

Cisco.com

Fundamentals of IP Multicast

Module 1

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

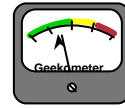
1

Module Objectives

Cisco.com

- **Recognize when to use IP Multicast**
- **Identify the fundamental concepts involved in IP Multicasting**
- **Characterize the differences in various IP Multicast routing protocols**

Agenda



Cisco.com

- **Why Multicast**
- **Multicast Applications**
- **Multicast Service Model**
- **Multicast Distribution Trees**
- **Multicast Forwarding**
- **Multicast Protocol Basics**
- **Multicast Protocol Review**

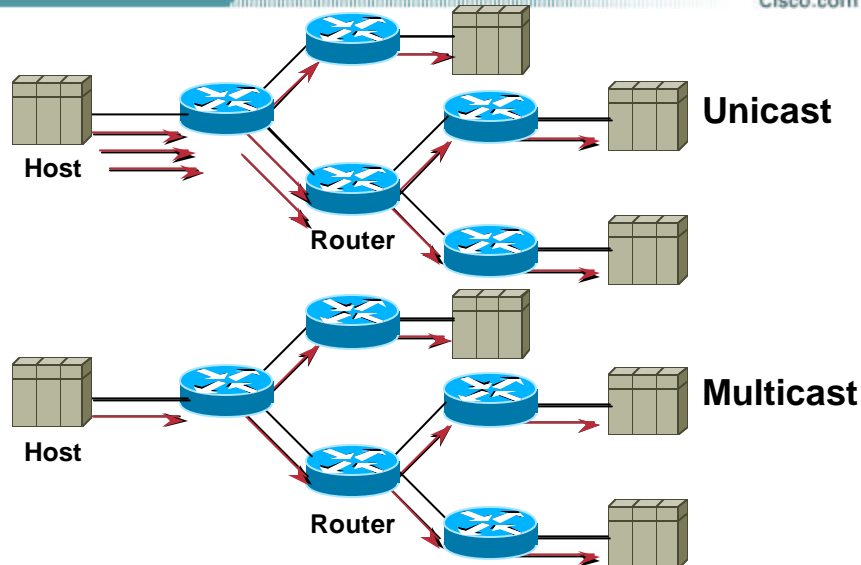
Why Multicast?

Cisco.com

- **When sending same data to multiple receivers**
- **Better bandwidth utilization**
- **Less host/router processing**
- **Receivers' addresses unknown**

Unicast vs Multicast

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

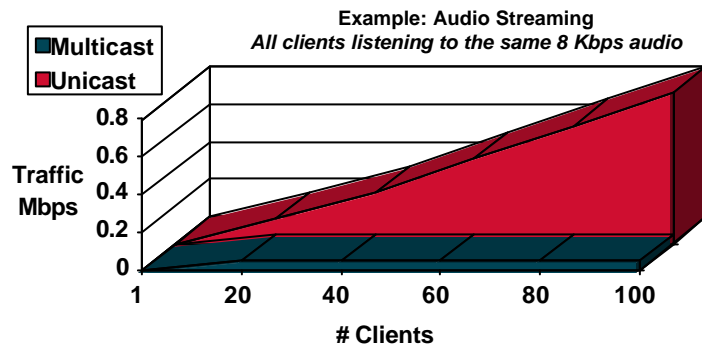
5

- **Unicast transmission sends multiple copies of data, one copy for each receiver**
 - Ex: host transmits 3 copies of data and network forwards each to 3 separate receivers
 - Ex: host can only send to one receiver at a time
- **Multicast transmission sends a single copy of data to multiple receivers**
 - Ex: host transmits 1 copy of data and network replicates at last possible hop for each receiver, each packet exists only one time on any given network
 - Ex: host can send to multiple receivers simultaneously

Multicast Advantages

Cisco.com

- **Enhanced Efficiency:** Controls network traffic and reduces server and CPU loads
- **Optimized Performance:** Eliminates traffic redundancy
- **Distributed Applications:** Makes multipoint applications possible



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

6

- **Multicast transmission affords many advantages over unicast transmission in a one-to-many or many-to-many environment**
 - Enhanced Efficiency: available network bandwidth is utilized more efficiently since multiple streams of data are replaced with a single transmission
 - Optimized Performance: less copies of data require forwarding and processing
 - Distributed Applications: multipoint applications will not be possible as demand and usage grows because unicast transmission will not scale
 - Ex: traffic level and clients increase at a 1:1 rate with unicast transmission
 - Ex: traffic level and clients do not increase at a greatly reduced rate with multicast transmission

Multicast Disadvantages

Cisco.com

Multicast is UDP Based!!!

- **Best Effort Delivery:** Drops are to be expected. Multicast applications should not expect reliable delivery of data and should be designed accordingly. Reliable Multicast is still an area for much research. Expect to see more developments in this area.
- **No Congestion Avoidance:** Lack of TCP windowing and “slow-start” mechanisms can result in network congestion. If possible, Multicast applications should attempt to detect and avoid congestion conditions.
- **Duplicates:** Some multicast protocol mechanisms (e.g. Asserts, Registers and Shortest-Path Tree Transitions) result in the occasional generation of duplicate packets. Multicast applications should be designed to expect occasional duplicate packets.
- **Out-of-Sequence Packets:** Various network events can result in packets arriving out of sequence. Multicast applications should be designed to handle packets that arrive in some other sequence than they were sent by the source.

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

7

• Multicast Disadvantages

- Most Multicast Applications are UDP based. This results in some undesirable side-effects when compared to similar unicast, TCP applications.
- Best Effort Delivery results in occasional packet drops. Many multicast applications that operate in real-time (e.g. Video, Audio) can be impacted by these losses. Also, requesting retransmission of the lost data at the application layer in these sort of real-time applications is not feasible.
 - Heavy drops on Voice applications result in jerky, missed speech patterns that can make the content unintelligible when the drop rate gets high enough.
 - Moderate to Heavy drops in Video is sometimes better tolerated by the human eye and appear as unusual “artifacts” on the picture. However, some compression algorithms can be severely impacted by even low drop rates; causing the picture to become jerky or freeze for several seconds while the decompression algorithm recovers.
- No Congestion Control can result in overall Network Degradation as the popularity of UDP based Multicast applications grow.
- Duplicate packets can occasionally be generated as multicast network topologies change.
 - Applications should expect occasional duplicate packets to arrive and should be designed accordingly.

IP Multicast Applications

Cisco.com

Live TV and Radio Broadcast
to the Desktop

Corporate Broadcasts

Multicast File Transfer
Data and File Replication

Distance Learning



Training

Whiteboard/Collaboration

Video Conferencing

Video-On-Demand

Real-Time Data Delivery—Financial

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

8

- **Many new multipoint applications are emerging as demand for them grows**

- Ex: Real-time applications include live broadcasts, financial data delivery, whiteboard collaboration, and video conferencing
- Ex: Non-real-time applications include file transfer, data and file replication, and video-on-demand
- Note also that the latest version of Novell Netware uses Ipmmc for file and print service announcements....see:
 - <http://developer.novell.com/research/appnotes/1999/march/02/index.htm>

Example Multicast Applications

Cisco.com

Mbone Multicast Applications

- **sdr—session directory**
 - Lists advertised sessions
 - Launches multicast application(s)
- **vat—audio conferencing**
 - PCM, DVI, GSM, and LPC4 compression
- **vic—video conferencing**
 - H.261 video compression
- **wb—white board**
 - Shared drawing tool
 - Can import PostScript images
 - Uses Reliable Multicast

Module1.ppt

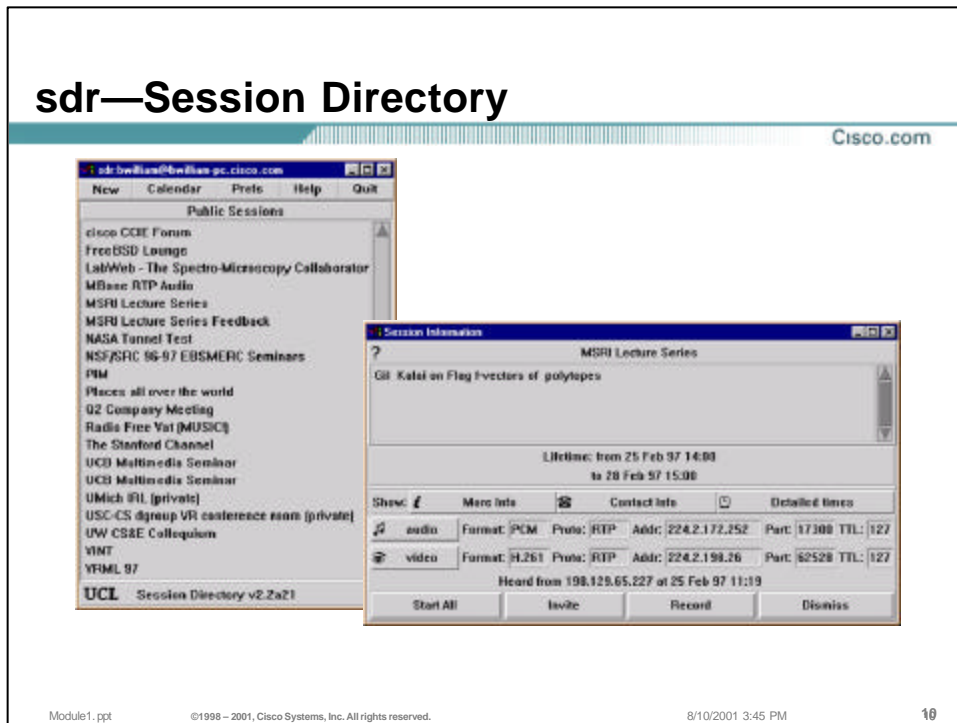
©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

9

- **Several MBONE multicast applications exist**
 - Ex: Session Directory is a tool that allows participants to view advertised multicast sessions and launch appropriate multicast applications to join an existing session
 - Ex: Audio Conferencing allows multiple participants to share audio interactively
 - Ex: Video Conferencing allows multiple participants to share video and audio interactively
 - Ex: White Boarding allows multiple participants to collaborate interactively in a text and graphical environment

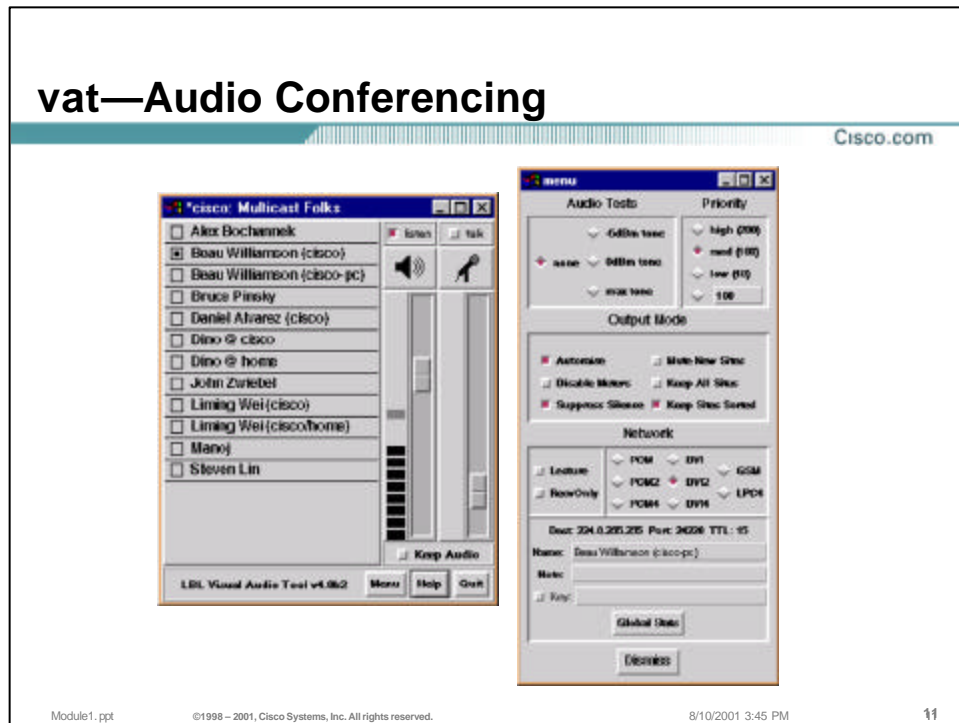
sdr—Session Directory



- **SDR - Session Directory (revised)**

- The SDR tool allows Multimedia multicast sessions to be created by other users in the network. These multimedia sessions (video, audio, etc.) are announced by the SDR application via well-known multicast groups.
- The window on the left shows an example of the SDR application in action. Each line is a multimedia session that has been created by some user in the network and is being announced (via multicast) by the creator's SDR application.
- By clicking on one of these sessions, the window on the right is brought up. This window displays various information about the multimedia session including:
 - General session information
 - Session schedule
 - Media type (in this case audio and video)
 - Media format
 - Multicast group and port numbers
- Using the window on the right, one can have SDR launch the appropriate multicast application(s) to join the session.

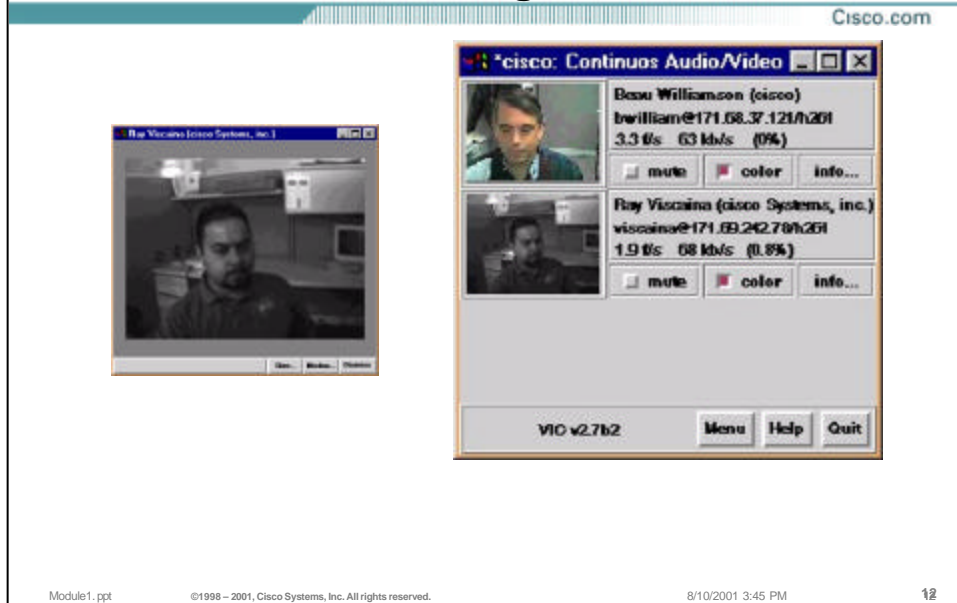
vat—Audio Conferencing



• Vat - Audio Conferencing Tool

- This is an example of the vat audio conferencing tool. The window on the left is the main window for the session. It contains a speaker gain slider widget and an output VU bar-graph meter along with a microphone gain slider widget and VU meter. When one wishes to address the conference, one usually presses the right mouse button on the workstation.
- The window on the right is a menu that can be brought up by pressing the “Menu” button on the main window. This menu allows various parameters about the session to be adjusted including encoding format.
- Notice that there are several members of this session listed in the main window even though only the second person is talking. (Indicated by the blackened square next to the name.) This points out that all members of the session are multicast sources even though they may never speak and only listed to the session. This is because vat uses the RTP/RTCP model to transport Real-Time audio data. In this model, all members of the session multicast member information and reception statistics to the entire group in an RTCP “back-channel” .
- Most all multimedia multicast applications use the RTP/RTCP model including:
 - vat (and its cousin application rat)
 - vic
 - wb - (Whiteboard)
 - IP/TV

vic—Video Conferencing



- **vic - Video Conferencing Tool**

- This is an example of the “vic” video conferencing tool. The window on the right is the main window for the video conferencing session. Notice that multiple video streams are being received, each with its own “thumbnail” image.
- The window on the left is a larger version of the thumbnail image from the main window.

wb—White Board

Cisco.com

Module1.ppt ©1998 – 2001, Cisco Systems, Inc. All rights reserved. 8/10/2001 3:45 PM 13

- **wb - Whiteboard**

- Just as its name implies, this is a form of electronic Whiteboard that can be shared by members of the multicast group.

- **“White Board” uses a form of Reliable Multicast**

- Reliable Multicasting is necessary to insure no loss of critical graphic information occurs.

- Most multimedia multicast applications simply use UDP, “Best Effort” datagram delivery mechanisms because of the time critical nature of the media. However, “wb” needs a reliable method to distribute the graphic images drawn on the electronic “Whiteboard”.

