

RMT Working Group
Internet Engineering Task Force
Category: Informational
December 2003
Expires June 2004

Brian Whetten, Consultant
Dah Ming Chiu, CUHK
Miriam Kadansky, Sun Microsystems
Seok Joo Koh, ETRI
Gursel Taskale, Tibco

Tree-Based ACK (TRACK) Building Block
for Reliable Multicast Transport

<[draft-chiu-rmt-bb-track-03.txt](#)>

Status of this Memo

This document is an Internet-Draft and is in full conformance with all provisions of [Section 10 of RFC 2026](#).

Internet-Drafts are 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 a "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>.

Abstract

This document defines the Tree-based ACK (TRACK) building block (BB) for reliable multicast transport (RMT) protocol instantiations. As an RMT building block, the TRACK BB is a coarse-grained modular component that may be common to multiple RMT protocols. The TRACK BB is designed to provide application-level confirmed delivery, local recovery, and enhanced flow and congestion control, and it assumes that the TREE BB (RFCyyyy) provides automatic tree configuration.

Table of Contents

1.	Introduction.....	3
2.	Terminology.....	4
3.	BB Rationale.....	5
4.	Functionality of TRACK BB.....	5
4.1	Hierarchical Session Creation and Maintenance.....	5
4.2	Data Sessions.....	6
4.3	TRACK Generation and Aggregation.....	7
4.4	Statistics Aggregation.....	7
4.5	Distributed RTT Calculations.....	7
5.	Applicability Statement.....	8
5.1	Application Types.....	9
5.2	Network Infrastructure.....	9
5.3	Manual vs. Automatic Controls.....	9
5.4	Heterogeneous Networks.....	9
5.5	Use of Network Infrastructure.....	10
5.6	Deployment Constraints.....	10
5.7	Target Scalability.....	10
5.8	Known Failure Modes.....	10
6.	TRACK Architecture.....	11
6.1	TRACK Entities.....	11
6.2	Basic Operation of the Protocol.....	13
7.	Details: TRACK Functionality.....	16
7.1	Session Creation and Maintenance.....	16
7.2	Data Sessions.....	22
7.3	Control Traffic Generation and Aggregation.....	27
7.4	Application Level Confirmed Delivery.....	30
7.5	Distributed RTT Calculations.....	32
7.6	SNMP Support.....	33
7.7	Late Join Semantics.....	33
8.	TRACK Message Types.....	34
9.	Global Configuration Parameters.....	38
9.1	Configuration Variables.....	38
9.2	Constants.....	39
9.3	Reason Codes.....	39
10.	Requirements from other Building Blocks.....	40
11.	Security Considerations.....	40
12.	References.....	41
13.	Acknowledgments.....	42
14.	Author's Addresses.....	42

1. Introduction

The Reliable Multicast Transport (RMT) working group was chartered to standardize IP multicast transport services [[RFC2887](#)]. Rather than create a set of monolithic protocol specifications, the RMT WG has chosen to break the reliable multicast protocols into Building Blocks (BB) and Protocol Instantiations (PI). A Building Block is a specification of the algorithms of a single component, with an abstract interface to other BBs and PIs. A PI combines a set of BBs, adds in the additional required functionality not specified in any BB, and specifies the specific instantiation of the protocol.

There are two primary reliability requirements for a transport protocol: ensuring goodput and confirming delivery. Other documents describe RMT building blocks to ensure goodput [[RFC3450](#), [RFC3451](#), [RFC3452](#), NORM-BB, NORM-PI], while this document describes the Tree-based ACK building block, or TRACK BB, which is concerned with confirming delivery. Specifically, the TRACK BB is designed to offer application-level confirmed delivery, aggregation of control traffic and sender statistics, local recovery, automatic tree building, and enhanced flow and congestion control.

The TRACK BB assumes that there is a Tree auto-configuration building block (e.g., the TREE BB [[RFCyyyy](#)]), which provides the list of parents to which each node joins. If receivers may serve as Repair Heads, the TRACK BB assumes the TREE BB is also responsible for selecting the role of each receiver as either receiver or Repair Head.

The TRACK BB is organized around a data channel and a control channel. The data channel is responsible for multicast data from the sender to all other nodes in a TRACK session. In order to integrate with goodput-ensuring transport protocols, these protocols are used as the data channel for a given data session. This data channel MAY also provide congestion control.

2. Terminology

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

In addition, the following terms are applied in this document as well as the TREE BB document [[RFCyyyy](#)].

Session

A session is used to distribute data over a multicast address. A Session Tree is used to provide reliability and feedback services for a session.

Session Identifier

A fixed-size number, chosen either by the application that creates the session or by the transport. Senders and receivers use the session Identifier to distinguish sessions.

Repair Head (RH)

A node within the tree which receives and retransmits data, and aggregates and forwards control information toward the sender. The sender operates as the root repair head in any session tree. An RH that has accepted children for a session is called a parent.

Session Tree (ST)

The session tree is a tree spanning all receivers of a multicast session. It is rooted at the sender, consisting of zero or more repair heads as interior nodes, and zero or more receivers as leaf nodes.

Parent

A parent is an RH or receiver's predecessor in the ST on the path toward the sender. Every RH or receiver on the tree except the sender itself has a parent. Each parent communicates with its children using either an assigned multicast address or through unicast.

Children

The set of receivers and RHs for which an RH or the sender is providing repair and feedback services.

3. BB Rationale

TRACK BB is primarily designed to run in conjunction with another transport protocol that is responsible for ensuring goodput.

The TRACK BB is responsible for specifying all of the TRACK-specific functionality. It interfaces with the TREE BB. The TRACK PI is then responsible for instantiating a complete protocol that includes all of the other components.

4. Functionality of TRACK BB

This TRACK BB is designed based on the following recommendations, as described in [Section 4.5 of RFC 3048](#):

It has been shown that the scalability of RM protocols can be greatly enhanced by the insertion of some kind of retransmission or feedback aggregation agents between the source and receivers. These agents are then used to form a tree with the source at (or near) the root, the receivers at the leaves of the tree, and the aggregation/local repair nodes in the middle. The internal nodes can either be dedicated software for this task, or they may be receivers that are performing dual duty.

The effectiveness of these agents to assist in the delivery of data is highly dependent upon how well the logical tree they use to communicate matches the underlying routing topology. The purpose of this building block would be to construct and manage the logical tree connecting the agents. Ideally, this building block would perform these functions in a manner that adapts to changes in session membership, routing topology, and network availability.

The TRACK BB provides the following detailed functionality.

4.1 Hierarchical Session Creation and Maintenance

This set of functionality is responsible for creating and maintaining a hierarchical tree of Repair Heads and receivers.

- o Bind. When a child knows the parent it wishes to join to for a given Data Session, it binds to that parent.
- o Unbind. When a child wishes to leave a data session, either because the session is over or because the application is finished with the session, it initiates an unbind operation with its parent.

- o Eject. A parent can also force a child to unbind. This happens if the parent needs to leave the session, if the child is not behaving correctly, or if the parent wants to move the child to another parent as part of tree configuration maintenance.
- o Fault Detection. In order to verify liveness, parents and children send regular heartbeat messages between themselves. The sender also sends regular null data messages to the group, if it has no data to send.
- o Fault Recovery. When a child detects that its parent is no longer reachable, it may switch to another parent. When a parent detects that one of its children is no longer reachable, it removes that child from its membership list and reports this up the tree to the sender of the Data Session.
- o Distributed Membership. Each parent is responsible for maintaining a local list of the children attached to it.

4.2 Data Sessions

This functionality is responsible for the reliable, ordered transmission of a set of data messages. These are initially transmitted using another transport protocol, the Data Channel Protocol, which has primary responsibility for ensuring goodput.

- o Data Transmission. The sender takes sequenced data messages from the application, and passes them to the data channel protocol for multicast transmission. It delays passing them to the data channel protocol if it is presently flow controlled.
- o Flow Control and Buffer Management. Senders and Repair Heads MAY maintain a set of buffers that are at least as large as the senders transmission window. The receivers pass their reception status up to the sender as part of their TRACK messages. This MAY be used to advance the buffer windows at each node and limit the senders window advancement to the speed of the slowest sender.
- o Retransmission Requests. While primary responsibility for goodput rests with the data channel protocol, receivers MAY request retransmission of lost messages from their parents.
- o Local Recovery. Repair Heads keep track of retransmission requests from their children, and provide repairs to them. If a Repair Head cannot fulfill a retransmission request, it forwards it up the tree.

- o End of Stream. When a data session is completed, this is signaled as an End of Stream condition.

