

Next Steps in Signaling
Internet-Draft
Expires: January 26, 2007

H. Schulzrinne
Columbia U.
R. Hancock
Siemens/RMR
July 25, 2006

GIST: General Internet Signaling Transport
draft-ietf-nsis-ntlp-10

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with [Section 6 of BCP 79](#).

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

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

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

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

This Internet-Draft will expire on January 26, 2007.

Copyright Notice

Copyright (C) The Internet Society (2006).

Abstract

This document specifies protocol stacks for the routing and transport of per-flow signaling messages along the path taken by that flow through the network. The design uses existing transport and security protocols under a common messaging layer, the General Internet Signaling Transport (GIST), which provides a universal service for diverse signaling applications. GIST does not handle signaling application state itself, but manages its own internal state and the

Internet-Draft

GIST

July 2006

configuration of the underlying transport and security protocols to enable the transfer of messages in both directions along the flow path. The combination of GIST and the lower layer transport and security protocols provides a solution for the base protocol component of the "Next Steps in Signaling" framework.

Table of Contents

1.	Introduction	4
1.1.	Restrictions on Scope	5
2.	Requirements Notation and Terminology	6
3.	Design Overview	9
3.1.	Overall Design Approach	9
3.2.	Modes and Messaging Associations	10
3.3.	Message Routing Methods	11
3.4.	GIST Messages	13
3.5.	Signaling Sessions	14
3.6.	Signaling Applications and NSLPIDs	15
3.7.	Example of Operation	15
4.	GIST Processing Overview	18
4.1.	GIST Service Interface	18
4.2.	GIST State	20
4.3.	Basic GIST Message Processing	21
4.4.	Routing State and Messaging Association Maintenance	27
5.	Message Formats and Transport	34
5.1.	GIST Messages	34
5.2.	Information Elements	36
5.3.	D-mode Transport	40
5.4.	C-mode Transport	43
5.5.	Message Type/Encapsulation Relationships	45
5.6.	Error Message Processing	46
5.7.	Messaging Association Setup	47
5.8.	Specific Message Routing Methods	51
6.	Formal Protocol Specification	56
6.1.	Node Processing	57
6.2.	Query Node Processing	59
6.3.	Responder Node Processing	62
6.4.	Messaging Association Processing	65
7.	Advanced Protocol Features	69
7.1.	Route Changes and Local Repair	69
7.2.	NAT Traversal	75
7.3.	Interaction with IP Tunnelling	79

7.4.	IPv4-IPv6 Transition and Interworking	79
8.	Security Considerations	81
8.1.	Message Confidentiality and Integrity	81
8.2.	Peer Node Authentication	82
8.3.	Routing State Integrity	82

Internet-Draft

GIST

July 2006

8.4.	Denial of Service Prevention	84
8.5.	Requirements on Cookie Mechanisms	85
8.6.	Security Protocol Selection Policy	87
8.7.	Residual Threats	88
9.	IANA Considerations	89
10.	Acknowledgements	94
11.	References	95
11.1.	Normative References	95
11.2.	Informative References	95
Appendix A.	Bit-Level Formats and Error Messages	98
A.1.	The GIST Common Header	98
A.2.	General Object Format	99
A.3.	GIST TLV Objects	100
A.4.	Errors	109
Appendix B.	API between GIST and Signaling Applications	117
B.1.	SendMessage	117
B.2.	RecvMessage	119
B.3.	MessageStatus	120
B.4.	NetworkNotification	121
B.5.	SetStateLifetime	121
B.6.	InvalidateRoutingState	122
Appendix C.	Example Routing State Table and Handshake Message Sequence	123
Appendix D.	Change History	125
D.1.	Changes In Version -10	125
D.2.	Changes In Version -09	128
D.3.	Changes In Version -08	129
D.4.	Changes In Version -07	131
D.5.	Changes In Version -06	132
D.6.	Changes In Version -05	133
D.7.	Changes In Version -04	134
D.8.	Changes In Version -03	135
D.9.	Changes In Version -02	136
D.10.	Changes In Version -01	137
Authors' Addresses	140
Intellectual Property and Copyright Statements	141

Internet-Draft

GIST

July 2006

1. Introduction

Signaling involves the manipulation of state held in network elements. 'Manipulation' could mean setting up, modifying and tearing down state; or it could simply mean the monitoring of state which is managed by other mechanisms.

This specification concentrates mainly on path-coupled signaling, which involves network elements which are located on the path taken by a particular data flow, possibly including but not limited to the flow endpoints. Indeed, there are almost always more than two participants in a path-coupled signaling session, although there is no need for every node on the path to participate. Path-coupled signaling thus excludes end-to-end higher-layer application signaling.

In the context of path-coupled signaling, examples of state management include network resource reservation, firewall configuration, and state used in active networking; examples of state monitoring are the discovery of instantaneous path properties, such as available bandwidth or cumulative queuing delay. Each of these different uses of signaling is referred to as a signaling application.

Every signaling application requires a set of state management rules, as well as protocol support to exchange messages along the data path. Several aspects of this protocol support are common to all or a large number of signaling applications, and hence can be developed as a common protocol. The NSIS framework given in [\[26\]](#) provides a

rationale for a function split between the common and application specific protocols, and gives outline requirements for the former, the 'NSIS Transport Layer Protocol' (NTLP). The application specific protocols are referred to as 'NSIS Signaling Layer Protocols' (NSLPs), and are defined in separate documents.

This specification provides a concrete solution for the NTLP. It is based on the use of existing transport and security protocols under a common messaging layer, the General Internet Signaling Transport (GIST). GIST does not handle signaling application state itself; in that crucial respect, it differs from application signaling protocols such as SIP, RTSP, and the control component of FTP. Instead, GIST manages its own internal state and the configuration of the underlying transport and security protocols to ensure the transfer of signaling messages on behalf of signaling applications in both directions along the flow path.

The structure of this specification is as follows. [Section 2](#) defines terminology, and [Section 3](#) gives an informal overview of the protocol

design principles and operation. The normative specification is contained mainly in [Section 4](#) to [Section 8](#). [Section 4](#) describes the message sequences and [Section 5](#) their format and contents. Note that the detailed bit formats are given in [Appendix A](#). The protocol operation is captured in the form of state machine language in [Section 6](#). [Section 7](#) describes some more advanced protocol features and security considerations are contained in [Section 8](#). In addition, [Section 9](#) gives the IANA considerations, [Appendix B](#) describes an abstract API for the service which GIST provides to signaling applications, and [Appendix C](#) provides an example message flow.

[1.1](#). Restrictions on Scope

This section briefly lists some important restrictions on GIST applicability and functionality. In some cases, these are implicit consequences of the functionality split developed in the NSIS framework; in others, they are restrictions on the types of scenario in which GIST can operate correctly.

Flow splitting: In some cases, e.g. where packet-level load sharing has been implemented, the path taken by a single flow in the network may not be well defined. If this is the case, GIST cannot

route signaling meaningfully. In some circumstances, GIST implementations could detect this condition, but even this cannot be guaranteed.

Multicast: GIST does not handle multicast flows. This includes classical IP multicast and any of the small group multicast schemes recently proposed.

Legacy NATs: GIST messages will generally pass through NATs, but unless the NAT is GIST-aware, any addressing data carried in the payload will not be handled correctly. There is a dual problem of whether the GIST peers either side of the boundary can work out how to address each other, and whether they can work out what translation to apply to the signaling packet payloads. The fundamental problem is that GIST messages contain three or four interdependent addresses which all have to be consistently translated, and existing generic NAT traversal techniques such as STUN [23] or TURN [24] can process only two. Appropriate behaviour for a GIST-aware NAT is discussed in [Section 7.2](#).

[2.](#) Requirements Notation and Terminology

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [RFC 2119](#) [2].

The terminology used in this specification is defined in this section. The basic entities relevant at the GIST level are shown in Figure 1. In particular, this diagram distinguishes the different address types as being associated with a flow (end-to-end addresses) or signaling (addresses of adjacent signaling peers).

Source	GIST (adjacent) peer nodes	Destination
IP address	IP addresses = Signaling	IP address

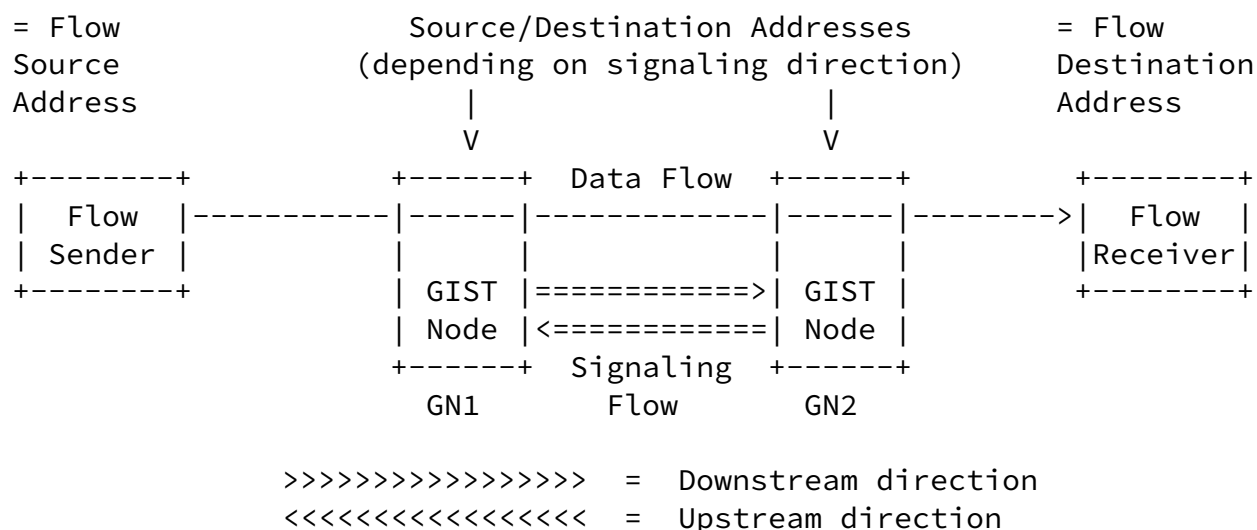


Figure 1: Basic Terminology

[Data] Flow: A set of packets identified by some fixed combination of header fields. Flows are unidirectional; a bidirectional communication is considered a pair of unidirectional flows.

Session: A single application layer flow of information for which some state information is to be manipulated or monitored. See [Section 3.5](#) for further detailed discussion.

Session Identifier (SID): A long, opaque identifier for a session.

[Flow] Sender: The node in the network which is the source of the packets in a flow. A sender could be a host, or a router if for example the flow is actually an aggregate.

[Flow] Receiver: The node in the network which is the sink for the packets in a flow.

Downstream: In the same direction as the data flow.

Upstream: In the opposite direction to the data flow.

GIST Node: Any node along the data path supporting GIST, regardless

of what signaling applications it supports.

[Adjacent] Peer: The next node along the data path, in the upstream or downstream direction, with which a GIST node explicitly interacts. The GIST peer discovery mechanisms implicitly determine whether two nodes will be adjacent. It is possible for adjacencies to skip over intermediate nodes which decide not to take part in the signaling exchange at the NTLP layer; even if such nodes process parts of the signaling messages, they store no state about the session and are never visible at the GIST level to the nodes on either side.

Datagram Mode (D-mode): A mode of sending GIST messages between nodes without using any transport layer state or security protection. Datagram mode uses UDP encapsulation, with source and destination IP addresses derived either from the flow definition or previously discovered adjacency information.

Connection Mode (C-mode): A mode of sending GIST messages directly between nodes using point-to-point messaging associations (see below). Connection mode allows the re-use of existing transport and security protocols where such functionality is required.

Messaging Association (MA): A single connection between two explicitly identified GIST adjacent peers, i.e. between a given signaling source and destination address. A messaging association may use a specific transport protocol and known ports. If security protection is required, it may use a specific network layer security association, or use a transport layer security association internally. A messaging association is bidirectional; signaling messages can be sent over it in either direction, and can refer to flows of either direction.

Message Routing Method (MRM): There can be different algorithms for discovering the route that signaling messages should take. These are referred to as message routing methods, and GIST supports alternatives within a common protocol framework. See [Section 3.3](#).

Message Routing Information (MRI): The set of data item values which

is used to route a signaling message according to a particular MRM; for example, for routing along a flow path, the MRI includes flow source and destination addresses, protocol and port numbers. See [Section 3.3](#).

Transfer Attributes: A description of the requirements which a signaling application has for the delivery of a particular message; for example, whether the message should be delivered reliably. See [Section 4.1.2](#).

Router Alert Option (RAO): An option that can be included in IP v4 and v6 headers to assist in the packet interception process; see [\[1\]](#) and [\[5\]](#).

[3.](#) Design Overview

[3.1.](#) Overall Design Approach

The generic requirements identified in the NSIS framework [\[26\]](#) for transport of signaling messages are essentially two-fold:

Routing: Determine how to reach the adjacent signaling node along each direction of the data path (the GIST peer), and if necessary explicitly establish addressing and identity information about that peer;

Transport: Deliver the signaling information to that peer.

To meet the routing requirement, one possibility is for the node to use local routing state information to determine the identity of the GIST peer explicitly. GIST defines a three-way handshake which probes the network to set up the necessary routing state between adjacent peers, during which signaling applications can also exchange data. Once the routing decision has been made, the node has to select a mechanism for transport of the message to the peer. GIST divides the transport problems into two categories, the easy and the difficult. It handles the easy cases internally, and uses well-understood transport protocols for the harder cases. Here, with details discussed later, "easy" messages are those that are sized well below the lowest maximum transmission unit (MTU) along a path, are infrequent enough not to cause concerns about congestion and flow control, and do not need security protection or guaranteed delivery.

In [\[26\]](#) all of these routing and transport requirements are assigned to a single notional protocol, the NSIS Transport Layer Protocol (NTLP). The strategy of splitting the transport problem leads to a layered structure for the NTLP, of a specialised GIST messaging layer running over standard transport and security protocols, as shown in Figure 2. This also shows GIST offering its services to upper layers at an abstract interface, the GIST API, further discussed in [Section 4.1](#).

July 2006

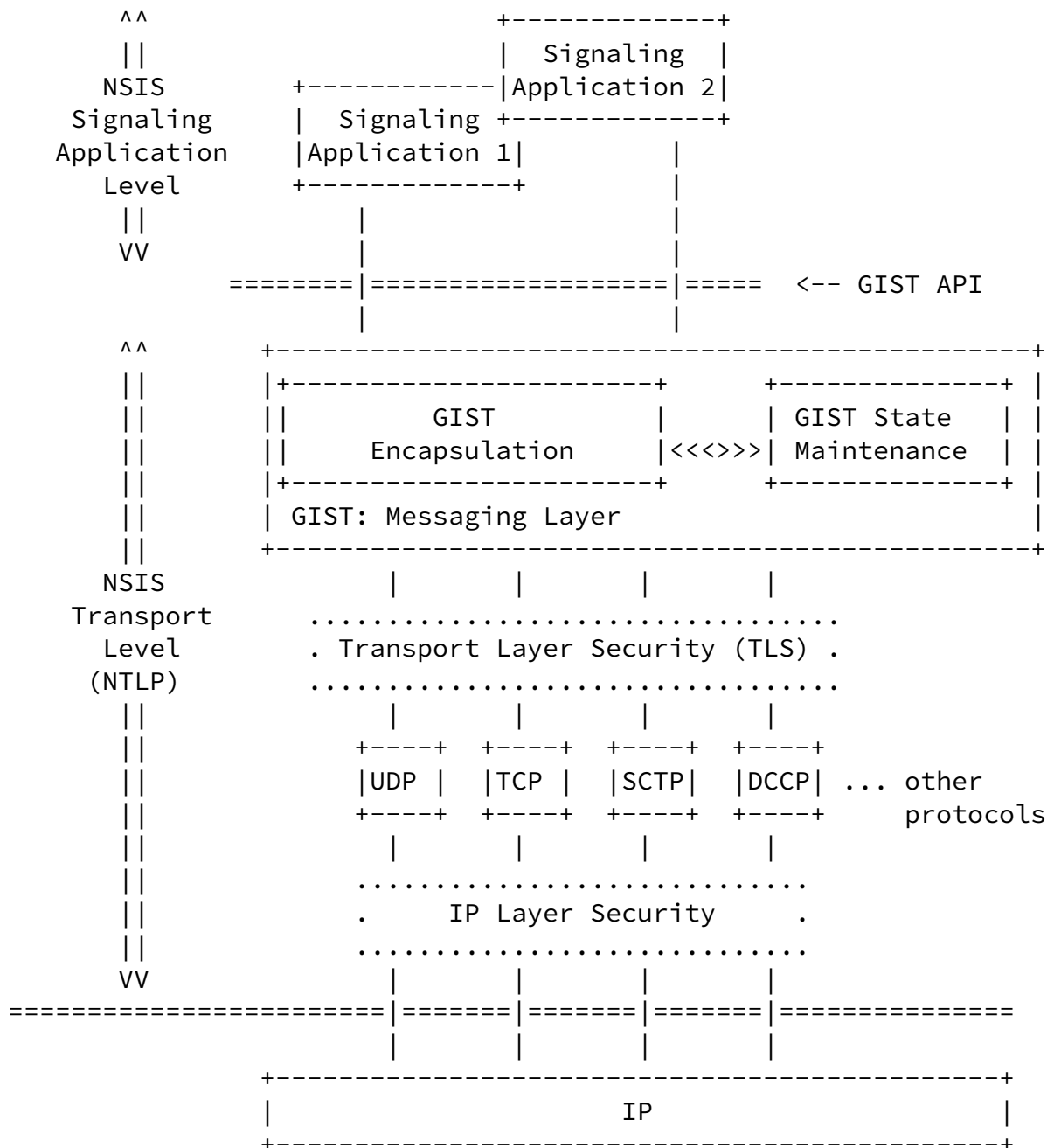


Figure 2: Protocol Stacks for Signaling Transport

3.2. Modes and Messaging Associations

Internally, GIST has two modes of operation:

Datagram mode (D-mode): used for small, infrequent messages with modest delay constraints and no security requirements; it must also be used when no routing state exists.

Connection mode (C-mode): used for larger data objects or where fast state setup in the face of packet loss is desirable, or where channel security is required.

D-mode uses UDP, as this is the only encapsulation which does not require per-message shared state to be maintained between the peers. C-mode can in principle use any stream or message-oriented transport protocol; this specification defines TCP as the initial choice. It can in principle employ specific network layer security associations, or an internal transport layer security association; this specification defines TLS as the initial choice. When GIST messages are carried in C-mode, they are treated just like any other traffic by intermediate routers between the GIST peers. Indeed, it would be impossible for intermediate routers to carry out any processing on the messages without terminating the transport and security protocols used.

It is possible to mix these two modes along a path. This allows, for example, the use of D-mode at the edges of the network and C-mode in the core of the network. Such combinations may make operation more efficient for mobile endpoints, while allowing multiplexing of signaling messages across shared security associations and transport connections between core routers.

It must be understood that the routing and transport decisions made by GIST are not independent. If the message transfer has requirements that require C-mode, for example if the message is so large that fragmentation is required, this can only be used between explicitly identified nodes. In such cases, GIST carries out the three-way handshake initially in D-mode to identify the peer and then sets up the necessary connections if they do not already exist. It must also be understood that the signaling application does not make

the D-mode/C-mode selection directly; rather, this decision is made by GIST on the basis of the message characteristics and the transfer attributes stated by the application. The distinction is not visible at the GIST service interface.

In general, the state associated with C-mode messaging to a particular peer (signaling destination address, protocol and port numbers, internal protocol configuration and state information) is referred to as a messaging association (MA). There may be any number of MAs between two GIST peers although the usual case is zero or one. They are set up and torn down by management actions within GIST itself.

[3.3.](#) Message Routing Methods

The baseline message routing functionality in GIST is that signaling

messages follow a route defined by an existing flow in the network, visiting a subset of the nodes through which it passes. This is the appropriate behaviour for application scenarios where the purpose of the signaling is to manipulate resources for that flow. However, there are scenarios for which other behaviours are applicable. Two examples are:

Predictive Routing: Here, the intent is to signal along a path that the data flow may follow in the future. Possible cases are pre-installation of state on the backup path that would be used in the event of a link failure, and predictive installation of state on the path that will be used after a mobile node handover.

NAT Address Reservations: This applies to the case where a node behind a NAT wishes to reserve an address at which it can be reached by a sender on the other side. This requires a message to be sent outbound from what will be the flow receiver although no reverse routing state for the flow yet exists.

Most of the details of GIST operation are independent of which alternative is being used. Therefore, the GIST design encapsulates the routing-dependent details as a message routing method (MRM), and allows multiple MRMs to be defined. The default is the path-coupled MRM, which corresponds to the baseline functionality described above; a second MRM for the NAT Address Reservation case is also defined.

The content of a MRM definition is as follows, using the path-coupled MRM as an example:

- o The format of the information that describes the path that the signaling should take, the Message Routing Information (MRI). For the path-coupled MRM, this is just the Flow Identifier (see [Section 5.8.1.1](#)) and some additional control information. Specifically, the MRI always includes a flag to distinguish between the two directions that signaling messages can take, denoted 'upstream' and 'downstream'.
- o A specification of the IP-level encapsulation of the messages which probe the network to discover the adjacent peers. A downstream encapsulation must be defined; an upstream encapsulation is optional. For the path-coupled MRM, this information is given in [Section 5.8.1.2](#) and [Section 5.8.1.3](#).
- o A specification of what validation checks GIST should apply to the probe messages, for example to protect against IP address spoofing attacks. The checks may be dependent on the direction (upstream or downstream) of the message. For the path-coupled MRM, the downstream validity check is basically a form of ingress

filtering, also discussed in [Section 5.8.1.2](#).

- o The mechanism(s) available for route change detection, i.e. any change in the neighbour relationships that the MRM discovers. The default case for any MRM is soft-state refresh, but additional supporting techniques may be possible; see [Section 7.1.2](#).

In addition, it should be noted that NAT traversal may require translation of fields in the MRI object carried in GIST messages (see [Section 7.2](#)). The generic MRI format includes a flag that must be given as part of the MRM definition, to indicate if some kind of translation is necessary. Development of a new MRM therefore includes updates to the GIST specification, and may include updates to specifications of NAT behaviour. These updates may be done in separate documents as is the case for the base GIST specification, as described in [Section 7.2.2](#).

The MRI is passed explicitly between signaling applications and GIST;

therefore, signaling application specifications must define which MRMs they require. Signaling applications may use fields in the MRI in their packet classifiers; if they use additional information for packet classification, this would be carried at the NSLP level and so would be invisible to GIST. Any node hosting a particular signaling application MUST use a GIST implementation that supports the corresponding MRMs. The GIST processing rules enforce that nodes which do not host the signaling application are not forced to handle messages for it at the GIST level, so it does not matter if they support the MRM or not.

[3.4.](#) GIST Messages

GIST has six message types: Query, Response, Confirm, Data, Error, and MA-Hello. Apart from the invocation of the messaging association protocols, all GIST communication consists of these messages. In addition, all signaling application data is carried as additional payloads in these messages, alongside the GIST information.

The first three messages implement the handshake that GIST uses to set up routing state and messaging associations. The handshake is initiated from the Querying node towards the Responding node. The first message is the Query, which is encapsulated in a special way depending on the message routing method, in order to probe the network infrastructure so that the correct peer will intercept it and become the Responding node. A Query always triggers a Response in the reverse direction as the second message of the handshake. As part of the defence against denial of service attacks, the Responding node can delay state installation until a return routability check, and require the Querying node to complete the handshake with the

Confirm. All of these three messages can optionally carry signaling application data.

The Data message is used purely to encapsulate signaling application data. Usually it is sent using pre-established routing state. However, if there are no security or transport requirements and no need for persistent reverse routing state, it can also be sent in the same way as the Query. Finally, Error messages are used to indicate error conditions at the GIST level, and the MA-Hello message can be used as a keepalive for the messaging association protocols.

[3.5.](#) Signaling Sessions

GIST allows signaling applications to associate each of their messages with a signaling session. Informally, given an application layer exchange of information for which some network control state information is to be manipulated or monitored, the corresponding signaling messages should be associated with the same session. Signaling applications provide the session identifier (SID) whenever they wish to send a message, and GIST reports the SID when a message is received.

Most GIST processing and state information is related to the flow (defined by the MRI, see above) and signaling application (given by the NSLP identifier, see below). There are several possible relationships between flows and sessions, for example:

- o The simplest case is that all messages for the same flow have the same SID.
- o Messages for more than one flow may use the same SID, for example because one flow is replacing another in a mobility or multihoming scenario.
- o A single flow may have messages for different SIDs, for example from independently operating signaling applications.

Because of this range of options, GIST does not perform any validation on how signaling applications map between flows and sessions, nor does it perform any validation on the properties of the SID itself. In particular, when a new SID is needed, logically it should be generated by the signaling application. NSIS implementations could provide common functionality to generate SIDs for use by any signaling application, but this is not part of GIST. GIST only defines the syntax of the SID as an opaque 128-bit identifier.

The SID assignment has the following impact on GIST processing:

- o Messages with the same SID that are to be delivered reliably between the same GIST peers are delivered in order.
- o All other messages are handled independently.

- o GIST identifies routing state (upstream and downstream peer) by the triplet (MRI, NSLP, SID).

Strictly, the routing state should not depend on the SID. However, if the routing state is keyed only by (MRI, NSLP), there is a trivial denial of service attack (see [Section 8.3](#)) where a malicious off-path node asserts that it is the peer for a particular flow. Instead, the routing state is also segregated between different SIDs, which means that the attacking node can only disrupt a signaling session if it can guess the corresponding SID. A consequence of this design is that signaling applications SHOULD choose SIDs so that they are cryptographically random, and SHOULD NOT use several SIDs for the same flow, to avoid additional load from routing state maintenance. Guidance on secure randomness generation can be found in [\[28\]](#).

[3.6.](#) Signaling Applications and NSLPIDs

The functionality for signaling applications is supported by NSIS signaling layer protocols (NSLPs). Each NSLP is identified by a 16 bit NSLP identifier (NSLPID), assigned by IANA ([Section 9](#)). A single signaling application, such as resource reservation, may define a family of NSLPs to implement its functionality, for example to carry out signaling operations at different levels in a hierarchy (cf. [\[19\]](#)). However, the interactions between the different NSLPs (for example, to relate aggregation levels or aggregation region boundaries in the resource management case) are handled at the signaling application level; the NSLPID is the only information visible to GIST about the signaling application being used.

[3.7.](#) Example of Operation

This section presents an example of GIST usage in a relatively simple (in particular, NAT-free) signaling scenario, to illustrate its main features.

Consider the case of an RSVP-like signaling application which allocates resources for a single unicast flow. We will consider how GIST transfers messages between two adjacent peers along the path, GN1 and GN2 (see Figure 1 in [Section 2](#)). In this example, the end-to-end exchange is initiated by the signaling application instance in the sender; we take up the story at the point where the first message is being processed above the GIST layer by the signaling application in GN1.

1. The signaling application in GN1 determines that this message is a simple description of resources that would be appropriate for the flow. It determines that it has no special security or transport requirements for the message, but simply that it should be transferred to the next downstream signaling application peer on the path that the flow will take.
2. The message payload is passed to the GIST layer in GN1, along with a definition of the flow and description of the message transfer attributes {unsecured, unreliable}. GIST determines that this particular message does not require fragmentation and that it has no knowledge of the next peer for this flow and signaling application; however, it also determines that this application is likely to require secured upstream and downstream transport of large messages in the future. This determination is a function of node-local policy interactions between GIST and the signaling application.
3. GN1 therefore constructs a GIST Query carrying the NSLP payload, and additional payloads at the GIST level to be used to initiate a messaging association. The Query is encapsulated in a UDP datagram and injected into the network, addressed towards the flow destination and with an IP Router Alert Option (RAO) included.
4. The Query passes through the network towards the flow receiver, and is seen by each router in turn. GIST-unaware routers will not recognise the RAO value and will forward the message unchanged; GIST-aware routers which do not support the NSLP in question will also forward the message basically unchanged, although they may need to process more of the message to decide this.
5. The message is intercepted at GN2. The GIST layer identifies the message as relevant to a local signaling application, and passes the NSLP payload and flow description upwards to it. The signaling application in GN2 indicates to GIST that it will peer with GN1 and so GIST should proceed to set up any routing state. In addition, the signaling application continues to process the message as in GN1 (compare step 1), forwarding the message downstream, and this will eventually result in the message reaching the flow receiver.
6. In parallel, the GIST instance in GN2 now knows that it should maintain routing state and a messaging association for future signaling with GN1. This is recognised because the message is a Query, and because the local signaling application has indicated

that it will peer with GN1. There are two possible cases for

Internet-Draft

GIST

July 2006

sending back the necessary GIST Response:

Association Exists: GN1 and GN2 already have an appropriate MA. GN2 simply records the identity of GN1 as its upstream peer for that flow and NSLP, and sends a Response back to GN1 over the MA identifying itself as the peer for this flow.

No Association: No MA exists. GN2 sends the Response in D-mode directly to GN1, identifying itself and agreeing to the association setup. The protocol exchanges needed to complete this will proceed in parallel with the following stages.