Downloading MBone Applications

Cisco.com

- **Multimedia conferencing application archive**
 - Contains sdr, vic, vat, rat, wb, nte, and other applications
 - URL:
 - <http://www-mice.cs.ucl.ac.uk/multimedia/software/>
 - Multiple platform support
 - SunOS, Solaris, HP, Linux, Windows 95/98/2000, Windows NT, etc.
 - Source code

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

14

- **Several multimedia applications for the MBONE are freely available**
 - Download the desired application for the appropriate platform
 - Source code and binaries are available

IP Multicast Service Model

Cisco.com

- **RFC 1112 (Host Ext. for Multicast Support)**
- **Each multicast group identified by a class-D IP address**
- **Members of the group could be present anywhere in the Internet**
- **Members join and leave the group and indicate this to the routers**
- **Senders and receivers are distinct:
i.e., a sender need not be a member**
- **Routers listen to all multicast addresses and use multicast routing protocols to manage groups**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

15

- **RFC 1112 is the Internet Group Management Protocol (IGMP)**
 - Allows hosts to join a group that receives multicast packets
 - Allows users to dynamically register (join/leave multicast groups) based on applications they execute
 - Uses IP datagrams to transmit data
- **Addressing**
 - Class D IP addresses (224-239) are dynamically allocated
 - Multicast IP addresses represent receiver groups, not individual receivers
- **Group Membership**
 - Receivers can be densely or sparsely distributed throughout the Internet
 - Receivers can dynamically join/leave a multicast session at any time using IGMP to manage group membership within the routers
 - Senders are not necessarily included in the multicast group they are sending to
 - Many applications have the characteristic of receivers also becoming senders eg RTCP streams from IP/TV clients and Tibco RV
- **Multicast Routing**
 - Group distribution requires packet distribution trees to efficiently forward data to multiple receivers
 - Multicast routing protocols effectively direct multicast traffic along network paths
 - Multicast Extension to Open Shortest Path First (MOSPF - 1584)
 - Core Based Tree (CBT)

IP Multicast Service Model

Cisco.com

- **IP group addresses**
 - Class D address—high-order 3 bits are set (224.0.0.0)
 - Range from 224.0.0.0 through 239.255.255.255
- **Well known addresses designated by IANA**
 - Reserved use: 224.0.0.0 through 224.0.0.255
 - 224.0.0.1—all multicast systems on subnet
 - 224.0.0.2—all routers on subnet
 - See “<ftp://ftp.isi.edu/in-notes/iana/assignments/multicast-addresses>”
- **Transient addresses, assigned and reclaimed dynamically**
 - Global scope: 224.0.1.0-238.255.255.255
 - Limited Scope: 239.0.0.0-239.255.255.255
 - Site-local scope: 239.253.0.0/16
 - Organization-local scope: 239.192.0.0/14

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

16

- **IP Addresses use the Class D address space**

- Class D addresses are denoted by the high 4 bits set to 1110.

- **Local Scope Addresses**

- Addresses 224.0.0.0 through 224.0.0.255
- Reserved by IANA for network protocol use

Examples:

224.0.0.1	All Hosts
224.0.0.2	All Multicast Routers
224.0.0.3	All DVMRP Routers
224.0.0.5	All OSPF Routers
224.0.0.6	All OSPF DR

- Multicasts in this range are never forwarded off the local network regardless of TTL
- Multicasts in this range are usually sent 'link local' with TTL = 1.

- **Global Scope Addresses**

- Addresses 224.0.1.0 through 238.255.255.255
- Allocated dynamically throughout the Internet

- **Administratively Scoped Addresses**

- Addresses 239.0.0.0 through 239.255.255.255
- Reserved for use inside of private Domains.

IP Multicast Addressing

Cisco.com

- **Dynamic Group Address Assignment**
 - **Historically accomplished using SDR application**
 - **Sessions/groups announced over well-known multicast groups**
 - **Address collisions detected and resolved at session creation time**
 - **Has problems scaling**
 - **Future dynamic techniques under consideration**
 - **Multicast Address Set-Claim (MASC)**
 - **Hierarchical, dynamic address allocation scheme**
 - **Extremely complex garbage-collection problem.**
 - **Long ways off**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

17

- **Dynamic Group Address Assignment**

- SDR
 - This was typically accomplished using the SDR application which would detect collisions in IP multicast group address assignment at the time new sessions were being created and pick an unused address.
 - While it was sufficient for use on the old Mbone when the total number of multicast sessions in the Internet was quite low, SDR has severe scaling problems that preclude it from continuing to be used as the number of sessions increase.
- Multicast Address Set-Claim (MASC)
 - MASC is new proposal for a dynamic multicast address allocation that is being developed by the “malloc” Working Group of the IETF.
 - This new proposal will provide for dynamic allocation of the global IP Multicast address space in a hierarchical manner.
 - In this proposal, domains “lease” IP multicast group address space from their parent domain. These leases are good for only a set period. It is possible that the parent domain may grant a completely different range at lease renewal time due to the need to reclaim address space for use elsewhere in the Internet.
 - As one can imagine, this is a very non-trivial mechanism and is a long ways from actual implementation.

IP Multicast Addressing

Cisco.com

- **Static Group Address Assignment (new)**
 - Temporary method to meet immediate needs
 - Group range: 233.0.0.0 - 233.255.255.255
 - Your AS number is inserted in middle two octets
 - Remaining low-order octet used for group assignment
 - Defined in IETF draft
 - “draft-ietf-mboned-glop-addressing-xx.txt”

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

18

- **Static Group Address Assignment**

- Until MASC has been fully specified and deployed, many content providers in the Internet require “something” to get going in terms of address allocation. This is being addressed with a temporary method of static multicast address allocation.
- This special allocation method is defined in:
 - “draft-ietf-mboned-glop-addressing-xx.txt”
- The basic concept behind this methodology is as follows:
 - Use the 233/8 address space for static address allocation
 - The middle two octets of the group address would contain your AS number
 - The final octet is available for group assignment.

Multicast Protocol Basics

Cisco.com

- **Multicast Distribution Trees**
- **Multicast Forwarding**
- **Types of Multicast Protocols**
 - Dense Mode Protocols
 - Sparse Mode Protocols

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

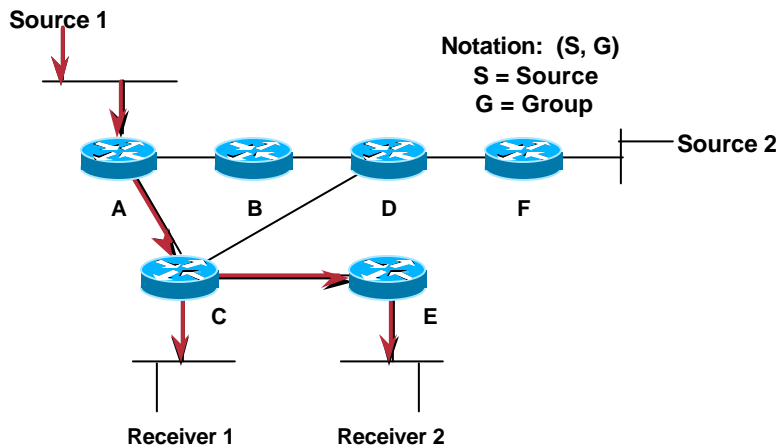
19

- **Multicast Distribution Trees**
 - Defines the path down which traffic flows from source to receiver(s).
- **Multicast Forwarding**
 - Unlike unicast forwarding which uses the “destination” address to make it’s forwarding decision, multicast forwarding uses the “source” address to make it’s forwarding decision.
- **Type of Multicast Protocols**
 - Dense Mode Protocols
 - Sparse Mode Protocols

Multicast Distribution Trees

Cisco.com

Shortest Path or Source Distribution Tree



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

20

- **Shortest Path Trees — aka Source Trees**

- A Shortest path or source distribution tree is a minimal spanning tree with the lowest cost from the source to all leaves of the tree.
- We forward packets on the Shortest Path Tree according to both the Source Address that the packets originated from and the Group address G that the packets are addressed to. For this reason we refer to the forwarding state on the SPT by the notation (S,G) (pronounced “S comma G”).

where:

- “S” is the IP address of the source.
- “G” is the multicast group address

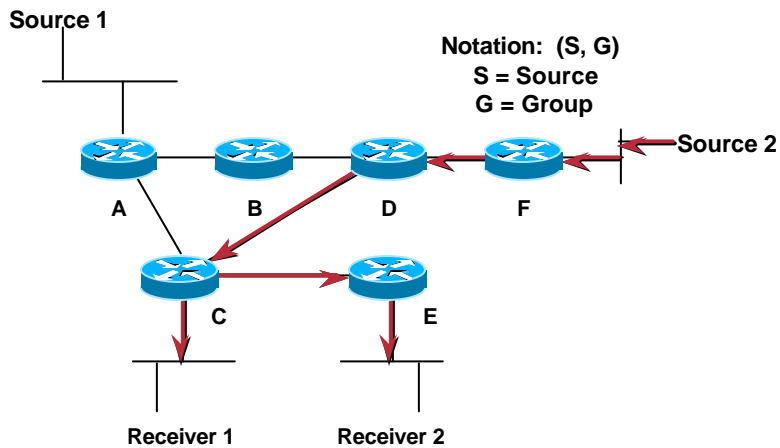
– Example 1:

- The shortest path between Source 1 and Receiver 1 is via Routers A and C, and shortest path to Receiver 2 is one additional hop via Router E.

Multicast Distribution Trees

Cisco.com

Shortest Path or Source Distribution Tree



Module1.ppt

©1998 - 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

21

- **Shortest Path Trees — aka Source Trees (cont.)**

- Every SPT is rooted at the source. This means that for every source sending to a group, there is a corresponding SPT.

- Example 2:

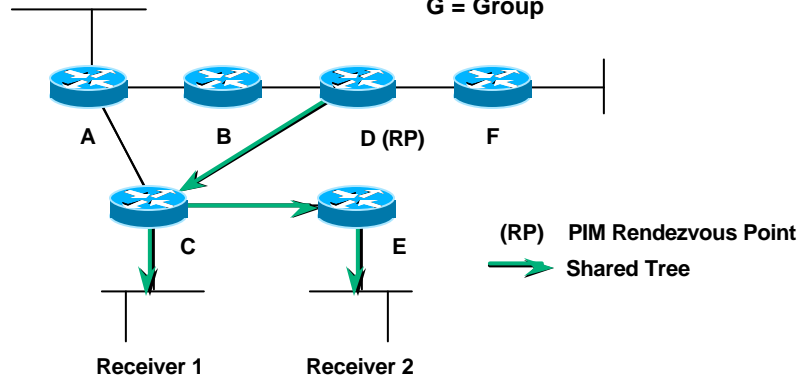
- The shortest path between Source 2 and Receiver 1 is via Routers D, F and C, and shortest path to Receiver 2 is one additional hop via Router E.

Multicast Distribution Trees

Cisco.com

Shared Distribution Tree

Notation: (*, G)
* = All Sources
G = Group



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

22

- **Shared Distribution Trees**

- Shared distribution tree whose root is a shared point in the network down which multicast data flows to reach the receivers in the network. In PIM-SM, this shared point is called the Rendezvous Point (RP).
- Multicast traffic is forwarded down the Shared Tree according to just the Group address G that the packets are addressed to, regardless of source address. For this reason we refer to the forwarding state on the shared tree by the notation (*,G) (pronounced “star comma G”)

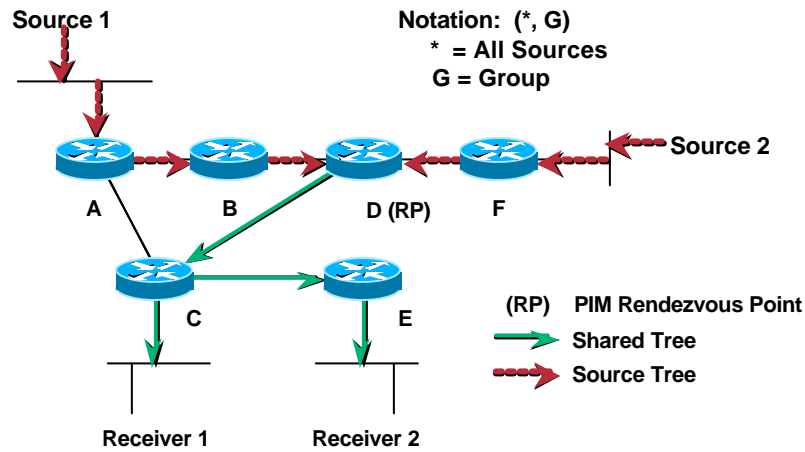
where:

- “*” means any source
- “G” is the group address

Multicast Distribution Trees

Cisco.com

Shared Distribution Tree



Module1.ppt

©1998 - 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

23

- **Shared Distribution Trees (cont.)**

- Before traffic can be sent down the Shared Tree it must somehow be sent to the Root of the Tree.
- In classic PIM-SM, this is accomplished by the RP joining the Shortest Path Tree back to each source so that the traffic can flow to the RP and from there down the shared tree. In order to trigger the RP to take this action, it must somehow be notified when a source goes active in the network.
 - In PIM-SM, this is accomplished by first-hop routers (i.e. the router directly connected to an active source) sending a special Register message to the RP to inform it of the active source.
- In the example above, the RP has been informed of Sources 1 and 2 being active and has subsequently joined the SPT to these sources.

Multicast Distribution Trees

Cisco.com

Characteristics of Distribution Trees

- **Source or Shortest Path trees**
Uses more memory $O(S \times G)$ but you get optimal paths from source to all receivers; minimizes delay
- **Shared trees**
Uses less memory $O(G)$ but you may get sub-optimal paths from source to all receivers; may introduce extra delay

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

24

- **Source or Shortest Path Tree Characteristics**
 - Provides optimal path (shortest distance and minimized delay) from source to all receivers, but requires more memory to maintain
- **Shared Tree Characteristics**
 - Provides sub-optimal path (may not be shortest distance and may introduce extra delay) from source to all receivers, but requires less memory to maintain

Multicast Distribution Trees

Cisco.com

How are Distribution Trees Built?

- **PIM**
 - Uses existing Unicast Routing Table plus Join/Prune/Graft mechanism to build tree.
- **DVMRP**
 - Uses DVMRP Routing Table plus special Poison-Reverse mechanism to build tree.
- **MOSPF**
 - Uses extension to OSPF's link state mechanism to build tree.
- **CBT**
 - Uses existing Unicast Routing Table plus Join/Prune/Graft mechanism to build tree.

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

25

- **Distribution trees are built in a variety of ways, depending upon the multicast routing protocol employed**

- PIM utilizes the underlying unicast routing table (any unicast routing protocol) plus:

Join: routers add their interfaces and/or send PIM-JOIN messages upstream to establish themselves as branches of the tree when they have interested receivers attached

Prune: routers prune their interfaces and/or send PIM-PRUNE messages upstream to remove themselves from the distribution tree when they no longer have interested receivers attached

Graft: routers send PIM-GRAFT messages upstream when they have a pruned interface and have already sent PIM-PRUNES upstream, but receive an IGMP host report for the group that was pruned; routers must reestablish themselves as branches of the distribution tree because of new interested receivers attached

- DVMRP utilizes a special RIP-like multicast routing table plus:

Poison-Reverse: a special metric of Infinity (32) plus the originally received metric, used to signal that the router should be placed on the distribution tree for the source network.

Prunes & Grafts: routers send Prunes and Grafts up the distribution similar to PIM-DM.

- MOSPF utilizes the underlying OSPF unicast routing protocol's link state advertisements to build (S,G) trees

Each router maintains an up-to-date image of the topology of the entire network

- CBT utilizes the underlying unicast routing table and the Join/Prune/Graft mechanisms (much like PIM-SM)

Multicast Forwarding

Cisco.com

- **Multicast Routing is backwards from Unicast Routing**
 - Unicast Routing is concerned about where the packet is going.
 - Multicast Routing is concerned about where the packet came from.
- **Multicast Routing uses “Reverse Path Forwarding”**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

26

- **Multicast Forwarding**
 - Routers must know packet origin, rather than destination (opposite of unicast)
 - ... origination IP address denotes known source
 - ... destination IP address denotes unknown group of receivers
 - Multicast routing utilizes Reverse Path Forwarding (RPF)
 - ... Broadcast: floods packets out all interfaces except incoming from source; initially assuming every host on network is part of multicast group
 - ... Prune: eliminates tree branches without multicast group members; cuts off transmission to LANs without interested receivers
 - ... Selective Forwarding: requires its own integrated unicast routing protocol

Reverse Path Forwarding (RPF)

Cisco.com

- **What is RPF?**

A router forwards a multicast datagram only if received on the up stream interface to the source (i.e. it follows the distribution tree).

- **The RPF Check**

- The routing table used for multicasting is checked against the “source” address in the multicast datagram.
- If the datagram arrived on the interface specified in the routing table for the source address; then the RPF check succeeds.
- Otherwise, the RPF Check fails.

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

27

- **Reverse Path Forwarding**

- Routers forward multicast datagrams received from incoming interface on distribution tree leading to source
- Routers check the source IP address against their multicast routing tables (RPF check); ensure that the multicast datagram was received on the specified incoming interface
- Note that changes in the unicast topology will not necessarily immediately reflect a change in RPF...this depends on how frequently the RPF check is performed on an lpmc stream - every 5 seconds is current Cisco default.

Reverse Path Forwarding (RPF)

Cisco.com

- If the RPF check succeeds, the datagram is forwarded
- If the RPF check fails, the datagram is typically silently discarded
- When a datagram is forwarded, it is sent out each interface in the outgoing interface list
- Packet is **never** forwarded back out the RPF interface!