4.3 TRACK Generation and Aggregation

This set of functionality is responsible for periodically generating TRACK messages from all receivers and aggregating them at Repair Heads. These messages provide updated flow control window information, roundtrip time measurements, and congestion control statistics. They OPTIONALLY acknowledge receipt of data, OPTIONALLY report missing messages, and OPTIONALLY provide group statistics.

The algorithms include:

- o TRACK Timing. In order to avoid ACK implosion, the senders and Repair Heads use timing algorithms to control the speed at which TRACK messages are sent.
- o TRACK Aggregation. In order to provide the highest levels of scalability and reliability, interior tree nodes provide aggregation of control traffic flowing up the tree. The aggregated feedback information includes that used for end-to-end confirmed delivery, flow control, congestion control, and group membership monitoring and management.
- o Statistics Request. A sender may prompt senders to generate and report a set of statistics back to the sender.

4.4 Statistics Aggregation

In addition to the predefined aggregation types, aggregation of self-describing data may also be performed on sender statistics flowing up the tree.

4.5 Distributed RTT Calculations

One of the primary challenges of congestion control is efficient RTT calculation. TRACK provides two methods to perform these calculations.

- o sender Per-Message RTT Calculations. On demand, a sender stamps outgoing messages with a timestamp. As each TRACK is passed up the tree, the amount of dally time spent waiting at each node is accumulated. The lowest measurements are passed up the tree, and the dally time is subtracted from the original measurement.

- o Local Per-Level RTT Calculations. Each parent measures the local RTT to each of its children as part of the keep-alive messages used for failure detection.

5. Applicability Statement

The primary objective of TRACK is to provide additional functionality in conjunction with a receiver reliable protocol. These functions MAY include application layer reliability, enhanced congestion control, flow control, statistics reporting, local recovery, and automatic tree building. It is designed to do this while still offering scalability in the range of 10,000s of receivers per data session. The primary corresponding design tradeoffs are additional complexity, and lower isolation of nodes in the face of network and host failures.

There is a fundamental tradeoff between reliability and real-time performance in the face of failures. There are two primary types of single layer reliability that have been proposed to deal with this: sender reliable and receiver reliable delivery.

Sender reliable delivery is similar to TCP, where the sender knows the identity of the senders in a Data Session, and is notified when any of them fails to receive all the data messages. Sender reliable delivery limits knowledge of group membership and failures to only the actual senders. Senders do not have any knowledge of the membership of a group, and do not require senders to explicitly join or leave a Data Session. Sender reliable protocols scale better in the face of networks that have frequent failures, and have very high isolation of failures between senders. This TRACK BB provides sender reliable delivery, typically in conjunction with a sender reliable system.

This BB is specified according to the guidelines in [RFC 3269](#), along with [RFC 2357](#) and [RFC 2887](#). It specifies all communication between entities in terms of messages, rather than packets. A message is an abstract communication unit, which may be part of, or all of, a given packet. It does not have a specific format, although it does contain a list of fields, some of which may be optional, and some of which may have fixed lengths associated with them. It is up to each protocol instantiation to combine the set of messages in this BB, with those in other components, and create the actual set of packet formats that will be used.

As mentioned in the introduction, this BB assumes the existence of a separate TREE BB [[RFCyyyy](#)].

5.1 Application Types

TRACK is designed to support a wide range of applications that require one to many bulk data transfer and application layer confirmed delivery. Examples of applications that fit into the one-to-many data dissemination model are: real time financial news and market data distribution, electronic software distribution, audio video streaming, distance learning, software updates and server replication.

Historically, financial applications have had the most stringent reliability requirements, while audio video streaming have had the least stringent. For applications that do not require this level of reliability, or that demand the lowest levels of latency and the highest levels of failure isolation, TRACK may be less applicable.

5.2 Network Infrastructure

TRACK is designed to work over almost all multicast and broadcast capable network infrastructures. It is specifically designed to be able to support both asymmetrical and single source multicast environments.

Asymmetric networks with very low upbound bandwidth and a very low loss Data Channel may be better served solely through NACK based protocols, particularly if high reliability is not required. A good example is some satellite networks.

5.3 Manual vs. Automatic Controls

Some networks can take advantage of manual or centralized tools for configuring and controlling the usage of a reliable multicast group. In public Internet the tools have to span multiple administrative domains where policies may be inconsistent. Hence, it is preferable to design tools that are fully distributed and automatic. To address these requirements, TRACK provides automatic configuration, but can also support manual configuration options.

5.4 Heterogeneous Networks

While the majority of controlled networks are symmetrical and support many-to-many multicast, in designing a protocol for the Internet, we must deal with virtually most network types. These include asymmetrical networks, satellite networks, networks where only a single node may send to a multicast group, and wireless networks. TRACK takes this into account by not requiring any many-to-many multicast services.

5.5 Use of Network Infrastructure

TRACK is designed to run in either single level or hierarchical configurations. In a single level, there is no need for specialized network infrastructure. In hierarchical configurations, special nodes called Repair Heads are defined, which may run either as part of a distributed application, or as part of dedicated server software. TRACK does not specifically support or require Generic Router Assist or other router level assist.

5.6 Deployment Constraints

The two primary tradeoffs TRACK has, for the functionality it provides, are additional complexity, and decreased failure isolation. Hence, if target applications are to be deployed in networks with high rates of persistent failures, and isolation of failed senders from affecting other senders is of high importance, TRACK may not be appropriate. Similarly, if simplicity is paramount, TRACK may not be appropriate.

5.7 Target Scalability

The target scalability of TRACK is tens of thousands of simultaneous senders per Data Session. Dedicated Repair Heads are targeted to be able to support thousands of simultaneous Data Sessions.

5.8 Known Failure Modes

If a hierarchical control tree is mis-configured, so that loop-free, contiguous connection is not provided, failure will occur. This failure is designed to occur gracefully, at the initialization of a Data Session.

If the configuration parameters on control traffic are poorly chosen on an asymmetrical network, where there is much less control channel bandwidth available than data channel bandwidth, there may be a very high rate of control traffic. This control traffic is not dynamically congestion controlled like the data traffic, and so could potentially cause congestion collapse. This potential control channel overload could be exacerbated by an application that makes overly heavy use of the application level confirmation or statistics gathering functions.

6. TRACK Architecture

6.1 TRACK Entities

6.1.1 Node Types

TRACK divides the operation of the protocol into three major entities: sender, sender, and Repair Head.

It is assumed that senders and senders typically run as part of an application on an end host client. Repair Heads MAY be components in the network infrastructure, managed by different network managers as part of different administrative domains, or MAY run on an end host client, in which case they function as both senders and Repair Heads. Absent of any automatic tree configuration, it is assumed that the Infrastructure Repair Heads have relatively static configurations, which consist of a list of nearby possible Repair Heads. Senders and receivers, on the other hand, are transient entities, which typically only exist for the duration of a single data session. In addition to these core components, applications that use TRACK are expected to interface with other services that reside in other network entities, such as multicast address allocation, session advertisement, network management consoles, DHCP, DNS, overlay networking, application level multicast, and multicast key management.

6.1.2 Multicast Group Address

A multicast group address is a logical address that is used to address a set of TRACK nodes. It is RECOMMENDED to consist of a pair consisting of an IP multicast address and a UDP port number. In this case, it may optionally have a Time To Live (TTL) value, although this value MUST only be used for providing a global scope to a data session, and not for scoping of local retransmissions. Data multicast addresses are multicast group addresses.

TRACK MAY be used with an overlay multicast or application layer multicast system. In this case, a Multicast Group Address MAY have a different format. The TRACK PI is responsible for specifying the format of a multicast group address.

6.1.3 Data Session

A data session is the unit of reliable delivery of TRACK. It consists of a sequence of sequentially numbered data messages, which are sent by a single sender over a single data multicast address.

They are delivered reliably, with acknowledgements and retransmissions occurring over the control tree. A data Session ID uniquely identifies it. A given Data Session is received by a set of zero or more senders, and a set of zero or more Repair Heads. One or more data sessions MAY share the same data multicast address (although this is NOT RECOMMENDED). Each TRACK node can simultaneously participate in multiple data sessions. A receiver MUST join all the data multicast addresses and control trees corresponding to the data sessions it wishes to receive.

6.1.4 Data Channel

A data Session is multicast over a data channel. The data channel is responsible for efficiently delivering the data messages to the members of a data Session, and providing statistical reliability guarantees on this delivery.

TRACK is then responsible for providing application level, sender based reliability, by confirming delivery to all senders, and optionally retransmitting lost messages that did not get correctly delivered by the data channel.

6.1.5 Data Multicast Address

This is the multicast group address used by the data channel protocol, to efficiently deliver data messages to all receivers and Repair Heads. All data multicast addresses used by TRACK are assumed to be unidirectional and only support a single sender.

