            Protocol Independent Multicast - Dense Mode (PIM-DM):
                      Protocol Specification (Revised)

Status of This Memo

Copyright Notice

Abstract

   This document specifies Protocol Independent Multicast - Dense Mode
   (PIM-DM).  PIM-DM is a multicast routing protocol that uses the
   underlying unicast routing information base to flood multicast
   datagrams to all multicast routers.  Prune messages are used to
   prevent future messages from propagating to routers without group
   membership information.

Table of Contents

1.  Introduction

   This specification defines a multicast routing algorithm for
   multicast groups that are densely distributed across a network.  This
   protocol does not have a topology discovery mechanism often used by a
   unicast routing protocol.  It employs the same packet formats sparse
   mode PIM (PIM-SM) uses.  This protocol is called PIM - Dense Mode.
   The foundation of this design was largely built on Deering's early
   work on IP multicast routing [12].

2.  Terminology

   The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT",
   "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" are to
   be interpreted as described in RFC 2119 [11] and indicate requirement
   levels for compliant PIM-DM implementations.

2.1.  Definitions

   Multicast Routing Information Base (MRIB)
      This is the multicast topology table, which is typically derived
      from the unicast routing table, or from routing protocols such as
      MBGP that carry multicast-specific topology information.  PIM-DM
      uses the MRIB to make decisions regarding RPF interfaces.

   Tree Information Base (TIB)
      This is the collection of state maintained by a PIM router and
      created by receiving PIM messages and IGMP information from local
      hosts.  It essentially stores the state of all multicast
      distribution trees at that router.

   Reverse Path Forwarding (RPF)
      RPF is a multicast forwarding mode in which a data packet is
      accepted for forwarding only if it is received on an interface used
      to reach the source in unicast.

   Upstream Interface
      Interface toward the source of the datagram.  Also known as the RPF
      Interface.

   Downstream Interface
      All interfaces that are not the upstream interface, including the
      router itself.

   (S,G) Pair
      Source S and destination group G associated with an IP packet.

2.2.  Pseudocode Notation

   We use set notation in several places in this specification.

   A (+) B
      is the union of two sets, A and B.

   A (-) B
      are the elements of set A that are not in set B.

   NULL
      is the empty set or list.

   Note that operations MUST be conducted in the order specified.  This
   is due to the fact that (-) is not a true difference operator,
   because B is not necessarily a subset of A.  That is, A (+) B (-) C =
   A (-) C (+) B is not a true statement unless C is a subset of both A
   and B.

   In addition, we use C-like syntax:

      =    denotes assignment of a variable.
      ==   denotes a comparison for equality.
      !=   denotes a comparison for inequality.

   Braces { and } are used for grouping.

3.  PIM-DM Protocol Overview

   This section provides an overview of PIM-DM behavior.  It is intended
   as an introduction to how PIM-DM works and is NOT definitive.  For
   the definitive specification, see Section 4, Protocol Specification.

   PIM-DM assumes that when a source starts sending, all downstream
   systems want to receive multicast datagrams.  Initially, multicast
   datagrams are flooded to all areas of the network.  PIM-DM uses RPF
   to prevent looping of multicast datagrams while flooding.  If some
   areas of the network do not have group members, PIM-DM will prune off
   the forwarding branch by instantiating prune state.

   Prune state has a finite lifetime.  When that lifetime expires, data
   will again be forwarded down the previously pruned branch.

   Prune state is associated with an (S,G) pair.  When a new member for
   a group G appears in a pruned area, a router can "graft" toward the
   source S for the group, thereby turning the pruned branch back into a
   forwarding branch.

The broadcast of datagrams followed by pruning of unwanted branches
is often referred to as a flood and prune cycle and is typical of
dense mode protocols.

To minimize repeated flooding of datagrams and subsequent pruning
associated with a particular (S,G) pair, PIM-DM uses a state refresh
message.  This message is sent by the router(s) directly connected to
the source and is propagated throughout the network.  When received
by a router on its RPF interface, the state refresh message causes an
existing prune state to be refreshed.

Compared with multicast routing protocols with built-in topology
discovery mechanisms (e.g., DVMRP [13]), PIM-DM has a simplified
design and is not hard-wired into a specific topology discovery
protocol.  However, this simplification does incur more overhead by
causing flooding and pruning to occur on some links that could be
avoided if sufficient topology information were available; i.e., to
decide whether an interface leads to any downstream members of a
particular group.  Additional overhead is chosen in favor of the
simplification and flexibility gained by not depending on a specific
topology discovery protocol.

PIM-DM differs from PIM-SM in two essential ways: 1) There are no
periodic joins transmitted, only explicitly triggered prunes and
grafts.  2) There is no Rendezvous Point (RP).  This is particularly
important in networks that cannot tolerate a single point of failure.
(An RP is the root of a shared multicast distribution tree.  For more
details, see [4]).

4.  Protocol Specification

The specification of PIM-DM is broken into several parts:

* Section 4.1 details the protocol state stored.
* Section 4.2 specifies the data packet forwarding rules.
* Section 4.3 specifies generation and processing of Hello messages.
* Section 4.4 specifies the Join, Prune, and Graft generation and
            processing rules.
* Section 4.5 specifies the State Refresh generation and forwarding
            rules.
* Section 4.6 specifies the Assert generation and processing rules.
* Section 4.7 gives details on PIM-DM Packet Formats.
* Section 4.8 summarizes PIM-DM timers and their defaults.

4.1.  PIM Protocol State

   This section specifies all the protocol states that a PIM-DM
   implementation should maintain to function correctly.  We term this
   state the Tree Information Base or TIB, as it holds the state of all
   the multicast distribution trees at this router.  In this
   specification, we define PIM-DM mechanisms in terms of the TIB.
   However, only a very simple implementation would actually implement
   packet forwarding operations in terms of this state.  Most
   implementations will use this state to build a multicast forwarding
   table, which would then be updated when the relevant state in the TIB
   changes.

   Unlike PIM-SM, PIM-DM does not maintain a keepalive timer associated
   with each (S,G) route.  Within PIM-DM, route and state information
   associated with an (S,G) entry MUST be maintained as long as any
   timer associated with that (S,G) entry is active.  When no timer
   associated with an (S,G) entry is active, all information concerning
   that (S,G) route may be discarded.

   Although we precisely specify the state to be kept, this does not
   mean that an implementation of PIM-DM has to hold the state in this
   form.  This is actually an abstract state definition, which is needed
   in order to specify the router's behavior.  A PIM-DM implementation
   is free to hold whatever internal state it requires and will still be
   conformant with this specification as long as it results in the same
   externally visible protocol behavior as an abstract router that holds
   the following state.

4.1.1.  General Purpose State

   A router stores the following non-group-specific state:

   For each interface:
     Hello Timer (HT)
     State Refresh Capable
     LAN Delay Enabled
     Propagation Delay (PD)
     Override Interval (OI)

     Neighbor State:
       For each neighbor:
         Information from neighbor's Hello
         Neighbor's Gen ID.
         Neighbor's LAN Prune Delay
         Neighbor's Override Interval
         Neighbor's State Refresh Capability
         Neighbor Liveness Timer (NLT)

4.1.2.  (S,G) State

   For every source/group pair (S,G), a router stores the following
   state:

   (S,G) state:
     For each interface:
       Local Membership:
         State: One of {"NoInfo", "Include"}

       PIM (S,G) Prune State:
         State: One of {"NoInfo" (NI), "Pruned" (P), "PrunePending"
                        (PP)}
                     Prune Pending Timer (PPT)
                     Prune Timer (PT)

       (S,G) Assert Winner State:
         State: One of {"NoInfo" (NI), "I lost Assert" (L), "I won
                        Assert" (W)}
         Assert Timer (AT)
         Assert winner's IP Address
         Assert winner's Assert Metric

     Upstream interface-specific:
       Graft/Prune State:
         State: One of {"NoInfo" (NI), "Pruned" (P), "Forwarding" (F),
                        "AckPending" (AP) }
         GraftRetry Timer (GRT)
         Override Timer (OT)
         Prune Limit Timer (PLT)

       Originator State:
         Source Active Timer (SAT)
         State Refresh Timer (SRT)

4.1.3.  State Summarization Macros

   Using the state defined above, the following "macros" are defined and
   will be used in the descriptions of the state machines and pseudocode
   in the following sections.

   The most important macros are those defining the outgoing interface
   list (or "olist") for the relevant state.

   immediate_olist(S,G) = pim_nbrs (-) prunes(S,G) (+)
                          (pim_include(*,G) (-) pim_exclude(S,G) ) (+)
                          pim_include(S,G) (-) lost_assert(S,G) (-)
                          boundary(G)

olist(S,G) = immediate_olist(S,G) (-) RPF_interface(S)

The macros pim_include(*,G) and pim_include(S,G) indicate the
interfaces to which traffic might or might not be forwarded because
of hosts that are local members on those interfaces.

```
pim_include(*,G) = {all interfaces I such that:
                   local_receiver_include(*,G,I)}
pim_include(S,G) = {all interfaces I such that:
                   local_receiver_include(S,G,I)}
pim_exclude(S,G) = {all interfaces I such that:
                   local_receiver_exclude(S,G,I)}
```

The macro RPF_interface(S) returns the RPF interface for source S.
That is to say, it returns the interface used to reach S as indicated
by the MRIB.

The macro local_receiver_include(S,G,I) is true if the IGMP module or
other local membership mechanism ([1], [2], [3], [6]) has determined
that there are local members on interface I that seek to receive
traffic sent specifically by S to G.

The macro local_receiver_include(*,G,I) is true if the IGMP module or
other local membership mechanism has determined that there are local
members on interface I that seek to receive all traffic sent to G.
Note that this determination is expected to account for membership
joins initiated on or by the router.

The macro local_receiver_exclude(S,G,I) is true if
local_receiver_include(*,G,I) is true but none of the local members
seek to receive traffic from S.

The set pim_nbrs is the set of all interfaces on which the router has
at least one active PIM neighbor.

The set prunes(S,G) is the set of all interfaces on which the router
has received Prune(S,G) messages:

```
prunes(S,G) = {all interfaces I such that
               DownstreamPState(S,G,I) is in Pruned state}
```

   The set lost_assert(S,G) is the set of all interfaces on which the
   router has lost an (S,G) Assert.

   lost_assert(S,G) = {all interfaces I such that
                       lost_assert(S,G,I) == TRUE}

   boundary(G) = {all interfaces I with an administratively scoped
                  boundary for group G}

   The following pseudocode macro definitions are also used in many
   places in the specification.  Basically RPF' is the RPF neighbor
   toward a source unless a PIM-DM Assert has overridden the normal
   choice of neighbor.

   neighbor RPF'(S,G) {
     if ( I_Am_Assert_loser(S, G, RPF_interface(S) )) {
       return AssertWinner(S, G, RPF_interface(S) )
     } else {
       return MRIB.next_hop( S )
     }
   }

   The macro I_Am_Assert_loser(S, G, I) is true if the Assert state
   machine (in Section 4.6) for (S,G) on interface I is in the "I am
   Assert Loser" state.

4.2.  Data Packet Forwarding Rules

   The PIM-DM packet forwarding rules are defined below in pseudocode.

   iif is the incoming interface of the packet.  S is the source address
   of the packet.  G is the destination address of the packet (group
   address).  RPF_interface(S) is the interface the MRIB indicates would
   be used to route packets to S.

   First, an RPF check MUST be performed to determine whether the packet
   should be accepted based on TIB state and the interface on which that
   the packet arrived.  Packets that fail the RPF check MUST NOT be
   forwarded, and the router will conduct an assert process for the
   (S,G) pair specified in the packet.  Packets for which a route to the
   source cannot be found MUST be discarded.

   If the RPF check has been passed, an outgoing interface list is
   constructed for the packet.  If this list is not empty, then the
   packet MUST be forwarded to all listed interfaces.  If the list is
   empty, then the router will conduct a prune process for the (S,G)
   pair specified in the packet.