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

28

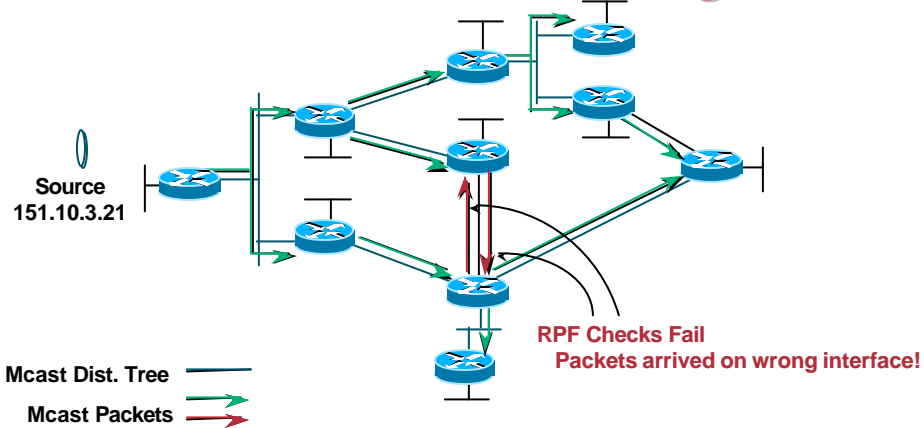
- **Multicast Forwarding**

- Successful RPF Checks allow the datagram to be forwarded
 - ... Datagram is forwarded out all outgoing interfaces, but not out the RPF interface the datagram was received on
- Unsuccessful RPF Checks cause the datagram to be silently dropped

Reverse Path Forwarding (RPF)

Cisco.com

Example: RPF Checking



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

29

• Multicast Forwarding: RPF Checking

- Source floods network with multicast data
- Each router has a designated incoming interface (RPF interface) on which multicast data can be received from a given source
- Each router receives multicast data on one or more interfaces, but performs RPF check to prevent duplicate forwarding

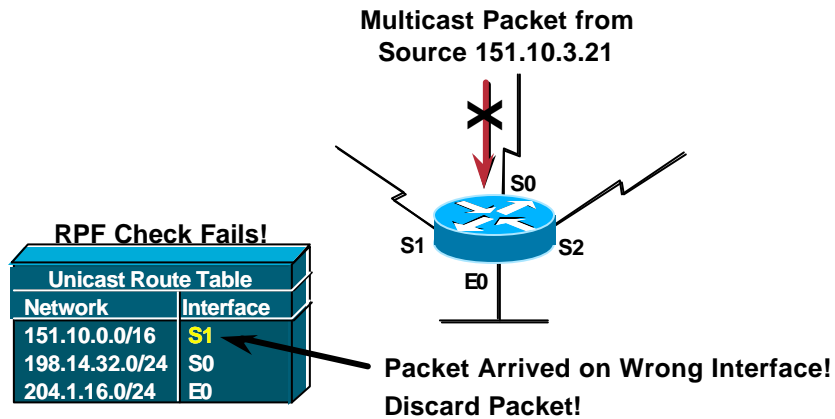
Example: Router receives multicast data on two interfaces

- 1) performs RPF Check on multicast data received on interface E0; RPF Check succeeds because data was received on specified incoming interface from source 151.10.3.21; data forwarded through all outgoing interfaces on the multicast distribution tree
- 2) performs RPF Check on multicast data received on interface E1; RPF Check fails because data was not received on specified incoming interface from source 151.10.3.21; data silently dropped

Reverse Path Forwarding (RPF)

Cisco.com

A closer look: RPF Check Fails



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

30

- **Multicast Forwarding: RPF Check Fails**

- Ex: Router can only accept multicast data from Source 151.10.3.21 on interface S1
... multicast data is silently dropped because it arrived on an interface not specified in the RPF check (S0)

Reverse Path Forwarding (RPF)

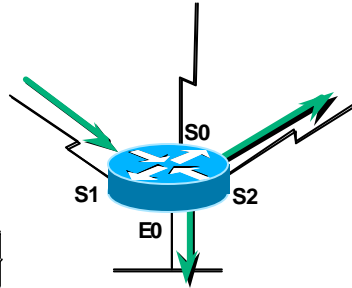
Cisco.com

A closer look: RPF Check Succeeds

Multicast Packet from Source 151.10.3.21

RPF Check Succeeds!

Unicast Route Table	
Network	Interface
151.10.0.0/16	S1
198.14.32.0/24	S0
204.1.16.0/24	E0



Packet Arrived on Correct Interface!
Forward out all outgoing interfaces.
(i. e. down the distribution tree)

Module1.ppt

©1998 - 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

31

- **Multicast Forwarding: RPF Check Succeeds**

- Ex: Router can only accept multicast data from Source 151.10.3.21 on interface S1
... multicast data is forwarded out all outgoing on the distribution tree because it arrive on the incoming interface specified in the RPF check (S1)

TTL Thresholds

Cisco.com

- **What is a TTL Threshold?**

A “TTL Threshold” may be set on a multicast router interface to limit the forwarding of multicast traffic to outgoing packets with TTLs greater than the Threshold.

- **The TTL Threshold Check**

- 1) All incoming IP packets first have their TTL decremented by one. If \leq Zero, they are dropped.
- 2) If a multicast packet is to be forwarded out an interface with a non-zero TTL Threshold; then its TTL is checked against the TTL Threshold. If the packet's TTL is $<$ the specified threshold, it is not forwarded out the interface.

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

32

- **TTL-Thresholds**

- Non-Zero, Multicast, TTL-Thresholds may be set on any multicast capable interface.
- IP multicast packets whose TTLs (after being decremented by one by normal router packet processing) are less than the TTL-Threshold on an outgoing interface, will be not be forwarded out that interface.
- Zero Multicast TTL implies NO threshold has been set.

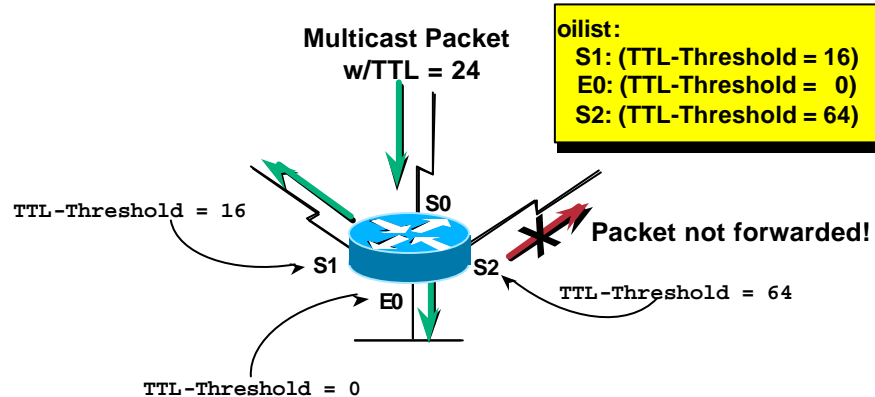
- **TTL-Threshold Application**

- Frequently used to set up multicast boundaries to prevent unwanted multicast traffic from entering/exiting the network.

TTL Thresholds

Cisco.com

A closer look: TTL-Thresholds



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

33

• TTL-Threshold Example

- In the above example, the interfaces have been configured with the following TTL-Thresholds:

S1: TTL-Threshold = 16
E0: TTL-Threshold = 0 (none)
S2: TTL-Threshold = 64

- An incoming Multicast packet is received on interface S0 with a TTL of 24.
- The TTL is decremented to 23 by the normal router IP packet processing.
- The outgoing interface list for this Group contains interfaces S1, E0 & S2.
- The TTL-Threshold check is performed on each outgoing interface as follows:

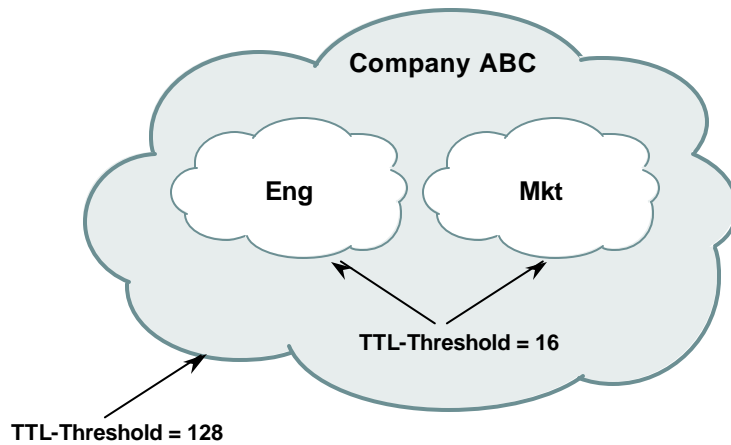
S1: TTL (23) > TTL-Threshold (16). FORWARD

E0: TTL (23) > TTL-Threshold (0). FORWARD

S2: TTL (23) < TTL-Threshold (64). DROP

TTL Threshold Boundaries

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

34

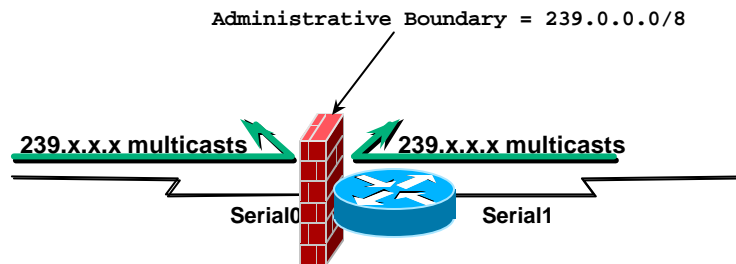
- **TTL-Threshold Boundaries**

TTL-Thresholds may be used as boundaries around portions of a network to prevent the entry/exit of unwanted multicast traffic. This requires multicast applications to transmit their multicast traffic with an initial TTL value set so as to not cross the TTL -Threshold boundaries.

In the example above, the Engineering or Marketing departments can prevent department related multicast traffic from leaving their network by using a TTL of 15 for their multicast sessions. Similarly, Company ABC can prevent private multicast traffic from leaving their network by using a TTL of 127 for their multicast sessions.

Administrative Boundaries

Cisco.com



- Configured using the `'ip multicast boundary <acl>'` interface command

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

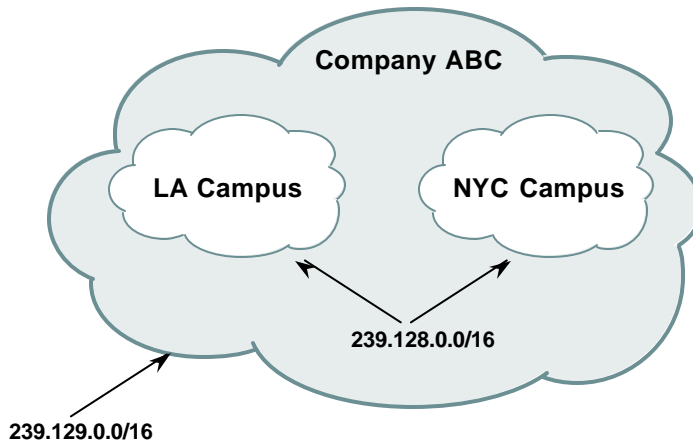
35

- **Administrative Boundaries**

- Administratively-scoped multicast address ranges may also be used as boundaries around portions of a network to prevent the entry/exit of unwanted multicast traffic. This requires multicast applications to transmit their multicast traffic with a group address that falls within the Administrative address range so that it will not cross the Administrative boundaries.
- In the example above, the entire Administratively-Scoped address range, (239.0.0.0/8) is being blocked from entering or leaving the router via interface Serial0. This is often done at the border of a network where it connects to the Internet so that potentially sensitive company Administratively-Scoped multicast traffic can leave the network. (Nor can it enter the network from the outside.)
- Administrative multicast boundaries can be configured in Cisco IOS by the use of the `'ip multicast boundary <acl>'` interface command.

Administrative Boundaries

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

36

- **Administrative Boundaries**

- Administratively-scoped multicast address ranges generally used in more than one location.
- In the example above, the Administratively-Scoped address range, (239.128.0.0/16) is being used by both the LA campus and the NYC campus. Multicast traffic originated in these address ranges will remain within each respective campus and not onto the WAN that exists between the two campuses.
This is often sort of configuration is often used so that each campus can source high-rate multicasts on the local campus and not worry about it being accidentally “leaked” into the WAN and causing congestion on the slower WAN links.
- In addition to the 239.128.0.0/16 range, the entire company network has a Administrative boundary for the 239.129.0.0/16 multicast range. This is so that multicasts in these ranges do not leak into the Internet.
 - Note: The Admin.-Scoped address range (239..0.0/8) is similar to the 10.0.0.0 unicast address range in that it is reserved and is not assigned for use in the Internet.

Types of Multicast Protocols

Cisco.com

- **Dense-mode**
 - Uses “Push” Model
 - Traffic Flooded throughout network
 - Pruned back where it is unwanted
 - Flood & Prune behavior (typically every 3 minutes)
- **Sparse-mode**
 - Uses “Pull” Model
 - Traffic sent only to where it is requested
 - Explicit Join behavior

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

37

- **Dense-mode multicast protocols**
 - Initially flood/broadcast multicast data to entire network, then prune back paths that don't have interested receivers
- **Sparse-mode multicast protocols**
 - Assumes no receivers are interested unless they explicitly ask for it

Multicast Protocol Review

Cisco.com

- **Currently, there are 4 multicast routing protocols:**

- † **DVMRPv3 (Internet-draft)**

DVMRPv1 (RFC1075) is obsolete and was never used.

- † **MOSPF (RFC 1584) “Proposed Standard”**

- † **PIM-DM (Internet-draft)**

- † **CBT (Internet-draft)**

- † **PIM-SM (RFC 2362) “Proposed Standard”**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

38

- **IETF status of Multicast Protocols**

- DVMRPv1 is obsolete and was never used. DVMRPv2 is an old “Internet-Draft” and is the current implementation used through-out the Mbone. DVMRPv3 is the current “Internet-Draft” although it has not been completely implemented by most vendors.
- MOSPF is currently at “Proposed Standard” status. However, most members of the IETF IDMR working group doubt that MOSPF will scale to any degree and are therefore uncomfortable with declaring MOSPF as a standard for IP Multicasting. (Even the author of MOSPF, J. Moy, has been quoted in an RFC that, “more work needs to be done to determine the scalability of MOSPF.”)
- PIM-DM is in Internet Draft form and work continues to move into an RFC.
- CBT is also in Internet Draft form and while it has been through three different and incompatible revisions, it is not enjoying significant usage nor is it a primary focus of the IETF IDMR working group.
- PIM-SM moved to “Proposed Standard” in early 2000. Much of the effort in the IETF towards a working multicast protocol is focused on PIM-SM.

Dense-Mode Protocols

Cisco.com

- **DVMRP - Distance Vector Multicast Routing Protocol**
- **MOSPF - Multicast OSPF**
- **PIM DM - Protocol Independent Multicasting (Dense Mode)**

DVMRP Overview

Cisco.com

- **Dense Mode Protocol**
 - **Distance vector-based**
 - **Similar to RIP**
 - **Infinity = 32 hops**
 - **Subnet masks in route advertisements**
 - **DVMRP Routes used:**
 - **For RPF Check**
 - **To build Truncated Broadcast Trees (TBTs)**
 - **Uses special “Poison-Reverse” mechanism**
 - **Uses Flood and Prune operation**
 - **Traffic initially flooded down TBT's**
 - **TBT branches are pruned where traffic is unwanted.**
 - **Prunes periodically time-out causing reflooding.**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

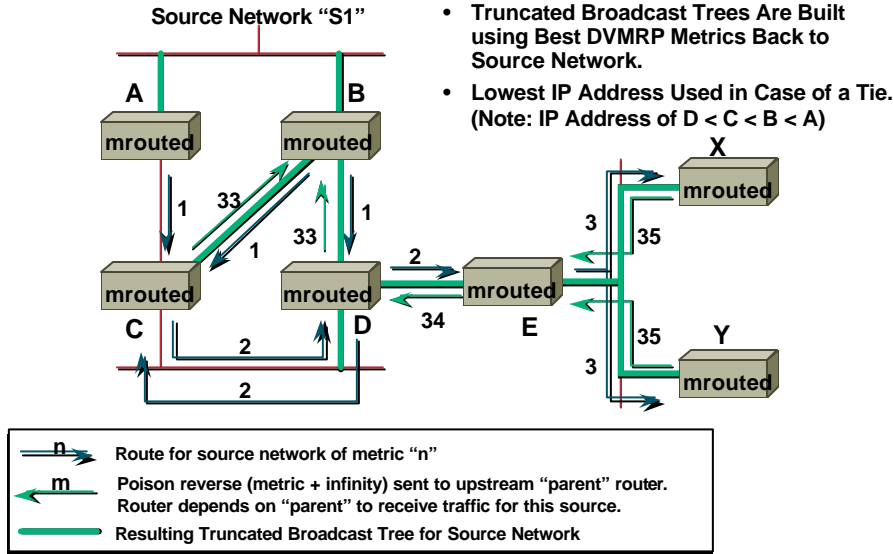
40

- **Distance Vector Multicast Routing Protocol**

- Builds a distribution tree per source network based on best metric (hop-count) back towards the source network.
 - Infinity = 32 hops
 - A “Poison Reverse” metric is used by DVMRP routers to signal their upstream neighbor that they are downstream and expect to receive traffic from a source network via the upstream router.
 - “Poison Reverse” is denoted by adding Infinity (32) to the received metric and then sending it back to the router from which it was originally received.
 - When a “Poison Reverse” is received for a source network, the interface over which it was received is placed on the outgoing interface list for the source network.
- Multicast traffic is “flooded” out all interfaces on the outgoing interface list (i.e. down the distribution tree for the source network).
- Downstream neighbors send Prunes up the distribution tree for multicast traffic for which they have no group members.