6.1.6 Control Tree or Session Tree

A control tree is a hierarchical communication path used to send control information from a set of receivers, through zero or more Repair Heads (RHs), to a sender. Information from lower nodes is aggregated as the information is relayed to higher nodes closer to the sender. Each data session uses a control tree. It is acceptable to have a degenerate control tree with no Repair Heads, which connects all of the receivers directly to the sender.

Each RH in the control tree uses a separate local control channel for communicating with its children. It is RECOMMENDED that each local control channel correspond to a separate multicast group address.

6.1.7 Local Control Channel

A local control channel is a unidirectional multicast path from a Repair Head or sender to its children. It uses a multicast group address for this communication.

[6.1.8](#) Host ID

With the widespread deployment of network address translators, creating a short globally unique ID for a host is a challenge. By default, TRACK uses a 48 bit long Host ID field, filled with the low-order 48 bits of the MD5 signature of the DNS name of the source. A TRACK PI, to match up with the goodput-ensuring protocol that TRACK PI uses as its Data Channel Protocol, MAY redefine the length and contents of this identifier.

[6.1.9](#) Data Session ID

A data Session ID is a globally unique identifier for a data session. It may either be selected by the data channel protocol (i.e. NORM) or by TRACK. By default, it is the combination of the Host ID for the sender, combined with the 16-bit port number used for the data session at the sender. This identifier is included in every TRACK message.

[6.1.10](#) Child ID

All members in a TRACK Data Session, besides the sender, are identified by the combination of their Host ID, and the port number with which they send IP packets to their parent.

[6.1.11](#) Message Sequence Numbers

A Message Sequence Number is a 32 bit number in the range from 1 through $2^{32} - 1$, which is used to specify the sequential order of a data message in a data stream. A sender node assigns consecutive Sequence Numbers to the data messages provided by the sender application. By default, zero is reserved to indicate that the data session has not yet started. A TRACK PI MAY redefine this. Message Sequence Numbers may wrap around, and so Sequence Number arithmetic MUST be used to compare any two Sequence Numbers.

[6.2](#) Basic Operation of the Protocol

For each data session, TRACK provides sequenced, reliable delivery of data from a single sender to up to tens of thousands of senders. A TRACK data session consists of a network that has exactly one sender node, zero or more receiver nodes and zero or more Repair Heads.

The figure below illustrates a TRACK Data Session with multiple Repair Heads.

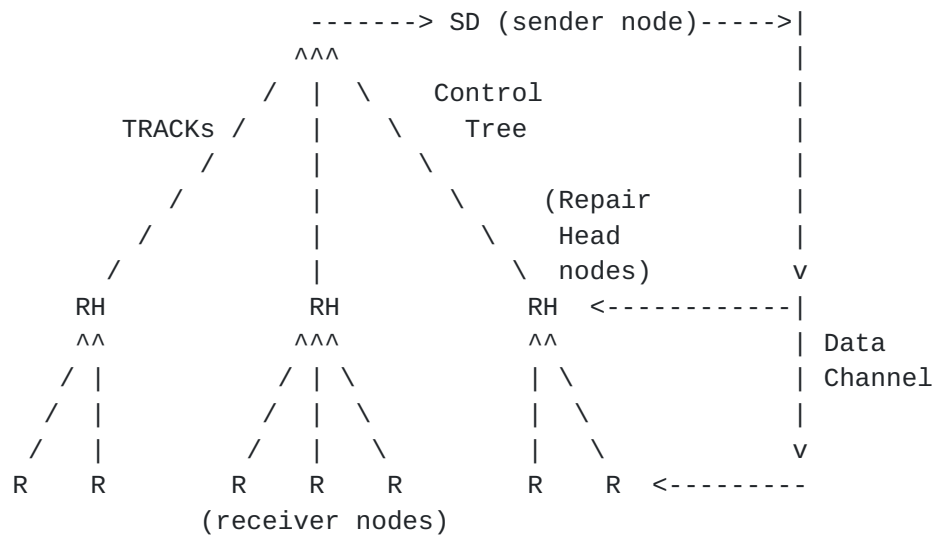


Figure 1. TRACK Session

Before a data session starts, a session advertisement MUST be received by all members of the Data Session, notifying them to join the group, and the appropriate configuration information for the data session. This MAY be provided directly by the application, by an external service, or by the TRACK PI.

A sender joins the control tree and a data channel protocol. It multicasts data messages on the data multicast address, using the data channel protocol. All of the nodes in the session subscribe to the data multicast address and join the data channel protocol.

There is no assumption of congruence between the topology of the data multicast address and the topology of the control tree.

A receiver joins the appropriate data channel, and the data multicast address used by that protocol, in order to receive data. A receiver periodically informs its parent about the messages that it has received by unicasting a TRACK message to the parent. It MAY also request retransmission of lost messages in this TRACK. Each parent node aggregates the TRACKs from its child nodes and (if it is not the sender) unicasts a single aggregated TRACK to its parent.

The sender and each Repair Head have a multicast local control channel to their children. This is used for transmitting Heartbeat messages that inform their child nodes that the parent node is still functioning. This channel is also used to perform local retransmission of lost Data messages to just these children. TRACK

MUST still provide correct operation even if multicast addresses are reused across multiple Data Sessions or multiple local control channels. It is NOT RECOMMENDED to use the same multicast address for multiple local control channels serving any given Data Session.

The communication path forms a loop from the sender to the receivers, through the Repair Heads back to the sender. Original data (ODATA), Retransmission (RDATA) and NullData messages regularly exercise the downward data direction. Heartbeat messages exercise the downward control direction. TRACK messages regularly exercise the Control Tree in the upward direction. This combination constantly checks that all of the nodes in the tree are still functioning correctly, and initiates fault recovery when required.

This hierarchical infrastructure allows TRACK to provide a number of functions in a scalable way. Application level confirmation of delivery and statistics aggregation both operate in a request-reply mode. A sender issues a request for application level confirmation or statistics reporting, and the receivers report back the appropriate information in their TRACK messages. This information is aggregated by the Repair Heads, and passed back up to the sender. Since TRACK messages are not delivered with the reliability of data messages, receivers and Repair Heads transmit this information redundantly.

TRACK also gathers control information that is useful for improving the performance of flow and congestion control algorithms, including scalable round trip time measurements.

7. Details: TRACK Functionality

7.1 Session Creation and Maintenance

7.1.1 Tree Configuration

Before a data session starts reliably delivering data, the tree for the data session needs to be created. This process binds each receiver to either a Repair Head or the sender, and binds the participating Repair Heads into a loop-free tree structure with the sender as the root of the tree. This process requires tree configuration knowledge, which can be provided with some combination of manual and/or automatic configuration. The algorithms for automatic tree configuration are part of the TREE BB [[RFCyyyy](#)]. They return to each node the address of the parent it should bind to, as well as zero or more backup parents to use if the primary parent fails.

7.1.2 Bind

In order to join a data session and bind to the tree, the following nodes need the following parameters.

A Repair Head requires the following parameters.

- Session: the unique identifier for the Data Session to join, received from the session advertisement algorithm specified in the PI.
- ParentAddress: the address and port of the parent node to which the node should connect, received from the TREE BB.
- UDPListenPort: the number of the port on which the node will listen for its childrens control messages. This parameter is configured by the application.
- RepairAddr: the multicast address, UDP port, and TTL on which this node sends control messages to its children. This parameter is configured by the application.

A sender requires the above parameters, except for the parentAddress. A sender requires the above parameters, except for the UDPListenPort and RepairAddr.

A Bind operation happens when a child wishes to join a parent in the distribution tree for a given Data Session. The receivers initiate the first Bind protocols to their parents, which then cause recursive binding by each parent, up to the sender. Each receiver sends a separate BindRequest message for each of the

streams that it would like to join. At the discretion of the PI, multiple BindRequest messages may be bundled together in a single message.

A node sends a BindRequest message to its automatically selected or manually configured parent node. The parent node sends either a BindConfirm message or a BindReject message. Reception of a BindConfirm message terminates the algorithm successfully, while receipt of a BindReject message causes the node to either retry the same parent or restart the Bind algorithm with its next parent candidate (depending on the BindReject reason code), or if it has none, to declare a REJECTED_BY_PARENT error. Once the node is accepted by a Repair Head, it informs the Tree BB using the setSN interface.

