

Routing over Large Clouds Working Group
Luciani
INTERNET-DRAFT
Networks)
<[draft-ietf-rolc-nhrp-08.txt](#)>
Katz

James V.
(Bay
Dave
(cisco
David
(Core Competence,
Bruce
(cisco
Expires December 1996

Systems)
Piscitello
Inc.)
Cole
Systems)

NBMA Next Hop Resolution Protocol (NHRP)

Status of this Memo

This document is an Internet-Draft. 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.''

To learn the current status of any Internet-Draft, please check the ``[id-abstracts.txt](#)'' listing contained in the Internet-Drafts Shadow

Directories on [ds.internic.net](#) (US East Coast), [nic.nordu.net](#) (Europe), [ftp.isi.edu](#) (US West Coast), or [munnari.oz.au](#) (Pacific Rim).

Abstract

This document describes the NBMA Next Hop Resolution Protocol (NHRP).

NHRP can be used by a source station (host or router) connected to a Non-Broadcast, Multi-Access (NBMA) subnetwork to determine the internetworking layer address and NBMA subnetwork addresses of the "NBMA next hop" towards a destination station. If the destination is

connected to the NBMA subnetwork, then the NBMA next hop is the destination station itself. Otherwise, the NBMA next hop is the egress router from the NBMA subnetwork that is "nearest" to the destination station. NHRP is intended for use in a multiprotocol internetworking layer environment over NBMA subnetworks.

This document is intended to be a functional superset of the NBMA Address Resolution Protocol (NARP) documented in [[1](#)].

Luciani, Katz, Piscitello, Cole
1]

[Page

Operation of NHRP as a means of establishing a transit path across an NBMA subnetwork between two routers will be addressed in a separate document.

1. Introduction

The NBMA Next Hop Resolution Protocol (NHRP) allows a source station (a host or router), wishing to communicate over a Non-Broadcast, Multi-Access (NBMA) subnetwork, to determine the internetworking layer addresses and NBMA addresses of suitable "NBMA next hops" toward a destination station. A subnetwork can be non-broadcast either because it technically doesn't support broadcasting (e.g., an X.25 subnetwork) or because broadcasting is not feasible for one reason or another (e.g., an SMDS multicast group or an extended Ethernet would be too large). If the destination is connected to the

NBMA subnetwork, then the NBMA next hop is the destination station itself. Otherwise, the NBMA next hop is the egress router from the NBMA subnetwork that is "nearest" to the destination station.

One way to model an NBMA network is by using the notion of logically independent IP subnets (LISs). LISs, as defined in [3] and [4], have the following properties:

- 1) All members of a LIS have the same IP network/subnet number and address mask.
- 2) All members within a LIS are directly connected to the same NBMA subnetwork.
- 3) All members outside of the LIS are accessed via a router.
- 4) All members within the LIS access each other directly (without routers)

Address resolution as described in [3] and [4] only resolves the next hop address if the destination station is a member of the same LIS as the source station; otherwise, the source station must forward packets to a router that is a member of multiple LIS's. In multi-LIS configurations, hop-by-hop address resolution may not be sufficient to resolve the "NBMA next hop" toward the destination station, and IP packets may have multiple IP hops through the NBMA subnetwork.

Another way to model NBMA is by using the notion of Local Address Groups (LAGs) [10]. The essential difference between the LIS and the LAG models is that while with the LIS model the outcome of the

"local/remote" forwarding decision is driven purely by addressing information, with the LAG model the outcome of this decision is

Luciani, Katz, Piscitello, Cole
2]

[Page

decoupled from the addressing information and is coupled with the Quality of Service and/or traffic characteristics. With the LAG model any two entities on a common NBMA network could establish a direct communication with each other, irrespective of the entities' addresses.

Support for the LAG model assumes the existence of a mechanism that allows any entity (i.e., host or router) connected to an NBMA network

to resolve an internetworking layer address to an NBMA address for any other entity connected to the same NBMA network. This

resolution

would take place regardless of the address assignments to these entities. NHRP describes such a mechanism. For example, when the internetworking layer address is of type IP, once the NBMA next hop has been resolved, the source may either start sending IP packets to the destination (in a connectionless NBMA subnetwork such as SMDS)

or

may first establish a connection to the destination with the desired bandwidth and QOS characteristics (in a connection-oriented NBMA subnetwork such as ATM).

Use of NHRP may be sufficient for hosts doing address resolution when

those hosts are directly connected to an NBMA subnetwork, allowing for straightforward implementations in NBMA stations. NHRP also has the capability of determining the egress point from an NBMA subnetwork when the destination is not directly connected to the

NBMA

subnetwork and the identity of the egress router is not learned by other methods (such as routing protocols). Optional extensions to NHRP provide additional robustness and diagnosability.

Address resolution techniques such as those described in [3] and [4] may be in use when NHRP is deployed. ARP servers and services over NBMA subnetworks may be required to support hosts that are not capable of dealing with any model for communication other than the LIS model, and deployed hosts may not implement NHRP but may

continue

to support ARP variants such as those described in [3] and [4].

NHRP

is intended to reduce or eliminate the extra router hops required by the LIS model, and can be deployed in a non-interfering manner alongside existing ARP services.

The operation of NHRP to establish transit paths across NBMA subnetworks between two routers requires additional mechanisms to avoid stable routing loops, and will be described in a separate document.

Luciani, Katz, Piscitello, Cole
3]

[Page

2. Overview

2.1 Terminology

The term "network" is highly overloaded, and is especially confusing in the context of NHRP. We use the following terms:

Internetwork layer--the media-independent layer (IP in the case of TCP/IP networks).

Subnetwork layer--the media-dependent layer underlying the internetwork layer, including the NBMA technology (ATM, X.25, SMDS, etc.)

The term "server", unless explicitly stated to the contrary, refers

to an Next Hop Server (NHS). An NHS is an entity performing the Next Hop Resolution Protocol service within the NBMA cloud. An NHS

is always tightly coupled with a routing entity (router, route server or edge device) although the converse is not yet guaranteed until ubiquitous deployment of this functionality occurs.

The term "client", unless explicitly stated to the contrary, refers

to an Next Hop Resolution Protocol client (NHC). An NHC is an entity which initiates NHRP requests of various types in order to obtain access to the NHRP service.

The term "station" generally refers to a host or router which contains an NHRP entity. Occasionally, the term station will describe a "user" of the NHRP client or service functionality; the difference in usage is largely semantic.

2.2 Protocol Overview

In this section, we briefly describe how a source S (which potentially can be either a router or a host) uses NHRP to determine the "NBMA next hop" to destination D.

For administrative and policy reasons, a physical NBMA subnetwork may

be partitioned into several, disjoint "Logical NBMA subnetworks". A Logical NBMA subnetwork is defined as a collection of hosts and routers that share unfiltered subnetwork connectivity over an NBMA subnetwork. "Unfiltered subnetwork connectivity" refers to the absence of closed user groups, address screening or similar features that may be used to prevent direct communication between stations connected to the same NBMA subnetwork. (Hereafter, unless otherwise specified, we use the term "NBMA subnetwork" to mean *logical* NBMA subnetwork.)

Placed within the NBMA subnetwork are one or more entities that implement the NHRP protocol. Such stations which are capable of answering Next Hop Resolution Requests are known as "Next Hop Servers" (NHSs). Each NHS serves a set of destination hosts, which may or may not be directly connected to the NBMA subnetwork. NHSs cooperatively resolve the NBMA next hop within their logical NBMA subnetwork. In addition to NHRP, NHSs may participate in protocols used to disseminate routing information across (and beyond the boundaries of) the NBMA subnetwork, and may support "classical" ARP service as well.

An NHS maintains a "next-hop resolution" cache, which is a table of address mappings (internetwork layer address to NBMA subnetwork layer address). This table can be constructed from information gleaned from NHRP Register packets (see [Section 5.2.3](#) and 5.2.4), extracted from Next Hop Resolution Requests/Replies that traverse the NHS as they are forwarded, or through mechanisms outside the scope of this document (examples of such mechanisms include ARP [[2](#), [3](#), [4](#)] and pre-configured tables). [Section 6.2](#) further describes cache management issues.

A host or router that is not an NHRP server must be configured with the identity of the NHS which serves it (see Configuration, [Section 4](#)).

[Note: for NBMA subnetworks that offer group or multicast addressing features, it may be desirable to configure stations with a group identity for NHSs, i.e., addressing information that would solicit a response from "all NHSs". The means whereby a group of NHSs divide responsibilities for next hop resolution are not described here.]

Whether or not a particular station within the NBMA subnetwork which is making use of the NHRP protocol needs to be able to act as an NHS is a local matter. For a station to avoid providing NHS functionality, there must be one or more NHSs within the NBMA subnetwork which are providing authoritative NBMA information on its behalf. If NHRP is to be able to resolve the NBMA address for stations that lack NHS functionality, these serving NHSs must exist along all routed paths between Next Hop Resolution Requesters and the station which cannot answer Next Hop Resolution Requests.

The protocol proceeds as follows. An event occurs triggering station

S to want to resolve the NBMA address of a path to D. This is most likely to be when a data packet addressed to station D is to be emitted from station S (either because station S is a host, or station S is a transit router), but the address resolution could also

be triggered by other means (a routing protocol update packet, for example). Station S first determines the next hop to station D

through normal routing processes (for a host, the next hop may simply be the default router; for routers, this is the "next hop" to the destination internetwork layer address). If the next hop is reachable through one of its NBMA interfaces, S constructs an Next Hop Resolution Request packet (see [Section 5.2.1](#)) containing station D's internetwork layer address as the (target) destination address, S's own internetwork layer address as the source address (Next Hop Resolution Request initiator), and station S's NBMA addressing information. Station S may also indicate that it prefers an authoritative Next Hop Resolution Reply (i.e., station S only wishes to receive a Next Hop Resolution Reply from the NHS-speaker that maintains the NBMA-to-internetwork layer address mapping for this destination). Station S emits the Next Hop Resolution Request packet towards the destination.

If the Next Hop Resolution Request is triggered by a data packet, station S may choose to dispose of the data packet while awaiting a Next Hop Resolution Reply in one of the following ways:

- (a) Drop the packet
- (b) Retain the packet until the Next Hop Resolution Reply arrives and a more optimal path is available
- (c) Forward the packet along the routed path toward D

The choice of which of the above to perform is a local policy matter, though option (c) is the recommended default, since it may allow data to flow to the destination while the NBMA address is being resolved. Note that an Next Hop Resolution Request for a given destination MUST NOT be triggered on every packet, though periodically retrying a Next Hop Resolution Request is permitted.

When the NHS receives an Next Hop Resolution Request, a check is made to see if it "serves" station D, i.e., the NHS checks to see if there is a "next hop" entry for D in its next-hop resolution cache. If the NHS does not serve D, the NHS forwards the Next Hop Resolution Request to another NHS. (Mechanisms for determining how to forward the Next Hop Resolution Request are discussed in [Section 3](#), Deployment.) Note that NHSs must be next hops to one another in order for forwarding of NHRP packets to be possible.

If this NHS serves D, the NHS resolves station D's NBMA address, and generates a positive Next Hop Resolution Reply (denoted by a 0 Code

in the reply) on D's behalf. (Next Hop Resolution Replies in this scenario are always marked as "authoritative".) The Next Hop Resolution Reply packet contains the next hop internetwork layer address and the NBMA address for station D and is sent back to S. (Note that if station D is not on the NBMA subnetwork, the next hop internetwork layer address will be that of the egress router through

which packets for station D are forwarded.)