DVMRP — Source Trees

Cisco.com



• DVMRP Source Trees

- DVMRP builds its Source Trees utilising the concept of “Truncated Broadcast Trees”. The basic definition of a Truncated Broadcast Tree (TBT) is as follows:
 - A Truncated Broadcast Tree (TBT) for source subnet “S1”, represent a shortest path spanning tree rooted at subnet “S1” to all other routers in the network.
- In DVMRP, the abstract notion of the TBT’s for all sub-networks are built by the exchange of periodic DVMRP routing updates between all DVMRP routers in the network. Just like its unicast cousin, RIPv2, DVMRP updates contain network prefixes/masks along with route metrics (in hop-counts) that describe the cost of reaching a particular subnets in the network.
- Unlike RIPv2, a downstream DVMRP router makes use of a special Poison-Reverse advertisement to signal an upstream router that this link is on the TBT for source subnet S1. This Poison-Reverse (PR) is created by adding 32 to the advertised metric and sending back to the upstream router.

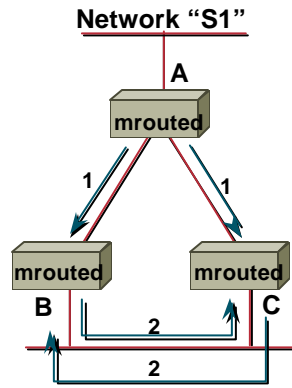
• Example DVMRP TBT for network “S1”:

- In the above example, DVMRP updates are being exchanged for source network “S1”. Routers A and B both advertise a metric of 1 (hop) to reach network S1 to routers C and D. In the case of router D, the advertisement from B is the best (only) route to source network S1 which causes router D to send back a PR advertisement (metric = 33) to B. This tells router B that router D is on the TBT for source network S1. In the case of router C, it received an advertisement from both A and B with the same metric. It breaks the tie using the lowest IP address and therefore sends a PR advertisement to router B. B now knows it has two branches of the TBT, one to router C and one to router D. These DVMRP updates flow throughout the entire network causing each router to send PR advertisements to its upstream DVMRP neighbor on the TBT for source network S1.

DVMRP — Source Trees

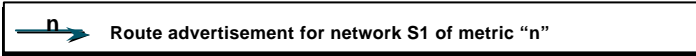
Cisco.com

Forwarding onto Multi-access Networks



- Both B & C have routes to network S1.
- To avoid duplicates, only one router can be "Designated Forwarder" for network S1.
- Router with best metric is elected as the "Designated Forwarder".
- Lowest IP address used as tie-breaker.
- Router C wins in this example.

(Note: IP Address of C < B)



Module1.ppt

©1998 - 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

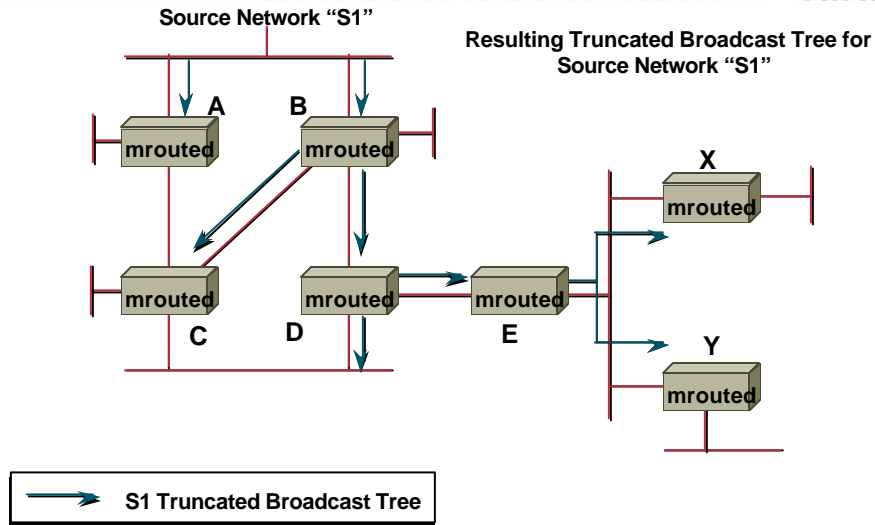
42

• Forwarding onto Multi-access Networks

- When two or more routers share a common Multi-access network, only one can be the "Designated Forwarder" which is responsible for forwarding a source network's traffic onto the Multi-access network; otherwise duplicate packets will be generated.
- The "Designated Forwarder" is selected based on the best route metric back to the source network (with the Lowest IP Address used as a tie-breaker).
- In the example above, both Router B and C share a common Multi-access network and each have routes to network S1. Since both have the same metric to network S1, the lowest IP address is used to break the tie (in this case, Router C wins).

DVMRP — Source Trees

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

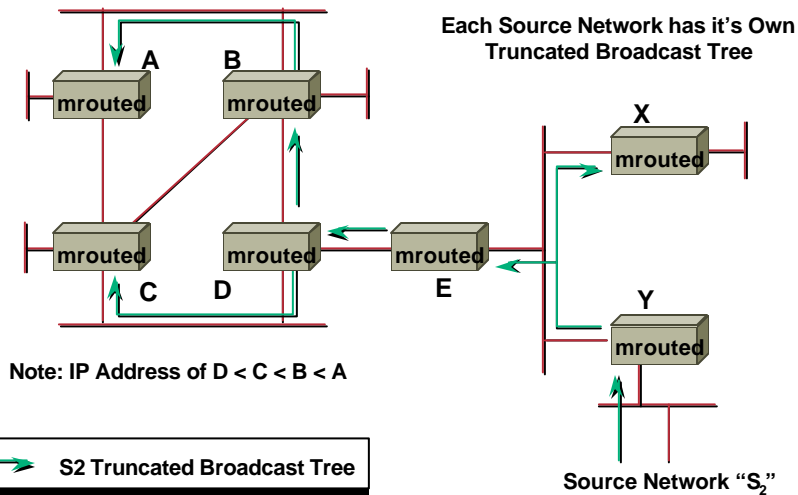
43

- **Example DVMRP TBT for network “S1” (cont.)**

- Once the DVMRP network has converged and all PR advertisements have been sent up the TBT toward source network “S1”, the S1 TBT has been built.
- The drawing above shows the S1 TBT that resulted in the DVMRP route update exchanges from the previous page. Notice that this is a minimum spanning tree that is rooted at source network “S1” and “spans” all routers in the network.
- If a multicast source were to now go active in network “S1”, the DVMRP routers in the network will initially “flood” this sources traffic down the S1 TBT.

DVMRP — Source Trees

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

44

- **Every source network has its own TBT**

- In the drawing above, the TBT for network S2 is shown. This TBT would also be created by the exchange of DVMRP route updates and by PR advertisements sent by all routers in the network toward network S2.
- It is important to remember that these TBT's simply exist in the form of PR advertisements in the DVMRP routing tables of the routers in the network and as such, there is one TBT for every source network in the DVMRP network.

- **Advantages of TBT's**

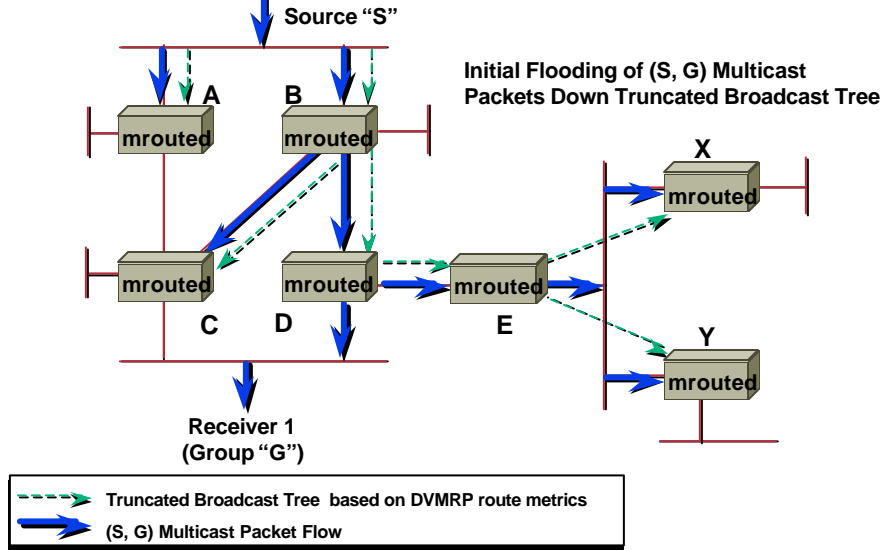
- The advantage of TBT's is that the initial flooding of multicast traffic throughout the DVMRP network is limited to flowing down the branches of the TBT. This insures that there are no duplicate packets sent as a result of parallel paths in the network.

- **Disadvantages of TBT's**

- The disadvantage of using TBT's is that it requires separate DVMRP routing information to be exchanged throughout the entire network. (Unlike other multicast protocols such as PIM that make use of the existing unicast routing table and do not have to exchange additional multicast routing data.
- Additionally, because DVMRP is based on a RIP model, it has all of the problems associated with a Distance-Vector protocol including, count-to-infinity, holddown, periodic updates.
 - One has to ask oneself, "Would I recommend someone build a unicast network based on RIP today?" The answer is of course not, protocols like OSPF, IS-IS, and EIGRP have long since superseded RIP in robustness and scalability. The same is true of DVMRP.

DVMRP — Flood & Prune

Cisco.com

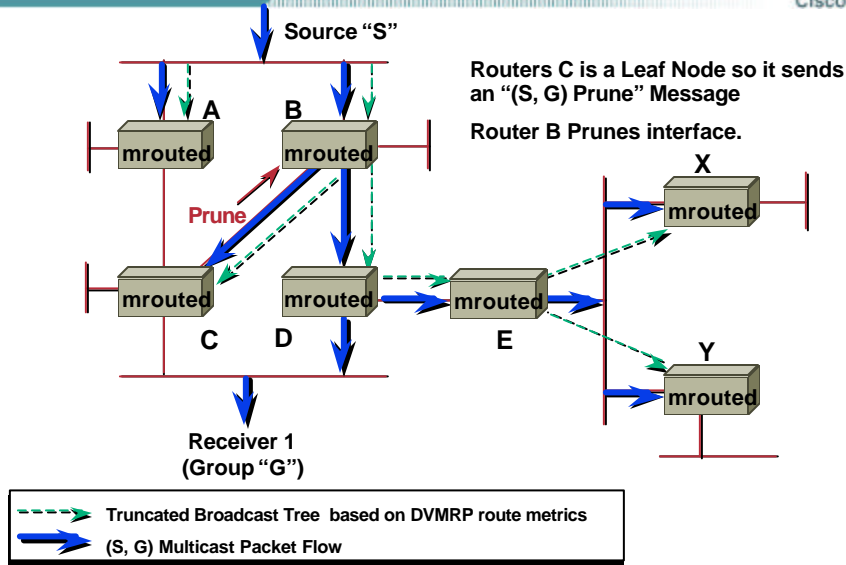


- **DVMRP Example**

- In this example we see source "S" has begun to transmit multicast traffic to group "G".
- Initially, the traffic (shown by the solid arrows) is flooded to all routers in the network down the Truncated Broadcast Tree (indicated by the dashed arrows) and is reaching Receiver 1.

DVMRP — Flood & Prune

Cisco.com

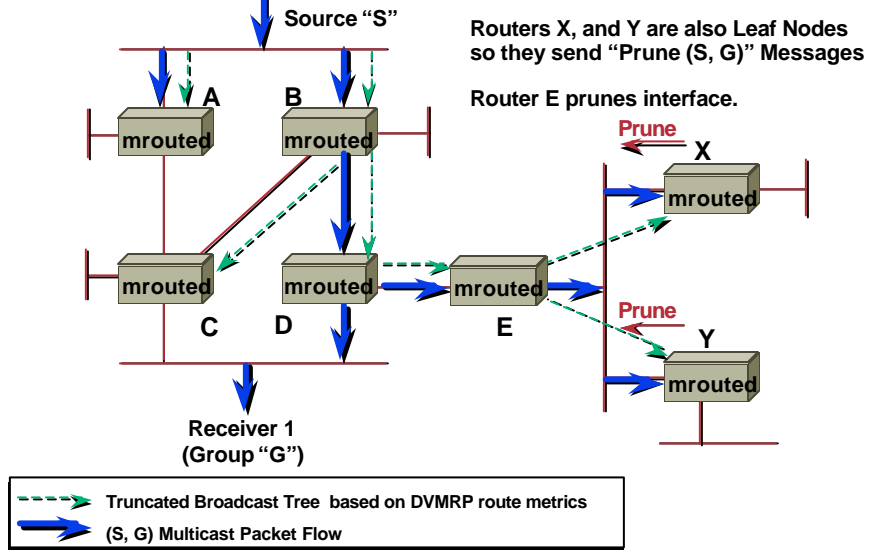


- **DVMRP Example (cont.)**

- At this point, we see that router C is a leaf node on the TBT and has no need for the traffic. Therefore, it sends a DVMRP (S, G) Prune message up the TBT to router B to shutoff the unwanted flow of traffic.
- Router B receives this (S, G) Prune message and shuts off the flow of (S, G) traffic to router C.

DVMRP — Flood & Prune

Cisco.com

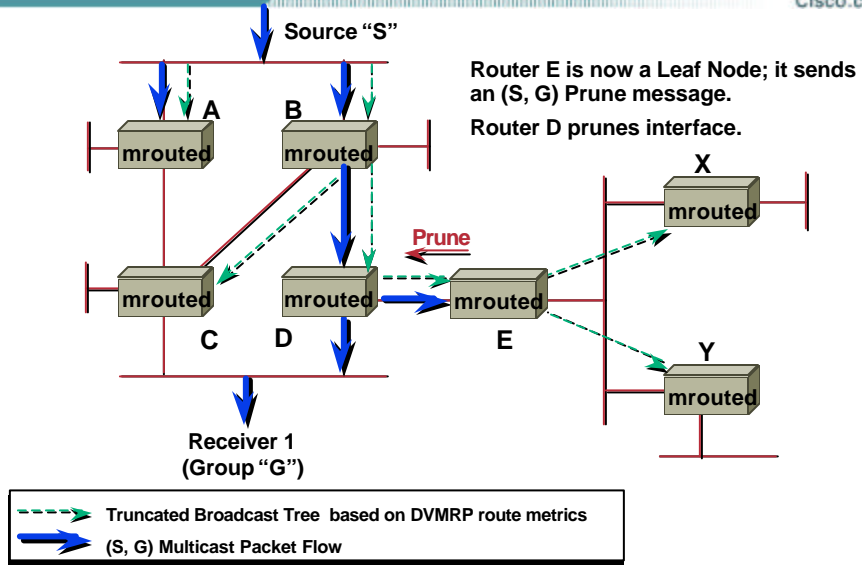


• DVMRP Example (cont.)

- Both routers X and Y are also leaf nodes that have no need for the (S, G) traffic (i.e. they have no directly connected receivers) and therefore send (S, G) Prunes up the TBT to router E.
- Once router E has received (S, G) Prunes messages *from all DVMRP neighbours on the subnet* it prunes the Ethernet interface connecting to router X and Y.

DVMRP — Flood & Prune

Cisco.com

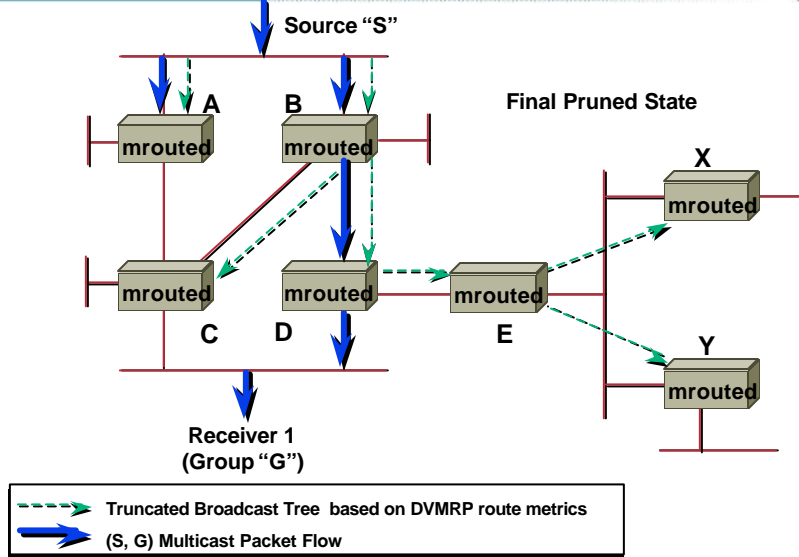


• DVMRP Example (cont.)

