



# *The Evolution of Multicast: From the MBone to Interdomain Multicast to Internet2 Deployment*

Kevin C. Almeroth, University of California

---

## Abstract

Multicast communication — the one-to-many or many-to-many delivery of data — is a hot topic. It is of interest in the research community, among standards groups, and to network service providers. For all the attention multicast has received, there are still issues that have not been completely resolved. One result is that protocols are still evolving, and some standards are not yet finished. From a deployment perspective, the lack of standards has slowed progress, but efforts to deploy multicast as an experimental service are in fact gaining momentum. The question now is how long it will be before multicast becomes a true Internet service. The goal of this article is to describe the past, present, and future of multicast. Starting with the Multicast Backbone (MBone), we describe how the emphasis has been on developing and refining intradomain multicast routing protocols. Starting in the middle to late 1990s, particular emphasis has been placed on developing interdomain multicast routing protocols. We provide a functional overview of the currently deployed solution. The future of multicast may hinge on several research efforts that are working to make the provision of multicast less complex by fundamentally changing the multicast model. We briefly survey these efforts. Finally, attempts are being made to deploy native multicast routing in both Internet2 networks and the commodity Internet. We examine how multicast is being deployed in these networks.

---

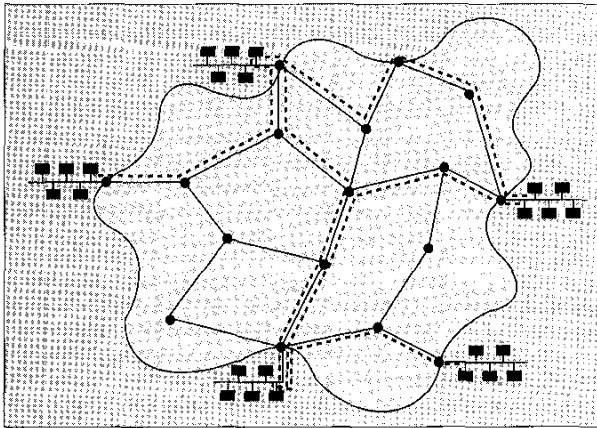


Without a doubt, multicast communication — the one-to-many or many-to-many delivery of data — has become a hot topic. It is the focus of intense study in the research community. It has become a highly desired feature of many vendors' network products. It is growing into a true deployment challenge for Internet engineers. It is evolving into a highly touted service being offered by some Internet service providers (ISPs). And finally, it is starting to be used by a number of companies offering large-scale Internet applications and services. From almost all perspectives, multicast is developing into one of the most interesting Internet services.

For all the potential multicast has, and for all the advocacy multicast has received, there are still some concerns. First, by Internet standards multicast is an old concept; yet by most measures, deployment has been very slow. To put deployment in perspective, compare multicast to the World Wide Web and HyperText Transfer Protocol (HTTP). IP multicast was first introduced in Steve Deering's Ph.D. dissertation in 1988 and tested on a wide scale during an "audiocast" at the 1992 Internet Engineering Task Force (IETF) meeting in San Diego [1]. The first Web browser was written in 1990, and in

1993 there were about 100 sites on the Web. So while multicast and the Web are roughly the same age, multicast is considered to be in the early stages of evolution [2], while the Web's success, influence, and use seem totally pervasive. Second, IP multicast is one of the first services to be deployed which requires additional "intelligence" in the network. Multicast requires a nontrivial amount of state and complexity in both core and edge routers. These requirements are at odds with the longstanding belief that intelligence should be pushed to the edges of the network. While many in the Internet community realize that the new generation of network services will put demands on the network, the difficulty is in deploying and managing these services in an infrastructure that has a lengthy history of only offering best-effort unicast service.

With these concerns in mind, the image of multicast may seem somewhat tarnished. Is multicast then more trouble than its efficiency gains and economies of scale are worth? This question is especially relevant if multicast is to be used as a money-making enterprise for commercial companies. The challenges are to define elegant protocols, to support an infrastructure on top of which new applications can be developed, and to continue to investigate new ways of increasing



■ Figure 1. A generic tunnel-based topology representative of the early MBone.

efficiency and reducing complexity. Doing multicast “the right way” is a noble endeavor and an appropriate long-term research topic, but the demand for working multicast has created an environment in which even short-term functional solutions are very attractive.

In this article we attempt to describe the past, present, and future of multicast. The history of multicast should help the reader understand how multicast has evolved into its current state. Relevant topics include a description of the Multicast Backbone (MBone) and an overview of the common *intradomain* multicast routing protocols. More recently, multicast evolution has been primarily focused in the area of *interdomain* protocol development. Multicast in the present can be characterized as an effort to deploy multicast on a wide scale using a triumvirate of routing protocols. These deployments have been carried out in the two Internet2 backbone networks — very-high-speed Backbone Network Service (vBNS) and Abilene — as well as in the commodity Internet (so designated in order to distinguish it from Internet2 networks). The future of multicast is rooted in the continued development, evaluation, and standardization of new protocols. However, unlike current efforts, which are focused primarily on routing, future efforts are likely to include other issues such as address allocation, management, and billing [3]. We are already starting to see some efforts in these areas.

The remainder of this article is organized as follows. We describe the early evolution of multicast, in particular the development of intradomain multicast. We then focus on interdomain multicast, including the best current practices and several of the efforts to define the next generation of protocols. We supply details on inter-domain deployment efforts in the commodity Internet and in Internet2 networks, and conclude the article.

### The Evolution of Intradomain Multicast

From the first Internet-wide experiments in 1992 to the middle of 1997, standardization and deployment in multicast focused on a single flat topology. This topology is in contrast to the Internet topology, which is based on a hierarchical routing structure. The initial multicast protocol research and standardization efforts were aimed at developing routing protocols for this flat topology. Beginning in 1997, when the multicast community realized the need for a hierarchical multicast infrastructure and interdomain routing, the existing protocols were categorized as intradomain protocols, and work began on standardizing an interdomain solution. In this section we describe the standard IP multicast model, and the evolution and characterization of intradomain multicast protocols.

### The Standard IP Multicast Model

Stephen Deering is responsible for describing the standard multicast model for IP networks [4]. This model describes how end systems are to send and receive multicast packets. The model includes both an explicit set of requirements and several implicit requirements. An understanding of the model will help the reader understand part of the evolutionary path multicast has taken. The model is as follows [5]:

- IP-style semantics. A source can send multicast packets at any time, with no need to register or to schedule transmission. IP multicast is based on User Datagram Protocol (UDP) (not TCP), so packets are delivered using a best-effort policy.
- Open groups. Sources only need to know a multicast address. They do not need to know group membership, and they do not need to be a member of the multicast group to which they are sending. A group can have any number of sources.
- Dynamic groups. Multicast group members can join or leave a multicast group at will. There is no need to register, synchronize, or negotiate with a centralized group management entity.

The standard IP multicast model is an end-system specification and does not discuss requirements for how the network should perform routing. The model also does not specify any mechanisms for providing quality of service, security, or address allocation.

### The Birth of the Multicast Backbone

Interest in building a multicast-capable Internet, motivated by Deering's work [4], began to achieve critical mass in the late 1980s. This work led to the creation of multicast in the Internet [6] and the creation of the Multicast Backbone (MBone) [7, 8]. In March 1992, the MBone carried its first worldwide event when 20 sites received audio from the meeting of the IETF [1] in San Diego. While the conferencing software itself represented a considerable accomplishment, the most significant achievement here was the deployment of a virtual multicast network. The multicast routing function was provided by workstations running a daemon process called *mrouted* (pronounced m-route-d), which received unicast-encapsulated multicast packets on an incoming interface and then forwarded packets over the appropriate set of outgoing interfaces. Connectivity among these machines was provided using point-to-point, IP-encapsulated *tunnels*. Each tunnel connected two endpoints via one logical link, but could cross several Internet routers. Once a packet is received, it can be sent to other tunnel endpoints or broadcast to local members. Routing decisions were made using the Distance Vector Multicast Routing Protocol (DVMRP) [9]. An example of connectivity provided via a virtual topology is shown in Fig. 1. In this earliest phase of the MBone, all tunnels were terminated on workstations, and the MBone topology was such that sometimes multiple tunnels ran over a common physical link. Multicast routing in the early MBone was actually a controlled form of flooding. The first versions of *mrouted* did not implement pruning. It was not until several years later that pruning was deployed.