An NHS receiving a Next Hop Resolution Reply may cache the NBMA next hop information contained therein. To a subsequent Next Hop Resolution Request, this NHS may respond with the cached, non-authoritative, NBMA next hop information or with cached negative information, if the NHS is allowed to do so, see [section 5.2.2](#) and 6.2. Non-authoritative Next Hop Resolution Replies are distinguished

from authoritative Next Hop Resolution Replies so that if a communication attempt based on non-authoritative information fails, a source station can choose to send an authoritative Next Hop Resolution Request. NHSs MUST NOT respond to authoritative Next Hop Resolution Requests with cached information.

[Note: An Next Hop Resolution Reply can be returned directly to the Next Hop Resolution Request initiator, i.e., without traversing the list of NHSs that forwarded the Next Hop Resolution Request, if all of the following criteria are satisfied:

- (a) Direct communication is available via datagram transfer (e.g., SMDS) or the NHS has an existing virtual circuit connection to the Next Hop Resolution Request initiator or is permitted to open one.
- (b) The Next Hop Resolution Request initiator has not included the NHRP Reverse NHS record Extension (see [Section 5.3.5](#)).
- (c) The authentication policy in force permits direct communication between the NHS and the Next Hop Resolution Request initiator.

The purpose of allowing an NHS to send a Next Hop Resolution Reply directly is to reduce response time. A consequence of allowing a direct Next Hop Resolution Reply is that NHSs that would under normal circumstances be traversed by the Next Hop Resolution Reply would not cache next hop information contained therein.]

The process of forwarding the Next Hop Resolution Request is repeated until the Next Hop Resolution Request is satisfied, or an error occurs (e.g., no NHS in the NBMA subnetwork can resolve the Next Hop Resolution Request.) If the determination is made that station D's next hop cannot be resolved, a negative Next Hop Resolution Reply (NAK) is returned. This occurs when (a) no next-hop resolution information is available for station D from any NHS, or (b) an NHS is unable to forward the Next Hop Resolution Request (e.g.,

connectivity
is lost).

NHRP Registration Requests, NHRP Purge Requests, NHRP Purge Replies,
and NHRP Error Indications follow the routed path from sender to
receiver in the same fashion that Next Hop Resolution Requests and

Luciani, Katz, Piscitello, Cole
7]

[Page

Next Hop Resolution Replies do. That is, "requests" and "indications" follow the routed path from Source Protocol Address (which is the address of the station initiating the communication) to the Destination Protocol Address. "Replies", on the other hand, follow the routed path from the Destination Protocol Address back to the Source Protocol Address with the exceptions mentioned above where a direct VC may be created. In the case of a NHRP Registration Reply, the packet is always returned via a direct VC (see [Section 5.2.4](#)).

NHRP Requests and NHRP Replies MUST NOT cross the borders of a logical NBMA subnetwork (an explicit NBMA subnetwork identifier may be included as an extension in the Next Hop Resolution Request, see [section 5.3.2](#)). Thus, the internetwork layer traffic out of and into a logical NBMA subnetwork always traverses an internetwork layer router at its border. Internetwork layer filtering can then be implemented at these border routers.

NHRP optionally provides a mechanism to send a Next Hop Resolution Reply which contains aggregated NBMA next hop information. Suppose that router X is the NBMA next hop from station S to station D. Suppose further that X is an egress router for all stations sharing an internetwork layer address prefix with station D. When a Next Hop Resolution Reply is generated in response to a Next Hop Resolution Request, the responder may augment the internetwork layer address of station D with a prefix length (see [Section 5.2.0.1](#)). A subsequent (non-authoritative) Next Hop Resolution Request for some destination that shares an internetwork layer address prefix (for the number of bits specified in the prefix length) with D may be satisfied with this cached information. See [section 6.2](#) regarding caching issues.

To dynamically detect subnetwork-layer filtering in NBMA subnetworks (e.g., X.25 closed user group facility, or SMDS address screens), to trace the routed path that an NHRP packet takes, or to provide loop detection and diagnostic capabilities, a "Route Record" may be included in NHRP packets (see Sections [5.3.4](#) and [5.3.5](#)). The Route Record extensions contain the internetwork (and subnetwork layer) addresses of all intermediate NHSs between source and destination (in the forward direction) and between destination and source (in the reverse direction). When a source station is unable to communicate with the responder (e.g., an attempt to open an SVC fails), it may attempt to do so successively with other subnetwork layer addresses in the Route Record until it succeeds (if authentication policy permits such action). This approach can find a suitable egress point in the presence of subnetwork-layer filtering (which may be

source/destination sensitive, for instance, without necessarily creating separate logical NBMA subnetworks) or subnetwork-layer

Luciani, Katz, Piscitello, Cole
8]

[Page

congestion (especially in connection-oriented media).

3. Deployment

Next Hop Resolution Requests traverse one or more hops within an NBMA subnetwork before reaching the station that is expected to generate a

response. Each station, including the source station, chooses a neighboring NHS to which it will forward the Next Hop Resolution Request. The NHS selection procedure typically involves applying a destination protocol layer address to the protocol layer routing table which causes a routing decision to be returned. This routing decision is then used to forward the Next Hop Resolution Request to the downstream NHS. The destination protocol layer address

previously

mentioned is carried within the Next Hop Resolution Request packet.

Note that even though a protocol layer address was used to acquire a routing decision, NHRP packets are not encapsulated within a protocol

layer header but rather are carried at the NBMA layer using the encapsulation described in [Section 5](#).

Each NHS/router examines the Next Hop Resolution Request packet on its way toward the destination. Each NHS which the NHRP packet traverses on the way to the packet's destination might modify the packet (e.g., updating the Forward Record extension). Ignoring error

situations, the Next Hop Resolution Request eventually arrives at a station that is to generate an Next Hop Resolution Reply. This responding station "serves" the destination. The responding station generates a Next Hop Resolution Reply using the source protocol address from within the NHRP packet to determine where the Next Hop Resolution Reply should be sent.

Rather than use routing to determine the next hop for an NHRP packet,

an NHS may use static configuration information (or other applicable means) in order to determine to which neighboring NHSs to forward the

Next Hop Resolution Request packet. The use of static configuration information for this purpose is beyond the scope of this document.

In order to forward NHRP packets to a neighboring NHS, NHRP clients must nominally be configured with the NBMA address of at least one NHS. In practice, a client's default router should also be its NHS in that way a client may be able to know the NBMA address of its NHS from the configuration which was already required for the client to be able to communicate.

The NHS serving a particular destination must lie along the routed

path to that destination. In practice, this means that all egress routers must double as NHSs serving the destinations beyond them, and that hosts on the NBMA subnetwork are served by routers that double

as NHSs. Also, this implies that forwarding of NHRP packets within an NBMA subnetwork requires a contiguous deployment of NHRP capable routers. During migration to NHRP, it cannot be expected that all routers within the NBMA subnetwork are NHRP capable. Thus, NHRP traffic which would otherwise need to be forwarded through such routers can be expected to be dropped due to the NHRP packet not being recognized. In this case, NHRP will be unable to establish

any

transit paths whose discovery requires the traversal of the non-NHRP speaking routers. If the client has tried and failed to acquire a cut through path then the client should use the network layer routed path as a default.

If a subnetwork offers a link layer group addressing or multicast feature, the client (station) may be configured with a group address assigned to the group of next-hop servers. The client might then submit Next Hop Resolution Requests to the group address, eliciting

a

response from one or more NHSs, depending on the response strategy selected. Note that the constraints described in [Section 2](#)

regarding

directly sending Next Hop Resolution Reply may apply.

4. Configuration

Clients

To participate in NHRP, a client connected to an NBMA subnetwork should be configured with the NBMA address(es) of its NHS(s) (alternatively, it should be configured with a means of acquiring them, i.e., the group address that members of a NHS group use for the purpose of address or next-hop resolution.) The NHS(s) will likely also represent the client's default or peer routers, so their NBMA addresses may be obtained from the client's existing configuration. If the client is attached to several subnetworks (including logical NBMA subnetworks), the client should also be configured to receive routing information from its NHS(s) and peer routers so that it can determine which internetwork layer networks are reachable through which subnetworks.

Next Hop Servers

An NHS is configured with knowledge of its own internetwork layer and NBMA addresses and a logical NBMA subnetwork identifier (see [Section 5.3.2](#)). An NHS MAY also be configured with a set of internetwork layer address prefixes that correspond to the internetwork layer addresses of the stations it serves. If a

served

client is attached to several subnetworks, the NHS may also need

to

be configured to advertise routing information to such client.

If an NHS acts as an egress router for stations connected to other subnetworks than the NBMA subnetwork, the NHS must, in addition to the above, be configured to exchange routing information between the NBMA subnetwork and these other subnetworks.

In all cases, routing information is exchanged using conventional intra-domain and/or inter-domain routing protocols.

The NBMA addresses of the stations served by the NHS may be learned via NHRP Register packets or manual configuration.

5. NHRP Packet Formats

This section describes the format of NHRP packets. In the following, unless otherwise stated explicitly, the unqualified term "request" refers generically to any of the NHRP packet types which are "requests". Further, unless otherwise stated explicitly, the unqualified term "reply" refers generically to any of the NHRP packet types which are "replies".

An NHRP packet consists of a Fixed Part, a Mandatory Part, and an Extensions Part. The Fixed Part is common to all NHRP packet types. The Mandatory Part **MUST** be present, but varies depending on packet type. The Extensions Part also varies depending on packet type, and need not be present.

The length of the Fixed Part is fixed at 20 octets. The length of the Mandatory Part is determined by the contents of the extensions offset field (ar\$extoff). If ar\$extoff=0x0 then the mandatory part length is equal to total packet length (ar\$pktsz) minus 20 otherwise the mandatory part length is equal to ar\$extoff minus 20. The length of the Extensions Part is implied by ar\$pktsz minus ar\$extoff. NHSS may increase the size of an NHRP packet as a result of extension processing, but not beyond the offered maximum SDU size of the NBMA network.

NHRP packets are encapsulated using the native formats used on the particular NBMA network over which NHRP is carried. For example, SMDS networks always use LLC/SNAP encapsulation at the NBMA layer, and an NHRP packet is preceded by the following LLC/SNAP encapsulation:

```
[0xAA-AA-03] [0x00-00-5E] [0x00-03]
```

The first three octets are LLC, indicating that SNAP follows. The SNAP OUI portion is the IANA's OUI, and the SNAP PID portion identifies NHRP (see [4]).

ar\$afn

Defines the type of "link layer" addresses being carried. This number is taken from the 'address family number' list specified in [6]. This field has implications to the coding of ar\$shtl and ar\$sstl as described below.

ar\$pro.type

field is a 16 bit unsigned integer representing the following number space:

0x0000 to 0x00FF Protocols defined by the equivalent NLPIDs.
0x0100 to 0x03FF Reserved for future use by the IETF.

0x0400 to 0x04FF Allocated for use by the ATM Forum.
0x0500 to 0x05FF Experimental/Local use.
0x0600 to 0xFFFF Protocols defined by the equivalent

Ethertypes.

(based on the observations that valid Ethertypes are never smaller than 0x600, and NLPIDs never larger than 0xFF.)

ar\$pro.snap

When ar\$pro.type has a value of 0x0080, a SNAP encoded extension is being used to encode the protocol type. This snap extension is placed in the ar\$pro.snap field. This is termed the 'long form' protocol ID. If ar\$pro != 0x0080 then the ar\$pro.snap field MUST be zero on transmit and ignored on receive. The ar\$pro.type field itself identifies the protocol being referred to. This is termed the 'short form' protocol ID.