- At this point, all of router E's downstream interfaces on the TBT have been pruned and it no longer has any need for the (S, G) traffic. As a result, it too sends an (S,G) Prune up the TBT to router D.
- When router D receives this (S, G) Prune, it prunes the interface and shuts off the flow of (S, G) traffic to router E.

DVMRP — Flood & Prune

Cisco.com



Module1.ppt

©1998 - 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

49

• DVMRP Example (cont.)

- In the drawing above, we see the final pruned state of the TBT which leaves traffic flowing to the receiver.
- However, because DVMRP is a "flood and prune" protocol, these pruned branches of the TBT will time out (typically after 2 minutes) and (S, G) traffic will once again flood down all branches of the TBT. This will again trigger the sending (S, G) Prune messages up the TBT to prune of unwanted traffic.

DVMRP — Evaluation

Cisco.com

- **Widely used on the MBONE (being phased out)**
- **Significant scaling problems**
 - Slow Convergence—RIP-like behavior
 - Significant amount of multicast routing state information stored in routers—(S,G) everywhere
 - No support for shared trees
 - Maximum number of hops < 32
- **Not appropriate** for large scale production networks
 - Due to flood and prune behavior
 - Due to its poor scalability

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

50

- **Appropriate for large number of densely distributed receivers located in close proximity to source**
- **Widely used, oldest multicast routing protocol**
- **Significant scaling problems**
 - Protocol limits maximum number of hops to 32 and requires a great deal of multicast routing state information to be retained
- **Not appropriate for...**
 - Few interested receivers (assumes everyone wants data initially)
 - Groups sparsely represented over WAN (floods frequently)

MOSPF (RFC 1584)

Cisco.com

- **Extension to OSPF unicast routing protocol**
 - OSPF: Routers use link state advertisements to understand all available links in the network (route messages along least-cost paths)
 - MOSPF: Includes multicast information in OSPF link state advertisements to construct multicast distribution trees (each router maintains an up-to-date image of the topology of the entire network)
- **Group membership LSAs are flooded throughout the OSPF routing domain so MOSPF routers can compute outgoing interface lists**
- **Uses Dijkstra algorithm to compute shortest-path tree**
 - Separate calculation is required for each (S_{Net} , G) pair

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

51

- **Multicast Extension to OSPF (RFC 1584)**

- Extension to OSPF unicast routing protocol; requires OSPF as underlying unicast routing protocol.

- **Group Membership LSAs**

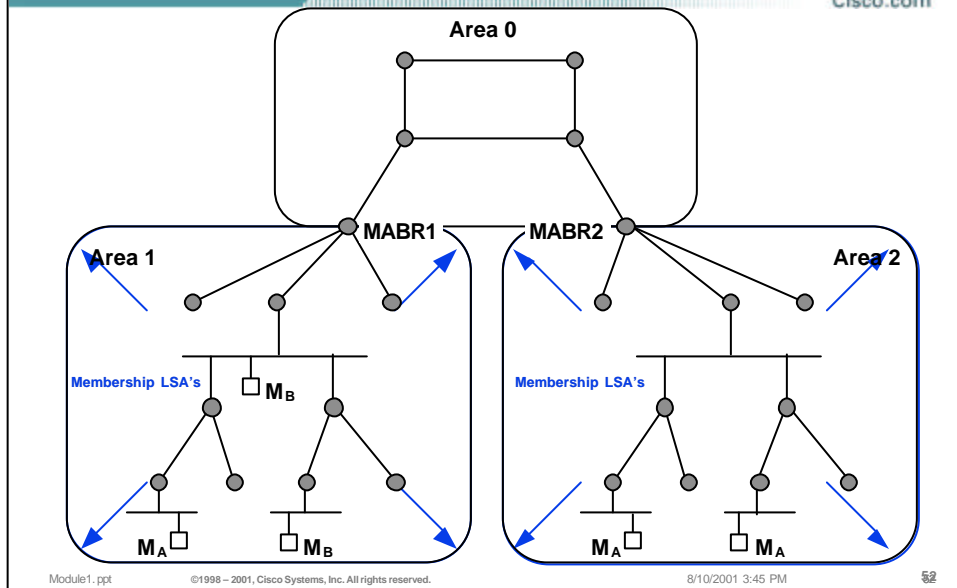
- MOSPF uses a new type of OSPF LSA called “Group-Membership LSA” to advertise the existence of Group members on networks.
- Group-Membership LSA’s are periodically flooded throughout an area in the same fashion as other OSPF LSAs.

- **Dijkstra Algorithm**

- Uses Dijkstra algorithm to compute shortest-path tree for every source-network/group pair!!!

MOSPF Membership LSA's

Cisco.com

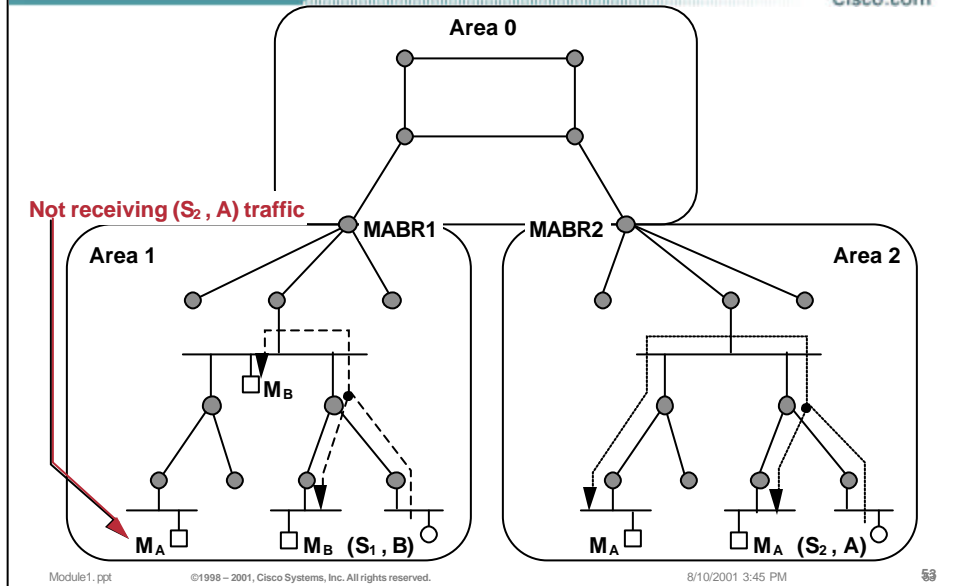


• Membership LSA Flooding Example

- In this example, Area 1 has members of both Group "A" and "B" while Area 2 has members of Group "A" only.
- Routers with directly connect members originate Membership LSA's announcing the existence of these members on their networks. These LSA's are flooded throughout the area.
- Notice that these Group Membership LSA's do not travel between Area 1 and Area 2. (This will be addressed later.)

MOSPF Intra-Area Traffic

Cisco.com



- **Intra-Area Multicast**

- Once all routers within the area have learned where all members are in the network topology, it is possible to construct Source-network trees for multicast traffic forwarding.

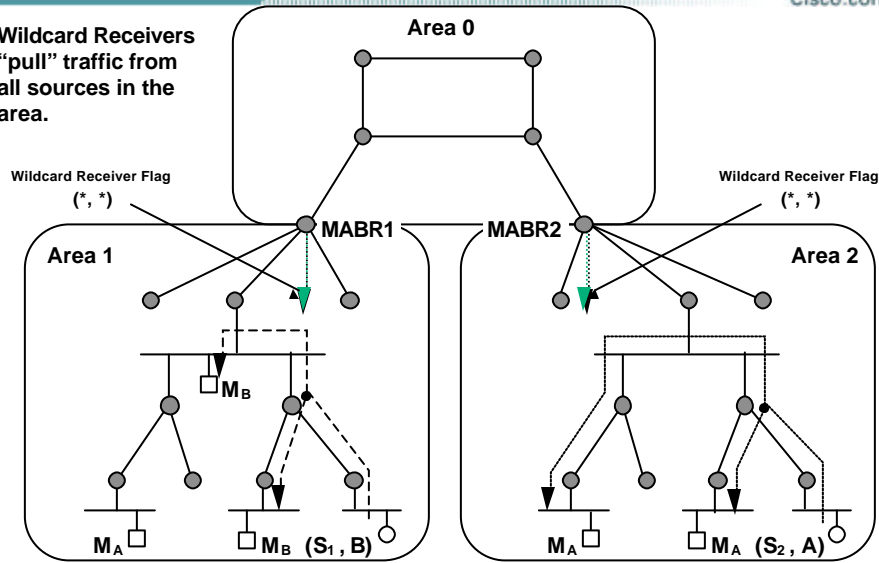
- **Example**

- In the above example, Source “S1” in Area 1 begins sending multicast traffic to Group “B”. As this data reaches the the routers in the area, each runs a Dijkstra calculation and computes a Shortest Path Tree rooted at the network for “S1” and that spans all the members of Group “B”. The results of these calculations are used to forward the (S1, B) traffic as seen in Area 1 above.
- In Area 2, Source “S2” begins sending multicast traffic to Group “A”. Again, the routers in the area use the Group-Membership information in their MOSPF database to run a Dijkstra calculation for the source network where “S2” resides and create a Shortest Path Tree for this traffic flow. The results are then used to forward (S2, A) traffic as shown.
- Notice that the routers in Area 2 are not aware of the member of Group “A” in Area 1 because Membership LSA’s are not flooded between these two areas. This Inter-Area traffic flow is handled by another mechanism that is described in the next few pages.

MOSPF Inter-Area Traffic

Cisco.com

Wildcard Receivers
“pull” traffic from
all sources in the
area.



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

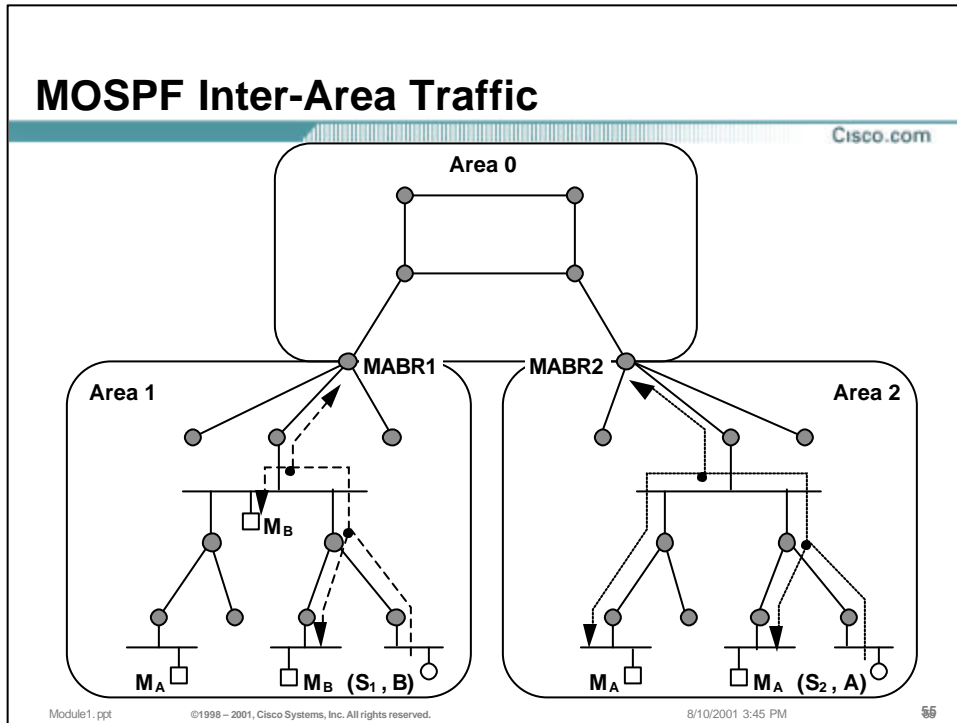
54

- **Wildcard Receivers**

- In order to get multicast traffic to flow between Areas, the concept of “Wildcard Receivers” is used by MOSPF Area Border Routers (MABR).
- Wildcard Receivers set the “Wildcard Receiver” flag in the Router LSA’s that they inject into the Area. This flag is equivalent to a “wildcard” Group Membership LSA that effectively says, “I have a directly connected member for every group.”

MOSPF Inter-Area Traffic

Cisco.com



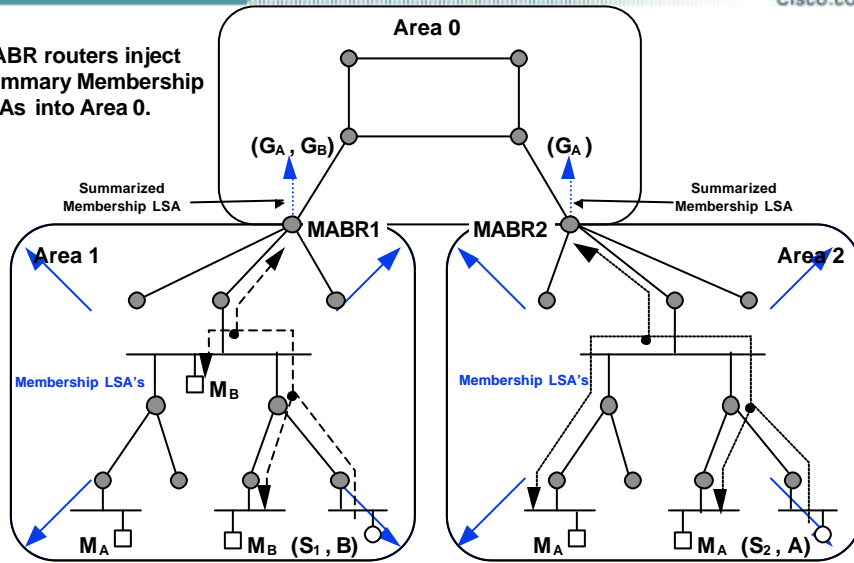
- **Multicast Area Border Routers (MABR)**

- Multicast Area Border routers (i.e. routers that connect an area to the backbone area, Area 0), always set the “Wildcard Receiver” flag in their Router LSA’s that they are injecting into a non-backbone area.
- This causes the MABR to be always be added as a branch of the Shortest Path Tree of any active source in the non-backbone area.
- In the above example, this has resulted in MABR1 being added to the SPT for (S1,B) traffic and MABR2 being added to the SPT for (S2, A) traffic. This “pulls” the source traffic in the area to the border router so that it can be sent into the backbone area.

MOSPF Inter-Area Traffic

Cisco.com

MABR routers inject Summary Membership LSAs into Area 0.



- **Summary Membership LSA's**

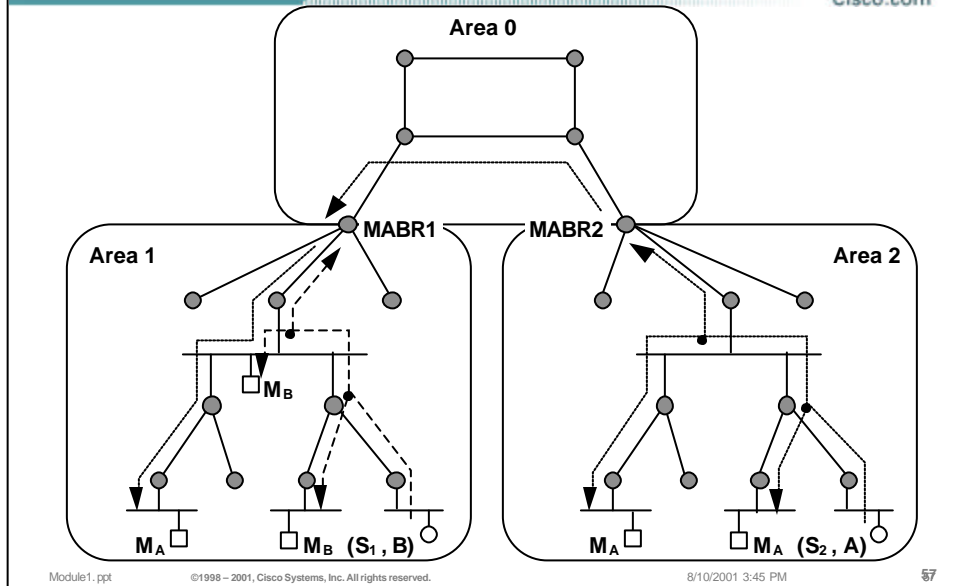
- In addition to Group Membership LSA's MOSPF also defines a new Summary Membership LSA that is used to summarise an area's group membership information.
- Summary Membership LSA's are injected into the backbone area, Area 0 so that routers in the backbone area are made aware of the existence of members in other areas.

- **Inter-Area Traffic Example**

- In the above example, the existence of members of groups "A" and "B" in Area 1 is being injected into the backbone area by MABR1 via Summary Membership LSA.
- In addition, MABR2 is injecting a Summary Membership LSA into the backbone area that indicates that Area 2 has members of group "A".
- Routers in the backbone area now use the information in these Summary Membership LSA's in their Dijkstra calculations to know which MABR's to include in the backbone SPT for which sources. (See next drawing.)