The original multicast routing protocol, DVMRP, creates multicast trees using a technique known as *broadcast-and-prune*. Because of the way the tree is constructed by DVMRP, it is called a *reverse shortest path tree*. The steps to creating this type of tree are as follows:

- The source broadcasts each packet on its local network. An attached router receives the packet and sends it on all outgoing interfaces.
- Each router that receives a packet performs a reverse path forwarding (RPF) check. That is, each router checks to see if the incoming interface on which a multicast packet is received is the interface the router would use as an outgoing interface to reach the source. In this way, a router will

choose to only receive packets on the one interface that it believes is the most efficient path back to the source. All packets received on the proper interface are forwarded on all outgoing interfaces. All others are discarded silently.<sup>1</sup>

- Eventually a packet will reach a router with some number of attached hosts. This *leaf router* will check to see if it knows of any group members on any of its attached subnets. A router discovers the existence of group members by periodically issuing Internet Group Management Protocol (IGMP) [5, 10, 11] queries. If there are members, the leaf router forwards the multicast packet on the subnet. Otherwise, the leaf router will send a *prune message* toward the source on the RPF interface, that is, the interface the leaf router would use to forward packets to the source.
- Prune packets are forwarded back toward the source, and routers along the way create prune state for the interface on which the prune message is received. If prune messages are received on all interfaces except the RPF interface, the router will send a prune message of its own toward the source.

In this way, reverse shortest path trees are created. These trees can be constructed even on a virtual topology like the Mbone. Broadcast-and-prune protocols are also known as *dense mode* protocols, because they are designed to perform best when the topology is densely populated with group members. Routers assume there are group members downstream, and thus forward packets. Only when explicit prune messages are received does a router not forward multicast traffic. If a group is densely populated, routers are unlikely to ever need to prune. The key disadvantage of dense mode protocols is that state information must be kept for *each* source at *every* router in the network, regardless of whether downstream group members exist. If a group is not densely populated, significant state must be stored in the network, and a significant amount of bandwidth may be wasted.

#### The Evolution of Intradomain Multicast

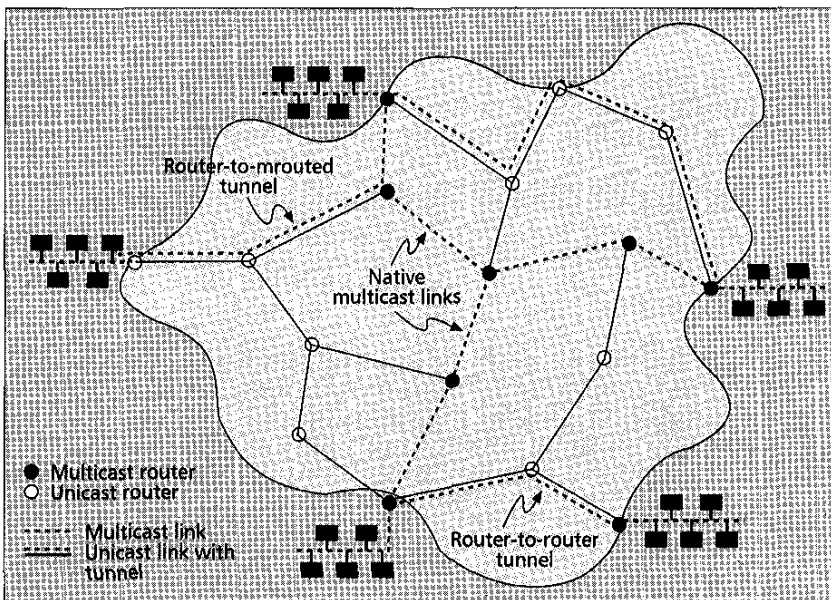
Since 1992, the Mbone has grown tremendously. It is no longer a simple virtual network sitting on top of the Internet, but is

rapidly becoming integrated into the Internet itself. In addition to simple DVMRP tunnels between workstations, the Mbone now has *native* multicast capability; that is, routers are capable of handling multicast packets (Fig. 2). Furthermore, ongoing research has led to the development and deployment of two additional dense mode protocols. These are described below.

**MOSPF** — Multicast Extensions to OSPF (MOSPF) [12] uses the Open Shortest Path First (OSPF) [13] protocol to provide multicast. Basically, MOSPF routers flood an OSPF area with information about group receivers. This allows all MOSPF routers in an area to have the same view of group membership. In the same way that each OSPF router independently constructs the unicast routing topology, each MOSPF router can construct the shortest-path tree for each source and group. While group membership reports are flooded throughout the OSPF area, data is not. MOSPF is something of an oddity in terms of classification. It is considered a dense mode protocol because membership information is broadcast to each MOSPF router, but it is also considered an explicit join protocol because data is only sent to those receivers that specifically request it. The key to understanding MOSPF is to realize that it is heavily dependent on OSPF and its link state routing paradigm.

**PIM-DM** — Protocol Independent Multicast (PIM) [14] has been split into two protocols, a dense mode version called PIM-DM [15] and a sparse mode version called PIM-SM [16]. PIM-DM is very similar to DVMRP; there are only two major differences. The first is that PIM (both dense mode and sparse mode) uses the unicast routing table to perform RPF checks. While DVMRP maintains its own routing table, PIM uses whatever unicast table is available. The name PIM is derived from the fact that the unicast table can be built using any unicast routing algorithm. PIM simply requires the unicast routing table to exist, and thus is *independent* of the algorithm used to build it. The second difference between PIM-DM and DVMRP is that DVMRP tries to avoid sending unnecessary packets to neighbors who will then generate prune messages based on a failed RPF check. The set of outgoing interfaces built by a given DVMRP router will include only those downstream routers that use the given router to reach the source (successful RPF check). PIM-DM avoids this complexity, but the trade-off is that packets are forwarded on all outgoing interfaces. Unnecessary packets are often forwarded to routers which must then generate prune messages because of the resulting RPF failure.

The next evolutionary step in intradomain routing was to develop protocols that addressed the disadvantages of dense mode protocols. A new class of protocols, called *sparse mode* protocols, was created. Instead of optimizing only for the case when a group has many members, sparse mode protocols are designed to work more efficiently when there are only a few



■ Figure 2. An example multicast topology with a combination of tunnels and native multicast links.

<sup>1</sup> In reality, the action for a packet that fails an RPF check depends on the protocol. Some protocols tell all upstream routers except the RPF router to stop forwarding packets.

widely distributed group members. Instead of broadcasting traffic and triggering prune messages, receivers are expected to send explicit join messages. These join messages are sent to a router acting as a core. Sources are expected to send their data traffic to this same node. The use of a core as a "meeting place" for sources and receivers facilitates creation of the multicast tree. Two of the most popular sparse mode protocols are described below.

**CBT** — The Core Based Trees (CBT) protocol was first discussed in the research community [17] and is now being standardized by the IETF [18]. CBT uses the basic sparse mode paradigm to create a single *shared tree* used by all sources. The tree is rooted at a core. All sources send their data to the core, and all receivers send explicit join messages to the core. There are two differences between CBT and PIM-SM. First, CBT uses only a shared tree, and is not designed to use shortest path trees. Second, CBT uses *bidirectional* shared trees, but PIM-SM uses *unidirectional* shared trees. Bidirectional shared trees involve slightly more complexity, but are more efficient when packets traveling from a source to the core cross branches of the multicast tree. In this case, instead of only sending traffic "up" to the core, packets can also be sent "down" the tree. While CBT has significant technical merits and is on par technically with PIM-SM, few routing vendors provide support for CBT.

**PIM-SM** — PIM-SM [16] is much more widely used than CBT. It is similar to PIM-DM in that routing decisions are based on whatever underlying unicast routing table exists, but the tree construction mechanism is quite different. PIM-SM's tree construction algorithm is actually more similar to that used by CBT than to that used by PIM-DM. In the following description of sparse mode protocol operation, we use PIM-SM as our example.

