SAM Research Group Internet-Draft Intended Status: Informational Expires: August 18, 2008

# Hybrid Overlay Multicast Framework draft-irtf-sam-hybrid-overlay-framework-02

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with <u>Section 6 of</u> <u>BCP 79</u>. This document may not be modified, and derivative works of it may not be created, except to publish it as an RFC and to translate it into languages other than English.

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

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

The list of current Internet-Drafts can be accessed at <a href="http://www.ietf.org/ietf/lid-abstracts.txt">http://www.ietf.org/ietf/lid-abstracts.txt</a>

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

This Internet-Draft will expire on August 18, 2008.

Copyright Notice

Copyright (C) The IETF Trust (2008).

Buford

Expires August 25, 2008

[Page 1]

# Abstract

We describe an experimental framework for constructing SAM sessions using hybrid combinations of Application Layer Multicast, native multicast, and multicast tunnels. We leverage AMT relay and gateway elements for interoperation between native regions and ALM regions. The framework allows different overlay algorithms and different ALM control algorithms to be used.

# Conventions used in this document

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in <u>RFC-2119</u> [1].

# Table of Contents

<u>2</u> . Definitions <u>2.1</u> . Overlay Network	. <u>4</u> . <u>4</u> . <u>4</u>
2.1. Overlay Network	. <u>4</u> . <u>4</u>
	.4
2.2. Overlay Multicast	
<u>2.3</u> . Peer	.4
2.4. Multi-Destination Routing	.5
<u>3</u> . Assumptions	. <u>5</u>
3.1. Overlay	.5
3.2. Overlay Multicast	.5
3.3. NAT	.6
3.4. Regions	.6
3.5. AMT	.6
4. ALM Tree Operations	.7
5. Hvbrid Connectivity	.8
6. Scenarios	.9
6.1. ALM-Only Tree - Algorithm 1	.9
6.2. ALM tree with peer at AMT site (AMT-GW)	10
6.3. ALM tree with NM peer using AMT-R	10
6.4. ALM tree with NM peer with P-AMT-R	11
6.5. Mixed Region Scenarios	11
7. Open Issues and Further Work	13
8. Security Considerations	13
9. TANA Considerations	13
10. References	13
10.1. Normative References	13
10.2 Informative References	14
Author's Address	15
Full Convright Statement	16
Intellectual Property	16

Buford Expires August 25, 2008

[Page 2]

## **1**. Introduction

The concept of scalable adaptive multicast [BUF2007] includes both scaling properties and adaptability properties. Scalability is intended to cover:

- o large group size
- o large numbers of small groups
- o rate of group membership change
- o admission control for QoS
- o use with network layer QoS mechanisms
- o varying degrees of reliability
- o trees connect nodes over global internet

Adaptability includes

- o use of different control mechanisms for different multicast trees depending on initial application parameters or application class
- o changing multicast tree structure depending on changes in application requirements, network conditions, and membership
- o use of different control mechanisms and tree structure in different regions of network depending on native multicast support, network characteristics, and node behavior

In this document we describe an experimental framework for constructing SAM sessions using hybrid combinations of Application Layer Multicast, native multicast, and multicast tunnels.

### 2. Definitions

#### 2.1. Overlay Network

P P P P P ..+...+...+...+...+... . +P P+ .. ..+...+...+...+... P P P P P

Overlay network - An application layer virtual or logical network in which end points are addressable and that provides connectivity, routing, and messaging between end points. Overlay networks are frequently used as a substrate for deploying new network services, or for providing a routing topology not available from the underlying physical network. Many peer-to-peer systems are overlay networks that run on top of the Internet.

In the above figure, "P" indicates overlay peers, and peers are connected in a logical address space. The links shown in the figure represent predecessor/successor links. Depending on the overlay routing model, additional or different links may be present.

## 2.2. Overlay Multicast

Overlay Multicast (OM): Hosts participating in a multicast session form an overlay network and utilize unicast connections among pairs of hosts for data dissemination. The hosts in overlay multicast exclusively handle group management, routing, and tree construction, without any support from Internet routers. This is also commonly known as Application Layer Multicast (ALM) or End System Multicast (ESM).

We call systems which use proxies connected in an overlay multicast backbone "proxied overlay multicast" or POM.