Reliability is achieved through the use of a standard request-response protocol. At the beginning of the algorithm, the child initializes TimeMaxBindResponse to the constant TIMEOUT_PARENT_RESPONSE and initializes NumBindResponseFailures to 0. Every time it sends a BindRequest message, it waits TimeMaxBindResponse for a response from the parent node. If no response is received, the node doubles its value for TimeMaxBindResponse, but limits TimeMaxBindResponse to be no larger than MAX_TIMEOUT_PARENT_RESPONSE. It also increments NumBindResponseFailures, and retransmits the BindRequest message. If NumBindResponseFailures reaches NUM_MAX_PARENT_ATTEMPTS, it reports a PARENT_UNREACHABLE error.

When a parent receives a BindRequest message, it first consults the TREE BB for approval (using the acceptchild Tree BB interface), for instance to ensure that accepting the BindRequest will not cause a loop in the tree. Then the parent checks to be sure that it does not have more than Maxchildren children already bound to it for this session. If it can accept the child, it sends back a BindConfirm message. Otherwise, it sends the node a BindReject message. Then the parent checks to see if it is already a member of this Data Session. If it is not yet a member of this session, it attempts to join the tree itself.

The BindConfirm message contains the lowest Sequence Number that the Repair Head has available. If this number is 0, then the Repair Head has all of the data available from the start of the session. Otherwise, the requesting node is attempting a late join, and can only use this Repair Head if late join was allowed by the PI. If late join is not allowed, the node may try another Repair Head, or give up.

Similarly, if a failure recovery occurs, when a node tries to bind to a new Repair Head, it must follow the same rules as for a late join. See Fault Recovery, below.

7.1.3 Unbind

A child may decide to leave a Data Session for the following reasons. 1) It detects that the Data Session is finished. 2) The application requests to leave the Data Session. 3) It is not able to keep up with the data rate of the Data Session. When any of these conditions occurs, it initiates an Unbind process.

An Unbind is, like the Bind function, a simple request-reply protocol. Unlike the Bind function, it only has a single response, UnbindConfirm. With this exception, the Unbind operation uses the same state variables and reliability algorithms as the Bind function.

When a child receives an UnbindConfirm message from its parent, it reports a LEFT_DATA_SESSION_GRACEFULLY event. If it does not receive this message after NUM_MAX_PARENT_ATTEMPTS, then it reports a LEFT_DATA_SESSION_ABNORMALLY event. Unbinds are reported to the Tree BB using the lostSN interface.

7.1.4 Eject

A parent may decide to remove one or more of its children from a data stream for the following reasons. 1) The parent needs to leave the group due to application reasons. 2) The Repair Head detects an unrecoverable failure with either its parent or the sender. 3) The parent detects that the child is not able to keep up with the speed of the data stream. 4) The parent is not able to handle the load of its children and needs some of them to move to another parent. In the first two cases, the parent needs to multicast the advertisement of the termination of one or more Data Sessions to all of its children. In the second two cases, it needs to send one or more unicast notifications to one or more of its children.

Consequently, an Eject can be done either with a repeated multicast advertisement message to all children, or a set of unicast request-reply messages to the subset of children that it needs to go to.

For the multicast version of Eject, the parent sends a multicast UnbindRequest message to all of its children for a given Data Session, on its Local Multicast Channel. It is only necessary to provide statistical reliability on this message, since children will detect the parents failure even if the message is not received.

Therefore, the UnbindRequest message is sent FAILURE_DETECTION_REDUNDANCY times.

For the unicast version of Eject, the parent sends a unicast UnbindRequest message to all of its children. Each of them responds with an EjectConfirm. Reliability is ensured through the same request-reply mechanism as the Bind operation.

Ejections are reported to the Tree BB using the removechild interface.

7.1.5 Fault Detection

There are three cases where fault detection is needed. 1) Detection (by a child) that a parent has failed. 2) Detection (by a parent) that a child has failed. 3) Detection (by either a Repair Head or sender) that a sender has failed.

In order to be scalable and efficient, fault detection is primarily accomplished by periodic keep-alive messages, combined with the existing TRACK messages. nodes expect to see keep-alive messages every set period of time. If more than a fixed number of periods go by, and no keep-alive messages of a given type are received, the node declares a preliminary failure. The detecting node may then ping the potentially failed node before declaring it failed, or it can just declare it failed.

Failures are detected through three keep-alive messages: Heartbeat, TRACK, and NullData. The Heartbeat message is multicast periodically from a parent to its children on its Local Control Channel. NullData messages are multicast by a sender on the Data Control Channel when it has no data to send. TRACK messages are generated periodically, even if no data is being sent to a Data Session.

Heartbeat messages are multicast every HeartbeatPeriod seconds, from a parent to its children. Every time that a parent sends a Retransmission message or a Heartbeat message (as well as at initialization time), it resets a timer for HeartbeatPeriod seconds. If the timer goes off, a Heartbeat is sent. The HeartbeatPeriod is dynamically computed as follows:

$$\text{interval} = \text{AckWindow} / \text{MessageRate}$$

$$\text{HeartbeatPeriod} = 2 * \text{interval}$$

Global configuration parameters ConstantHeartbeatPeriod and MinimumHeartbeatPeriod can be used to either set HeartbeatPeriod to a constant, or give HeartbeatPeriod a lower bound, globally.

Similarly, a NullData message is multicast by the sender to all data session members, every NULL_DATA_PERIOD. The NullData timer is set to NULL_DATA_PERIOD, and is reset every time that a Data or NullData message is sent by the sender.

The key parameter for failure detection is the global tree parameter FAILURE_DETECTION_REDUNDANCY. The higher the value for this parameter, the more keep-alive messages that must be missed before a failure is declared.

A major goal of failure detection is for children to detect parent failures fast enough that there is a high probability they can rejoin the stream at another parent, before flow control has advanced the buffer window to a point where the child can not recover all lost messages in the stream. In order to attempt to do this, children detect a failure of a parent if $\text{FAILURE_DETECTION_REDUNDANCY} * \text{HeartbeatPeriod}$ time goes by without any heartbeats. As part of buffer window advancement, all parents MAY choose to buffer all messages for a minimum of $\text{FAILURE_DETECTION_REDUNDANCY} * 2 * \text{HeartbeatPeriod}$ seconds, which gives children a period of time to find a new parent before the buffers are freed. Children report parent failures to the Tree BB using the lostSN interface.

A parent detects a preliminary failure of one of its children if it does not receive any TRACK messages from that child in $\text{FAILURE_DETECTION_REDUNDANCY} * \text{TrackTimeout}$ seconds (see discussion of how TrackTimeout is computed below). Because a failed child can slow down the groups progress, it is very important that a parent resolve the childs status quickly. Once a parent declares a preliminary failure of a child, it issues a set of up to $\text{FAILURE_DETECTION_REDUNDANCY}$ Heartbeat messages that are unicast (or multicast) to the failed sender(s). These messages are spaced apart by $2 * \text{LocalRTT}$, where LocalRTT is the round trip time that has been measured to the child in question (see below for description of how LocalRTT is measured). These Heartbeat messages contain a childrenList field that contains the children who are requested to send a TRACK immediately.

Whenever a child receives a Heartbeat message where the child is identified in the childrenList field, it immediately sends a TRACK to its parent. If a parent does not receive a TRACK message from a child after waiting a period of $2 * \text{LocalRTT}$ after the last Heartbeat message to that child, it declares the child failed, and removes it from the parents child membership list. It informs the Tree BB using the removechild interface.

A child or a Repair Head detects the failure of a sender if it does not receive a Data or NullData message from a sender in $\text{FAILURE_DETECTION_REDUNDANCY} * \text{NULL_DATA_PERIOD}$.

Note that the more senders there are in a tree, and the higher the loss rate, the larger $\text{FAILURE_DETECTION_REDUNDANCY}$ must be, in order to give the same probability that erroneous failures won't be declared.

7.1.6 Fault Notification

When a parent detects the failure of a child, it adds a failure notification field to the next $\text{TRANSMISSION_REDUNDANCY TRACK}$ messages that it sends up the tree. It sends this notification multiple times because TRACKs are not delivered reliably. A failure notification field includes the failure code, as well as a list of one or more failed nodes. Failure notifications are aggregated up the tree and delivered to the sender. A failure notification is not a definitive report of a node failure, as the child may have detected a communication failure with its parent and moved to a different Repair Head.

7.1.7 Fault Recovery

The Fault Recovery algorithms require a list of one or more addresses of alternate parents that can be bound to, and that still provide loop free operation.

If a child detects the failure of its parent, it then re-runs the Bind operation to a new parent candidate, in order to rejoin the tree. A node may perform a late join, i.e. binding with a Repair Head which cannot provide all the necessary repair data, only if allowed by the PI.