MOSPF Inter-Area Traffic

Cisco.com

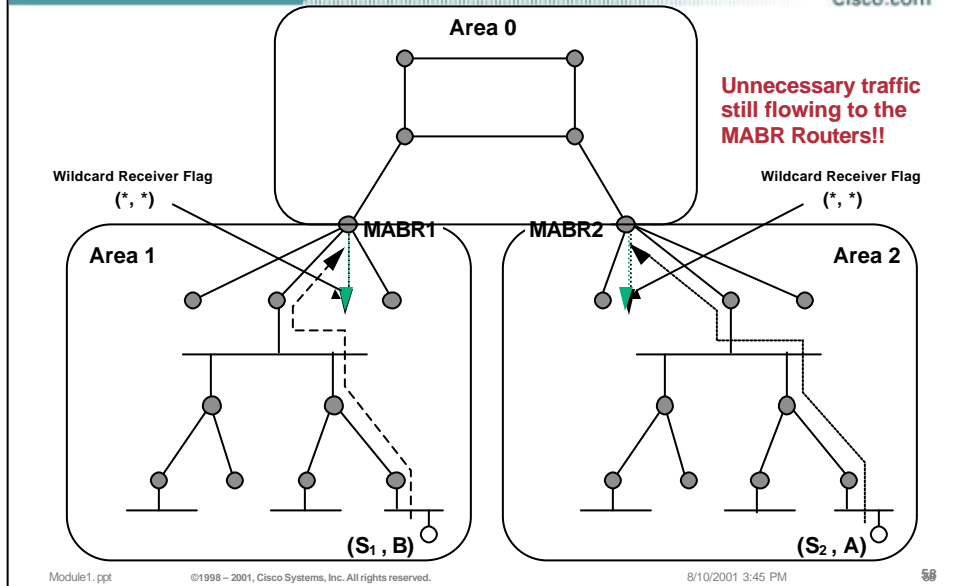


• Inter-Area Traffic Example

- The combination of the “Wildcard Receiver” mechanism and the injection of Summary Membership LSA's into the backbone area permits the SPT for (S2,A) traffic to be extended across the backbone area.
- (S2, A) traffic is now flowing from Area 2 and into the backbone area (Area 0) via MABR2. The routers in the backbone are forwarding this traffic to MABR1 who is sending the traffic into Area 1. Routers inside of Area 1 run the Dijkstra calculation on (S2, A) traffic and construct an (S2, A) SPT inside of Area 1 to route the traffic to members of group “A” as shown above.

MOSPF Inter-Area Traffic

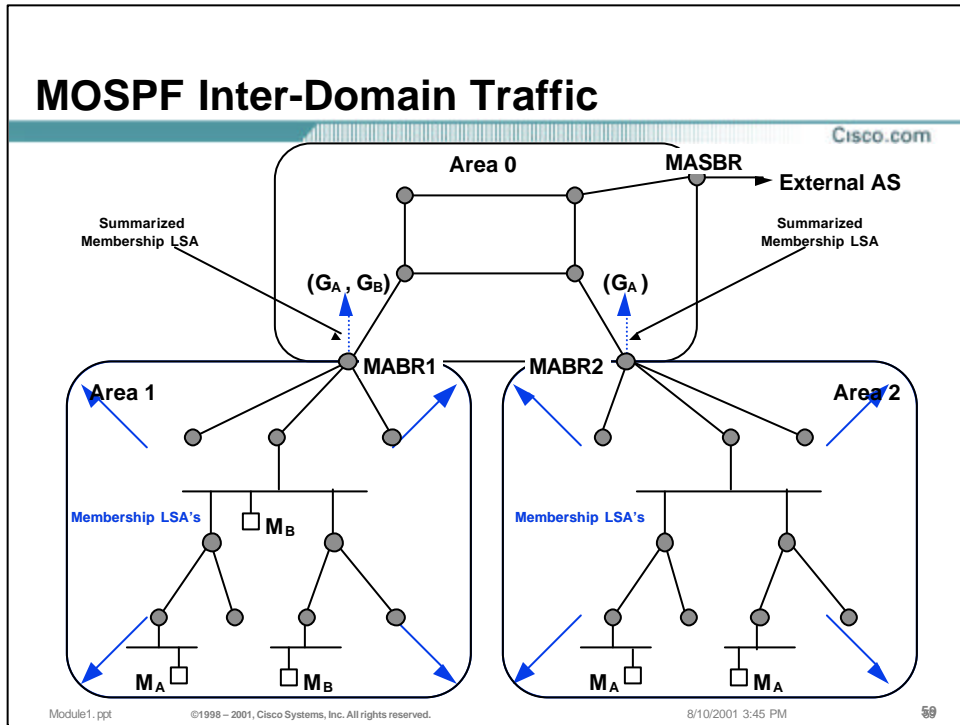
Cisco.com



- **Unnecessary Traffic Flows**

- In the case where there are no members for a multicast group, traffic is still “pulled” to the MABR’s as a result of the “Wildcard Receiver” mechanisms. This can result in bandwidth being consumed inside of the area unnecessarily.

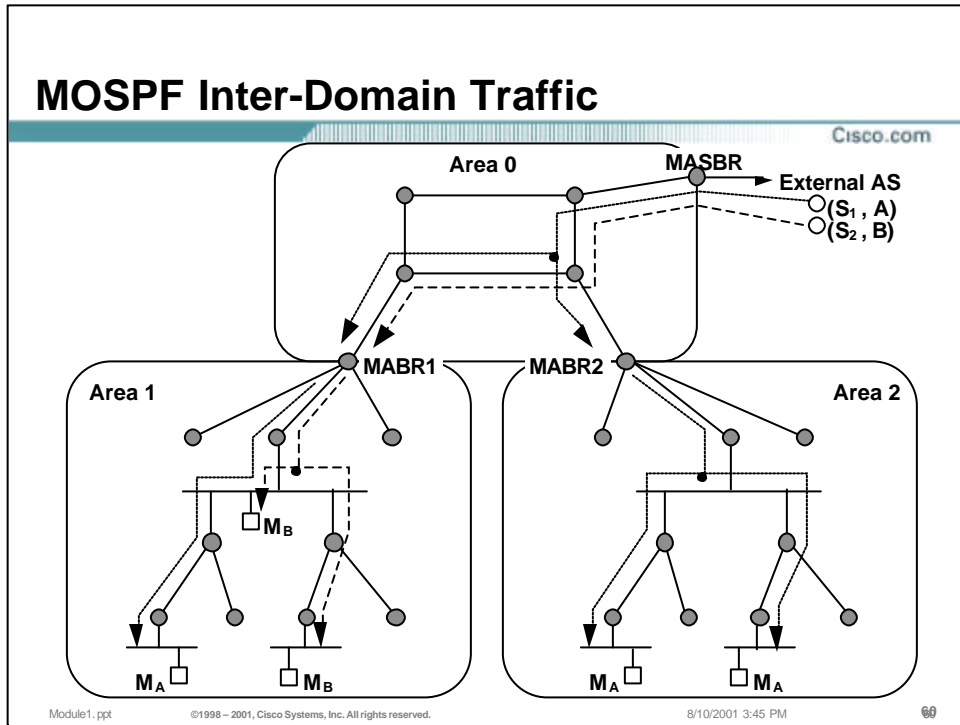
MOSPF Inter-Domain Traffic



- **Inter-Domain Traffic**

- Inter-domain multicast traffic flow basically follows the same mechanisms that were used for Inter-Area traffic flows.
- Summary Membership LSA's inform the routers in the backbone of which MABR's has members of which groups.

MOSPF Inter-Domain Traffic



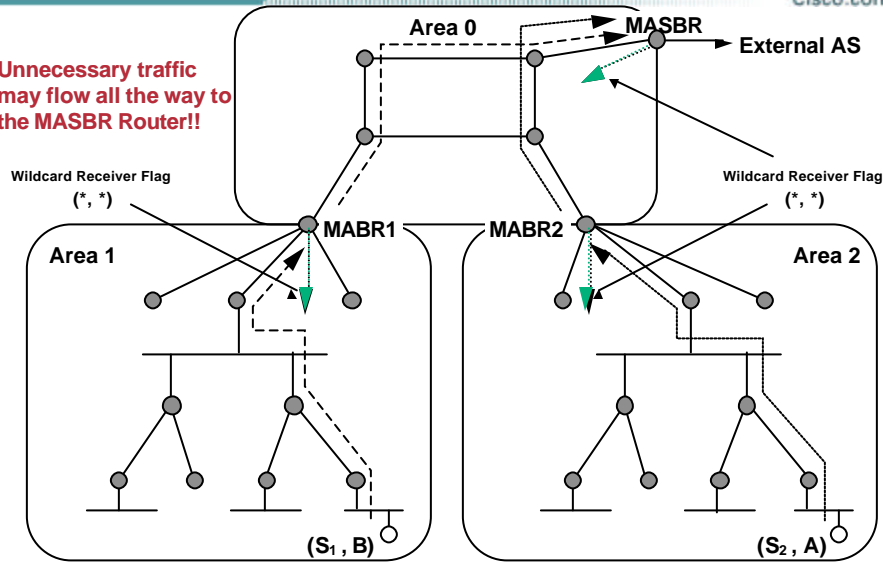
- **Inter-domain Traffic (cont.)**

- When traffic arrives from outside the domain via the Multicast AS Border Router (MASBR), this traffic is forwarded across the backbone to the MABR's as necessary based on the Summary Membership LSA's that they have injected into the area.
- This causes the multicast traffic for group "A" and "B" arriving from outside the AS to be forwarded as shown above.

MOSPF Inter-Domain Traffic

Cisco.com

Unnecessary traffic may flow all the way to the MASBR Router!!



Module1.ppt

©1998 - 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

61

• Inter-Domain Traffic (cont.)

- MASBR's also use the "Wildcard Receiver" mechanism to automatically "pull" all source traffic in the area to them so that they can forward this traffic as needed to the outside world.
- In the example above, the "Wildcard Receiver" mechanism is causing the (S1,B) and (S2,A) traffic to be "pulled" into the backbone area and from there to the MASBR so that it can be forwarded to the outside world.

MOSPF — Evaluation

Cisco.com

- Does not flood multicast traffic everywhere to create state, Uses LSAs and the link-state database
- Protocol dependent—works only in OSPF-based networks
- Significant **scaling problems**
 - Dijkstra algorithm run for EVERY multicast (S_{Net}, G) pair!
 - Dijkstra algorithm rerun when:
 - Group Membership changes
 - Line-flaps
 - Does not support **shared-trees**
- **Not appropriate for...**
 - General purpose multicast networks where the number of senders may be quite large.
 - IP/TV - (Every IP/TV client is a multicast source)

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

62

- **Appropriate for use within single routing domain**
- **Requires OSPF as underlying unicast routing protocol**
- **Significant scaling problems**
 - Frequent flooding of link-state/membership information hinders performance
 - Router CPU demands grow rapidly to keep track of current network topology (source-group pairs)
 - Dijkstra algorithm must be run for every single multicast source
 - Volatility of multicast groups can be lethal
- **Not appropriate for...**
 - Networks with unstable links (too much Dijkstra algorithm computing required for each source)
 - Many simultaneous active source-network/group pairs (routers must maintain too much information relating to the entire network topology)
 - Ubiquitous Multicast Applications permit any user in the network to create a new source-group pair.
 - There is no way for Network Administrator to control the number of source-network/group pairs in the network!!!
 - Therefore, the Network Administrator has little control to prevent MOSPF from “melting down” his/her network as multicast applications become popular with the Users!!!

PIM-DM

Cisco.com

- **Protocol Independent**
 - Supports all underlying unicast routing protocols including: static, RIP, IGRP, EIGRP, IS-IS, BGP, and OSPF
- **Uses reverse path forwarding**
 - Floods network and prunes back based on multicast group membership
 - Assert mechanism used to prune off redundant flows
- **Appropriate for...**
 - Smaller implementations and pilot networks

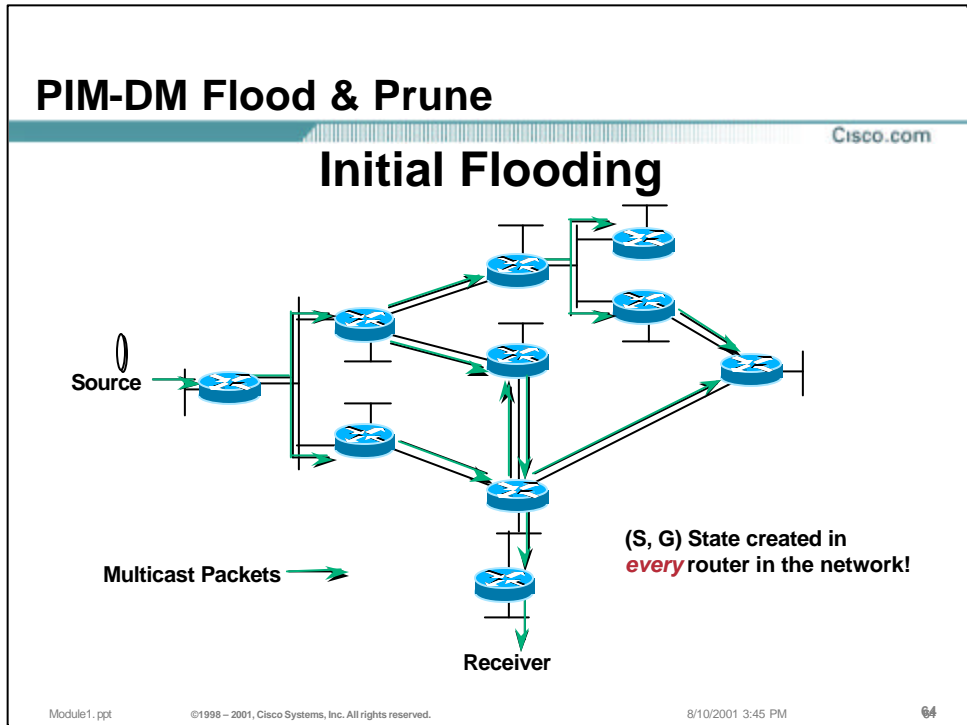
Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

63

- **Protocol Independent Multicast (PIM) Dense-mode (Internet-draft)**
 - Uses Reverse Path Forwarding (RPF) to flood the network with multicast data, then prune back paths based on uninterested receivers
 - Interoperates with DVMRP
- **Appropriate for**
 - Small implementations and pilot networks



- **PIM-DM Initial Flooding**

- PIM-DM is similar to DVMRP in that it initially floods multicast traffic to all parts of the network.
- However unlike DVMRP, which pre-builds a “Truncated Broadcast Tree” that is used for initial flooding, PIM-DM initially floods traffic out ALL non RPF interfaces where there is:
 - Another PIM-DM neighbor or
 - A directly connected member of the group

The reason that PIM-DM does not use “Truncated Broadcast Trees” to pre-build a spanning tree for each source network is that this would require running a separate routing protocol as does DVMRP. (At the very least, some sort of Poison-Reverse messages would have to be sent to build the TBT.) Instead, PIM-DM uses other mechanisms to prune back the traffic flows and build Source Trees.

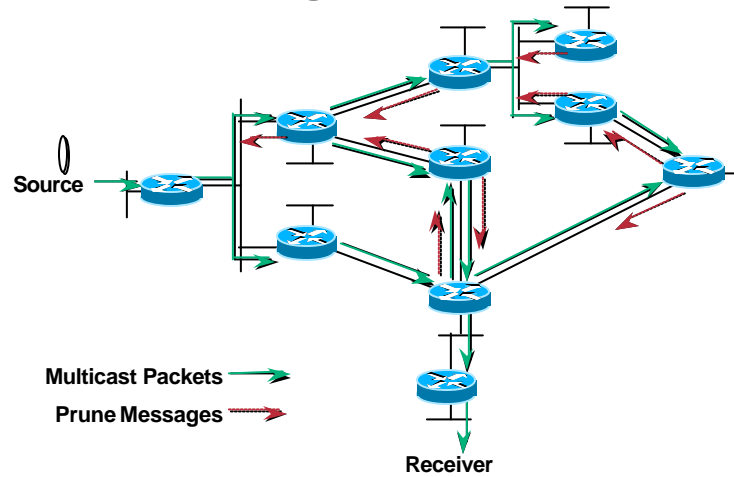
- **Initial Flooding Example**

- In this example, multicast traffic being sent by the source is flooded throughout the entire network.
- As each router receives the multicast traffic via its RPF interface (the interface in the direction of the source), it forwards the multicast traffic to all of its PIM-DM neighbors.
- Note that this results in some traffic arriving via a non-RPF interface such as the case of the two routers in the center of the drawing. (Packets arriving via the non-RPF interface are discarded.) These non-RPF flows are normal for the initial flooding of data and will be corrected by the normal PIM-DM pruning mechanism.