- A core, called a rendezvous point (RP) in PIM terminology, must be configured.<sup>2</sup> Different groups may use different routers for RPs, but a group can only have a single RP.
  - Information about which routers in the network are RPs, and the mappings of multicast groups to RPs, must be discovered by all routers.
  - RP discovery is done using a bootstrap protocol. However, because the RP discovery mechanism is not included in the PIM-SMv1 specification, each vendor implementation of PIM-SMv1 has its own RP discovery mechanism. For PIM-SMv2, the bootstrap protocol is included in the protocol specification.
  - The basic function of the bootstrap protocol, in addition to RP discovery, is to provide robustness in case of RP failure. The bootstrap protocol includes mechanisms to select an alternate RP if the primary RP goes down.
- Receivers send explicit *join messages* to the RP. Forwarding state is created in each router along the path from the receiver to the RP. A single shared tree, rooted at the RP, is formed for each group. As with other multicast protocols, the tree is a reverse shortest path tree — join messages follow a reverse path from receivers to the RP.
- Each source sends multicast data packets, encapsulated in unicast packets, to the RP. When an RP receives one of these *register packets*, a number of actions are possible. First, if the RP has forwarding state for the group (i.e., there are receivers who have joined the group), the encapsulation is stripped off the packet, and it is sent on the shared tree. However, if the RP does not have forwarding state for the group, it sends a *register-stop message* to the RP. This avoids wasting bandwidth between the source and the RP. Second, the RP may wish to send a join message toward the source. By establishing multicast forwarding state between the source and the RP, the RP can receive the source's traffic as multicast and avoid the overhead of encapsulation.

These steps describe the basic mechanism used by sparse mode protocols in general and PIM-SM in particular. In summary, the basic goal is to use the RP as a "meeting place" for sources and receivers. Receivers explicitly join the shared tree, and sources register with the RP.

Sparse mode protocols have a number of advantages over dense mode protocols. First, sparse mode protocols typically offer better scalability in terms of routing state. Only routers on the path between a source and a group member must keep state. Dense mode protocols require state in all routers in the network. Second, sparse mode protocols are more efficient because the use of explicit join messages means multicast traffic only flows across links that have been explicitly added to the tree.

Sparse mode protocols do have a few disadvantages. These are mostly related to the use of RPs. First, the RP can be a single point of failure. Second, the RP can become a hot spot for multicast traffic. Third, having traffic forwarded from a source to the RP and then to receivers means that nonoptimal paths may exist in the multicast tree. The first problem is mostly solved with the bootstrap router protocol. The second and third problems are solved in CBT by using bidirectional trees. PIM-SM solves these problems by providing a mechanism to switch from a shared tree to a shortest path tree. This change occurs when a leaf router sends a special message toward the source. Forwarding state is changed so that traffic flows directly to the receiver instead of first through the RP. This action occurs when a traffic rate threshold is violated.

Finally, not only has progress been made in protocol development, but Mbone growth has led to increased user awareness of multicast, which in turn has led to demand for new applications and better support for real-time data. Improvements have been made in transport layer protocols. For example, the Real-Time Protocol (RTP) [20] assists loss- and delay-sensitive applications in adapting to the Internet's best-effort service model. With respect to applications, the Mbone has seen an increasingly diverse set of media types. Originally, the Mbone was considered a research effort, and its evolution was overseen by members of the Mbone community. Coordination of events was handled almost exclusively through the use of a global session directory tool, originally called *sd*, but now called *sdr*. As multicast deployment has continued, and as multicast has been integrated into the Internet as a native service, the informal use agreements and guidelines have faded. Even though *sdr*-based sessions remain at the core of Internet multicast events, their percentage of the total is shrinking. Other applications are being deployed that do not coordinate sessions through *sdr* or use RTP. This potpourri of tools has enriched the diversity of applications available, but has stressed the ability of the network to provide multicast according to the standard IP multicast model.

For clarity, it is worth summarizing the key multicast terminology. Multicast protocols use either a *broadcast-and-prune* or an *explicit join* mechanism. Broadcast-and-prune protocols are commonly called **dense mode protocols** and always use a **reverse shortest path tree** rooted at a source. Explicit join protocols, commonly called **sparse mode protocols**, can use either a reverse shortest path tree or a *shared tree*. A shared tree uses a **core** or **rendezvous point** to bring sources and receivers together.

<sup>2</sup> Deciding how many RPs to have and where to place them in the network is a network planning issue and is beyond the scope of this article. A recent book offers some discussion on this topic [19].

## Problems with Multicast

As the MBone has grown, it has suffered from an increasing number of problems, and these problems have been occurring with increasing frequency. The most important reason for this is the growing difficulty of managing a flat virtual topology. The same problems experienced with class-based unicast routing have manifested themselves in the MBone. As the MBone has grown, its size has become a problem, in terms of both routing state and susceptibility to misconfigurations. As a result, the multicast community has realized the need to deploy hierarchical interdomain routing. In particular, the MBone faces problems of scalability and manageability.

**Scalability** — Large, flat networks are inherently unstable. Exacerbating this problem are organizational mechanisms which do not provide significant route aggregation. For these two reasons, the MBone has experienced substantial scalability problems. At its peak, the MBone had almost 10,000 routes. Unfortunately, most of these routes had long prefixes (between /28 and /32), which meant that very few hosts could be represented in each routing table entry. These scalability problems are not new. As the Internet has grown, unicast routing had to be fundamentally changed to enable continued growth and stability. The solutions — route aggregation and hierarchical routing — have proven successful, and the issue now is how to apply them to multicast.

**Manageability** — As the MBone has grown, it has become harder to manage. The MBone has no central management, and most tasks have been handled on a per-site basis. Most coordination takes place via the MBone mailing list. Because the MBone is a virtual topology and new sites can be connected anywhere, there should be a formal procedure for adding new sites. Because no such mechanism exists, the MBone has grown randomly, and there are many inefficiencies. Two types of inefficiency commonly observed are:

- **Virtual topology (tunnel) management.** The MBone is characterized as a set of multicast-capable islands connected by tunnels. The goal has always been to connect these islands in the most efficient manner, but over time suboptimal tunnels have been created. Tunnels are often set up in very inefficient ways (Fig. 1 for several examples). This behavior was observed very early in the history of the MBone, especially with regard to the MCI Backbone. To avoid the growing tangle of tunnels, engineers at MCI undertook the difficult task of enforcing a policy that tunnels through or into the MCI network would have to be terminated at designated border points. The goal was to resolve the observed problem of single physical links being crossed by several (up to 10) tunnels. The work of the MCI engineers set an example that helped keep the MBone reasonably efficient for a number of years.
- **Interdomain policy management.** Domain boundaries are another source of problems when trying to manage a flat topology. The model in today's Internet is to establish autonomous system (AS) boundaries between Internet domains. ASes are commonly managed or owned by different organizations. Entities in one AS are typically not trusted by entities in another AS. As a result, exchange of routing information across AS boundaries is handled very carefully. Peering relationships among ASes are provisioned using the Border Gateway Protocol (BGP), which provides routing abstraction and policy control [21–23]. As a result of widescale use of BGP, there is a commonly accepted procedure when two ASes wish to communicate. Because the MBone does not provide such an interdomain protocol, it offers no protection across domain boundaries. When

there is a single flat topology connected using tunnels, routing problems can easily spread throughout the topology.

To summarize, the first problem is the complexity and instability of a large flat topology. The second problem is that there are no protocol mechanisms to build a hierarchical multicast routing topology. The need to solve these two problems created the first attempts to deploy interdomain multicast.

## The Evolution of Interdomain Multicast

Interdomain multicast has evolved out of the need to provide scalable, hierarchical, Internet-wide multicast. Protocols that provide the necessary functionality have been developed, but the technology is relatively immature. These protocols are being considered by the IETF, while simultaneously being evaluated through extensive deployment. The particular interdomain solution in use is considered near-term, and is possibly only an interim solution. While the solution is functional, it lacks elegance and long-term scalability. As a result, additional work is underway to find long-term solutions. Some of these proposals are based on the standard IP multicast model. Others attempt to refine the service model in hopes of making the problem easier.