   Upon receipt of a data packet from S addressed to G on interface iif:

   if (iif == RPF_interface(S) AND UpstreamPState(S,G) != Pruned) {
       oiflist = olist(S,G)
   } else {
       oiflist = NULL
   }
   forward packet on all interfaces in oiflist

   This pseudocode employs the following "macro" definition:

   UpstreamPState(S,G) is the state of the Upstream(S,G) state machine
   in Section 4.4.1.

4.3.  Hello Messages

   This section describes the generation and processing of Hello
   messages.

4.3.1.  Sending Hello Messages

   PIM-DM uses Hello messages to detect other PIM routers.  Hello
   messages are sent periodically on each PIM enabled interface.  Hello
   messages are multicast to the ALL-PIM-ROUTERS group.  When PIM is
   enabled on an interface or when a router first starts, the Hello
   Timer (HT) MUST be set to random value between 0 and
   Triggered_Hello_Delay.  This prevents synchronization of Hello
   messages if multiple routers are powered on simultaneously.

   After the initial Hello message, a Hello message MUST be sent every
   Hello_Period.  A single Hello timer MAY be used to trigger sending
   Hello messages on all active interfaces.  The Hello Timer SHOULD NOT
   be reset except when it expires.

4.3.2.  Receiving Hello Messages

   When a Hello message is received, the receiving router SHALL record
   the receiving interface, the sender, and any information contained in
   recognized options.  This information is retained for a number of
   seconds in the Hold Time field of the Hello Message.  If a new Hello
   message is received from a particular neighbor N, the Neighbor
   Liveness Timer (NLT(N,I)) MUST be reset to the newly received Hello
   Holdtime.  If a Hello message is received from a new neighbor, the
   receiving router SHOULD send its own Hello message after a random
   delay between 0 and Triggered_Hello_Delay.

4.3.3.  Hello Message Hold Time

   The Hold Time in the Hello Message should be set to a value that can
   reasonably be expected to keep the Hello active until a new Hello
   message is received.  On most links, this will be 3.5 times the value
   of Hello_Period.

   If the Hold Time is set to '0xffff', the receiving router MUST NOT
   time out that Hello message.  This feature might be used for on-
   demand links to avoid keeping the link up with periodic Hello
   messages.

   If a Hold Time of '0' is received, the corresponding neighbor state
   expires immediately.  When a PIM router takes an interface down or
   changes IP address, a Hello message with a zero Hold Time SHOULD be
   sent immediately (with the old IP address if the IP address is
   changed) to cause any PIM neighbors to remove the old information
   immediately.

4.3.4.  Handling Router Failures

   If a Hello message is received from an active neighbor with a
   different Generation ID (GenID), the neighbor has restarted and may
   not contain the correct (S,G) state.  A Hello message SHOULD be sent
   after a random delay between 0 and Triggered_Hello_Delay (see 4.8)
   before any other messages are sent.  If the neighbor is downstream,
   the router MAY replay the last State Refresh message for any (S,G)
   pairs for which it is the Assert Winner indicating Prune and Assert
   status to the downstream router.  These State Refresh messages SHOULD
   be sent out immediately after the Hello message.  If the neighbor is
   the upstream neighbor for an (S,G) entry, the router MAY cancel its
   Prune Limit Timer to permit sending a prune and reestablishing a
   Pruned state in the upstream router.

   Upon startup, a router MAY use any State Refresh messages received
   within Hello_Period of its first Hello message on an interface to
   establish state information.  The State Refresh source will be the
   RPF'(S), and Prune status for all interfaces will be set according to
   the Prune Indicator bit in the State Refresh message.  If the Prune
   Indicator is set, the router SHOULD set the PruneLimitTimer to
   Prune_Holdtime and set the PruneTimer on all downstream interfaces to
   the State Refresh's Interval times two.  The router SHOULD then
   propagate the State Refresh as described in Section 4.5.1.

4.3.5.  Reducing Prune Propagation Delay on LANs

   If all routers on a LAN support the LAN Prune Delay option, then the
   PIM routers on that LAN will use the values received to adjust their
   J/P_Override_Interval on that interface and the interface is LAN
   Delay Enabled.  Briefly, to avoid synchronization of Prune Override
   (Join) messages when multiple downstream routers share a multi-access
   link, sending of these messages is delayed by a small random amount
   of time.  The period of randomization is configurable and has a
   default value of 3 seconds.

   Each router on the LAN expresses its view of the amount of
   randomization necessary in the Override Interval field of the LAN
   Prune Delay option.  When all routers on a LAN use the LAN Prune
   Delay Option, all routers on the LAN MUST set their Override_Interval
   to the largest Override value on the LAN.

   The LAN Delay inserted by a router in the LAN Prune Delay option
   expresses the expected message propagation delay on the link and
   SHOULD be configurable by the system administrator.  When all routers
   on a link use the LAN Prune Delay Option, all routers on the LAN MUST
   set Propagation Delay to the largest LAN Delay on the LAN.

   PIM implementers should enforce a lower bound on the permitted values
   for this delay to allow for scheduling and processing delays within
   their router.  Such delays may cause received messages to be
   processed later and triggered messages to be sent later than
   intended.  Setting this LAN Prune Delay to too low a value may result
   in temporary forwarding outages, because a downstream router will not
   be able to override a neighbor's prune message before the upstream
   neighbor stops forwarding.

4.4.  PIM-DM Prune, Join, and Graft Messages

   This section describes the generation and processing of PIM-DM Join,
   Prune, and Graft messages.  Prune messages are sent toward the
   upstream neighbor for S to indicate that traffic from S addressed to
   group G is not desired.  In the case of downstream routers A and B,
   where A wishes to continue receiving data and B does not, A will send
   a Join in response to B's Prune to override the Prune.  This is the
   only situation in PIM-DM in which a Join message is used.  Finally, a
   Graft message is used to re-join a previously pruned branch to the
   delivery tree.

4.4.1.  Upstream Prune, Join, and Graft Messages

   The Upstream(S,G) state machine for sending Prune, Graft, and Join
   messages is given below.  There are three states.

      Forwarding (F)
        This is the starting state of the Upsteam(S,G) state machine.
        The state machine is in this state if it just started or if
        oiflist(S,G) != NULL.

      Pruned (P)
        The set, olist(S,G), is empty.  The router will not forward data
        from S addressed to group G.

      AckPending (AP)
        The router was in the Pruned(P) state, but a transition has
        occurred in the Downstream(S,G) state machine for one of this
        (S,G) entry's outgoing interfaces, indicating that traffic from S
        addressed to G should again be forwarded.  A Graft message has
        been sent to RPF'(S), but a Graft Ack message has not yet been
        received.

   In addition, there are three state-machine-specific timers:

      GraftRetry Timer (GRT(S,G))
        This timer is set when a Graft is sent upstream.  If a
        corresponding GraftAck is not received before the timer expires,
        then another Graft is sent, and the GraftRetry Timer is reset.
        The timer is stopped when a Graft Ack message is received.  This
        timer is normally set to Graft_Retry_Period (see 4.8).

      Override Timer (OT(S,G))
        This timer is set when a Prune(S,G) is received on the upstream
        interface where olist(S,G) != NULL.  When the timer expires, a
        Join(S,G) message is sent on the upstream interface.  This timer
        is normally set to t_override (see 4.8).

      Prune Limit Timer (PLT(S,G))
        This timer is used to rate-limit Prunes on a LAN.  It is only
        used when the Upstream(S,G) state machine is in the Pruned state.
        A Prune cannot be sent if this timer is running.  This timer is
        normally set to t_limit (see 4.8).

```
  +------------+                             +------------+
  |            |       olist == NULL         |            |
  |  Forward   |---------------------------->|   Pruned   |
  |            |                             |            |
  +------------+                             +------------+
      ^    |                                     ^    |
      |    |                                     |    |
      |    |RPF'(S) Changes      olist == NULL|    |
      |    |                                     |    |
      |    |              +------------+         |    |
      |    +-------->|            |---------+   |
      |              | AckPending |             |
  +------------|              |<------------+
 Rcv GraftAck OR +------------+ olist != NULL
  Rcv State Refresh
     With (P==0) OR
  S Directly Connect
```

                 Figure 1: Upstream Interface State Machine

   In tabular form, the state machine is defined as follows:

| | Previous State | | |
| Event | Forwarding | Pruned | AckPending |
|---|---|---|---|
| Data packet arrives on RPF_Interface(S) AND olist(S,G) == NULL AND PLT(S,G) not running | ->P Send Prune(S,G) Set PLT(S,G) | ->P Send Prune(S,G) Set PLT(S,G) | N/A |
| State Refresh(S,G) received from RPF'(S) AND Prune Indicator == 1 | ->F  Set OT(S,G) | ->P Reset PLT(S,G) | ->AP  Set OT(S,G) |
| State Refresh(S,G) received from RPF'(S) AND Prune Indicator == 0 AND PLT(S,G) not running | ->F | ->P Send Prune(S,G) Set PLT(S,G) | ->F Cancel GRT(S,G) |

| | Previous State | | |
|---|---|---|---|
| Event | Forwarding | Pruned | AckPending |
| See Join(S,G) to RPF'(S) | ->F Cancel OT(S,G) | ->P | ->AP Cancel OT(S,G) |
| See Prune(S,G) | ->F Set OT(S,G) | ->P | ->AP Set OT(S,G) |
| OT(S,G) Expires | ->F Send Join(S,G) | N/A | ->AP Send Join(S,G) |
| olist(S,G)->NULL | ->P Send Prune(S,G) Set PLT(S,G) | N/A | ->P Send Prune(S,G) Set PLT(S,G) Cancel GRT(S,G) |
| olist(S,G)->non-NULL | N/A | ->AP Send Graft(S,G) Set GRT(S,G) | N/A |
| RPF'(S) Changes AND olist(S,G) != NULL | ->AP Send Graft(S,G) Set GRT(S,G) | ->AP Send Graft(S,G) Set GRT(S,G) | ->AP Send Graft(S,G) Set GRT(S,G) |
| RPF'(S) Changes AND olist(S,G) == NULL | ->P | ->P Cancel PLT(S,G) | ->P Cancel GRT(S,G) |
| S becomes directly connected | ->F | ->P | ->F Cancel GRT(S,G) |
| GRT(S,G) Expires | N/A | N/A | ->AP Send Graft(S,G) Set GRT(S,G) |
| Receive GraftAck(S,G) from RPF'(S) | ->F | ->P | ->F Cancel GRT(S,G) |

   The transition event "RcvGraftAck(S,G)" implies receiving a Graft Ack
   message targeted to this router's address on the incoming interface
   for the (S,G) entry.  If the destination address is not correct, the
   state transitions in this state machine must not occur.

4.4.1.1.  Transitions from the Forwarding (F) State

   When the Upstream(S,G) state machine is in the Forwarding (F) state,
   the following events may trigger a transition:

      Data Packet arrives on RPF_Interface(S) AND olist(S,G) == NULL AND
      S NOT directly connected
         The Upstream(S,G) state machine MUST transition to the Pruned (P)
         state, send a Prune(S,G) to RPF'(S), and set PLT(S,G) to t_limit
         seconds.

      State Refresh(S,G) Received from RPF'(S)
         The Upstream(S,G) state machine remains in a Forwarding state.
         If the received State Refresh has the Prune Indicator bit set to
         one, this router must override the upstream router's Prune state
         after a short random interval.  If OT(S,G) is not running and the
         Prune Indicator bit equals one, the router MUST set OT(S,G) to
         t_override seconds.

      See Join(S,G) to RPF'(S)
         This event is only relevant if RPF_interface(S) is a shared
         medium.  This router sees another router on RPF_interface(S) send
         a Join(S,G) to RPF'(S,G).  If the OT(S,G) is running, then it
         means that the router had scheduled a Join to override a
         previously received Prune.  Another router has responded more
         quickly with a Join, so the local router SHOULD cancel its
         OT(S,G), if it is running.  The Upstream(S,G) state machine
         remains in the Forwarding (F) state.

      See Prune(S,G) AND S NOT directly connected
         This event is only relevant if RPF_interface(S) is a shared
         medium.  This router sees another router on RPF_interface(S) send
         a Prune(S,G).  As this router is in Forwarding state, it must
         override the Prune after a short random interval.  If OT(S,G) is
         not running, the router MUST set OT(S,G) to t_override seconds.
         The Upstream(S,G) state machine remains in Forwarding (F) state.