In all cases, where a protocol has an assigned number in the ar\$pro.type space (excluding 0x0080) the short form MUST be used when transmitting NHRP messages. Additionally, where a protocol

has valid short and long forms of identification, receivers MAY choose to recognize the long form.

ar\$hopcnt

The Hop count indicates the maximum number of NHSS that an NHRP packet is allowed to traverse before being discarded.

ar\$pktsz

The total length of the NHRP packet, in octets (excluding link layer encapsulation).

ar\$chksum

The standard IP checksum over the entire NHRP packet (starting with the fixed header). If only the hop count field is changed, the checksum is adjusted without full recomputation. The checksum is completely recomputed when other header fields are changed.

ar\$extoff

This field identifies the existence and location of NHRP extensions. If this field is 0 then no extensions exist otherwise this field represents the offset from the beginning of the NHRP packet (i.e., starting from the ar\$afn field) of the first extension.

ar\$op.version

This field is set to 0x01 for NHRP version 1.

ar\$op.type

This is the NHRP packet type: NHRP Next Hop Resolution Request(1),

Luciani, Katz, Piscitello, Cole
13]

[Page

NHRP Next Hop Resolution Reply(2), NHRP Registration Request(3), NHRP Registration Reply(4), NHRP Purge Request(5), NHRP Purge Reply(6), or NHRP Error Indication(7). Use of NHRP packet Types

in

the range 128 to 255 are reserved for research or use in other protocol development and will be administered by IANA.

ar\$shtl

Type & length of source NBMA address interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0003 for NSAP, ar\$afn=8 for E.164). When ar\$afn=0x000F (E.164 address plus NSAP subaddress) then both ar\$shtl and ar\$sstl must be coded appropriately (see below).

ar\$sstl

Type & length of source NBMA subaddress interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x000F for NSAP). When an NBMA technology has no concept

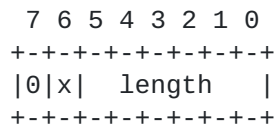
of

a subaddress, the subaddress length is always coded ar\$sstl = 0

and

no storage is allocated for the subaddress in the appropriate mandatory part.

ar\$shtl, ar\$sstl, subnetwork layer addresses, and subnetwork layer subaddresses fields are coded as follows:



The most significant bit is reserved and MUST be set to zero. The second most significant bit (x) is a flag indicating whether the address being referred to is in:

- NSAP format (x = 0).
- Native E.164 format (x = 1).

For NBMA technologies that use neither NSAP nor E.164 format addresses, x = 0 SHALL be used to indicate the native form for the particular NBMA technology.

In the case where the NBMA is ATM, if a subaddress is to be included then ar\$afn MUST be set to 0x000F which means that if a subaddress exists then it is of type NSAP.

The bottom 6 bits is an unsigned integer value indicating the length of the associated NBMA address in octets. If this value is zero the flag x is ignored.

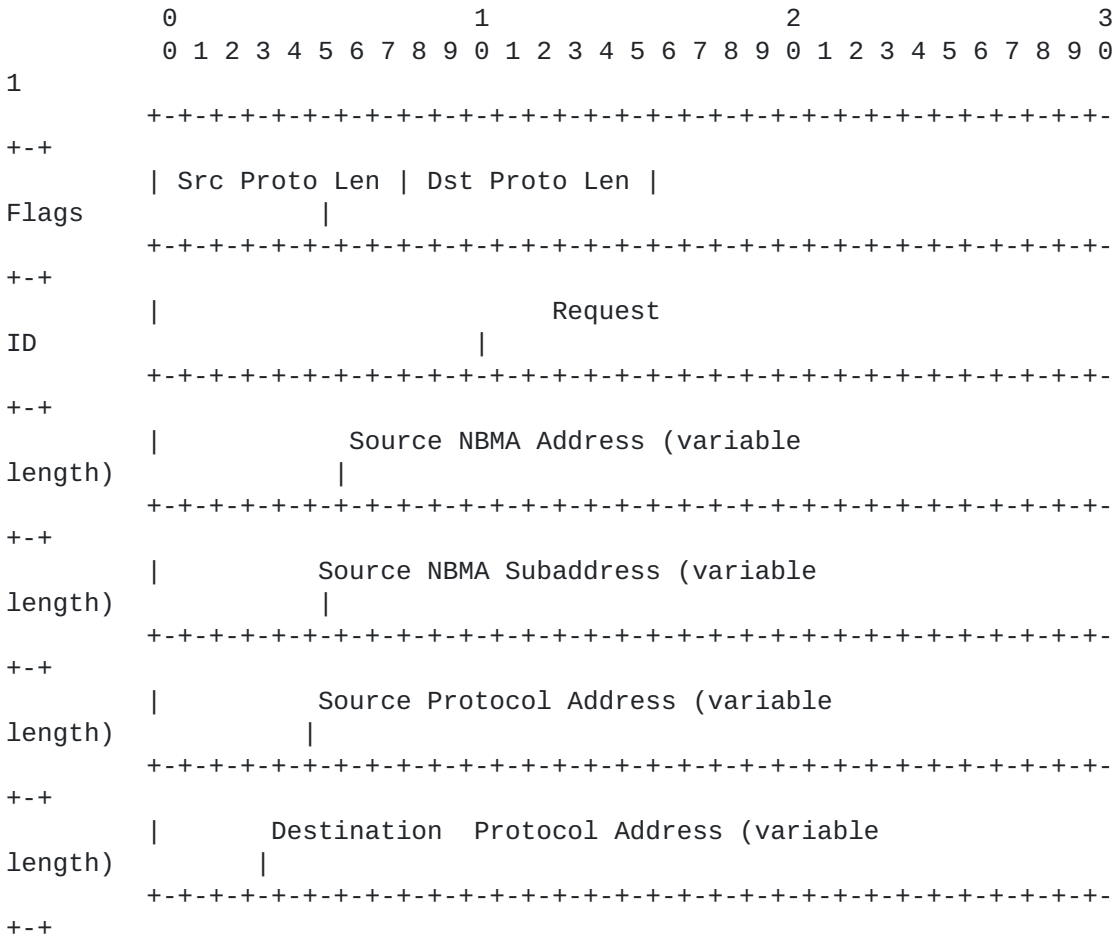
5.2.0 Mandatory Part

The Mandatory Part of the NHRP packet contains the operation specific information (e.g., Next Hop Resolution Request/Reply, etc.) and variable length data which is pertinent to the packet type.

5.2.0.1 Mandatory Part Format

Sections 5.2.1 through 5.2.6 have a very similar mandatory part. This mandatory part includes a common header and zero or more Client Information Entries (CIEs). Section 5.2.7 has a different format which is specified in that section.

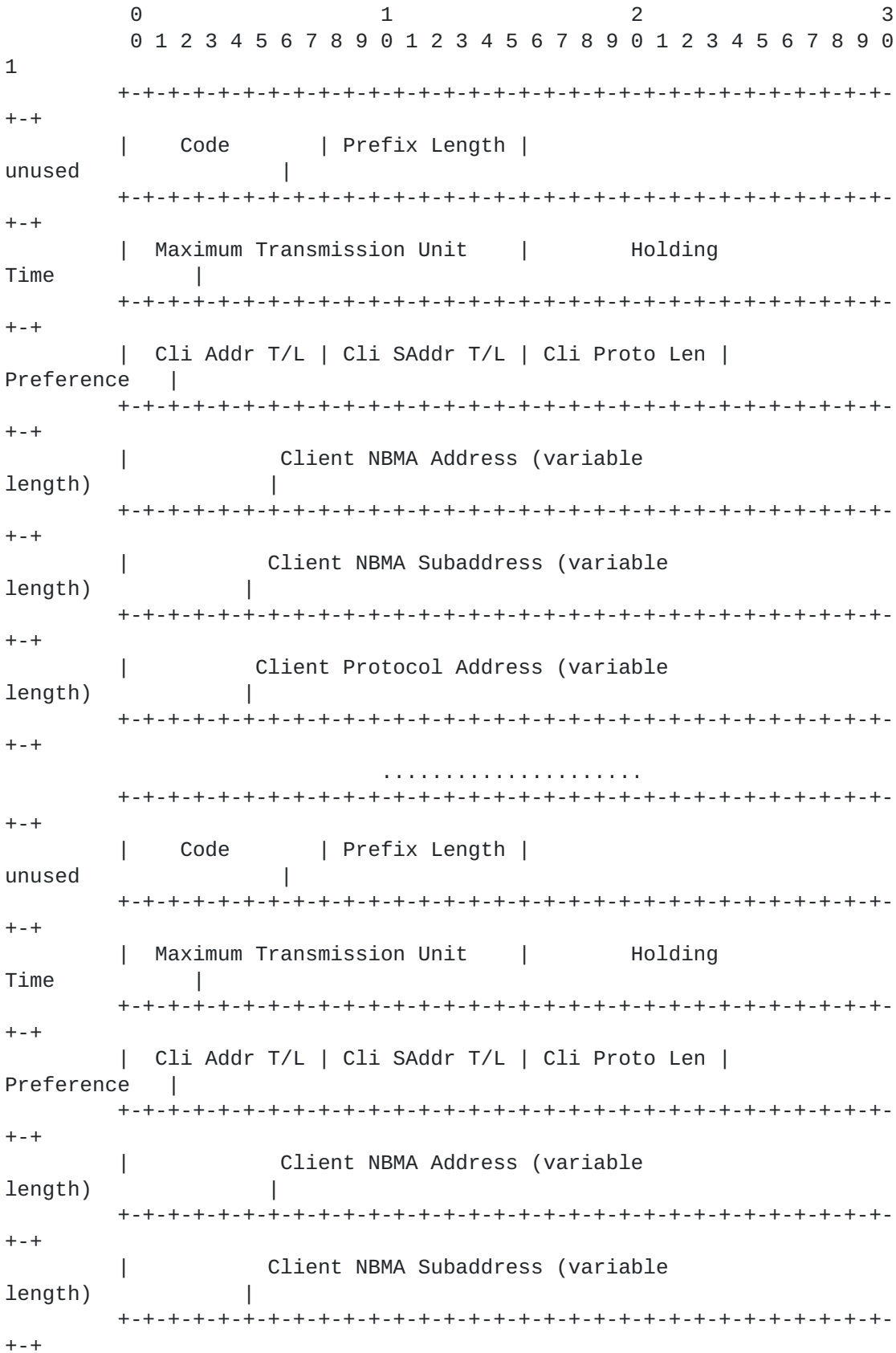
The common header looks like the following:



And the CIEs have the following format:

Luciani, Katz, Piscitello, Cole
15]

[Page



packet into the associated "reply". When a sender of a "request" receives "reply", it will compare the Request ID and source address information in the received "reply" against that found in its outstanding "request" list. When a match is found then the "request" is considered to be acknowledged.

The value is taken from a 32 bit counter that is incremented each time a new "request" is transmitted. The same value MUST be used when resending a "request", i.e., when a "reply" has not been received for a "request" and a retry is sent after an appropriate interval.

The NBMA address/subaddress form specified below allows combined E.164/NSAPA form of NBMA addressing. For NBMA technologies without a subaddress concept, the subaddress field is always ZERO length and ar\$stl = 0.

Source NBMA Address

The Source NBMA address field is the address of the source station which is sending the "request". If the field's length as specified in ar\$stl is 0 then no storage is allocated for this address at all.

Source NBMA SubAddress

The Source NBMA subaddress field is the address of the source station which is sending the "request". If the field's length as specified in ar\$stl is 0 then no storage is allocated for this address at all.

Source Protocol Address

This is the protocol address of the station which is sending the "request". This is also the protocol address of the station toward which a "reply" packet is sent.