PIM-DM Flood & Prune

Cisco.com

Pruning unwanted traffic



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

65

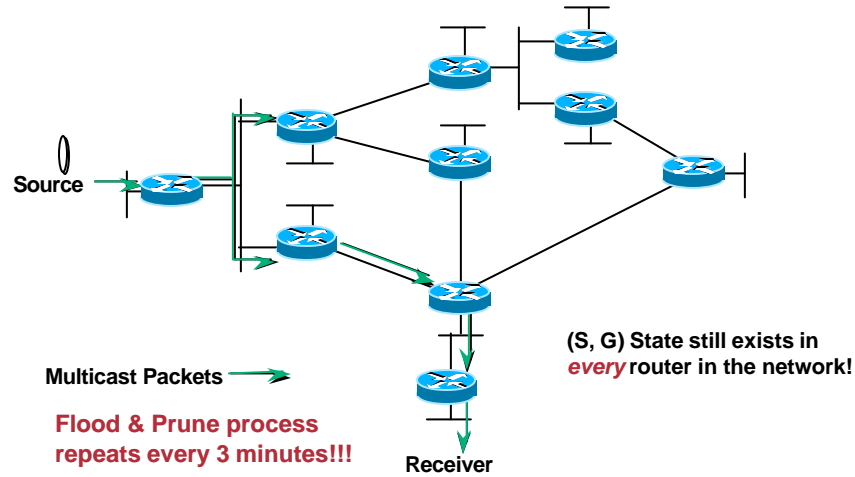
- **Pruning unwanted traffic**

- In the example above, PIM Prunes (denoted by the dashed arrows) are sent to stop the flow of unwanted traffic.
- Prunes are sent on the RPF interface when the router has no downstream members that need the multicast traffic.
- Prunes are also sent on non-RPF interfaces to shutoff the flow of multicast traffic that is arriving via the wrong interface (i.e. traffic arriving via an interface that is not in the shortest path to the source.)
 - An example of this can be seen at the second router from the receiver near the center of the drawing. Multicast traffic is arriving via a non-RPF interface from the router above (in the center of the network) which results in a Prune message.

PIM-DM Flood & Prune

Cisco.com

Results after Pruning



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

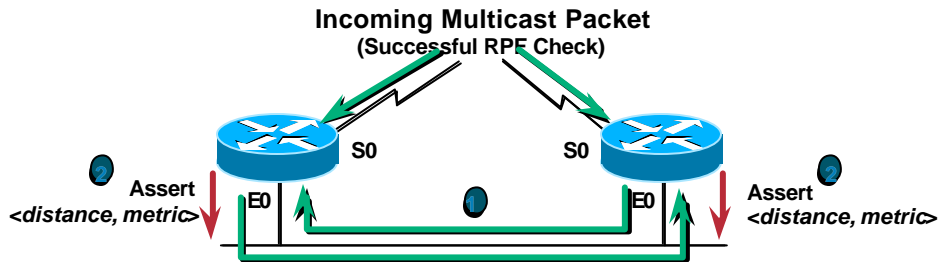
66

- **Results after Pruning**

- In the final drawing in our example shown above, multicast traffic has been pruned off of all links except where it is necessary. This results in a Shortest Path Tree (SPT) being built from the Source to the Receiver.
- Even though the flow of multicast traffic is no longer reaching most of the routers in the network, (S, G) state still remains in ALL routers in the network. This (S, G) state will remain until the source stops transmitting.
- In PIM-DM, Prunes expire after three minutes. This causes the multicast traffic to be re-flooded to all routers just as was done in the “Initial Flooding” drawing. This periodic (every 3 minutes) “Flood and Prune” behavior is normal and must be taken into account when the network is designed to use PIM-DM.

PIM-DM Assert Mechanism

Cisco.com



- 1 Routers **receive** packet on an interface in their “*oilst*”!!
 - Only one router should continue sending to avoid duplicate packets.
- 2 Routers send “PIM Assert” messages
 - Compare *distance* and *metric* values
 - Router with best route to source wins
 - If *metric* & *distance* equal, highest IP adr wins
 - Losing router stops sending (prunes interface)

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

67

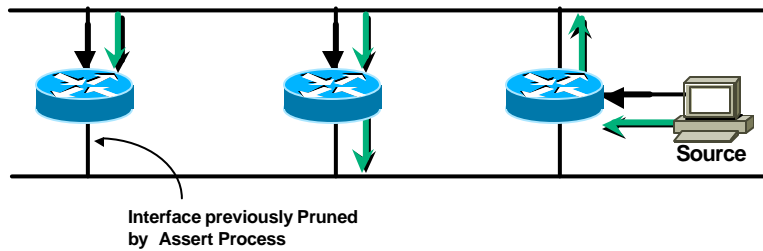
• PIM Assert Mechanism

- The PIM Assert mechanism is used to shutoff duplicate flows onto the same multi-access network.
 - Routers detect this condition when they receive an (S, G) packet via a multi-access interface that it is in the (S, G) OIL.
 - This causes the routers to send Assert Messages.
- Assert messages containing the Admin. Distance and metric to the source combined into a single assert value. (The Admin. Distance is the high-order part of this assert value.)
- Routers compare these values to determine who has the best path (lowest value) to the source. (If both values are the same, the highest IP address is used as the tie breaker.)
- The Losing routers (the ones with the higher value) Prunes its interface while the winning router continues to forward multicast traffic onto the LAN segment.

Potential PIM-DM Route Loop

Cisco.com

Normal Steady-State Traffic Flow



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

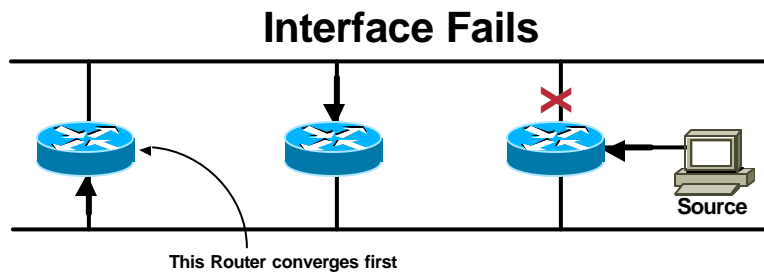
68

• Potential PIM-DM Route Loops

- The non-deterministic behavior of PIM-DM along with its flood-and-prune mechanism can sometimes result in serious network outages including “blackholes” and multicast route loops.
- The network in the above example is a simplified version of a frequently used network design whereby multiple routers are used to provide redundancy in the network.
- Under normal steady-state conditions, traffic flows from the source via the RPF interfaces as shown.
 - Note that the routers have performed the Assert process and one interface on one router is in the pruned state.

Potential PIM-DM Route Loop

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

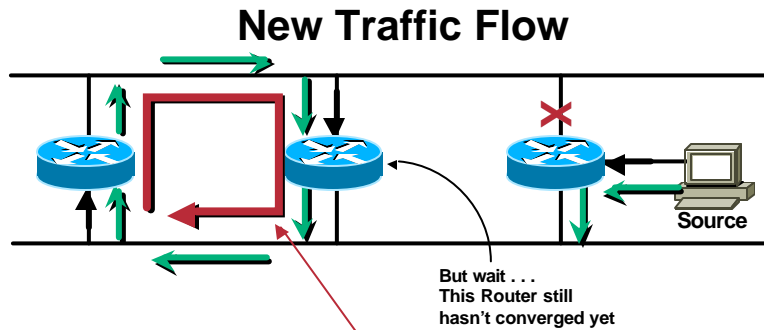
69

- **Potential PIM-DM Route Loops**

- Now let's assume that the forwarding interface of the first-hop router fails as shown above.
- Let's also assume that the unicast routing of router on the left converges first and PIM computes the new RPF interface as shown.

Potential PIM-DM Route Loop

Cisco.com



Multicast Route Loop !!



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

70

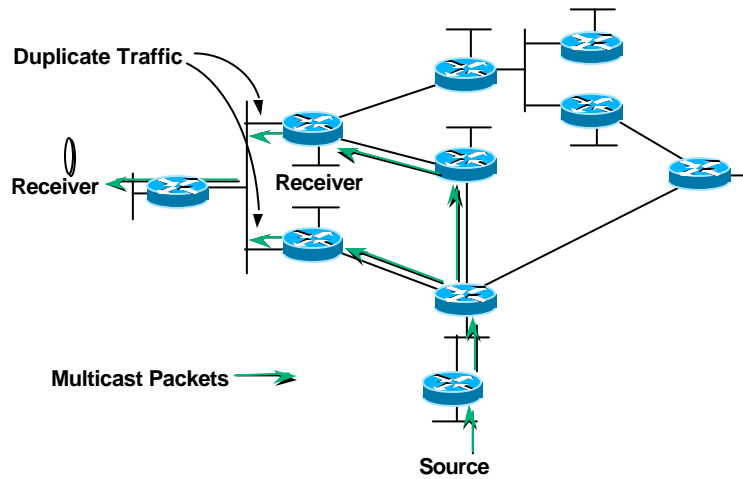
• Potential PIM-DM Route Loops

- Unfortunately, the middle router has not yet converged and is still forwarding multicast traffic using the old RPF interface.
- At this point, a multicast route loop exists in the network due to the transient condition of the two routers having opposite RPF interfaces.
- During the time that this route loop exists, virtually all of the bandwidth on the network segments can be consumed. This situation will continue until the router in the middle of the picture finally converges and the new “correct” RPF interface is calculated.
- Unfortunately, if the router needs some bandwidth to complete this convergence (as in the case when EIGRP goes active), then this condition will never be resolved!

PIM-DM Assert Problem

Cisco.com

Initial Flow



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

71

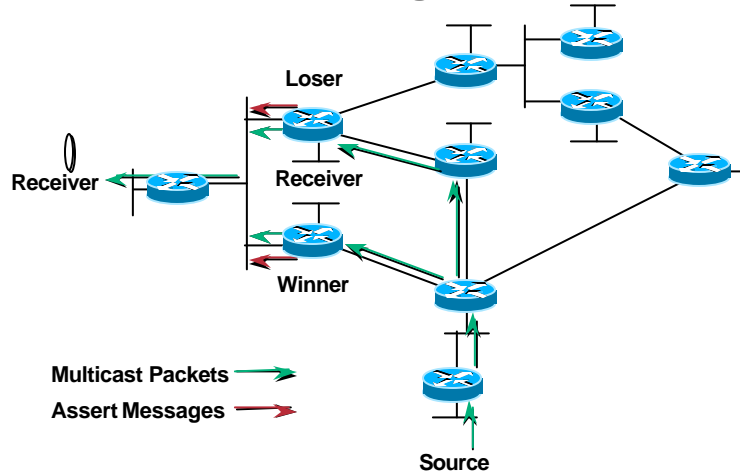
- **PIM-DM Assert Problem**

- While the PIM Assert mechanism is effective in pruning off duplicate traffic, it is not without its weaknesses.
- Consider the above example where duplicate traffic is flowing onto a LAN segment.

PIM-DM Assert Problem

Cisco.com

Sending Asserts



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

72

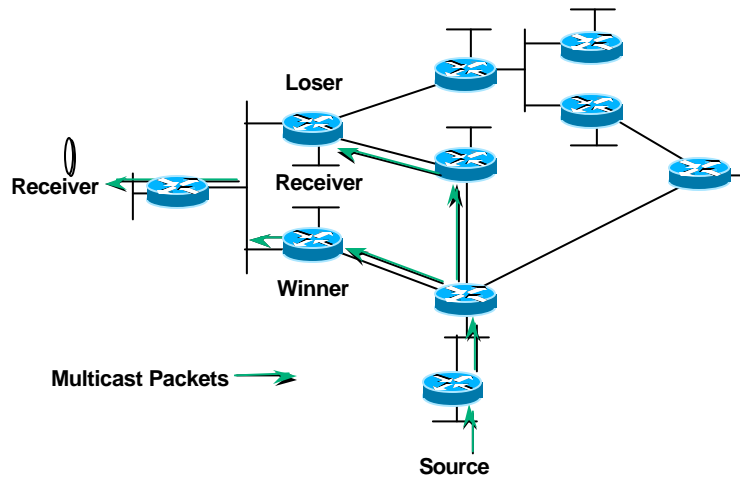
- **PIM-DM Assert Problem**

- The normal PIM Assert mechanism takes place and the two routers exchange routing metrics to determine which one has the best route to the source.
- In this case, the bottom router has the best metric and is the Assert Winner.

PIM-DM Assert Problem

Cisco.com

Assert Loser Prunes Interface



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

73

- **PIM-DM Assert Problem**

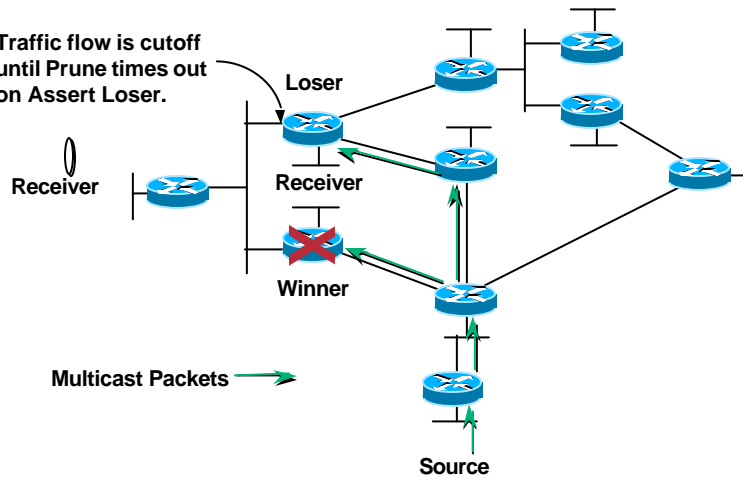
- The normal PIM Assert mechanism takes place and the Assert Winner continues forwarding while the Assert Loser prunes its interface and starts its prune timer.

PIM-DM Assert Problem

Cisco.com

Assert Winner Fails

Traffic flow is cutoff until Prune times out on Assert Loser.



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

74

- **PIM-DM Assert Problem**

- Let's now assume that the Assert Winner fails immediately after winning the Assert process.
- Unfortunately, the Assert Loser has no way of knowing that the Assert Winner has failed and will wait 3 minutes before timing out its pruned interface. This results in a 3 minute (worst-case) loss of traffic.

PIM-DM — Evaluation

Cisco.com

- **Most effective for small pilot networks**
- **Advantages:**
 - Easy to configure—two commands
 - Simple flood and prune mechanism
- **Potential issues...**
 - Inefficient flood and prune behavior
 - Complex Assert mechanism
 - Mixed control and data planes
 - Results in (S, G) state in every router in the network
 - Can result in non-deterministic topological behaviors
 - No support for shared trees

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

76

- **Evaluation: PIM Dense-mode**

- Most effective for small pilot networks.
- Advantages
 - Minimal number of commands required for configuration (two)
 - Simple mechanism for reaching all possible receivers and eliminating distribution to uninterested receivers
 - Simple behavior is easier to understand and therefore easier to debug
 - Interoperates with DVMRP
- Potential issues
 - Necessity to flood frequently because prunes expire after 3 minutes.

Sparse-Mode Protocols

Cisco.com

- **PIM SM- Protocol Independent Multicasting (Sparse Mode)**
- **CBT - Core Based Trees**

PIM-SM (RFC 2362)

Cisco.com

- Supports both source and shared trees
 - Assumes no hosts want multicast traffic unless they specifically ask for it
- Uses a **Rendezvous Point (RP)**
 - Senders and Receivers “rendezvous” at this point to learn of each others existence.
 - Senders are “registered” with RP by their first-hop router.
 - Receivers are “joined” to the Shared Tree (rooted at the RP) by their local Designated Router (DR).
- Appropriate for...
 - Wide scale deployment for **both** densely and sparsely populated groups in the enterprise
 - Optimal choice for all production networks regardless of size and membership density.

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

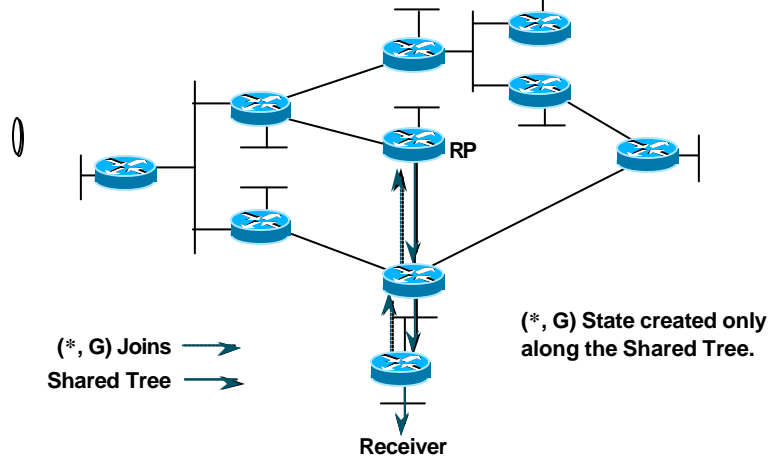
8/10/2001 3:45 PM

77

- **Protocol Independent Multicast (PIM) Sparse-mode (RFC 2362)**
 - Utilizes a rendezvous point (RP) to coordinate forwarding from source to receivers
 - Regardless of location/number of receivers, senders register with RP and send a single copy of multicast data through it to registered receivers
 - Regardless of location/number of sources, group members register to receive data and always receive it through the RP
 - Appropriate for
 - Wide scale deployment for both densely and sparsely populated groups in the Enterprise
 - Optimal choice for all production networks regardless of size and membership density.