7.1.8 Distributed Membership.

Each Repair Head is responsible for maintaining a set of state variables on the status of its children. Unlike the Generic Router Assist, this is hard state, that only is removed when a child leaves that Repair Head gracefully, or after the Repair Head detects that a child has failed. These variables MUST include, but are not necessarily limited to, the following:

- childID. This is the two-byte identifier assigned to the child by the Repair Head. This uniquely identifies this child to this Repair Head, but has no meaning outside that scope.
- GlobalchildIdentifier. This is the globally unique identifier for this child.
- childRTT. This is the weighted average of the local RTT to child.

- LastTRACK. This is the contents of the last TRACK message sent from this child, if any, not including options.
- LastApplicationLevelConfirmation. This is the content of the last Application Level Confirmation sent from this child, if any.
- Last Statistics. This is the contents of the last Statistics message sent from this child, if any.
- ChildLiveness. This is a set of variables that keep track of the liveness of each child. This includes the last time a TRACK message was received from this child, as well as the number of Heartbeat messages that have been directed at it, and the time at which the last Heartbeat message was sent to the child. Please see Fault Detection, above, for more details.

7.2 Data Sessions

7.2.1 Data Transmission and Retransmission

Data is multicast by a sender on the Data Multicast Address via the Data Channel Protocol. The Data Channel Protocol is responsible for taking care of as many retransmissions as possible, and for ensuring the goodput of the Data Session. TRACK is then responsible for providing OPTIONAL flow control and application level reliability. The mechanics of an application level confirmation of delivery are handled by TRACK, including keeping track of the distributed membership list of receivers and aggregating acknowledgements up the control tree. Please see below for more details on flow control and application level confirmation.

A common scenario for handling recovery of lost messages is to allow the data channel protocol to provide statistical reliability, and then allow TRACK to provide retransmissions for more persistent failure cases, such as if a sender is not able to receive any data messages for a few minutes.

Retransmissions of data messages may be multicast by the sender on the data multicast address or be multicast on a local control channel by a Repair Head.

A Repair Head joins all of the Data Multicast Addresses that any of its descendants have joined. A Repair Head is responsible for receiving and buffering all data messages using the reliability semantics configured for a stream. As a simple to implement option, a Repair Head MAY also function as a sender, and pass these data messages to an attached application.

For additional fault tolerance, a sender MAY subscribe to the multicast address associated with the Local Control Channel of one or more Repair Heads in addition to the multicast address of its

parent. In this case it does not bind to this Repair Head or sender, but will process Retransmission messages sent to this address. If the receivers Repair Head fails and it transfers to another Repair Head, this minimizes the number of data messages it needs to recover after binding to the new Repair Head.

7.2.2 Local Retransmission

If a Repair Head or sender determines from its child nodes TRACK messages that a Data message was missed, the Repair Head retransmits the Data message. The Repair Head or sender multicasts the Retransmission message on its multicast Local Control Channel. In the event that a Repair Head receives a retransmission and knows that its children need this repair, it re-multicasts the retransmission to its children.

The scope of retransmission (the multicast TTL) is considered part of the Control Channels multicast address, and is derived during tree configuration.

A Repair Head maintains the following state for each of its children, for the purpose of providing repair service to the local group:

- HighestConsecutivelyReceived. A Sequence Number indicating all Data messages up to this number (inclusive) that have been received by a given child.
- MissingMessages. A data structure to keep track of the reception status of the Data messages with Sequence Number higher than HighestConsecutivelyReceived.

The minimum HighestConsecutivelyReceived value of all its children is kept as the variable LocalStable.

A Repair Head also maintains a retransmission buffer. The size of the retransmission buffer MUST be greater than the maximum value of a sender transmission window. The retransmission buffer MUST keep all the data messages received by the Repair Head with Sequence Number higher than LocalStable, optionally some messages with Sequence Number lower than LocalStable if there is room (beyond the maximum value of senders transmission window). The latter messages are kept in the retransmission buffer in case a sender from another group losses its parent and needs to join this group.

As TRACK messages are received, the Repair Head updates the above state variables.

To perform local repair, a Repair Head implements a retransmission queue with memory. Each lost message is entered into the retransmission queue in increasing order according to its Sequence Number. If the same data message has already been retransmitted recently (recognized due to the queues memory) it is delayed by the local group RTT (see roundtrip time measurement) before retransmission.

Retransmissions MAY NOT be sent at a faster rate than the current TransmissionRate advertised by the sender.

7.2.3 Flow and Rate Control

TRACK offers the ability to limit the rate of Data traffic, through both flow control and rate limits.

When a sender sends a TRACK to its parent, the HighestAllowed field provides information on the status of the senders flow control window. The value of HighestAllowed is computed as follows:

$$\text{HighestAllowed} = \text{seqnum} + \text{senderWindow}$$

Where seqnum is the highest Sequence Number of consecutively received data messages at the sender. The size of the senderWindow may either be based on a parameter local to the sender or be a global parameter.

If flow control is enabled for a given Data Session, then a sender MUST NOT send any Data messages to the Data Channel Protocol that are higher than the current value for HighestAllowed that it has. On startup, HighestAllowed is initialized to senderWindow.

In addition, the sender application MAY provide minimum and maximum rate limits. Unless overridden by the Data Channel Protocol, a sender will not offer Data messages to the Data Channel Protocol at lower than MinimumDataRate (except possibly during short periods of time when certain slow senders are being ejected), or higher than MaximumDataRate. If a sender is not able to keep up with the minimum rate for a period of time, it SHOULD leave the group promptly. senders that leave the group MAY attempt to rejoin the group at a later time, but SHOULD NOT attempt an immediate reconnection.

7.2.4 Reliability Window

The sender and each Repair Head maintain a window of messages for possible retransmission. As messages are acknowledged by all of its children, they are released from the parents retransmission buffer, as described in 4.2.2. In addition, there are two global

parameters that can affect when a parent releases a data message from the retransmission buffer -- MinHoldTime, and MaxHoldTime.

MinHoldTime specifies a minimum length of time a message must be held for retransmission from when it was received. This parameter is useful to handle scenarios where one or more children have been disconnected from their parent, and have to reconnect to another. If, for example, MinHoldTime is set to `FAILURE_DETECTION_REDUNDANCY * 2 * ConstantHeartbeatPeriod`, then there is a high likelihood that any child will be able to recover any lost messages after reconnecting to another parent.

The sender continually advertises to the members of the Data Session both edges of its retransmission window. The higher value is the SeqNum field in each Data or NullData message, which specifies the highest Sequence Number of any data message sent. The trailing edge of the window is advertised in the HighestReleased field. This specifies the largest Sequence Number of any message sent that has subsequently been released from the sender retransmission window. If both values are the same then the window is presently empty. Zero is not a legitimate value for a data Sequence Number, so if either field has a value of zero, then no messages have yet reached that state. All Sequence Number fields use Sequence Number arithmetic so that a Data Session can continue after exhausting the Sequence Number space.

When a member of a Data Session receives an advertisement of a new HighestReleased value, it stores this, and is no longer allowed to ask for retransmission for any messages up to and including the HighestReleased value. If it has any outstanding missing messages that are less than or equal to HighestReleased, it MAY move forward and continue delivering the next data messages in the stream. It also SHOULD report an error for the messages that are no longer recoverable.

MaxHoldTime specifies the maximum length of time a message may be held for retransmission. This parameter is set at the sender which uses it to set the HighestReleased field in data message headers. This is particularly useful for real-time, semi-reliable streams such as live video, where retransmissions are only useful for up to a few seconds. When combined with Unordered delivery semantics, and application-level jitter control at the senders, this provides Time Bounded Reliability. MaxHoldTime MUST always be larger than MinHoldTime.

7.2.5 Ordering Semantics

TRACK offers two flavors of ordering semantics: Ordered or Unordered. One of these is selected on a per session basis as part of the Session Configuration Parameters.

Unordered service provides a reliable stream of messages, without duplicates, and delivers them to the application in the order received. This allows the lowest latency delivery for time sensitive applications. It may also be used by applications that wish to provide its own jitter control.

Ordered service provides TCP semantics on delivery. All messages are delivered in the order sent, without duplicates.

7.2.6 Retransmission Requests.

A sender detects that it has missed one or more Data messages by gaps in the sequence numbers of received messages. Each sender keeps track of HighestSequenceNumber, the highest sequence number known of for a Data Session, as observed from Data, RData, and NullData messages. Any sequence numbers between HighestReleased and HighestSequenceNumber that have not been received are assumed to be missing.

When a sender detects missing messages it MAY send off a request for retransmission, if local retransmission is enabled. It does this by sending a Retransmission Request message. The timing of this request is described below.

7.2.7 End Of Stream.

When an application signals that a Data Session is complete, the sender advertises this to its children by setting the End of Session option on the last Data Message in the Data Session, as well as all subsequent retransmissions of that Data Message, and all subsequent Null Data messages.