Destination Protocol Address

This is the protocol address of the station toward which a "request" packet is sent.

Code

This field is message specific. See the relevant message sections below. In general, this field is a NAK code; i.e., when the field is 0 in a reply then the packet is acknowledging a request and if it contains any other value the packet contains a negative acknowledgment.

Prefix Length

This field is message specific. See the relevant message sections below. In general, however, this field is used to indicate that

the information carried in an NHRP the message pertains to an equivalence class of internetwork layer addresses rather than just a single internetwork layer address specified. All internetwork layer addresses that match the first "Prefix Length" bit positions for the specific internetwork layer address are included in the equivalence class.

Maximum Transmission Unit

This field gives the maximum transmission unit for the relevant client station. If this value is 0 then either the default MTU is used or the MTU negotiated via signaling is used if such negotiation is possible for the given NBMA.

Holding Time

The Holding Time field specifies the number of seconds for which the Next Hop NBMA information specified in the CIE is considered to be valid. Cached information SHALL be discarded when the holding time expires. This field must be set to 0 on a NAK.

Cli Addr T/L

Type & length of next hop NBMA address specified in the CIE. This field is interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0003 for ATM).

Cli SAddr T/L

Type & length of next hop NBMA subaddress specified in the CIE. This field is interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0015 for ATM makes the address an E.164 and the subaddress an ATM Forum NSAP address).

When an NBMA technology has no concept of a subaddress, the subaddress is always null with a length of 0. When the address length is specified as 0 no storage is allocated for the address.

Cli Proto Len

This field holds the length in octets of the Client Protocol Address specified in the CIE.

Preference

This field specifies the preference for use of the specific CIE relative to other CIEs. Higher values indicate higher preference. Action taken when multiple CIEs have equal or highest preference value is a local matter.

Client NBMA Address

This is the client's NBMA address.

Client NBMA SubAddress

This is the client's NBMA subaddress.

Client Protocol Address

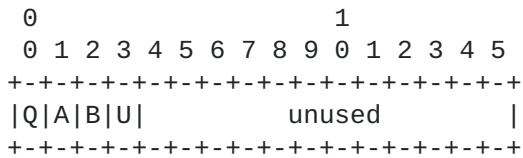
This is the client's internetworking layer address specified.

Note that an NHS SHOULD NOT cache source information which is in an NHRP message because this information could be inappropriately used to set up a cut-through without doing proper filtering along a routed path. Further, in the case where a distributed router exists in the network, incorrect or incomplete information may be included in the source information.

5.2.1 NHRP Next Hop Resolution Request

The NHRP Next Hop Resolution Request packet has a Type code of 1. Its mandatory part is coded as described in [Section 5.2.0.1](#) and the message specific meanings of the fields are as follows:

Flags - The flags field is coded as follows:



Q
Set if the station sending the Next Hop Resolution Request is a router; clear if the it is a host.

A
This bit is set in a Next Hop Resolution Request if only authoritative next hop information is desired and is clear otherwise. See the NHRP Next Hop Resolution Reply section below for further details on the "A" bit and its usage.

B
Unused (clear on transmit)

U
This is the Uniqueness bit. This bit aids in duplicate address detection. When this bit is set in an NHRP Resolution Request and one or more entries exist in the NHS cache which meet the requirements of the NHRP Resolution Request then only the CIE in the NHS's cache with this bit set will be returned. Note that even if this bit was set at registration time, there may still

be multiple CIEs that might fulfill the NHRP Resolution Request because an entire subnet can be registered through use of the Prefix Length in the CIE and the address of interest might be within such a subnet. If the "uniqueness" bit is set and the

responding NHS has one or more cache entries which match the request but no such cache entry has the "uniqueness" bit set, then the NHRP Resolution Reply returns with a NAK code of "13 - Binding Exists But Is Not Unique" and no CIE is included. If a client wishes to receive non-unique Next Hop Entries, then the client must have the "uniqueness" bit set to zero in its

NHRP

Resolution Request. Note that when this bit is set in an NHRP Registration Request, only a single CIE may be specified in the NHRP Registration Request and that CIE must have the Prefix Length field set to 0xFF.

Zero or one CIEs (see [Section 5.2.0.1](#)) may be specified in an NHRP Next Hop Resolution Request. If one is specified then that entry carries the pertinent information for the client sourcing the NHRP Next Hop Resolution Request. Usage of the CIE in the NHRP Next Hop Resolution Request is described below:

Prefix Length

If a CIE is specified in the NHRP Next Hop Resolution Request then the Prefix Length field may be used to qualify the widest acceptable prefix which may be used to satisfy the NHRP Next Hop Resolution Request. In the case of NHRP Next Hop Resolution Request/Reply, the Prefix Length specifies the equivalence class of addresses which match the first "Prefix Length" bit positions of the Destination Protocol Address. If this field is set to 0x00 then this field MUST be ignored. If the "U" bit is set in the common header then this field MUST be set to 0xFF.

Maximum Transmission Unit

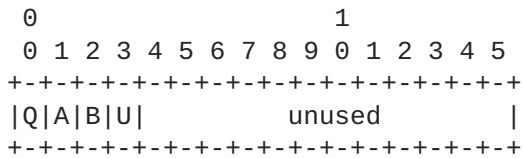
This field gives the maximum transmission unit for the source station. A possible use of this field in the Next Hop Resolution Request packet is for the Next Hop Resolution Requester to ask for a target MTU. In lieu of that usage, the CIE must be omitted.

All other fields in the CIE MUST be ignored and SHOULD be set to 0.

5.2.2 NHRP Next Hop Resolution Reply

The NHRP Next Hop Resolution Reply packet has a Type code of 2. CIEs correspond to Next Hop Entries in an NHS's cache which match the criteria in the NHRP Next Hop Resolution Request. Its mandatory part is coded as described in [Section 5.2.0.1](#). The message specific meanings of the fields are as follows:

Flags - The flags field is coded as follows:



Q

Hop

Copied from the Next Hop Resolution Request. Set if the Next Resolution Requester is a router; clear if it is a host.

A

Set if the next hop CIE in the Next Hop Resolution Reply is authoritative; clear if the Next Hop Resolution Reply is non-authoritative.

When an NHS receives a Next Hop Resolution Request for authoritative information for which it is the authoritative source, it MUST respond with a Next Hop Resolution Reply containing all and only those next hop CIEs which are contained in the NHS's cache which both match the criteria of the Next Hop Resolution Request and are authoritative cache entries. An NHS is an authoritative source for a Next Hop Resolution Request if the information in the NHS's cache matches the Next Hop Resolution Request criteria and that information was obtained through a NHRP Registration Request or through synchronization with an NHS which obtained this information through a NHRP Registration Request. An authoritative cache entry is one which is obtained through a NHRP Registration Request or through synchronization with an NHS which obtained this information through a NHRP Registration Request.

are

Hop

NHS's

An NHS obtains non-authoritative CIEs through promiscuous listening to NHRP packets other than NHRP Registrations which are directed at it. A Next Hop Resolution Request which indicates a request for non-authoritative information should cause a Next Hop Resolution Reply which contains all entries in the replying NHS's cache (i.e., both authoritative and non-authoritative) which match the criteria specified in the request.

B

Hop

Set if the association between the destination and the next hop information is guaranteed to be stable for the lifetime of the information (the holding time). This is the case if the Next Hop protocol address identifies the destination (though it may be different in value than the Destination address if the destination system has multiple addresses) or if the destination is not connected directly to the NBMA subnetwork but the egress

router to that destination is guaranteed to be stable (such as

Luciani, Katz, Piscitello, Cole
21]

[Page

when the destination is immediately adjacent to the egress router through a non-NBMA interface). This information affects caching strategies (see [section 6.2](#)).

U

This is the Uniqueness bit. See the NHRP Resolution Request section above for details. When this bit is set only, only one CIE is included since only one unique binding should exist in an NHS's cache.

One or more CIEs are specified in the NHRP Next Hop Resolution Reply.

Each CIE contains NHRP next hop information which the responding NHS has cached and which matches the parameters specified in the NHRP Next Hop Resolution Request. If no match is found by the NHS issuing

the NHRP Next Hop Resolution Reply then a single CIE is enclosed with

the a CIE Code set appropriately (see below) and all other fields MUST be ignored and SHOULD be set to 0. In order to facilitate the use of NHRP by minimal client implementations, the first CIE MUST contain the next hop with the highest preference value so that such an implementation need parse only a single CIE.

Code

If this field is set to zero then this packet contains a positively acknowledged NHRP Resolution Reply. If this field contains any other value then this message contains an NHRP Resolution Reply NAK which means that an appropriate internetworking layer to NBMA address binding was not available in the responding NHS's cache. If NHRP Resolution Reply contains

a Client Information Entry with a NAK Code other than 0 then it MUST NOT contain any other CIE. Currently defined NAK Codes are as follows:

12 - No Internetworking Layer Address to NBMA Address Binding Exists

This code states that there were absolutely no internetworking layer address to NBMA address bindings found in the responding NHS's cache.

13 - Binding Exists But Is Not Unique

This code states that there were one or more internetworking layer address to NBMA address bindings found in the responding NHS's cache, however none of them had the uniqueness bit set.

Prefix Length

In the case of NHRP Next Hop Resolution Reply, the Prefix Length

specifies the equivalence class of addresses which match the

Luciani, Katz, Piscitello, Cole
22]

[Page

first "Prefix Length" bit positions of the Destination Protocol Address.

Holding Time

The Holding Time specified in a CIE of an NHRP Resolution Reply is the amount of time remaining before the expiration of the client information which is cached at the replying NHS. It is not the value which was registered by the client.

The remainder of the fields for the CIE for each next hop are filled out as they were defined when the next hop was registered with the responding NHS (or one of the responding NHS's synchronized servers) via the NHRP Registration Request.

Load-splitting may be performed when more than one Client Information

Entry is returned to a requester when equal preference values are specified. Also, the alternative addresses may be used in case of connectivity failure in the NBMA subnetwork (such as a failed call attempt in connection-oriented NBMA subnetworks).

Any extensions present in the Next Hop Resolution Request packet MUST

be present in the NHRP Next Hop Resolution Reply even if the extension is non-Compulsory.

If an unsolicited NHRP Next Hop Resolution Reply packet is received, an Error Indication of type Invalid Next Hop Resolution Reply Received SHOULD be sent in response.

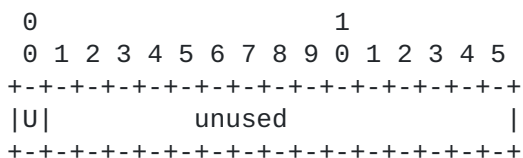
5.2.3 NHRP Registration Request

The NHRP Registration Request is sent from a station to an NHS to notify the NHS of the station's NBMA information. It has a Type code

of 3. Each CIE corresponds to Next Hop information which is to be cached at an NHS. The mandatory part of an NHRP Registration Request

is coded as described in [Section 5.2.0.1](#). The message specific meanings of the fields are as follows:

Flags - The flags field is coded as follows:



U

This is the Uniqueness bit. When set in an NHRP Registration Request, this bit indicates that the registration of the

protocol

Luciani, Katz, Piscitello, Cole
23]

[Page

address is unique within the confines of the set of synchronized NHSs. This "uniqueness" qualifier MUST be stored in the NHS/NHC cache. Any attempt to register a binding between the protocol address and an NBMA address when this bit is set MUST be rejected with a Code of "14 - Unique Internetworking Layer Address Already Registered" if the replying NHS already has a cache entry for the protocol address and the cache entry has the "uniqueness" bit set. A registration of a CIE's information is rejected when the CIE is returned with the Code field set to anything other than 0x00. See the description of the uniqueness bit in NHRP Resolution Request section above for further details. When this bit is set only, only one CIE MAY be included in the NHRP Registration Request.