## 2.3. Peer

Peer: an autonomous end system that is connected to the physical network and participates in and contributes resources to overlay construction, routing and maintenance. Some peers may also perform additional roles such as connection relays, super nodes, NAT traversal, and data storage.

Buford

Expires August 25, 2008

[Page 4]

### **<u>2.4</u>**. Multi-Destination Routing

Multi-Destination Routing (MDR): A type of multicast routing in which group member's addresses are explicitly listed in each packet transmitted from the sender [<u>AGU1984</u>]. XCAST [<u>RFC5058</u>] is an experimental MDR protocol. A hybrid host group and MDR design is described in [<u>HE2005</u>].

## 3. Assumptions

### <u>3.1</u>. Overlay

Peers connect in a large-scale overlay, which may be used for a variety of peer-to-peer applications in addition to multicast sessions.

Peers may assume additional roles in the overlay beyond participation in the overlay and in multicast trees.

We assume a single structured overlay routing algorithm is used. Any of a variety of multi-hop, one-hop, or variable-hop overlay algorithms could be used.

Castro et al. [CAS2003] compared multi-hop overlays and found that tree-based construction in a single overlay out-performed using separate overlays for each multicast session. We use a single overlay rather than separate overlays per multicast sessions. We defer federated and hierarchical multi-overlay designs to later versions of this document.

Peers may be distributed throughout the network, in regions where native multicast (NM) is available as well as regions where it is not available.

An overlay multicast algorithm may leverage the overlay's mechanism for maintaining overlay state in the face of churn. For example, a peer may hold a number of DHT (Distributed Hash Table) entries. When the peer gracefully leaves the overlay, it transfers those entries to the nearest peer. When another peers joins which is closer to some of the entries than the current peer which holds those entries, than those entries are migrated. Overlay churn affects multicast trees as well; remedies include automatic migration of the tree state and automatic re-join operations for dislocated children nodes.

### 3.2. Overlay Multicast

The overlay supports concurrent multiple multicast trees. The limit on number of concurrent trees depends on peer and network resources

Buford

Expires August 25, 2008

and is not an intrinsic property of the overlay. Some multicast trees will contain peers use ALM only, i.e., the peers do not have NM connectivity. Some multicast trees will contain peers with a combination of ALM and NM. Although the overlay could be used to form trees of NM-only peers, if such peers are all in the same region we expect native mechanisms to be used for such tree construction, and if such peers are in different regions we expect AMT to handle most cases of interest.

Peers are able to determine, through configuration or discovery: o Can they connect to a NM router

- o Is an AMT gateway accessible
- o Can the peer support the AMT-GW functionality locally
- o Is MDR supported in the region

#### <u>3.3</u>. NAT

Some peers in the overlay may be in a private address space and behind firewalls. We assume that mechanisms are available for the following, and that the mechanisms scale as the ratio of NATed peers to public address (public) peers grows, to a limit.

- o Connectivity establishment between NATed peers and public peers
- Routing of overlay control messages to/from NATed and public peers.
- o Routing of data messages over the topology of the tree

NAT traversal solutions developed elsewhere in IETF will be used, and new NAT traversal mechanisms are out of scope to this framework.

### 3.4. Regions

A region is a contiguous internetwork such that if native multicast is available, all routers and end systems can connect to native multicast groups available in that region.

A region may include end systems.

### 3.5. AMT

We use AMT [<u>THA2007</u>] to connect peers in ALM region with peers in NM region. AMT permits AMT-R and AMT-GW functionality to be embedded in

Expires August 25, 2008

[Page 6]

hosts or specially configured routers. We assume AMT-R and AMT-GW can be implemented in peers.

AMT has certain restrictions: 1) isolated sites/hosts can receive SSM, 2) isolated non-NAT sites/hosts can send SSM, 3) isolated sites/hosts can receive general multicast. AMT does not permit isolated sites/hosts to send general multicast.

## **<u>4</u>**. ALM Tree Operations

Peers use the overlay to support ALM operations such as:

- o Create tree
- o Join
- o Leave
- o Re-Form or optimize tree

There are a variety of algorithms for peers to form multicast trees in the overlay. We permit multiple such algorithms to be supported in the overlay, since different algorithms may be more suitable for certain application requirements, and since we wish to support experimentation. Therefore, overlay messaging corresponding to the set of overlay multicast operations must carry algorithm identification information.