      OT(S,G) Expires AND S NOT directly connected
         The OverrideTimer (OT(S,G)) expires.  The router MUST send a
         Join(S,G) to RPF'(S) to override a previously detected prune.
         The Upstream(S,G) state machine remains in the Forwarding (F)
         state.

      olist(S,G) -> NULL AND S NOT directly connected
         The Upstream(S,G) state machine MUST transition to the Pruned (P)
         state, send a Prune(S,G) to RPF'(S), and set PLT(S,G) to t_limit
         seconds.

   RPF'(S) Changes AND olist(S,G) is non-NULL AND S NOT directly
   connected
      Unicast routing or Assert state causes RPF'(S) to change,
      including changes to RPF_Interface(S).  The Upstream(S,G) state
      machine MUST transition to the AckPending (AP) state, unicast a
      Graft to the new RPF'(S), and set the GraftRetry Timer (GRT(S,G))
      to Graft_Retry_Period.

   RPF'(S) Changes AND olist(S,G) is NULL
      Unicast routing or Assert state causes RPF'(S) to change,
      including changes to RPF_Interface(S).  The Upstream(S,G) state
      machine MUST transition to the Pruned (P) state.

## 4.4.1.2.  Transitions from the Pruned (P) State

   When the Upstream(S,G) state machine is in the Pruned (P) state, the
   following events may trigger a transition:

   Data arrives on RPF_interface(S) AND PLT(S,G) not running AND S NOT
   directly connected
      Either another router on the LAN desires traffic from S addressed
      to G or a previous Prune was lost.  To prevent generating a
      Prune(S,G) in response to every data packet, the PruneLimit Timer
      (PLT(S,G)) is used.  Once the PLT(S,G) expires, the router needs
      to send another prune in response to a data packet not received
      directly from the source.  A Prune(S,G) MUST be sent to RPF'(S),
      and the PLT(S,G) MUST be set to t_limit.

   State Refresh(S,G) Received from RPF'(S)
      The Upstream(S,G) state machine remains in a Pruned state.  If
      the State Refresh has its Prune Indicator bit set to zero and
      PLT(S,G) is not running, a Prune(S,G) MUST be sent to RPF'(S),
      and the PLT(S,G) MUST be set to t_limit.  If the State Refresh
      has its Prune Indicator bit set to one, the router MUST reset
      PLT(S,G) to t_limit.

   See Prune(S,G) to RPF'(S)
      A Prune(S,G) is seen on RPF_interface(S) to RPF'(S).  The
      Upstream(S,G) state machine stays in the Pruned (P) state.  The
      router MAY reset its PLT(S,G) to the value in the Holdtime field
      of the received message if it is greater than the current value
      of the PLT(S,G).

   olist(S,G)->non-NULL AND S NOT directly connected
      The set of interfaces defined by the olist(S,G) macro becomes
      non-empty, indicating that traffic from S addressed to group G
      must be forwarded.  The Upstream(S,G) state machine MUST cancel
      PLT(S,G), transition to the AckPending (AP) state and unicast a

      Graft message to RPF'(S).  The Graft Retry Timer (GRT(S,G)) MUST
      be set to Graft_Retry_Period.

   RPF'(S) Changes AND olist(S,G) == non-NULL AND S NOT directly
   connected
      Unicast routing or Assert state causes RPF'(S) to change,
      including changes to RPF_Interface(S).  The Upstream(S,G) state
      machine MUST cancel PLT(S,G), transition to the AckPending (AP)
      state, send a Graft unicast to the new RPF'(S), and set the
      GraftRetry Timer (GRT(S,G)) to Graft_Retry_Period.

   RPF'(S) Changes AND olist(S,G) == NULL AND S NOT directly connected
      Unicast routing or Assert state causes RPF'(S) to change,
      including changes to RPF_Interface(S).  The Upstream(S,G) state
      machine stays in the Pruned (P) state and MUST cancel the
      PLT(S,G) timer.

   S becomes directly connected
      Unicast routing changed so that S is directly connected.  The
      Upstream(S,G) state machine remains in the Pruned (P) state.

4.4.1.3.  Transitions from the AckPending (AP) State

   When the Upstream(S,G) state machine is in the AckPending (AP) state,
   the following events may trigger a transition:

   State Refresh(S,G) Received from RPF'(S) with Prune Indicator == 1
      The Upstream(S,G) state machine remains in an AckPending state.
      The router must override the upstream router's Prune state after
      a short random interval.  If OT(S,G) is not running and the Prune
      Indicator bit equals one, the router MUST set OT(S,G) to
      t_override seconds.

   State Refresh(S,G) Received from RPF'(S) with Prune Indicator == 0
      The router MUST cancel its GraftRetry Timer (GRT(S,G)) and
      transition to the Forwarding (F) state.

   See Join(S,G) to RPF'(S,G)
      This event is only relevant if RPF_interface(S) is a shared
      medium.  This router sees another router on RPF_interface(S) send
      a Join(S,G) to RPF'(S,G).  If the OT(S,G) is running, then it
      means that the router had scheduled a Join to override a
      previously received Prune.  Another router has responded more
      quickly with a Join, so the local router SHOULD cancel its
      OT(S,G), if it is running.  The Upstream(S,G) state machine
      remains in the AckPending (AP) state.

See Prune(S,G)
   This event is only relevant if RPF_interface(S) is a shared
   medium.  This router sees another router on RPF_interface(S) send
   a Prune(S,G).  As this router is in AckPending (AP) state, it
   must override the Prune after a short random interval.  If
   OT(S,G) is not running, the router MUST set OT(S,G) to t_override
   seconds.  The Upstream(S,G) state machine remains in AckPending
   (AP) state.

OT(S,G) Expires
   The OverrideTimer (OT(S,G)) expires.  The router MUST send a
   Join(S,G) to RPF'(S).  The Upstream(S,G) state machine remains in
   the AckPending (AP) state.

olist(S,G) -> NULL
   The set of interfaces defined by the olist(S,G) macro becomes
   null, indicating that traffic from S addressed to group G should
   no longer be forwarded.  The Upstream(S,G) state machine MUST
   transition to the Pruned (P) state.  A Prune(S,G) MUST be
   multicast to the RPF_interface(S), with RPF'(S) named in the
   upstream neighbor field.  The GraftRetry Timer (GRT(S,G)) MUST be
   cancelled, and PLT(S,G) MUST be set to t_limit seconds.

RPF'(S) Changes AND olist(S,G) does not become NULL AND S NOT
directly connected
   Unicast routing or Assert state causes RPF'(S) to change,
   including changes to RPF_Interface(S).  The Upstream(S,G) state
   machine stays in the AckPending (AP) state.  A Graft MUST be
   unicast to the new RPF'(S) and the GraftRetry Timer (GRT(S,G))
   reset to Graft_Retry_Period.

RPF'(S) Changes AND olist(S,G) == NULL AND S NOT directly connected
   Unicast routing or Assert state causes RPF'(S) to change,
   including changes to RPF_Interface(S).  The Upstream(S,G) state
   machine MUST transition to the Pruned (P) state.  The GraftRetry
   Timer (GRT(S,G)) MUST be cancelled.

S becomes directly connected
   Unicast routing has changed so that S is directly connected.  The
   GraftRetry Timer MUST be cancelled, and the Upstream(S,G) state
   machine MUST transition to the Forwarding(F) state.

GRT(S,G) Expires
    The GraftRetry Timer (GRT(S,G)) expires for this (S,G) entry.
    The Upstream(S,G) state machine stays in the AckPending (AP)
    state.  Another Graft message for (S,G) SHOULD be unicast to
    RPF'(S) and the GraftRetry Timer (GRT(S,G)) reset to
    Graft_Retry_Period.  It is RECOMMENDED that the router retry a
    configured number of times before ceasing retries.

See GraftAck(S,G) from RPF'(S)
    A GraftAck is received from  RPF'(S).  The GraftRetry Timer MUST
    be cancelled, and the Upstream(S,G) state machine MUST transition
    to the Forwarding(F) state.

## 4.4.2.  Downstream Prune, Join, and Graft Messages

The Prune(S,G) Downstream state machine for receiving Prune, Join and
Graft messages on interface I is given below.  This state machine
MUST always be in the NoInfo state on the upstream interface.  It
contains three states.

NoInfo(NI)
    The interface has no (S,G) Prune state, and neither the Prune
    timer (PT(S,G,I)) nor the PrunePending timer ((PPT(S,G,I)) is
    running.

PrunePending(PP)
    The router has received a Prune(S,G) on this interface from a
    downstream neighbor and is waiting to see whether the prune will
    be overridden by another downstream router.  For forwarding
    purposes, the PrunePending state functions exactly like the
    NoInfo state.

Pruned(P)
    The router has received a Prune(S,G) on this interface from a
    downstream neighbor, and the Prune was not overridden.  Data from
    S addressed to group G is no longer being forwarded on this
    interface.

In addition, there are two timers:

PrunePending Timer (PPT(S,G,I))
    This timer is set when a valid Prune(S,G) is received.  Expiry of
    the PrunePending Timer (PPT(S,G,I)) causes the interface to
    transition to the Pruned state.

   Prune Timer (PT(S,G,I))
      This timer is set when the PrunePending Timer (PT(S,G,I))
      expires.  Expiry of the Prune Timer (PT(S,G,I)) causes the
      interface to transition to the NoInfo (NI) state, thereby
      allowing data from S addressed to group G to be forwarded on the
      interface.

```
      +------------+                          +------------+
      |            |      PPT Expires         |            |
      |PrunePending|------------------------->|  Pruned    |
      |            |                          |            |
      +------------+                          +------------+
          |   ^                                      |
          |   |                                      |
          |   |Rcv Prune                             |
          |   |                                      |
          |   |        +------------+                |
          |   +--------|            |                |
          |            |  NoInfo    |<-------------+
          +------------>|            |  Rcv Join/Graft OR
       Rcv Join/Graft OR +------------+ PT Expires OR
        RPF_Interface(S)->I          RPF_Interface(S)->I
```

               Figure 2: Downstream Interface State Machine

In tabular form, the state machine is as follows:

| Event | Previous State | | |
|-------|---------|-----------|--------|
|       | No Info | PrunePend | Pruned |
| Receive Prune(S,G) | ->PP  Set PPT(S,G,I) | ->PP | ->P Reset PT(S,G,I) |
| Receive Join(S,G) | ->NI | ->NI Cancel PPT(S,G,I) | ->NI Cancel PT(S,G,I) |
| Receive Graft(S,G) | ->NI Send GraftAck | ->NI Send GraftAck Cancel PPT(S,G,I) | ->NI Send GraftAck Cancel PT(S,G,I) |
| PPT(S,G) Expires | N/A | ->P Set PT(S,G,I) | N/A |
| PT(S,G) Expires | N/A | N/A | ->NI |
| RPF_Interface(S) becomes I | ->NI | ->NI Cancel PPT(S,G,I) | ->NI Cancel PT(S,G,I) |
| Send State Refresh(S,G) out I | ->NI | ->PP | ->P Reset PT(S,G,I) |

The transition events "Receive Graft(S,G)", "Receive Prune(S,G)", and
"Receive Join(S,G)" denote receiving a Graft, Prune, or Join message
in which this router's address on I is contained in the message's
upstream neighbor field.  If the upstream neighbor field does not
match this router's address on I, then these state transitions in
this state machine must not occur.

4.4.2.1.  Transitions from the NoInfo State

   When the Prune(S,G) Downstream state machine is in the NoInfo (NI)
   state, the following events may trigger a transition:

      Receive Prune(S,G)
         A Prune(S,G) is received on interface I with the upstream
         neighbor field set to the router's address on I.  The Prune(S,G)
         Downstream state machine on interface I MUST transition to the
         PrunePending (PP) state.  The PrunePending Timer (PPT(S,G,I))
         MUST be set to J/P_Override_Interval if the router has more than

one neighbor on I.  If the router has only one neighbor on
interface I, then it SHOULD set the PPT(S,G,I) to zero,
effectively transitioning immediately to the Pruned (P) state.