7. Eventually, another NSLP message works its way upstream from the receiver to GN2. This message contains a description of the actual resources requested, along with authorisation and other security information. The signaling application in GN2 passes this payload to the GIST level, along with the flow definition and transfer attributes {secured, reliable}.
8. The GIST layer in GN2 identifies the upstream peer for this flow and NSLP as GN1, and determines that it has an MA with the appropriate properties. The message is queued on the MA for transmission; this may incur some delay if the procedures begun in step 6.B have not yet completed.

Further messages can be passed in each direction in the same way. The GIST layer in each node can in parallel carry out maintenance operations such as route change detection (see [Section 7.1](#)).

It should be understood that several of these details of GIST operations can be varied, either by local policy or according to signaling application requirements. The authoritative details are contained in the remainder of this document.

Internet-Draft

GIST

July 2006

[4.](#) GIST Processing Overview

This section defines the basic structure and operation of GIST. [Section 4.1](#) describes the way in which GIST interacts with local signaling applications in the form of an abstract service interface. [Section 4.2](#) describes the per-flow and per-peer state that GIST maintains for the purpose of transferring messages. [Section 4.3](#) describes how messages are processed in the case where any necessary messaging associations and routing state already exist; this includes the simple scenario of pure D-mode operation, where no messaging associations are necessary in the first place. Finally, [Section 4.4](#) describes how routing state and messaging associations are created and managed.

[4.1.](#) GIST Service Interface

This section defines the service interface that GIST presents to signaling applications in terms of abstract properties of the message transfer. Note that the same service interface is presented at every GIST node; however, applications may invoke it differently at different nodes, depending for example on local policy. In addition, the service interface is defined independently of any specific transport protocol, or even the distinction between D-mode and C-mode. The initial version of this specification defines how to support the service interface using a C-mode based on TCP; if additional protocol support is added, this will support the same interface and so the change will be invisible to applications, except as a possible performance improvement. A more detailed description of this service interface is given in [Appendix B](#).

[4.1.1.](#) Message Handling

Fundamentally, GIST provides a simple message-by-message transfer

service for use by signaling applications: individual messages are sent, and individual messages are received. At the service interface, the NSLP payload, which is opaque to GIST, is accompanied by control information expressing the application's requirements about how the message should be routed, and the application also provides the session identifier (see [Section 3.5](#)). Additional message transfer attributes control the specific transport and security properties that the signaling application desires.

The distinction between GIST D- and C-mode is not visible at the service interface. In addition, the functionality to handle fragmentation and reassembly, bundling together of small messages for efficiency, and congestion control are not directly visible at the service interface; GIST will take whatever action is necessary based on the properties of the messages and local node state.

[4.1.2](#). Message Transfer Attributes

Message transfer attributes are used to define certain performance and security related aspects of message processing. The attributes available are as follows:

Reliability: This attribute may be 'true' or 'false'. When 'true', messages MUST be delivered to the signaling application in the peer exactly once or not at all; if there is a chance that the message was not delivered, an error MUST be indicated to the local signaling application identifying the routing information for the message in question. GIST implements reliability by using an appropriate transport protocol within a messaging association, so mechanisms for the detection of message loss depend on the protocol in question; for the current specification, the case of TCP is considered in [Section 5.7.2](#). Messages with the same SID to the same peer MUST be delivered in order. When 'false', a message may be delivered, once, several times or not at all, with no error indications in any case.

Security: This attribute defines the security properties that the signaling application requires for the message, including the type of protection required, and what authenticated identities should be used for the signaling source and destination. This information maps onto the corresponding properties of the security associations established between the peers in C-mode. It can be

specified explicitly by the signaling application, or reported by GIST to the signaling application. The latter can take place either on receiving a message, or just before sending a message but after configuring or selecting the messaging association to be used for it. This attribute can also be used to convey information about any address validation carried out by GIST, such as whether a return routability check has been carried out. Further details are discussed in [Appendix B](#).

Local Processing: An NSLP may provide hints to GIST to enable more efficient or appropriate processing. For example, the NSLP may select a priority from a range of locally defined values to influence the sequence in which messages leave a node. Any priority mechanism **MUST** respect the ordering requirements for reliable messages within a session, and priority values are not carried in the protocol or available at the signaling peer or intermediate nodes. An NSLP may also indicate that reverse path routing state will not be needed for this flow, to inhibit the node requesting its downstream peer to create it.

[4.2.](#) GIST State

[4.2.1.](#) Message Routing State

For each flow, the GIST layer can maintain message routing state to manage the processing of outgoing messages. This state is conceptually organised into a table with the following structure. The primary key (index) for the table is the combination of the information about how the message is to be routed, the session being signaled for, and the signaling application itself:

Message Routing Information (MRI): This defines the method to be used to route the message, the direction in which to send the message, and any associated addressing information; see [Section 3.3](#).

Session Identification (SID): The signaling session with which this message should be associated; see [Section 3.5](#).

NSLP Identification (NSLPID): This is an IANA-assigned identifier

associated with the NSLP which is generating messages for this flow. The inclusion of this identifier allows the routing state to be different for different NSLPs, for example because of different adjacencies.

The information associated with a given key consists of the routing state to reach the peer in the direction given by the MRI. For any flow, there will usually be two entries, one each for the upstream and downstream MRI. The routing state includes information about the peer identity (see [Section 4.4.2](#)), and a UDP port number for D-mode, or a reference to one or more MAs for C-mode. All of this information is learned from prior GIST exchanges.

It is also possible for the state information for either direction to be empty. There are several possible cases:

- o The signaling application has indicated that no messages will actually be sent in that direction.
- o The node is a flow endpoint, so there can be no signaling peer in one or other direction.
- o The node is the endpoint of the signaling path, for example because it is acting as a proxy, or because it has determined that there are no further signaling nodes in that direction.
- o The node is using other techniques to route the message. For example, it can encapsulate it the same way as a Query and rely on the peer to intercept it.

Each entry in the routing state table has an associated validity timer for how long it can be considered accurate; when this timer expires, the entry **MUST** be purged if it has not been refreshed. Installation and maintenance of routing state is described in more detail in [Section 4.4](#).

Note also that the routing state is described as a table of per-flow entries, but that there is no implied constraint on how the information is stored. However, in general, and especially if GIST peers are several IP hops away, there is no way to identify the correct downstream peer for a flow and signaling application from the local forwarding table using prefix matching, and the same applies

always to upstream peer state because of the possibility of asymmetric routing: prefix aggregation is not possible, and per-flow state has to be stored, just as for RSVP [12].

[4.2.2.](#) Peer-Peer Messaging Association State

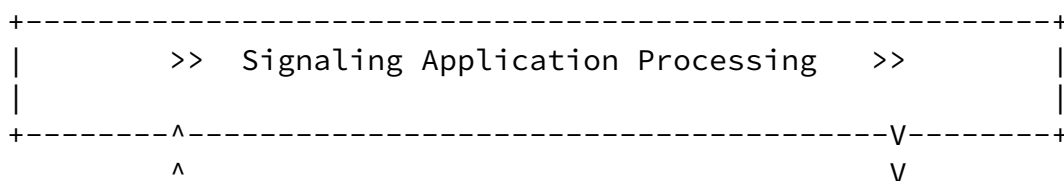
The per-flow message routing state is not the only state stored by GIST. There is also the state required to manage the MAs. Since these are not per-flow, they are stored separately from the routing state, including the following per-MA information:

- o a queue of messages pending transmission while an MA is being established;
- o a timer for how long since the peer re-stated its desire to keep the MA open (see [Section 4.4.3](#)).

In addition, per-MA state is held in the messaging association protocols themselves. However, the details of this state are not directly visible to GIST, and they do not affect the rest of the protocol description.

[4.3.](#) Basic GIST Message Processing

This section describes how signaling application messages are processed in the case where any necessary messaging associations and routing state are already in place. The description is divided into several parts. Firstly, message reception, local processing and message transmission are described for the case where the node hosts the NSLPID identified in the message. Secondly, the case where the message is handled directly in the IP or GIST layer (because there is no matching signaling application on the node) is given. An overview is given in Figure 3.



Reception in D-mode depends on the message type.

Normal encapsulation: Normal messages arrive UDP-encapsulated and addressed directly to the receiving signaling node, at an address and port learned previously. Each datagram contains a single message which is passed to the GIST layer for further processing, just as in the C-mode case.

Q-mode encapsulation: Where GIST is sending messages to be intercepted by the appropriate peer rather than directly addressed to it (in particular, Query messages), these are UDP encapsulated with an IP router alert option. Each signaling node will therefore see all such messages. Several RAO values may be used by NSIS; the assignment rationale is discussed in [11]. The case where the NSLPID does not match a local signaling application at all is considered below in [Section 4.3.4](#); otherwise, the message is again passed up to the GIST layer for further processing.

[4.3.2](#). Local Processing and Validation

Once a message has been received, it is processed locally within the GIST layer. The GIST hop count, which every message contains to prevent looping, MUST be checked and decremented immediately the message has been received. Further processing depends on the message type and payloads carried; most of the GIST payloads are associated with state maintenance and details are covered in [Section 4.4](#).

In the case of a Query, there is an interaction with signaling application policy to determine which of two courses to follow:

1. The receiving signaling application wishes to become a signaling peer with the Querying node. GIST MUST continue with the handshake process to set up message routing state, as described in [Section 4.4.1](#). The application MAY provide an NSLP payload for the same NSLPID, which GIST will transfer in the Response.
2. The signaling application does not wish to set up state with the Querying node and become its peer. GIST MUST propagate the Query, similar to the case described in [Section 4.3.4](#). No message is sent back to the Querying node. The application MAY provide an updated NSLP payload for the same NSLPID, which will be used in the Query forwarded by GIST.

This interaction with the signaling application, including the generation or update of an NSLP payload, SHOULD take place synchronously as part of the Query processing. In terms of the GIST service interface, this can be implemented by providing appropriate return values for the primitive that is triggered when such a message

Internet-Draft

GIST

July 2006

is received; see [Appendix B.2](#) for further discussion.

For all other message types, the GIST payloads are processed as described in [Section 4.4](#). The remainder of the GIST message consists of an NSLP payload, which is delivered locally to the signaling application identified by the NSLPID. The format of the payload is not constrained by GIST, and the content is not interpreted. Delivery is subject to the following validation checks:

- o if the message was explicitly routed (see [Section 7.1.4](#)) or is a Data message delivered without routing state (see [Section 5.3.2](#)), the payload is delivered but flagged to the receiving NSLP to indicate that routing state was not validated;
- o else, if there is no routing state for the MRI/SID/NSLPID the message MUST be rejected with a "No Routing State" error message ([Appendix A.4.4.5](#));
- o else, if the message arrived on an association which is not associated with the MRI/NSLPID/SID combination given in the message, the message MUST be rejected with an "Incorrectly Delivered Message" error message ([Appendix A.4.4.4](#));
- o else, the payload is delivered as normal.

[4.3.3](#). Message Transmission

Signaling applications can generate their messages for transmission, either asynchronously, or in response to a normal input message, and GIST can also generate messages autonomously. When a message is available for transmission, GIST uses internal policy and the stored routing state to determine how to handle it. The following processing applies equally to locally generated messages and messages forwarded from within the GIST or signaling application levels. However, see [Section 5.6](#) for special rules applying to the transmission of error messages by GIST.

The main decision is whether the message must be sent in C-mode or D-mode. Reasons for using the former are:

- o signaling application requirements: for example, it has requested channel secured delivery, or reliable delivery;

- o protocol specification: a message that requires fragmentation MUST be sent over a messaging association;
- o local policy: for example, a node MAY send messages over a messaging association to benefit from adaptive congestion control.

In principle, as well as determining that some messaging association must be used, GIST MAY select between a set of alternatives, e.g. for load sharing or because different messaging associations provide different transport or security attributes.

If the use of a messaging association is selected, the message is queued on the association found from the routing state table, and further output processing is carried out according to the details of the protocol stacks used. If no appropriate association exists, the message is queued while one is created (see [Section 4.4.1](#)). If no association can be created, this is an error condition, and should be indicated back to the local signaling application.

If a messaging association is not required, the message is sent in D-mode. The processing in this case depends on the message type and whether routing state exists or not.

- o If the message is not a Query, and routing state exists, it is UDP encapsulated and sent directly to the address from the routing state table.
- o If the message is a Query, then it is UDP encapsulated with IP address and router alert option determined from the MRI and NSLPID; further details depend on the message routing method.
- o If no routing state exists, GIST can attempt to use the same Q-mode encapsulation as in the Query case. If this is not possible, e.g. because the encapsulation for the MRM is only defined for one message direction, then this is an error condition which is reported back to the local signaling application.

[4.3.4](#). Nodes not Hosting the NSLP

A node may receive messages where it has no signaling application corresponding to the message NSLPID. There are several possible

cases depending mainly on the encapsulation:

1. A Q-mode encapsulated message contains an RAO value which is relevant to NSIS but not to the specific node, but the IP layer is unable to recognise whether it needs to be passed to GIST for further processing or whether the packet should be forwarded just like a normal IP datagram.
2. A Q-mode encapsulated message contains an RAO value which is relevant to the node, but the specific signaling application for the NSLPID in the message is not processed there.

3. A directly addressed message (in D-mode or C-mode) is delivered to a node for which there is no corresponding signaling application. With the current specification, this should never happen. While future versions might find a use for such a feature, currently this MUST cause an "Unknown NSLPID" error message, [Appendix A.4.4.6](#).
4. A Q-mode encapsulated message arrives at the end-system which does not handle the signaling application. This is possible in normal operation, and MUST be indicated to the sender with an "Endpoint Found" informational message (Appendix A.4.4.7). The end-system includes the MRI and SID from the original message in the error message without interpreting them.
5. The node is GIST-aware NAT. See [Section 7.2](#).

In cases (1) and (2), the role of GIST is to forward the message essentially unchanged, and it will not become a peer to the node sending the message. Forwarding with modified NSLP payloads is covered above in [Section 4.3.2](#). However, a GIST implementation must ensure that the IP-layer TTL field and GIST hop count are managed correctly to prevent message looping, and this should be done consistently independently of whether the processing takes place on the fast path or in GIST-specific code. The rules are that in cases (1) and (2), the IP-layer TTL MUST be decremented just as if the message was a normal IP forwarded packet; in case (2) the GIST hop count MUST be decremented as in the case of normal input processing, which indeed applies to cases (3) and (4).

A GIST node processing Q-mode encapsulated messages in this way SHOULD make the routing decision based on the full contents of the MRI and not only the IP destination address. It MAY also apply a restricted set of sanity checks and under certain conditions return an error message rather than forward the message. These conditions are:

1. The message is so large that it would be fragmented on downstream links, for example because the downstream MTU is very small. The error "Message Too Large" (Appendix A.4.4.8) SHOULD be returned to the sender, which SHOULD begin messaging association setup.
2. The GIST hop count has reached zero. The error "Hop Limit Exceeded" (Appendix A.4.4.2) SHOULD be returned to the sender, which MAY retry with a larger initial hop count if it is clear that a loop has not been formed.
3. The MRI represents a flow definition which is too general to be forwarded along a unique path (e.g. the destination address

prefix is too short). The error "MRI Validation Failure" (Appendix A.4.4.12) with subcode 0 ("MRI Too Wild") SHOULD be returned to the sender, which MAY retry with restricted MRIs, possibly starting additional signaling sessions to do so. If the GIST node does not understand the MRM in question it MUST NOT apply this check, instead forwarding the message transparently.

In the first two cases, only the common header of the GIST message is examined; in the third case, the MRI is also examined. The rest of the message MUST never be inspected or modified.

Note that the GIST hop count is only intended to prevent message looping at the GIST level, and by default NSLPs must take their own measures to prevent looping at the application level. However, the GIST API (Appendix B) provides the incoming hop count to the NSLPs, which can preserve it on outgoing messages as they are forwarded further along the path. This provides a lightweight loop-prevention mechanism for NSLPs which do not define anything more sophisticated.

[4.4.](#) Routing State and Messaging Association Maintenance

The main responsibility of GIST is to manage the routing state and messaging associations which are used in the message processing described above. Routing state is installed and maintained by specific GIST messages. Messaging associations depend on the existence of routing state, but are actually set up by the normal procedures of the transport and security protocols that comprise them. Timers control routing state and messaging association refresh and expiration.

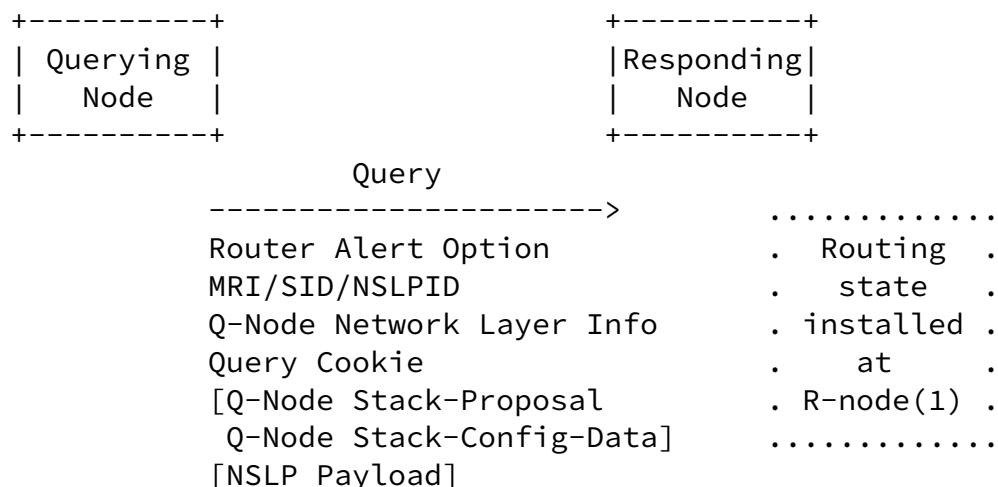
There are two different cases for state installation and refresh:

1. Where routing state is being discovered or a new association is to be established; and
2. Where an existing association can be re-used, including the case where routing state for the flow is being refreshed.

These cases are now considered in turn, followed by the case of background general management procedures.

[4.4.1.](#) State Setup

The complete sequence of possible messages for GIST state setup between adjacent peers is shown in Figure 4 and described in detail in the following text. The figure informally summarises the contents of each message, including optional elements in []. An example is given in [Appendix C](#).



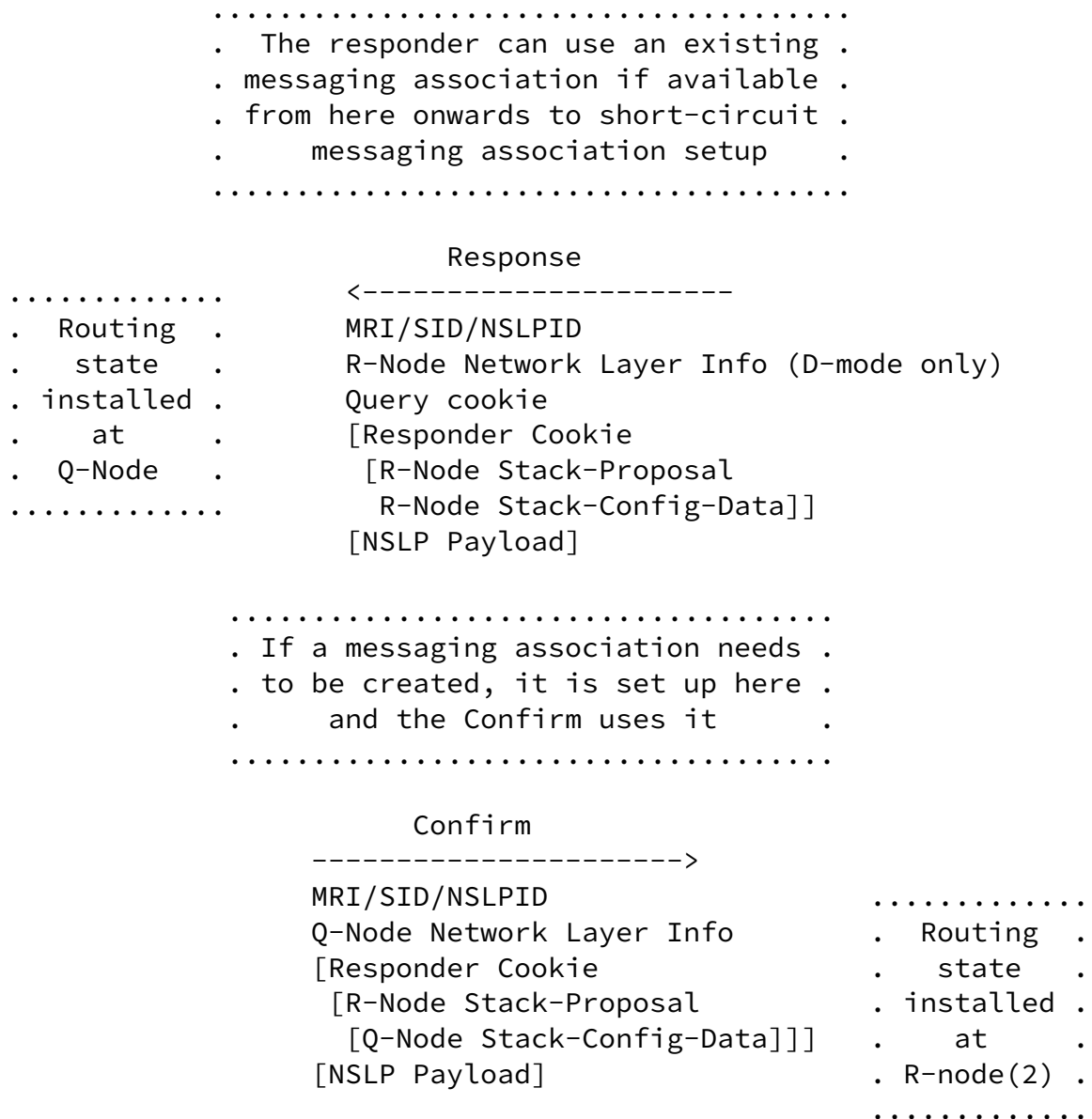


Figure 4: Message Sequence at State Setup

The initial message in any routing state maintenance operation is a Query, sent from the querying node and intercepted at the responding node. This message has addressing and other identifiers appropriate for the flow and signaling application that state maintenance is being done for, addressing information about the node itself, and it MAY contain an NSLP payload. It also includes a Query Cookie, and optionally capability information about messaging association

protocol stacks. The role of the cookies in this and subsequent messages is to protect against certain denial of service attacks and to correlate the various events in the message sequence (see [Section 8.5](#) for further details).

Provided that the signaling application has indicated that message routing state should be set up (see [Section 4.3.2](#)), reception of a Query MUST elicit a Response. This is a normally encapsulated D-mode message with additional payloads. It contains network layer information about the responding node, echoes the Query Cookie, and MAY contain an NSLP payload, possibly a response to the NSLP payload in the initial message. In case a messaging association was requested, it MUST also contain a Responder Cookie and its own capability information about messaging association protocol stacks. Even if a messaging association is not requested, the Response MAY still include a Responder Cookie if the node's routing state setup policy requires it (see below).

Setup of a new messaging association begins when peer addressing information is available and a new messaging association is actually needed. Any setup MUST take place immediately after the specific Query/Response exchange, because the addressing information used may have a limited lifetime, either because it depends on limited lifetime NAT bindings or because it refers to agile destination ports for the transport protocols. The Stack-Proposal and Stack-Configuration-Data objects carried in the exchange carry capability information about what messaging association protocols can be used, and the processing of these objects is described in more detail in [Section 5.7](#). With the protocol options currently defined, setup of the messaging association always starts from the Querying node, although more flexible configurations are possible within the overall GIST design. In any case, once set up, the association itself can be used equally in both directions.

Finally, after any necessary messaging association setup has completed, a Confirm MUST be sent if the Response requested it. If a messaging association is being used, the Confirm MUST be sent over it before any other messages for the same flow, and it echoes the Responder Cookie and Stack-Proposal from the Response. The former is used to allow the receiver to validate the contents of the message (see [Section 8.5](#)), and the latter is to prevent certain bidding-down

attacks on messaging association security (see [Section 8.6](#)). The Confirm MAY also contain an abbreviated form of the original Stack-Configuration-Data to finalise details of the messaging association configuration. The association can be used in the upstream direction for the MRI and NSLPID carried in the Confirm, after the Confirm has been received.

The querying node MUST install the responder address as routing state information after verifying the Query Cookie in the Response. The responding node MAY install the querying address as peer state information at two points in time:

1. after the receipt of the initial Query, or
2. after a Confirm containing the Responder Cookie.

The precise constraints on when state information is installed are a matter of security policy considerations on prevention of denial-of-service attacks and state poisoning attacks, which are discussed further in [Section 8](#). Because the responding node MAY choose to delay state installation as in case (2), the Confirm must contain sufficient information to allow it to be processed in the same way as the original Query. This places some special requirements on NAT traversal and cookie functionality, which are discussed in [Section 7.2](#) and [Section 8](#) respectively.

[4.4.2](#). Messaging Association Re-use

It is a design goal of GIST that, as far as possible, messaging associations should be re-used for multiple flows and sessions, rather than setting up a new MA for each. This is to ensure that the MA cost scales only with the number of peers, and to avoid the latency of new MA setup where possible.

However, re-use requires the identification of an existing MA which matches the same routing state and desired properties that would be the result of a full handshake in D-mode, and this identification must be done as reliably and securely as continuing with the full procedure. Note that this requirement is complicated by the fact that NATs may remap the node addresses in D-mode messages, and also interacts with the fact that some nodes may peer over multiple interfaces (and thus with different addresses).

MA re-use is controlled by the Network-Layer-Information (NLI) object, which is carried in Query/Confirm and optionally Response messages. The NLI object includes:

Internet-Draft

GIST

July 2006

Peer-Identity: For a given node, this is an interface independent value with opaque syntax. It MUST be chosen so as to have a high probability of uniqueness between peers, and SHOULD be stable at least until the next node restart. Note that there is no cryptographic protection of this identity; attempting to provide this would essentially duplicate the functionality in the messaging association security protocols.

Interface-Address: This is an IP address through which the signaling node can be reached. There may be several choices available for the Interface-Address, and further discussion of this is contained in [Section 5.2.2](#).

By default, a messaging association is associated with the NLI object that was provided by the peer in the Query/Response/Confirm at the time the association was set up. There may be more than one association for a given NLI object, for example with different security or transport properties.

MA re-use is controlled by matching the NLI provided in a GIST message with those associated with existing MAs. This can be done on receiving either a Query or Response, although the former is more likely:

- o If there is a perfect match to the NLI of an existing association, that association SHOULD be re-used, provided it has the appropriate properties in other respects. This is indicated by sending the remaining messages in the handshake over that association. This will only fail, that is, lead to re-use of an association to the wrong node, if signaling nodes have colliding Peer-Identities and one is reachable at the same Interface-Address as another. This could be done by an on-path attacker.
- o In all other cases, the full handshake MUST be executed in D-mode as usual. There are in fact four possibilities:
 1. Nothing matches: this is clearly a new peer.
 2. Only the Peer-Identity matches: this may be either a new interface on an existing peer, or a changed address mapping behind a NAT, or an attacker attempting to hijack the Peer-Identity. These should be rare events, so the expense of a new association setup is acceptable.

3. Only the Interface-Address matches: this is probably a new peer behind the same NAT as an existing one. A new association setup is required.

4. The full NLI object matches: this is a degenerate case, where one node recognises an existing peer, but wishes to allow the option to set up a new association in any case, for example to create an association with different properties.

[4.4.3.](#) State Maintenance Procedures

Refresh and expiration of all types of GIST state is controlled by timers.

Each item of routing state expires after a lifetime which is negotiated during the Query/Response/Confirm handshake. The Network Layer Info (NLI) object in the Query contains a proposal for the lifetime value, and the NLI in the Response contains the value the Responding node requires. A default timer value of 30 seconds is RECOMMENDED. Nodes which can exploit alternative, more powerful, route change detection methods such as those described in [Section 7.1.2](#) MAY choose to use much longer times. Nodes MAY use shorter times to provide more rapid change detection, but MUST take into account the fact that the Query messages generated may stress the rate limits applied to D-mode traffic ([Section 5.3.3](#)).

The Querying node MUST generate a Query before this timer expires, if it believes that the signaling session is still active; otherwise, the Responding node MAY delete the state. Receipt of the message at the Responding node will refresh peer addressing state for one direction, and receipt of a Response at the querying node will refresh it for the other. There is no mechanism at the GIST level for explicit teardown of routing state. However, GIST MUST NOT refresh routing state if a signaling session is known to be inactive, either because upstream state has expired, or because the signaling application has indicated via the GIST API (Appendix B.5) that the state is no longer required, because this would prevent correct state repair in the case of network rerouting.

Unneeded MAs are torn down by GIST, using the teardown mechanisms of

the underlying transport or security protocols if available, for example by simply closing a TCP connection. The teardown can be initiated by either end. Whether an MA is needed is a combination of two factors:

- o local policy, which could take into account the cost of keeping the messaging association open, the level of past activity on the association, and the likelihood of future activity, e.g. if there is routing state still in place which might generate messages to use it.

- o whether the peer still wants the MA in place. During MA setup, each node indicates its own MA-Hold-Time as part of the Stack-Configuration-Data. A node MUST NOT tear down the MA if it has received traffic from its peer over that period. A peer which has generated no traffic but still wants the MA retained may use a special null message (MA-Hello) to indicate the fact. A default value for MA-Hold-Time of 30 seconds is RECOMMENDED. Nodes MAY use shorter times to achieve more rapid peer failure detection, but need to take into account the load on the network created by the MA-Hello messages. Nodes MAY use longer times, but need to take into account the cost of retaining idle MAs for extended periods.