For example, for small groups, the join point might be directly assigned by the rendezvous point, while for large trees the join request might be propagated down the tree with candidate parents forwarding their position directly to the new node.

In addition to these overlay level tree operations, some peers may implement additional operations to map tree operations to native multicast and/or AMT [THA2007] connections.

+----+ +----+ AMT Site | P P P P P | Native MCast | 1 ++--+ +P . +---++ P+ AMT | AMT | . L GW | |RLY | +P . +--++ ++---+ . . +----+ +----+ . +----+ . . Native | | . MDR | . ....+...+P | P+...+P . | P . +----+ +----+ . | Native . MCast| . . . +----+ | P-AMT-R+ |Native Mcast | P+ . . ++---+ | P-AMT-R+ | P-AMT-GW+===|AMT | ...+...+.. . |RLY | P | .+...+....+...+ ++---+ +----+ P P P +-----+

# 5. Hybrid Connectivity

In the above figure we show the hybrid architecture in six regions of the network. All peers are connected in an overlay, and the figure shows the predecessor/successor links between peers. The peers may have other connections in the overlay.

o No native multicast: Peers (P) in this region connect to the overlay

Buford

Expires August 25, 2008

[Page 8]

- o Native multicast (NM) with a local AMT gateway (AMT GW). There are one or more peers (P) connected to the overlay in this region.
- o Native multicast with a local AMT relay (AMT RLY). There are one or more peers (P) connected to the overlay in this region.
- Native multicast with one or more peers which emulate the AMT relay behavior (P-AMT-R) which also connect to the overlay. There may be other peers (P) which also connect to the overlay.
- o Native MDR is a native multicast region using multi-destination routing, in which one or more peers reside in the region.
- o Native multicast with no peers that connect to the overlay, but for which there is at least one peer in the unicast-only part of the network which can behave as an AMT-GW (P-AMT-GW) to connect to multicast sources through an AMT-R for that region. It may be feasible to also allow non-peer hosts in such a region to participate as receivers of overlay multicast; for this version, we prefer to require all hosts to join the overlay as peers.

## 6. Scenarios

#### 6.1. ALM-Only Tree - Algorithm 1

Here is a simplistic algorithm for forming a multicast tree in the overlay. Its main advantage is use of the overlay routing mechanism for routing both control and data messages. The group creator doesn't have to be the root of the tree or even in the tree. It doesn't consider per node load, admission control, or alternative paths.

As stated earlier, multiple algorithms will co-exist in the overlay.

1. Peer which initiates multicast group:

groupID = create(); // allocate a unique groupId

// the root is the nearest peer in the overlay

// out of band advertisement/distribution of groupID, perhaps by
publishing in DHT

2. Any joining peer:

// out of band discovery of groupID, perhaps by lookup in DHT

joinTree(groupID); // sends "join groupID" message

Buford Expires August 25, 2008

[Page 9]

The overlay routes the join request using the overlay routing mechanism toward the peer with the nearest id to the groupID. This peer is the root. Peers on the path to the root join the tree as forwarding points.

3. Leave Tree:

leaveTree(groupID) // removes this node from the tree

Propagates a leave message to each child node and to the parent node. If the parent node is a forwarding node and this is its last child, then it propagates a leave message to its parent. A child node receiving a leave message from a parent sends a join message to the groupID.

4. Message forwarding:

multicastMsg(groupID, msg);

- o SSM tree The creator of the tree is the source. It sends data messages to the tree root which are forwarded down the tree.
- o ASM tree A node sending a data message sends the message to its parent and its children. Each node receiving a data message from one edge forwards it to remaining tree edges it is connected to.

#### 6.2. ALM tree with peer at AMT site (AMT-GW)

The joining peer connects to the tree using the ALM protocol, or, if the tree includes a peer in an NM region, then the peer can use the AMT GW to connect to the NM peer through the AMT relay. The peer can choose the delivery path based on latency and throughput.