Receive Graft(S,G)
A Graft(S,G) is received on the interface I with the upstream
neighbor field set to the router's address on I.  The Prune(S,G)
Downstream state machine on interface I stays in the NoInfo (NI)
state.  A GraftAck(S,G) MUST be unicast to the originator of the
Graft(S,G) message.

4.4.2.2.  Transitions from the PrunePending (PP) State

When the Prune(S,G) downstream state machine is in the PrunePending
(PP) state, the following events may trigger a transition.

Receive Join(S,G)
A Join(S,G) is received on interface I with the upstream neighbor
field set to the router's address on I.  The Prune(S,G)
Downstream state machine on interface I MUST transition to the
NoInfo (NI) state.  The PrunePending Timer (PPT(S,G,I)) MUST be
cancelled.

Receive Graft(S,G)
A Graft(S,G) is received on interface I with the upstream
neighbor field set to the router's address on I.  The Prune(S,G)
Downstream state machine on interface I MUST transition to the
NoInfo (NI) state and MUST unicast a Graft Ack message to the
Graft originator.  The PrunePending Timer (PPT(S,G,I)) MUST be
cancelled.

PPT(S,G,I) Expires
The PrunePending Timer (PPT(S,G,I)) expires, indicating that no
neighbors have overridden the previous Prune(S,G) message.  The
Prune(S,G) Downstream state machine on interface I MUST
transition to the Pruned (P) state.  The Prune Timer (PT(S,G,I))
is started and MUST be initialized to the received
Prune_Hold_Time minus J/P_Override_Interval.  A PruneEcho(S,G)
MUST be sent on I if I has more than one PIM neighbor.  A
PruneEcho(S,G) is simply a Prune(S,G) message multicast by the
upstream router to a LAN, with itself as the Upstream Neighbor.
Its purpose is to add additional reliability so that if a Join
that should have overridden the Prune is lost locally on the LAN,
the PruneEcho(S,G) may be received and trigger a new Join
message.  A PruneEcho(S,G) is OPTIONAL on an interface with only
one PIM neighbor.  In addition, the router MUST evaluate any
possible transitions in the Upstream(S,G) state machine.

      RPF_Interface(S) becomes interface I
        The upstream interface for S has changed.  The Prune(S,G)
        Downstream state machine on interface I MUST transition to the
        NoInfo (NI) state.  The PrunePending Timer (PPT(S,G,I)) MUST be
        cancelled.

4.4.2.3.  Transitions from the Prune (P) State

   When the Prune(S,G) Downstream state machine is in the Pruned (P)
   state, the following events may trigger a transition.

      Receive Prune(S,G)
        A Prune(S,G) is received on the interface I with the upstream
        neighbor field set to the router's address on I.  The Prune(S,G)
        Downstream state machine on interface I remains in the Pruned (P)
        state.  The Prune Timer (PT(S,G,I)) SHOULD be reset to the
        holdtime contained in the Prune(S,G) message if it is greater
        than the current value.

      Receive Join(S,G)
        A Join(S,G) is received on the interface I with the upstream
        neighbor field set to the router's address on I.  The Prune(S,G)
        downstream state machine on interface I MUST transition to the
        NoInfo (NI) state.  The Prune Timer (PT(S,G,I)) MUST be
        cancelled.  The router MUST evaluate any possible transitions in
        the Upstream(S,G) state machine.

      Receive Graft(S,G)
        A Graft(S,G) is received on interface I with the upstream
        neighbor field set to the router's address on I.  The Prune(S,G)
        Downstream state machine on interface I MUST transition to the
        NoInfo (NI) state and send a Graft Ack back to the Graft's
        source.  The Prune Timer (PT(S,G,I)) MUST be cancelled.  The
        router MUST evaluate any possible transitions in the
        Upstream(S,G) state machine.

      PT(S,G,I) Expires
        The Prune Timer (PT(S,G,I)) expires, indicating that it is again
        time to flood data from S addressed to group G onto interface I.
        The Prune(S,G) Downstream state machine on interface I MUST
        transition to the NoInfo (NI) state.  The router MUST evaluate
        any possible transitions in the Upstream(S,G) state machine.

      RPF_Interface(S) becomes interface I
        The upstream interface for S has changed.  The Prune(S,G)
        Downstream state machine on interface I MUST transition to the
        NoInfo (NI) state.  The PruneTimer (PT(S,G,I)) MUST be cancelled.

      Send State Refresh(S,G) out interface I
         The router has refreshed the Prune(S,G) state on interface I.
         The router MUST reset the Prune Timer (PT(S,G,I)) to the Holdtime
         from an active Prune received on interface I.  The Holdtime used
         SHOULD be the largest active one but MAY be the most recently
         received active Prune Holdtime.

## 4.5.  State Refresh

   This section describes the major portions of the state refresh
   mechanism.

## 4.5.1.  Forwarding of State Refresh Messages

   When a State Refresh message, SRM, is received, it is forwarded
   according to the following pseudo-code.

```
   if (iif != RPF_interface(S))
     return;
   if (RPF'(S) != srcaddr(SRM))
     return;
   if (StateRefreshRateLimit(S,G) == TRUE)
     return;

   for each interface I in pim_nbrs {
     if (TTL(SRM) == 0 OR (TTL(SRM) - 1) < Threshold(I))
       continue;      /* Out of TTL, skip this interface */
     if (boundary(I,G))
       continue;      /* This interface is scope boundary, skip it */
     if (I == iif)
       continue;      /* This is the incoming interface, skip it */
     if (lost_assert(S,G,I) == TRUE)
       continue;      /* Let the Assert Winner do State Refresh */

     Copy SRM to SRM';   /* Make a copy of SRM to forward */

     if (I contained in prunes(S,G)) {
       set Prune Indicator bit of SRM' to 1;

       if StateRefreshCapable(I) == TRUE
         set PT(S,G) to largest active holdtime read from a Prune
         message accepted on I;
```

```
   } else {
     set Prune Indicator bit of SRM' to 0;
   }

   set srcaddr(SRM') to my_addr(I);
   set TTL of SRM' to TTL(SRM) - 1;
   set metric of SRM' to metric of unicast route used to reach S;
   set pref of SRM' to preference of unicast route used to reach S;
   set mask of SRM' to mask of route used to reach S;

   if (AssertState == NoInfo) {
     set Assert Override of SRM' to 1;
   } else {
     set Assert Override of SRM' to 0;
   }

   transmit SRM' on I;
}
```

   The pseudocode above employs the following macro definitions.

   Boundary(I,G) is TRUE if an administratively scoped boundary for
   group G is configured on interface I.

   StateRefreshCapable(I) is TRUE if all neighbors on an interface use
   the State Refresh option.

   StateRefreshRateLimit(S,G) is TRUE if the time elapsed since the last
   received StateRefresh(S,G) is less than the configured
   RefreshLimitInterval.

   TTL(SRM) returns the TTL contained in the State Refresh Message, SRM.
   This is different from the TTL contained in the IP header.

   Threshold(I) returns the minimum TTL that a packet must have before
   it can be transmitted on interface I.

   srcaddr(SRM) returns the source address contained in the network
   protocol (e.g., IPv4) header of the State Refresh Message, SRM.

   my_addr(I) returns this node's network (e.g., IPv4) address on
   interface I.

4.5.2.  State Refresh Message Origination

   This section describes the origination of State Refresh messages.
   These messages are generated periodically by the PIM-DM router
   directly connected to a source.  One Origination(S,G) state machine
   exists per (S,G) entry in a PIM-DM router.

   The Origination(S,G) state machine has the following states:

     NotOriginator(NO)
       This is the starting state of the Origination(S,G) state machine.
       While in this state, a router will not originate State Refresh
       messages for the (S,G) pair.

     Originator(O)
       When in this state the router will periodically originate State
       Refresh messages.  Only routers directly connected to S may
       transition to this state.

   In addition, there are two state machine specific timers:

     State Refresh Timer (SRT(S,G))
       This timer controls when State Refresh messages are generated.
       The timer is initially set when that Origination(S,G) state
       machine transitions to the O state.  It is cancelled when the
       Origination(S,G) state machine transitions to the NO state.  This
       timer is normally set to StateRefreshInterval (see 4.8).

     Source Active Timer (SAT(S,G))
       This timer is first set when the Origination(S,G) state machine
       transitions to the O state and is reset on the receipt of every
       data packet from S addressed to group G.  When it expires, the
       Origination(S,G) state machine transitions to the NO state.  This
       timer is normally set to SourceLifetime (see 4.8).

```
         +-------------+   Rcv Directly From S   +------------+
         |             |------------------------->|            |
         |NotOriginator|                          | Originator |
         |             |<-------------------------|            |
         +------------+        SAT Expires OR     +------------+
                           S NOT Direct Connect
```

                   Figure 3: State Refresh State Machine

In tabular form, the state machine is defined as follows:

```
+--------------------------------+---------------+------------------+
|                                |         Previous State           |
|                                +---------------+------------------+
|             Event              | NotOriginator |    Originator    |
+--------------------------------+---------------+------------------+
| Receive Data from S AND        | ->O           | ->O Reset        |
| S directly connected           | Set SRT(S,G)  |    SAT(S,G)       |
|                                | Set SAT(S,G)  |                  |
+--------------------------------+---------------+------------------+
| SRT(S,G) Expires               | N/A           | ->O    Send      |
|                                |               | StateRefresh(S,G) |
|                                |               | Reset SRT(S,G)   |
+--------------------------------+---------------+------------------+
| SAT(S,G) Expires               | N/A           | ->NO  Cancel     |
|                                |               |     SRT(S,G)     |
+--------------------------------+---------------+------------------+
| S no longer directly connected | ->NO          | ->NO             |
|                                |               |  Cancel SRT(S,G) |
|                                |               |  Cancel SAT(S,G) |
+--------------------------------+---------------+------------------+
```

4.5.2.1.  Transitions from the NotOriginator (NO) State

   When the Originating(S,G) state machine is in the NotOriginator (NO)
   state, the following event may trigger a transition:

      Data Packet received from directly connected Source S addressed to
      group G
         The router MUST transition to an Originator (O) state, set
         SAT(S,G) to SourceLifetime, and set SRT(S,G) to
         StateRefreshInterval.  The router SHOULD record the TTL of the
         packet for use in State Refresh messages.

4.5.2.2.  Transitions from the Originator (O) State

   When the Originating(S,G) state machine is in the Originator (O)
   state, the following events may trigger a transition:

      Receive Data Packet from S addressed to G
         The router remains in the Originator (O) state and MUST reset
         SAT(S,G) to SourceLifetime.  The router SHOULD increase its
         recorded TTL to match the TTL of the packet, if the packet's TTL
         is larger than the previously recorded TTL.  A router MAY record
         the TTL based on an implementation specific sampling policy to
         avoid examining the TTL of every multicast packet it handles.

SRT(S,G) Expires
   The router remains in the Originator (O) state and MUST reset
   SRT(S,G) to StateRefreshInterval.  The router MUST also generate
   State Refresh messages for transmission, as described in the
   State Refresh Forwarding rules (Section 4.5.1), except for the
   TTL.  If the TTL of data packets from S to G are being recorded,
   then the TTL of each State Refresh message is set to the highest
   recorded TTL.  Otherwise, the TTL is set to the configured State
   Refresh TTL.  Let I denote the interface over which a State
   Refresh message is being sent.  If the Prune(S,G) Downstream
   state machine is in the Pruned (P) state, then the Prune-
   Indicator bit MUST be set to 1 in the State Refresh message being
   sent over I. Otherwise, the Prune-Indicator bit MUST be set to 0.

SAT(S,G) Expires
   The router MUST cancel the SRT(S,G) timer and transition to the
   NotOriginator (NO) state.

S is no longer directly connected
   The router MUST transition to the NotOriginator (NO) state and
   cancel both the SAT(S,G) and SRT(S,G).

4.6.  PIM Assert Messages

4.6.1.  Assert Metrics

   Assert metrics are defined as follows:

   struct assert_metric {
     metric_preference;
     route_metric;
     ip_address;
   };

   When assert_metrics are compared, the metric_preference and
   route_metric field are compared in order, where the first lower value
   wins.  If all fields are equal, the IP address of the router that
   sourced the Assert message is used as a tie-breaker, with the highest
   IP address winning.

