source address conforming to Figure 6(b), when using a mapping that was created with a 4rd-tag destination.

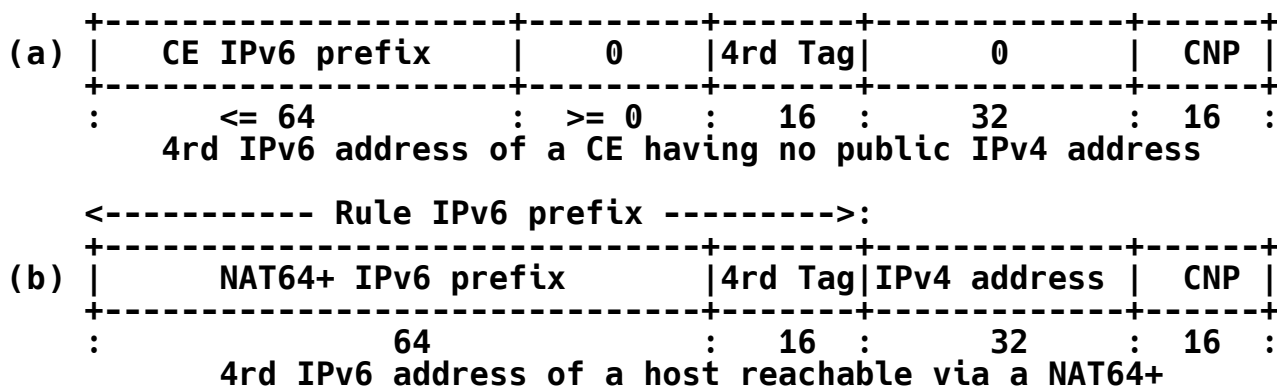


Figure 6: Rules for CE and NAT64+

- R-12: For anti-spoofing protection, CEs and BRs **MUST** check that the IPv6 source address of each received Tunnel packet is that which, according to R-9, is derived from the source 4rd IPv4 address. For this, the IPv4 address used to obtain the source 4rd IPv4 address is that embedded in the IPv6 source address (in its bits 80-111). (This verification is needed because IPv6 ingress filtering [RFC3704] applies only to IPv6 prefixes, without any guarantee that Tunnel packets are built as specified in R-9.)
- R-13: For additional protection against packet corruption at a link layer that might be undetected at this layer during Domain traversal, CEs and BRs **SHOULD** verify that source and destination IPv6 addresses have not been modified. This can be done by checking that they remain checksum neutral (see the Note above regarding the CNP).

4.6. Fragmentation Processing

4.6.1. Fragmentation at Domain Entry

- R-14: If an IPv4 packet enters a CE or BR with a size such that the derived Tunnel packet would be longer than the Domain PMTU, the packet has to be either discarded or fragmented. The Domain-entry node **MUST** discard it if the packet has DF = 1, with an ICMP error message returned to the source. It **MUST** fragment it otherwise, with the payload of each fragment not exceeding PMTU - 48. The first fragment has its offset equal to the received offset. Subsequent fragments have offsets increased by the lengths of the payloads of previous fragments. Functionally, fragmentation is supposed to be done in IPv4 before applying reversible header translation to each fragment; see Section 4.3.

4.6.2. Ports of Fragments Addressed to Shared-Address CEs

Because ports are available only in the first fragments of IPv4 fragmented packets, a BR needs a mechanism to send to the right shared-address CEs all fragments of fragmented packets.

For this, a BR MAY systematically reassemble fragmented IPv4 packets before tunneling them. But this consumes large memory space, creates opportunities for denial-of-service-attacks, and can significantly increase forwarding delays. This is the reason for the following requirement:

R-15: BRs SHOULD support an algorithm whereby received IPv4 packets can be forwarded on the fly. The following is an example of such an algorithm:

- (1) At BR initialization, if at least one CE Mapping rule deals with one or more shared public IPv4 addresses (i.e., length of Rule IPv4 prefix + EA-bits length > 32), the BR initializes an empty "IPv4 packet table" whose entries have the following items:
 - IPv4 source
 - IPv4 destination
 - IPv4 identification
 - Destination port
- (2) When the BR receives an IPv4 packet whose matching Mapping rule deals with one or more shared public IPv4 addresses (i.e., length of Rule IPv4 prefix + EA-bits length > 32), the BR searches the table for an entry whose IPv4 source, IPv4 destination, and IPv4 identification are those of the received packet. The BR then performs actions as detailed in Table 5, depending on which conditions hold.

+-----+-----+-----+-----+-----+-----+-----+-----+-----+								
- CONDITIONS -								
First Fragment (offset = 0)	Y	Y	Y	Y	N	N	N	N
Last fragment (MF = 0)	Y	Y	N	N	Y	Y	N	N
An entry has been found	Y	N	Y	N	Y	N	Y	N

- RESULTING ACTIONS -								
Create a new entry	-	-	-	X	-	-	-	-
Use port of the entry	-	-	-	-	X	-	X	-
Update port of the entry	-	-	X	-	-	-	-	-
Delete the entry	X	-	-	-	X	-	-	-
Forward the packet	X	X	X	X	X	-	X	-
+-----+-----+-----+-----+-----+-----+-----+-----+-----+								

Table 5: BR Actions

- (3) The BR performs garbage collection for table entries that remain unchanged for longer than some limit. This limit, which is normally longer than the maximum time normally needed to reassemble a packet, is not critical. It should not, however, be longer than 15 seconds [RFC791].