The sender SHOULD NOT leave the Data Session until it has a report from the TRACK reports that all group members have left the Data Session, or it has waited a period of at least `FAILURE_DETECTION_REDUNDANCY * TrackTimeout` seconds.

7.3 Control Traffic Generation and Aggregation.

One of the largest challenges for scalable reliable multicast protocols has been that of controlling the potential explosion of control traffic. There is a fundamental tradeoff between the latency with which losses can be detected and repaired, and the amount of control traffic generated by the protocol.

TRACK messages are the primary form of control traffic in this BB. They are sent from senders and Repair Heads to their parents.

TRACK messages may be sent for the following purposes:

- to request retransmission of messages
- to advance the senders transmission window for flow control purposes
- to deliver application level confirmation of data reception
- to propagate other relevant feedback information up through the session (such as RTT and loss reports, for congestion control)

7.3.1 TRACK Generation with the Rotating TRACK Algorithm

Each receiver sends a TRACK message to its parent once per AckWindow of data messages received. A sender uses an offset from the boundary of each AckWindow to send its TRACK, in order to reduce burstiness of control traffic at the parents. Each parent has a maximum number of children, Maxchildren. When a child binds to the parent, the parent assigns a locally unique childID to that child, between 0 and Maxchildren-1.

Each child in a tree generates a TRACK message at least once every AckWindow of data messages, when the most recent data messages Sequence Number, modulo AckWindow, is equal to MemberID. If the message that would have triggered a given TRACK for a given node is missed, the node will generate the TRACK as soon as it learns that it has missed the message, typically through receipt of a higher numbered data message.

Together, AckWindow and Maxchildren determine the maximum ratio of control messages to data messages seen by each parent, given a constant load of data messages. In each data message, the sender advertises the current MessageRate (measured in messages per second) it is sending data at. This rate is generated by the congestion control algorithms in use at the sender.

At the time a node sends a regular TRACK, it also computes a TRACKTimeout value:

$$\text{interval} = \text{AckWindow} / \text{MessageRate}$$
$$\text{TRACKTimeout} = 2 * \text{interval}$$

If no TRACKs are sent within TRACKTimeout interval, a TRACK is generated, and TRACKTimeout is increased by a factor of 2, up to a value of MAX_TRACK_TIMEOUT.

This timer mechanism is used by a sender to ensure timely repair of lost messages and regular feedback propagation up the tree even when the sender is not sending data continuously. This mechanism complements the AckWindow-based regular TRACK generation mechanism.

7.3.2 TRACK Aggregation

There are many reasons for providing feedback from all the receivers to the sender in an aggregated form. The major ones are listed below:

- 1) End-to-end delivery confirmation. This confirmation tells the sender that all the senders (in the entire tree) have received data messages up to a certain Sequence Number. This is carried in an Application Level Confirmation message.
- 2) Flow control. The aggregated information is carried in the field HighestAllowed. It tells the sender the highest Sequence Number that all the senders (in the entire tree) are prepared to receive.
- 3) Congestion control feedback. Information about the state of the tree can be passed up to help control the congestion control algorithms for the group.
- 4) Counting current membership in the group. This information is carried in the field SubTreeCount. This lets the sender know the number of senders currently connected to the repair tree.
- 5) Measuring the round-trip time from the sender to the "worst" sender.

A Repair Head maintains state for each child. Each time a TRACK (from a child) is received, the corresponding states for that child are updated based on the information in the TRACK message. When a Repair Head sends a TRACK message to its parent, the following fields of its TRACK message are derived from the aggregation of the corresponding states for its children. The following rules describe how the aggregation is performed:

- WorstLossRate. Take the maximum value of the WorstLossRate from all children.
- SubTreeCount. Take the sum of the SubTreeCount from all children.
- HighestAllowed. Take the minimum of the HighestAllowed value from all children.

- WorstEdgeThroughput. Take the minimum value of the WorstEdgeThroughput field from all children.
- UnicastCost. Take the sum of the UnicastCost from all children.
- MulticastCost. Take the sum of MulticastCost from all children.
- senderDallyTime: take the minimum value, for all of the children, of (childs reported senderDallyTime + childs local dally time).
- FailureCount: take the sum of the FailureCount for all children.
- FailureList: concatenate the FailureList fields for all children, up to a maximum list size of MaxFailureListSize.

Note, the senderTimeStamp, parentTimestamp, and parentDallyTime fields are not aggregated. The sender will derive the roundtrip time to the worst sender by doing its local aggregation for senderDallyTime.

Application level confirmations (ALCs) are handled as follows. For a set of ALC requests from receivers, the ones with the highest value for HighConfirmationSequenceNumber are considered, and all others are discarded.

For the ConfirmationStatus field, the following rules apply. Note that ConfirmationStatus of SomesendersAcknowledge can correspond to a ConfirmationCount of zero.

```

If all children report AllsendersAcknowledge Then
    ConfirmationStatus = AllsendersAcknowledge
Else If at least one child reports (ListOfFailures OR
    FailuresExceedMaximumListSize) Then
    If the count of all reported failures >
        MaximumFailureListSize Then
        ConfirmationStatus = FailuresExceedMaximumListSize
    Else
        ConfirmationStatus = ListOfFailures
Else
    ConfirmationStatus = SomesendersAcknowledge

```

The ConfirmationCount field is equal to the sum of the ConfirmationCount for the aggregated ALC reports of all children. The PendingCount field is equal to the sum of the PendingCount fields of all children. The FailureList field is the concatenation of the FailureList fields of all aggregated ALC reports of all children, up to a maximum length of MaximumFailureListSize.

In addition to these fields with fixed aggregation rules, TRACK supports a set of user defined aggregation statistics. These statistics are self-describing in terms of their data type and aggregation method. Statistics reports are numbered, and only the most recent statistics report request is aggregated to the sender. Statistics are aggregated over the set of child statistics reports

that have been received with that number. Aggregation methods include minimum, maximum, sum, product, and concatenation.

7.3.3 Statistics Reporting.

A sender can request a list of aggregated statistics from all senders in the group. There are a set of predefined statistics, such as loss rate and average throughput. There is also the capacity to request a set of other TRACK statistics, as well as application-defined statistics.

The format of each statistic is self-describing, both in terms of data type, size, and aggregation method. A sender reliably sends out a statistics request by attaching it as an option to a Data message. When a sender gets a request for a statistic, it fills in the data fields and forwards it up the tree in the next TRACK message. Since TRACKs are not reliable, multiple copies are sent in a total of NumReplies consecutive TRACK messages from each sender. Each statistics report is aggregated according to the method described in the statistic and the result is delivered to the sender.

Most aggregation options have fixed length no matter how many senders there are. The one exception is concatenation, which creates a list of values from some or all senders, up to a length of MaximumStatisticsListSize entries. It is NOT RECOMMENDED to use this to create group-wide lists, unless the group size is carefully controlled.

7.4 Application Level Confirmed Delivery.

Flow control and the reliability window are concerned with goodput, of delivering data with a high probability that it is delivered at all senders. However, neither mechanism provides explicit confirmation to the sender as to the list of recipients for each message. Application level confirmed delivery allows applications to determine the set of applications that have received a set of data messages.

There are three primary factors that determine the reliability semantics of a message: the senders knowledge of the sender list, the application level actions that must be performed in order to consider a message delivered, and the response to persistent failure conditions at senders. For example, an extremely strong distributed guarantee would consist of the following. First, the full sender membership list is known at the sender, and verified to make sure no receivers have left the group. Second, the application at each receiver must write the data to persistent

store before it can be acknowledged. Third, receivers are given a very long period of time to recover all lost data messages, before they are ejected from the data session. In the meantime, transmission of data messages is flow controlled by the slowest receivers.

A weaker form of reliability would include the following. First, that the sender gets a count of receivers, and otherwise depends on the distributed group membership algorithms to maintain the membership list. Second, that data messages are considered reliably delivered as soon as the application receives the data from TRACK. Third, that retransmissions are limited to only 30 seconds, and receivers must choose to leave the Data Session or continue with missing data messages, if a failure takes longer than this period to recover from.

TRACK provides the functionality to easily implement a wide range of application level confirmation semantics, based on how these three items are configured. It is the applications responsibility to then select the configurations it desires for a given data session.

The primary mechanism for application level confirmation (ALC) of delivery is the ALC report. To check for ALC of delivery, a sender issues an Application Level Confirmation Request, by attaching this message as an option to a Data message, and reliably transmitting it to all senders. Each ALC Request includes a specified level of reliability, a reply redundancy factor, and the range of Data message sequence numbers that the ALC Confirmation covers.