An Assert metric for (S,G) to include in (or compare against) an
Assert message sent on interface I should be computed by using the
following pseudocode:

```
assert_metric
my_assert_metric(S,G,I) {
  if (CouldAssert(S,G,I) == TRUE) {
    return spt_assert_metric(S,G,I)
  } else {
    return infinite_assert_metric()
  }
}
```

spt_assert_metric(S,I) gives the Assert metric we use if we're
sending an Assert based on active (S,G) forwarding state:

```
assert_metric
spt_assert_metric(S,I) {
  return {0,MRIB.pref(S),MRIB.metric(S),my_addr(I)}
}
```

MRIB.pref(X) and MRIB.metric(X) are the routing preference and
routing metrics associated with the route to a particular (unicast)
destination X, as determined by the MRIB.  my_addr(I) is simply the
router's network (e.g., IP) address associated with the local
interface I.

infinite_assert_metric() gives the Assert metric we need to send an
Assert but doesn't match (S,G) forwarding state:

```
assert_metric
infinite_assert_metric() {
  return {1,infinity,infinity,0}
}
```

4.6.2.  AssertCancel Messages

   An AssertCancel(S,G) message is simply an Assert message for (S,G)
   with infinite metric.  The Assert winner sends this message when it
   changes its upstream interface to this interface.  Other routers will
   see this metric, causing those with forwarding state to send their
   own Asserts and re-establish an Assert winner.

   AssertCancel messages are simply an optimization.  The original
   Assert timeout mechanism will eventually allow a subnet to become
   consistent; the AssertCancel mechanism simply causes faster
   convergence.  No special processing is required for an AssertCancel
   message, as it is simply an Assert message from the current winner.

4.6.3.  Assert State Macros

   The macro lost_assert(S,G,I), is used in the olist computations of
   Section 4.1.3, and is defined as follows:

```
bool lost_assert(S,G,I) {
  if ( RPF_interface(S) == I ) {
    return FALSE
  } else {
    return (AssertWinner(S,G,I) != me  AND
            (AssertWinnerMetric(S,G,I) is better than
             spt_assert_metric(S,G,I)))
  }
}
```

   AssertWinner(S,G,I) defaults to NULL, and AssertWinnerMetric(S,G,I)
   defaults to Infinity when in the NoInfo state.

4.6.4.  (S,G) Assert Message State Machine

   The (S,G) Assert state machine for interface I is shown in Figure 4.
   There are three states:

     NoInfo (NI)
       This router has no (S,G) Assert state on interface I.

     I am Assert Winner (W)
       This router has won an (S,G) Assert on interface I.  It is now
       responsible for forwarding traffic from S destined for G via
       interface I.

     I am Assert Loser (L)
       This router has lost an (S,G) Assert on interface I.  It must not
       forward packets from S destined for G onto interface I.

   In addition, an Assert Timer (AT(S,G,I)) is used to time out the
   Assert state.

```
     +------------+                          +------------+
     |            | Rcv Pref Assert or SR    |            |
     |   Winner   |------------------------->|   Loser    |
     |            |                          |            |
     +------------+                          +------------+
        ^  |                                    ^  |
        |  |                                    |  |
        |  |AT Expires OR         Rcv Pref Assert or|  |
        |  |CouldAssert->FALSE       State Refresh|  |
        |  |                                    |  |
        |  |            +------------+          |  |
        |  +-------->|            |----------+  |
        |            | No Info    |             |
     +------------|            |<------------+
     Rcv Data from dnstrm +------------+ Rcv Inf Assert from Win OR
     OR Rcv Inferior Assert            Rcv Inf SR from Winner OR
        OR Rcv Inferior SR             AT Expires OR
                                       CouldAssert Changes OR
                                       Winner's NLT Expires
```

Figure 4: Assert State Machine

In tabular form, the state machine is defined as follows:

| Event | Previous State | | |
| | No Info | Winner | Loser |
|---|---|---|---|
| An (S,G) Data packet received on downstream interface | ->W Send Assert(S,G) Set AT(S,G,I) | ->W Send Assert(S,G) Set AT(S,G,I) | ->L |
| Receive Inferior (Assert OR State Refresh) from Assert Winner | N/A | N/A | ->NI Cancel AT(S,G,I) |
| Receive Inferior (Assert OR State Refresh) from non-Assert Winner AND CouldAssert==TRUE | ->W Send Assert(S,G) Set AT(S,G,I) | ->W Send Assert(S,G) Set AT(S,G,I) | ->L |

| Event | Previous State | | |
| :--- | :--- | :--- | :--- |
| | No Info | Winner | Loser |
| Receive Preferred Assert OR State Refresh | ->L Send Prune(S,G) Set AT(S,G,I) | ->L Send Prune(S,G) Set AT(S,G,I) | ->L  Set AT(S,G,I) |
| Send State Refresh | ->NI | ->W Reset AT(S,G,I) | N/A |
| AT(S,G) Expires | N/A | ->NI | ->NI |
| CouldAssert -> FALSE | ->NI | ->NI Cancel AT(S,G,I) | ->NI Cancel AT(S,G,I) |
| CouldAssert -> TRUE | ->NI | N/A | ->NI Cancel AT(S,G,I) |
| Winner's NLT(N,I) Expires | N/A | N/A | ->NI Cancel AT(S,G,I) |
| Receive Prune(S,G), Join(S,G) or Graft(S,G) | ->NI | ->W | ->L Send Assert(S,G) |

   Terminology: A "preferred assert" is one with a better metric than
   the current winner.  An "inferior assert" is one with a worse metric
   than my_assert_metric(S,G,I).

   The state machine uses the following macro:

   CouldAssert(S,G,I) = (RPF_interface(S) != I)

## 4.6.4.1.  Transitions from NoInfo State

   In the NoInfo state, the following events may trigger transitions:

     An (S,G) data packet arrives on downstream interface I
        An (S,G) data packet arrived on a downstream interface.  It is
        optimistically assumed that this router will be the Assert winner
        for this (S,G).  The Assert state machine MUST transition to the
        "I am Assert Winner" state, send an Assert(S,G) to interface I,
        store its own address and metric as the Assert Winner, and set
        the Assert_Timer (AT(S,G,I) to Assert_Time, thereby initiating
        the Assert negotiation for (S,G).

      Receive Inferior (Assert OR State Refresh) AND
      CouldAssert(S,G,I)==TRUE
         An Assert or State Refresh is received for (S,G) that is inferior
         to our own assert metric on interface I. The Assert state machine
         MUST transition to the "I am Assert Winner" state, send an
         Assert(S,G) to interface I, store its own address and metric as
         the Assert Winner, and set the Assert Timer (AT(S,G,I)) to
         Assert_Time.

      Receive Preferred Assert or State Refresh
         The received Assert or State Refresh has a better metric than
         this router's, and therefore the Assert state machine MUST
         transition to the "I am Assert Loser" state and store the Assert
         Winner's address and metric.  If the metric was received in an
         Assert, the router MUST set the Assert Timer (AT(S,G,I)) to
         Assert_Time.  If the metric was received in a State Refresh, the
         router MUST set the Assert Timer (AT(S,G,I)) to three times the
         received State Refresh Interval.  If CouldAssert(S,G,I) == TRUE,
         the router MUST also multicast a Prune(S,G) to the Assert winner
         with a Prune Hold Time equal to the Assert Timer and evaluate any
         changes in its Upstream(S,G) state machine.

4.6.4.2.  Transitions from Winner State

   When in "I am Assert Winner" state, the following events trigger
   transitions:

      An (S,G) data packet arrives on downstream interface I
         An (S,G) data packet arrived on a downstream interface.  The
         Assert state machine remains in the "I am Assert Winner" state.
         The router MUST send an Assert(S,G) to interface I and set the
         Assert Timer (AT(S,G,I) to Assert_Time.

      Receive Inferior Assert or State Refresh
         An (S,G) Assert is received containing a metric for S that is
         worse than this router's metric for S.  Whoever sent the Assert
         is in error.  The router MUST send an Assert(S,G) to interface I
         and reset the Assert Timer (AT(S,G,I)) to Assert_Time.

      Receive Preferred Assert or State Refresh
         An (S,G) Assert or State Refresh is received that has a better
         metric than this router's metric for S on interface I.  The
         Assert state machine MUST transition to "I am Assert Loser" state
         and store the new Assert Winner's address and metric.  If the
         metric was received in an Assert, the router MUST set the Assert
         Timer (AT(S,G,I)) to Assert_Time.  If the metric was received in
         a State Refresh, the router MUST set the Assert Timer (AT(S,G,I))
         to three times the State Refresh Interval.  The router MUST also

        multicast a Prune(S,G) to the Assert winner, with a Prune Hold
        Time equal to the Assert Timer, and evaluate any changes in its
        Upstream(S,G) state machine.

     Send State Refresh
        The router is sending a State Refresh(S,G) message on interface
        I.  The router MUST set the Assert Timer (AT(S,G,I)) to three
        times the State Refresh Interval contained in the State
        Refresh(S,G) message.

     AT(S,G,I) Expires
        The (S,G) Assert Timer (AT(S,G,I)) expires.  The Assert state
        machine MUST transition to the NoInfo (NI) state.

     CouldAssert(S,G,I) -> FALSE
        This router's RPF interface changed, making CouldAssert(S,G,I)
        false.  This router can no longer perform the actions of the
        Assert winner, so the Assert state machine MUST transition to
        NoInfo (NI) state, send an AssertCancel(S,G) to interface I,
        cancel the Assert Timer (AT(S,G,I)), and remove itself as the
        Assert Winner.

4.6.4.3.  Transitions from Loser State

   When in "I am Assert Loser" state, the following transitions can
   occur:

     Receive Inferior Assert or State Refresh from Current Winner
        An Assert or State Refresh is received from the current Assert
        winner that is worse than this router's metric for S (typically,
        the winner's metric became worse).  The Assert state machine MUST
        transition to NoInfo (NI) state and cancel AT(S,G,I).  The router
        MUST delete the previous Assert Winner's address and metric and
        evaluate any possible transitions to its Upstream(S,G) state
        machine.  Usually this router will eventually re-assert and win
        when data packets from S have started flowing again.

     Receive Preferred Assert or State Refresh
        An Assert or State Refresh is received that has a metric better
        than or equal to that of the current Assert winner.  The Assert
        state machine remains in Loser (L) state.  If the metric was
        received in an Assert, the router MUST set the Assert Timer
        (AT(S,G,I)) to Assert_Time.  If the metric was received in a
        State Refresh, the router MUST set the Assert Timer (AT(S,G,I))
        to three times the received State Refresh Interval.  If the
        metric is better than the current Assert Winner, the router MUST

store the address and metric of the new Assert Winner, and if
CouldAssert(S,G,I) == TRUE, the router MUST multicast a
Prune(S,G) to the new Assert winner.

AT(S,G,I) Expires
  The (S,G) Assert Timer (AT(S,G,I)) expires.  The Assert state
  machine MUST transition to NoInfo (NI) state.  The router MUST
  delete the Assert Winner's address and metric.  If CouldAssert ==
  TRUE, the router MUST evaluate any possible transitions to its
  Upstream(S,G) state machine.

CouldAssert -> FALSE
  CouldAssert has become FALSE because interface I has become the
  RPF interface for S.  The Assert state machine MUST transition to
  NoInfo (NI) state, cancel AT(S,G,I), and delete information
  concerning the Assert Winner on I.

CouldAssert -> TRUE
  CouldAssert has become TRUE because interface I used to be the
  RPF interface for S, and now it is not.  The Assert state machine
  MUST transition to NoInfo (NI) state, cancel AT(S,G,I), and
  delete information concerning the Assert Winner on I.

Current Assert Winner's NeighborLiveness Timer Expires
  The current Assert winner's NeighborLiveness Timer (NLT(N,I)) has
  expired.  The Assert state machine MUST transition to the NoInfo
  (NI) state, delete the Assert Winner's address and metric, and
  evaluate any possible transitions to its Upstream(S,G) state
  machine.