- R-16: For the above algorithm to be effective, CEs that are assigned shared public IPv4 addresses MUST NOT interleave fragments of several fragmented packets.
- R-17: CEs that are assigned IPv4 prefixes and are in nodes that route public IPv4 addresses rather than only using NAT44s MUST have the same behavior as that described just above for BRs.

4.6.3. Packet Identifications from Shared-Address CEs

When packets go from CEs that share the same IPv4 address to a common destination, a precaution is needed to guarantee that packet identifications set by sources are different. Otherwise, packet reassembly at the destination could be confused because it is based only on source IPv4 address and Identification. The probability of such confusing situations may in theory be very low, but a safe solution is needed in order to avoid creating new attack opportunities.

R-18: A CE that is assigned a shared public IPv4 address **MUST** only use packet identifications that have the CE PSID in their bits 0 to PSID length - 1.

R-19: A BR or a CE that receives a packet from a shared-address CE **MUST** check that bits 0 to PSID length - 1 of their packet identifications are equal to the PSID found in the source 4rd IPv4 address.

4.7. TOS and Traffic Class Processing

IPv4 TOS and IPv6 traffic class have the same semantic, that of the differentiated services field, or DS field, specified in [RFC2474] and [RFC6040]. Their first 6 bits contain a differentiated services codepoint (DSCP), and their last 2 bits can convey explicit congestion notifications (ECNs), which both may evolve during Domain traversal. [RFC2983] discusses how the DSCP can be handled by tunnel endpoints. The Tunnel Traffic Class option provides permission to ignore DS-field evolutions occurring during Domain traversal, if the desired behavior is that of generic tunnels conforming to [RFC2473].

R-20: Unless the Tunnel Traffic Class option is configured for the Domain, BRs and CEs **MUST** copy the IPv4 TOS into the IPv6 traffic class at Domain entry and copy back the IPv6 traffic class into the IPv4 TOS at Domain exit.

R-21: If the Tunnel Traffic Class option is configured for a Domain, BRs and CEs **MUST** at Domain entry take the configured Tunnel Traffic Class as the IPv6 traffic class and copy the received IPv4 TOS into the IPv4_TOS of the fragment header (Figure 3). At Domain exit, they **MUST** copy back the IPv4_TOS of the fragment header into the IPv4 TOS.

4.8. Tunnel-Generated ICMPv6 Error Messages

If a Tunnel packet is discarded on its way across a 4rd domain because of an unreachable destination, an ICMPv6 error message is returned to the IPv6 source. For the IPv4 source of the discarded packet to be informed of packet loss, the ICMPv6 message has to be converted into an ICMPv4 message.

R-22: If a CE or BR receives an ICMPv6 error message [RFC4443], it MUST synthesize an ICMPv4 error packet [RFC792]. This packet MUST contain the first 8 octets of the discarded packet's IP payload. The reserved IPv4 dummy address (192.0.0.8/32; see Section 6) MUST be used as its source address.

As described in [RFC6145], ICMPv6 Type = 1 and Code = 0 (Destination Unreachable, No route to destination) MUST be translated into ICMPv4 Type = 3 and Code = 0 (Destination Unreachable, Net unreachable), and ICMPv6 Type = 3 and Code = 0 (Time Exceeded, Hop limit exceeded in transit) MUST be translated into ICMPv4 Type = 11 and Code = 0 (Time Exceeded, time to live exceeded in transit).

4.9. Provisioning 4rd Parameters to CEs

Domain parameters listed in Section 4.2 are subject to the following constraints:

R-23: Each Domain MUST have a BR Mapping rule and/or a NAT64+ Mapping rule. The BR Mapping rule is only used by CEs that are assigned public IPv4 addresses, shared or not. The NAT64+ Mapping rule is only used by CEs that are assigned the unspecified IPv4 address (Section 4.4) and therefore need an ISP NAT64 to reach IPv4 destinations.

R-24: Each CE and each BR MUST support up to 32 Mapping rules.

Support for up to 32 Mapping rules ensures that independently acquired CEs and BR nodes can always interwork.

ISPs that need Mapping rules for more IPv4 prefixes than this number SHOULD split their networks into multiple Domains. Communication between these domains can be done in IPv4 or by some other implementation-dependent, but equivalent, means.

- R-25:** For mesh topologies, where CE-CE paths don't go via BRs, all Mapping rules of the Domain **MUST** be sent to all CEs. For hub-and-spoke topologies, where all CE-CE paths go via BRs, each CE **MAY** be sent only the BR Mapping rule of the Domain plus, if different, the CE Mapping rule that applies to its CE IPv6 prefix.
- R-26:** In a Domain where the chosen topology is hub-and-spoke, all CEs **MUST** have IPv6 prefixes that match a CE Mapping rule. (Otherwise, packets sent to CEs whose IPv6 prefixes would match only the BR Mapping rule would, with longest-match selected routes, be routed directly to these CEs. This would be contrary to the hub-and-spoke requirement.)
- R-27:** CEs **MUST** be able to acquire parameters of 4rd domains (Section 4.2) in DHCPv6 [RFC3315]. Formats of DHCPv6 options to be used are detailed in Figures 7, 8, and 9, with field values specified after each figure.