When a sender gets an ALC Request, it checks to see if the application has delivered the specified range of Data Messages, including both the Low Confirmation Sequence Number and the High Confirmation Sequence Number. When it sends the next TRACK out, it sets the ConfirmationStatus field to either SomesendersAcknowledge if it is still pending confirmation, AllsendersAcknowledge if it has application level confirmation, ListOfFailures if it has a failure and MaximumFailureListSize > 0, or FailuresExceedsMaximumListSize otherwise. It also sets the ConfirmCount to 1 if it has a confirmation, and PendingCount to 1 if it is still pending. If the Immediate ACK bit is set in the ALC Request, the sender generates an ACK immediately.

One example of how an application can implicitly signal confirmation of delivery is through the freeing of buffers passed to it by the transport. The API could specify that whenever an application has freed up a buffer containing one or more data messages, then these messages are considered acknowledged by the

application. Alternatively, the application could be required to explicitly acknowledge each message.

7.5 Distributed RTT Calculations.

This TRACK BB provides two algorithms for distributed RTT calculations: LocalRTT measurements and senderRTT measurements. LocalRTT measurements are only between a parent and its children. senderRTT measurements are end-to-end RTT measurements, measuring the RTT to the worst sender as selected by the congestion control algorithms.

The senderRTT is useful for congestion control. It can be used to set the data rate based on the TCP response function, which is being proposed for the congestion control building blocks.

The LocalRTT can be used to (a) quickly detect faulty children (as described under fault detection) or (b) avoid sending unnecessary retransmissions (as described in the local repair algorithm).

In the case of LocalRTT measurements, a parent initiates measurement by including a parentTimestamp field in a Heartbeat message sent to its children. When a child receives a Heartbeat message with this field set, it notes the time of receipt using its local system clock, and stores this with the message as HeartbeatReceiveTime. When the child next generates a TRACK, just before sending it, it measures its system clock again as TRACKSendTime, and calculates the LocalDallyTime.

$$\text{LocalDallyTime} = \text{TRACKSendTime} - \text{HeartbeatReceiveTime}.$$

The child includes this value, along with the parentTimestamp field, as fields in the next TRACK message sent. Every heartbeat message that is multicast to all children SHOULD include a parentTimestamp field.

The senderRTT algorithm is similar. A sender initiates the process by including a senderTimestamp field in a data message. When a sender gets a message with this field set, it keeps track of the DataReceiveTime for that message, and when it generates the next TRACK message, includes the senderTimestamp and senderDallyTime value. These values are aggregated by Repair Heads.

Each node only keeps track of the most recent value for {senderTimestamp, DataReceiveTime} and {parentTimestamp, HeartbeatReceiveTime}, replacing any older values any time that a new message is received with these values set. As long as it has non-zero values to report, each node sends up both a

{senderTimestamp, senderDallyTime} and a {parentTimestamp, LocalDallyTime} set of fields in each TRACK message generated.

Unless redefined by the TRACK PI, these RTT measurements are averaged using an exponentially weighted moving average, where the first RTT measurement, `RTT_measurement`, initializes the average `RTT_average`, and then each successive measurement is averaged in according to the following formula. The RECOMMENDED value for alpha is 1/8.

$$\text{RTT_average} = \text{RTT_measurement} * \alpha + \text{RTT_average} (1-\alpha)$$

[7.6](#) **SNMP Support**

The Repair Heads and the sender are designed to interact with SNMP management tools. This allows network managers to easily monitor and control the sessions being transmitted. All TRACK nodes MAY have SNMP MIBs defined in a separate document. SNMP support is OPTIONAL for sender nodes, but is RECOMMENDED for all other nodes.

[7.7](#) **Late Join Semantics**

TRACK offers three flavors of late join support:

a) No Recovery

A sender binds to a Repair Head after the session has started and agrees to the reliability service starting from the Sequence Number in the current data message received from the sender.

b) Continuation

This semantic is used when a sender has lost its Repair Head and needs to re-affiliate. In this case, the sender must indicate the oldest Sequence Number it needs to repair in order to continue the reliability service it had from the previous Repair Head. The binding occurs if this is possible.

c) No Late Join

For some applications, it is important that a sender receives either all data or no data (e.g. software distribution). In this case option (c) is used.

These are specified by the `LateJoinSemantics` session parameter, and enforced by a parent when a child attempts to bind to it.

8. TRACK Message Types

The following table summarizes the messages and their fields used by the TRACK BB. All messages contain the session identifier.

Table 1. TRACK Messages

Message	From	To	Mcast?	Fields
BindRequest	child	parent	no	Scope, Level, Role, Rejoin BindSequenceNumber, SubTreeCount
BindConfirm	parent	child	no	RepairAddr, BindSequenceNumber LowestAvailableRepair Level, childIndex, Role
BindReject	parent	child	no	Reason, BindSequenceNumber
UnbindRequest	child	parent	no	Reason, childIndex
UnbindConfirm	parent	child	no	
EjectRequest	parent	child	either	Reason, Alternateparent
EjectConfirm	child	parent	no	
Heartbeat	parent	child	either	Level, parentTimestamp childrenList, SeqNum HighestReleased
NullData, OData	sender	all	yes	senderTimeStamp, DataLength HighestReleased, SeqNum EndOfStream, TransmissionRate
Rdata	parent	child	yes	senderTimeStamp, DataLength HighestReleased, SeqNum EndOfStream, TransmissionRate
Track	child	parent	no	BitMask, SubTreeCount Slowest, HighestAllowed parentThere, parentTimeStamp parentDallyTime, senderTimeStamp senderDallyTime, CongestionControl, FailureList
StatsRequest	sender	sender	yes	Immediate, StatsSeqNum NumReplies, StatsList
StatsReply	child	parent	yes	StatsSeqNum, StatsList

The various fields of the messages are described as follows:

- BindSequenceNumber: This is a monotonically increasing sequence number for each bind request from a given sender for a given Data Session.
- Scope: an integer to indicate how far a repair message travels. This is optional.
- Rejoin: a flag as to whether this sender was previously a member of this Data Session.
- Level: an integer that indicates the level in the repair tree. This value is used to keep loops in the tree from forming, in addition to indicating the distance from the sender. Any changes in a nodes level are passed down to the Tree BB using the treeLevelUpdate interface.
- Role: This indicates if the bind requestor is a sender or Repair Head.
- SubTreeCount: This is an integer indicating the current number of senders below the node.
- RepairAddr: This field in the BindConfirm message is used to tell the sender which multicast address the Repair Head will be sending retransmissions on. If this field is null, then the sender should expect retransmissions to be sent on the senders data multicast address.
- Alternateparent: This is an optional field that specifies another parent a child may attempt to bind to.
- SeqNum: an integer indicating the Sequence Number of a data message within a given Data Session. For a Heartbeat, it is the highest sequence number the parent knows about.
- ChildIndex: This is an integer the Repair Head assigns to a particular child. The child sender uses this value to implement the rotating TRACK Generation algorithm.
- LowestRepairAvailable: This is the lowest sequence number that a Repair Head will provide repairs for.
- Reason: a code indicating the reason for the BindReject, UnbindRequest, or EjectRequest message.
- ParentTimestamp: This field is included in Heartbeat messages to

signal the need to do a local RTT measurement from a parent. It is the time when the parent sent the message.

- childrenList: This field contains the identifiers for a list of children. As part of the keepalive message, this field together with the SeqNum field is used to urge those listed senders to send a TRACK (for the provided SeqNum). The Repair Head sending this must have been missing the regular TRACKs from these children for an extended period of time.
- senderTimestamp: This field is included in Data messages to signal the need to do a roundtrip time measurement from the sender, through the tree, and back to the sender. It is the time (measured by the senders local clock) when it sent the message.
- ApplicationSynch: a Sequence Number signaling a request for confirmed delivery by the application.
- EndOfStream: indicates that this message is the end of the data for this session.
- TransmissionRate: This field is used by the sender to tell the senders its sending rate, in messages per second. It is part of the data or nulldata messages.
- HighestReleased: This field contains a Sequence Number, corresponding to the trailing edge of the senders retransmission window. It is used (as part of the data, nulldata or retransmission headers) to inform the senders that they should no longer attempt to recover those messages with a smaller (or same) Sequence Number.
- HighestAllowed: a Sequence Number, used for flow control from the senders. It signals the highest Sequence Number the sender is allowed to send that will not overrun the senders buffer pools.
- BitMask: an array of 1s and 0s. Together with a Sequence Number it is used to indicate lost data messages. If the *i*th element is a 1, it indicates the message SeqNum+*i* is lost.
- Slowest: This field contains a field that characterizes the slowest sender in the subtree beneath (and including) the node sending the TRACK. This is used to provide information for the congestion control BB.
- SenderDallyTime: This field is associated with a senderTimestamp field. It contains the sum of the waiting time that should be subtracted from the RTT measurement at the sender.