Because the Responding node can choose not to retain state until a Confirm, an abbreviated Stack-Configuration-Data object containing just this information MUST be repeated by the Querying node in the first Confirm sent on a new MA.

Messaging associations can always be set up on demand, and messaging association status is not made directly visible outside the GIST layer. Therefore, even if GIST tears down and later re-establishes a messaging association, signaling applications cannot distinguish this from the case where the MA is kept permanently open. To maintain the transport semantics described in [Section 4.1](#), GIST MUST close transport connections carrying reliable messages gracefully or report an error condition, and MUST NOT open a new association for a given session and peer while messages on a previous association may still be outstanding.

This specification defines precisely only the time at which routing state or messaging associations expire; it does not define when refresh handshakes or keepalives should be initiated. Implementations MUST select timer settings which take at least the following into account:

- o The transmission latency between source and destination;
- o The need for retransmissions, either explicitly or within the messaging association protocols;
- o The need to avoid network synchronisation of control traffic (cf. [\[39\]](#)).

In most cases, a reasonable policy is to initiate the refresh process when between 1/2 and 3/4 of the appropriate validity time has elapsed since the last successful refresh. The actual moment is chosen randomly within this interval to avoid synchronisation effects.

[5.](#) Message Formats and Transport

[5.1.](#) GIST Messages

All GIST messages begin with a common header, followed by a sequence of type-length-value (TLV) objects. This subsection describes the various GIST messages and their contents at a high level in ABNF [\[9\]](#); a more detailed description of the header and each object is given in [Section 5.2](#). Note that the NAT traversal mechanism for GIST involves the insertion of an additional NAT-Traversal-Object in Query, Response and Error messages; the rules for this are given in [Section 7.2](#).

GIST-Message: The primary messages are either one of the stages in the three-way handshake, or a simple message carrying NSLP data. Additional types are defined for errors and keeping messaging associations alive.

GIST-Message = Query / Response / Confirm /
Data / Error / MA-Hello

The common header includes a version number, message type and size,

and NSLPID. It also carries a hop count to prevent message looping and various control flags, including one (the R flag) to indicate if a reply of some sort is requested. The objects following the common header MUST be carried in a fixed order, depending on message type. Messages with missing, duplicate or invalid objects for the message type MUST be rejected with an "Object Type Error" error message with the appropriate subcode (Appendix A.4.4.9).

Query: A Query MUST be sent in D-mode, in fact with the special Q-mode encapsulation. In addition to the common header, it contains certain mandatory control objects, and MAY contain a signaling application payload. A stack proposal and configuration data MUST be included if the message exchange relates to setup of a messaging association. The R flag MUST always be set (R=1) in a Query, since this message always elicits a Response.

```
Query = Common-Header
       [ NAT-Traversal-Object ]
       Message-Routing-Information
       Session-Identification
       Network-Layer-Information
       Query-Cookie
       [ Stack-Proposal Stack-Configuration-Data ]
       [ NSLP-Data ]
```

Response: A Response may be sent in D-mode, or C-mode if a messaging

association is being re-used. It MUST echo the MRI SID and Query-Cookie of the Query, and in D-mode carries its own Network-Layer-Information; if the message exchange relates to setup of a messaging association, which can only take place in D-mode, a Responder cookie MUST be included, as well as its own stack proposal and configuration data. The R flag MUST be set (R=1) if a Responder cookie is present but otherwise is optional; if the R flag is set, a Confirm MUST be sent as a reply. Note that the direction of this MRI will be inverted compared to that in the Query, that is, an upstream MRI becomes downstream and vice versa (see [Section 3.3](#)).

```
Response = Common-Header
          [ NAT-Traversal-Object ]
          Message-Routing-Information
          Session-Identification
```

```

[ Network-Layer-Information ]
Query-Cookie
[ Responder-Cookie
  [ Stack-Proposal Stack-Configuration-Data ] ]
[ NSLP-Data ]

```

Confirm: A Confirm may be sent in D-mode, or C-mode if a messaging association has been re-used. It MUST echo the MRI (with inverted direction), SID, and Responder-Cookie if the Response carried one; if the message exchange relates to setup of a new messaging association or re-use of an existing one (which can only take place in C-mode), the message MUST also echo the Stack-Proposal from the Response so it can be verified that this has not been tampered with. The first message on an association MUST also repeat the Stack-Configuration-Data from the original Query in an abbreviated form, just containing the MA-Hold-Time.

```

Confirm = Common-Header
        Message-Routing-Information
        Session-Identification
        Network-Layer-Information
        [ Responder-Cookie
          [ Stack-Proposal
            [ Stack-Configuration-Data ] ] ]
        [ NSLP-Data ]

```

Data: The Data message is used to transport NSLP data without modifying GIST state. It contains no control objects, but only the MRI and SID associated with the NSLP data being transferred. Network-Layer-Information (NLI) MUST be carried in the D-mode case, but MUST NOT be included otherwise.

```

Data = Common-Header
      Message-Routing-Information
      Session-Identification
      [ Network-Layer-Information ]
      NSLP-Data

```

Error: An Error message reports a problem determined at the GIST level. (Errors generated by signaling applications are reported in

NSLP-Data payloads and are not treated specially by GIST.) The message includes a Network-Layer-Information object for the originator of the error message if it is being sent in D-mode; all other information related to the error is carried in a GIST-Error-Data object.

```
Error = Common-Header
      [ NAT-Traversal-Object ]
      [ Network-Layer-Information ]
      GIST-Error-Data
```

MA-Hello: This message MUST be sent only in C-mode to indicate that a node wishes to keep a messaging association open. It contains only the common header, with a NSLPID of zero. The R flag MAY be set (R=1); if so, the peer MUST send another message back along the messaging association. This allows a node to test the liveness of the peer.

```
MA-Hello = Common-Header
```

[5.2.](#) Information Elements

This section describes the content of the various objects that can be present in each GIST message, both the common header, and the individual TLVs. The bit formats are provided in [Appendix A](#).

[5.2.1.](#) The Common Header

Each message begins with a fixed format common header, which contains the following information:

Version: The version number of the GIST protocol. This specification defines GIST version 1.

Length: The number of 32 bit words in the message following the common header.

Upper layer identifier (NSLPID): This gives the specific NSLP that

this message is used for.

GIST hop count: A hop count to prevent a message from looping.

Message type: The message type (Query, Response, etc.)

Source addressing mode: If set (S=1), this indicates that the IP source address of the message is the same as the IP address of the signaling peer, so replies to this message can be sent safely to this address. S is always set in C-mode. It is cleared (S=0) if the IP source address was derived from the message routing information in the payload and this is different from the signaling source address.

Response requested: A flag which if set (R=1) indicates that a GIST message should be sent in response to this message. The appropriate message type for the response depends on the type of the initial message.

Explicit routing: A flag which if set (E=1) indicates that the message was explicitly routed (see [Section 7.1.4](#)).

[5.2.2](#). TLV Objects

All data following the common header is encoded as a sequence of type-length-value objects. Currently, each object can occur at most once; the set of required and permitted objects is determined by the message type and encapsulation (D-mode or C-mode).

Message-Routing-Information (MRI): Information sufficient to define how the signaling message should be routed through the network.

Message-Routing-Information = message-routing-method
method-specific-information

The format of the method-specific-information depends on the message-routing-method requested by the signaling application. Note that it always includes a flag defining the direction as either 'upstream' or 'downstream' (see [Section 3.3](#)). It is provided by the NSLP in the message sender and used by GIST to select the message routing.

Session-Identification (SID): The GIST session identifier is a 128 bit, cryptographically random identifier chosen by the node which originates the signaling exchange. See [Section 3.5](#).

Network-Layer-Information (NLI): This object carries information about the network layer attributes of the node sending the message, including data related to the management of routing state. This includes a peer identity and IP address for the sending node. It also includes IP-TTL information to allow the IP hop count between GIST peers to be measured and reported, and a validity time (RS-validity-time) for the routing state.

```
Network-Layer-Information = peer-identity
                           interface-address
                           RS-validity-time
                           IP-TTL
```

The use of the RS-validity-time field is described in [Section 4.4.3](#). The peer-identity and interface-address are used for matching existing associations, as discussed in [Section 4.4.2](#).

The interface-address must be routable, i.e. it MUST be usable as a destination IP address for packets to be sent back to the node generating the signaling message, whether in D-mode or C-mode. If this object is carried in a Query or Confirm, the interface-address MUST specifically be set to an address bound to the interface associated with the MRI, to allow its use in route change handling as discussed in [Section 7.1](#). A suitable choice is the interface that is carrying the outbound flow. A node may have several choices for which of its addresses to use as the interface-address. For example, there may be a choice of IP versions, or addresses of limited scope (e.g. link-local), or addresses bound to different interfaces in the case of a router or multi-homed host. However, some of these interface addresses may not be usable by the peer. A node MUST follow a policy of using a global address of the same IP version as in the MRI, unless it can establish that an alternative address would also be usable.

The setting and interpretation of the IP-TTL field depends on the message direction (upstream/downstream as determined from the MRI as described above) and encapsulation.

- * If the message is sent downstream, if the TTL that will be set in the IP header for the message can be determined, the IP-TTL value MUST be set to this value, or else set to 0.
- * On receiving a downstream message in D-mode, a non-zero IP-TTL is compared to the TTL in the IP header, and the difference is stored as the IP-hop-count-to-peer for the upstream peer in the routing state table for that flow. Otherwise, the field is

ignored.

Internet-Draft

GIST

July 2006

- * If the message is sent upstream, the IP-TTL MUST be set to the value of the IP-hop-count-to-peer stored in the routing state table, or 0 if there is no value yet stored.
- * On receiving an upstream message, the IP-TTL is stored as the IP-hop-count-to-peer for the downstream peer.

In all cases, the IP-TTL value reported to signaling applications is the one stored with the routing state for that flow, after it has been updated if necessary from processing the message in question.

Stack-Proposal: This field contains information about which combinations of transport and security protocols are available for use in messaging associations, and is also discussed further in [Section 5.7](#).

Stack-Proposal = 1*stack-profile

stack-profile = 1*protocol-layer

Each protocol-layer field identifies a protocol with a unique tag; any additional data, such as higher-layer addressing or other options data associated with the protocol, will be carried in a MA-protocol-options field in the Stack-Configuration-Data TLV (see below).

Stack-Configuration-Data (SCD): This object carries information about the overall configuration of a messaging association.

Stack-Configuration-Data = MA-Hold-Time
0*MA-protocol-options

The MA-Hold-Time field indicates how long a node will hold open an inactive association; see [Section 4.4.3](#) for more discussion. The MA-protocol-options fields give the configuration of the protocols (e.g. TCP, TLS) to be used for new messaging associations, and they are described in more detail in [Section 5.7](#).

Query-Cookie/Responder-Cookie: A Query-Cookie is contained in a Query and MUST be echoed in a Response; a Responder-Cookie MAY be sent in a Response, and if present MUST be echoed in the following Confirm. Cookies are variable length bit strings, chosen by the cookie generator. See [Section 8.5](#) for further details on requirements and mechanisms for cookie generation.

NSLP-Data: The NSLP payload to be delivered to the signaling application. GIST does not interpret the payload content.

GIST-Error-Data: This contains all the information to determine the cause and context of an error.

```
GIST-Error-Data = error-class error-code error-subcode
                  common-error-header
                  [ Message-Routing-Information-content ]
                  [ Session-Identification-content ]
                  0*additional-information
                  [ comment ]
```

The error-class indicates the severity level, and the error-code and error-subcode identify the specific error itself. A full list of GIST errors and their severity levels is given in [Appendix A.4](#). The common-error-header carries the Common-Header from the original message, and contents of the Message-Routing-Information (MRI) and Session-Identification (SID) objects are also included if they were successfully decoded. For some errors, additional information fields must be included according to a fixed format; finally, an optional free-text comment may be added.

[5.3](#). D-mode Transport

This section describes the various encapsulation options for D-mode messages. Although there are several possibilities, depending on message type, MRM, and local policy, the general design principle is that the sole purpose of the encapsulation is to ensure that the message is delivered to or intercepted at the correct peer. Beyond that, minimal significance is attached to the type of encapsulation or the values of addresses or ports used for it. This allows new

options to be developed in the future to handle particular deployment requirements without modifying the overall protocol specification.

[5.3.1.](#) Normal Encapsulation

Normal encapsulation MUST be used for all D-mode messages where the signaling peer is already known from previous signaling. This includes Response and Confirm messages, and Data messages except if these are being sent without using local routing state. Normal encapsulation is simple: the complete set of GIST payloads is concatenated together with the common header, and placed in the data field of a UDP datagram. UDP checksums MUST be enabled. The message is IP addressed directly to the adjacent peer; the UDP port numbering MUST be compatible with that used on Query messages (see below), that is, the same for messages in the same direction and with port numbers swapped for messages in the opposite direction.

[5.3.2.](#) Q-mode Encapsulation

Q-mode encapsulation MUST be used for messages where no routing state is available or where the routing state is being refreshed, in particular for Query messages. Q-mode encapsulation is similar to normal encapsulation, with changes in IP address selection, IP options, and a defined method for selecting UDP ports.

In general, the IP addresses are derived from information in the MRI; the exact rules depend on the MRM. In addition, the IP header is given a Router Alert Option ([\[1\]](#) and [\[5\]](#)) to assist the peer in intercepting the message depending on the NSLPID. Each NSLPID maps to a unique RAO value, but one RAO value may refer to multiple NSLPIDs; further details are discussed in [\[11\]](#).

The source UDP port is selected by the message sender as the port at which it is prepared to receive UDP messages in reply, and a destination UDP port is allocated for GIST by IANA (see [Section 9](#)). Note that for some MRMs, GIST nodes anywhere along the path can generate GIST packets with source addresses that spoof the source address of the data flow. Therefore, destinations cannot distinguish these packets from genuine end-to-end data purely on address analysis. Instead, it must be possible to distinguish such GIST packets by port analysis; furthermore, the mechanism to do so must remain valid even if the destination is GIST-unaware. GIST solves

this problem by using a fixed destination UDP port from the "well known" space for the Q-mode encapsulation. This port should never be allocated on a GIST-unaware host, and therefore Q-mode encapsulated messages will always be rejected with an ICMP error.

5.3.3. Retransmission and Rate Control

D-mode uses UDP, and hence has no automatic reliability or congestion control capabilities. Signaling applications requiring reliability should be serviced using C-mode, which should also carry the bulk of signaling traffic. However, some form of messaging reliability is required for the GIST control messages themselves, as is rate control to handle retransmissions and also bursts of unreliable signaling or state setup requests from the signaling applications.

Query messages which do not receive Responses MAY be retransmitted; retransmissions MUST use a binary exponential backoff. The initial timer value is T1, which the backoff process can increase up to a maximum value of T2 seconds. The default value for T1 is 500 ms. T1 is an estimate of the round-trip time between the querying and responding nodes. Elements MAY use smaller values of T1 if it is known that the Query should be answered within the local network. T1 MAY be chosen larger, and this is RECOMMENDED if it is known in

advance (such as on high latency access links) that the round-trip time is larger. The default value of T2 is $64 \times T1$.

Note that Queries may go unanswered either because of message loss (in either direction), or because there is no reachable GIST peer. Therefore, implementations MAY trade off reliability (large T2) against promptness of error feedback to applications (small T2). If the NSLP has indicated a timeout on the validity of this payload (see [Appendix B.1](#)), T2 MUST be chosen so that the process terminates within this timeout.

Retransmitted Queries MUST use different Query-Cookie values. If the Query carries NSLP data, it may be delivered multiple times to the signaling application. These rules apply equally to the message that first creates routing state, and those that refresh it.

This algorithm is sufficient to handle lost Queries and Responses. The case of a lost Confirm is more subtle. Notionally, we can

distinguish between two cases:

1. Where the Responding node is already prepared to store per-flow state after receiving a single (Query) message. This would include any cases where the node has NSLP data queued to send. Here, the Responding node MAY run a retransmission timer to resend the Response until a Confirm is received, since the node is already managing state for that flow. The problem of an amplification attack stimulated by a malicious Query is handled by requiring the cookie mechanism to enable the node receiving the Response to discard it efficiently if it does not match a previously sent Query.
2. Where the responding node is not prepared to store per-flow state until receiving a properly formed Confirm.

Case (2) should be handled without requiring a retransmission timer, since this would require per-flow state at the Responding node. However, we can assume that the next signaling message will be in the direction Querying Node -> Responding Node (if there is no next signaling message, the fact that the Confirm has been lost is moot). In this case, the responding node will start to receive messages at the GIST level for a MRI/NSLP combination for which there is no stored routing state, since this state is only created on receipt of a Confirm.

Therefore, the error condition is detected at the Responding node when such a message arrives, without the need for a specific timer. Recovery requires a Confirm to be transmitted and successfully received. The mechanism to cause this is that the Responding node

MUST reject the incoming message with a "No Routing State" error message (Appendix A.4.4.5) back to the Querying node, which MUST interpret this as caused by a lost Confirm; the Querying node MUST regenerate the Confirm purely from local state. In particular, it needs to remember a valid Responder Cookie.

In all cases, Responses MUST be sent promptly to avoid spurious retransmissions. Nodes generating any type of retransmission MUST be prepared to receive and match a reply to any of them, not just the one most recently sent.

The basic rate-control requirements for D-mode traffic are deliberately minimal. A single rate limiter applies to all traffic, for all interfaces and message types. It applies to retransmissions as well as new messages, although an implementation MAY choose to prioritise one over the other. Rate-control applies only to locally generated D-mode messages, not to messages which are being forwarded. When the rate limiter is in effect, D-mode messages MUST be queued until transmission is re-enabled, or an error condition MAY be indicated back to local signaling applications. The rate limiting mechanism is implementation-defined, but it is RECOMMENDED that a token bucket limiter as described in [31] be used. The token bucket MUST be sized to ensure that a node cannot saturate the network with D-mode traffic, for example when re-probing the network for multiple flows after a route change. A suitable approach is to restrict the token bucket parameters so that the mean output rate is a small fraction, such as 5%, of the node's lowest-speed interface.

[5.4.](#) C-mode Transport

Encapsulation in C-mode is more complex, because of the variation in available transport functionality. This issue is treated in [Section 5.4.1](#). The actual encapsulation is given in [Section 5.4.2](#).

[5.4.1.](#) Choice of Transport Protocol

It is a general requirement of the NTLP defined in [26] that it should be able to support bundling of small messages, fragmentation of large messages, and message boundary delineation. Not all transport protocols natively support all these features.

TCP provides both bundling and fragmentation, but not message boundaries. However, the length information in the GIST common header allows the message boundary to be discovered during parsing.

SCTP [16] satisfies all requirements.

DCCP [30] is message-based but does not provide bundling or fragmentation, nor does it provide reliability. Bundling can be

carried out by the GIST layer sending multiple messages in a single datagram; because the common header includes length information, the message boundaries within the datagram can be discovered during parsing. Fragmentation of GIST messages over multiple datagrams should be avoided, because of amplification of message loss rates that this would cause.

The bundling together of small messages is either built into the transport protocol or can be carried out by the GIST layer during message construction. Either way, two approaches can be distinguished:

1. As messages arrive for transmission they are gathered into a bundle until a size limit is reached or a timeout expires (cf. the Nagle algorithm of TCP or similar optional functionality in SCTP). This provides maximal efficiency at the cost of some latency.
2. Messages awaiting transmission are gathered together while the node is not allowed to send them, for example because it is congestion controlled.

The second type of bundling is always appropriate. For GIST, the first type MUST NOT be used for trigger messages (i.e. messages that update GIST or signaling application state), but may be appropriate for refresh messages (i.e. messages that just extend timers). These distinctions are known only to the signaling applications, but MAY be indicated (as an implementation issue) by setting the priority transfer attribute ([Section 4.1.2](#)).

It can be seen that all of these transport protocol options can be supported by the basic GIST message format already presented. GIST messages requiring fragmentation must be carried using a reliable transport protocol, TCP or SCTP. This specification defines only the use of TCP, but other possibilities could be included without additional work on message formatting.

[5.4.2](#). Encapsulation Format

The GIST message, consisting of common header and TLVs, is carried directly in the transport protocol, possibly incorporating transport layer security protection. Further messages can be carried in a continuous stream (for TCP), or up to the next transport layer message boundary (for SCTP, DCCP, or UDP). This is shown in

Internet-Draft

GIST

July 2006

Message	Normal D-mode	Query D-mode (Q-mode)	C-mode
Query	Never	Always	Never
Response	Unless a messaging association is being re-used	Never	If a messaging association is being re-used
Confirm	Unless a messaging association has been set up or is being re-used	Never	If a messaging association has been set up or is being re-used
Data	If routing state exists for the flow but no messaging association	If no routing state exists and the MRI can be used to derive the Q-mode encapsulation	If a messaging association exists

[5.6.](#) Error Message Processing

Special rules apply to the encapsulation and transmission of error messages.

GIST only generates error messages in response to incoming messages. Error messages **MUST NOT** be generated in response to incoming error messages. The routing and encapsulation of the error message is derived from that of the message that caused the error; in particular, local routing state is not consulted. Routing state and messaging association state **MUST NOT** be created to handle the error, and error messages **MUST NOT** be retransmitted explicitly by GIST, although they are subject to the same rate control as other messages.

- o If the incoming message was received in D-mode, the error MUST be sent in D-mode using the normal encapsulation, using the addressing information from the NLI object in the incoming message. If the NLI could not be determined, the error MUST be sent to the IP source of the incoming message if the S flag was set in it. The NLI object in the Error message reports information about the originator of the error.

- o If the incoming message was received over a messaging association, the error MUST be sent back over the same messaging association.

The NSLPID in the common header of the Error message has the value zero. If for any reason the message cannot be sent, for example, because it is too large to send in D-mode, an error SHOULD be logged locally.

[5.7.](#) Messaging Association Setup

[5.7.1.](#) Overview

A key attribute of GIST is that it is flexible in its ability to use existing transport and security protocols. Different transport protocols may have performance attributes appropriate to different environments; different security protocols may fit appropriately with different authentication infrastructures. Even given an initial default mandatory protocol set for GIST, the need to support new protocols in the future cannot be ruled out, and secure feature negotiation cannot be added to an existing protocol in a backwards-compatible way. Therefore, some sort of capability discovery is required.

Capability discovery is carried out in Query and Response messages, using Stack-Proposal and Stack-Configuration-Data objects. If a new messaging association is required it is then set up, followed by a Confirm. Messaging association re-use is achieved by short-circuiting this exchange by sending the Response or Confirm messages on an existing association ([Section 4.4.2](#)); whether to do this is a matter of local policy. The end result of this process is a messaging association which is a stack of protocols. If multiple associations exist, it is a matter of local policy how to distribute messages over them, subject to respecting the transfer attributes

requested for each message.

Every possible protocol for a messaging association has the following attributes:

- o MA-Protocol-ID, a 1-byte IANA assigned value (see [Section 9](#)).
- o A specification of the (non-negotiable) policies about how the protocol should be used; for example, in which direction a connection should be opened.
- o [Depending on the specific protocol:] Formats for an MA-protocol-options field to carry the protocol addressing and other configuration information in the Stack-Configuration-Data object. The format may differ depending on whether the field is present in

the Query or Response. Some protocols do not require the definition of such additional data, in which case no corresponding MA-protocol-options field will occur in the SCD object.

A Stack-Proposal object is simply a list of profiles; each profile is a sequence of MA-Protocol-IDs. A profile lists the protocols in 'top to bottom' order (e.g. TLS over TCP, or TCP over IPsec). A Stack-Proposal is generally accompanied by a Stack-Configuration-Data object which carries an MA-protocol-options field for any protocol listed in the Stack-Proposal which needs it. An MA-protocol-options field may apply globally, to all instances of the protocol in the Stack-Proposal; or it can be tagged as applying to a specific instance. The latter approach can be used to carry different port numbers for TCP depending on whether it is to be used with or without TLS. An MA-protocol-options field may also be flagged as not usable; for example, a NAT which could not handle SCTP would set this in an MA-protocol-options field about SCTP. A protocol flagged this way MUST NOT be used for a messaging association. If the Stack-Proposal and Stack-Configuration-Data are both present but not consistent, for example, if they refer to different protocols, or an MA-protocol-options field refers to a non-existent profile, an "Object Value Error" error message (Appendix A.4.4.10) with subcode 5 ("Stack-Proposal - Stack-Configuration-Data Mismatch") MUST be returned and the message dropped.

A node generating a Stack-Configuration-Data object MUST honour the

implied protocol configurations for the period during which a messaging association might be set up; in particular, it MUST be immediately prepared to accept incoming datagrams or connections at the protocol/port combinations advertised. This MAY require the creation of listening endpoints for the transport and security protocols in question, or a node MAY keep a pool of such endpoints open for extended periods. However, the received object contents MUST be retained only for the duration of the Query/Response exchange and to allow any necessary association setup to complete. They may become invalid because of expired bindings at intermediate NATs, or because the advertising node is using agile ports. Once the setup is complete, or if it is not necessary, or fails for some reason, the object contents MUST be discarded. A default time of 30 seconds to keep the contents is RECOMMENDED.

A Query requesting association setup always contains a Stack-Proposal and Stack-Configuration-Data object. The Stack-Proposal MUST only include protocol configurations that are suitable for the transfer attributes of the messages that the Querying node wishes use the messaging association for. For example, it should not simply include all configurations that the Querying node is capable of supporting.

The Response always contains a Stack-Proposal and Stack-Configuration-Data object, unless re-use (where the Responder decides to use an existing association) occurs. For such a Response, the Stack-Proposal MUST NOT depend on the Query. A node MAY make different proposals depending on the combination of interface and NSLPID. If re-use does occur, which is indicated by sending the Response over an existing messaging association, the following rules apply:

- o The re-used messaging association MUST NOT have weaker security properties than would have been offered in the full Response that would have been sent without re-use.
- o The re-used messaging association MUST have equivalent or better transport and security characteristics as at least one of the protocol configurations that was offered in the Query.

Once the messaging association is set up, the querying node repeats the responder's Stack-Proposal over it in the Confirm. The

responding node MUST verify that this has not been changed as part of bidding-down attack prevention. If a difference is detected, the responding node MUST terminate the messaging association and SHOULD log an error condition locally. See [Section 8.6](#) for further discussion.

[5.7.2.](#) Protocol Definition: Forwards-TCP

This MA-Protocol-ID denotes a basic use of TCP between peers. Support for this protocol is REQUIRED. If this protocol is offered, MA-protocol-options data MUST also be carried in the SCD object. The MA-protocol-options field formats are:

- o in a Query: no information apart from the field header.
- o in a Response: 2 byte port number at which the connection will be accepted, followed by 2 pad bytes.

The connection is opened in the forwards direction, from the querying node towards the responder. The querying node MAY use any source address and source port. The destination information MUST be derived from information in the Response: the address from the interface-address from the Network-Layer-Information object and the port from the SCD object as described above.

Associations using Forwards-TCP can carry messages with the transfer attribute Reliable=True. If an error occurs on the TCP connection such as a reset, as can be detected for example by a socket exception condition, GIST MUST report this to NSLPs as discussed in

[Section 4.1.2.](#)

[5.7.3.](#) Protocol Definition: Transport Layer Security

This MA-Protocol-ID denotes a basic use of transport layer channel security. Support for this protocol is mandatory; associations using it can carry messages with the transfer attribute Secure=True. For use with TCP, implementation of TLS1.0 [\[7\]](#) is REQUIRED and implementation of TLS1.1 [\[10\]](#) is RECOMMENDED. (If an unreliable transport such as DCCP or UDP is defined for GIST in the future, this MA-Protocol-ID would be implemented for it using DTLS [\[41\]](#).) GIST nodes supporting TLS1.0 or TLS1.1 MUST be able to negotiate the TLS

ciphersuite TLS_RSA_WITH_3DES_EDE_CBC_SHA and SHOULD be able to negotiate the TLS ciphersuite TLS_RSA_WITH_AES_128_CBC_SHA.

The default mode of TLS authentication, which applies in particular to the above ciphersuites, uses a client/server certificate exchange. The Querying node acts as a TLS client, and the Responding node acts as a TLS server. Where one of the above ciphersuites is negotiated, the GIST node acting as a server MUST provide a certificate, and MUST request one from the GIST node acting as a TLS client. This allows either server-only or mutual authentication, depending on the certificates available to the client and the policy applied at the server.

GIST nodes MAY negotiate other TLS ciphersuites. In some cases, the negotiation of alternative ciphersuites is used to trigger alternative authentication procedures, such as the use of pre-shared keys [29]. The use of other authentication procedures may require additional specification work to define how they can be used as part of TLS within the GIST framework, and may or may not require the definition of additional MA-Protocol-IDs.

No MA-protocol-options field is required for this use of TLS.

[5.7.4.](#) Additional Protocol Options

Further protocols or configurations could be defined in the future for additional performance or flexibility. Examples are:

- o SCTP or DCCP as alternatives to TCP, with essentially the same configuration.
- o SigComp [21] for message compression.
- o IPsec [35], ssh [36], or HIP/IPsec [37] for channel security.

- o Alternative modes of TCP operation, for example where it is set up from the responder to the querying node.