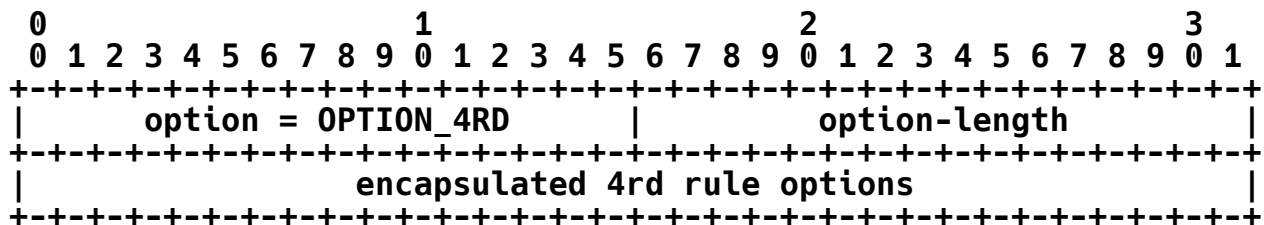


Figure 7: DHCPv6 Option for 4rd

- o option code: 97, OPTION_4RD (see Section 6)
- o option-length: the length of encapsulated options, in octets

- o encapsulated 4rd rule options: The OPTION_4RD DHCPv6 option contains at least one encapsulated OPTION_4RD_MAP_RULE option and a maximum of one encapsulated OPTION_4RD_NON_MAP_RULE option. Since DHCP servers normally send whatever options the operator configures, operators are advised to configure these options appropriately. DHCP servers MAY check to see that the configuration follows these rules and notify the operator in an implementation-dependent manner if the settings for these options aren't valid. The length of encapsulated options is in octets.



Figure 8: Encapsulated Option for Mapping-Rule Parameters

- o option code: 98, encapsulated OPTION_4RD_MAP_RULE option (see Section 6)
- o option-length: 20
- o prefix4-len: number of bits of the Rule IPv4 prefix
- o prefix6-len: number of bits of the Rule IPv6 prefix
- o ea-len: EA-bits length
- o W: WKP authorized, = 1 if set
- o rule-ipv4-prefix: Rule IPv4 prefix, left-aligned
- o rule-ipv6-prefix: Rule IPv6 prefix, left-aligned

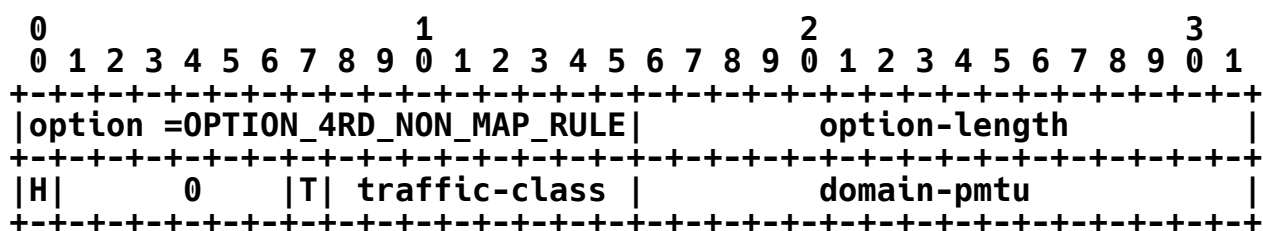


Figure 9: Encapsulated Option for Non-Mapping-Rule Parameters of 4rd Domains

- o option code: 99, encapsulated OPTION_4RD_NON_MAP_RULE option (see Section 6)
- o option-length: 4
- o H: Hub-and-spoke topology (= 1 if Yes)
- o T: Traffic Class flag (= 1 if a Tunnel Traffic Class is provided)
- o traffic-class: Tunnel Traffic Class
- o domain-pmtu: Domain PMTU (at least 1280 octets)

Means other than DHCPv6 that may prove useful to provide 4rd parameters to CEs are off-scope for this document. The same or similar parameter formats would, however, be recommended to facilitate training and operation.

5. Security Considerations

Spoofing attacks

With IPv6 ingress filtering in effect in the Domain [RFC3704], as required in Section 3 (Figure 1 in particular), and with consistency checks between 4rd IPv4 and IPv6 addresses (Section 4.5), no spoofing opportunity in IPv4 is introduced by 4rd: being able to use as source IPv6 address only one that has been allocated to him, a customer can only provide as source 4rd IPv4 address that which derives this IPv6 address according to Section 4.5, i.e., one that his ISP has allocated to him.

Routing loop attacks

Routing loop attacks that may exist in some "automatic tunneling" scenarios are documented in [RFC6324]. No opportunities for routing loop attacks have been identified with 4rd.

Fragmentation-related attacks

As discussed in Section 4.6, each BR of a Domain that assigns shared public IPv4 addresses should maintain a dynamic table of fragmented packets that go to these shared-address CEs.

This leaves BRs vulnerable to denial-of-service attacks from hosts that would send very large numbers of first fragments and would never send last fragments having the same packet identifications. This vulnerability is inherent in IPv4 address sharing, be it static or dynamic. Compared to what it is with algorithms that reassemble IPv4 packets in BRs, it is, however, significantly mitigated by the algorithm provided in Section 4.6.2, as that algorithm uses much less memory space.