### The Near-Term Solution

The near-term solution for interdomain multicast routing has three parts. The first is a straightforward extension of the interdomain unicast route exchange protocol, BGP. The second and third are additional protocols needed to build and interconnect trees across domain boundaries.

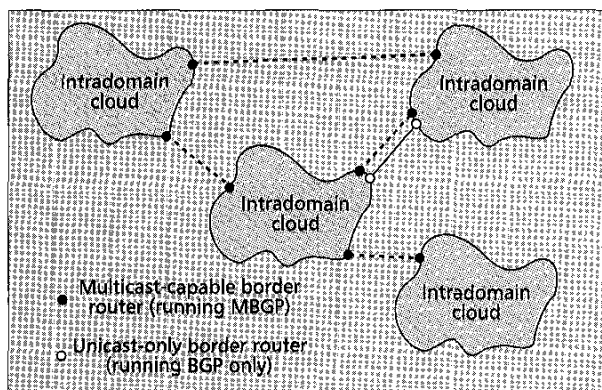
**Carrying Multicast Routes in BGP** — The first requirement follows from the need to make multicast routing hierarchical in the same manner as unicast routing. Route aggregation and abstraction, as well as hop-by-hop policy routing, are provided in unicast using BGP [22]. BGP offers substantial abstraction and control among domains. Within a domain, a network administrator can run any routing protocol desired. Routing to hosts in an external domain is simply a matter of choosing the best external link.

BGP supports interdomain routing by reliably exchanging network reachability information. This information is used to compute an end-to-end distance-vector-style path of AS numbers. Each AS advertises the set of routes it can reach and an associated cost. Each border router can then compute the set of ASes that should be traversed to reach any network. The use of a distance vector algorithm together with full path information allows BGP to overcome many of the limitations of traditional distance vector algorithms. Packets are still routed on a hop-by-hop basis, but less information is needed and better routing decisions can be made.

The functionality provided by BGP, and its well-understood paradigm for connecting ASes, are important catalysts for supporting interdomain multicast. A version of BGP capable of carrying multicast routes would not only provide hierarchical routing and policy decisions, but would also allow a service provider to use different topologies for unicast and multicast traffic.

The mechanism by which BGP has been extended to carry multicast routes is called Multiprotocol Extensions to BGP4 (MBGP) [24].<sup>3</sup> MBGP is able to carry multiprotocol routes by adding the Subsequent Address Family Identifier (SAFI) to two BGP4 messages: MP\_REACH\_NLRI and MP\_UNREACH\_NLRI. Specifically for multicast, the SAFI field can specify unicast, multicast, or unicast/multicast. With MBGP, instead of

<sup>3</sup> There is some ambiguity over terminology here. First, multiprotocol BGP4 is sometimes also referred to as BGP4+. Second, some think that MBGP stands for Multicast BGP. All three terms refer to the same protocol.



■ Figure 3. Example inter-domain multicast topology running BGP and/or MBGP.

every router needing to know the entire flat multicast topology, each router only needs to know the topology of its own domain and the paths to reach each of the other domains. Figure 3 shows an example of several domains connected together by MBGP sessions. In one case, two domains are connected together using different connections for unicast and multicast.

There is some confusion over exactly what functionality MBGP provides. To be clear, we offer the following example. If one domain advertises reachability for multicast, the message will say, "I have a path to sources on the networks listed in this message." MBGP messages do not carry information about multicast groups (i.e., class D addresses are never carried in an MBGP message). Recall that multicast trees are constructed using a reverse path back to the source. Therefore, MBGP information is used when a join message is sent from an RP or receiver toward the source. This join message needs to know the best reverse path toward the source. MBGP provides this next-hop information between domains. If all unicast and multicast topologies were assumed to be the same, the reverse path join could simply follow the same next hop that any unicast traffic would follow. MBGP allows a network administrator to specify a different reverse path for the join to follow, and (subsequently) a different forward path when data is sent.

While MBGP is the first step toward providing interdomain multicast, it alone is not a complete solution. MBGP is capable of determining the next hop to a host, but not of providing multicast tree construction functions. More specifically, what is the format of the join message? When should join messages be sent, and how often? Support for this functionality is not provided by MBGP; a true interdomain multicast routing protocol is needed. Furthermore, conventional wisdom suggests that this protocol should not use the broadcast-and-prune method of tree construction. The near-term solution being advocated is to use PIM-SM, to establish a multicast tree between domains containing group members.

*The Multicast Source Discovery Protocol* — To summarize: various intradomain routing protocols exist, there is a route exchange protocol to support multicast, and PIM-SM is to be used to connect receivers and sources across domain boundaries. But there is still one function missing from the near-term solution. This function is needed when trying to connect sparse mode domains together.

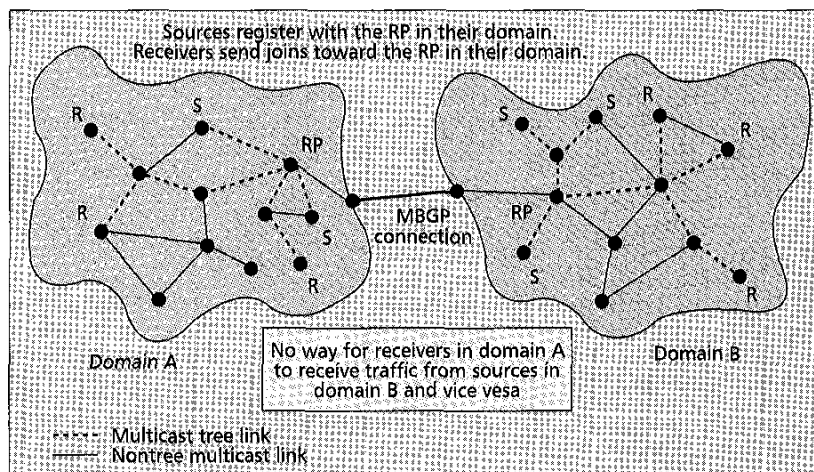
er. Given that PIM-SM is the only sparse mode protocol that has seen significant deployment, this function tends to be heavily influenced by PIM-SM. The problem is basically how to inform an RP in one domain that there are sources in other domains. The underlying assumption here is that a group can now have multiple RPs. However, the reality is that there is still only one RP per domain, but now multiple domains may be involved. The approach adopted is largely motivated by the perceived needs of the ISP community. In fact, the decision to have multiple RPs rather than a single root is what differentiates the near-term solution from other proposed solutions.

A problem arises when group members are spread over multiple domains. There is no mechanism to connect the various intradomain multicast trees together. While traffic from all the sources for a particular group *within a particular domain* will reach the group's receivers, any sources outside the domain will remain disjoint. Why is this the case? Within a domain, receivers send join messages toward one RP, and sources send register messages to the same RP. However, there is no way for an RP in one domain to find out about sources in other domains using different RPs. There is no mechanism for RPs to communicate with each other when one receives a source register message. This problem is summarized in Fig. 4.