[5.8.](#) Specific Message Routing Methods

Each message routing method (see [Section 3.3](#)) requires the definition of the format of the message routing information (MRI) and Q-mode encapsulation rules. These are given in the following subsections for the various possible MRMs.

[5.8.1](#). The Path-Coupled MRM

[5.8.1.1](#). Message Routing Information

For the path-coupled MRM, this is conceptually the Flow Identifier as in the NSIS Framework [26]. Minimally, this could just be the flow destination address; however, to account for policy based forwarding and other issues a more complete set of header fields SHOULD be specified if possible (see [Section 4.3.4](#) and [Section 7.2](#) for further discussion).

```
MRI = network-layer-version
      source-address prefix-length
      destination-address prefix-length
      IP-protocol
      diffserv-codepoint
      [ flow-label ]
      [ ipsec-SPI / L4-ports]
```

Additional control information defines whether the flow-label, IPsec Security Parameters Index (SPI), and port information are present, and whether the IP-protocol and diffserv-codepoint fields should be interpreted as significant. The source and destination addresses MUST be real node addresses, but prefix lengths other than 32/128 (for IPv4/6) MAY be used to implement address wildcarding, allowing the MRI to refer to traffic to or from a wider address range.

The MRI format allows a potentially very large number of different flag and field combinations. A GIST implementation that cannot interpret the MRI in a message MUST return an "Object Value Error" message (Appendix A.4.4.10) with subcodes 1 ("Value Not Supported") or 2 ("Invalid Flag-Field Combination") and drop the message.

[5.8.1.2](#). Downstream Q-mode Encapsulation

Where the signaling message is travelling in the same ('downstream') direction as the flow defined by the MRI, the IP addressing for Q-mode encapsulated messages is as follows. Support for this

encapsulation is REQUIRED.

- o The destination IP address MUST be the flow destination address as given in the MRI of the message payload.
- o By default, the source address is the flow source address, again from the MRI. This provides the best likelihood that the message will be correctly routed through any region performing per-packet policy-based forwarding or load balancing which takes the source address into account. However, there may be circumstances where the use of the signaling source address is preferable, such as:
 - * In order to receive ICMP error messages about the signaling message (such as unreachable port or address). If these are delivered to the flow source rather than the signaling source, it will be very difficult for the querying node to detect that it is the last GIST node on the path.
 - * In order to receive GIST Error messages where the error message sender could not interpret the NLI in the original message.
 - * In order to attempt to run GIST through an unmodified NAT, which will only process and translate IP addresses in the IP header.

Because of these considerations, use of the signaling source address is allowed as an option, with use based on local policy. A node SHOULD use the flow source address for initial Query messages, but SHOULD transition to the signaling source address for some retransmissions or as a matter of static configuration, for example if a NAT is known to be in the path out of a certain interface. The S-flag in the common header tells the message receiver which option was used.

It is vital that the Query mimics the actual data flow as closely as possible, since this is the basis of how the signaling message is attached to the data path. To this end, GIST SHOULD set the DiffServ codepoint and (for IPv6) flow label to match the values in the MRI.

Any message sent in D-mode MUST have a size below a conservative estimate of the path MTU, for which this specification takes the value 512 bytes as a default. It is possible that fragmented datagrams including an RAO will not be correctly handled in the network, so the sender MAY set the Don't Fragment (DF) bit in the IPv4 header in order to detect that a message has encountered a link with an unusually low MTU. If the sender sets DF for any reason, it SHOULD use the signaling source address for the IP source address in order to receive the ICMP error.

Internet-Draft

GIST

July 2006

A GIST implementation SHOULD apply validation checks to the MRI, to reject Query messages that are being injected by nodes with no legitimate interest in the flow being signalled for. In general, if the GIST node can detect that no flow could arrive over the same interface as the Query, it MUST be rejected with an appropriate error message. Such checks apply only to messages with the Q-mode encapsulation, since only those messages are required to track the flow path. The main checks are that the IP version should match the version(s) used on that interface, and that the full range of source addresses (the source-address masked with its prefix-length) would pass ingress filtering checks. For these cases, the error message is "MRI Validation Failure" (Appendix A.4.4.12) with subcodes 1 or 2 ("IP Version Mismatch" or "Ingress Filter Failure") respectively.

[5.8.1.3](#). Upstream Q-mode Encapsulation

In some deployment scenarios it is desirable to set up routing state in the upstream direction, (i.e. from flow receiver towards the sender). This could be used to support firewall signaling to control traffic from an un-cooperative sender, or signaling in general where the flow sender was not NSIS-capable. This capability is incorporated into GIST by defining an encapsulation and processing rules for sending Query messages upstream.

In general, it is not possible to determine the hop-by-hop route upstream because of asymmetric routing. However, in particular cases, the upstream peer can be discovered with a high degree of confidence, for example:

- o The upstream GIST peer is 1 IP hop away, and can be reached by tracing back through the interface on which the flow arrives.
- o The upstream peer is a border router of a single-homed (stub) network.

This section defines an upstream Q-mode encapsulation and validation checks for when it can be used. The functionality to generate upstream Queries is OPTIONAL, but if received they MUST be processed in the normal way. No special functionality is needed for this.

It is possible for routing state at a given node, for a specific MRI and NSLPID, to be created by both an upstream Query exchange

(initiated by the node itself), and a downstream Query exchange (where the node is the responder). If the SIDs are different, these items of routing state **MUST** be considered as independent; if the SIDs match, the routing state installed by the downstream exchange **MUST** take precedence, provided that the downstream Query passed ingress filtering checks. The rationale for this is that the downstream

Query is in general a more reliable way to install state, since it directly probes the routing infrastructure along the flow path, whereas use of the upstream Query depends on the correctness of the Querying node's understanding of the topology.

The details of the encapsulation are as follows:

- o The destination address **SHOULD** be the flow source address as given in the MRI of the message payload. An implementation with more detailed knowledge of local routing **MAY** use an alternative destination address (e.g. the address of its default router).
- o The source address **SHOULD** be the signaling node address.
- o The DiffServ codepoint and (for IPv6) flow label **MAY** be set to match the values from the MRI, as in the downstream case. The same considerations about message size and fragmentation also apply as in the downstream case, and RAO setting and UDP port selection are also the same.
- o The IP layer TTL of the message **MUST** be set to 255.

The sending GIST implementation **SHOULD** attempt to send the Query via the same interface and to the same link layer neighbour from which the data packets of the flow are arriving.

The receiving GIST node **MAY** apply validation checks to the message and MRI, to reject Query messages which have reached a node at which they can no longer be trusted. In particular, a node **SHOULD** reject a message which has been propagated more than one IP hop, with an "Invalid IP layer TTL" error message (Appendix A.4.4.11). This can be determined by examining the received IP layer TTL, similar to the generalised IP TTL security mechanism described in [\[25\]](#). Alternatively, receipt of an upstream Query at the flow source **MAY** be used to trigger setup of GIST state in the downstream direction.

These restrictions may be relaxed in a future version.

[5.8.2.](#) The Loose-End MRM

The Loose-End MRM is used to discover GIST nodes with particular properties in the direction of a given address, for example to discover a NAT along the upstream data path as in [\[32\]](#).

[5.8.2.1.](#) Message Routing Information

For the loose-end MRM, only a simplified version of the Flow Identifier is needed.

```
MRI = network-layer-version
      source-address
      destination-address
```

The source address is the address of the node initiating the discovery process, for example the node that will be the data receiver in the NAT discovery case. The destination address is the address of a node which is expected to be the other side of the node to be discovered. Additional control information defines the direction of the message relative to this flow as in the path-coupled case.

[5.8.2.2.](#) Downstream Q-mode Encapsulation

Only one encapsulation is defined for the loose-end MRM; by convention, this is referred to as the downstream encapsulation, and is defined as follows:

- o The IP destination address **MUST** be the destination address as given in the MRI of the message payload.
- o By default, the IP source address is the source address, again from the MRI. However, the use of the signaling source address is allowed as in the case of the path-coupled MRM.

There are no special requirements on the setting of the DiffServ codepoint, IP layer TTL, or (for IPv6) the flow label. Nor are any special validation checks applied. Restrictions on message size and

setting of the Don't Fragment (DF) bit apply as in the case of the path-coupled MRM.

[6.](#) Formal Protocol Specification

This section provides a more formal specification of the operation of GIST processing, in terms of rules for transitions between states of a set of communicating state machines within a node. The following description captures only the basic protocol specification; additional mechanisms can be used by an implementation to accelerate route change processing, and these are captured in [Section 7.1](#).

Conceptually, GIST processing at a node may be seen in terms of four types of cooperating state machine:

1. There is a top-level state machine which represents the node itself (Node-SM). It is responsible for the processing of events which cannot be directed towards a more specific state machine, for example, inbound messages for which no routing state currently exists. This machine exists permanently, and is responsible for creating per-MRI state machines to manage the GIST handshake and routing state maintenance procedures.
2. For each flow and signaling direction where the node is

responsible for the creation of routing state, there is an instance of a Query-Node state machine (Querying-SM). This machine sends Query and Confirm messages and waits for Responses, according to the requirements from local API commands or timer processing, such as message repetition or routing state refresh.

3. For each flow and signaling direction where the node has accepted the creation of routing state by a peer, there is an instance of a Responding-Node state machine (Responding-SM). This machine is responsible for managing the status of the routing state for that flow. Depending on policy, it MAY be responsible for [re]transmission of Response messages, or this MAY be handled by the Node-SM, and a Responding-SM is not even created for a flow until a properly formatted Confirm has been accepted.
4. Messaging associations have their own lifecycle, represented by MA-SM, from when they are first created (in an incomplete state, listening for an inbound connection or waiting for outbound connections to complete), to when they are active and available for use.

Apart from the fact that the various machines can be created and destroyed by each other, there is almost no interaction between them. The machines for different flows do not interact; the Querying-SM and Responding-SM for a single flow and signaling direction do not interact. That is, the Responding-SM which accepts the creation of routing state for a flow on one interface has no direct interaction

with the Querying-SM which sets up routing state on the next interface along the path. This interaction is mediated instead through the NSLP.

The state machine descriptions use the terminology rx_MMMM, tg_TTTT and er_EEEE for incoming messages, API/lower layer triggers and error conditions respectively. The possible events of these types are given in the table below. In addition, timeout events denoted to_TTTT may also occur; the various timers are listed independently for each type of state machine in the following subsections.

+-----+-----+	
Name	Meaning
+-----+-----+	

rx_Query	A Query has been received.
rx_Response	A Response has been received.
rx_Confirm	A Confirm has been received.
rx_Data	A Data message has been received.
rx_Message	rx_Query rx_Response rx_Confirm rx_Data.
rx_MA-Hello	A MA-Hello message has been received.
tg_NSLPData	A signaling application has requested data transfer (via API SendMessage).
tg_Connected	The protocol stack for a messaging association has completed connecting.
tg_RawData	GIST wishes to transfer data over a particular messaging association.
er_NoRSM	A "No Routing State" error was received.
er_MACconnect	A messaging association protocol failed to complete a connection.
er_MAFailure	A messaging association failed.

[6.1.](#) Node Processing

The Node level state machine is responsible for processing events for which no more appropriate messaging association state or routing state exists. Its structure is trivial: there is a single state

('Idle'); all events cause a transition back to Idle. Some events cause the creation of other state machines. The only events that are processed by this state machine are incoming GIST messages (Query/Response/Confirm/Data) and API requests to send data; no other events are possible. In addition to this event processing, the Node level machine is responsible for managing listening endpoints for messaging associations. Although these relate to Responding node operation,

they cannot be handled by the Responder state machine since they are not created per flow. The processing rules for each event are as follows:

Rule 1 (rx_Query):

```
use the GIST service interface to determine the signaling application
  policy relating to this peer
if (the signaling application indicates that routing state should
  be created) then
  if (routing state can be created without a 3-way handshake) then
    create Responding-SM and transfer control to it
  else
    send Response
else
  propagate the Query with any updated NSLP payload provided
```

Rule 2 (rx_Response):

```
// should already have a Querying-SM to handle this
discard message
send "No Routing State" error message
```

Rule 3 (rx_Confirm):

```
if (routing state can be created before receiving a Confirm) then
  // we should already have Responding-SM for it,
  // which would handle this message
  discard message
  send "No Routing State" error message
else
  create Responding-SM and pass message to it
```

Rule 4 (rx_Data):

```
if (node policy will only process Data messages with matching
  routing state) then
  send "No Routing State" error message
else
  pass directly to NSLP
```

Rule 5 (tg_NSLPData):

```
if Q-mode encapsulation is not possible for this MRI
    reject message with an error
else
    if (local policy & transfer attributes say routing
        state is not needed) then
        send message statelessly
    else
        create Querying-SM and pass message to it
```

[6.2.](#) Query Node Processing

The Querying-Node state machine (Querying-SM) has three states:

- o Awaiting Response
- o Established
- o Awaiting Refresh

The Querying-SM is created by the Node-SM machine as a result of a request to send a message for a flow in a signaling direction where the appropriate state does not exist. The Query is generated immediately and the No_Response timer is started. The NSLP data MAY be carried in the Query if local policy and the transfer attributes allow it, otherwise it MUST be queued locally pending MA establishment. Then the machine transitions to the Awaiting Response state, in which timeout-based retransmissions are handled. Data messages (rx_Data events) should not occur in this state; if they do, this may indicate a lost Response and a node MAY also retransmit a Query for this reason.

Once a Response has been successfully received and routing state created, the machine transitions to Established, during which NSLP data can be sent and received normally. Further Responses received in this state (which may be the result of a lost Confirm) MUST be treated the same way. The Awaiting Refresh state can be considered as a substate of Established, where a new Query has been generated to refresh the routing state (as in Awaiting Response) but NSLP data can be handled normally.

The timers relevant to this state machine are as follows:

Refresh_QNode: Indicates when the routing state stored by this state machine must be refreshed. It is reset whenever a Response is received indicating that the routing state is still valid. Implementations MUST set the period of this timer based on the

Internet-Draft

GIST

July 2006

value in the RS-validity-time field of a Response to ensure that a Query is generated before the peer's routing state expires.

No_Response: Indicates that a Response has not been received in answer to a Query. This is started whenever a Query is sent and stopped when a Response is received.

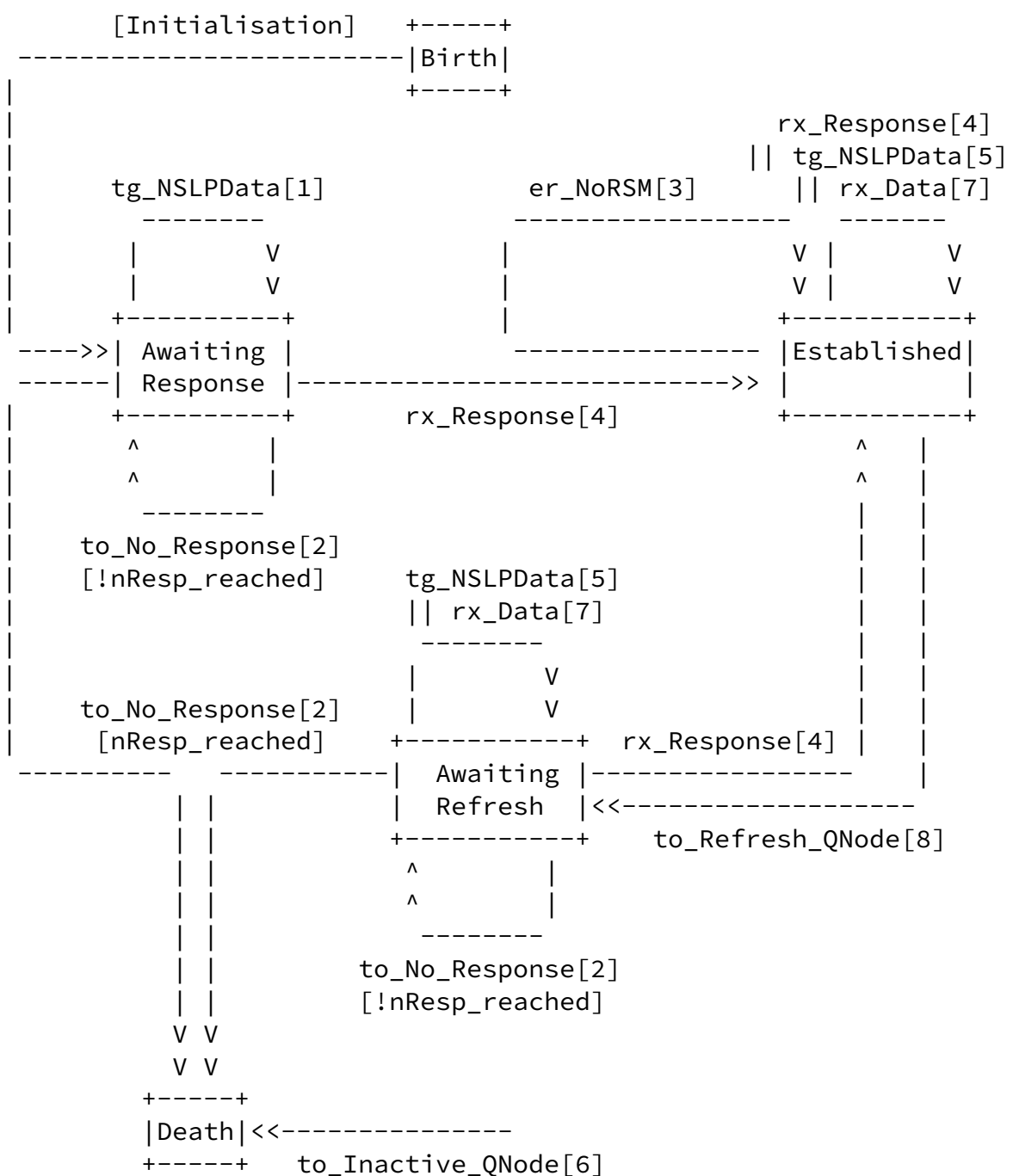
Inactive_QNode: Indicates that no traffic is currently being handled by this state machine. This is reset whenever the state machine handles NSLP data, in either direction. When it expires, the state machine MAY be deleted. The period of the timer can be set at any time via the API (SetStateLifetime), and if the period is reset in this way the timer itself MUST be restarted.

The main events (including all those that cause state transitions) are shown in the figure below, tagged with the number of the processing rule that is used to handle the event. These rules are listed after the diagram. All events not shown or described in the text above are assumed to be impossible in a correct implementation.

Internet-Draft

GIST

July 2006



(from all states)

Figure 6: Query Node State Machine

The processing rules are as follows:

Rule 1: Store the message for later transmission

Rule 2:

```
if number of Queries sent has reached the threshold
  // nQuery_isMax is true
  indicate No Response error to NSLP
  destroy self
else
  send Query
  start No_Response timer with new value
```

Rule 3:

```
// Assume the Confirm was lost in transit so resend it
// for the last Response we received
send Confirm
restart Refresh_QNode and Inactive_QNode timers
```

Rule 4:

```
if a new MA-SM is needed create one
if the R flag was set send a Confirm
pass any NSLP data to the NSLP
send any stored Data messages
stop No_Response timer
start Refresh_QNode and Inactive_QNode timers
```

Rule 5:

```
send Data message
restart Inactive_QNode timer
```

Rule 6: Terminate

Rule 7:

pass any data to the NSLP
(re)start Inactive_QNode timer

Rule 8:

send Query
start No_Response timer
stop Refresh_QNode timer

[6.3.](#) Responder Node Processing

The Responding-Node state machine (Responding-SM) has three states:

- o Awaiting Confirm

- o Established
- o Awaiting Refresh

The policy governing the creation of the Responding-SM has three cases:

1. It is created on receiving a Query, no Confirm is requested.
2. It is created on receiving a Query, but a Confirm is requested. A timer is used to retransmit Response messages and the Responding-SM is destroyed if no valid Confirm is received.
3. It cannot be created until a valid Confirm is received; the initial Query will have been handled by the Node level machine.

In case 2 the Responding-SM is created in the Awaiting Confirm state, and remains there until a Confirm is received, at which point it transitions to Established. In cases 1 and 3 the Responding-SM is created directly in the Established state. Note that if the machine is created on receiving a Query, some of the message processing will already have been performed in the Node state machine. In the

Established state the NSLP can send and receive data normally, and any additional rx_Confirm events MUST be silently ignored. The Awaiting Refresh state can be considered a substate of Established, where a Query has been received to begin the routing state refresh.

In the Awaiting Refresh state the Responding-SM behaves as in the Awaiting Confirm state, except that the NSLP can still send and receive data. In particular, in both states there is timer-based retransmission of Response messages until a Confirm is received; additional rx_Query events in these states MUST also generate a response and restart the no_Confirm timer.

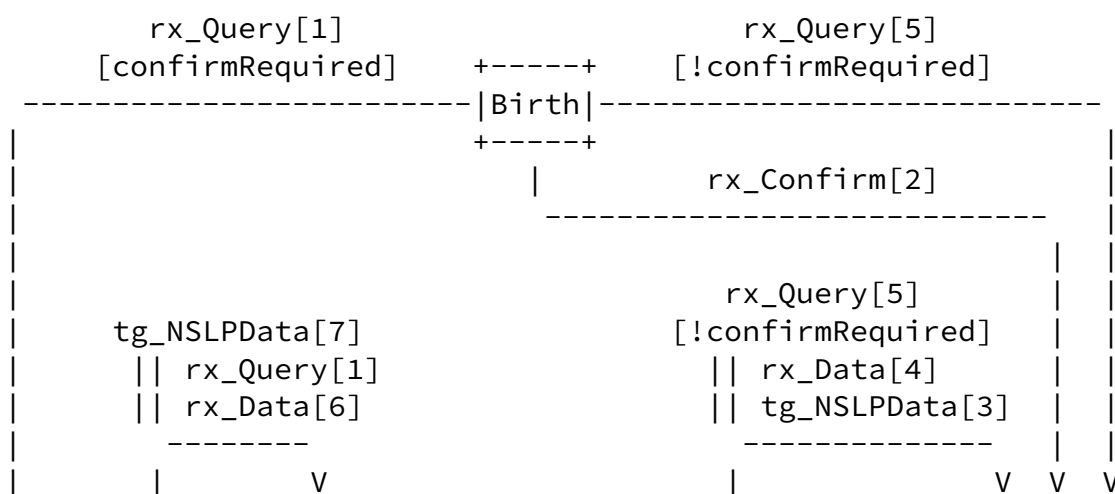
The timers relevant to the operation of this state machine are as follows:

Expire_RNode: Indicates when the routing state stored by this state machine needs to be expired. It is reset whenever a Query or Confirm (depending on local policy) is received indicating that the routing state is still valid. Note that state cannot be refreshed from the R-Node.

No_Confirm: Indicates that a Confirm has not been received in answer to a Response. This is started/reset whenever a Response is sent and stopped when a Confirm is received.

The detailed state transitions and processing rules are described

below as in the Query node case.



Rule 2:

pass any piggybacked data to the NSLP
start Expire_RNode timer

Rule 3: send the Data message

Rule 4: pass data to NSLP

Rule 5:

// no Confirm is required
send Response
start Expire_RNode timer

Rule 6: send "No Routing State" error message

Rule 7: store Data message

Rule 8:

pass any piggybacked data to the NSLP
send any stored Data messages
stop No_Confirm timer
start Expire_RNode timer

Rule 9:

if number of Responses sent has reached threshold
 // nResp_isMax is true
 destroy self
else
 send Response
 start No_Response timer

Rule 10: destroy self

[6.4.](#) Messaging Association Processing

Messaging associations (MAs) are modelled for use within GIST with a simple three-state process. The Awaiting Connection state indicates that the MA is waiting for the connection process(es) for every

protocol in the messaging association to complete; this might involve

creating listening endpoints or attempting active connects. Timers may also be necessary to detect connection failure (e.g. no incoming connection within a certain period), but these are not modelled explicitly. The Connected state indicates that the MA is open and ready to use. In addition there is an Idle state in which the local node no longer requires the messaging association but the remote node still wants it to be kept open.

Clearly, many internal details of the messaging association protocols are hidden in this model, especially where the messaging association uses multiple protocol layers. Note also that although the existence of messaging associations is not directly visible to signaling applications, there is some interaction between the two because security-related information becomes available during the open process, and this may be indicated to signaling applications if they have requested it.

The timers relevant to the operation of this state machine are as follows:

SendHello: Indicates that an MA-Hello message should be sent to the remote node. The period of this timer is determined by the MA-Hold-Time sent by the remote node during the Query/Response/Confirm exchange.

NoHello: Indicates that no MA-Hello has been received from the remote node for a period of time. The period of this timer is sent to the remote node as the MA-Hold-Time during the Query/Response exchange.

NoActivity: Indicates that the link has been inactive for a period of time. The period of this timer is implementation-specific but is likely to be related to the period of the NoHello timer.

The detailed state transitions and processing rules are described below as in the Query node case.

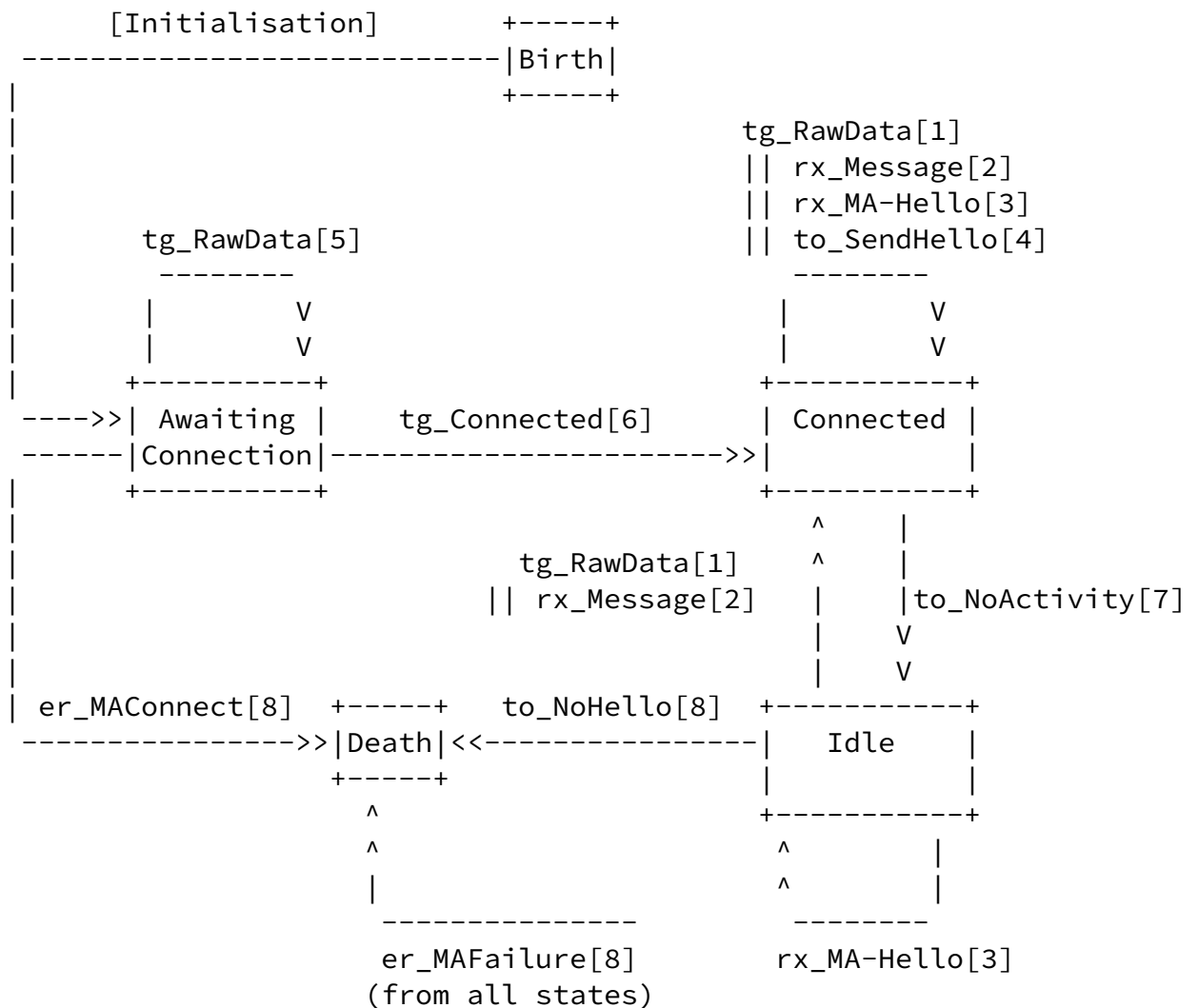


Figure 8: Messaging Association State Machine

The processing rules are as follows:

Rule 1:

```

pass message to transport layer
(re)start NoActivity timer
(re)start SendHello

```

Rule 2:

```

pass message to Node-SM
(re)start NoActivity timer

```

Internet-Draft

GIST

July 2006

Rule 3:

```
if reply requested
  send MA-Hello
restart SendHello timer
```

Rule 4:

```
send MA-Hello message
restart NoHello timer
```

Rule 5: queue message for later transmission

Rule 6:

```
pass outstanding queued messages to transport layer
stop any timers controlling connection establishment
start NoActivity timer
start SendHello timer
```

Rule 7:

```
stop NoActivity timer
stop sendHello timer
start NoHello timer
```

Rule 8: destroy self

Internet-Draft

GIST

July 2006

7. Advanced Protocol Features

7.1. Route Changes and Local Repair

7.1.1. Introduction

When re-routing takes place in the network, GIST and signaling application state need to be updated for all flows whose paths have changed. The updates to signaling application state depend mainly on the signaling application: for example, if the path characteristics have actually changed, simply moving state from the old to the new path is not sufficient. Therefore, GIST cannot carry out the complete path update processing. Its responsibilities are to detect the route change, update its local routing state consistently, and inform interested signaling applications at affected nodes.

```

                                xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
                                x  +--+      +--+      +--+  x      Initial
                                x  .|C1|_.....|D1|_.....|E1|  x      Configuration
                                x  . +--+ .      .+--+ .      .+--+ \.  x
>>xxxxxxxxxxxxxxxxx .      .      .      .      .      . xxxxxx>>
    +-+      +-+ .      ..      ..      .      +-+
...|A|_.....|B|/      ..      ..      .|F|_.....
    +-+      +-+ .      .      .      .      .      +-+
                                .
                                . +--+      +--+      +--+ .
                                .|C2|_.....|D2|_.....|E2|/
                                +--+      +--+      +--+

                                +--+      +--+      +--+
                                .|C1|_.....|D1|_.....|E1|
                                . +--+      .+--+      +--+
                                Configuration
                                after failure
                                of E1-F link

```

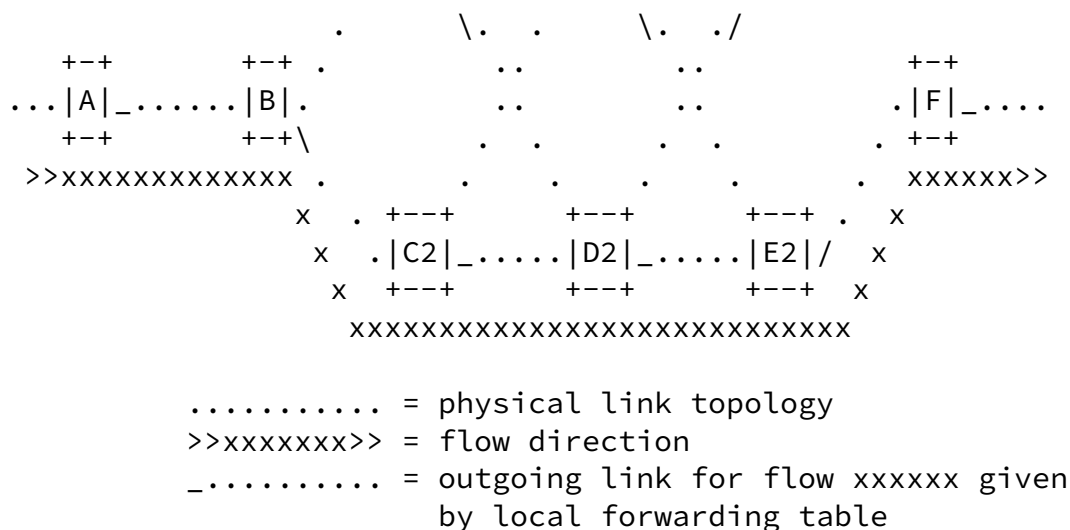


Figure 9: A Re-Routing Event

Route change management is complicated by the distributed nature of the problem. Consider the re-routing event shown in Figure 9. An external observer can tell that the main responsibility for controlling the updates will probably lie with nodes B and F; however, E1 is best placed to detect the event quickly at the GIST level, and C1 and D1 could also attempt to initiate the repair.