6. IANA Considerations

IANA has allocated the following:

- o Encapsulated options `OPTION_4RD` (97), `OPTION_4RD_MAP_RULE` (98), and `OPTION_4RD_NON_MAP_RULE` (99). These options have been recorded in the option code space of the "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)" registry. See Section 4.9 of this document and Section 24.3 of [RFC3315]).

Value	Description	Reference
97	<code>OPTION_4RD</code>	this document
98	<code>OPTION_4RD_MAP_RULE</code>	this document
99	<code>OPTION_4RD_NON_MAP_RULE</code>	this document

- o Reserved IPv4 address 192.0.0.8/32 to be used as the "IPv4 dummy address" (Section 4.8).

7. Relationship with Previous Works

The present specification has been influenced by many previous IETF drafts, in particular those accessible at <http://tools.ietf.org/html/draft-xxxx>, where "xxxx" refers to the following (listed in order, by date of their first versions):

- o bagnulo-behave-nat64 (2008-06-10)
- o xli-behave-ivi (2008-07-06)
- o despres-sam-scenarios (2008-09-28)
- o boucadair-port-range (2008-10-23)

- o ymbk-aplusp (2008-10-27)
- o xli-behave-divi (2009-10-19)
- o thaler-port-restricted-ip-issues (2010-02-28)
- o cui-softwire-host-4over6 (2010-07-06)
- o dec-stateless-4v6 (2011-03-05)
- o matsushima-v6ops-transition-experience (2011-03-07)
- o despres-intarea-4rd (2011-03-07)
- o deng-aplusp-experiment-results (2011-03-07)
- o operators-softwire-stateless-4v6-motivation (2011-05-05)
- o xli-behave-divi-pd (2011-07-04)
- o murakami-softwire-4rd (2011-07-04)
- o murakami-softwire-4v6-translation (2011-07-04)
- o despres-softwire-4rd-addmapping (2011-08-19)
- o boucadair-softwire-stateless-requirements (2011-09-08)
- o chen-softwire-4v6-add-format (2011-10-12)
- o mawatari-softwire-464xlat (2011-10-16)
- o mdt-softwire-map-dhcp-option (2011-10-24)
- o mdt-softwire-mapping-address-and-port (2011-10-24)
- o mdt-softwire-map-translation (2012-01-10)
- o mdt-softwire-map-encapsulation (2012-01-27)

8. References

8.1. Normative References

- [RFC791] Postel, J., "Internet Protocol", STD 5, RFC 791, DOI 10.17487/RFC0791, September 1981, <<http://www.rfc-editor.org/info/rfc791>>.
- [RFC792] Postel, J., "Internet Control Message Protocol", STD 5, RFC 792, DOI 10.17487/RFC0792, September 1981, <<http://www.rfc-editor.org/info/rfc792>>.
- [RFC793] Postel, J., "Transmission Control Protocol", STD 7, RFC 793, DOI 10.17487/RFC0793, September 1981, <<http://www.rfc-editor.org/info/rfc793>>.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", BCP 14, RFC 2119, DOI 10.17487/RFC2119, March 1997, <<http://www.rfc-editor.org/info/rfc2119>>.
- [RFC2460] Deering, S. and R. Hinden, "Internet Protocol, Version 6 (IPv6) Specification", RFC 2460, DOI 10.17487/RFC2460, December 1998, <<http://www.rfc-editor.org/info/rfc2460>>.
- [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, DOI 10.17487/RFC2474, December 1998, <<http://www.rfc-editor.org/info/rfc2474>>.
- [RFC2983] Black, D., "Differentiated Services and Tunnels", RFC 2983, DOI 10.17487/RFC2983, October 2000, <<http://www.rfc-editor.org/info/rfc2983>>.
- [RFC3315] Droms, R., Ed., Bound, J., Volz, B., Lemon, T., Perkins, C., and M. Carney, "Dynamic Host Configuration Protocol for IPv6 (DHCPv6)", RFC 3315, DOI 10.17487/RFC3315, July 2003, <<http://www.rfc-editor.org/info/rfc3315>>.
- [RFC4291] Hinden, R. and S. Deering, "IP Version 6 Addressing Architecture", RFC 4291, DOI 10.17487/RFC4291, February 2006, <<http://www.rfc-editor.org/info/rfc4291>>.

- [RFC4443] Conta, A., Deering, S., and M. Gupta, Ed., "Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification", RFC 4443, DOI 10.17487/RFC4443, March 2006, <<http://www.rfc-editor.org/info/rfc4443>>.
- [RFC4862] Thomson, S., Narten, T., and T. Jinmei, "IPv6 Stateless Address Autoconfiguration", RFC 4862, DOI 10.17487/RFC4862, September 2007, <<http://www.rfc-editor.org/info/rfc4862>>.
- [RFC4925] Li, X., Ed., Dawkins, S., Ed., Ward, D., Ed., and A. Durand, Ed., "Softwire Problem Statement", RFC 4925, DOI 10.17487/RFC4925, July 2007, <<http://www.rfc-editor.org/info/rfc4925>>.
- [RFC5082] Gill, V., Heasley, J., Meyer, D., Savola, P., Ed., and C. Pignataro, "The Generalized TTL Security Mechanism (GTSM)", RFC 5082, DOI 10.17487/RFC5082, October 2007, <<http://www.rfc-editor.org/info/rfc5082>>.
- [RFC6040] Briscoe, B., "Tunnelling of Explicit Congestion Notification", RFC 6040, DOI 10.17487/RFC6040, November 2010, <<http://www.rfc-editor.org/info/rfc6040>>.