The decision to maintain a separate multicast tree and RP for each domain is driven by the need to reduce administrative dependencies between domains. Two potential problems are avoided this way:

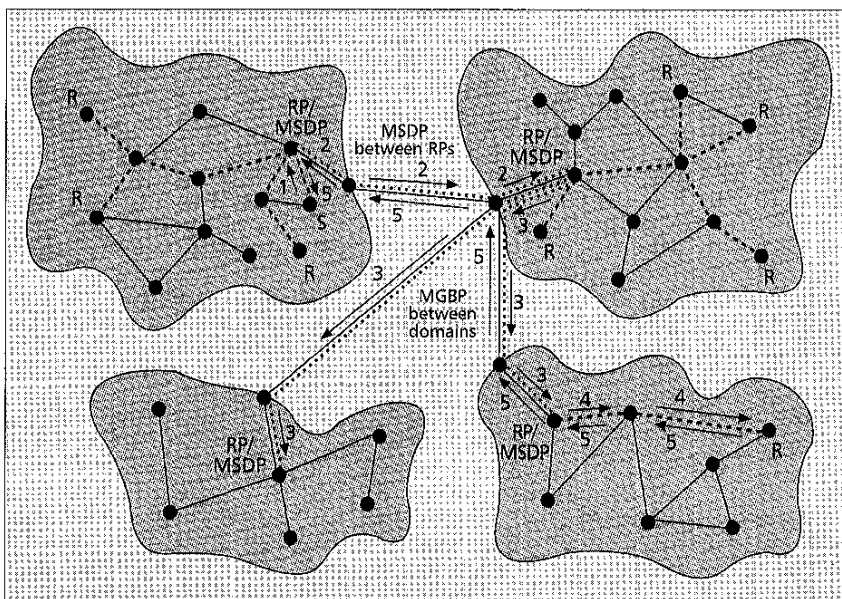
- It is not necessary for two domains to co-administer a single sparse mode cloud. Relevant administrative functions include identifying candidate RPs and establishing the group-RP mapping.
- It becomes possible to avoid second- and third-party dependencies, in which multicast delivery for sources and groups in one or more domains is dependent on another domain whose only function is to provide the RP. Dependencies can occur when all sources and receivers in the RP's domain leave or become inactive. The domain with the RP has no group members, but is still providing the RP service. Depending on how multicast and interdomain traffic billing is handled, this could be particularly undesirable.

The near-term solution adopted for this problem is a new protocol, appropriately named the Multicast Source Discovery Protocol (MSDP) [25]. This protocol works by having representatives in each domain announce to other domains the existence of active sources. MSDP is run in the same router as



■ Figure 4. The problem of connecting sources and receivers across two sparse mode domains.





■ Figure 5. MSDP operation, including the flow of Source Active messages.

a domain's RP (or one of the RPs). MSDP's operation is similar to that of MBGP, in that MSDP sessions are configured between domains and TCP is used for reliable session message exchange. MSDP operation is described below, with each step shown in Fig. 5:

- 1) When a new source for a group becomes active it will register with the domain's RP.
- 2) The MSDP peer in the domain will detect the existence of the new source and send a Source Active (SA) message to all directly connected MSDP peers.
- 3) MSDP message flooding:
  - MSDP peers that receive an SA message will perform a *peer-RPF check*. The MSDP peer that received the SA message will check to see if the MSDP peer that sent the message is along the "correct" MSDP-peer path. These peer-RPF checks are necessary to prevent SA message looping.
  - If an MSDP peer receives an SA message on the correct interface, the message is forwarded to all MSDP peers except the one from which the message was received. This is called *peer-RPF flooding*.
- 4) Within a domain, an MSDP peer (also the RP) will check to see if it has state for any group members in the domain. If state does exist, the RP will send a PIM join message to the source address advertised in the SA message.
- 5) If data is contained in the message, the RP then forwards it on the multicast tree. Once group members receive data, they may choose to switch to a shortest path tree using PIM-SM conventions.
- 6) Steps 3–5 are repeated until all MSDP peers have received the SA message and all group members are receiving data from the source.

This ends the description of the short-term interdomain multicast routing solution. The solution is referred to with the abbreviations for the three relevant protocols: MBGP/PIM-SM/MSDP. However, while the given description is relatively complete, there are a number of details which are not discussed. And as with any system, most of the complexity is in the details. Furthermore, we have not yet discussed the limitations of the current solution in any detail. In particular, a qualitative assessment of the scalability, complexity, and overall quality of the protocols would be valuable.

The MBGP/PIM-SM/MSDP solution is relatively straight-

forward once a person understands all the abbreviations and understands the motivating factors that drove the design of the protocols. While some argue that the current set of protocols is not simple, it really is no more complex than many other Internet services, such as unicast routing. The key advantage of MBGP/PIM-SM/MSDP is that it is a functional solution largely built on existing protocols. Furthermore, it is already being deployed with a fair amount of success. The key disadvantage is that, as a long-term solution, the MBGP/PIM-SM/MSDP protocol suite may be susceptible to scalability problems. Further discussion of two particular problems follows.

**MSDP and Dynamic Groups** — When multicast sources begin to transmit, the network is required to create some type of routing state to control packet

flow. We have already discussed how different types of multicast routing protocols accomplish this function. However, in the case of MSDP, information about the existence of sources must first be transmitted before routing state can be created. This extra complexity increases the overhead of managing groups. When groups are dynamic, due to either bursty sources or frequent group member join/leave events, the overhead of managing the group can be significant.<sup>4</sup> A formidable task would be created for networks that must establish and remove information for thousands of sources and receivers scattered around the world. Two specific problems related to dynamic groups/sources are:

- **Join latency.** Because SA messages are only sent periodically, there may be a significant delay between when new receivers join and when they hear the next SA message. To solve this problem, MSDP peers may be configured to cache SA messages. A noncaching MSDP peer can send an SA-Request message to an MSDP peer that does perform caching. This gives MSDP peers a mechanism to actively determine source, thereby reducing join latency. The trade-off is the extra state and complexity of maintaining the cache.
- **Bursty sources.** This type of source can be characterized as sending short packet bursts separated by silent periods on the order of several minutes. One example is when a tool like *sdr* is used to periodically advertise a session. A problem occurs when trying to establish a multicast tree for this kind of source. The problem begins when one or a few packets are sent to the RP. The RP will hear the packet and flood an SA message, and RPs in other domains will send join messages back to the source. However, because no multicast forwarding state existed when the packet was originally sent, and because it takes time to forward SA messages and have other RPs establish forwarding state, the original burst will not reach new receivers. Once state is established, all subsequent packets should reach these receivers. The problem occurs when the period of silence between packet bursts exceeds the forwarding state timeout value (typically 3 min). Because

<sup>4</sup> Again, it should be noted that because no formal study of MBGP/PIM-SM/MSDP performance has been conducted, many of these statements are hypothetical.

no packets are sent, the forwarding state is discarded. When another session announcement is sent, the same process of establishing state but losing the initial burst is repeated. In this way, no packets from bursty sources ever reach group members. The solution, specified in the MSDP protocol, is to have SA messages carry the first  $n$  data packets. This is not a particularly elegant solution, but it does solve the problem. The lack of elegance is making the protocol harder to standardize. Because data packets are delivered via SA messages, which are delivered over TCP connections, some in the multicast community wonder if this will have undesirable side effects or break assumptions of higher-layer protocols. As a result, recent discussions in the MSDP working group have generated proposals which allow data to be carried in either GRE or UDP packets. The final decision on which data delivery options to support has not been made.

**MSDP Scalability** — The issue of scalability is an important one to consider for MSDP. Because of the way MSDP operates, if multicast becomes tremendously successful, the overhead of MSDP may become too large. The limitation occurs if multicast use grows to the point where there are thousands of multicast sources. The number of SA messages (plus data) being flooded around the network could become very large. The generally-agreed-upon conclusion is that MSDP is not a particularly scalable solution, and will likely be insufficient for the long term. But, given that long-term solutions are not ready to be deployed, MSDP is seen as an immediate solution to an immediate need.

#### *Long-Term Proposals*

While MBGP/PIM-SM/MSDP is a recognized near-term solution, there is still a need to develop long-term solutions. Numerous efforts are being undertaken in this direction. These efforts can be broken down into two groups: those based on the standard IP multicast philosophy, and those that look to change this model in hopes of simplifying the problem. Efforts in each of these areas are described next.

**Border Gateway Multicast Protocol** — The Border Gateway Multicast Protocol (BGMP) was first proposed as a long-term solution for Internet-wide interdomain multicast [25].<sup>5</sup> The key idea of BGMP is to construct bidirectional shared trees between domains using a single root. One of the functions of BGMP is then to decide in which particular domain to root the shared tree. BGMP relies on the belief that interdomain dependencies can be avoided by using a strict address allocation scheme. Such an address allocation scheme allows domains to own specific addresses or specific ranges of addresses. The belief is that if a particular domain owns the address for a particular group, the domain will be significantly involved in the multicast service. Finally, this means that dependency problems, even though there is a single root, should be highly unlikely. For example, a video-on-demand application will likely be rooted at the server; a video conference group will be rooted at the primary source or at a session coordinator. The belief is that no matter the type of session, one domain will always be the logical choice for the root domain.