Request ID

The request ID has the same meaning as described in [Section 5.2.0.1](#). However, the request ID for NHRP Registrations which is maintained at each client MUST be kept in non-volatile memory so that when a client crashes and reregisters there will be no inconsistency in the NHS's database. In order to reduce the overhead associated with updating non-volatile memory, the actual updating need not be done with every increment of the Request ID but could be done, for example, every 50 or 100 increments. In this scenario, when a client crashes and reregisters it knows to add 100 to the value of the Request ID in the non-volatile memory before using the Request ID for subsequent registrations.

One or more CIEs are specified in the NHRP Registration Request. Each CIE contains next hop information which a client is attempting to register with its servers. Generally, all fields in CIEs enclosed in NHRP Registration Requests are coded as described in [Section 5.2.0.1](#). However, if a station is only registering itself with the NHRP Registration Request then it MAY code the Cli Addr T/L, Cli SAddr T/L, and Cli Proto Len as zero which signifies that the client address information is to be taken from the source information in the common header (see [Section 5.2.0.1](#)). Below, further clarification is given for some fields in a CIE in the context of a NHRP Registration Request.

Code

This field is set to 0x00 in NHRP Registration Requests.

Prefix Length

This field may be used in a NHRP Registration Request to register equivalence information for the Client Protocol Address specified

Luciani, Katz, Piscitello, Cole
24]

[Page

in the CIE of an NHRP Registration Request In the case of NHRP Registration Request, the Prefix Length specifies the equivalence class of addresses which match the first "Prefix Length" bit positions of the Client Protocol Address. If this field is set to 0x00 then this field MUST be ignored and no equivalence information is assumed (i.e., only Client Protocol Address is bound to the NBMA information). If the "U" bit is set in the common header then this field MUST be set to 0xFF.

This packet is used to register a station's NHRP information with its

NHSs, as configured or known through conventional routing means. NHSs may also be configured with the identities of stations that they

serve. If an NHS receives an NHRP Registration Request packet which has the Destination Protocol Address field set to an address other than the NHS's own protocol address then the NHS MUST forward the packet along the routed path toward the Destination Protocol Address.

It is possible that a misconfigured station will attempt to register with the wrong NHS (i.e., one that cannot serve it due to policy constraints or routing state). If this is the case, the NHS MUST reply with a NAK-ed Registration Reply of type Can't Serve This Address.

If an NHS cannot serve a station due to a lack of resources, the NHS MUST reply with a NAK-ed Registration Reply of type Registration Overflow.

In order to keep the registration entry from being discarded, the station MUST re-send the NHRP Registration Request packet often enough to refresh the registration, even in the face of occasional packet loss. It is recommended that the NHRP Registration Request packet be sent at an interval equal to one-third of the Holding Time specified therein.

5.2.4 NHRP Registration Reply

The NHRP Registration Reply is sent by an NHS to a client in response

to that client's NHRP Registration Request. If the Code field of a CIE in the NHRP Registration Reply has anything other than 0 zero in it then the NHRP Registration Reply is a NAK otherwise the reply is an ACK. The NHRP Registration Reply has a Type code of 4.

An NHRP Registration Reply is formed from an NHRP Registration Request by changing the type code to 4, updating the CIE Code field, and filling in the appropriate extensions if they exist. The message

specific meanings of the fields are as follows:

Attempts to register the information in the CIEs of an NHRP

Luciani, Katz, Piscitello, Cole
25]

[Page

Registration Request may fail for various reasons. If this is the case then each failed attempt to register the information in a CIE of an NHRP Registration Request is logged in the associated NHRP Registration Reply by setting the CIE Code field to the appropriate error code as shown below:

CIE Code

0 - Successful Registration

The information in the CIE was successfully registered with the NHS.

4 - Can't Serve This Address

An NHS may refuse an NHRP Registration Request attempt for administrative reasons (due to policy constraints or routing state). If so, the NHS MUST send an NHRP Registration Reply which contains a NAK code of 4.

5 - Registration Overflow

If an NHS cannot serve a station due to a lack of resources, the NHS MUST reply with a NAKed NHRP Registration Reply which contains a NAK code of 5.

14 - Unique Internetworking Layer Address Already Registered

If a client tries to register a protocol address to NBMA address binding with the uniqueness bit on and the protocol address already exists in the NHS's cache then if that cache entry also has the uniqueness bit on then this NAK Code is returned in the CIE in the NHRP Registration Reply.

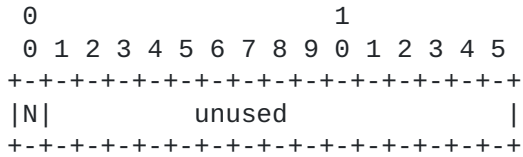
Due to the possible existence of asymmetric routing, an NHRP Registration Reply may not be able to merely follow the routed path back to the source protocol address specified in the common header of the NHRP Registration Reply. As a result, there MUST exist a direct NBMA level connection between the client and its NHS on which to send the NHRP Registration Reply before NHRP Registration Reply may be returned to the client. If such a connection does not exist then the NHS must setup such a connection to the client by using the source NBMA information supplied in the common header of the NHRP Registration Request.

5.2.5 NHRP Purge Request

The NHRP Purge Request packet is sent in order to invalidate cached information in a station. The NHRP Purge Request packet has a type

code of 5. The mandatory part of an NHRP Purge Request is coded as described in [Section 5.2.0.1](#). The message specific meanings of the fields are as follows:

Flags - The flags field is coded as follows:



N

When set, this bit tells the receiver of the NHRP Purge Request that the requester does not expect to receive an NHRP Purge Reply. If an unsolicited NHRP Purge Reply is received by a station where that station is identified in the Source Protocol Address of the packet then that packet must be ignored.

One or more CIEs are specified in the NHRP Purge Request. Each CIE contains next hop information which is to be purged from an NHS/NHC cache. Generally, all fields in CIEs enclosed in NHRP Purge Requests are coded as described in [Section 5.2.0.1](#). Below, further clarification is given for some fields in a CIE in the context of a NHRP Purge Request.

Code

This field is set to 0x00 in NHRP Purge Requests.

Prefix Length

In the case of NHRP Purge Requests, the Prefix Length specifies the equivalence class of addresses which match the first "Prefix Length" bit positions of the Client Protocol Address specified in the CIE. All next hop information which contains a protocol address which matches an element of this equivalence class is to be purged from the receivers cache. If this field is set to 0x00 then this field MUST be ignored and no equivalence information is assumed.

The Maximum Transmission Unit and Preference fields of the CIE are coded as zero. The Holding Time should be coded as zero but there may be some utility in supplying a "short" holding time to be applied to the matching next hop information before that information would be purged; this usage is for further study. The Client Protocol Address field and the Cli Proto Len field MUST be filled in. The Client Protocol Address is filled in with the protocol address to be purged from the receiving station's cache

while the Cli Proto Len is set the length of the purged client's protocol address. All remaining fields in the CIE MAY be set to zero although the client NBMA information (and associated length fields) MAY be specified to narrow the scope of the NHRP Purge Request if requester desires. However, the receiver of an NHRP Purge Request may choose to ignore the Client NBMA information if it is supplied.

An NHRP Purge Request packet is sent from an NHS to a station to cause it to delete previously cached information. This is done when the information may be no longer valid (typically when the NHS has previously provided next hop information for a station that is not directly connected to the NBMA subnetwork, and the egress point to that station may have changed).

An NHRP Purge Request packet may also be sent from a client to an NHS

with which the client had previously registered. This allows for a client to invalidate its registration with NHRP before it would otherwise expire via the holding timer.

The station sending the NHRP Purge Request MAY periodically retransmit the NHRP Purge Request until either NHRP Purge Request is acknowledged or until the holding time of the information being purged has expired. Retransmission strategies for NHRP Purge Requests are a local matter.

When a station receives an NHRP Purge Request, it MUST discard any previously cached information that matches the information in the CIEs.

An NHRP Purge Reply MUST be returned for the NHRP Purge Request even if the station does not have a matching cache entry assuming that the "N" bit is off in the NHRP Purge Request.

If the station wishes to reestablish communication with the destination shortly after receiving an NHRP Purge Request, it should make an authoritative Next Hop Resolution Request in order to avoid any stale cache entries that might be present in intermediate NHSS (See [section 6.2.2](#)). It is recommended that authoritative Next Hop Resolution Requests be made for the duration of the holding time of the old information.

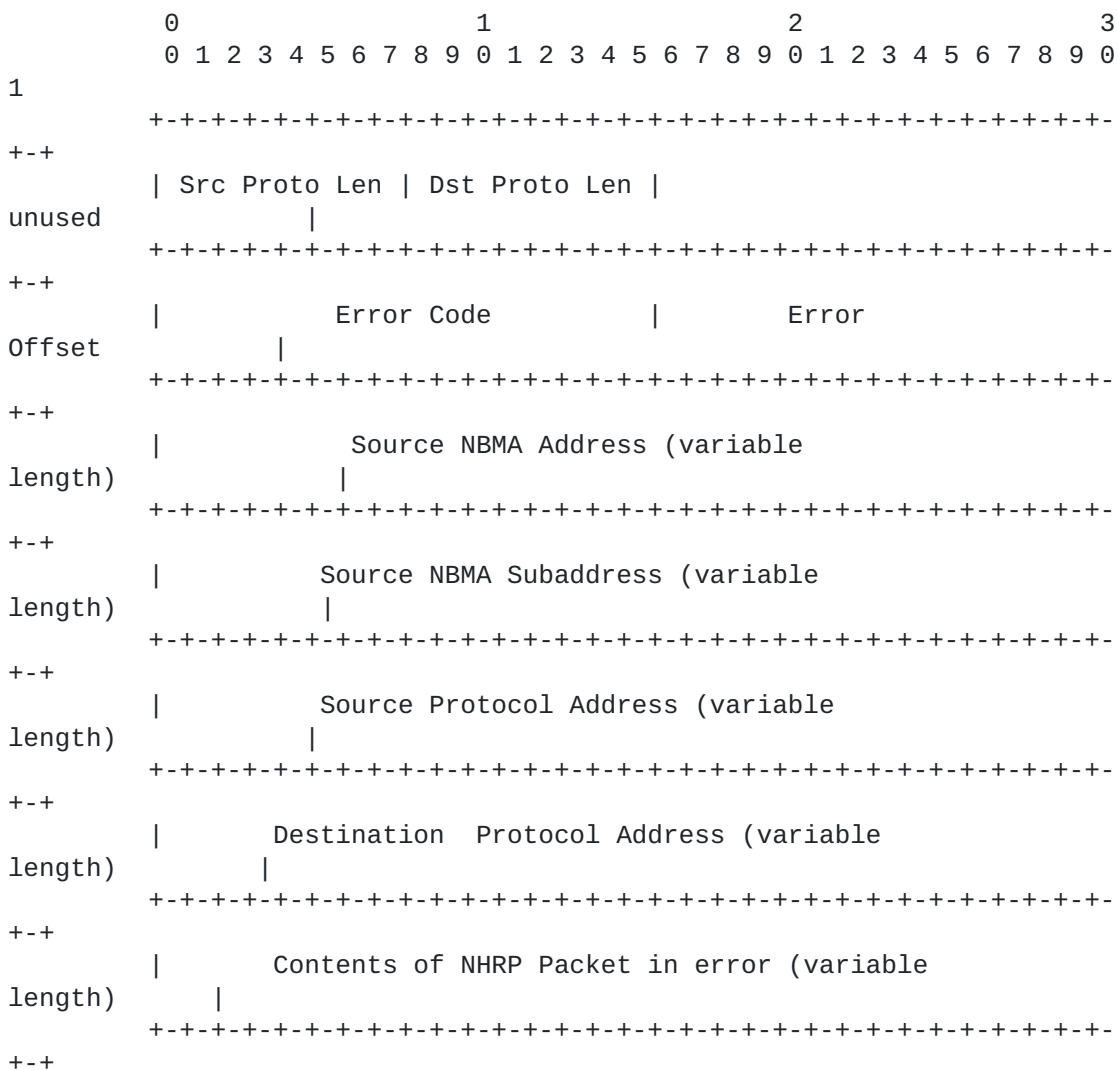
5.2.6 NHRP Purge Reply

The NHRP Purge Reply packet is sent in order to assure the sender of an NHRP Purge Request that all cached information of the specified type has been purged from the station sending the reply. The NHRP Purge Reply has a type code of 6.