8.2. Informative References

- [NAT444] Yamagata, I., Shirasaki, Y., Nakagawa, A., Yamaguchi, J., and H. Ashida, "NAT444", Work in Progress, draft-shirasaki-nat444-06, July 2012.
- [RFC768] Postel, J., "User Datagram Protocol", STD 6, RFC 768, DOI 10.17487/RFC0768, August 1980, <<http://www.rfc-editor.org/info/rfc768>>.
- [RFC1191] Mogul, J. and S. Deering, "Path MTU discovery", RFC 1191, DOI 10.17487/RFC1191, November 1990, <<http://www.rfc-editor.org/info/rfc1191>>.
- [RFC1918] Rekhter, Y., Moskowitz, B., Karrenberg, D., de Groot, G., and E. Lear, "Address Allocation for Private Internets", BCP 5, RFC 1918, DOI 10.17487/RFC1918, February 1996, <<http://www.rfc-editor.org/info/rfc1918>>.

- [RFC2473] Conta, A. and S. Deering, "Generic Packet Tunneling in IPv6 Specification", RFC 2473, DOI 10.17487/RFC2473, December 1998, <<http://www.rfc-editor.org/info/rfc2473>>.
- [RFC3022] Srisuresh, P. and K. Egevang, "Traditional IP Network Address Translator (Traditional NAT)", RFC 3022, DOI 10.17487/RFC3022, January 2001, <<http://www.rfc-editor.org/info/rfc3022>>.
- [RFC3704] Baker, F. and P. Savola, "Ingress Filtering for Multihomed Networks", BCP 84, RFC 3704, DOI 10.17487/RFC3704, March 2004, <<http://www.rfc-editor.org/info/rfc3704>>.
- [RFC3828] Larzon, L-A., Degermark, M., Pink, S., Jonsson, L-E., Ed., and G. Fairhurst, Ed., "The Lightweight User Datagram Protocol (UDP-Lite)", RFC 3828, DOI 10.17487/RFC3828, July 2004, <<http://www.rfc-editor.org/info/rfc3828>>.
- [RFC4271] Rekhter, Y., Ed., Li, T., Ed., and S. Hares, Ed., "A Border Gateway Protocol 4 (BGP-4)", RFC 4271, DOI 10.17487/RFC4271, January 2006, <<http://www.rfc-editor.org/info/rfc4271>>.
- [RFC4821] Mathis, M. and J. Heffner, "Packetization Layer Path MTU Discovery", RFC 4821, DOI 10.17487/RFC4821, March 2007, <<http://www.rfc-editor.org/info/rfc4821>>.
- [RFC4961] Wing, D., "Symmetric RTP / RTP Control Protocol (RTCP)", BCP 131, RFC 4961, DOI 10.17487/RFC4961, July 2007, <<http://www.rfc-editor.org/info/rfc4961>>.
- [RFC5595] Fairhurst, G., "The Datagram Congestion Control Protocol (DCCP) Service Codes", RFC 5595, DOI 10.17487/RFC5595, September 2009, <<http://www.rfc-editor.org/info/rfc5595>>.
- [RFC5969] Townsley, W. and O. Troan, "IPv6 Rapid Deployment on IPv4 Infrastructures (6rd) -- Protocol Specification", RFC 5969, DOI 10.17487/RFC5969, August 2010, <<http://www.rfc-editor.org/info/rfc5969>>.
- [RFC6145] Li, X., Bao, C., and F. Baker, "IP/ICMP Translation Algorithm", RFC 6145, DOI 10.17487/RFC6145, April 2011, <<http://www.rfc-editor.org/info/rfc6145>>.