On the assumption that signaling applications are soft-state based and operate end to end, and because GIST also periodically updates its picture of routing state, route changes will eventually be repaired automatically. The specification as already given includes this functionality. However, especially if upper layer refresh times are extended to reduce signaling load, the duration of inconsistent state may be very long indeed. Therefore, GIST includes logic to exchange prompt notifications with signaling applications, to allow local repair if possible. The additional mechanisms to achieve this are described in the following subsections. To a large extent, these additions can be seen as implementation issues; the protocol messages and their significance are not changed, but there are extra interactions through the API between GIST and signaling applications, and additional triggers for transitions between the various GIST states.

[7.1.2.](#) Route Change Detection Mechanisms

There are two aspects to detecting a route change at a single node:

- o Detecting that the outgoing path, in the direction of the Query, has or may have changed.
- o Detecting that the incoming path, in the direction of the Response, has (or may have) changed, in which case the node may no longer be on the path at all.

At a single node, these processes are largely independent, although clearly a change in one direction at a node corresponds to a change the opposite direction at its peer. Note that there are two possible forms for a route change: the interface through which a flow leaves or enters a node may change, and the adjacent peer may change. In general, a route change can include one or the other or both (or indeed neither, although such changes are very hard to detect).

The route change detection mechanisms available to a node depend on the MRM in use and the role the node played in setting up the routing state in the first place (i.e. as Querying or Responding node). The following discussion is specific to the case of the path-coupled MRM using downstream Queries only; other scenarios may require other methods. However, the repair logic described in the subsequent

subsections is intended to be universal.

There are five mechanisms for a node to detect that a route change has occurred, which are listed below. They apply differently depending on whether the change is in the Query or Response direction, and these differences are summarised in the following table.

Local Trigger: In local trigger mode, GIST finds out from the local forwarding table that the next hop has changed. This only works if the routing change is local, not if it happens a few routing hops away, including the case that it happens at a GIST-unaware node.

Extended Trigger: Here, GIST checks a link-state topology database to discover that the path has changed. This makes certain assumptions on consistency of route computation and only works within a single area for OSPF [\[13\]](#) and similar link-state

protocols. Where available, this offers the most accurate and rapid indication of route changes, but requires more access to the routing internals than a typical operating system may provide.

GIST C-mode Monitoring: GIST may find that C-mode packets are arriving (from either peer) with a different IP layer TTL or on a different interface. This provides no direct information about the new flow path, but indicates that routing has changed and that rediscovery may be required.

Data Plane Monitoring: The signaling application on a node may detect a change in behaviour of the flow, such as IP layer TTL change, arrival on a different interface, or loss of the flow altogether. The signaling application on the node is allowed to notify this information locally to GIST (Appendix B.6).

GIST Probing: According to the specification, each GIST node **MUST** periodically repeat the discovery (Query/Response) operation. Values for the probe frequency are discussed in [Section 4.4.3](#). The querying node will discover the route change by a modification in the Network-Layer-Information in the Response. The period can be negotiated independently for each GIST hop, so nodes that have access to the other techniques listed above **MAY** use long periods for the probing operation.

Method	Query direction	Response direction
Local Trigger	Discovers new interface (and peer if local)	Not applicable
Extended Trigger	Discovers new interface and may determine new peer	May determine that route from peer will have changed
C-mode	Provides hint that	Provides hint that change

Monitoring	change has occurred	has occurred	
Data Plane	Not applicable	NSLP informs GIST that a	
Monitoring		change may have occurred	
Probing	Discovers changed NLI in	Discovers changed NLI in	
	Response	Query	
+-----+-----+-----+			

[7.1.3.](#) GIST Behaviour Supporting Re-Routing

The GIST behaviour necessary to support re-routing can be modelled using a 3-level classification of the validity of each item of routing state. (This classification applies separately to the Querying and Responding node for each pair of GIST peers.) The levels are:

Bad: The routing state is either missing altogether, or not safe to use to send data.

Tentative: The routing state may have changed, but it is still usable for sending NSLP data pending verification.

Good: The routing state has been established and no events affecting it have since been detected.

These classifications are not identical to the states described in [Section 6](#), but there are dependencies between them. Specifically, routing state is considered Bad until the machine first enters the Established state, at which point it becomes Good. Thereafter, the status may be invalidated for any of the reasons discussed above; it is an implementation issue to decide which techniques to implement in any given node, and how to reclassify routing state (as Bad or Tentative) for each. The status returns to Good, either when the state machine re-enters the Established state, or if GIST can determine from direct examination of the routing or forwarding tables that the peer has not changed. When the status returns to Good, GIST

MUST if necessary update its routing state table so that the relationships between MRI/SID/NSLPID tuples and messaging associations are up to date.

When classification of the routing state for the downstream direction changes to Bad/Tentative because of local routing indications, GIST MAY automatically change the classification in the upstream direction to Tentative unless local routing indicates that this is not necessary. This SHOULD NOT be done in the case where the initial change was indicated by the signaling application. This mechanism accounts for the fact that a routing change may affect several nodes, and so can be an indication that upstream routing may also have changed. In any case, whenever GIST updates the routing status, it informs the signaling application with the NetworkNotification API (Appendix B.4), unless the change was caused via the API in the first place.

The GIST behaviour for state repair is different for the Querying and Responding node. At the Responding node, there is no additional behaviour, since the Responding node cannot initiate protocol transitions autonomously, it can only react to the Querying node. The Querying node has three options, depending on how the transition from 'Good' was initially caused:

1. To inspect the routing/forwarding table and verifying that the next peer has not changed. This technique MUST NOT be used if the transition was caused by a signaling application, but SHOULD be used otherwise if available.
2. To move to the 'Awaiting Refresh' state. This technique MUST NOT be used if the current status is 'Bad', since data is being incorrectly delivered.
3. To move to the 'Awaiting Response' state. This technique may be used at any time, but has the effect of freezing NSLP communication while GIST state is being repaired.

The second and third techniques trigger the execution of a GIST handshake to carry out the repair. It may be desirable to delay the start of the handshake process, either to wait for the network to stabilise, to avoid flooding the network with Query traffic for a large number of affected flows, or to wait for confirmation that the node is still on the path from the upstream peer. One approach is to delay the handshake until there is NSLP data to be transmitted. Implementation of such delays is a matter of local policy; however, GIST MUST begin the handshake immediately if the status change was caused by an InvalidateRoutingState API call marked as 'Urgent', and SHOULD begin it if the upstream routing state is still known to be

Good.

[7.1.4.](#) Signaling Application Operation

Signaling applications can use these functions as provided by GIST to carry out rapid local repair following re-routing events. The signaling application instances carry out the multi-hop aspects of the procedure, including crossover node detection, and tear-down/reinstallation of signaling application state; they also trigger GIST to carry out the local routing state maintenance operations over each individual hop. The local repair procedures depend heavily on the fact that stateful NSLP nodes are a single GIST hop apart; this is enforced by the details of the GIST peer discovery process.

The following outline description of a possible set of NSLP actions takes the scenario of Figure 9 as an example.

1. The signaling application at node E1 is notified by GIST of route changes affecting the downstream and upstream directions. The downstream status was updated to Bad because of a trigger from the local forwarding tables, and the upstream status changed automatically to Tentative as a consequence. The signaling application at E1 MAY begin local repair immediately, or MAY propagate a notification upstream to D1 that re-routing has occurred.
2. The signaling application at node D1 is notified of the route change, either by signaling application notifications or from the GIST level (e.g. by a trigger from a link-state topology database). If the information propagates faster within the routing protocol, GIST will change the upstream/downstream routing state to Tentative/Bad automatically, and this will cause the signaling application to propagate the notification further upstream.
3. This process continues until the notification reaches node A. Here, there is no downstream routing change, so GIST only learns of the update via the signaling application trigger. Since the upstream status is still Good, it therefore begins the repair handshake immediately.
4. The handshake initiated by node A causes its downstream routing state to be confirmed as Good and unchanged there; it also confirms the (Tentative) upstream routing state at B as Good. This is enough to identify B as the crossover router, and the signaling application and GIST can begin the local repair process.

Internet-Draft

GIST

July 2006

An alternative way to reach step (4) is that node B is able to determine autonomously that there is no likelihood of an upstream route change. For example, it could be an area border router and the route change is only intra-area. In this case, the signaling application and GIST will see that the upstream state is Good and can begin the local repair directly.

After a route change, a signaling application may wish to remove state at another node which is no longer on the path. However, since it is no longer on the path, in principle GIST can no longer send messages to it. In general, provided this state is soft, it will time out anyway; however, the timeouts involved may have been set to be very long to reduce signaling load. The requirement to remove state in a specific peer node is identified in [\[33\]](#).

This requirement can be met provided that GIST is able to use the old path to the signaling application peer for some period while the signaling application still needs it. Since NSLP peers are a single GIST hop apart, the necessary information is just the old entry in the node's routing state table for that flow. Rather than requiring GIST to maintain multiple generations of this information, it can just be provided to the signaling application in the same node in an opaque form for each message that is received. The signaling application can store it if necessary and provide it back to the GIST layer in case it needs to be used. Because this is a reference to information about the source of a prior signaling message, it is denoted 'SII-Handle' (for Source Identification Information) in the abstract API of [Appendix B](#). Note that GIST if possible SHOULD use the same SII-Handle for multiple sessions to the same peer, since this then allows signaling applications to aggregate some signaling, such as summary refreshes or bulk teardowns.

Messages sent this way MUST bypass the GIST routing state tables at the sender, and this is indicated by setting the E flag in the common header (Appendix A.1); at the receiver, GIST MUST NOT validate the MRI/SID/NSLPID against local routing state and instead indicates the mode of reception to signaling applications through the API (Appendix B.2). Signaling applications should validate the source and effect of the message themselves, and if appropriate should in particular indicate to GIST (see [Appendix B.5](#)) that routing state is

no longer required for this flow. This is necessary to prevent GIST in nodes on the old path initiating routing state refresh and thus causing state conflicts at the crossover router.

[7.2.](#) NAT Traversal

[7.2.1.](#) Overview

GIST messages must carry packet addressing and higher layer information as payload data in order to define the flow signalled for. (This applies to all GIST messages, regardless of how they are encapsulated or which direction they are travelling in.) At an addressing boundary the data flow packets will have their headers translated; if the signaling payloads are not translated consistently, the signaling messages will refer to incorrect (and probably meaningless) flows after passing through the boundary. In addition, GIST handshake messages carry additional addressing information about the GIST nodes themselves, and this must also be processed appropriately when traversing a NAT.

The simplest solution to this problem is to require that a NAT is GIST-aware, and to allow it to modify messages based on the contents of the MRI. This makes the assumption that NATs only rewrite the header fields included in this payload, and not other higher layer identifiers. Provided this is done consistently with the data flow header translation, signaling messages will be valid each side of the boundary, without requiring the NAT to be signaling application aware. Note, however, that if the NAT does not understand the MRI, and the N-flag in the MRI is clear (see [Appendix A.3.1](#)), it should reject the message with an "Object Type Error" message (Appendix A.4.4.9) with subcode 4 ("Untranslated Object").

This specification defines an additional object that a NAT can insert into Q-mode encapsulated messages and which is echoed back in any replies, i.e. Response or Error messages. NATs apply GIST-specific processing only to Q-mode encapsulated messages or replies carrying the NAT traversal object. All other GIST messages, either in C-mode, or D-mode messages with no NAT-Traversal object, should be treated as normal data traffic by the NAT, i.e. with IP and transport layer

header translation but no GIST-specific processing.

The new object, the NAT-Traversal object (Appendix A.3.8), carries the translation between the MRIs which are appropriate for the internal and external sides of the NAT. It also carries a list of which other objects in the message have been translated. This should always include the NLI, and the Stack-Configuration-Data if present; if GIST is extended with further objects that carry addressing data, this list allows a message receiver to know if the new objects were supported by the NAT. Finally, the NAT-Traversal object MAY be used to carry data to be used in back-translating D-mode responses; this could be the original NLI or SCD, or opaque equivalents in the case of topology hiding.

A consequence of this approach is that the routing state tables at

the signaling application peers each side of the NAT are no longer directly compatible. In particular, the values for Message-Routing-Information are different, which is why the unmodified MRI is propagated in the NAT-Traversal object to allow subsequent C-mode messages to be interpreted correctly.

7.2.2. Message Processing Rules

This specification normatively defines the behaviour of a GIST node receiving a message containing a NAT-Traversal object. However, it does not define normative behaviour for a NAT translating GIST messages, since much of this will depend on NAT implementation and policy about allocating bindings. In addition, it is not necessary for a GIST implementation itself. Therefore, those aspects of the following description are informative; full details of NAT behaviour for handling GIST messages can be found in [\[40\]](#).

A possible set of operations for a NAT to process a Q-mode encapsulated message is as follows:

1. Verify that bindings for any data flow are actually in place.
2. Create a new Message-Routing-Information object with fields modified according to the data flow bindings.
3. Create bindings for subsequent C-mode signaling based on the

information in the Network-Layer-Information and Stack-Configuration-Data objects.

4. Create new Network-Layer-Information and if necessary Stack-Configuration-Data objects with fields to force D-mode response messages through the NAT, and to allow C-mode exchanges using the C-mode signaling bindings.
5. Add a NAT-Traversal object, listing the objects which have been modified and including the unmodified MRI and any other data needed to interpret the response. If a NAT-Traversal object is already present, in the case of a sequence of NATs, the list of modified objects may be updated and further opaque data added, but the MRI contained in it is left unchanged.
6. Encapsulate the message according to the normal rules of this specification for the Q-mode encapsulation. If the S-flag was set in the original message, the same IP source address selection policy should be applied to the forwarded message.
7. Forward the message with these new payloads.

A GIST node receiving such a message MUST verify that all mandatory objects containing addressing have been translated correctly, or else reject the message with an "Object Type Error" message (Appendix A.4.4.9) with subcode 4 ("Untranslated Object"). The error message MUST include the NAT-Traversal object as the first TLV after the common header, and this is also true for any other error message generated as a response. Otherwise, the message is processed essentially as normal. If no state needs to be updated for the message, the NAT-Traversal object can be effectively ignored. The other possibility is that a Response must be returned, either because the message is the beginning of a handshake for a new flow, or it is a refresh for existing state. In both cases, the GIST node MUST create the Response in the normal way using the local form of the MRI, and its own NLI and (if necessary) SCD. It MUST also include the NAT-Traversal object as the first object in the Response after the common header.

A NAT will intercept D-mode messages with the normal encapsulation containing such echoed NAT-Traversal objects. The NAT processing is

a subset of the processing for the Q-mode encapsulated case:

1. Verify the existence of bindings for the data flow.
2. Leave the Message-Routing-Information object unchanged.
3. Modify the NLI and SCD objects for the Responding node if necessary, and create or update any bindings for C-mode signaling traffic.
4. Forward the message.

A GIST node receiving such a message MUST use the MRI from the NAT-Traversal object as the key to index its internal routing state; it MAY also store the translated MRI for additional (e.g. diagnostic) information, but this is not used in the GIST processing. The remainder of GIST processing is unchanged.

Note that Confirm messages are not given GIST-specific processing by the NAT. Thus, a Responding node which has delayed state installation until receiving the Confirm, only has available the untranslated MRI describing the flow, and the untranslated NLI as peer routing state. This would prevent the correct interpretation of the signaling messages; also, subsequent Query (refresh) messages would always be seen as route changes because of the NLI change. Therefore, a Responding node that wishes to delay state installation until receiving a Confirm must somehow reconstruct the translations when the Confirm arrives. How to do this is an implementation issue; one approach is to carry the translated objects as part of the

Responder cookie which is echoed in the Confirm. Indeed, for one of the cookie constructions in [Section 8.5](#) this is automatic.

[7.3.](#) Interaction with IP Tunnelling

The interaction between GIST and IP tunnelling is very simple. An IP packet carrying a GIST message is treated exactly the same as any other packet with the same source and destination addresses: in other words, it is given the tunnel encapsulation and forwarded with the other data packets.

Tunnelled packets will not be identifiable as GIST messages until

they leave the tunnel, since any router alert option and the standard GIST protocol encapsulation (e.g. port numbers) will be hidden within the standard tunnel encapsulation. If signaling is needed for the tunnel itself, this has to be initiated as a separate signaling session by one of the tunnel endpoints - that is, the tunnel counts as a new flow. Because the relationship between signaling for the microflow and signaling for the tunnel as a whole will depend on the signaling application in question, it is a signaling application responsibility to be aware of the fact that tunnelling is taking place and to carry out additional signaling if necessary; in other words, at least one tunnel endpoint must be signaling application aware.

In some cases, it is the tunnel exit point (i.e. the node where tunnelled data and downstream signaling packets leave the tunnel) that will wish to carry out the tunnel signaling, but this node will not have knowledge or control of how the tunnel entry point is carrying out the data flow encapsulation. The information about how the inner MRI/SID relate to the tunnel MRI/SID needs to be carried in the signaling data from the tunnel entry point; this functionality is the equivalent to the RSVP SESSION_ASSOC object of [14]. In the NSIS protocol suite, these bindings are managed by the signaling applications, either implicitly (e.g. by SID re-use) or explicitly by carrying objects that bind the inner and outer SIDs as part of the NSLP payload.

[7.4.](#) IPv4-IPv6 Transition and Interworking

GIST itself is essentially IP version neutral: version dependencies are isolated in the formats of the Message-Routing-Information, Network-Layer-Information and Stack-Configuration-Data objects, and GIST also depends on the version independence of the protocols that support messaging associations. In mixed environments, GIST operation will be influenced by the IP transition mechanisms in use. This section provides a high level overview of how GIST is affected, considering only the currently predominant mechanisms.

Dual Stack: (As described in [34].) In mixed environments, GIST MUST use the same IP version for Q-mode encapsulated messages as the flow it is signaling for, and SHOULD do so for other signaling also (see [Section 5.2.2](#)). The IP version used in D-mode is closely tied to the IP version used by the data flow, so it is

intrinsically impossible for an IPv4-only or IPv6-only GIST node to support signaling for flows using the other IP version. Hosts which are dual stack for applications and routers which are dual stack for forwarding need GIST implementations which can support both IP versions. Applications with a choice of IP versions might select a version based on which could be supported in the network by GIST, which could be established by invoking parallel discovery procedures.

Packet Translation: (Applicable to SIIT [\[6\]](#) and NAT-PT [\[15\]](#).) Some transition mechanisms allow IPv4 and IPv6 nodes to communicate by placing packet translators between them. From the GIST perspective, this should be treated essentially the same way as any other NAT operation (e.g. between internal and external addresses) as described in [Section 7.2](#). The translating node needs to be GIST-aware; it will have to translate the addressing payloads between IPv4 and IPv6 formats for flows which cross between the two. The translation rules for the fields in the MRI payload (including e.g. DiffServ-codepoint and flow-label) are as defined in [\[6\]](#).

Tunnelling: (Applicable to 6to4 [\[17\]](#).) Many transition mechanisms handle the problem of how an end to end IPv6 (or IPv4) flow can be carried over intermediate IPv4 (or IPv6) regions by tunnelling; the methods tend to focus on minimising the tunnel administration overhead.

From the GIST perspective, the treatment should be as similar as possible to any other IP tunnelling mechanism, as described in [Section 7.3](#). In particular, the end to end flow signaling will pass transparently through the tunnel, and signaling for the tunnel itself will have to be managed by the tunnel endpoints. However, additional considerations may arise because of special features of the tunnel management procedures. In particular, [\[18\]](#) is based on using an anycast address as the destination tunnel endpoint. GIST MAY use anycast destination addresses in the Q-mode encapsulation of D-mode messages if necessary, but MUST NOT use them in the Network-Layer-Information addressing field; normal unicast addresses MUST be used instead. Note that the addresses from the IP header are not used by GIST in matching requests and responses, so there is no requirement to use anycast source addresses.

8. Security Considerations

The security requirement for GIST is to protect the signaling plane against identified security threats. For the signaling problem as a whole, these threats have been outlined in [\[27\]](#); the NSIS framework [\[26\]](#) assigns a subset of the responsibilities to the NTLP. The main issues to be handled can be summarised as:

Message Protection: Signaling message content can be protected against eavesdropping, modification, injection and replay while in transit. This applies both to GIST payloads, and GIST should also provide such protection as a service to signaling applications between adjacent peers.

Routing State Integrity Protection: It is important that signaling messages are delivered to the correct nodes, and nowhere else. Here, 'correct' is defined as 'the appropriate nodes for the signaling given the Message-Routing-Information'. In the case where the MRI is based on the Flow Identification for path-coupled signaling, 'appropriate' means 'the same nodes that the infrastructure will route data flow packets through'. GIST has no role in deciding whether the data flow itself is being routed correctly; all it can do is ensure the signaling is routed consistently with it. GIST uses internal state to decide how to route signaling messages, and this state needs to be protected against corruption.

Prevention of Denial of Service Attacks: GIST nodes and the network have finite resources (state storage, processing power, bandwidth). The protocol tries to minimise exhaustion attacks against these resources and not allow GIST nodes to be used to launch attacks on other network elements.

The main additional issue is handling authorisation for executing signaling operations (e.g. allocating resources). This is assumed to be done in each signaling application.

In many cases, GIST relies on the security mechanisms available in messaging associations to handle these issues, rather than introducing new security measures. Obviously, this requires the interaction of these mechanisms with the rest of the GIST protocol to be understood and verified, and some aspects of this are discussed in [Section 5.7](#).

8.1. Message Confidentiality and Integrity

GIST can use messaging association functionality, specifically in this version TLS ([Section 5.7.3](#)), to ensure message confidentiality

Internet-Draft

GIST

July 2006

and integrity. Implementation of this functionality is REQUIRED but its use for any given flow or signaling application is OPTIONAL. In some cases, confidentiality of GIST information itself is not likely to be a prime concern, in particular since messages are often sent to parties which are unknown ahead of time, although the content visible even at the GIST level gives significant opportunities for traffic analysis. Signaling applications may have their own mechanism for securing content as necessary; however, they may find it convenient to rely on protection provided by messaging associations, since it runs unbroken between signaling application peers.

[8.2.](#) Peer Node Authentication

Cryptographic protection (of confidentiality or integrity) requires a security association with session keys. These can be established by an authentication and key exchange protocol based on shared secrets, public key techniques or a combination of both. Authentication and key agreement is possible using the protocols associated with the messaging association being secured. TLS incorporates this functionality directly; IKE, IKEv2 or KINK could provide it for IPsec. GIST nodes rely on these protocols to authenticate the identity of the next hop, and GIST has no authentication capability of its own.

With routing state discovery, there are few effective ways to know what is the legitimate next or previous hop as opposed to an impostor. In other words, cryptographic authentication here only provides assurance that a node is 'who' it is (i.e. the legitimate owner of identity in some namespace), not 'what' it is (i.e. a node which is genuinely on the flow path and therefore can carry out signaling for a particular flow). Authentication provides only limited protection, in that a known peer is unlikely to lie about its role. Additional methods of protection against this type of attack are considered in [Section 8.3](#) below.

It is an implementation issue whether peer node authentication should be made signaling application dependent; for example, whether successful authentication could be made dependent on presenting credentials related to a particular signaling role (e.g. signaling for QoS). The abstract API of [Appendix B](#) leaves open such policy and authentication interactions between GIST and the NSLP it is serving. However, it does allow applications to inspect the authenticated

identity of the peer to which a message will be sent before transmission.

[8.3.](#) Routing State Integrity

Internal state in a node (see [Section 4.2](#)) is used to route messages.

If this state is corrupted, signaling messages may be misdirected.

In the case where the MRM is path-coupled, the messages need to be routed identically to the data flow described by the MRI, and the routing state table is the GIST view of how these flows are being routed through the network in the immediate neighbourhood of the node. Routes are only weakly secured (e.g. there is no cryptographic binding of a flow to a route), and there is no authoritative information about flow routes other than the current state of the network itself. Therefore, consistency between GIST and network routing state has to be ensured by directly interacting with the routing mechanisms to ensure that the signaling peers are the appropriate ones for any given flow. An overview of security issues and techniques in this context is provided in [\[38\]](#).

In one direction, peer identification is installed and refreshed only on receiving a Response (compare Figure 4). This MUST echo the cookie from a previous Query, which will have been sent along the flow path with the Q-mode encapsulation, i.e. end-to-end addressed. Hence, only the true next peer or an on-path attacker will be able to generate such a message, provided freshness of the cookie can be checked at the querying node.

In the other direction, peer identification MAY be installed directly on receiving a Query containing addressing information for the signaling source. However, any node in the network could generate such a message; indeed, many nodes in the network could be the genuine upstream peer for a given flow. To protect against this, three strategies are used:

Filtering: the receiving node MAY reject signaling messages which claim to be for flows with flow source addresses which could be ruled out by ingress filtering. An extension of this technique would be for the receiving node to monitor the data plane and to check explicitly that the flow packets are arriving over the same

interface and if possible from the same link layer neighbour as the D-mode signaling packets. If they are not, it is likely that at least one of the signaling or flow packets is being spoofed.

Authentication (weak or strong): the receiving node MAY refuse to install upstream state until it has completed a Confirm handshake with the peer. This echoes the Response cookie of the Response, and discourages nodes from using forged source addresses. This also plays a role in denial of service prevention, see below. A stronger approach is to require full peer authentication within the messaging association, the reasoning being that an authenticated peer can be trusted not to pretend that it is on path when it is not.

SID segregation: The routing state lookup for a given MRI and NSLPID MUST also take the SID into account. A malicious node can only overwrite existing routing state if it can guess the corresponding SID; it can insert state with random SID values, but generally this will not be used to route messages for which state has already been legitimately established.