As a result of a protocol like BGMP, there is a need for a strict address allocation scheme. Strict means that ownership

must be clearly defined, and that there cannot be collisions. Therefore, the *sdr* mechanism of randomly choosing an address is not sufficient. Because of BGMP, as well as demands from ISPs and application writers, work is being conducted to develop the necessary address allocation schemes. Before discussing two of the proposals for address allocation, it is worthwhile to make two points. First, BGMP is relatively flexible, and can use any scheme as long as it provides strict address allocation. Second, independent of BGMP, there is a need for better address allocation. The *sdr* mechanism is not particularly scalable and is no longer sufficient, even for the current MBGP/PIM-SM/MSDP solution. Proposals, usable in both the current model and with BGMP, are being considered by the IETF. They are described below.

**MASC** — The Multicast Address-Set Claim (MASC) protocol supports address allocation between domains [26]. MASC includes mechanisms to guarantee that address collisions are immediately resolved. From a more abstract perspective, MASC provides the functionality required at the highest layer of a more general addressing scheme called the Multicast Address Allocation Architecture (MAAA) [27]. MASC and its supporting protocols are specific instances of protocols that meet the requirements of the MAAA specification. In MAAA, there are three levels of address allocation: at the domain level, within a domain, and between hosts and the network. Work to develop protocols at each level is underway in the IETF. MASC would act as a top-level address allocation protocol and operate between domains; the multicast Address Allocation Protocol (AAP) [28] would allocate addresses within a domain; and the Multicast Address Dynamic Client Allocation Protocol (MADCAP) [29] would be used by hosts to request addresses from a multicast address allocation server (MAAS).

**GLOP** — Another, much simpler, proposal is to statically allocate multicast addresses to each AS. A "glop" of addresses is assigned to each AS. The AS number is encoded as part of the address [30]. The first version of GLOP is being evaluated with only part of the 224/4 address range. Only the 233/8 address range is being used. As a result, the first octet is static, the next two octets encode the AS number, and the final octet provides a range of addresses to be allocated. This proposal is gaining in popularity, but it has two limitations. First, because only 8 bits, or 256 addresses, are available to each AS, there is likely to be an insufficient number of addresses per AS. This problem could be solved by using more of the class D address space, or switching to IPv6 addressing. The second problem is that GLOP does not specify a mechanism by which addresses are allocated within the domain. This problem could be solved by using a simple administrative procedure, a dynamic protocol like AAP/MADCAP, or a modified intradomain version of *sdr*.

**The Root Addressed Multicast Architecture** — In response to the perceived complexity of MBGP/PIM-SM/MSDP and BGMP, and the need to address additional multicast-related issues such as security, billing, and management [3], some members of the multicast community are looking to make fundamental changes in the multicast model. One class of proposals being offered is called the Root Addressed Multicast Architecture (RAMA) [31]. The premise for RAMA-style protocols is that most multicast applications are single-source or have an easily identifiable primary source. By making this source the root of the tree, the complexity of core placement in other multicast routing protocols can be eliminated. This trade-off raises a

<sup>5</sup> BGMP should not be confused with MBGP. After reading this article the differences should be obvious, but the similarity of the names and abbreviations has led to constant confusion. Furthermore, BGMP was recently renamed. It was previously known as Grand Unified Multicast (GUM).



number of important issues which are described at the end of this section. There are two primary RAMA-style protocols being discussed: Express Multicast [32] and Simple Multicast [33]. The key aspects of these two protocols are discussed below.

**Express Multicast** — Express is designed specifically as a single-source protocol. The root of the tree is placed at the source, and group members send join messages along the reverse path to the source. Express also provides mechanisms to efficiently collect information about subscribers. The protocol is specifically designed for subscriber-based systems that use logical channels. Representative applications include TV broadcasts,

file distribution, and any single-source multimedia application. The key advantages of Express are that routing complexity can be reduced and that *closed groups* can be offered.

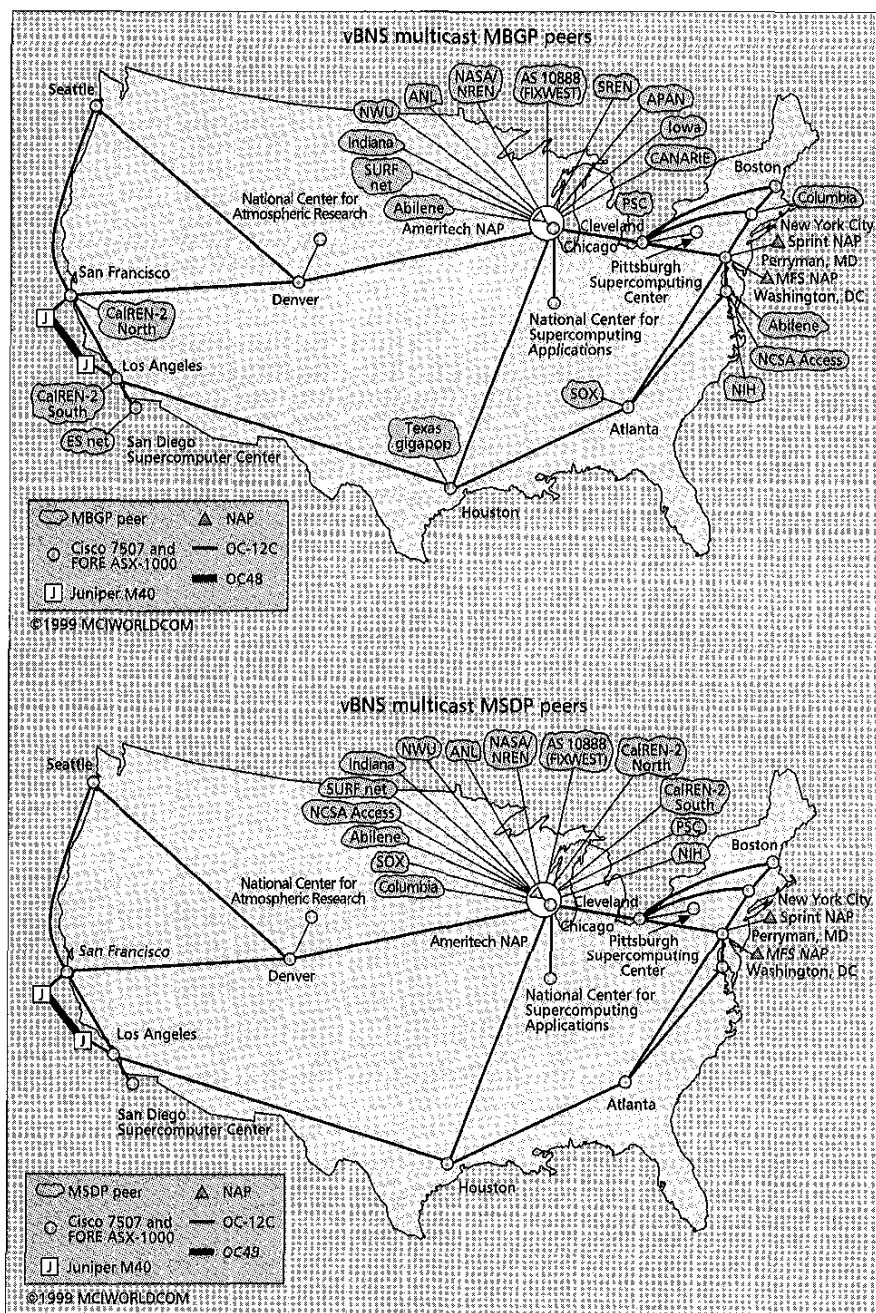
**Simple Multicast** — Simple Multicast and Express Multicast are similar, but Simple Multicast has the added flexibility of allowing multiple sources per group. A particular source must be chosen as the primary, and the tree is rooted at this node's first-hop router. Receivers send join messages to the source, and a bidirectional tree is constructed. Additional sources send packets to the primary source. Because the tree is bidirectional, as soon as packets reach a router in the tree they are forwarded both downstream to

receivers and upstream to the core. The advantages and disadvantages of this proposal are being heavily debated, but the proposal's authors believe that it eliminates the address allocation problem and the need to place and locate RPs. Address allocation is done by using the core address and multicast group address to uniquely identify a group. By routing on this pair of addresses, each root/core/source can allocate, without collision, up to  $2^{32}$  addresses.