- [RFC6146] Bagnulo, M., Matthews, P., and I. van Beijnum, "Stateful NAT64: Network Address and Protocol Translation from IPv6 Clients to IPv4 Servers", RFC 6146, DOI 10.17487/RFC6146, April 2011, <<http://www.rfc-editor.org/info/rfc6146>>.
- [RFC6324] Nakibly, G. and F. Templin, "Routing Loop Attack Using IPv6 Automatic Tunnels: Problem Statement and Proposed Mitigations", RFC 6324, DOI 10.17487/RFC6324, August 2011, <<http://www.rfc-editor.org/info/rfc6324>>.
- [RFC6346] Bush, R., Ed., "The Address plus Port (A+P) Approach to the IPv4 Address Shortage", RFC 6346, DOI 10.17487/RFC6346, August 2011, <<http://www.rfc-editor.org/info/rfc6346>>.
- [RFC6437] Amante, S., Carpenter, B., Jiang, S., and J. Rajahalme, "IPv6 Flow Label Specification", RFC 6437, DOI 10.17487/RFC6437, November 2011, <<http://www.rfc-editor.org/info/rfc6437>>.
- [RFC6535] Huang, B., Deng, H., and T. Savolainen, "Dual-Stack Hosts Using "Bump-in-the-Host" (BIH)", RFC 6535, DOI 10.17487/RFC6535, February 2012, <<http://www.rfc-editor.org/info/rfc6535>>.
- [RFC6887] Wing, D., Ed., Cheshire, S., Boucadair, M., Penno, R., and P. Selkirk, "Port Control Protocol (PCP)", RFC 6887, DOI 10.17487/RFC6887, April 2013, <<http://www.rfc-editor.org/info/rfc6887>>.
- [RFC7136] Carpenter, B. and S. Jiang, "Significance of IPv6 Interface Identifiers", RFC 7136, DOI 10.17487/RFC7136, February 2014, <<http://www.rfc-editor.org/info/rfc7136>>.
- [Solutions-4v6] Boucadair, M., Ed., Matsushima, S., Lee, Y., Bonness, O., Borges, I., and G. Chen, "Motivations for Carrier-side Stateless IPv4 over IPv6 Migration Solutions", Work in Progress, draft-ietf-softwire-stateless-4v6-motivation-05, November 2012.

Appendix A. Textual Representation of Mapping Rules

In the sections that follow, each Mapping rule will be represented as follows, using 0bXXX to represent binary number XXX; square brackets ("[]") indicate optional items:

```
{Rule IPv4 prefix, EA-bits length, Rule IPv6 prefix
  [, WKPs authorized]}
```

EXAMPLES:

```
{0.0.0.0/0, 32, 2001:db8:0:1:300::/80}
    a BR Mapping rule
{198.16.0.0/14, 22, 2001:db8:4000::/34}
    a CE Mapping rule
{0.0.0.0/0, 32, 2001:db8:0:1::/80}
    a NAT64+ Mapping rule
{198.16.0.0/14, 22, 2001:db8:4000::/34, Yes}
    a CE Mapping rule
    and hub-and-spoke topology
```

Appendix B. Configuring Multiple Mapping Rules

As far as Mapping rules are concerned, the simplest deployment model is that in which the Domain has only one rule (the BR Mapping rule). To assign an IPv4 address to a CE in this model, an IPv6 /112 is assigned to it, comprising the BR /64 prefix, the 4rd Tag, and the IPv4 address. However, this model has the following limitations: (1) shared IPv4 addresses are not supported; (2) IPv6 prefixes used for 4rd are too long to also be used for native IPv6 addresses; (3) if the IPv4 address space of the ISP is split with many disjoint IPv4 prefixes, the IPv6 routing plan must be as complex as an IPv4 routing plan based on these prefixes.

With more Mapping rules, CE prefixes used for 4rd can be those used for native IPv6. How to choose CE Mapping rules for a particular deployment does not need to be standardized.

The following is only a particular pragmatic approach that can be used for various deployment scenarios. It is applied in some of the use cases that follow.

- (1) Select a "Common_IPv6_prefix" that will appear at the beginning of all 4rd CE IPv6 prefixes.
- (2) Choose all IPv4 prefixes to be used, and assign one of them to each CE Mapping rule *i*.

(3) For each CE Mapping rule i , do the following:

- A. Choose the length of its Rule IPv6 prefix (possibly the same for all CE Mapping rules).
- B. Determine its $\text{PSID_length}(i)$. A CE Mapping rule that assigns shared addresses with a sharing ratio of 2^{K_i} has $\text{PSID_length} = K_i$. A CE Mapping rule that assigns IPv4 prefixes of length $L < 32$ is considered to have a negative PSID_length ($\text{PSID_length} = L - 32$).
- C. Derive EA-bits $\text{length}(i) = 32 - L(\text{Rule IPv4 prefix}(i)) + \text{PSID_length}(i)$.
- D. Derive the length of $\text{Rule_code}(i)$, the prefix to be appended to the common prefix to get the Rule IPv6 prefix of rule i :
$$\begin{aligned} L(\text{Rule_code}(i)) = & L(\text{CE IPv6 prefix}(i)) \\ & - L(\text{Common IPv6 prefix}) \\ & - (32 - L(\text{Rule IPv4 prefix}(i))) \\ & - \text{PSID_length}(i) \end{aligned}$$
- E. Derive $\text{Rule_code}(i)$ with the following constraints: (1) its length is $L(\text{Rule_code}(i))$; (2) it does not overlap with any of the previously obtained Rule_codes (for instance, 010 and 01011 do overlap, while 00, 011, and 010 do not); (3) it has the lowest possible value as a fractional binary number (for instance, $0100 < 10 < 11011 < 111$). Thus, rules whose Rule_code lengths are 4, 3, 5, and 2 give Rule_codes 0000, 001, 00010, and 01.