- ParentDallyTime: This is the same as the senderDallyTime, but is associated with a parentTimestamp instead of a senderTimestamp.
- DataLength: This is the length of the Data payload.
- CongestionControl: This includes any additional congestion control variables for aggregation, such as WorstLossRate, WorstEdgeThroughput, UnicastCost, and MulticastCost.
- ApplicationConfirms: This is the SeqNum value for which delivery has been confirmed by all children at or below this parent.
- Failedchildren: This is a list of all children that have recently been dropped from the repair tree.
- Immediate: If set to 1, a sender should immediately send a TRACK on receipt of this packet.
- Reliability: The level of reliability required in order to consider the set of data packets reliably delivered.
- NumReplies: The number of consecutive TRACK messages that should be sent with this message attached
- SeqNumRange: The set of data messages that the ALC request applies to.
- ConfirmStatus: The acknowledgement status of the senders in the subtree up to the node that sends this message.
- ConfirmCount: The number of senders in the subtree up to the node that sends this message, that have acknowledged the ALC request.
- PendingCount: The number of senders in this subtree that are still pending in their decision as to acknowledging this ALC request.
- StatsSeqNum: The number of this request for statistics.
- StatsList: The list of statistics to be filled in by senders, and aggregated by the control tree.

9. Global Configuration Parameters

9.1 Configuration Variables

These are variables that control the session and are advertised to all participants. Some of them MAY be configured as constants.

- TimeMaxBindResponse: the time, in seconds, to wait for a response to a BindRequest. Initial value is TIMEOUT_PARENT_RESPONSE (recommended value is 3). Maximum value is MAX_TIMEOUT_PARENT_RESPONSE.
- Maxchildren: The maximum number of children a Repair Head is allowed to handle. Recommended value: 32.
- ConstantHeartbeatPeriod: Instead of dynamically calculating the HeartbeatPeriod, a constant period may be used instead. Recommended value: 3 seconds.
- MinimumHeartbeatPeriod: The minimum value for the dynamically calculated HeartbeatPeriod. Recommended value: 1 second.
- MinHoldTime: The minimum amount of time a Repair Head holds on to data messages.
- MaxHoldTime: The maximum amount of time a Repair Head holds on to data messages.
- AckWindow: The number of messages seen before a sender issues an acknowledgement. Recommended value: 32.
- LateJoinSemantics: The options available to a sender who wishes to join a Data Session that is already in progress.
- MaximumFailureListSize: The maximum number of entries that can be in a failure list. This MUST be small enough that the FailureList does not ever cause a TRACK to exceed the size of a maximum UDP packet. Recommended value: 800.
- MaximumStatisticsListSize: The maximum number of entries that can be in a statistics list. This MUST be small enough that the FailureList does not ever cause a TRACK to exceed the size of a maximum UDP packet. Recommended value: 100.
- MaximumDataRate: The maximum admission rate for data messages from the application to the Data Channel Protocol.
- MinimumDataRate: The minimum admission rate for data messages from the application to the Data Channel Protocol.

9.2 Constants

- NUM_MAX_PARENT_ATTEMPTS: The number of times to try to bind to a Repair Head before declaring a PARENT_UNREACHABLE error. Recommended value is 5.
- TIMEOUT_PARENT_RESPONSE: The minimum value, in seconds, between attempts to contact a parent. Recommended value is 1 second.
- MAX_TIMEOUT_PARENT_RESPONSE: The maximum value, in seconds, between attempts to contact a parent. Recommended value is 16.
- NULL_DATA_PERIOD: The time between transmission of NullData Messages. Recommended value is 1.
- FAILURE_DETECTION_REDUNDANCY: The number of times a message is sent without receiving a response before declaring an error. Recommended value is 3.
- MAX_TRACK_TIMEOUT: The maximum value for TRACKTimeout. Recommended value is 5 seconds.
- TRANSMISSION_REDUNDANCY: The number of times a failure notification is redundantly sent up the tree in a TRACK message. Recommended value is 3.

9.3 Reason Codes

- o BindReject reason codes
 - LOOP_DETECTED
 - MAX_CHILDREN_EXCEEDED
- o UnbindRequest reason codes
 - SESSION_DONE
 - APPLICATION_REQUEST
 - RECEIVER_TOO_SLOW
- o EjectRequest reason codes
 - PARENT_LEAVING
 - PARENT_FAILURE
 - CHILD_TOO_SLOW
 - PARENT_OVERLOADED

10. Requirements from other Building Blocks

This TRACK BB can be interfaced to any other BB or PI wishing to use a tree structure. To actually use this BB's features, the PI needs to include the messages described in this BB in its packets.

11. Security Considerations

Basically, this document is for informational and security issues are not applied. The following considerations are given just for information:

- a. The primary security requirement for a TRACK protocol is protection of the transport infrastructure. This is accomplished through the use of lightweight group authentication of the control and, optionally, the data messages sent to the group. These algorithms use IPsec and shared symmetric keys.
- b. For TRACK, it is recommended that there be one shared key for the Data Session and one for each Local Control Channel. These keys are distributed through a separate key manager component, which may be either centralized or distributed. Each member of the group is responsible for contacting the key manager, establishing a pair-wise security association with the key manager, and obtaining the appropriate keys.
- c. The exact algorithms for this BB are presently the subject of research within standardization within the IETF Multicast Security (MSEC) working group.

12. References

Normative:

- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels," [BCP 14](#), [RFC 2119](#), March 1997
- [RFC3048] Whetten, B., Vicisano, L., Kermode, R., Handley, M., Floyd, S. and M. Luby, "Reliable Multicast Transport Building Blocks for One-to-Many Bulk-Data Transfer," [RFC 3048](#), January 2001.
- [RFCyyyy] Chiu, D., Koh, S., Kadansky, M., Whetten, B. and G. Taskale, "Tree Auto-Configuration (TREE) Building Block for Reliable Multicast Transport," RFC yyyy, 2004.

Informative:

- [RFC3269] Kermode, R., Vicisano, L., "Author Guidelines for Reliable Multicast Transport (RMT) Building Blocks and Protocol Instantiation documents," [RFC 3269](#), April 2002.
- [RFC2887] Handley, M., Whetten, B., Kermode, R., Floyd, S., Vicisano, L., and Luby, M., "The Reliable Multicast Design Space for Bulk Data Transfer," [RFC 2887](#), August 2000.
- [RFC2357] Mankin, A., Romanow, A., Bradner, S. and V. Paxson, "IETF Criteria for Evaluating Reliable Multicast Transport and Application Protocols," [RFC 2357](#), June 1998.
- [RFC3450] Luby, M., Gemmell, J., Vicisano, L., Rizzo, L. and J. Crowcroft, "Asynchronous Layered Coding (ALC) Protocol Instantiation," [RFC 3450](#), December 2002.
- [RFC3451] Luby, M., Gemmell, J., Vicisano, L., Rizzo, L., Handley, M. and J. Crowcroft, "Layered Coding Transport (LCT) Building Block," [RFC 3451](#), December 2002.
- [RFC3452] Luby, M., Vicisano, L., Gemmell, J., Rizzo, L., Handley, M., and J. Crowcroft, "Forward Error Correction (FEC) Building Block," [RFC 3452](#), December 2002.
- [NORM-BB] Adamson, B., Bormann, C., Handley M., Macker J. "NACK-Oriented Reliable Multicast (NORM) Building Blocks," Internet Draft, December 2003.
- [NORM-PI] Adamson, B., Bormann, C., Handley M., Macker J. "NACK-Oriented Reliable Multicast Protocol (NORM)," Internet Draft, December 2003.

13. Acknowledgments

The authors would like to give special thanks to Sanjoy Paul, Joe Wesley and Juyoung Park for their valuable comments.

14. Author's Addresses

Brian Whetten
brian@whetten.net
2430 20th Street #D, Santa Monica, CA 90405

Dah Ming Chiu
dmchiu@ie.cuhk.edu.hk
Information Engineering Department,
The Chinese University of Hong Kong Shatin, N.T. Hong Kong

Miriam Kadansky
miriam.kadansky@sun.com
Sun Microsystems Laboratories 1 Network Drive
Burlington, MA 01803

Seok Joo Koh
sjkoh@etri.re.kr
Protocol Engineering Center,
ETRI, 161 Kajung-Dong Yusong-Gu, TAEJON, 305-350, KOREA

Gursel Taskale
gursel@tibco.com
TIBCO
3303 Hillview Ave. Palo Alto, CA. 94304-1213

Full Copyright Statement

Copyright (C) The Internet Society (2003). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS 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."