An NHRP Purge Reply is formed from an NHRP Purge Request by merely changing the type code in the request to 6. The packet is then returned to the requester after filling in the appropriate extensions if they exist.

5.2.7 NHRP Error Indication

The NHRP Error Indication is used to convey error indications to the sender of an NHRP packet. It has a type code of 7. The Mandatory Part has the following format:



Src Proto Len
This field holds the length in octets of the Source Protocol Address.

Dst Proto Len

This field holds the length in octets of the Destination Protocol Address.

Error Code

An error code indicating the type of error detected, chosen from the following list:

1 - Unrecognized Extension

When the Compulsory bit of an extension in NHRP packet is set, the NHRP packet cannot be processed unless the extension has been processed. The responder MUST return an NHRP Error Indication of type Unrecognized Extension if it is incapable

of

processing the extension. However, if a transit NHS (one which is not going to generate a reply) detects an unrecognized extension, it SHALL ignore the extension.

2 - Subnetwork ID Mismatch

This error occurs when the current station receives an NHRP packet whose NBMA subnetwork identifier matches none of the locally known identifiers for the NBMA subnetwork on which the packet is received.

3 - NHRP Loop Detected

A Loop Detected error is generated when it is determined that an NHRP packet is being forwarded in a loop.

6 - Protocol Address Unreachable

This error occurs when a packet is moving along the routed path and it reaches a point such that the protocol address of interest is not reachable.

7 - Protocol Error

A generic packet processing error has occurred (e.g., invalid version number, invalid protocol type, failed checksum, etc.)

8 - NHRP SDU Size Exceeded

If the SDU size of the NHRP packet exceeds the MTU size of the NBMA network then this error is returned.

9 - Invalid Extension

If an NHS finds an extension in a packet which is inappropriate for the packet type, an error is sent back to the sender with Invalid Extension as the code.

10- Invalid Next Hop Resolution Reply Received

If a client receives a Next Hop Resolution Reply for a Next Hop Resolution Request which it believes it did not make then an error packet is sent to the station making the reply with an error code of Invalid Reply Received.

11- Authentication Failure

If a received packet fails an authentication test then this

error is returned.

14- Hop Count Exceeded

The hop count which was specified in the Fixed Header of an NHRP message has been exceeded.

Error Offset

The offset in octets into the NHRP packet, starting at the NHRP Fixed Header, at which the error was detected.

Source NBMA Address

The Source NBMA address field is the address of the station which observed the error.

Source NBMA SubAddress

The Source NBMA subaddress field is the address of the station which observed the error. If the field's length as specified in ar\$sst1 is 0 then no storage is allocated for this address at all.

Source Protocol Address

This is the protocol address of the station which issued the Error packet.

Destination Protocol Address

This is the protocol address of the station which sent the packet which was found to be in error.

An NHRP Error Indication packet SHALL NEVER be generated in response to another NHRP Error Indication packet. When an NHRP Error Indication packet is generated, the offending NHRP packet SHALL be discarded. In no case should more than one NHRP Error Indication packet be generated for a single NHRP packet.

If an NHS sees its own Protocol and NBMA Addresses in the Source NBMA

and Source Protocol address fields of a transiting NHRP Error Indication packet then the NHS will quietly drop the packet and do nothing (this scenario would occur when the NHRP Error Indication packet was itself in a loop).

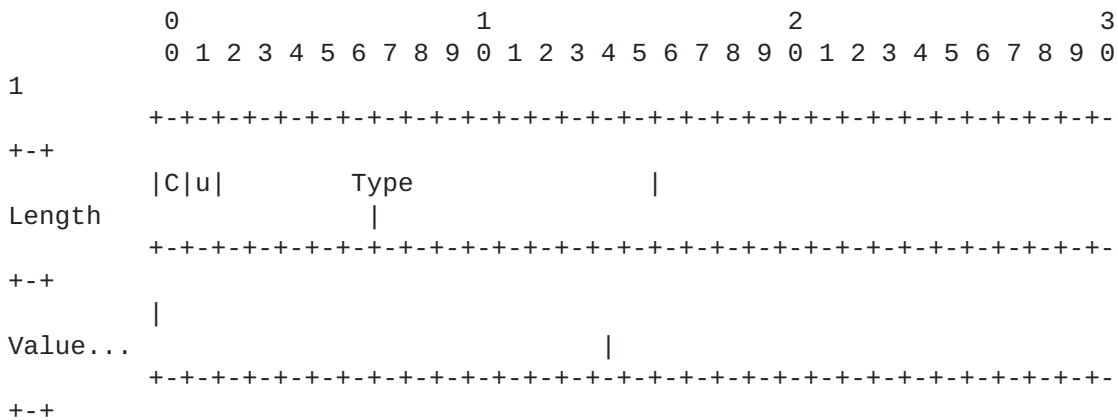
Note that no extensions may be added to an NHRP Error Indication.

5.3 Extensions Part

The Extensions Part, if present, carries one or more extensions in {Type, Length, Value} triplets. Extensions are only present in a "reply" if they were present in the corresponding "request";

therefore, minimal NHRP client implementations which do not also act as an NHS and do not transmit extensions need not be able to receive extensions. The previous statement is not intended to preclude the creation of NHS-only extensions which might be added to and removed from NHRP packets by the same NHS; such extensions MUST not be propagated to clients. An implementation that is incapable of processing extensions SHALL return an NHRP Error Indication of type Unrecognized Extension when it receives an NHRP packet containing extensions.

Extensions have the following format:



C
"Compulsory." If clear, and the NHS does not recognize the type code, the extension may safely be ignored. If set, and the NHS does not recognize the type code, the NHRP "request" is considered to be in error. (See below for details.)

u
Unused and must be set to zero.

Type
The extension type code (see below). The extension type is not qualified by the Compulsory bit, but is orthogonal to it.

Length
The length in octets of the value (not including the Type and Length fields; a null extension will have only an extension header and a length of zero).

When extensions exist, the extensions list is terminated by the Null TLV, having Type = 0 and Length = 0.

Extensions may occur in any order, but any particular extension type (except for the vendor-private extension) may occur only once in an NHRP packet. The vendor-private extension may occur multiple times

in a packet in order to allow for extensions which do not share the same vendor ID to be represented.

Luciani, Katz, Piscitello, Cole
32]

[Page

The Compulsory bit provides for a means to add to the extension set. If the bit is set, the NHRP message cannot be properly processed by the station responding to the message (e.g., the station that would issue a Next Hop Resolution Reply in response to a Next Hop Resolution Request) without processing the extension. As a result, the responder MUST return an NHRP Error Indication of type Unrecognized Extension. If the Compulsory bit is clear then the extension can be safely ignored; however, if an ignored extension is in a "request" then it MUST be returned, unchanged, in the corresponding "reply" packet type.

If a transit NHS (one which is not going to generate a "reply") detects an unrecognized extension, it SHALL ignore the extension.

If

the Compulsory bit is set, the transit NHS MUST NOT cache the information contained in the packet and MUST NOT identify itself as an egress router (in the Forward Record or Reverse Record extensions). Effectively, this means, if a transit NHS encounters

an

extension which it cannot process and which has the Compulsory bit set then that NHS MUST NOT participate in any way in the protocol exchange other than acting as a forwarding agent.

Use of NHRP extension Types in the range 8192 to 16383 are reserved for research or use in other protocol development and will be administered by IANA.

5.3.0 The End Of Extensions

Compulsory = 1
Type = 0
Length = 0

When extensions exist, the extensions list is terminated by the End Of Extensions/Null TLV.

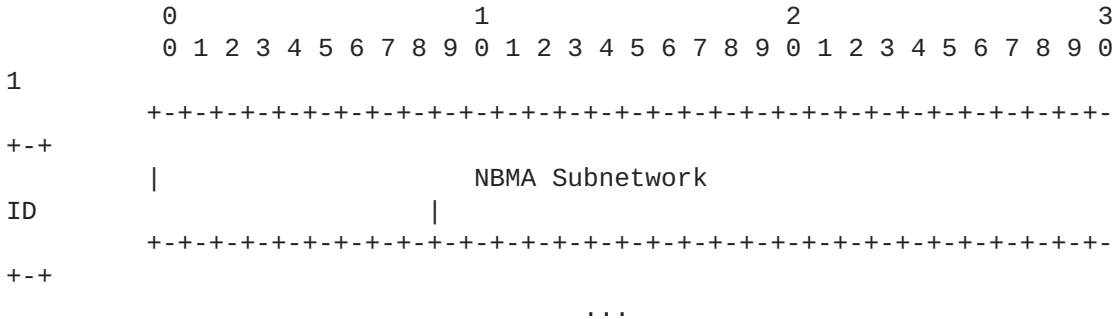
5.3.1 Extension with Type 1 not assigned.

5.3.2 NBMA Subnetwork ID Extension

Compulsory = 1
Type = 2
Length = variable

This extension is used to carry one or more identifiers for the NBMA subnetwork. This can be used as a validity check to ensure that an NHRP packet does not leave a particular NBMA subnetwork. The extension is placed in a "request" packet with an ID value of zero. The first NHS along the routed path fills in the field with the identifier(s) for the NBMA subnetwork.

Multiple NBMA Subnetwork IDs may be used as a transition mechanism while NBMA Subnetworks are being split or merged.



Each identifier consists of a 32 bit globally unique value assigned to the NBMA subnetwork. This value may be chosen from the internetworking layer address space administered by the operators of the NBMA subnetwork if such an address can fit into a 32 bit field. This value is used for identification only, not for routing or any other purpose.

Each NHS processing a "request" or "reply" SHALL verify these values.