Applying the principles of Appendix B with $L(\text{Common_IPv6_prefix}) = 36$, $L(\text{PSID}) = 2$ for all rules, and $L(\text{CE IPv6 prefix}(i)) = 56$ for all rules, Rule_codes and Rule IPv6 prefixes are as follows:

CE Rule IPv4 prefix	EA bits length	Rule-Code length	Code (binary)	CE Rule IPv6 prefix
192.8.0.0/15	19	1	0	2001:db8:0::/37
192.4.0.0/16	18	2	10	2001:db8:800::/38
192.2.0.0/16	18	2	11	2001:db8:c00::/38

Mapping rules are then the following:

```
{192.8.0.0/15, 19, 2001:0db8:0000::/37}
{192.4.0.0/16, 18, 2001:0db8:0800::/38}
{192.2.0.0/16, 18, 2001:0db8:0c00::/38}
{0.0.0.0/0, 32, 2001:0db8:0000:0001:300::/80}
```

The CE whose IPv6 prefix is, for example, 2001:db8:0bbb:bb00::/56 derives its IPv4 address and its port set as follows (Section 4.4):

```
CE IPv6 prefix      : 2001:0db8:0bbb:bb00::/56
Rule IPv6 prefix(i) : 2001:0db8:0800::/38 (longest match)
EA-bits length(i)   : 18
EA bits              : 0b11 1011 1011 1011 1011
Rule IPv4 prefix(i) : 0b1100 0000 0000 0100 (192.4.0.0/16)
IPv4 address         : 0b1100 0000 0000 0100 1110 1110 1110 1110
                     : 192.4.238.238
PSID                 : 0b11
Ports                : 0bYYYY 11XX XXXX XXXX
                     : with YYYY > 0, and X...X any value
```


An IPv4 packet sent to address 192.4.238.238 and port 7777 is tunneled to the IPv6 address obtained as follows (Section 4.5):

```

IPv4 address      : 192.4.238.238 (0xc004 eeee)
                  : 0b1100 0000 0000 0100 1110 1110 1110 1110
Rule IPv4 prefix(i): 192.4.0.0/16  (longest match)
                  : 0b1100 0000 0000 0100
IPv4 suffix(i)    : 0b1110 1110 1110 1110
EA-bits length(i) : 18
PSID length(i)    : 2  (= 16 + 18 - 32)
Port field        : 0b 0001 1110 0110 0001 (7777)
PSID              : 0b11
Rule IPv6 prefix(i): 2001:0db8:0800::/38
CE IPv6 prefix    : 2001:0db8:0bbb:bb00::/56
IPv6 address      : 2001:0db8:0bbb:bb00:300:c004:eeee:YYYY
                  : with YYYY = the computed CNP

```

C.2. With Some CEs behind Third-Party Router CPEs

We now consider an ISP that has the same need as the ISP described in the previous section, except that (1) instead of using only its own IPv6 infrastructure, it uses that of a third-party provider and (2) some of its customers use this provider's Customer Premises Equipment (CPEs) so that they can use specific services offered by the provider. In these CPEs, a non-zero index is used to route IPv6 packets to the physical port to which CEs are attached, say 0x2. Each such CPE delegates to the CE nodes the customer-site IPv6 prefix followed by this index.

The ISP is supposed to have the same IPv4 prefixes as those in the previous use case -- 192.8.0.0/15, 192.4.0.0/16, and 192.2.0.0/16 -- and to use the same Common_IPv6_prefix, 2001:db8:0::/36.

We also assume that only a minority of customers use third-party CPEs, so that it is sufficient to use only one of the two /16s for them.

Mapping rules are then (see Appendix C.1):

```

{192.8.0.0/15, 19, 2001:0db8:0000::/37}
{192.4.0.0/16, 18, 2001:0db8:0800::/38}
{192.2.0.0/16, 18, 2001:0db8:0c00::/38}
{0.0.0.0/0,    32, 2001:0db8:0000:0001:300::/80}

```

CEs that are behind third-party CPEs derive their own IPv4 addresses and port sets as described in Appendix C.1.

In a BR, and also in a CE if the topology is mesh, the IPv6 address that is derived from IPv4 address 192.4.238.238 and port 7777 is obtained as described in the previous section, except for the last two steps, which are modified as follows:

```

IPv4 address      : 192.4.238.238 (0xc004 eeee)
                  : 0b1100 0000 0000 0100 1110 1110 1110 1110
Rule IPv4 prefix(i): 192.4.0.0/16  (longest match)
                  : 0b1100 0000 0000 0100
IPv4 suffix(i)    : 0b1110 1110 1110 1110
EA-bits length(i) : 18
PSID length(i)    : 2  (= 16 + 18 - 32)
Port field        : 0b 0001 1110 0110 0001 (7777)
PSID              : 0b11
Rule IPv6 prefix(i): 2001:0db8:0800::/38
CE IPv6 prefix    : 2001:0db8:0bbb:bb00::/60
IPv6 address      : 2001:0db8:0bbb:bb00:300:192.4.238.238:YYYY
                  with YYYY = the computed CNP

```

Appendix D. Replacing Dual-Stack Routing with IPv6-Only Routing

In this use case, we consider an ISP that offers IPv4 service with public addresses individually assigned to its customers. It also offers IPv6 service, as it has deployed dual-stack routing. Because it provides its own CPEs to customers, it can upgrade all of its CPEs to support 4rd. It wishes to take advantage of this capability to replace dual-stack routing with IPv6-only routing, without changing any IPv4 address or IPv6 prefix.