If the peer is not a joining peer and is on the overlay path of a join request:

- o If its next hop is a peer in an NM region with AMT-R, then it can select either overlay routed multicast messages or AMT delivered multicast messages.
- o If its next hop is a peer outside of an NM region, then it could use either ALM only or use AMT delivery as an alternate path

#### 6.3. ALM tree with NM peer using AMT-R

There are these cases:

Buford Expires August 25, 2008 [Page 10]

o There is no peer in the tree which has an AMT-GW

The NM peer uses ALM routing

o There is at least one peer in the tree which can function as P-AMT-GW

The NM peer can join the tree using ALM routing and/or connecting to the P-AMT-GW.

o There is at least one peer in the tree which is in an AMT-GW region

The NM peer can join the tree using ALM routing and/or connecting to the AMT-GW.

#### 6.4. ALM tree with NM peer with P-AMT-R

Either the NM peer supports P-AMT-R or another peer in the multcast tree in the same region is P-AMT-R capable.

The three cases above apply here, replacing AMT-R with P-AMT-R.

### 6.5. Mixed Region Scenarios

In version 2 of this document we elaborate on:

- o ALM tree topology vs NM topology and NM-ALM edges
- o Single NM-ALM edge nodes vs multi NM peers from same region in the tree
- o Initial tree membership is ALM vs initial tree membership is NM

For ALM tree topology vs NM topology, all peers belong to the overlay, but only P-ALM peers use overlay routing for multicast data transmission. As a default behavior, a P-NM peer should generally prefer to join the tree via an AMT-GW node. But there may be special cases (small trees, short multicast sessions, trees where most of the members are known to be P-ALM) in which the peer can override this to specify an ALM-only join. A P-NM peer may also accept P-ALM children which don't use the AMT tunnel path to participate in the multicast tree.

Consider 3 types of tree links: P-ALM to P-ALM, P-NM to P-NM and P-ALM to/from P-NM:

Buford Expires August 25, 2008 [Page 11]

- o P-ALM to P-ALM This is a normal ALM tree path with management strictly in the overlay
- o P-NM to P-NM If the peers are in the same region, then the data path use native multicast capability in that region, and control occurs in ALM layer for ALM tree coordination and NM layer for native multicast purposes. If the peers are in different NM regions, then, if AMT gateways are available and configured to support an AMT tunnel between the regions, a tunnel is created using the AMT protocol (or already exists for this multicast group). The peers connect to their respective AMT gateways using the AMT procedure.
- o P-ALM to/from P-NM The connection can be either ALM or AMT tunnel depending on the context.

We expect two new functions are needed to build hybrid trees:

o joinViaAMTGateway(peer, AMT-GW, group\_id) where 'Peer' is the peer requesting to join the ALM group identified by group\_id, and AMT-GW is the ip address of the AMT gateway that the peer uses in its native multicast region. Request is transmitted to one or more parent peer candiates and/or rendezvous peers for the specified group id, according to the usual join protocol in this overlay. If the parent peer is a P-AMT-GW, then a tunnel is formed using the AMT protocol from the P-AMT-GW to the specified AMT-GW. If parent peer is a peer P-NM in native multicast region, then the

```
tunnel is created between P-NM's AMT-GW and the specified AMT-GW, using
```

the AMT protocol. If parent peer is a P-ALM, then the requested is propagated to other peers in the tree according to the join rules.

o leaveViaAMTGateway(peer, AMT-GW, group\_id)where 'Peer' is the peer requesting to leave the ALM group identified by group\_id, and AMT-GW is the ip address of the AMT gateway that the peer uses in its native multicast region. Request is transmitted the parent peer which is associated with the AMT-GW or provides that role. If the parent peer is a P-AMT-GW, then it removes the child from its AMT children list and may tear down the AMT tunnel P-AMT-GW to

the

specified AMT-GW if no other children are using it. If parent peer is a peer P-NM in native multicast region, then the tunnel is created between P-NM's AMT-GW and the specified AMT-GW, using the

AMT protocol.

Regarding initial tree membership being either P-NM or P-ALM node(s),

we expect the general case should be that hybrid tree formation is supported transparently regardless.

Buford

Expires August 25, 2008 [Page 12]

## 7. Open Issues and Further Work

- AMT [THA2007] has some restrictions on connecting isolated sites/hosts as SSM/ASM sources and receivers. Further analysis is needed to insure that OM data path is consistent with these constraints and whether additional operating restrictions between the overlay and AMT need be specified.
- o For NM regions with no AMT support, specifics of how peers selfselect as P-AMT-GW and P-AMT-RLY, and what additional behavior if any is needed beyond that specified in [THA2007].
- o We expect that the evolution of this document will lead to protocol specification related to the interopation points of the hybrid interfaces of the network.