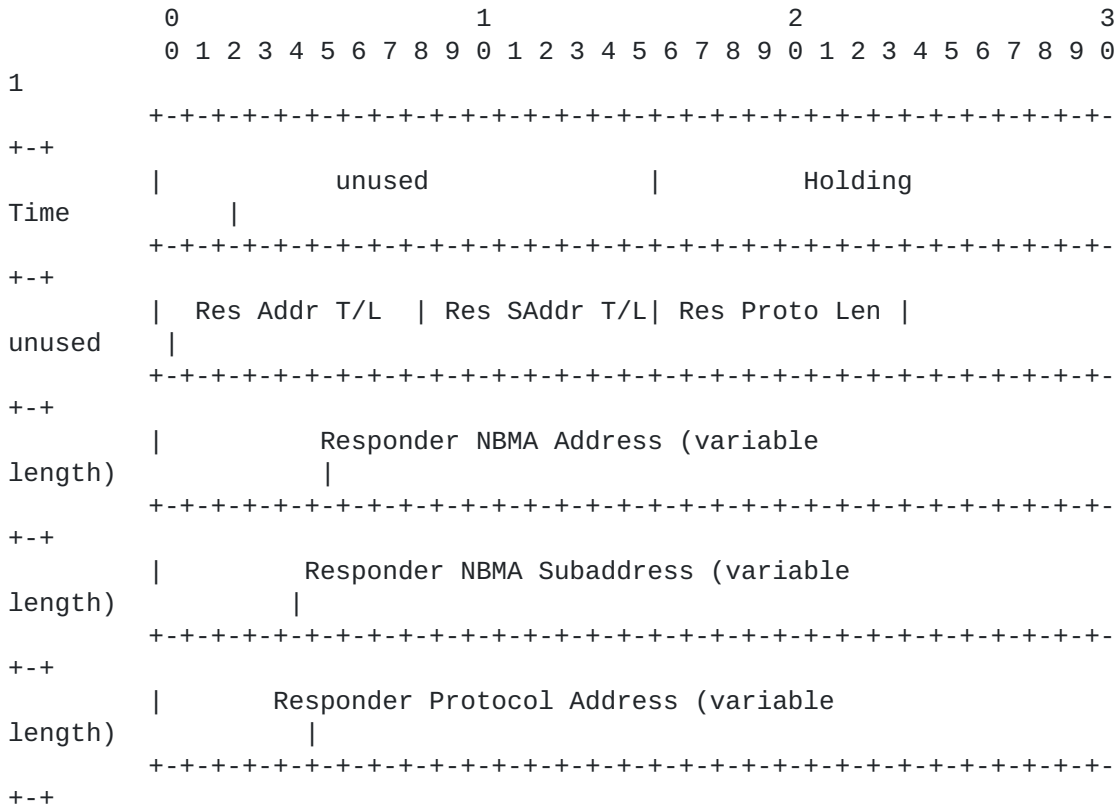
If the value is not zero and none of the values matches the NHS's NBMA Subnetwork ID, the NHS SHALL return an NHRP Error Indication to the entity identified in Source Protocol Address if the packet type is a "request" and to the Destination Protocol Address if the packet type is a "reply". The error indicated in this case is "Subnetwork ID Mismatch". The packet is discarded by the station sending the NHRP Error Indication.

When an NHS is building a "reply" and the NBMA Subnetwork ID extension is present in the correspond "request" then the NBMA Subnetwork ID extension SHALL be copied from the "request" to the "reply".

5.3.3 Responder Address Extension

- Compulsory = 1
- Type = 3
- Length = variable

This extension is used to determine the address of the NHRP responder; i.e., the entity that generates the appropriate "reply" packet for a given "request" packet. In the case of an Next Hop Resolution Request, the station responding may be different (in the case of cached replies) than the system identified in the Next Hop field of the Next Hop Resolution Reply. Further, this extension may aid in detecting loops in the NHRP forwarding path.



Holding Time

The Holding Time field specifies the number of seconds for which the NBMA information is considered to be valid. Cached information SHALL be discarded when the holding time expires.

Res Addr T/L

Type & length of the responder NHS's NBMA address interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0003 for ATM). When the address length is specified as 0 no storage is allocated for the address.

Res SAddr T/L

Type & length of responder NHS's NBMA subaddress interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0015 for ATM makes the address an E.164 and the subaddress an ATM Forum NSAP address). When an NBMA technology has no concept of a subaddress, the subaddress is always null with a length of 0. When the address length is specified as 0 no storage is allocated for the address.

Res Proto Len

This field holds the length in octets of responding NHS's Protocol Address.

Responder NBMA Address

This is the NBMA address of the responding NHS.

Responder NBMA SubAddress

This is the NBMA subaddress of the responding NHS.

Responder Protocol Address

This is the Protocol Address of responding NHS.

If a "requester" desires this information, the "requester" SHALL

include this extension with a value of zero. Note that this implies that no storage is allocated for the Holding Time and Type/Length fields until the "Value" portion of the extension is filled out.

If an NHS is generating a "reply" packet in response to a "request" containing this extension, the NHS SHALL include this extension, containing its protocol address in the "reply". If an NHS has more than one protocol address, it SHALL use the same protocol address consistently in all of the Responder Address, Forward NHS Record, and Reverse NHS Record extensions. The choice of which of several protocol address to include in this extension is a local matter.

If an NHRP Next Hop Resolution Reply packet being forwarded by an NHS contains a protocol address of that NHS in the Responder Address Extension then that NHS SHALL generate an NHRP Error Indication of type "NHRP Loop Detected" and discard the Next Hop Resolution Reply.

If an NHRP Next Hop Resolution Reply packet is being returned by an intermediate NHS based on cached data, it SHALL place its own address in this extension (differentiating it from the address in the Next Hop field).

5.3.4 NHRP Forward Transit NHS Record Extension

Compulsory = 1
Type = 4
Length = variable

The NHRP Forward Transit NHS record contains a list of transit NHSs through which a "request" has traversed. Each NHS SHALL append to the extension a Forward Transit NHS element (as specified below) containing its Protocol address The extension length field and the ar\$chksum fields SHALL be adjusted appropriately.

The responding NHS, as described in [Section 5.3.3](#), SHALL NOT update this extension.

In addition, NHSs that are willing to act as egress routers for packets from the source to the destination SHALL include information about their NBMA Address.

The Forward Transit NHS element has the following form:

NHS NBMA Address

This is the NBMA address of the transit NHS.

NHS NBMA SubAddress

This is the NBMA subaddress of the transit NHS.

NHS Protocol Address

This is the Protocol Address of the transit NHS.

If a "requester" wishes to obtain this information, it SHALL include

Holding Time

The Holding Time field specifies the number of seconds for which the NBMA information is considered to be valid. Cached information

SHALL be discarded when the holding time expires.

NHS Addr T/L

Type & length of the responding NHS's NBMA address interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0003 for ATM). When the address length is specified as 0 no storage is allocated for the address.

NHS SAddr T/L

Type & length of the responding NHS's NBMA subaddress interpreted in the context of the 'address family number'[6] indicated by ar\$afn (e.g., ar\$afn=0x0015 for ATM makes the address an E.164 and the subaddress an ATM Forum NSAP address). When an NBMA

technology

has no concept of a subaddress the subaddress is always null with a length of 0. When the address length is specified as 0 no storage is allocated for the address.

NHS Proto Len

This field holds the length in octets of the transit NHS's Protocol Address.

NHS NBMA Address

This is the NBMA address of the transit NHS.

NHS NBMA SubAddress

This is the NBMA subaddress of the transit NHS.

NHS Protocol Address

This is the Protocol Address of the transit NHS.

If a "requester" wishes to obtain this information, it SHALL include this extension with a length of zero. Note that this implies that no storage is allocated for the Holding Time and Type/Length fields until the "Value" portion of the extension is filled out.

If an NHS has more than one Protocol address, it SHALL use the same Protocol address consistently in all of the Responder Address, Forward NHS Record, and Reverse NHS Record extensions. The choice of which of several Protocol addresses to include in this extension is a local matter.

If a "reply" that is being forwarded by an NHS contains the Protocol

Address of that NHS in one of the Reverse Transit NHS elements then the NHS SHALL generate an NHRP Error Indication of type "NHRP Loop Detected" and discard the "reply".

Note that this information may be cached at intermediate NHSS; if so, the cached value SHALL be used when generating a reply.

5.3.6 NHRP QoS Extension

Compulsory = 0
Type = 6
Length = variable

The NHRP QoS Extension is carried in Next Hop Resolution Request packets to indicate the desired QoS of the path to the indicated destination. This information may be used to help select the appropriate NBMA Next Hop.

It may also be carried in NHRP Register Request packets to indicate the QoS to which the registration applies.

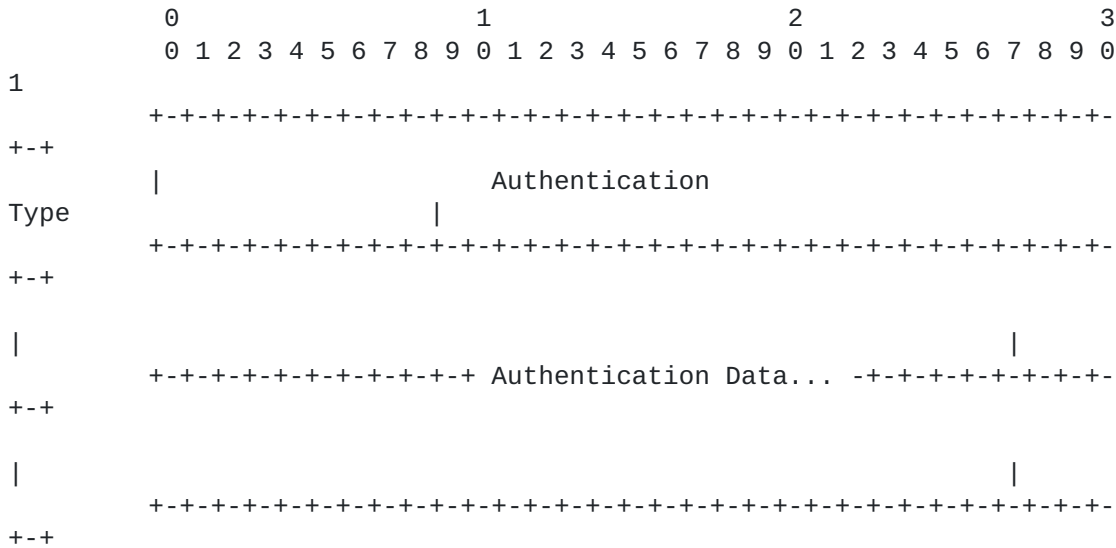
The syntax and semantics of this extension are TBD; alignment with resource reservation may be useful.

5.3.7 NHRP Authentication Extension

Compulsory = 1
Type = 7
Length = variable

The NHRP Authentication Extension is carried in NHRP packets to convey authentication information between NHRP speakers. The Authentication Extension may be included in any NHRP "request" or "reply".

Except in the case of an NHRP Registration Request/Reply Authentication is done pairwise on an NHRP hop-by-hop basis; i.e., the authentication extension is regenerated at each hop. In the case of an NHRP Registration Request/Reply, the Authentication is checked on an end-to-end basis rather than hop-by-hop. If a received packet fails the authentication test, the station SHALL generate an Error Indication of type "Authentication Failure" and discard the packet. Note that one possible authentication failure is the lack of an Authentication Extension; the presence or absence of the Authentication Extension is a local matter.



The Authentication Type field identifies the authentication method in use. Currently assigned values are:

- 1 - Cleartext Password
- 2 - Keyed MD5

All other values are reserved.

The Authentication Data field contains the type-specific authentication information.

In the case of Cleartext Password Authentication, the Authentication Data consists of a variable length password.

In the case of Keyed MD5 Authentication, the Authentication Data contains the 16 byte MD5 digest of the entire NHRP packet, including the encapsulated protocol's header, with the authentication key appended to the end of the packet. The authentication key is not transmitted with the packet.

Distribution of authentication keys is outside the scope of this document.

5.3.8 NHRP Vendor-Private Extension

Compulsory = 0
Type = 8
Length = variable

The NHRP Vendor-Private Extension is carried in NHRP packets to convey vendor-private information or NHRP extensions between NHRP speakers.

```

      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
  +-+-+-+-+-+-+-+-+
+-+
|           Vendor ID           |
Data....  |
  +-+-+-+-+-+-+-+-+
+-+

```

Vendor ID

802 Vendor ID as assigned by the IEEE [6]

Data

The remaining octets after the Vendor ID in the payload are vendor-dependent data.