For this, the ISP can use the single-rule model described at the beginning of Appendix B. If the prefix routed to BRs is chosen to start with 2001:db8:0:1::/64, this rule is:

```
{0.0.0.0/0, 32, 2001:db8:0:1:300::/80}
```

All that is needed in the network before disabling IPv4 routing is the following:

- o In all routers, where there is an IPv4 route toward x.x.x.x/n, add a parallel route toward 2001:db8:0:1:300:x.x.x.x::/(80+n).
- o Where IPv4 address x.x.x.x was assigned to a CPE, now delegate IPv6 prefix 2001:db8:0:1:300:x.x.x.x::/112.

NOTE: In parallel with this deployment, or after it, shared IPv4 addresses can be assigned to IPv6 customers. It is sufficient that IPv4 prefixes used for this be different from those used for exclusive-address assignments. Under this constraint, Mapping rules can be set up according to the same principles as those described in Appendix C.

Appendix E. Adding IPv6 and 4rd Service to a Net-10 Network

In this use case, we consider an ISP that has only deployed IPv4, possibly because some of its network devices are not yet IPv6 capable. Because it did not have enough IPv4 addresses, it has assigned private IPv4 addresses [RFC1918] to customers, say 10.x.x.x. It thus supports up to 2^{24} customers (a "Net-10" network, using the NAT444 model [NAT444]).

Now, it wishes to offer IPv6 service without further delay, using 6rd [RFC5969]. It also wishes to offer incoming IPv4 connectivity to its customers with a simpler solution than that provided by the Port Control Protocol (PCP) [RFC6887].

This appendix describes an example that adds IPv6 (using 6rd) and 4rd services to the "Net-10" private IPv4 network.

The IPv6 prefix to be used for 6rd is supposed to be 2001:db8::/32, and the public IPv4 prefix to be used for shared addresses is supposed to be 198.16.0.0/16 (0xc610). The resulting sharing ratio is $2^{24} / 2^{(32 - 16)} = 256$, giving a PSID length of 8.

The ISP installs one or several BRs at its border to the public IPv4 Internet. They support 6rd, and 4rd above it. The BR prefix /64 is supposed to be that which is derived from IPv4 address 10.0.0.1 (i.e., 2001:db8:0:100:/64).

In accordance with [RFC5969], 6rd BRs are configured with the following parameters: IPv4MaskLen = 8; 6rdPrefix = 2001:db8::/32; 6rdBRIPv4Address = 192.168.0.1 (0xc0a80001).

4rd Mapping rules are then the following:

```
{198.16.0.0/16, 24, 2001:db8:0:0:300::/80}
{0.0.0.0/0,      32, 2001:db8:0:100:300:/80,}
```

Any customer device that supports 4rd in addition to 6rd can then use its assigned shared IPv4 address with 240 assigned ports.

If its NAT44 supports port forwarding to provide incoming IPv4 connectivity (statically, or dynamically with Universal Plug and Play (UPnP) and/or the NAT Port Mapping Protocol (NAT-PMP)), it can use it with ports of the assigned port set (a possibility that does not exist in Net-10 networks without 4rd/6rd).

Acknowledgements

This specification has benefited over several years from independent proposals, questions, comments, constructive suggestions, and useful criticisms from numerous IETF contributors. The authors would like to express recognition of all of these contributors, and especially the following, in alphabetical order by their first names: Behcet Sarikaya, Bing Liu, Brian Carpenter, Cameron Byrne, Congxiao Bao, Dan Wing, Derek Atkins, Erik Kline, Francis Dupont, Gabor Bajko, Hui Deng, Jacni Quin (who was an active coauthor of some earlier versions of this specification), James Huang, Jan Zorz, Jari Arkko, Kathleen Moriarty, Laurent Toutain, Leaf Yeh, Lorenzo Colitti, Marcello Bagnulo, Mark Townsley, Mohamed Boucadair, Nejc Skoberne, Olaf Maennel, Ole Troan, Olivier Vautrin, Peng Wu, Qiong Sun, Rajiv Asati, Ralph Droms, Randy Bush, Satoru Matsushima, Simon Perreault, Stuart Cheshire, Suresh Krishnan, Ted Lemon, Teemu Savolainen, Tetsuya Murakami, Tina Tsou, Tomek Mrugalski, Washam Fan, Wojciech Dec, Xiaohong Deng, Xing Li, and Yu Fu.

Authors' Addresses

Remi Despres
RD-IPtech
3 rue du President Wilson
Levallois
France

Email: despres.remi@laposte.net

Sheng Jiang (editor)
Huawei Technologies Co., Ltd
Q14, Huawei Campus, No. 156 BeiQing Road
Hai-Dian District, Beijing 100095
China

Email: jiangsheng@huawei.com

Reinaldo Penno
Cisco Systems, Inc.
170 West Tasman Drive
San Jose, CA 95134
United States

Email: repenno@cisco.com

Yiu Lee
Comcast
One Comcast Center
Philadelphia, PA 19103
United States

Email: yiulee@cable.comcast.com

Gang Chen
China Mobile
29, Jinrong Avenue
Xicheng District, Beijing 100033
China

Email: phdgang@gmail.com, chengang@chinamobile.com

Maoke Chen (a.k.a. Noriyuki Arai)
BBIX, Inc.
Tokyo Shiodome Building, Higashi-Shimbashi 1-9-1
Minato-ku, Tokyo 105-7310
Japan

Email: maoke@bbix.net