## **<u>8</u>**. Security Considerations

Overlays are vulnerable to DOS and collusion attacks. We are not solving overlay security issues. For this version we assume centralized peer authentication model similar to what is proposed for P2P-SIP.

# 9. IANA Considerations

This document has no actions for IANA.

## **10**. References

## <u>**10.1</u>**. Normative References</u>

- [1] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, March 199
- [RFC0792] Postel, J., "Internet Control Message Protocol", STD 5, RFC 792, September 1981.
- [RFC3376] Cain, B., Deering, S., Kouvelas, I., Fenner, B., and A. Thyagarajan, "Internet Group Management Protocol, Version 3", <u>RFC 3376</u>, October 2002.
- [RFC3810] Vida, R. and L. Costa, "Multicast Listener Discovery Version 2 (MLDv2) for IPv6", <u>RFC 3810</u>, June 2004.

Buford

Expires August 25, 2008

- [RFC4605] Fenner, B., He, H., Haberman, B., and H. Sandick, "Internet Group Management Protocol (IGMP) / Multicast Listener Discovery (MLD)-Based Multicast Forwarding ("IGMP/MLD Proxying")", <u>RFC 4605</u>, August 2006.
- [RFC4607] Holbrook, H. and B. Cain, "Source-Specific Multicast for IP", <u>RFC 4607</u>, August 2006.
- [RFC5058] R. Boivie, N. Feldman , Y. Imai , W. Livens , D. Ooms, "Explicit Multicast (Xcast) Concepts and Options", IETF <u>RFC</u> <u>5058</u>. November 2007.

#### <u>10.2</u>. Informative References

- [AGU1984] L. Aguilar, Datagram Routing for Internet Multicasting, Sigcomm 84, March 1984.
- [BUF2007] J. Buford, S. Kadadi. SAM Problem Statement. Dec 2006. Internet Draft draft-irtf-sam-problem-statement-01.txt, work in progress.
- [CAS2002] M. Castro, P. Druschel, A.-M. Kermarrec, An. Rowstron, Scribe: A large-scale and decentralized application-level multicast infrastructure IEEE Journal on Selected Areas in Communications, Vol.20, No.8. October 2002.
- [CAS2003] M. Castro, M. Jones, A. Kermarrec, A. Rowstron, M. Theimer, H. Wang and A. Wolman, "An Evaluation of Scalable Application-level Multicast Built Using Peer-to-peer overlays," in Proceedings of IEEE INFOCOM 2003, April 2003.
- [HE2005] Q. He, M. Ammar. Dynamic Host-Group/Multi-Destination Routing for Multicast Sessions. J. of Telecommunication Systems, vol. 28, pp. 409-433, 2005.
- [MUR2006] E. Muramoto, Y. Imai, N. Kawaguchi. Requirements for Scalable Adaptive Multicast Framework in Non-GIG Networks. November 2006. Internet Draft <u>draft-muramoto-irtf-sam-</u> <u>generic-require-01.txt</u>, work in progress.
- [THA2007] D. Thale, M. Talwar, A. Aggarwal, L. Vicisano, T. Pusateri. Automatic IP Multicast Without Explicit Tunnels (AMT). Internet Draft <u>draft-ietf-mboned-auto-multicast-08</u>, Work in progress. Oct. 2007.

Buford Expires August 25, 2008 [Page 14]

Author's Address

John Buford Avaya Labs 307 Middletown-Lincroft Road Lincroft, NJ 07738 USA

Email: buford@samrg.org

Full Copyright Statement

Copyright (C) The IETF Trust (2008). This document is subject to the rights, licenses and restrictions contained in <u>BCP 78</u>, 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, THE IETF TRUST 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 procedures with respect to rights in RFC documents can be found in BCP 78 and BCP 79.

Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <a href="http://www.ietf.org/ipr">http://www.ietf.org/ipr</a>.

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.

## Acknowledgment

Funding for the RFC Editor function is provided by the IETF Administrative Support Activity (IASA).

Buford