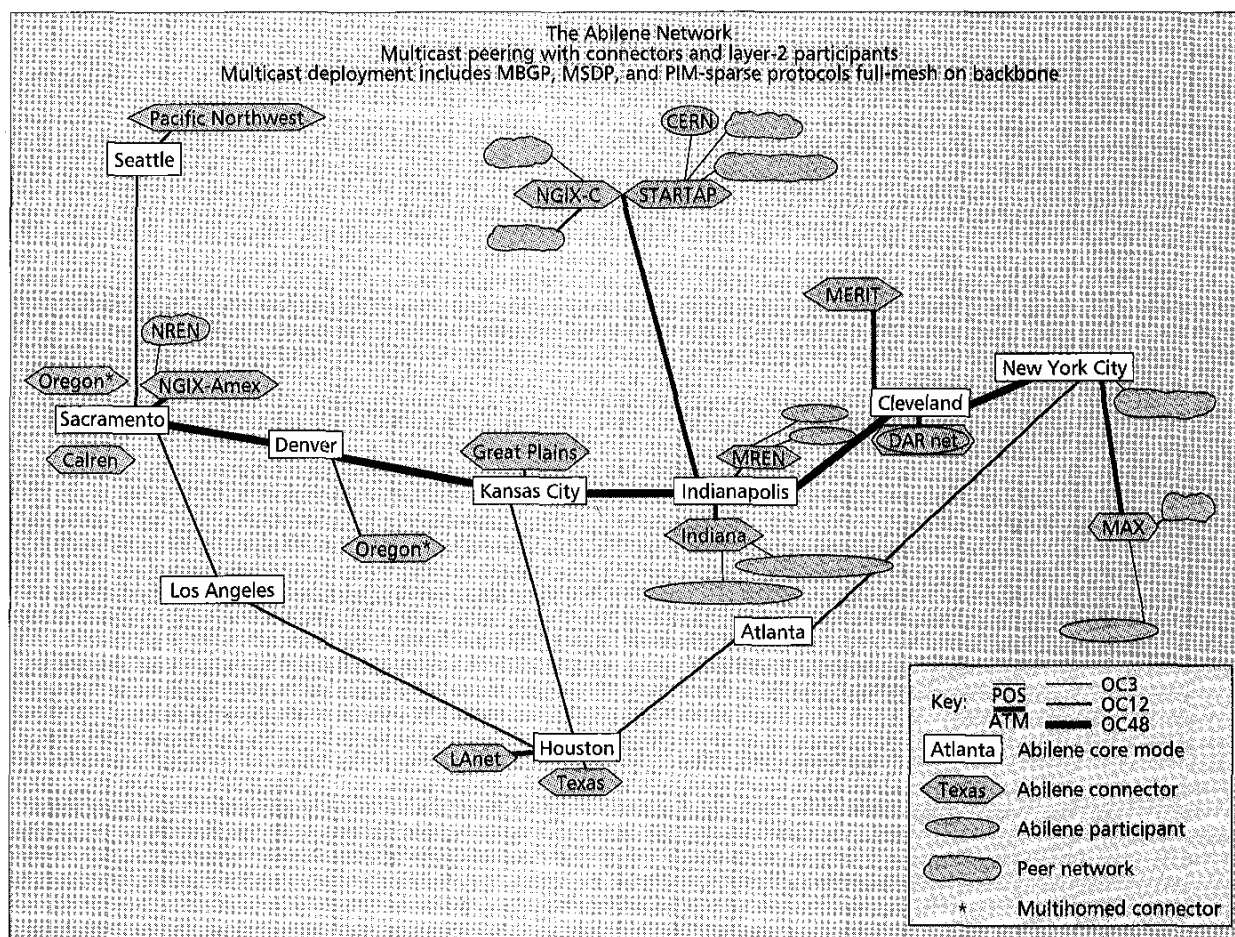
The Express and Simple multicast proposals have received significant attention in both the research community and the IETF. There is another question in addition to that of the merits of these new protocols. If these protocols are standardized, will they be expected to replace all existing protocols, or will they work in parallel with the existing multicast infrastructure? If the RAMA-style protocols are expected to work in cooperation with existing protocols, there will be yet another set of protocols to deploy, evaluate, and interoperate with. This does not make the provision of Internet-wide multicast easier. If RAMA-style protocols are expected to replace the current set of protocols, the question becomes whether they have enough flexibility to support all types of multicast applications. The bottom line is that these new protocols are still proposals, and it is uncertain what their future will be.

### Interdomain Multicast Deployment

The successful deployment of multicast, or lack thereof, was one of the original motivations for developing interdomain routing protocols. In this section we describe efforts to deploy these protocols. Our description is divided into two parts: a discussion of the commodity Internet, and one of the Internet2 architecture.



■ Figure 6. The vBNS MBGP and MSDP peering topology.



■ Figure 7. The Abilene multicast map.

### Deployment in the Commodity Internet

Measuring the success of interdomain deployment, either from a qualitative point of view or by taking a count of connected hosts, is a difficult problem. Published studies have so far only dealt with the MBone, although several studies that distinguish between the MBone and interdomain multicast are currently underway. It is beyond the scope of this article to offer any quantitative results. However, it is possible to describe the plan, now being implemented, to transition from the MBone's flat virtual topology to a true interdomain multicast infrastructure.

Now that interdomain multicast routing is possible, the issue is how to deal with the MBone. While the rest of the Internet is working to deploy interdomain multicast, the challenge is how to bring MBone users into the new infrastructure. The solution has been to make the MBone its own AS, called AS10888. All MBone tunnels and sites connected by tunnels are relegated to AS10888. Connectivity between AS10888 and other multicast-capable ASes is provided at the NASA Ames Multicast-Friendly Internet Exchange (MIX) [34]. The NASA Ames MIX provides connectivity between the MBone (AS10888) and all other ASes that have deployed MBGP/PIM-SM/MSDP. The deployment of interdomain multicast can continue to grow while the flat routing topology that is the MBone is eliminated. Sites on the MBone will hopefully transition to native multicast by deploying whatever interdomain solution is appropriate. When this occurs, these sites will no longer need their old MBone tunnels. Observational analysis suggests that this transition process is indeed occurring. Because of the differences in

route aggregation between MBGP routes and MBone routes, it is difficult to quantify this assertion. However, the number of routes in the MBone has decreased dramatically, and the number of MBGP routes has increased dramatically.

### Deployment in Internet2

For Internet2, the plan has always been to try and do multicast "the right way" to the extent possible given the currently available set of protocols. As a result, Internet2 multicast deployment is following guidelines set forth by the Internet2 Multicast Working Group. Briefly, these guidelines require all multicast deployed in Internet2 to be native and sparse mode. No tunnels are allowed, and all routers must support interdomain multicast routing using MBGP/MSDP. To date, Internet2 has experienced a reasonable amount of success in deploying multicast. This success includes backbone deployment, connecting other high-speed networks, connecting member institutions, and running several high-bandwidth (on the order of 30 Mb/s) multicast applications.

There are two Internet2 backbones in the United States. One is vBNS [35, 36] and the other is Abilene. vBNS has been in existence since 1995, and from a very early stage has had basic dense mode capability. During the 1998 Internet2 Member Meeting in San Francisco, the inherent problems of dense mode protocols were painfully realized when tens of megabits of traffic were flooded across the network. As a result, vBNS engineers worked hard to transition the network to PIM-SM and MBGP/MSDP. As of mid-1999, the network had successfully deployed interdomain multicast, and was in the process of establishing MBGP and MSDP peering relationships with other

networks. Figure 6 shows the topology of vBNS, including the existing MBGP and MSDP peering relationships. As vBNS engineers and other network operators gain experience in using MBGP and MSDP, the rate at which new MBGP/MSDP peerings are added will increase. A number of additional networks, including several international high-speed networks, are planning to connect to vBNS in the very near future.