This extension may be added to any "request" or "reply" packet and it

is the only extension that may be included multiple times. If the receiver does not handle this extension, or does not match the

Vendor

ID in the extension then the extension may be completely ignored by the receiver. If a Vendor Private Extension is included in a "request" then it must be copied in the corresponding "reply".

5.3.9 Extension with Type 9 not assigned.

6. Protocol Operation

In this section, we discuss certain operational considerations of NHRP.

6.1 Router-to-Router Operation

In practice, the initiating and responding stations may be either hosts or routers. However, there is a possibility under certain conditions that a stable routing loop may occur if NHRP is used between two routers. In particular, attempting to establish an NHRP path across a boundary where information used in route selection is lost may result in a routing loop. Such situations include the loss of BGP path vector information, the interworking of multiple routing protocols with dissimilar metrics (e.g, RIP and OSPF), etc. In such circumstances, NHRP should not be used. This situation can be avoided if there are no "back door" paths between the entry and egress router outside of the NBMA subnetwork. Protocol mechanisms to relax these restrictions are under investigation.

In general it is preferable to use mechanisms, if they exist, in routing protocols to resolve the egress point when the destination lies outside of the NBMA subnetwork, since such mechanisms will be more tightly coupled to the state of the routing system and will probably be less likely to create loops.

6.2 Cache Management Issues

The management of NHRP caches in the source station, the NHS serving the destination, and any intermediate NHSs is dependent on a number of factors.

6.2.1 Caching Requirements

Source Stations

Source stations MUST cache all received Next Hop Resolution Replies

that they are actively using. They also must cache "incomplete" entries, i.e., those for which a Next Hop Resolution Request has been sent but which a Next Hop Resolution Reply has not been received. This is necessary in order to preserve the Request ID for retries, and provides the state necessary to avoid triggering Next Hop Resolution Requests for every data packet sent to the destination.

Source stations MUST purge expired information from their caches. Source stations MUST purge the appropriate cached information upon receipt of an NHRP Purge Request packet.

When a station has a co-resident client and NHS, the station may reply to Next Hop Resolution Requests with information which the station cached as a result of the station making its own Next Hop Resolution Requests and receiving its own Next Hop Resolution Replies as long as the station follows the rules for Transit NHSs as seen below.

Serving NHSs

The NHS serving the destination (the one which responds authoritatively to Next Hop Resolution Requests) SHOULD cache information about all Next Hop Resolution Requests to which it has responded if the information in the Next Hop Resolution Reply has the possibility of changing during its lifetime (so that an NHRP Purge Request packet can be sent). The NBMA information provided by the source station in the Next Hop Resolution Request may be cached for the duration of its holding time. This information is considered to be stable, since it identifies a station directly attached to the NBMA subnetwork. An example of unstable information is NBMA information derived from a routing table, where that routing table information has not been guaranteed to be stable through administrative means.

Transit NHSs

A Transit NHS (lying along the NHRP path between the source station and the responding NHS) may cache information contained in Next Hop Resolution Request packets that it forwards. A Transit NHS may cache information contained in Next Hop Resolution Reply packets that it forwards only if that Next Hop Resolution Reply has the

Stable (B) bit set. It MUST discard any cached information whose

Luciani, Katz, Piscitello, Cole
43]

[Page

holding time has expired. It may return cached information in response to non-authoritative Next Hop Resolution Requests only.

6.2.2 Dynamics of Cached Information

NBMA-Connected Destinations

NHRP's most basic function is that of simple NBMA address resolution of stations directly attached to the NBMA subnetwork. These mappings are typically very static, and appropriately chosen holding times will minimize problems in the event that the NBMA address of a station must be changed. Stale information will

cause

a loss of connectivity, which may be used to trigger an authoritative Next Hop Resolution Request and bypass the old data. In the worst case, connectivity will fail until the cache entry times out.

This applies equally to information marked in Next Hop Resolution Replies as being "stable" (via the "B" bit).

This also applies equally well to source stations that are routers as well as those which are hosts.

Note that the information carried in the Next Hop Resolution Request packet is always considered "stable" because it represents a station that is directly connected to the NBMA subnetwork.

Destinations Off of the NBMA Subnetwork

If the source of a Next Hop Resolution Request is a host and the destination is not directly attached to the NBMA subnetwork, and the route to that destination is not considered to be "stable,"

the

destination mapping may be very dynamic (except in the case of a subnetwork where each destination is only singly homed to the NBMA subnetwork). As such the cached information may very likely

become

stale. The consequence of stale information in this case will be

a

suboptimal path (unless the internetwork has partitioned or some other routing failure has occurred).

6.3 Use of the Prefix Length field of a CIE

A certain amount of care needs to be taken when using the Prefix Length field of a CIE, in particular with regard to the prefix length

advertised (and thus the size of the equivalence class specified by it). Assuming that the routers on the NBMA subnetwork are exchanging

routing information, it should not be possible for an NHS to create

a

black hole by advertising too large of a set of destinations, but suboptimal routing (e.g., extra internetwork layer hops through the

NBMA) can result. To avoid this situation an NHS that wants to send the Prefix Length MUST obey the following rule:

The NHS examines the Network Layer Reachability Information (NLRI) associated with the route that the NHS would use to forward towards the destination (as specified by the Destination internetnetwork layer address in the Next Hop Resolution Request), and extracts from this NLRI the shortest address prefix such that: (a) the Destination internetnetwork layer address (from the Next Hop Resolution Request) is covered by the prefix, (b) the NHS does not have any routes with NLRI that forms a subset of what is covered by the prefix. The prefix may then be used in the CIE.

The Prefix Length field of the CIE should be used with restraint, in order to avoid NHRP stations choosing suboptimal transit paths when overlapping prefixes are available. This document specifies the use of the prefix length only when all the destinations covered by the prefix are "stable". That is, either:

- (a) All destinations covered by the prefix are on the NBMA network, or
- (b) All destinations covered by the prefix are directly attached to the NHRP responding station.

Use of the Prefix Length field of the CIE in other circumstances is outside the scope of this document.

6.4 Domino Effect

One could easily imagine a situation where a router, acting as an ingress station to the NBMA subnetwork, receives a data packet, such that this packet triggers an Next Hop Resolution Request. If the router forwards this data packet without waiting for an NHRP transit path to be established, then when the next router along the path receives the packet, the next router may do exactly the same - originate its own Next Hop Resolution Request (as well as forward the packet). In fact such a data packet may trigger Next Hop Resolution Request generation at every router along the path through an NBMA subnetwork. We refer to this phenomena as the NHRP "domino" effect.

The NHRP domino effect is clearly undesirable. At best it may result in excessive NHRP traffic. At worst it may result in an excessive number of virtual circuits being established unnecessarily. Therefore, it is important to take certain measures to avoid or

suppress this behavior. NHRP implementations for NHSs MUST provide
a mechanism to address this problem. One possible strategy to address
this problem would be to configure a router in such a way that Next
Hop Resolution Request generation by the router would be driven only

by the traffic the router receives over its non-NBMA interfaces (interfaces that are not attached to an NBMA subnetwork). Traffic received by the router over its NBMA-attached interfaces would not trigger NHRP Next Hop Resolution Requests. Such a router avoids the NHRP domino effect through administrative means.

7. NHRP over Legacy BMA Networks

There would appear to be no significant impediment to running NHRP over legacy broadcast subnetworks. There may be issues around running NHRP across multiple subnetworks. Running NHRP on broadcast media has some interesting possibilities; especially when setting up a cut-through for inter-ELAN inter-LIS/LAG traffic when one or both end stations are legacy attached. This use for NHRP requires further research.

8. Security Considerations

As in any routing protocol, there are a number of potential security attacks possible. Plausible examples include denial-of-service attacks, and masquerade attacks using register and purge packets. The use of authentication on all packets is recommended to avoid such attacks.

The authentication schemes described in this document are intended to allow the receiver of a packet to validate the identity of the sender; they do not provide privacy or protection against replay attacks.

Detailed security analysis of this protocol is for further study.

9. Discussion

The result of an Next Hop Resolution Request depends on how routing is configured among the NHSS of an NBMA subnetwork. If the destination station is directly connected to the NBMA subnetwork and the the routed path to it lies entirely within the NBMA subnetwork, the Next Hop Resolution Replies always return the NBMA address of the destination station itself rather than the NBMA address of some egress router. On the other hand, if the routed path exits the NBMA subnetwork, NHRP will be unable to resolve the NBMA address of the destination, but rather will return the address of the egress router.

For destinations outside the NBMA subnetwork, egress routers and routers in the other subnetworks should exchange routing information

so that the optimal egress router may be found.

Luciani, Katz, Piscitello, Cole
46]

[Page

In addition to NHSs, an NBMA station could also be associated with one or more regular routers that could act as "connectionless servers" for the station. The station could then choose to resolve the NBMA next hop or just send the packets to one of its connectionless servers. The latter option may be desirable if communication with the destination is short-lived and/or doesn't require much network resources. The connectionless servers could, of course, be physically integrated in the NHSs by augmenting them with internetwork layer switching functionality.

References

- [1] NBMA Address Resolution Protocol (NARP), Juha Heinanen and Ramesh Govindan, [RFC1735](#).
- [2] Address Resolution Protocol, David C. Plummer, [RFC 826](#).
- [3] Classical IP and ARP over ATM, Mark Laubach, [RFC 1577](#).
- [4] Transmission of IP datagrams over the SMDS service, J. Lawrence and D. Piscitello, [RFC 1209](#).
- [5] Protocol Identification in the Network Layer, ISO/IEC TR 9577:1990.
- [6] Assigned Numbers, J. Reynolds and J. Postel, [RFC 1700](#).
- [7] Multiprotocol Encapsulation over ATM Adaptation Layer 5, J. Heinanen, [RFC1483](#).
- [8] Multiprotocol Interconnect on X.25 and ISDN in the Packet Mode, A. Malis, D. Robinson, and R. Ullmann, [RFC1356](#).
- [9] Multiprotocol Interconnect over Frame Relay, T. Bradley, C. Brown, and A. Malis, [RFC1490](#).
- [10] "Local/Remote" Forwarding Decision in Switched Data Link Subnetworks, Yakov Rekhter, Dilip Kandlur, RFCxxxx.

Acknowledgments

We would like to thank Juha Heinanen of Telecom Finland and Ramesh Govindan of ISI for their work on NBMA ARP and the original NHRP draft, which served as the basis for this work. Russell Gardo of

IBM, John Burnett of Adaptive, Dennis Ferguson of ANS, Joel Halpern

Luciani, Katz, Piscitello, Cole
47]

[Page

of Newbridge, Paul Francis of NTT, Tony Li and Yakov Rekhter of cisco, and Grenville Armitage of Bellcore should also be acknowledged

for comments and suggestions that improved this work substantially. We would also like to thank the members of the ION working group of the IETF, whose review and discussion of this document have been invaluable.

Authors' Addresses

James V. Luciani
Bay Networks
3 Federal Street
Mail Stop: BL3-04
Billerica, MA 01821
Phone: +1 508 439 4737
Email: luciani@baynetworks.com

Dave Katz
cisco Systems
170 W. Tasman Dr.
San Jose, CA 95134 USA
Phone: +1 408 526 8284
Email: dkatz@cisco.com

David Piscitello
Core Competence
1620 Tuckerstown Road
Dresher, PA 19025 USA
Phone: +1 215 830 0692
Email: dave@corecom.com

Bruce Cole
cisco Systems
170 W. Tasman Dr.
San Jose, CA 95134 USA
Phone: +1 408 526 4000
Email: bcole@cisco.com