Receive Prune(S,G), Join(S,G), or Graft(S,G)
  A Prune(S,G), Join(S,G), or Graft(S,G) message was received on
  interface I with its upstream neighbor address set to the
  router's address on I.  The router MUST send an Assert(S,G) on
  the receiving interface I to initiate an Assert negotiation.  The
  Assert state machine remains in the Assert Loser(L) state.  If a
  Graft(S,G) was received, the router MUST respond with a
  GraftAck(S,G).

4.6.5.  Rationale for Assert Rules

   The following is a summary of the rules for generating and processing
   Assert messages.  It is not intended to be definitive (the state
   machines and pseudocode provide the definitive behavior).  Instead,
   it provides some rationale for the behavior.

   1. The Assert winner for (S,G) must act as the local forwarder for
      (S,G) on behalf of all downstream members.
   2. PIM messages are directed to the RPF' neighbor and not to the
      regular RPF neighbor.
   3. An Assert loser that receives a Prune(S,G), Join(S,G), or
      Graft(S,G) directed to it initiates a new Assert negotiation so
      that the downstream router can correct its RPF'(S).
   4. An Assert winner for (S,G) sends a cancelling assert when it is
      about to stop forwarding on an (S,G) entry.  Example: If a router
      is being taken down, then a canceling assert is sent.

4.7.  PIM Packet Formats

   All PIM-DM packets use the same format as PIM-SM packets.  In the
   event of a discrepancy, PIM-SM [4] should be considered the
   definitive specification.  All PIM control messages have IP protocol
   number 103.  All PIM-DM messages MUST be sent with a TTL of 1.  All
   PIM-DM messages except Graft and Graft Ack messages MUST be sent to
   the ALL-PIM-ROUTERS group.  Graft messages SHOULD be unicast to the
   RPF'(S).  Graft Ack messages MUST be unicast to the sender of the
   Graft.

   The IPv4 ALL-PIM-ROUTERS group is 224.0.0.13.  The IPv6 ALL-PIM-
   ROUTERS group is 'ff02::d'.

4.7.1.  PIM Header

   All PIM control messages have the following header:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|PIM Ver| Type  |   Reserved    |           Checksum            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   PIM Ver PIM version number is 2.

   Type
     Types for specific PIM messages.  Available types are as follows:
     0 = Hello
     1 = Register (PIM-SM only)
     2 = Register Stop (PIM-SM only)
     3 = Join/Prune
     4 = Bootstrap (PIM-SM only)
     5 = Assert
     6 = Graft
     7 = Graft Ack
     8 = Candidate RP Advertisement (PIM-SM only)
     9 = State Refresh

   Reserved
     Set to zero on transmission.  Ignored upon receipt.

   Checksum
     The checksum is the standard IP checksum; i.e., the 16 bit one's
     complement of the one's complement sum of the entire PIM message.
     For computing checksum, the checksum field is zeroed.

     For IPv6, the checksum also includes the IPv6 "pseudo-header", as
     specified in RFC 2460, Section 8.1 [5].

4.7.2.  Encoded Unicast Address

   An Encoded Unicast Address has the following format:

    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
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   | Addr Family   | Encoding Type |      Unicast Address
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+...

   Addr Family
     The PIM Address Family of the 'Unicast Address' field of this
     address.  Values 0 - 127 are as assigned by the IANA for Internet
     Address Families in [9].  Values 128 - 250 are reserved to be
     assigned by the IANA for PIM specific Address Families.  Values 251
     - 255 are designated for private use.  As there is no assignment
     authority for this space; collisions should be expected.

   Encoding Type
     The type of encoding used with a specific Address Family.  The
     value '0' is reserved for this field and represents the native
     encoding of the Address Family.

   Unicast Address
      The unicast address as represented by the given Address Family and
      Encoding Type.

4.7.3.  Encoded Group Address

   An Encoded Group address has the following format:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|  Addr Family  | Encoding Type |B| Reserved  |Z|  Mask Len     |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                     Group Multicast Address
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+...
```

   Addr Family
      As described above.

   Encoding Type
      As described above.

   B
      Indicates that the group range should use Bidirectional PIM [16].
      Transmitted as zero; ignored upon receipt.

   Reserved
      Transmitted as zero.  Ignored upon receipt.

   Z
      Indicates that the group range is an admin scope zone.  This is
      used in the Bootstrap Router Mechanism [18] only.  For all other
      purposes, this bit is set to zero and ignored on receipt.

   Mask Len
      The mask length field is 8 bits.  The value is the number of
      contiguous left justified one bits used as a mask, which, combined
      with the address, describes a range of addresses.  It is less than
      or equal to the address length in bits for the given Address Family
      and Encoding Type.  If the message is sent for a single address
      then the mask length MUST equal the address length.  PIM-DM routers
      MUST only send for a single address.

   Group Multicast Address
      The address of the multicast group.

4.7.4.  Encoded Source Address

   An Encoded Source address has the following format:

```
    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
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |  Addr Family  | Encoding Type |   Rsrvd   |S|W|R|  Mask Len   |
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
   |                        Source Address
   +-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+...
```

   Addr Family
     As described above.

   Encoding Type
     As described above.

   Rsrvd
     Reserved.  Transmitted as zero.  Ignored upon receipt.

   S
     The Sparse Bit.  Set to 0 for PIM-DM.  Ignored upon receipt.

   W
     The Wild Card Bit.  Set to 0 for PIM-DM.  Ignored upon receipt.

   R
     The Rendezvous Point Tree bit.  Set to 0 for PIM-DM.  Ignored upon
     receipt.

   Mask Len
     As described above.  PIM-DM routers MUST only send for a single
     source address.

   Source Address
     The source address.

4.7.5.  Hello Message Format

   The PIM Hello message, as defined by PIM-SM [4], has the following
   format:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|PIM Ver| Type  |   Reserved    |           Checksum            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Option Type          |         Option Length         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         Option Value                          |
|                             ...                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                               .                               |
|                               .                               |
|                               .                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Option Type          |         Option Length         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                         Option Value                          |
|                             ...                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   PIM Ver, Type, Reserved, Checksum
     Described above.

   Option Type
     The type of option given in the Option Value field.  Available
     types are as follows:

```
      0               Reserved
      1               Hello Hold Time
      2               LAN Prune Delay
      3 - 16          Reserved
      17              To be assigned by IANA
      18              Deprecated and SHOULD NOT be used
      19              DR Priority (PIM-SM Only)
      20              Generation ID
      21              State Refresh Capable
      22              Bidir Capable
      23 - 65000      To be assigned by IANA
      65001 - 65535   Reserved for Private Use [9]
```

   Unknown options SHOULD be ignored.

4.7.5.1.  Hello Hold Time Option

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Type = 1           |           Length = 2          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|           Hold Time            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

Hold Time is the number of seconds a receiver MUST keep the neighbor
reachable.  If the Hold Time is set to '0xffff', the receiver of this
message never times out the neighbor.  This may be used with dial-
on-demand links to avoid keeping the link up with periodic Hello
messages.  Furthermore, if the Holdtime is set to '0', the
information is timed out immediately.  The Hello Hold Time option
MUST be used by PIM-DM routers.

4.7.5.2.  LAN Prune Delay Option

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Type = 2           |           Length = 4          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|T|      LAN Prune Delay         |        Override Interval       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

The LAN_Prune_Delay option is used to tune the prune propagation
delay on multi-access LANs.  The T bit is used by PIM-SM and SHOULD
be set to 0 by PIM-DM routers and ignored upon receipt.  The LAN
Delay and Override Interval fields are time intervals in units of
milliseconds and are used to tune the value of the J/P Override
Interval and its derived timer values.  Section 4.3.5 describes how
these values affect the behavior of a router.  The LAN Prune Delay
SHOULD be used by PIM-DM routers.

4.7.5.3.  Generation ID Option

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|            Type = 20          |           Length = 4          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                        Generation ID                          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   Generation ID is a random value for the interface on which the Hello
   message is sent.  The Generation ID is regenerated whenever PIM
   forwarding is started or restarted on the interface.  The Generation
   ID option MAY be used by PIM-DM routers.

4.7.5.4.  State Refresh Capable Option

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|            Type = 21          |           Length = 4          |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|  Version = 1  |    Interval   |            Reserved           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   The Interval field is the router's configured State Refresh Interval
   in seconds.  The Reserved field is set to zero and ignored upon
   receipt.  The State Refresh Capable option MUST be used by State
   Refresh capable PIM-DM routers.

4.7.6.  Join/Prune Message Format

   PIM Join/Prune messages, as defined in PIM-SM [4], have the following
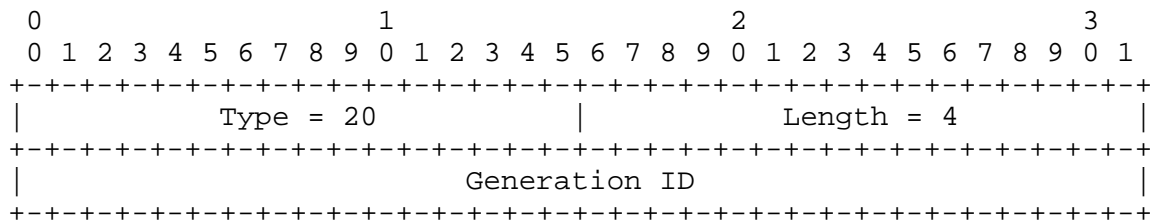   format:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|PIM Ver| Type  |   Reserved    |            Checksum           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Upstream Neighbor Address (Encoded Unicast Format)     |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Reserved    |  Num Groups   |          Hold Time            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+


+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Multicast Group Address 1 (Encoded Group Format)       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Number of Joined Sources    |   Number of Pruned Sources    |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Joined Source Address 1 (Encoded Source Format)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                              .                                |
|                              .                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Joined Source Address n (Encoded Source Format)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Pruned Source Address 1 (Encoded Source Format)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                              .                                |
|                              .                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Pruned Source Address n (Encoded Source Format)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                              .                                |
|                              .                                |
|                              .                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Multicast Group Address m (Encoded Group Format)       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Number of Joined Sources    |   Number of Pruned Sources    |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|        Joined Source Address 1 (Encoded Source Format)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                              .                                |
|                              .                                |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Joined Source Address n (Encoded Source Format)      |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Pruned Source Address 1 (Encoded Source Format)      |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                               .                               |
|                               .                               |
|                               .                               |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Pruned Source Address n (Encoded Source Format)      |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

PIM Ver, Type, Reserved, Checksum
  Described above.

Upstream Neighbor Address
  The address of the upstream neighbor.  The format for this address
  is given in the Encoded Unicast address in Section 4.7.2.  PIM-DM
  routers MUST set this field to the RPF next hop.

Reserved
  Transmitted as zero.  Ignored upon receipt.

Hold Time
  The number of seconds a receiving PIM-DM router MUST keep a Prune
  state alive, unless removed by a Join or Graft message.  If the
  Hold Time is '0xffff', the receiver MUST NOT remove the Prune state
  unless a corresponding Join or Graft message is received.  The Hold
  Time is ignored in Join messages.

Number of Groups
  Number of multicast group sets contained in the message.

Multicast Group Address
  The multicast group address in the Encoded Multicast address format
  given in Section 4.7.3.

Number of Joined Sources
  Number of Join source addresses listed for a given group.

Number of Pruned Sources
  Number of Prune source addresses listed for a given group.

Join Source Address 1..n
  This list contains the sources from which the sending router wishes
  to continue to receive multicast messages for the given group on
  this interface.  The addresses use the Encoded Source address
  format given in Section 4.7.4.

Prune Source Address 1..n
  This list contains the sources from which the sending router does
  not wish to receive multicast messages for the given group on this
  interface.  The addresses use the Encoded Source address format
  given in Section 4.7.4.

4.7.7.  Assert Message Format

   PIM Assert Messages, as defined in PIM-SM [4], have the following
   format:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|PIM Ver| Type  |   Reserved    |            Checksum           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Multicast Group Address (Encoded Group Format)       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|             Source Address (Encoded Unicast Format)           |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|R|                     Metric Preference                       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                             Metric                            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   PIM Ver, Type, Reserved, Checksum
     Described above.

   Multicast Group Address
     The multicast group address in the Encoded Multicast address format
     given in Section 4.7.3.

   Source Address
     The source address in the Encoded Unicast address format given in
     Section 4.7.2.

   R
     The Rendezvous Point Tree bit.  Set to 0 for PIM-DM.  Ignored upon
     receipt.

   Metric Preference
     The preference value assigned to the unicast routing protocol that
     provided the route to the source.

   Metric
     The cost metric of the unicast route to the source.  The metric is
     in units applicable to the unicast routing protocol used.
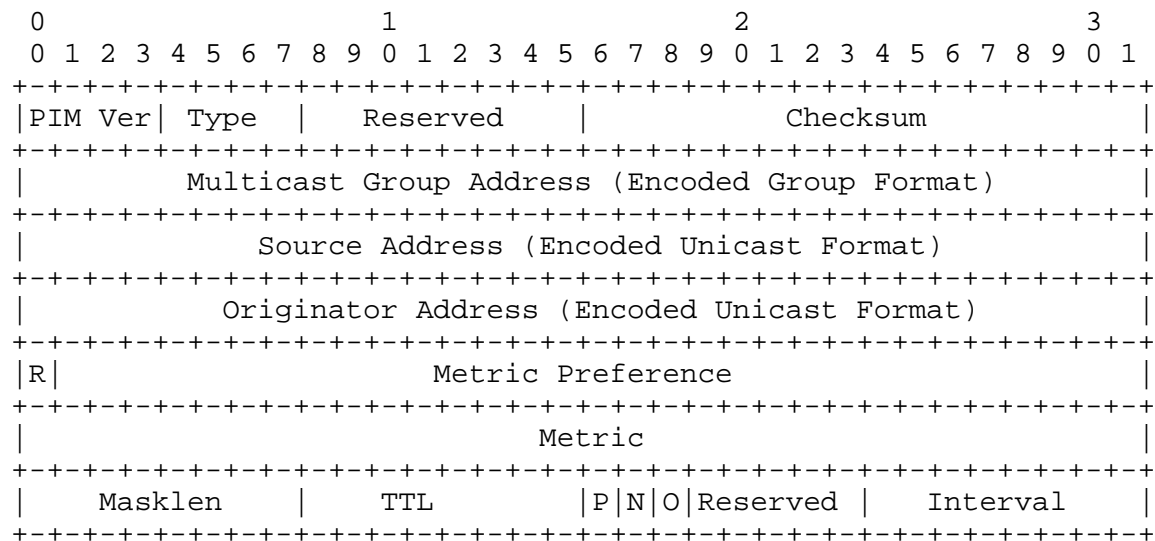
4.7.8.  Graft Message Format

   PIM Graft messages use the same format as Join/Prune messages, except
   that the Type field is set to 6.  The source address MUST be in the
   Join section of the message.  The Hold Time field SHOULD be zero and
   SHOULD be ignored when a Graft is received.

4.7.9.  Graft Ack Message Format

   PIM Graft Ack messages are identical in format to the received Graft
   message, except that the Type field is set to 7.  The Upstream
   Neighbor Address field SHOULD be set to the sender of the Graft
   message and SHOULD be ignored upon receipt.

4.7.10.  State Refresh Message Format

   PIM State Refresh Messages have the following format:

```
 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
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|PIM Ver| Type  |   Reserved    |           Checksum            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|         Multicast Group Address (Encoded Group Format)        |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|           Source Address (Encoded Unicast Format)            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|          Originator Address (Encoded Unicast Format)         |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|R|                     Metric Preference                       |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|                            Metric                            |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
|   Masklen     |     TTL       |P|N|O|Reserved |   Interval    |
+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+-+
```

   PIM Ver, Type, Reserved, Checksum
     Described above.

Multicast Group Address
   The multicast group address in the Encoded Multicast address format
   given in Section 4.7.3.

Source Address
   The address of the data source in the Encoded Unicast address
   format given in Section 4.7.2.

Originator Address
   The address of the first hop router in the Encoded Unicast address
   format given in Section 4.7.2.

R
   The Rendezvous Point Tree bit.  Set to 0 for PIM-DM.  Ignored upon
   receipt.

Metric Preference
   The preference value assigned to the unicast routing protocol that
   provided the route to the source.

Metric
   The cost metric of the unicast route to the source.  The metric is
   in units applicable to the unicast routing protocol used.

Masklen
   The length of the address mask of the unicast route to the source.

TTL
   Time To Live of the State Refresh message.  Decremented each time
   the message is forwarded.  Note that this is different from the IP
   Header TTL, which is always set to 1.

P
   Prune indicator flag.  This MUST be set to 1 if the State Refresh
   is to be sent on a Pruned interface.  Otherwise, it MUST be set to
   0.

N
   Prune Now flag.  This SHOULD be set to 1 by the State Refresh
   originator on every third State Refresh message and SHOULD be
   ignored upon receipt.  This is for compatibility with earlier
   versions of state refresh.

O
   Assert Override flag.  This SHOULD be set to 1 by upstream routers
   on a LAN if the Assert Timer (AT(S,G)) is not running and SHOULD be
   ignored upon receipt.  This is for compatibility with earlier
   versions of state refresh.

   Reserved
     Set to zero and ignored upon receipt.

   Interval
     Set by the originating router to the interval (in seconds) between
     consecutive State Refresh messages for this (S,G) pair.

4.8.  PIM-DM Timers

   PIM-DM maintains the following timers.  All timers are countdown
   timers -- they are set to a value and count down to zero, at which
   point they typically trigger an action.  Of course they can just as
   easily be implemented as count-up timers, where the absolute expiry
   time is stored and compared against a real-time clock, but the
   language in this specification assumes that they count downward
   towards zero.

   Global Timers
     Hello Timer: HT

     Per interface (I):
       Per neighbor (N):
         Neighbor Liveness Timer: NLT(N,I)

       Per (S,G) Pair:
         (S,G) Assert Timer: AT(S,G,I)
         (S,G) Prune Timer: PT(S,G,I)
         (S,G) PrunePending Timer: PPT(S,G,I)

       Per (S,G) Pair:
         (S,G) Graft Retry Timer: GRT(S,G)
         (S,G) Upstream Override Timer: OT(S,G)
         (S,G) Prune Limit Timer: PLT(S,G)
         (S,G) Source Active Timer: SAT(S,G)
         (S,G) State Refresh Timer: SRT(S,G)

When timer values are started or restarted, they are set to default
values.  The following tables summarize those default values.

Timer Name: Hello Timer (HT)
+----------------------+--------+--------------------------------------+
| Value Name           | Value  | Explanation                          |
+----------------------+--------+--------------------------------------+
|Hello_Period          | 30 sec | Periodic interval for hello messages |
+----------------------+--------+--------------------------------------+
|Triggered_Hello_Delay | 5 sec  | Random interval for initial Hello    |
|                      |        | message on bootup or triggered Hello |
|                      |        | message to a rebooting neighbor      |
+----------------------+--------+--------------------------------------+

Hello messages are sent on every active interface once every
Hello_Period seconds.  At system power-up, the timer is initialized
to rand(0,Triggered_Hello_Delay) to prevent synchronization.  When a
new or rebooting neighbor is detected, a responding Hello is sent
within rand(0,Triggered_Hello_Delay).

Timer Name: Neighbor Liveness Timer (NLT(N,I))
+-------------------+----------------+-----------------------------+
| Value Name        | Value          | Explanation                 |
+-------------------+----------------+-----------------------------+
| Hello Holdtime    | From message   | Hold Time from Hello Message|
+-------------------+----------------+-----------------------------+

Timer Name: PrunePending Timer (PPT(S,G,I))
+----------------------+-------------+----------------------------+
| Value Name           | Value       | Explanation                |
+----------------------+-------------+----------------------------+
| J/P_Override_Interval | OI(I) + PD(I) | Short time after a Prune to |
|                      |             | allow other routers on the |
|                      |             | LAN to send a Join         |
+----------------------+-------------+----------------------------+

The J/P_Override_Interval is the sum of the interface's
Override_Interval (OI(I)) and Propagation_Delay (PD(I)).  If all
routers on a LAN are using the LAN Prune Delay option, both
parameters MUST be set to the largest value on the LAN.  Otherwise,
the Override_Interval (OI(I)) MUST be set to 2.5 seconds, and the
Propagation_Delay (PD(I)) MUST be set to 0.5 seconds.

Timer Name: Prune Timer (PT(S,G,I))
+----------------+----------------+----------------------------------+
| Value Name     | Value          | Explanation                      |
+----------------+----------------+----------------------------------+
| Prune Holdtime | From message   | Hold Time read from Prune Message |
+----------------+----------------+----------------------------------+

Timer Name: Assert Timer (AT(S,G,I))
+------------------------+---------+------------------------------+
| Value Name             | Value   | Explanation                  |
+------------------------+---------+------------------------------+
| Assert Time            | 180 sec | Period after last assert before |
|                        |         | assert state is timed out    |
+------------------------+---------+------------------------------+

   Note that, for historical reasons, the Assert message lacks a
   Holdtime field.  Thus, changing the Assert Time from the default
   value is not recommended.  If all members of a LAN are state refresh
   enabled, the Assert Time will be three times the received
   RefreshInterval(S,G).

Timer Name: Graft Retry Timer (GRT(S,G))
+-------------------+-------+----------------------------------------+
| Value Name        | Value | Explanation                            |
+-------------------+-------+----------------------------------------+
| Graft_Retry_Period | 3 sec | In the absence of receipt of a GraftAck |
|                   |       | message, the time before retransmission |
|                   |       | of a Graft message                     |
+-------------------+-------+----------------------------------------+

Timer Name: Upstream Override Timer (OT(S,G))
+------------+---------------+----------------------------------------+
| Value Name | Value         | Explanation                            |
+------------+---------------+----------------------------------------|
| t_override | rand(0, OI(I)) | Randomized delay to prevent response   |
|            |               | implosion when sending a join message  |
|            |               | to override someone else's prune       |
+------------+---------------+----------------------------------------+

   t_override is a random value between 0 and the interface's
   Override_Interval (OI(I)).  If all routers on a LAN are using the LAN
   Prune Delay option, the Override_Interval (OI(I)) MUST be set to the
   largest value on the LAN.  Otherwise, the Override_Interval (OI(I))
   MUST be set to 2.5 seconds.

Timer Name: Prune Limit Timer (PLT(S,G))
+------------+-------------------+------------------------------------+
| Value Name | Value             | Explanation                        |
+------------+-------------------+------------------------------------|
| t_limit    | Default: 210 secs | Used to prevent Prune storms on a  |
|            |                   | LAN                                |
+------------+-------------------+------------------------------------+

Timer Name: Source Active Timer (SAT(S,G))
+----------------+-------------------+------------------------------------+
| Value Name     | Value             | Explanation                        |
+----------------+-------------------+------------------------------------+
| SourceLifetime | Default: 210 secs | Period of time after receiving     |
|                |                   | a multicast message a directly     |
|                |                   | attached router will continue      |
|                |                   | to send State Refresh messages     |
+----------------+-------------------+------------------------------------+

Timer Name: State Refresh Timer (SRT(S,G))
+----------------+-----------------+------------------------------------+
| Value Name     | Value           | Explanation                        |
+----------------+-----------------+------------------------------------+
| RefreshInterval | Default: 60 secs | Interval between successive        |
|                |                 | state refresh messages             |
+----------------+-----------------+------------------------------------+

5.  Protocol Interaction Considerations

   PIM-DM is designed to be independent of underlying unicast routing
   protocols and will interact only to the extent needed to perform RPF
   checks.  It is generally assumed that multicast area and autonomous
   system boundaries will correspond to the same boundaries for unicast
   routing, though a deployment that does not follow this assumption is
   not precluded by this specification.

   In general, PIM-DM interactions with other multicast routing
   protocols should be in compliance with RFC 2715 [7].  Other specific
   interactions are noted below.

5.1.  PIM-SM Interactions

   PIM-DM is not intended to interact directly with PIM-SM, even though
   they share a common packet format.  It is particularly important to
   note that a router cannot differentiate between a PIM-DM neighbor and
   a PIM-SM neighbor based on Hello messages.

   In the event that a PIM-DM router becomes a neighbor of a PIM-SM
   router, the two will effectively form a simplex link, with the PIM-DM
   router sending all multicast messages to the PIM-SM router while the
   PIM-SM router sends no multicast messages to the PIM-DM router.

   The common packet format permits a hybrid PIM-SM/DM implementation
   that would use PIM-SM when a rendezvous point is known and PIM-DM
   when one is not.  Such an implementation is outside the scope of this
   document.

5.2.  IGMP Interactions

   PIM-DM will forward received multicast data packets to neighboring
   host group members in all cases except when the PIM-DM router is in
   an Assert Loser state on that interface.  Note that a PIM Prune
   message is not permitted to prevent the delivery of messages to a
   network with group members.

   A PIM-DM Router MAY use the DR Priority option described in PIM-SM
   [14] to elect an IGMP v1 querier.

5.3.  Source Specific Multicast (SSM) Interactions

   PIM-DM makes no special considerations for SSM [15].  All Prunes and
   Grafts within the protocol are for a specific source, so no
   additional checks have to be made.

5.4.  Multicast Group Scope Boundary Interactions

   Although multicast group scope boundaries are generally identical to