The other Internet2 backbone is the Abilene network. Because Abilene is a newer network and only recently (February 1999) became operational, the state of interdomain multicast in Abilene is not nearly as advanced as in vBNS. However, Abilene has been running PIM-SM since mid-1999, and has begun to establish its first set of interdomain peering relationships. The challenge has been to climb the learning curve and establish multicast capability in the backbone. Now that the first MBGP/PIM-SM/MSDP peering relationships have been established, additional peerings are being added rapidly. The current topology is shown in Fig. 7.

## Conclusions

In this article we present a tutorial-style overview of multicast. We cover the early development of intradomain routing protocols, the evolution of the MBone, the needs and current solutions for interdomain multicast, the set of next-generation protocols currently under investigation, and the current state of deployment in the Internet and Internet2. Whatever the future holds for multicast, it is likely to present major challenges for both research and deployment.

## Acknowledgments

This article would not have been possible without the expertise of many, many people. It would be quite impossible for any list to do justice to them all. However, specific technical and qualitative suggestions were offered by Dave Meyer, Dave Thaler, Mostafa Ammar, Christophe Diot, Brian Levine, and Ben Chinowsky. Also, engineers from both the vBNS and Abilene contributed topologies for their respective networks. These individuals included Kevin Thompson, Steven Wallace, and Brent Sweeny.

## References

- [1] S. Casner and S. Deering, "First IETF Internet Audiocast," *ACM Comp. Commun. Rev.*, July 1992, pp. 92-97.
- [2] K. Almeroth, "A Long-Term Analysis of Growth and Usage Patterns in the Multicast Backbone (MBone)," to appear, *IEEE Infocom*, Tel Aviv, Israel, Mar. 2000.
- [3] C. Diot et al., "Deployment Issues for the IP Multicast Service and Architecture," *IEEE Network*, this issue.
- [4] S. Deering, "Multicast Routing in a Datagram Internetwork," Ph.D. dissertation, 1991.
- [5] S. Deering and D. Cheriton, "Multicast Routing in Datagram Internetworks and Extended LANs," *ACM Trans. Comp. Sys.*, May 1990, pp. 85-111.
- [6] S. Deering, "Host Extensions for IP Multicasting," RFC 1112, Aug. 1989.
- [7] S. Casner, "Frequently Asked Questions (FAQ) on the Multicast Backbone (MBone)," USC/ISI, Dec. 1994; [ftp://ftp.isi.edu/mbone/faq.txt](http://ftp.isi.edu/mbone/faq.txt)
- [8] H. Eriksson, "The Multicast Backbone," *Commun. ACM*, vol. 8, 1994, pp. 54-60.
- [9] D. Waitzman, C. Partridge, and S. Deering, "Distance Vector Multicast Routing Protocol (DVMRP)," RFC 1075, Nov. 1988.

- [10] W. Fenner, "Internet Group Management Protocol, Version 2," RFC 2236, Nov. 1997.
- [11] B. Cain, S. Deering, and A. Thyagarajan, "Internet Group Management Protocol, Version 3," Internet draft, draft-ietf-idmr-igmp-v3\*.txt, Feb. 1999.
- [12] J. Moy, "Multicast Extensions to OSPF," RFC 1584, Mar. 1994.
- [13] J. Moy, "OSPF version 2," RFC 2178, Apr. 1998.
- [14] S. Deering et al., "PIM Architecture for Wide-Area Multicast Routing," *IEEE/ACM Trans. Net.*, Apr. 1996, pp. 153-62.
- [15] S. Deering et al., "Protocol Independent Multicast Version 2 Dense Mode Specification," Internet draft, draft-ietf-pim-v2-dm\*.txt, Nov. 1998.
- [16] D. Estrin et al., "Protocol Independent Multicast Sparse-Mode (PIM-SM): Protocol Specification," RFC 2362, June 1998.
- [17] T. Ballardie, P. Francis, and J. Crowcroft, "Core Based Trees (CBT): An Architecture for Scalable Multicast Routing," *ACM SIGCOMM*, San Francisco, CA, Sept. 1995, pp. 85-95.
- [18] A. Ballardie, "Core Based Trees (CBT Version 2) Multicast Routing," RFC 2189, Sept. 1997.
- [19] B. Williamson, *Developing IP Multicast Networks, Volume 1 (Fundamentals)*, Indianapolis, IN: Cisco Press, 1999.
- [20] H. Schulzrinne et al., "RTP: A Transport Protocol for Real-Time Applications," RFC 1889, Jan. 1996.
- [21] C. Huitema, *Routing in the Internet*, Englewood Cliffs, NJ: Prentice Hall, 1995.
- [22] Y. Rekhter and T. Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 1771, Mar. 1995.
- [23] P. Traina, "Experience with the BGP-4 Protocol," RFC 1773, Mar. 1995.
- [24] T. Bates et al., "Multiprotocol Extensions for BGP-4," RFC 2283, Feb. 1998.
- [25] D. Farinacci et al., "Multicast Source Discovery Protocol (MSDP)," Internet draft, draft-farinacci-msdp\*.txt, June 1998.
- [26] S. Kumar et al., "The MASC/BGMP Architecture for Inter-Domain Multicast Routing," *ACM SIGCOMM*, Vancouver, Canada, Aug. 1998.
- [27] M. Handley, D. Thaler, and D. Estrin, "The Internet Multicast Address Allocation Architecture," Internet draft, draft-ietf-malloc-arch\*.txt, Dec. 1997.
- [28] M. Handley, "Multicast Address Allocation Protocol (AAP)," Internet draft, draft-ietf-malloc-aap\*.txt, Aug. 1998.
- [29] B. Patel, M. Shah, and S. Hanna, "Multicast Address Dynamic Client Allocation Protocol (MADCAP)," Internet draft, draft-ietf-malloc-madcap\*.txt, Feb. 1999.
- [30] D. Meyer and P. Lothberg, "Static Allocations in 233/8," Internet draft, draft-ietf-mbone-static-allocation-00.txt, May 1999.
- [31] T. Ballardie et al., "On Extending THE Standard IP Multicast Architecture," tech. rep., Univ. College London res. rep. RN/99/21, Oct. 1999.
- [32] H. Holbrook and D. Cheriton, "IP Multicast Channels: EXPRESS Support for Large-Scale Single-Source Applications," *ACM SIGCOMM*, Cambridge, MA, Aug. 1999.
- [33] R. Perlman et al., "Simple Multicast: A Design for Simple, Low-Overhead Multicast," Internet draft, draft-perlman-simple-multicast\*.txt, Feb. 1999.
- [34] H. LaMaster et al., "Multicast-Friendly Internet Exchange (MIX)," Internet draft, draft-ietf-mbone-mix\*.txt, June 1999.
- [35] J. Jamison and R. Wilder, "vBNS: The Internet Fast Lane for Research and Education," *IEEE Commun.*, Jan. 1997.
- [36] J. Jamison et al., "vBNS: Not your Father's Internet," *IEEE Spectrum*, July 1999.

## Biography

KEVIN C. ALMEROOTH [M] ([almeroth@cs.ucsb.edu](mailto:almeroth@cs.ucsb.edu)) earned his Ph.D. in computer science from the Georgia Institute of Technology in 1997. He is currently an assistant professor at the University of California in Santa Barbara where his main research interests include computer networks and protocols, multicast communication, large-scale multimedia systems, and performance evaluation. At UCSB, he is a founding member of the Media Arts and Technology Program (MATP), associate director of the Center for Information Technology and Society (CITS), and on the Executive Committee for the University of California Digital Media Innovation (DIMI) program. In the research community, he is on the Editorial Board of *IEEE Network*, has served as tutorial chair for various conferences (ICNP 1999 and ACM MM 2000), and has been on the program committees of several conferences (ICNP 1999, ICNP 2000, and NGC 1999). He is currently serving as chair of the Internet2 Working Group on Multicast, and is a member of the IETF Multicast Directorate (MADDOGS). He has been a member of both the ACM and IEEE since 1993.