[8.4.](#) Denial of Service Prevention

GIST is designed so that in general each Query only generates at most one Response, so that a GIST node cannot become the source of a denial of service amplification attack. (There is a special case of retransmitted Response messages, see [Section 5.3.3.](#))

However, GIST can still be subjected to denial-of-service attacks where an attacker using forged source addresses forces a node to establish state without return routability, causing a problem similar to TCP SYN flood attacks. Furthermore, an adversary might use modified or replayed unprotected signaling messages as part of such an attack. There are two types of state attacks and one computational resource attack. In the first state attack, an attacker floods a node with messages that the node has to store until it can determine the next hop. If the destination address is chosen so that there is no GIST-capable next hop, the node would accumulate messages for several seconds until the discovery retransmission attempt times out. The second type of state-based attack causes GIST state to be established by bogus messages. A related computational/network-resource attack uses unverified messages to cause a node

query an authentication or authorisation infrastructure, or attempt to cryptographically verify a digital signature.

We use a combination of two defences against these attacks:

1. The responding node need not establish a session or discover its next hop on receiving the Query, but MAY wait for a Confirm, possibly on a secure channel. If the channel exists, the additional delay is one one-way delay and the total is no more than the minimal theoretically possible delay of a three-way handshake, i.e., 1.5 node-to-node round-trip times. The delay gets significantly larger if a new connection needs to be established first.
2. The Response to the Query contains a cookie, which is repeated in the Confirm. State is only established for messages that contain a valid cookie. The setup delay is also 1.5 round-trip times. This mechanism is similar to that in SCTP [[16](#)] and other modern protocols.

Once a node has decided to establish routing state, there may still be transport and security state to be established between peers. This state setup is also vulnerable to denial of service attacks. GIST relies on the lower layer protocols that make up messaging associations to mitigate such attacks. In the current specification, the querying node is always the one wishing to establish a messaging association, so it is the responding node that needs to be protected.

Signaling applications can use the services provided by GIST to defend against certain (e.g. flooding) denial of service attacks. In particular, they can elect to process only messages from peers that have passed a return routability check or been authenticated at the messaging association level (see [Appendix B.2](#)). Signaling applications that accept messages under other circumstances (in particular, before routing state has been fully established at the GIST level) need to take this into account when designing their denial of service prevention mechanisms, for example by not creating local state as a result of processing such messages.

[8.5](#). Requirements on Cookie Mechanisms

The requirements on the Query cookie can be summarised as follows:

Liveness: The cookie must be live, that is, it must change from one handshake to the next. To prevent replay attacks.

Unpredictability: The cookie must not be guessable e.g. from a sequence or timestamp. To prevent direct forgery based on seeing a history of captured messages.

Easily validated: It must be efficient for the Q-Node to validate that a particular cookie matches an in-progress handshake, for a routing state machine which already exists. To discard spoofed responses, or responses to spoofed queries.

Uniqueness: The cookie must be unique to a given handshake since it is actually used to match the Response to a handshake anyway, e.g. during messaging association re-use.

Likewise, the requirements on the Responder cookie can be summarised as follows:

Liveness: The cookie must be live as above. To prevent replay attacks.

Creation simplicity: The cookie must be lightweight to generate. To avoid resource exhaustion at the responding node.

Validation simplicity: It must be simple for the R-node to validate that an R-cookie was generated by itself and no-one else, without storing state about the handshake it was generated for.

Binding: The cookie must be bound to the routing state that will be installed. To prevent use with different routing state e.g. in a modified Confirm. The routing state here includes the NLI of the Query, the MRI/NSLPID for the messaging, and the interface on which the Query was received.

A suitable implementation for the Q-Cookie is a cryptographically

strong random number which is unique for this routing state machine handshake. A node MUST implement this or an equivalently strong mechanism. Guidance on random number generation can be found in [28].

A suitable implementation for the R-Cookie is as follows:

R-Cookie = liveness data + hash (locally known secret,
Q-Node NLI, MRI, NSLPID,
reception interface,
liveness data)

A node MUST implement this or an equivalently strong mechanism. There are several alternatives for the liveness data. One is to use a timestamp like SCTP. Another is to give the local secret a (rapid) rollover, with the liveness data as the generation number of the secret, like IKEv2. In both cases, the liveness data has to be carried outside the hash, to allow the hash to be verified at the Responder. Another approach is to replace the hash with encryption under a locally known secret, in which case the liveness data does not need to be carried in the clear. Any symmetric cipher immune to known plaintext attacks can be used.

To support the validation simplicity requirement, the Responder can check the liveness data to filter out some blind (flooding) attacks before beginning any cryptographic cookie verification. To support this usage, the liveness data must be carried in the clear and not be easily guessable; this rules out the timestamp approach, and suggests the use of sequence of secrets with the liveness data identifying the position in the sequence. The secret strength and rollover frequency must be high enough that the secret cannot be brute-forced during its lifetime. Note that any node can use a Query to discover the current liveness data, so it remains hard to defend against sophisticated attacks which disguise such probes within a flood of Queries from

forged source addresses. Therefore, it remains important to use an efficient hashing mechanism or equivalent.

If a node receives a message for which cookie validation fails, it MAY return an "Object Value Error" error message (Appendix A.4.4.10) with subcode 4 ("Invalid Cookie") to the sender, as well as dropping the message. However, doing so in general makes a node a source of

backscatter. Therefore, this MUST only be enabled selectively, e.g. during initial deployment or debugging.

[8.6.](#) Security Protocol Selection Policy

This specification defines a single mandatory-to-implement security protocol (TLS, [Section 5.7.3](#)). However, it is possible to define additional security protocols in the future, for example to allow re-use with other types of credentials, or migrate towards protocols with stronger security properties. In addition, use of any security protocol for a messaging association is optional. Security protocol selection is carried out as part of the GIST handshake mechanism ([Section 4.4.1](#)).

The selection process may be vulnerable to downgrade attacks, where a man in the middle modifies the capabilities offered in the Query or Response to mislead the peers into accepting a lower level of protection than is achievable. There is a two part defence against such attacks (the following is based the same concepts as [\[22\]](#)):

1. The Response does not depend on the Stack-Proposal in the Query (see [Section 5.7.1](#)). Therefore, tampering with the Query has no effect on the resulting messaging association configuration.
2. The Responding node's Stack-Proposal is echoed in the Confirm. The Responding node checks this to validate that the proposal it made in the Response is the same as the one received by the Querying node. Note that as a consequence of the previous point, the Responding node does not have to remember the proposal explicitly, since it is a static function of local policy.

The validity of the second part depends on the strength of the security protection provided for the Confirm. If the Querying node is prepared to create messaging associations with null security properties (e.g. TCP only), the defence is ineffective, since the man in the middle can re-insert the original Responder's Stack-Proposal, and the Responding node will assume that the minimal protection is a consequence of Querying node limitations. However, if the messaging association provides at least integrity protection that cannot be broken in real-time, the Confirm cannot be modified in this way. Therefore, if the Querying node does not apply a security

policy to the messaging association protocols to be created that ensures at least this minimal level of protection is met, it remains open to the threat that a downgrade has occurred. Applying such a policy ensures capability discovery process will result in the setup of a messaging association with the correct security properties as appropriate for the two peers involved.

8.7. Residual Threats

Taking the above security mechanisms into account, the main residual threats against NSIS are three types of on-path attack.

An on-path attacker who can intercept the initial Query can do most things it wants to the subsequent signaling. It is very hard to protect against this at the GIST level; the only defence is to use strong messaging association security to see whether the Responding node is authorised to take part in NSLP signaling exchanges. To some extent, this behaviour is logically indistinguishable from correct operation, so it is easy to see why defence is difficult. Note that an on-path attacker of this sort can do anything to the traffic as well as the signaling. Therefore, the additional threat induced by the signaling weakness seems tolerable.

At the NSLP level, there is a concern about transitivity of trust of correctness of routing along the signaling chain. The NSLP at the querying node can have good assurance that it is communicating with an on-path peer or a node delegated by the on-path node. However, it has no assurance that the node beyond the responder is also on-path, or that the MRI (in particular) is not being modified by the responder to refer to a different flow. Therefore, if it sends signaling messages with payloads (e.g. authorisation tokens) which are valuable to nodes beyond the adjacent hop, it is up to the NSLP to ensure that the appropriate chain of trust exists, which must in general use strong messaging association security.

There is a further residual attack by a node which is not on the path of the Query, but is on the path of the Response, or is able to use a Response from one handshake to interfere with another. The attacker modifies the Response to cause the Querying node to form an adjacency with it rather than the true peer. In principle, this attack could be prevented by including an additional cryptographic object in the Response which ties the Response to the initial Query and the routing state and can be verified by the Querying node.

Internet-Draft

GIST

July 2006

9. IANA Considerations

This section defines the registries and initial codepoint assignments for GIST. It also defines the procedural requirements to be followed by IANA in allocating new codepoints. Note that the guidelines on the technical criteria to be followed in evaluating requests for new codepoint assignments are covered normatively in a separate document which considers the NSIS protocol suite in a unified way. That document discusses the general issue of NSIS extensibility, as well as the technical criteria for particular registries; see [[11](#)] for further details.

The registry definitions that follow leave large blocks of codes reserved as unused. This is to allow a future revision of this specification to modify the allocation policies without having to retrospectively change the initial rules if they turn out to have been suboptimal, e.g. if the space for one particular policy is exhausted too quickly.

The allocation policies used in this section follow the guidance given in [3]. In addition, for a number of the GIST registries, this specification also defines private/experimental ranges as discussed in [8]. Note that the only environment in which these codepoints can validly be used is a closed one in which the experimenter knows all the experiments in progress.

This specification allocates the following codepoints in existing registries:

Well-known UDP port XXX as the destination port for Q-mode encapsulated GIST messages ([Section 5.3](#)).

This specification creates the following registries with the structures as defined below:

NSLP Identifiers: Each signaling application requires the assignment of one of more NSLPIDs. The following NSLPID is allocated by this specification:

NSLPID	Application
0	Used for GIST messages not related to any signaling

| application. |
+-----+

Every other NSLPID MUST be associated with a specific RAO value; multiple NSLPIDs MAY be associated with the same value. The

NSLPID is a 16 bit integer, and allocation policies for further values are as follows:

1-32703: IESG Approval

32704-32767: Private/Experimental Use

32768-65536: Reserved

GIST Message Type: The GIST common header (Appendix A.1) contains a 1 byte message type field. The following values are allocated by this specification:

MType	Message
0	Query
1	Response
2	Confirm
3	Data
4	Error
5	MA-Hello

Allocation policies for further values are as follows:

6-63: Standards Action

64-119: Expert Review

120-127: Private/Experimental Use

128-255: Reserved

Object Types: There is a 12-bit field in the object header (Appendix A.2). The following values for object type are defined by this specification:

OType	Object Type
0	Message Routing Information
1	Session ID
2	Network Layer Information
3	Stack Proposal
4	Stack Configuration Data
5	Query Cookie
6	Responder Cookie
7	NAT Traversal
8	NSLP Data
9	Error

Allocation policies for further values are as follows:

10-1023: Standards Action

1024-1999: Specification Required

2000-2047: Private/Experimental Use

2048-4095: Reserved

When a new object type is defined, the object format MUST be provided, and the setting of the extensibility bits (A/B, see [Appendix A.2.1](#)) MUST also be defined.

Message Routing Methods: GIST allows multiple message routing methods (see [Section 3.3](#)). The MRM is indicated in the leading byte of the MRI object (Appendix A.3.1). This specification defines the following values:

+-----+-----+	
MRM-ID	Message Routing Method
+-----+-----+	
0	Path Coupled MRM
1	Loose End MRM
+-----+-----+	

Allocation policies for further values are as follows:

2-63: Standards Action

64-119: Expert Review

120-127: Private/Experimental Use

128-255: Reserved

When a new MRM is defined, the specification MUST provide the information described in [Section 3.3](#).

MA-Protocol-IDs: Each upper layer protocol that can be used in a

messaging association is identified by a 1-byte MA-Protocol-ID ([Section 5.7](#)). This is used as a tag in the Stack-Proposal and Stack-Configuration-Data objects (Appendix A.3.4 and [Appendix A.3.5](#)). The following values are defined by this specification:

MA-Protocol-ID	Higher Layer Protocol
1	TCP opened in the forwards direction
2	TLS initiated in the forwards direction

Allocation policies for further values are as follows:

3-63: Standards Action

64-119: Expert Review

120-127: Private/Experimental Use

128-255: Reserved

Allocation of a new MA-Protocol-ID MUST define the format for the MA-protocol-options field (if any) in the Stack-Configuration-Data

object that is needed to define its configuration. If a protocol is to be used for reliable message transfer, it MUST be described how delivery errors are to be detected by GIST. Note that the MA-Protocol-ID is not an IP Protocol number; indeed, some of the messaging association protocols - such as TLS - do not have an IP Protocol number.

Error Codes/Subcodes: There is a 2 byte error code and 1 byte subcode in the Value field of the Error object (Appendix A.4.1). Error codes 1-12 are defined in [Appendix A.4.4](#) together with subcodes 0-3 for code 1, 0-5 for code 9, 0-5 for code 10, and 0-2 for code 12. Additional codes and subcodes are allocated on a first-come, first served basis. When a new error code/subcode combination is allocated, the Error Class and the format of any associated error-specific information MUST also be defined.

[10](#). Acknowledgements

This document is based on the discussions within the IETF NSIS working group. It has been informed by prior work and formal and informal inputs from: Cedric Aoun, Attila Bader, Roland Bless, Bob Braden, Marcus Brunner, Benoit Campedel, Luis Cordeiro, Elwyn Davies, Christian Dickmann, Pasi Eronen, Alan Ford, Xiaoming Fu, Bo Gao, Ruediger Geib, Eleanor Hepworth, Thomas Herzog, Cheng Hong, Jia Jia,

Cornelia Kappler, Georgios Karagiannis, Ruud Klaver, Chris Lang, John Loughney, Allison Mankin, Jukka Manner, Pete McCann, Andrew McDonald, Glenn Morrow, Dave Oran, Andreas Pashalidis, Henning Peters, Tom Phelan, Takako Sanda, Charles Shen, Melinda Shore, Martin Stiernerling, Martijn Swanink, Mike Thomas, Hannes Tschofenig, Sven van den Bosch, Michael Welzl, Lars Westberg, and Mayi Zoumaro-djayoon. Parts of the TLS usage description ([Section 5.7.3](#)) were derived from the Diameter base protocol specification, [RFC3588](#). In addition, Hannes Tschofenig provided a detailed set of review comments on the security section, and Andrew McDonald provided the formal description for the initial packet formats. Chris Lang's implementation work provided objective feedback on the clarity and feasibility of the specification, and he also provided the state machine description and the initial error catalogue and formats. Finally, Magnus Westerlund carried out a detailed AD review which identified a number of issues and led to significant clarifications.

11.1. Normative References

- [1] Katz, D., "IP Router Alert Option", [RFC 2113](#), February 1997.
- [2] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [3] Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA Considerations Section in RFCs", [BCP 26](#), [RFC 2434](#), October 1998.
- [4] Nichols, K., Blake, S., Baker, F., and D. Black, "Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers", [RFC 2474](#), December 1998.
- [5] Partridge, C. and A. Jackson, "IPv6 Router Alert Option", [RFC 2711](#), October 1999.
- [6] Nordmark, E., "Stateless IP/ICMP Translation Algorithm (SIIT)", [RFC 2765](#), February 2000.
- [7] Dierks, T. and C. Allen, "The TLS Protocol Version 1.0", [RFC 2246](#), January 1999.
- [8] Narten, T., "Assigning Experimental and Testing Numbers Considered Useful", [BCP 82](#), [RFC 3692](#), January 2004.
- [9] Crocker, D., Ed. and P. Overell, "Augmented BNF for Syntax Specifications: ABNF", [RFC 4234](#), October 2005.
- [10] Dierks, T. and E. Rescorla, "The Transport Layer Security (TLS) Protocol Version 1.1", [RFC 4346](#), April 2006.
- [11] Loughney, J., "NSIS Extensibility Model", [draft-loughney-nsis-ext-02](#) (work in progress), March 2006.

11.2. Informative References

- [12] Braden, B., Zhang, L., Berson, S., Herzog, S., and S. Jamin, "Resource ReSerVation Protocol (RSVP) -- Version 1 Functional Specification", [RFC 2205](#), September 1997.
- [13] Moy, J., "OSPF Version 2", STD 54, [RFC 2328](#), April 1998.
- [14] Terzis, A., Krawczyk, J., Wroclawski, J., and L. Zhang, "RSVP Operation Over IP Tunnels", [RFC 2746](#), January 2000.

-
- [15] Tsirtsis, G. and P. Srisuresh, "Network Address Translation - Protocol Translation (NAT-PT)", [RFC 2766](#), February 2000.
 - [16] Stewart, R., Xie, Q., Morneault, K., Sharp, C., Schwarzbauer, H., Taylor, T., Rytina, I., Kalla, M., Zhang, L., and V. Paxson, "Stream Control Transmission Protocol", [RFC 2960](#), October 2000.
 - [17] Carpenter, B. and K. Moore, "Connection of IPv6 Domains via IPv4 Clouds", [RFC 3056](#), February 2001.
 - [18] Huitema, C., "An Anycast Prefix for 6to4 Relay Routers", [RFC 3068](#), June 2001.
 - [19] Baker, F., Iturralde, C., Le Faucheur, F., and B. Davie, "Aggregation of RSVP for IPv4 and IPv6 Reservations", [RFC 3175](#), September 2001.
 - [20] Grossman, D., "New Terminology and Clarifications for Diffserv", [RFC 3260](#), April 2002.
 - [21] Price, R., Bormann, C., Christoffersson, J., Hannu, H., Liu, Z., and J. Rosenberg, "Signaling Compression (SigComp)", [RFC 3320](#), January 2003.
 - [22] Arkko, J., Torvinen, V., Camarillo, G., Niemi, A., and T. Haukka, "Security Mechanism Agreement for the Session Initiation Protocol (SIP)", [RFC 3329](#), January 2003.
 - [23] Rosenberg, J., Weinberger, J., Huitema, C., and R. Mahy, "STUN - Simple Traversal of User Datagram Protocol (UDP) Through Network Address Translators (NATs)", [RFC 3489](#), March 2003.
 - [24] Rosenberg, J., "Obtaining Relay Addresses from Simple Traversal of UDP Through NAT (STUN)", [draft-ietf-behave-turn-01](#) (work in progress), June 2006.
 - [25] Gill, V., Heasley, J., and D. Meyer, "The Generalized TTL Security Mechanism (GTSM)", [RFC 3682](#), February 2004.
 - [26] Hancock, R., Karagiannis, G., Loughney, J., and S. Van den Bosch, "Next Steps in Signaling (NSIS): Framework", [RFC 4080](#), June 2005.
 - [27] Tschofenig, H. and D. Kroesenberg, "Security Threats for Next Steps in Signaling (NSIS)", [RFC 4081](#), June 2005.

- [28] Eastlake, D., Schiller, J., and S. Crocker, "Randomness

- Requirements for Security", [BCP 106](#), [RFC 4086](#), June 2005.
- [29] Eronen, P. and H. Tschofenig, "Pre-Shared Key Ciphersuites for Transport Layer Security (TLS)", [RFC 4279](#), December 2005.
- [30] Kohler, E., Handley, M., and S. Floyd, "Datagram Congestion Control Protocol (DCCP)", [RFC 4340](#), March 2006.
- [31] Conta, A., Deering, S., and M. Gupta, "Internet Control Message Protocol (ICMPv6) for the Internet Protocol Version 6 (IPv6) Specification", [RFC 4443](#), March 2006.
- [32] Stiemerling, M., "NAT/Firewall NSIS Signaling Layer Protocol (NSLP)", [draft-ietf-nsis-nslp-natfw-12](#) (work in progress), June 2006.
- [33] Manner, J., "NSLP for Quality-of-Service Signaling", [draft-ietf-nsis-qos-nslp-11](#) (work in progress), June 2006.
- [34] Nordmark, E. and R. Gilligan, "Basic Transition Mechanisms for IPv6 Hosts and Routers", [RFC 4213](#), October 2005.
- [35] Kent, S. and K. Seo, "Security Architecture for the Internet Protocol", [RFC 4301](#), December 2005.
- [36] Ylonen, T. and C. Lonvick, "The Secure Shell (SSH) Protocol Architecture", [RFC 4251](#), January 2006.
- [37] Moskowitz, R., "Host Identity Protocol", [draft-ietf-hip-base-06](#) (work in progress), June 2006.
- [38] Nikander, P., Arkko, J., Aura, T., Montenegro, G., and E. Nordmark, "Mobile IP Version 6 Route Optimization Security Design Background", [RFC 4225](#), December 2005.
- [39] Floyd, S. and V. Jacobson, "The Synchronisation of Periodic Routing Messages", SIGCOMM Symposium on Communications Architectures and Protocols pp. 33--44, September 1993.

- [40] Pashalidis, A. and H. Tschofenig, "GIST NAT Traversal", [draft-pashalidis-nsis-gimps-nattraversal-03](#) (work in progress), June 2006.
- [41] Rescorla, E. and N. Modadugu, "Datagram Transport Layer Security", [RFC 4347](#), April 2006.

Appendix A. Bit-Level Formats and Error Messages

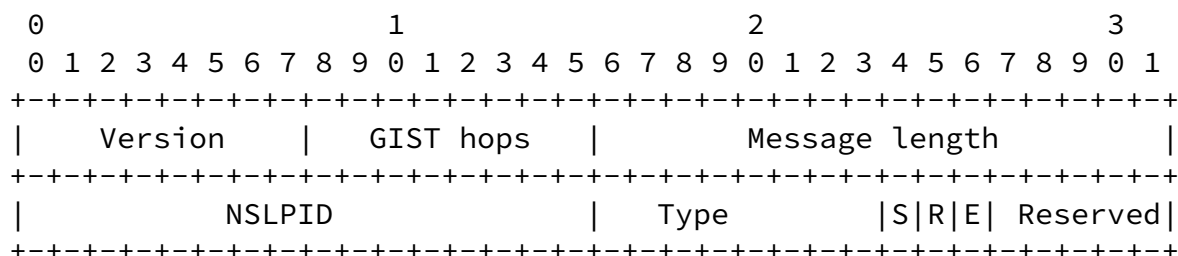
This appendix provides formats for the various component parts of the GIST messages defined abstractly in [Section 5.2](#).

Each GIST message consists of a header and a sequence of objects. The GIST header has a specific format, described in more detail in [Appendix A.1](#) below. An NSLP message is one object within a GIST message. Note that GIST itself provides the NSLP message length information and signaling application identification. General object formatting guidelines are provided in [Appendix A.2](#) below, followed in [Appendix A.3](#) by the format for each object. Finally, [Appendix A.4](#) provides the formats used for error reporting.

In the following object diagrams, '//' is used to indicate a variable sized field and ':' is used to indicate a field that is optionally present.

A.1. The GIST Common Header

This header begins all GIST messages. It has a fixed format, as shown below.



Version (8 bits): The GIST protocol version number.

GIST hops (8 bits): A hop count for the number of GIST-aware nodes this message can still pass through.

Message Length (16 bits): The total number of words in the message after the common header itself.

NSLPID (16 bits): IANA assigned identifier of the signaling application the message refers to.

Type (8 bits): The GIST message type (Query, Response, etc.).

S flag: S=1 if the IP source address is the same as the signaling source address, S=0 if it is different.

R flag: R=1 if a reply to this message is explicitly requested.

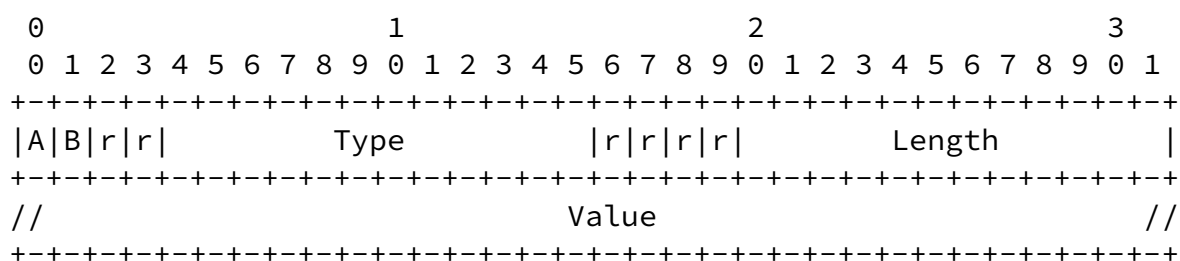
E flag: E=1 if the message was explicitly routed [Section 7.1.4](#).

The rules governing the use of the R-flag depend on the GIST message type. It MUST always be set (R=1) in Query messages, since these always elicit a Response, and never in Confirm, Data or Error messages. It is optional in an MA-Hello; if set, another MA-Hello is sent in reply. It is optional in a Response, but MUST be set if the Response contains a Responder cookie; if set, a Confirm is sent in reply.

Parsing failures may be caused by unknown Version or Type values, inconsistent R flag setting, or a Message Length inconsistent with the set of objects carried. In all cases the receiver MUST if possible return a "Common Header Parse Error" message (Appendix A.4.4.1) with the appropriate subcode, and not process the message further.

[A.2](#). General Object Format

Each object begins with a fixed header giving the object Type and object Length. This is followed by the object Value, which is a whole number of 32-bit words long.



A/B flags: The bits marked 'A' and 'B' are extensibility flags which are defined in [Appendix A.2.1](#) below; the remaining bits marked 'r' are reserved.

Type (12 bits): An IANA-assigned identifier for the type of object.

Length (12 bits): Length has the units of 32 bit words, and measures the length of Value. If there is no Value, Length=0. If the Length is not consistent with the contents of the object, an "Object Value Error" message (Appendix A.4.4.10) with subcode 0 "Incorrect Length" MUST be returned and the message dropped.

Value (variable): Value is (therefore) a whole number of 32 bit words. If there is any padding required, the length and location must be defined by the object-specific format information; objects which contain variable length (e.g. string) types may need to include additional length subfields to do so.

Any part of the object used for padding or defined as reserved (marked 'Reserved' or 'Rsv' or, in the case of individual bits, 'r' in the diagrams below) MUST be set to 0 on transmission and MUST be ignored on reception.

[A.2.1](#). Object Extensibility

The leading two bits of the TLV header are used to signal the desired treatment for objects whose Type field is unknown at the receiver. The following three categories of object have been identified, and are described here.

AB=00 ("Mandatory"): If the object is not understood, the entire message containing it MUST be rejected with an "Object Type Error" message (Appendix A.4.4.9) with subcode 1 ("Unrecognised Object").

AB=01 ("Ignore"): If the object is not understood, it MUST be deleted and the rest of the message processed as usual.

AB=10 ("Forward"): If the object is not understood, it MUST be retained unchanged in any message forwarded as a result of message processing, but not stored locally.

The combination AB=11 is reserved. If a message is received containing an object with AB=11, it MUST be rejected with an "Object Type Error" message (Appendix A.4.4.9) with subcode 5 ("Invalid Extensibility Flags").

[A.3.](#) GIST TLV Objects

[A.3.1.](#) Message-Routing-Information

Type: Message-Routing-Information

Length: Variable (depends on MRM)

```

      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|      MRM-ID      |N|  Reserved  |                               |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//      Method-specific addressing information (variable)      //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

MRM-ID (8 bits): An IANA-assigned identifier for the message routing method.

N flag: If set (N=1), this means that NATs do not need to translate this MRM; if clear (N=0) it means that the method-specific information contains network or transport layer information that a NAT must process.

The remainder of the object contains method-specific addressing information, which is described below.

[A.3.1.1.](#) Path-Coupled MRM

In the case of basic path-coupled routing, the addressing information takes the following format. The N-flag N=0 for this MRM.

```

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
                                     +---+---+---+---+---+---+---+---+
                                     |IP-Ver |P|T|F|S|A|B|D|Reserved |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                               Source Address                               //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                               Destination Address                           //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Source Prefix | Dest Prefix | Protocol   | DS-field |Rsv|
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:      Reserved      |      Flow Label      :
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:                               SPI                               :
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:      Source Port      :      Destination Port      :
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

IP-Ver (4 bits): The IP version number, 4 or 6.

Source/Destination address (variable): The source and destination addresses are always present and of the same type; their length depends on the value in the IP-Ver field.

Source/Dest Prefix (each 8 bits): The length of the mask to be applied to the source and destination addresses for address wildcarding. In the normal case where the MRI refers only to traffic between specific host addresses, the Source/Dest Prefix

values would both be 32/128 for IPv4/6 respectively.

P flag: P=1 means that IP Protocol should be interpreted.

Protocol (8 bits): The IP protocol number. Ignored if P=0. In the case of IPv6, the Protocol field refers to the true upper layer protocol carried by the packets, i.e. excluding any IP option headers. This is therefore not necessarily the same as the Next Header value from the base IPv6 header.

T flag: T=1 means that DiffServ field (DS-field) should be interpreted.

DS-field (6 bits): The DiffServ field. See [\[4\]](#) and [\[20\]](#).

F flag: F=1 means that flow label is present and should be interpreted.

Flow Label (20 bits): The flow label; only present if F=1. If F=0, the entire 32 bit word containing the Flow Label is absent. F may only be set if IP-Ver is 6.

S flag: S=1 means that the SPI field is present. Can only be set if P=1.

SPI field (32 bits): The SPI field; see [\[35\]](#). Only present if S=1.