   routing area boundaries, it is conceivable that a routing area might
   be partitioned for a particular multicast group.  PIM-DM routers MUST
   NOT send any messages concerning a particular group across that
   group's scope boundary.

6.  IANA Considerations

6.1.  PIM Address Family

   The PIM Address Family field was chosen to be 8 bits as a tradeoff
   between packet format and use of the IANA assigned numbers.  When the
   PIM packet format was designed, only 15 values were assigned for
   Address Families, and large numbers of new Address Families were not
   envisioned; 8 bits seemed large enough.  However, the IANA assigns
   Address Families in a 16 bit value.  Therefore, the PIM Address
   Family is allocated as follows:

Values 0 - 127 are designated to have the same meaning as IANA
assigned Address Family Numbers [9].

Values 128 - 250 are designated to be assigned by the IANA based on
IESG approval, as defined in [8].

Values 251 - 255 are designated for Private Use, as defined in [8].

## 6.2.  PIM Hello Options

Values 17 - 65000 are to be assigned by the IANA.  Since the space is
large, they may be assigned as First Come First Served, as defined in
[8].  Assignments are valid for one year and may be renewed.
Permanent assignments require a specification, as defined in [8].

## 7.  Security Considerations

The IPsec authentication header [10] MAY be used to provide data
integrity protection and groupwise data origin authentication of PIM
protocol messages.  Authentication of PIM messages can protect
against unwanted behaviors caused by unauthorized or altered PIM
messages.  In any case, a PIM router SHOULD NOT accept and process
PIM messages from neighbors unless a valid Hello message has been
received from that neighbor.

Note that PIM-DM has no rendezvous point, and therefore no single
point of failure that may be vulnerable.  Because PIM-DM uses unicast
routes provided by an unknown routing protocol, it may suffer
collateral effects if the unicast routing protocol is attacked.

## 7.1.  Attacks Based on Forged Messages

The extent of possible damage depends on the type of counterfeit
messages accepted.  We next consider the impact of possible
forgeries. A forged PIM-DM message is link local and can only reach a
LAN if it was sent by a local host or if it was allowed onto the LAN
by a compromised or non-compliant router.

1. A forged Hello message can cause multicast traffic to be delivered
   to links where there are no legitimate requestors, potentially
   wasting bandwidth on that link.  On a multi-access LAN, the
   effects are limited without the capability to forge a Join
   message, as other routers will Prune the link if the traffic is
   not desired.

2. A forged Join/Prune message can cause multicast traffic to be
   delivered to links where there are no legitimate requestors,
   potentially wasting bandwidth on that link.  A forged Prune

message on a multi-access LAN is generally not a significant
attack in PIM, because any legitimately joined router on the LAN
would override the Prune with a Join before the upstream router
stops forwarding data to the LAN.

3. A forged Graft message can cause multicast traffic to be delivered
   to links where there are no legitimate requestors, potentially
   wasting bandwidth on that link.  In principle, Graft messages
   could be sent multiple hops because they are unicast to the
   upstream router.  This should not be a problem, as the remote
   forger should have no way to get a Hello message to the target of
   the attack.  Without a valid Hello message, the receiving router
   SHOULD NOT accept the Graft.

4. A forged GraftAck message has no impact, as it will be ignored
   unless the router has recently sent a Graft to its upstream
   router.

5. By forging an Assert message on a multi-access LAN, an attacker
   could cause the legitimate forwarder to stop forwarding traffic to
   the LAN.  Such a forgery would prevent any hosts downstream of
   that LAN from receiving traffic.

6. A forged State Refresh message on a multi-access LAN would have
   the same impact as a forged Assert message, having the same
   general functions.  In addition, forged State Refresh messages
   would be propagated downstream and might be used in a denial of
   service attack.  Therefore, a PIM-DM router SHOULD rate limit
   State Refresh messages propagated.

7.2.  Non-cryptographic Authentication Mechanisms

   A PIM-DM router SHOULD provide an option to limit the set of
   neighbors from which it will accept PIM-DM messages.  Either static
   configuration of IP addresses or an IPSec security association may be
   used.  All options that restrict the range of addresses from which
   packets are accepted MUST default to allowing all packets.

   Furthermore, a PIM router SHOULD NOT accept protocol messages from a
   router from which it has not yet received a valid Hello message.

7.3.  Authentication Using IPsec

   The IPSec [10] transport mode using the Authentication Header (AH) is
   the recommended method to prevent the above attacks in PIM.  The
   specific AH authentication algorithm and parameters, including the
   choice of authentication algorithm and the choice of key, are
   configured by the network administrator.  The Encapsulating Security

Payload (ESP) MAY also be used to provide both encryption and
authentication of PIM protocol messages.  When IPsec authentication
is used, a PIM router SHOULD reject (drop without processing) any
unauthorized PIM protocol messages.

To use IPSec, the administrator of a PIM network configures each PIM
router with one or more Security Associations and associated Security
Parameters Indices that are used by senders to authenticate PIM
protocol messages and are used by receivers to authenticate received
PIM protocol messages.  This document does not describe protocols for
establishing Security Associations.  It assumes that manual
configuration of Security Associations is performed, but it does not
preclude the use of some future negotiation protocol such as GDOI
[17] to establish Security Associations.

The network administrator defines a Security Association (SA) and
Security Parameters Index (SPI) to be used to authenticate all PIM-DM
protocol messages from each router on each link in a PIM-DM domain.

In order to avoid the problem of allocating individual keys for each
neighbor on a link to each individual router, it is acceptable to
establish only one authentication key for all PIM-DM routers on a
link.  This will not specifically authenticate the individual router
sending the message, but will ensure that the sender is a PIM-DM
router on that link.  If this method is used, the receiver of the
message MUST ignore the received sequence number, thus disabling
anti-replay mechanisms.  The effects of disabling anti-replay
mechanisms are essentially the same as the effects of forged
messages, described in Section 7.1, with the additional protection
that the forger can only reuse legitimate messages.

The Security Policy Database at a PIM-DM router should be configured
to ensure that all incoming and outgoing PIM-DM packets use the SA
associated with the interface to which the packet is sent.  Note
that, according to [10], there is nominally a different Security
Association Database (SAD) for each router interface.  Thus, the
selected Security Association for an inbound PIM-DM packet can vary
depending on the interface on which the packet arrived.  This fact
allows the network administrator to use different authentication
methods for each link, even though the destination address is the
same for most PIM-DM packets, regardless of interface.

7.4.  Denial of Service Attacks

   There are a number of possible denial of service attacks against PIM
   that can be caused by generating false PIM protocol messages or even
   by generating false data traffic.  Authenticating PIM protocol
   traffic prevents some, but not all, of these attacks.  The possible
   attacks include the following:

   *  Sending packets to many different group addresses quickly can
      amount to a denial of service attack in and of itself.  These
      messages will initially be flooded throughout the network before
      they are pruned back.  The maintenance of state machines and State
      Refresh messages will be a continual drain on network resources.

   *  Forged State Refresh messages sent quickly could be propagated by
      downstream routers, creating a potential denial of service attack.
      Therefore, a PIM-DM router SHOULD limit the rate of State Refresh
      messages propagated.

8.  Acknowledgments

   The major features of PIM-DM were originally designed by Stephen
   Deering, Deborah Estrin, Dino Farinacci, Van Jacobson, Ahmed Helmy,
   David Meyer, and Liming Wei.  Additional features for state refresh
   were designed by Dino Farinacci, Isidor Kouvelas, and Kurt Windisch.
   This revision was undertaken to incorporate some of the lessons
   learned during the evolution of the PIM-SM specification and early
   deployments of PIM-DM.

   Thanks the PIM Working Group for their comments.

9.  References

9.1.  Normative References

   [1]  Deering, S., "Host extensions for IP multicasting", STD 5, RFC
        1112, August 1989.

   [2]  Fenner, W., "Internet Group Management Protocol, Version 2", RFC
        2236, November 1997.

   [3]  Cain, B., Deering, S., Kouvelas, I., Fenner, B., and A.
        Thyagarajan, "Internet Group Management Protocol, Version 3",
        RFC 3376, October 2002.

   [4]  Estrin, D., Farinacci, D., Helmy, A., Thaler, D., Deering, S.,
        Handley, M., Jacobson, V., Liu, C., Sharma, P., and L. Wei,
        "Protocol Independent Multicast-Sparse Mode (PIM-SM): Protocol
        Specification", RFC 2362, June 1998.

   [5]  Deering, S. and R. Hinden, "Internet Protocol, Version 6 (IPv6)
        Specification", RFC 2460, December 1998.

   [6]  Deering, S., Fenner, W., and B. Haberman, "Multicast Listener
        Discovery (MLD) for IPv6", RFC 2710, October 1999.

   [7]  Thaler, D., "Interoperability Rules for Multicast Routing
        Protocols", RFC 2715, October 1999.

   [8]  Narten, T. and H. Alvestrand, "Guidelines for Writing an IANA
        Considerations Section in RFCs", BCP 26, RFC 2434, October 1998.

   [9]  IANA, "Address Family Numbers", linked from
        http://www.iana.org/numbers.html.

   [10] Kent, S. and R. Atkinson, "Security Architecture for the
        Internet Protocol", RFC 2401, November 1998.

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

9.2.  Informative References

   [12] Deering, S.E., "Multicast Routing in a Datagram Internetwork",
        Ph.D. Thesis, Electrical Engineering Dept., Stanford University,
        December 1991.

   [13] Waitzman, D., Partridge, C., and S. Deering, "Distance Vector
        Multicast Routing Protocol", RFC 1075, November 1988.

   [14] Fenner,  W., Handley, M., Holbrook, H., and I. Kouvelas,
        "Protocol Independent Multicast - Sparse Mode (PIM-SM): Protocol
        Specification (Revised)", Work in Progress.

   [15] Holbrook, H. and B. Cain, "Source Specific Multicast for IP",
        Work in Progress.

   [16] Handley, M., Kouvelas, I., Speakman, T., and L. Vicisano, "Bi-
        directional Protocol Independent Multicast", Work in Progress.

   [17] Baugher, M., Weis, B., Hardjono, T., and H. Harney, "The Group
        Domain of Interpretation", RFC 3547, July 2003.

   [18] Fenner, W., Handley, M., Kermode, R., and D. Thaler, "Bootstrap
        Router (BSR) Mechanism for PIM Sparse Mode", Work in Progress.

Authors' Addresses

   Andrew Adams
   NextHop Technologies
   825 Victors Way, Suite 100
   Ann Arbor, MI 48108-2738

   EMail: ala@nexthop.com


   Jonathan Nicholas
   ITT Industries
   Aerospace/Communications Division
   100 Kingsland Rd
   Clifton, NJ   07014

   EMail: jonathan.nicholas@itt.com


   William Siadak
   NextHop Technologies
   825 Victors Way, Suite 100
   Ann Arbor, MI 48108-2738

   EMail: wfs@nexthop.com

Full Copyright Statement

   Copyright (C) The Internet Society (2005).

   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.

   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.

Intellectual Property

   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 IETF's procedures with respect to rights in IETF 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.

Acknowledgement

   Funding for the RFC Editor function is currently provided by the
   Internet Society.