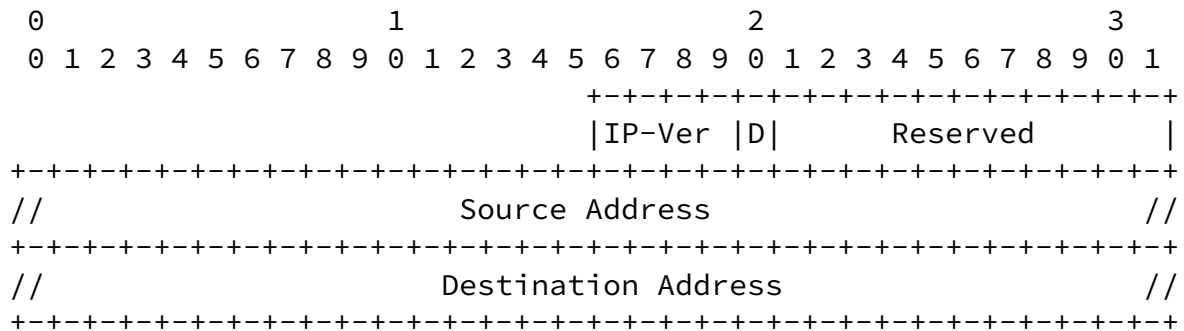
A/B flags: These can only be set if P=1. If either is set, the port fields are also present.

Source/Destination Port (each 16 bits): If either of A, B is set the word containing the port numbers is included in the object. However, the contents of each field is only significant if the corresponding flag is set; otherwise, the contents of the field is regarded as padding, and the MRI refers to all ports (i.e. acts as a wildcard). If the flag is set and Port=0x0000, the MRI will apply to a specific port, whose value is not yet known. If neither of A or B is set, the word is absent.

D flag: The Direction flag has the following meaning: the value 0 means 'in the same direction as the flow' (i.e. downstream), and the value 1 means 'in the opposite direction to the flow' (i.e. upstream).

A.3.1.2. Loose-End MRM

In the case of the loose-end MRM, the addressing information takes the following format. The N-flag N=0 for this MRM.



IP-Ver (4 bits): The IP version number, 4 or 6.

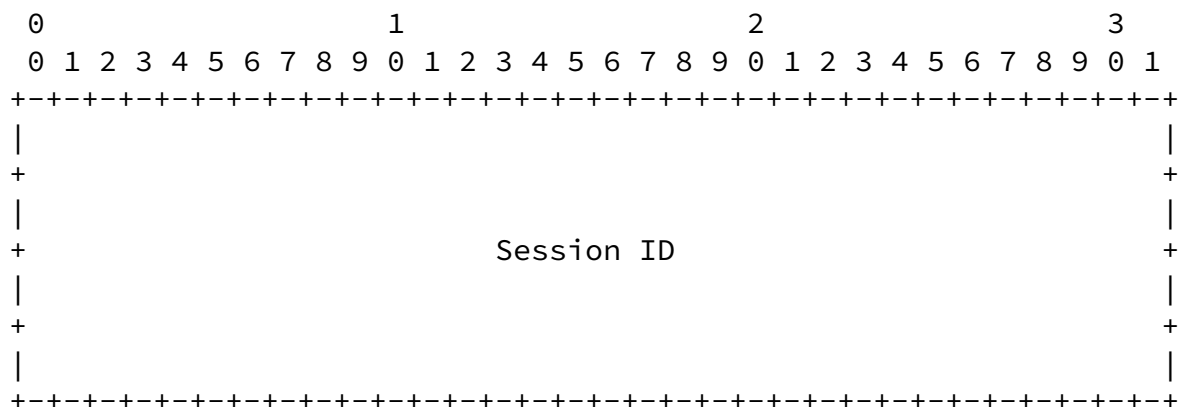
Source/Destination address (variable): The source and destination addresses are always present and of the same type; their length depends on the value in the IP-Ver field.

D flag: The direction flag. The only valid value is D=0 (see [Section 5.8.2](#)).

A.3.2. Session Identification

Type: Session-Identification

Length: Fixed (4 32-bit words)



A.3.3. Network-Layer-Information

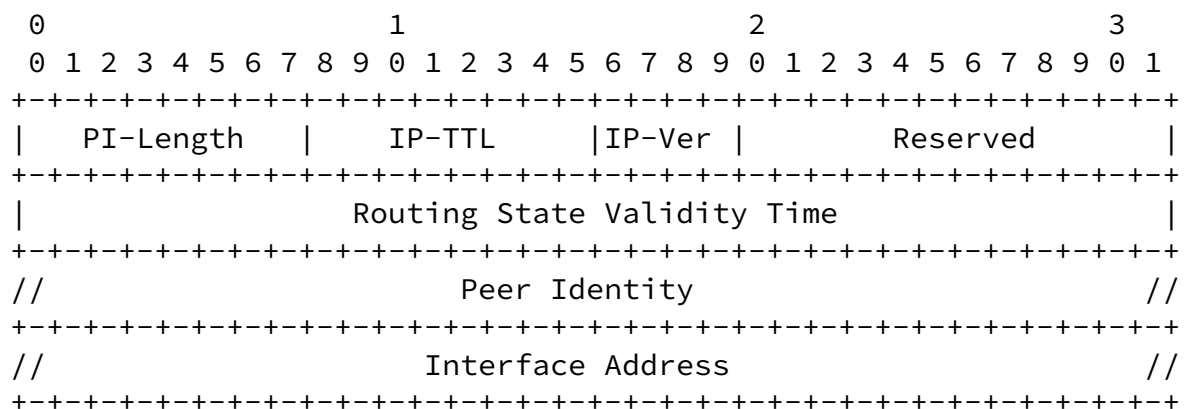
Internet-Draft

GIST

July 2006

Type: Network-Layer-Information

Length: Variable (depends on length of Peer-Identity and IP version)



PI-Length (8 bits): The byte length of the Peer Identity field.

Peer Identity (variable): The Peer Identity field. Note that the Peer-Identity field itself is padded to a whole number of words.

IP-TTL (8 bits): Initial or reported IP layer TTL.

IP-Ver (4 bits): The IP version for the Interface Address field.

Interface Address (variable): The IP address allocated to the interface, matching the IP-Ver field.

Routing State Validity Time (32 bits): The time for which the routing state for this flow can be considered correct without a refresh. Given in milliseconds.

[A.3.4.](#) Stack Proposal

Type: Stack-Proposal

Length: Variable (depends on number of profiles and size of each profile)

Internet-Draft

GIST

July 2006

```

      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Prof-Count   |      Reserved      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//              Profile 1              //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:                                                         :
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//              Profile 2              //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

Prof-Count (8 bits): The number of profiles listed. MUST be > 0.

Each profile is itself a sequence of protocol layers, and the profile is formatted as a list as follows:

- o The first byte is a count of the number of layers in the profile.
- o This is followed by a sequence of 1-byte MA-Protocol-IDs as described in [Section 5.7](#).
- o The profile is padded to a word boundary with 0, 1, 2 or 3 zero bytes.

If there are no profiles (i.e. all bytes are null), then an "Object Value Error" message (Appendix A.4.4.10) with subcode 3 ("Empty List") MUST be returned and the message dropped.

[A.3.5](#). Stack-Configuration-Data

Type: Stack-Configuration-Data

Length: Variable (depends on number of protocols and size of each MA-protocol-options field)

Internet-Draft

GIST

July 2006

```

      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|  MPO-Count   |      Reserved      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
|                                     MA-Hold-Time                                     |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                     MA-protocol-options 1                                     //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:                                                                                       :
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                     MA-protocol-options N                                     //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

MPO-Count (8 bits): The number of MA-protocol-options fields present (these contain their own length information).

MA-Hold-Time (32 bits): The time for which the messaging association will be held open without traffic or a hello message. Given in milliseconds.

The MA-protocol-options fields are formatted as follows:

```

      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```



```

|MA-Protocol-ID |      Profile      |      Length      |D| Reserved  |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                         Options Data                                         //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

MA-Protocol-ID (8 bits): Protocol identifier as described in [Section 5.7](#).

Profile (8 bits): Tag indicating which profile from the accompanying Stack-Proposal object this applies to. Profiles are numbered from 1 upwards; the special value 0 indicates 'applies to all profiles'.

Length (8 bits): The byte length of MA-protocol-options field that follows. This will be zero-padded up to the next word boundary.

D flag: If set (D=1), this protocol MUST NOT be used for a messaging association.

Options Data (variable): Any options data for this protocol. Note that the format of the options data may differ depending on whether the field is in a Query or Response.

[A.3.6](#). Query Cookie

Type: Query-Cookie

Length: Variable (selected by querying node)

```

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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                         Query Cookie                                         //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

The contents are implementation defined. See [Section 8.5](#) for further discussion.

[A.3.7.](#) Responder Cookie

Type: Responder-Cookie

Length: Variable (selected by responding node)

```

      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                               Responder Cookie                               //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
```

The contents are implementation defined. See [Section 8.5](#) for further discussion.

[A.3.8.](#) NAT Traversal

Type: NAT-Traversal

Length: Variable (depends on length of contained fields)

This object is used to support the NAT traversal mechanisms described in [Section 7.2](#).

```

      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
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| MRI-Length      | Type-Count      | NAT-Count      | Reserved      |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                Original Message-Routing-Information                //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                List of translated objects                          //
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Length of opaque information |                                         |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                Information replaced by NAT #1                      |
```

```

+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
:                                                                                   :
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
| Length of opaque information |                                                    |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+
//                                                                                   //
//                                     Information replaced by NAT #N                 |
+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+---+

```

MRI-Length (8 bits): The word length of the included MRI payload.

Original Message-Routing-Information (variable): The MRI data from when the message was first sent, not including the object header.

Type-Count (8 bits): The number of objects in the 'List of translated objects' field.

List of translated objects (variable): This field lists the object types of the objects that were translated by every NAT through which the message has passed. It is initialised by the first NAT on the path; subsequent NATs may delete elements in the list. Padded with 2 null bytes if necessary.

NAT-Count (8 bits): The number of NATs traversed by the message, and the number of opaque payloads at the end of the object. The length fields for each opaque payload are byte counts, not including the 2 bytes of the length field itself. Note that each opaque information field is zero-padded to the next 32-bit word boundary if necessary.

[A.3.9.](#) NSLP Data

Type: NSLP-Data

Length: Variable (depends on NSLP)

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 2 3 4 5 6 7 8 9	0 1		

```

+-----+
//                               NSLP Data                               //
+-----+

```

A.4. Errors

A.4.1. Error Object

Type: Error

Length: Variable (depends on error)

```

      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
+-----+-----+-----+-----+-----+-----+-----+-----+
| Error Class |           Error Code           | Error Subcode |
+-----+-----+-----+-----+-----+-----+-----+-----+
|S|M|C|D|Q|      Reserved      | MRI Length  | Info Count  |
+-----+-----+-----+-----+-----+-----+-----+-----+
|
+                               +
|                               |
+                               +
|                               |
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Session Id                               :
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Message Routing Information              :
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Additional Information                   :
+-----+-----+-----+-----+-----+-----+-----+-----+
:                               Debugging Comment                       :
+-----+-----+-----+-----+-----+-----+-----+-----+

```

The flags are:

S - S=1 means the Session ID object is present

M - M=1 means MRI object is present

C - C=1 means a debug Comment is present after header.

D - D=1 means the original message was received in D-mode

Q - Q=1 means the original message was received Q-mode encapsulated
(can't be set if D=0).

A GIST Error object contains an 8 bit error-class (see

[Appendix A.4.3](#)), a 16 bit error-code, an 8 bit error-subcode, and as much information about the message which triggered the error as is available. This information MUST include the Common header of the original message and MUST also include the Session Id and MRI objects if these could be decoded correctly. These objects are included in their entirety, except for their TLV Headers.

The Info Count field contains the number of Additional Information fields in the object, and the possible formats for these fields are given in [Appendix A.4.2](#). The precise set of fields to include depends on the error code/subcode. For every error description in the error catalogue [Appendix A.4.4](#), the line "Additional Info:" states what fields MUST be included, and in what order if there can be more than one; if this line is given as 'None' then additional information MUST NOT be included. The Debugging Comment is a null-terminated UTF-8 string, padded if necessary to a whole number of 32-bit words with more null characters.

[A.4.2](#). Additional Information Fields

The Common Error Header may be followed by some Additional Information objects. The possible formats of these objects are shown below.

Message Length Info:

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	2	3	4	5	6	7	8	9																								
Calculated Length																Reserved																																															

Calculated Length (16 bits): the length of the original message calculated by adding up all the objects in the message.

MTU Info:

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	2	3	4	5	6	7	8	9																								
Link MTU																Reserved																																															

Link MTU (16 bits): the MTU for a link along which a message could not be sent.

Internet-Draft

GIST

July 2006

Object Type Info:

```

      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
+-----+-----+-----+-----+-----+-----+-----+-----+
|               Object Type               |               Reserved               |
+-----+-----+-----+-----+-----+-----+-----+-----+

```

Object type (16 bits): This provides information about the type of object which caused the error.

Object Value Info:

```

      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
+-----+-----+-----+-----+-----+-----+-----+-----+
| Rsv | Real Object Length |               Offset               |
+-----+-----+-----+-----+-----+-----+-----+-----+
//                               Object                               //
+-----+-----+-----+-----+-----+-----+-----+-----+

```

This object carries information about a TLV object which was found to be invalid in the original message. An error message may contain more than one Object Value Info object.

Real Object Length (12 bits) Since the length in the original TLV header may be inaccurate, this field provides the actual length of the object (including the TLV Header) included in the error message.

Offset (16 bits): The byte in the object at which the GIST node found the error. The first byte in the object has offset=0.

Object (variable): The invalid TLV object (including the TLV Header).

A.4.3. Error Classes

The first byte of the error object, "Error Class", indicates the severity level. The currently defined severity levels are:

0 (Informational): response data which should not be thought of as changing the condition of the protocol state machine.

- 1 (Success): response data which indicates that the message being responded to has been processed successfully in some sense.

- 2 (Protocol-Error): the message has been rejected because of a protocol error (e.g. an error in message format).
- 3 (Transient-Failure): the message has been rejected because of a particular local node status which may be transient (i.e. it may be worthwhile to retry after some delay).
- 4 (Permanent-Failure): the message has been rejected because of local node status which will not change without additional out of band (e.g. management) operations.

Additional error class values are reserved.

The allocation of error classes to particular errors is not precise; the above descriptions are deliberately informal. Actual error processing should take into account the specific error in question; the error class may be useful supporting information (e.g. in network debugging).

[A.4.4.](#) Error Catalogue

This section lists all the possible GIST errors, including when they are raised and what additional information fields should be carried in the error object.

[A.4.4.1.](#) Common Header Parse Error

Class:	Protocol-Error
Code:	1
Additional Info:	For subcode 3 only, Message Length Info carries the calculated message length.

This message is sent if a GIST node receives a message where the common header cannot be parsed correctly, or where an error in the overall message format is detected. Note that in this case the

original MRI and Session ID are not included in the Error Object.
This error code is split into subcodes as follows:

- 0: Unknown Version: The GIST version is unknown. The (highest) supported version supported by the node can be inferred from the Common Header of the Error message itself.
- 1: Unknown Type: The GIST message type is unknown.
- 2: Invalid R-flag: The R flag in the header is inconsistent with the message type.

- 3: Incorrect Message Length: The overall message length is not consistent with the set of objects carried.

[A.4.4.2.](#) Hop Limit Exceeded

Class: Permanent-Failure
Code: 2
Additional Info: None

This message is sent if a GIST node receives a message with a GIST hop count of zero, or a GIST node decrements a packet's GIST hop count to zero on reception. This message indicates either a routing loop or too small an initial hop count value.

[A.4.4.3.](#) Incorrect Encapsulation

Class: Protocol-Error
Code: 3
Additional Info: None

This message is sent if a GIST node receives a message which uses an incorrect encapsulation method (e.g. a Query arrives over an MA).

[A.4.4.4.](#) Incorrectly Delivered Message

Class: Protocol-Error
Code: 4
Additional Info: None

This message is sent if a GIST node receives a message over an MA which is not associated with the MRI/NSLPID/SID combination in the message.

[A.4.4.5.](#) No Routing State

Class: Protocol-Error
Code: 5
Additional Info: None

This message is sent if a node receives a message for which routing state should exist, but has not yet been created and thus there is no appropriate Querying-SM or Responding-SM. This can occur either on receiving a Response to an unknown Query, or on receiving a Data message at a node whose policy requires routing state to exist before such messages can be accepted. See also [Section 6.1](#) and [Section 6.3](#).

[A.4.4.6.](#) Unknown NSLPID

Class: Permanent-Failure
Code: 6
Additional Info: None

This message is sent if a router receives a directly addressed message for an NSLP which it does not support.

[A.4.4.7.](#) Endpoint Found

Class: Informational
Code: 7
Additional Info: None

This message is sent if a GIST node at a flow endpoint receives a Query message for an NSLP which it does not support.

[A.4.4.8.](#) Message Too Large

Class: Permanent-Failure

Code: 8
Additional Info: MTU Info

A router receives a message which it can't forward because it exceeds the MTU on the next or subsequent hops.

[A.4.4.9.](#) Object Type Error

Class: Protocol-Error
Code: 9
Additional Info: Object Type Info

This message is sent if a GIST node receives a message containing a TLV object with an invalid type. The message includes the object at fault. This error code is split into subcodes as follows:

- 0: Duplicate Object: This subcode is used if a GIST node receives a message containing multiple instances of an object which may only appear once in a message. In the current specification, this applies to all objects.
- 1: Unrecognised Object: This subcode is used if a GIST node receives a message containing an object which it does not support, and the extensibility flags AB=00.

- 2: Missing Object: This subcode is used if a GIST node receives a message which is missing one or more mandatory objects. This message is also sent if a Stack-Proposal is sent without a matching Stack-Configuration-Data object when one was necessary, or vice versa.
- 3: Invalid Object Type: This subcode is used if the object type is known, but it is not valid for this particular GIST message type.
- 4: Untranslated Object: This subcode is used if the object type is known and is mandatory to interpret, but it contains addressing data which has not been translated by an intervening NAT.
- 5: Invalid Extensibility Flags: This subcode is used if an object is

received with the extensibility flags AB=11.

A.4.4.10. Object Value Error

Class: Protocol-Error
Code: 10
Additional Info: 1 or 2 Object Value Info fields as given below

This message is sent if a router receives a packet containing an object which cannot be properly parsed. The message contains a single Object Value Info object, except for subcode 5 as stated below. This error code is split into subcodes as follows:

- 0: Incorrect Length: The overall length does not match the object length calculated from the object contents.
- 1: Value Not Supported: The value of a field is not supported by the GIST node.
- 2: Invalid Flag-Field Combination: An object contains an invalid combination of flags and/or fields. At the moment this only relates to the Path-Coupled MRM object, but in future there may be more.
- 3: Empty List: At the moment this only relates to Stack-Proposals. The error message is sent if a stack proposal with a length > 0 (a length of 0 is handled as "Value Not Supported") contains only null bytes.
- 4: Invalid Cookie: The message contains a cookie which could not be verified by the node.

- 5: Stack-Proposal - Stack-Configuration-Data Mismatch: This subcode is used if a GIST node receives a message in which the data in the Stack-Proposal object is inconsistent with the information in the Stack Configuration Data object. In this case, both the Stack-Proposal object and Stack-Configuration-Data object MUST be included in separate Object Value Info fields in that order.

[A.4.4.11.](#) Invalid IP layer TTL

Class: Permanent-Failure
Code: 11
Additional Info: None

This error indicates that a message was received with an IP layer TTL outside an acceptable range; for example, that an upstream Query was received with an IP layer TTL of less than 254 (i.e. more than one IP hop from the sender). The actual IP distance can be derived from the IP-TTL information in the NLI object carried in the same message.

[A.4.4.12.](#) MRI Validation Failure

Class: Permanent-Failure
Code: 12
Additional Info: Object Value Info

This error indicates that a message was received with an MRI that could not be accepted, e.g. because of too much wildcarding or failing some validation check (cf. [Section 5.8.1.2](#)). The Object Value Info includes the MRI so the error originator can indicate the part of the MRI which caused the problem. The error code is divided into subcodes as follows:

- 0: MRI Too Wild: The MRI contained too much wildcarding (e.g. too short a destination address prefix) to be forwarded correctly down a single path.
- 1: IP Version Mismatch: The MRI in a path-coupled Query message uses an IP version which is not implemented on the interface used.
- 2: Ingress Filter Failure: The MRI in a path-coupled Query message describes a flow which would not pass ingress filtering on the interface used.

[Appendix B.](#) API between GIST and Signaling Applications

This appendix provides an abstract API between GIST and signaling applications. It should not constrain implementors, but rather help clarify the interface between the different layers of the NSIS protocol suite. In addition, although some of the data types carry the information from GIST information elements, this does not imply that the format of that data as sent over the API has to be the same.

Conceptually the API has similarities to the sockets API, particularly that for unconnected UDP sockets. An extension for an API like that for UDP connected sockets could be considered. In this case, for example, the only information needed in a `SendMessage` primitive would be `NSLP-Data`, `NSLP-Data-Size`, and `NSLP-Message-Handle` (which can be null). Other information which was persistent for a group of messages could be configured once for the socket. Such extensions may make a concrete implementation more efficient but do not change the API semantics, and so are not considered further here.

[B.1.](#) `SendMessage`

This primitive is passed from a signaling application to GIST. It is used whenever the signaling application wants to initiate sending a message.

```
SendMessage ( NSLP-Data, NSLP-Data-Size, NSLP-Message-Handle,  
              NSLPID, Session-ID, MRI, SII-Handle,  
              Transfer-Attributes, Timeout, IP-TTL, GIST-Hop-Count )
```

The following arguments are mandatory.

`NSLP-Data`: The NSLP message itself.

`NSLP-Data-Size`: The length of `NSLP-Data`.

`NSLP-Message-Handle`: A handle for this message, that can be used by GIST as a reference in subsequent `MessageStatus` notifications (Appendix B.3). Notifications could be about error conditions or about the security attributes that will be used for the message. A NULL handle may be supplied if the NSLP is not interested in such notifications.

`NSLPID`: An identifier indicating which NSLP this is.

`Session-ID`: The NSIS session identifier. Note that it is assumed that the signaling application provides this to GIST rather than GIST providing a value itself.

Internet-Draft

GIST

July 2006

MRI: Message routing information for use by GIST in determining the correct next GIST hop for this message. The MRI implies the message routing method to be used and the message direction.

The following arguments are optional:

SII-Handle: A handle, previously supplied by GIST, to a data structure that should be used to route the message explicitly to a particular GIST next hop.

Transfer-Attributes: Attributes defining how the message should be handled (see [Section 4.1.2](#)). The following attributes can be considered:

Reliability: Values 'unreliable' or 'reliable'.

Security: This attribute allows the NSLP to specify what level of security protection is requested for the message (selected from 'integrity' and 'confidentiality'), and can also be used to specify what authenticated signaling source and destination identities should be used to send the message. The possibilities can be learned by the signaling application from prior MessageStatus or RecvMessage notifications. If an NSLP-Message-Handle is provided, GIST will inform the signaling application of what values it has actually chosen for this attribute via a MessageStatus callback. This might take place either synchronously (where GIST is selecting from available messaging associations), or asynchronously (when a new messaging association needs to be created).

Local Processing: This attribute contains hints from the signaling application about what local policy should be applied to the message; in particular, its transmission priority relative to other messages, or whether GIST should attempt to set up or maintain forward routing state.

Timeout: Length of time GIST should attempt to send this message before indicating an error.

IP-TTL: The value of the IP layer TTL that should be used when sending this message (may be overridden by GIST for particular messages).

GIST-Hop-Count: The value for the hop count when sending the message.

[B.2.](#) RecvMessage

This primitive is passed from GIST to a signaling application. It is used whenever GIST receives a message from the network, including the case of null messages (zero length NSLP payload), typically initial Query messages. For Queries, the results of invoking this primitive are used by GIST to check whether message routing state should be created (see the discussion of the 'Routing-State-Check' argument below).

```
RecvMessage ( NSLP-Data, NSLP-Data-Size, NSLPID, Session-ID, MRI,  
              Routing-State-Check, SII-Handle, Transfer-Attributes,  
              IP-TTL, IP-Distance, GIST-Hop-Count, Inbound-Interface )
```

NSLP-Data: The NSLP message itself (may be empty).

NSLP-Data-Size: The length of NSLP-Data (may be zero).

NSLPID: An identifier indicating which NSLP this is message is for.

Session-ID: The NSIS session identifier.

MRI: Message routing information that was used by GIST in forwarding this message. Implicitly defines the message routing method that was used and the direction of the message relative to the MRI.

Routing-State-Check: This boolean is True if GIST is checking with the signaling application to see if routing state should be created with the peer or the message should be forwarded further (see [Section 4.3.2](#)). If True, the signaling application should return the following values via the RecvMessage call:

- A boolean indicating whether to set up the state.

- Optionally, an NSLP-Payload to carry in the generated Response or forwarded Query respectively.

This mechanism could be extended to enable the signaling application to indicate to GIST whether state installation should be immediate or deferred (see [Section 5.3.3](#) and [Section 6.3](#) for further discussion).

SII-Handle: A handle to a data structure, identifying a peer address and interface. Can be used to identify route changes and for explicit routing to a particular GIST next hop.

Transfer-Attributes: The reliability and security attributes that were associated with the reception of this particular message. As well as the attributes associated with SendMessage, GIST may indicate the level of verification of the addresses in the MRI. Three attributes can be indicated:

- * Whether the signaling source address is one of the flow endpoints (i.e. whether this is the first or last GIST hop);
- * Whether the signaling source address has been validated by a return routability check.
- * Whether the message was explicitly routed (and so has not been validated by GIST as delivered consistently with local routing state).

IP-TTL: The value of the IP layer TTL this message was received with (if available).

IP-Distance: The number of IP hops from the peer signaling node which sent this message along the path, or 0 if this information is not available.

GIST-Hop-Count: The value of the hop count the message was received with, after being decremented in the GIST receive-side processing.

Inbound-Interface: Attributes of the interface on which the message was received, such as whether it lies on the internal or external side of a NAT. These attributes have only local significance and

are implementation defined.

[B.3.](#) MessageStatus

This primitive is passed from GIST to a signaling application. It is used to notify the signaling application that a message that it requested to be sent could not be dispatched, or to inform the signaling application about the transfer attributes that have been selected for the message (specifically, security attributes). The signaling application can respond to this message with a return code to abort the sending of the message if the attributes are not acceptable.

MessageStatus (NSLP-Message-Handle, Transfer-Attributes, Error-Type)

NSLP-Message-Handle: A handle for the message provided by the signaling application in SendMessage.

Transfer-Attributes: The reliability and security attributes that will be used to transmit this particular message.

Error-Type: Indicates the type of error that occurred. For example, 'no next node found'.

[B.4.](#) NetworkNotification

This primitive is passed from GIST to a signaling application. It indicates that a network event of possible interest to the signaling application occurred.

NetworkNotification (NSLPID, MRI, Network-Notification-Type)

NSLPID: An identifier indicating which NSLP this is message is for.

MRI: Provides the message routing information to which the network notification applies.

Network-Notification-Type: Indicates the type of event that caused the notification and associated additional data. Two events have been identified:

Last Node: GIST has detected that this is the last NSLP-aware node in the path. See [Section 4.3.4](#).

Routing Status Change: GIST has installed new routing state, or has detected that the routing state may no longer be valid, or has re-established the routing state. See [Section 7.1.3](#). The new status is reported; if the status is Good, the SII-Handle of the peer is also reported, as for RecvMessage.

[B.5](#). SetStateLifetime

This primitive is passed from a signaling application to GIST. It indicates the duration for which the signaling application would like GIST to retain its routing state. It can also give a hint that the signaling application is no longer interested in the state.

SetStateLifetime (NSLPID, MRI, State-Lifetime)

NSLPID: Provides the NSLPID to which the routing state lifetime applies.

MRI: Provides the message routing information to which the routing state lifetime applies; includes the direction (in the D flag).

State-Lifetime: Indicates the lifetime for which the signaling application wishes GIST to retain its routing state (may be zero, indicating that the signaling application has no further interest in the GIST state).

[B.6](#). InvalidateRoutingState

This primitive is passed from a signaling application to GIST. It indicates that the signaling application has knowledge that the next

signaling hop known to GIST may no longer be valid, either because of changes in the network routing or the processing capabilities of signaling application nodes. See [Section 7.1](#).

InvalidateRoutingState (NSLPID, MRI, Status, Urgent)

NSLPID: The NSLP originating the message. May be null (in which case the invalidation applies to all signaling applications).

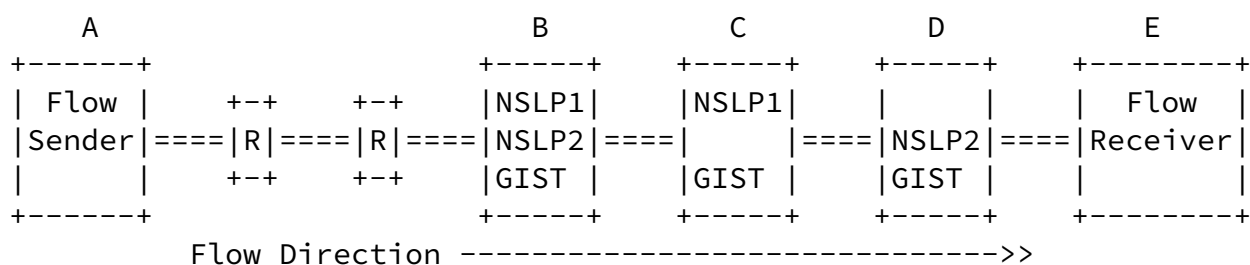
MRI: The flow for which routing state should be invalidated; includes the direction of the change (in the D flag).

Status: The new status that should be assumed for the routing state, one of Bad or Tentative (see [Section 7.1.3](#)).

Urgent: A hint as to whether rediscovery should take place immediately, or only with the next signaling message.

[Appendix C](#). Example Routing State Table and Handshake Message Sequence

Figure 10 shows a signaling scenario for a single flow being managed by two signaling applications using the path-coupled message routing method. The flow sender and receiver and one router support both, two other routers support one each. The figure also shows the routing state table at node B.



Message Routing Information	Session ID	NSLPID	Routing State
MRM = Path Coupled; Flow ID = {IP-A, IP-E, proto/ports}; D=up	0xABCD	NSLP1	IP-A
MRM = Path Coupled; Flow ID = {IP-A, IP-E, proto/ports}; D=down	0xABCD	NSLP1	(null)
MRM = Path Coupled; Flow ID = {IP-A, IP-E, proto/ports}; D=up	0x1234	NSLP2	IP-A
MRM = Path Coupled; Flow ID = {IP-A, IP-E, proto/ports}; D=down	0x1234	NSLP2	Points to B-D MA

Figure 10: A Signaling Scenario

The upstream state is just the same address for each application. For the downstream direction, NSLP1 only requires D-mode messages and so no explicit routing state towards C is needed. NSLP2 requires a messaging association for its messages towards node D, and node C does not process NSLP2 at all, so the peer state for NSLP2 is a pointer to a messaging association that runs directly from B to D. Note that E is not visible in the state table (except implicitly in the address in the message routing information); routing state is stored only for adjacent peers. (In addition to the peer identification, IP hop counts are stored for each peer where the state itself if not null; this is not shown in the table.)

Figure 11 shows the message sequence for a GIST handshake that sets

up the messaging association for B-D signaling. It shows the

exchange of Stack Proposals and MA-protocol-options data in each direction. Then the Querying node selects TLS/TCP as the stack configuration to use and sets up the messaging association over which it sends the Confirm.

```

----- Query ----->
IP(Src=IP#A; Dst=IP#E; RAO for NSLP2); UDP(Src=6789; Dst=GIST)
GIST(Header(Type=Query; NSLPID=NSLP2; R=1; S=0)
  MRI(MRM=Path-Coupled; Flow=F; Direction=down)
  SessionID(0x1234)
  NLI(Peer='string1'; IA=IP#B)
  QueryCookie(0x139471239471923526)
  StackProposal(#Proposals=3;1=TLS/TCP; 2=TLS/SCTP; 3=TCP)
  StackConfigurationData(#MPO=2;
    TCP(Applicable: all; Data: null)
    SCTP(Applicable: all; Data: null)))

<----- Response -----
IP(Src=IP#D; Dst=IP#B); UDP(Src=GIST; Dst=6789)
GIST(Header(Type=Response; NSLPID=NSLP2; R=1; S=1)
  MRI(MRM=Path-Coupled; Flow=F; Direction=up)
  SessionID(0x1234)
  NLI(Peer='stringr2', IA=IP#D)
  QueryCookie(0x139471239471923526)
  ResponderCookie(0xacdefedcdfaeeded)
  StackProposal(#Proposals=3; 1=TCP; 2=SCTP; 3=TLS/TCP)
  StackConfigurationData(#MPO=3;
    TCP(Applicable: 3; Data: port=6123)
    TCP(Applicable: 1; Data: port=5438)
    SCTP(Applicable: all; Data: port=3333)))

-----TCP SYN----->
<-----TCP SYN/ACK-----
-----TCP ACK----->
TCP connect(IP Src=IP#B; IP Dst=IP#D; Src Port=9166; Dst Port=6123)
<-----TLS INIT----->

----- Confirm ----->
[Sent within messaging association]
GIST(Header(Type=Confirm; NSLPID=NSLP2; R=0; S=1)
  MRI(MRM=Path-Coupled; Flow=F; Direction=down)
  SessionID(0x1234)
  NLI(Peer='string1'; IA=IP#B)
  ResponderCookie(0xacdefedcdfaeeded)
  StackProposal(#Proposals=3; 1=TCP; 2=SCTP; 3=TLS/TCP))

```

Figure 11: GIST Handshake Message Sequence

[Appendix D](#). Change History

[D.1](#). Changes In Version -10

1. Added further guidance on parameter setting for initial backoff and rate control for D-mode to [Section 5.3.3](#) [AD review comment M1].
2. Rephrased the end of [Section 8.6](#) to highlight the threat left open when the Querying node does not apply a strong security policy to offered Stack-Proposal [AD review comment M2].
3. Clarified in [Section 7.2](#) that although NAT behaviour is only informatively described in this specification, it is being defined in a separate document [AD review M3].
4. Strengthened and clarified the reference to the extensibility document for technical guidance on codepoint allocation, and made the reference normative. Added rationale for the 'Reserved' blocks in the various registries, and added further notes on what information must be provided to support an allocation request [AD review comment M4].
5. Fixed an identifier collision in the ABNF for the GIST messages in [Section 5.2.2](#) (Common-Header in the message header and common-header as a payload in error messages) and re-verified the ABNF [AD review comment L1].
6. Clarified the text in [Section 3.3](#) about the impact on NATs of defining a new MRM, referring to the specification split described in [Section 7.2.2](#). Also added a flag to the MRM format (Appendix A.3.1) to denote MRIs which do not contain network or transport addresses, and made more specific the error message to be returned if a NAT does not understand an MRM in [Section 7.2.1](#) [AD review comment L2].
7. Added discussion in [Section 4.1.2](#) on delivery failure detected for reliable messaging in general, and for the case of Forwards-TCP in particular in [Section 5.7.2](#). Also noted that this needs to be considered for future MA-Protocol-IDs used for reliable messaging ([Section 9](#)) [AD review comment L3].
8. Added clarifying text to [Section 5.1](#) on what it means to invert the direction of an MRI [AD review comment L5].
9. Enhanced the format descriptions in [Appendix A](#) to include

descriptions of all message and object fields and also field lengths [AD review comment L6].

Internet-Draft

GIST

July 2006

10. Added more explanation in [Section 5.2.2](#) of how a message direction is defined, in particular in the context of TTL measurement [AD review comment L7].
11. Added a new explanation of why a well-known port is needed for the query encapsulation in [Section 5.3.2](#) [AD review comment L8].
12. Added a note that DCCP does not provide reliability in [Section 5.4.1](#) [AD review comment L9].
13. Clarified the rules on how long to retain stack configuration data in [Section 5.7.1](#) and included a default timer value [AD review comment L10].
14. Modified the text about stack-proposal verification as part of downgrade protection in [Section 5.7.1](#), to clarify that the MUST applies directly to the object verification itself; also noted the action to be taken in case of a failed verification [AD review comment L11].
15. Added further information on the addressing used in opening a forwards-TCP connection in [Section 5.7.2](#) [AD review comment L12].
16. Modified the text in [Section 5.8.1.2](#) to say that using the signaling source address is a consequence of setting DF itself rather than why DF was set in the first place; also weakened the instruction from MUST to SHOULD [AD review comment L14].
17. Added further clarification of why routing state installed by a downstream Query should supersede that from an upstream Query in [Section 5.8.1.3](#) [AD review comment L15].
18. Corrected a timer in the Messaging Association state machine ([Section 6.4](#)) from NoHello to SendHello. Also, added default values for MA-Hold-Time and route change probe frequency, and explanatory text for each, to [Section 4.4.3](#) [AD review comment L16].

19. Re-arranged the text in [Section 7.2](#) to highlight the rules about precisely which messages are and are not translated in a GIST-specific way by NATs [AD review comment L19].
20. Explicitly noted that 'r' bits are also reserved in [Appendix A.2](#) [AD review comment L20].
21. Added an error condition for processing messages which have the extensibility flags AB set to 11 in [Appendix A.2.1](#) [AD review

- comment L21].
22. Fixed the table of MRM identifiers in [Section 9](#) so the field name matches that in [Appendix A.3.1](#) [AD review comment L22].
 23. Clarified why only D=0 is valid for the loose-end MRM in [Appendix A.3.1.2](#) [AD review comment L23].
 24. Clarified the rules about processing the NAT traversal object in [Appendix A.3.8](#) to cover the case where there are several NATs along the path with different capabilities [AD review comment L24].
 25. Strengthened the text in [Appendix A.4.1](#) to be clearer about what additional information fields must be included in error messages [AD review comment L25].
 26. Tidied up the use of acronyms throughout the document, including adding some to the terminology list in [Section 2](#) [AD review comment N1].
 27. Added references to [RFC4086](#) and updated 2119 language for cryptographic randomness of SIDs and cookies in [Section 3.5](#) and [Section 8.5](#) respectively [AD review comment N2].
 28. Modified the transition labelling in Figure 7 to make it clearer that in the Established-Established transition, the [!confirmRequired] qualification applies only to the rx_Query case [AD review comment N4].
 29. Added a reference for OSPF in [Section 7.1.2](#) [AD review comment N5].

30. Changed NAT terminology from public/private to external/internal to match BEHAVE usage in [Section 7.2](#) and [Section 7.4](#) [AD review comment N6].
31. Updated a number of i-d references to published RFCs or working group documents [AD review comment N7 partial].
32. Fixed [rfc2119](#) capitalisation of MUST not in [Appendix A.3.5](#) [AD review comment N8].
33. Fixed an error subcode name from 'Invalid Object' to 'Invalid Object Type' in [Appendix A.4.4.9](#) [AD review comment N9].
34. Added the NTO to the GIST message ABNF in [Section 5.1](#) and updated the forward reference to the NAT traversal section

[tracker issue 104].

35. Removed a spurious rule about creating listening MAs in [Section 6.3](#) and strengthened the rules about needing to have these available but with an open policy on when to create and destroy them in [Section 5.7.1](#) [tracker issue 105].
36. Added text that limits the applicability of the private/experimental space to closed network environments [tracker issue 106].
37. Added text in [Section 7.1.4](#) encouraging GIST to use a single SII across multiple sessions if possible to allow signaling application aggregation [tracker issue 107].
38. Specified that this document would define GIST version 1 in [Section 5.2.1](#) [tracker issue 108].
39. Added the ability for RecvMessage to pass up interface attributes in [Appendix B.2](#) [tracker issue 110].
40. Added additional text on rules for selecting stack proposals and MA re-use in [Section 5.7.1](#) to ensure that re-used associations have properties that the Querying node actually needs [tracker issue 111].

41. Added a brief introduction to the GIST message types in a new [Section 3.4](#).

In addition, the following AD review comments did not lead to text changes. See the mailing list discussion at <http://www1.ietf.org/mail-archive/web/nsis/current/msg06307.html>.

L4: Direct use of PMTUD by GIST.

L13: Use of TLS 1.0 rather than 1.1.

L17: Guidance on NSLP behaviour during rerouting,

L18: Behaviour of GIST-unaware NATs.

N3: Node state machine logic.

[D.2](#). Changes In Version -09

1. Added a new [Section 3.6](#) clarifying the relationship between signaling applications and NSLPIDs; modified terminology in the remainder of the document likewise.

2. Added a new [Section 8.6](#) explaining the rationale behind the downgrade attack prevention mechanism.
3. Re-wrote parts of [Section 4.3.2](#), [Section 6.1](#) and [Appendix B.2](#) to clarify the way that GIST is assumed to interact with signaling applications to exercise policy control over whether or not two nodes become signaling peers during a GIST handshake.
4. Generalised an error message [Appendix A.4.4.12](#) to cover additional MRI validation checks in [Section 4.3.4](#) and [Section 5.8.1.2](#).
5. Allowed an optional Stack-Configuration-Data object in Confirm messages to allow messaging association lifetime to be negotiated even in the case of late state installation at the Responding node (see [Section 4.4.1](#) and [Section 4.4.3](#)).
6. Removed the option in [Section 4.4.2](#) of allowing a node to treat

messaging associations with the same authenticated end points as equivalent.

7. Include additional guidance in [Section 4.4.3](#) to prevent routing state being erroneously refreshed in the case of rerouting events; also included general guidance notes on timer setting.
8. Clarified that the Stack-Proposal lists protocols in top-to-bottom order (see [Section 5.7.1](#)).
9. Enhanced the definition of TLS usage in [Section 5.7.3](#) with details on ciphersuite requirements and authentication methods.
10. Tidied up terminology and discussion of how protocol options data is carried in the SCD; renamed higher-layer-addressing to MA-protocol-options.

[D.3](#). Changes In Version -08

1. Changed the protocol name from GIMPS to GIST (everywhere).
2. Inserted [RFC2119](#) language (MUST etc.) in the appropriate places.
3. Added references to the actions to be taken in various error conditions, including the error messages to be send (throughout).
4. Added legacy NAT traversal to the list of excluded functions in [Section 1.1](#).

5. Included some text at the end of [Section 3.3](#) analysing the case of a GIST node which does not support a particular MRM.
6. Added a flag to mark when messages have been explicitly routed, so they can bypass validation against current routing state (see [Section 4.3.1](#), TBD).
7. Re-wrote the discussion in [Section 4.3.4](#) to cover all cases of nodes not hosting an NSLP (including end systems), in particular the validations that can be performed at intermediate GIST nodes (this replaces the old [section 7.2](#)).

8. Clarified the rules about R and S flag setting in the common header and D flag in the MRI ([Section 5](#)).
9. Included discussion of how a node with a choice of interfaces or IP versions should select one to use in the NLI ([Section 5.2.2](#)).
10. Modified the description of messaging association protocol selections ([Section 5.7](#) and elsewhere) to clarify that this is essentially capability discovery rather than an open ended protocol negotiation.
11. Modified the description of how higher layer addressing information is carried ([Section 5.7.1](#) and [Appendix A.3.5](#)) to allow the data to be tagged against a specific profile if necessary, or omitted if the protocol does not need it.
12. Added a higher layer protocol definition for TLS in [Section 5.7.3](#).
13. Simplified and restructured the state machine presentation in [Section 6](#), in particular using a single list for the events and eliminating the transition tables. Also modified the operation of the Responder machine to handle retransmitted Query messages correctly.
14. Re-wrote the route change handling text in [Section 7.1](#) to clarify the relative responsibilities of GIST and NSLPs and their interaction through the API. Notifications are now assumed to be a signaling application responsibility, and GIST behaviour is defined in terms of handling changes in a 3-state model of the correctness of the routing state for each direction.
15. Updated the NAT traversal description in [Section 7.2](#), including normative text about how GIST nodes should handle messages containing NAT-Traversal objects.

16. Likewise, clarified that the responsibility for session/flow binding in the case of tunnelling is handled by NSLPs ([Section 7.3](#)).

17. Formalised the IANA considerations ([Section 9](#)).
18. Extended the routing state example (Appendix C) to include a message sequence for association setup.
19. Re-arranged the sequence of sections, including placing this change history at the end.

[D.4](#). Changes In Version -07

1. The open issues section has finally been removed in favour of the authoritative list of open issues in an online issue tracker at <http://nsis.srmr.co.uk/cgi-bin/roundup.cgi/nsis-ntlp-issues/index>.
2. Clarified terminology on peering and adjacencies that there may be NSIS nodes between GIMPS peers that do some message processing, but that are not explicitly visible in the peer state tables.
3. Added a description of the loose-end MRM ([Section 5.8.2](#) and [Appendix A.3.1.2](#)).
4. Added a description of an upstream Query encapsulation for the path-coupled MRM, [Section 5.8.1.3](#), including rationale for and restrictions on its use.
5. The formal description of the protocol in [Section 6](#) has been significantly updated and extended in terms of detail.
6. Modified the description of the interaction between NSLPs and GIMPS for handling inbound messages for which no routing state exists, to allow the NSLP to indicate whether state setup should proceed and to provide NSLP payloads for the Response or forwarded message ([Section 3.7](#), [Section 4.3.2](#) and [Appendix B](#)).
7. Included new text, [Section 5.6](#), on the processing and encapsulation of error messages. Also added formats and an error message catalogue in [Appendix A.4](#), including a modified format for the overall GIMPS-Error message and the GIMPS-Error-Data object.
8. Removed the old [section 5.3.3](#) on NSLPID/RA0 setting on the assumption that this will be covered in the extensibility document.

9. Included a number of other minor corrections and clarifications.

[D.5.](#) Changes In Version -06

Version -06 does not introduce any major structural changes to the protocol definition, although it does clarify a number of details and resolve some outstanding open issues. The primary changes are as follows:

1. Added a new high level [Section 3.3](#) which gathers together the various aspects of the message routing method concept.
2. Added a new high level [Section 3.5](#) which explains the concept and significance of the session identifier. Also clarified that the routing state always depends on the session identifier.
3. Added notes about the level of address validation performed by GIMPS in [Section 4.1.2](#) and extensions to the API in [Appendix B](#).
4. Split the old Node-Addressing object into a Network-Layer-Information object and Stack-Configuration-Data object. The former refers to basic information about a node, and the latter carries information about messaging association configuration. Redefined the content of the various handshake messages accordingly in [Section 4.4.1](#) and [Section 5.1](#).
5. Re-wrote [Section 4.4.3](#) to clarify the rules on refresh and purge of routing state and messaging associations. Also, moved the routing state lifetime into the Network-Layer-Information object and added a messaging association lifetime to the Stack-Configuration-Data object ([Section 5.2](#)).
6. Added specific message types for errors and MA-Refresh in [Section 5.1](#). The error object is now GIMPS-specific ([Appendix A.4.1](#)).
7. Moved the Flow-Identifier information about the message routing method from the general description of the object to the path-coupled MRM section ([Section 5.8.1.1](#)), and made a number of clarifications to the bit format ([Appendix A.3.1.1](#)).
8. Removed text about assumptions on the version numbering of NSLPs, and restricted the scope of the description of TLV object formats and extensibility flags to GIMPS rather than the whole of NSIS ([Appendix A](#)).
9. Added a new [Section 5.5](#) explaining the possible relationships

between message types and encapsulation formats.

Internet-Draft

GIST

July 2006

10. Added a new [Section 6](#) in outline form, to capture the formal specification of the protocol operation.
11. Added new security sections on cookie requirements ([Section 8.5](#)) and residual threats ([Section 8.7](#)).

[D.6](#). Changes In Version -05

Version -05 reformulates the specification, to describe routing state maintenance in terms of exchanging explicitly identified Query/Response/Confirm messages, leaving the upstream/downstream distinction as a specific detail of how Query messages are encapsulated. This necessitated widespread changes in the specification text, especially [Section 4.2.1](#), [Section 4.4](#), [Section 5.1](#) and [Section 5.3](#) (although the actual message sequences are unchanged). A number of other issues, especially in the area of message encapsulation, have also been closed. The main changes are the following:

1. Added a reference to an individual draft on the Loose End MRM as a concrete example of an alternative message routing method.
2. Added further text (particularly in [Section 2](#)) on what GIMPS means by the concept of 'session'.
3. Firmed up the selection of UDP as the encapsulation choice for D-mode, removing the open issue on this topic.
4. Defined the interaction between GIMPS and signaling applications for communicating about the cryptographic security properties of how a message will be sent or has been received (see [Section 4.1.2](#) and [Appendix B](#)).
5. Closed the issue on whether Query messages should use the signaling or flow source address in the IP header; both options are allowed by local policy and a flag in the common header indicates which was used. (See [Section 5.8.1.2](#).)
6. Added the necessary information elements to allow the IP hop count between adjacent GIMPS peers to be measures and reported.

(See [Section 5.2.2](#) and [Appendix A.3.3](#).)

7. The old open-issue text on selection of IP router alert option values has been moved into the main specification to capture the technical considerations that should be used in assigning such values (in old [section 5.3.3](#)).

8. Resolved the open issue on lost Confirm messages by allowing a choice of timer-based retransmission of the Response, or an error message from the responding node which causes the retransmission of the Confirm (see [Section 5.3.3](#)).
9. Closed the open issue on support for message scoping (this is now assumed to be a NSLP function).
10. Moved the authoritative text for most of the remaining open issues to an online issue tracker.

[D.7](#). Changes In Version -04

Version -04 includes mainly clarifications of detail and extensions in particular technical areas, in part to support ongoing implementation work. The main details are as follows:

1. Substantially updated [Section 4](#), in particular clarifying the rules on what messages are sent when and with what payloads during routing and messaging association setup, and also adding some further text on message transfer attributes.
2. The description of messaging association protocol setup including the related object formats has been centralised in a new [Section 5.7](#), removing the old [Section 6.6](#) and also closing old open issues 8.5 and 8.6.
3. Made a number of detailed changes in the message format definitions (Appendix A), as well as incorporating initial rules for encoding message extensibility information. Also included explicit formats for a general purpose Error object, and the objects used to discover supported messaging association protocols. Updated the corresponding open issues section (old

[section 9.3](#)) with a new item on NSLP versioning.

4. Updated the GIMPS API (Appendix B), including more precision on message transfer attributes, making the NSLP hint about storing reverse path state a return value rather than a separate primitive, and adding a new primitive to allow signaling applications to invalidate GIMPS routing state. Also, added a new parameter to SendMessage to allow signaling applications to 'bypass' a message statelessly, preserving the source of an input message.
5. Added an outline for the future content of an IANA considerations section ([Section 9](#)). Currently, this is restricted to identifying the registries and allocations required, without defining the allocation policies and other

considerations involved.

6. Shortened the background design discussion in [Section 3](#).
7. Made some clarifications in the terminology section relating to how the use of C-mode does and does not mandate the use of transport or security protection.
8. The ABNF for message formats in [Section 5.1](#) has been re-written with a grammar structured around message purpose rather than message direction, and additional explanation added to the information element descriptions in [Section 5.2](#).
9. The description of the D-mode transport in [Section 5.3](#) has been updated. The encapsulation rules (covering IP addressing and UDP port allocation) have been corrected, and a new subsection on message retransmission and rate limiting has been added, superseding the old open issue on the same subject ([section 8.10](#)).
10. A new open issue on IP TTL measurement to detect non-GIMPS capable hops has been added (old [section 9.5](#)).

[D.8](#). Changes In Version -03

Version -03 includes a number of minor clarifications and extensions

compared to version -02, including more details of the GIMPS API and messaging association setup and the node addressing object. The full list of changes is as follows:

1. Added a new section pinning down more formally the interaction between GIMPS and signaling applications ([Section 4.1](#)), in particular the message transfer attributes that signaling applications can use to control GIMPS ([Section 4.1.2](#)).
2. Added a new open issue identifying where the interaction between the security properties of GIMPS and the security requirements of signaling applications should be identified (old [section 9.10](#)).
3. Added some more text in [Section 4.2.1](#) to clarify that GIMPS has the (sole) responsibility for generating the messages that refresh message routing state.
4. Added more clarifying text and table to GHC and IP TTL handling discussion of [Section 4.3.4](#).
5. Split [Section 4.4](#) into subsections for different scenarios, and added more detail on Node-Addressing object content and use to

handle the case where association re-use is possible in [Section 4.4.2](#).

6. Added strawman object formats for Node-Addressing and Stack-Proposal objects in [Section 5.1](#) and [Appendix A](#).
7. Added more detail on the bundling possibilities and appropriate configurations for various transport protocols in [Section 5.4.1](#).
8. Included some more details on NAT traversal in [Section 7.2](#), including a new object to carry the untranslated address-bearing payloads, the NAT-Traversal object.
9. Expanded the open issue discussion in old [section 9.3](#) to include an outline set of extensibility flags.

[D.9](#). Changes In Version -02

Version -02 does not represent any radical change in design or

structure from version -01; the emphasis has been on adding details in some specific areas and incorporation of comments, including early review comments. The full list of changes is as follows:

1. Added a new [Section 1.1](#) which summarises restrictions on scope and applicability; some corresponding changes in terminology in [Section 2](#).
2. Closed the open issue on including explicit GIMPS state teardown functionality. On balance, it seems that the difficulty of specifying this correctly (especially taking account of the security issues in all scenarios) is not matched by the saving of state enabled.
3. Removed the option of a special class of message transfer for reliable delivery of a single message. This can be implemented (inefficiently) as a degenerate case of C-mode if required.
4. Extended [Appendix A](#) with a general discussion of rules for message and object formats across GIMPS and other NSLPs. Some remaining open issues are noted in old [section 9.3](#) (since removed).
5. Updated the discussion of RAO/NSLPID relationships to take into account the proposed message formats and rules for allocation of NSLP id, and propose considerations for allocation of RAO values.

6. Modified the description of the information used to route messages (first given in [Section 4.2.1](#) but also throughout the document). Previously this was related directly to the flow identification and described as the Flow-Routing-Information. Now, this has been renamed Message-Routing-Information, and identifies a message routing method and any associated addressing.
7. Modified the text in [Section 4.3](#) and elsewhere to impose sanity checks on the Message-Routing-Information carried in C-mode messages, including the case where these messages are part of a GIMPS-Query/Response exchange.

8. Added rules for message forwarding to prevent message looping in a new [Section 4.3.4](#), including rules on IP TTL and GIMPS hop count processing. These take into account the new RAO considerations described above.
9. Added an outline mechanism for messaging association protocol stack setup, with the details in a new [Section 6.6](#) and other changes in [Section 4.4](#) and the various sections on message formats.
10. Removed the open issue on whether storing reverse routing state is mandatory or optional. This is now explicit in the API (under the control of the local NSLP).
11. Added an informative annex describing an abstract API between GIMPS and NSLPs in [Appendix B](#).

[D.10](#). Changes In Version -01

The major change in version -01 is the elimination of 'intermediaries', i.e. imposing the constraint that signaling application peers are also GIMPS peers. This has the consequence that if a signaling application wishes to use two classes of signaling transport for a given flow, maybe reaching different subsets of nodes, it must do so by running different signaling sessions; and it also means that signaling adaptations for passing through NATs which are not signaling application aware must be carried out in D-mode. On the other hand, it allows the elimination of significant complexity in the C-mode handling and also various other protocol features (such as general route recording).

The full set of changes is as follows:

1. Added a worked example in [Section 3.7](#).

2. Stated that nodes which do not implement the signaling application should bypass the message ([Section 4.3](#)).
3. Decoupled the state handling logic for routing state and messaging association state in [Section 4.4](#). Also, allow

messaging associations to be used immediately in both directions once they are opened.

4. Added simple ABNF for the various GIMPS message types in a new [Section 5.1](#), and more details of the common header and each object in [Section 5.2](#), including bit formats in [Appendix A](#). The common header format means that the encapsulation is now the same for all transport types ([Section 5.4.1](#)).
5. Added some further details on D-mode encapsulation in [Section 5.3](#), including more explanation of why a well known port is needed.
6. Removed the possibility for fragmentation over DCCP ([Section 5.4.1](#)), mainly in the interests of simplicity and loss amplification.
7. Removed all the tunnel mode encapsulations (old sections [5.3.3](#) and [5.3.4](#)).
8. Fully re-wrote the route change handling description ([Section 7.1](#)), including some additional detection mechanisms and more clearly distinguishing between upstream and downstream route changes. Included further details on GIMPS/NSLP interactions, including where notifications are delivered and how local repair storms could be avoided. Removed old discussion of propagating notifications through signaling application unaware nodes (since these are now bypassed automatically). Added discussion on how to route messages for local state removal on the old path.
9. Revised discussion of policy-based forwarding (old [Section 7.2](#)) to account for actual Flow-Routing-Information definition, and also how wilddarding should be allowed and handled.
10. Removed old route recording section (old [Section 6.3](#)).
11. Extended the discussion of NAT handling ([Section 7.2](#)) with an extended outline on processing rules at a GIMPS-aware NAT and a pointer to implications for C-mode processing and state management.

12. Clarified the definition of 'correct routing' of signaling messages in [Section 8](#) and GIMPS role in enforcing this. Also, opened the possibility that peer node authentication could be signaling application dependent.
13. Removed old open issues on C-mode Encapsulation ([section 8.7](#)); added new open issues on Message Routing (old [Section 9.3](#) of version -05, later moved to [Section 3.3](#)) and D-mode congestion control.
14. Added this change history.

Internet-Draft

GIST

July 2006

Authors' Addresses

Henning Schulzrinne
Columbia University
Department of Computer Science
450 Computer Science Building
New York, NY 10027
US

Phone: +1 212 939 7042
Email: hgs+nsis@cs.columbia.edu
URI: <http://www.cs.columbia.edu>

Robert Hancock
Siemens/Roke Manor Research
Old Salisbury Lane
Romsey, Hampshire S051 0ZN
UK

Email: robert.hancock@roke.co.uk
URI: <http://www.roke.co.uk>

Internet-Draft

GIST

July 2006

Intellectual Property Statement

The IETF takes no position regarding the validity or scope of any Intellectual Property Rights or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; nor does it represent that it has made any independent effort to identify any such rights. Information on the procedures with respect to rights in RFC documents can be found in [BCP 78](#) and [BCP 79](#).

Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <http://www.ietf.org/ipr>.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights that may cover technology that may be required to implement this standard. Please address the information to the IETF at ietf-ipr@ietf.org.

Disclaimer of Validity

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Copyright Statement

Copyright (C) The Internet Society (2006). This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights.

Acknowledgment

Funding for the RFC Editor function is currently provided by the Internet Society.