PIM-SM Shared Tree Joins

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

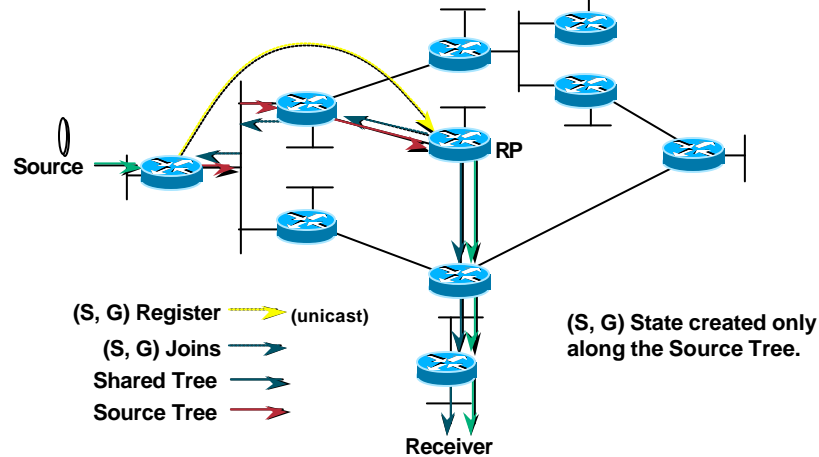
78

- PIM-SM Shared Tree Joins

- In this example, there is an active receiver (attached to leaf router at the bottom of the drawing) has joined multicast group “G”.
- The leaf router knows the IP address of the Rendezvous Point (RP) for group G and when it sends a (*,G) Join for this group towards the RP.
- This (*, G) Join travels hop-by-hop to the RP building a branch of the Shared Tree that extends from the RP to the last-hop router directly connected to the receiver.
- At this point, group “G” traffic can flow down the Shared Tree to the receiver.

PIM-SM Sender Registration

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

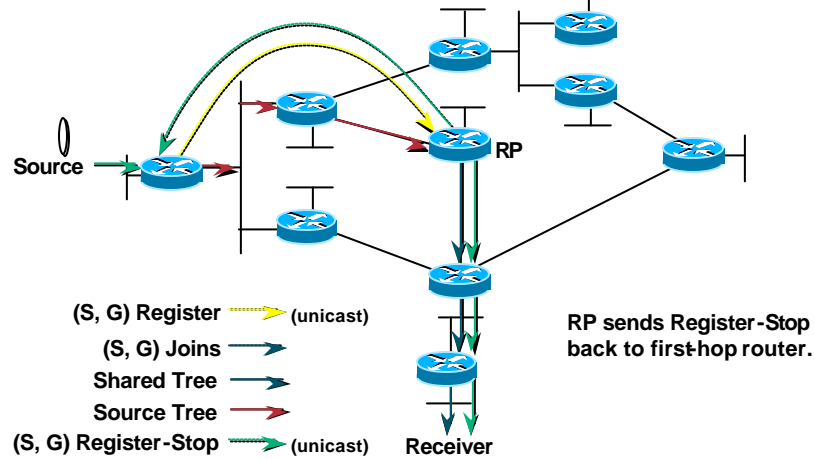
79

• PIM-SM Sender Registration

- As soon as an active source for group G sends a packet the leaf router that is attached to this source is responsible for “Registering” this source with the RP and requesting the RP to build a tree back to that router.
- The source router encapsulates the multicast data from the source in a special PIM SM message called the Register message and unicasts that data to the RP.
- When the RP receives the Register message it does two things
 - It de-encapsulates the multicast data packet inside of the Register message and forwards it down the Shared Tree.
 - The RP also sends an (S,G) Join back towards the source network S to create a branch of an (S, G) Shortest-Path Tree. This results in (S, G) state being created in all the router along the SPT, including the RP.

PIM-SM Sender Registration

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

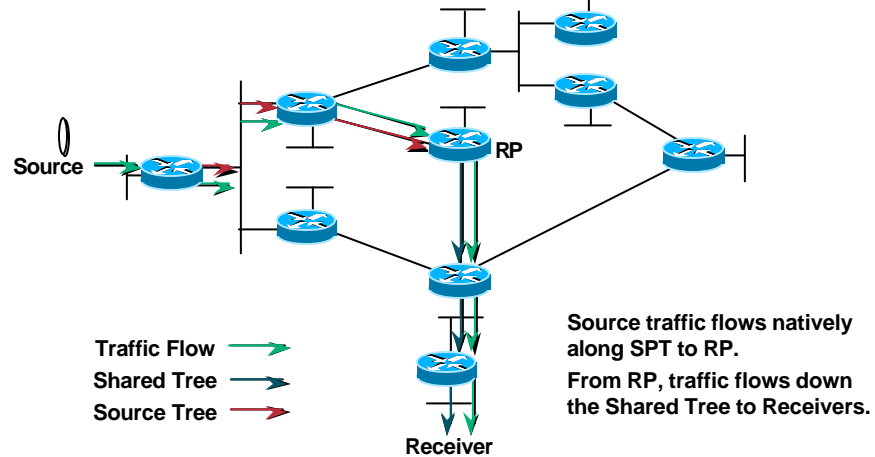
80

- **PIM-SM Sender Registration (cont.)**

- As soon as the SPT is built from the Source router to the RP, multicast traffic begins to flow natively from source S to the RP.
- Once the RP begins receiving data natively (i.e. down the SPT) from source S it sends a 'Register Stop' to the source's first hop router to inform it that it can stop sending the unicast Register messages.

PIM-SM Sender Registration

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

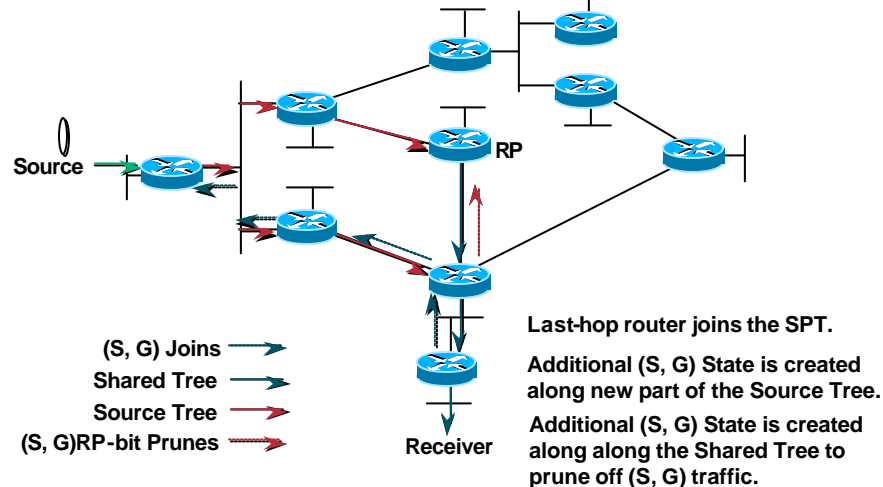
81

- **PIM-SM Sender Registration (cont.)**

- At this point, multicast traffic from the source is flowing down the SPT to the RP and from there, down the Shared Tree to the receiver.

PIM-SM SPT Switchover

Cisco.com



Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

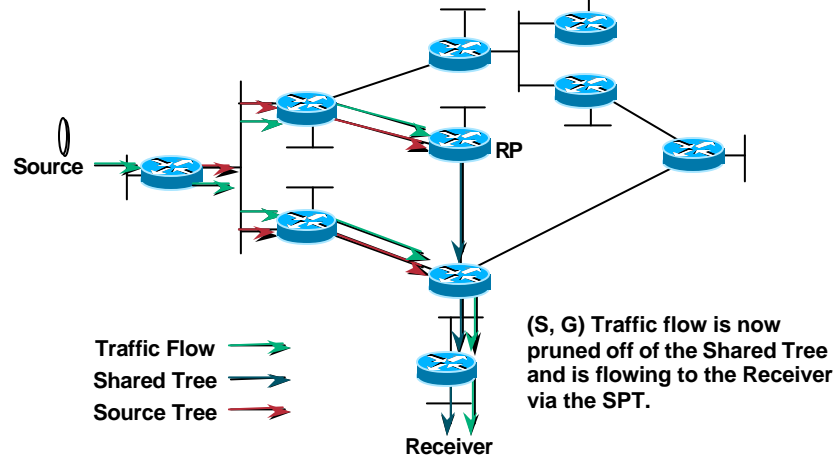
82

• PIM-SM Shortest-Path Tree Switchover

- PIM-SM has the capability for last-hop routers (i.e. routers with directly connected members) to switch to the Shortest-Path Tree and bypass the RP if the traffic rate is above a set threshold called the “SPT-Threshold”.
 - The default value of the SPT-Threshold in Cisco routers is zero. This means that the default behaviour for PIM-SM leaf routers attached to active receivers is to immediately join the SPT to the source as soon as the first packet arrives via the (*,G) shared tree.
- In the above example, the last-hop router (at the bottom of the drawing) sends an (S, G) Join message toward the source to join the SPT and bypass the RP.
- This (S, G) Join messages travels hop-by-hop to the first-hop router (i.e. the router connected directly to the source) thereby creating another branch of the SPT. This also creates (S, G) state in all the routers along this branch of the SPT.
- Finally, special (S, G)RP-bit Prune messages are sent up the Shared Tree to prune off this (S,G) traffic from the Shared Tree.
 - If this were not done, (S, G) traffic would continue flowing down the Shared Tree resulting in duplicate (S, G) packets arriving at the receiver.

PIM-SM SPT Switchover

Cisco.com



Module1.ppt

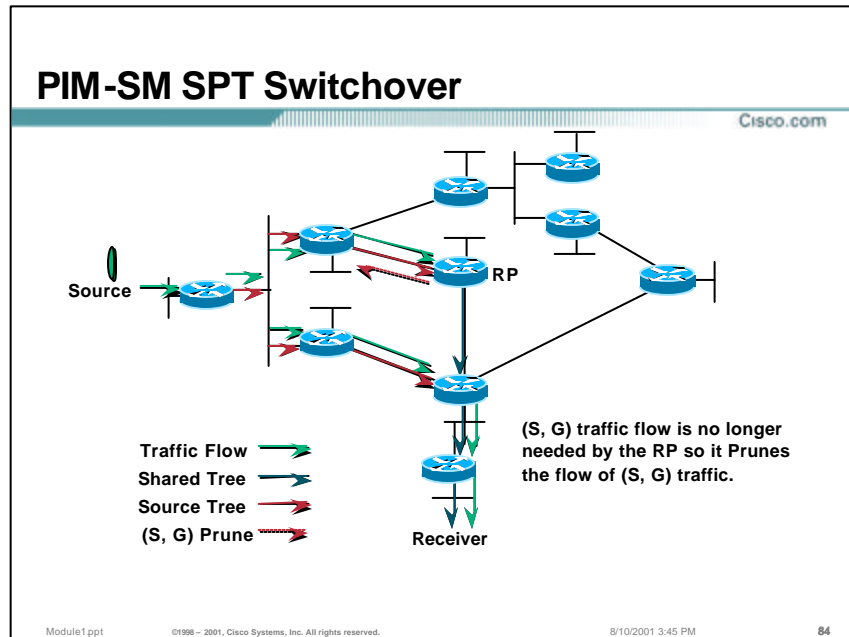
©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

83

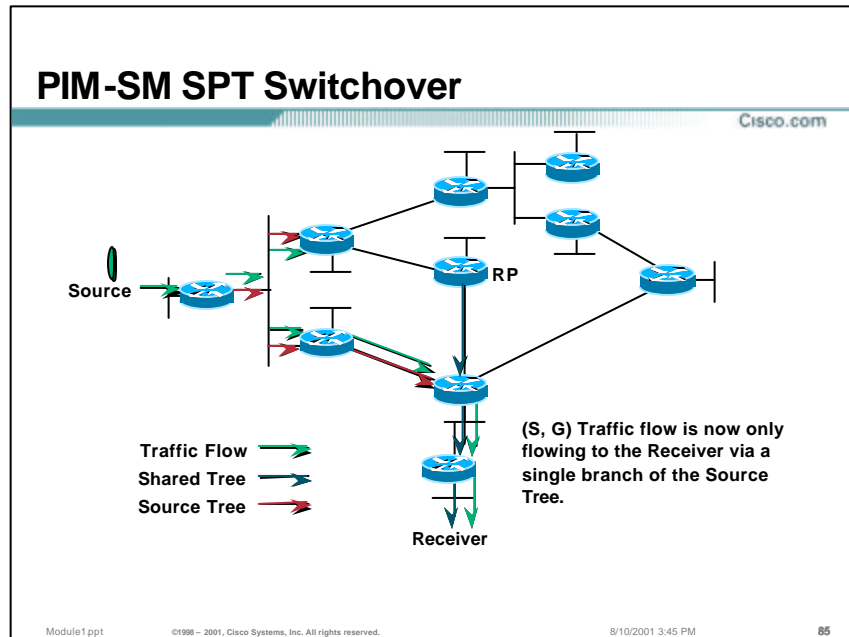
- **PIM-SM Shortest-Path Tree Switchover**

- At this point, (S, G) traffic is now flowing directly from the first-hop router to the last-hop router and from there to the receiver.
 - Note: The RP will normally send (S, G) Prunes back toward the source to shutoff the flow of now unnecessary (S, G) traffic to the RP *IFF* it has received an (S, G)RP-bit Prune on all interfaces on the Shared Tree. (This step has been omitted from the example above.)
- As a result of this SPT-Switchover mechanism, PIM SM also supports the construction and use of SPT (S,G) trees but in a much more economical fashion than PIM DM in terms of forwarding state.



- **PIM-SM Shortest-Path Tree Switchover**

- At this point, the RP no longer needs the flow of (S, G) traffic since all branches of the Shared Tree (in this case there is only one) have pruned off the flow of (S, G) traffic.
- As a result, the RP will send (S, G) Prunes back toward the source to shutoff the flow of the now unnecessary (S, G) traffic to the RP
 - Note: This will occur *IFF* the RP has received an (S, G)RP-bit Prune on all interfaces on the Shared Tree.



- **PIM-SM Shortest-Path Tree Switchover**

- As a result of the SPT-Switchover, (S, G) traffic is now only flowing from the first-hop router to the last-hop router and from there to the receiver. Notice that traffic is no longer flowing to the RP.
- As a result of this SPT-Switchover mechanism, it is clear that PIM SM also supports the construction and use of SPT (S,G) trees but in a much more economical fashion than PIM DM in terms of forwarding state.

PIM-SM Frequently Forgotten Fact

“The default behavior of PIM-SM is that routers with directly connected members will join the Shortest Path Tree as soon as they detect a new multicast source.”

- **Frequently Forgotten Fact**

- Unless configured otherwise, the default behaviour of Cisco routers running PIM-SM is for last-hop routers to immediately switch to the SPT for any new source.

PIM-SM — Evaluation

Cisco.com

- **Effective for **sparse or dense** distribution of multicast receivers**
- **Advantages:**
 - Traffic only sent down “joined” branches
 - Can switch to optimal source-trees for high traffic sources dynamically
 - Unicast routing protocol-independent
 - Basis for inter-domain multicast routing
 - When used with MBGP and MSDP

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

87

- **Evaluation: PIM Sparse-mode**

- Can be used for sparse or dense distribution of multicast receivers (no necessity to flood)
- Advantages
 - Traffic sent only to registered receivers that have explicitly joined the multicast group
 - RP can be switched to optimal shortest-path-tree when high-traffic sources are forwarding to a sparsely distributed receiver group
 - Interoperates with DVMRP
- Potential issues
 - Requires RP during initial setup of distribution tree (can switch to shortest-path-tree once RP is established and determined suboptimal)

CBT Overview

Cisco.com

- **Constructs single, shared delivery tree (not source-based) for multicast group members**
 - Traffic is sent and received over same tree, regardless of source(s)
 - Reduced amount of multicast state information stored in routers
- **Uses core router to construct shared tree**
 - Routers send join message to core and form branch of tree, suppressing downstream join messages
 - Downstream routers connect to shared tree through on-tree routers
 - Source unicasts data to core, then multicasts using group ID
 - Aggregates traffic onto smaller subset of links
- **No Commercial implementation available**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

88

- **Core Based Trees (Internet-draft)**

- Utilizes shared delivery tree constructed around core router (much like PIM's RP)
 - Unlike PIM, the Shared tree is bi-directional.
 - If the first-hop router for a source is already on the tree, it forwards the multicast packets out all branches of the tree.
 - If the first-hop router for a source is not on the Shared tree, a single copy of multicast data is sent through the core router to receivers.
 - Regardless of location/number of sources, group members always receive multicast data through through the Shared tree.
- Key benefits
 - Reduced amount of multicast state information stored in routers (always send and receive over same distribution tree)
 - Traffic is aggregated onto smaller subset of links

CBT — Evaluation

Cisco.com

- **Academic work-in-progress**
- **Runs primarily on MS Power Point**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

89

- **Evaluation: CBT**

- Appropriate for inter- and intra-domain multicast routing (no necessity to flood)
- Current Deployment
 - New protocol that is not widely deployed in production environments (no commercial implementation available)
 - Improves scalability of some existing multicast algorithms to support sparse distribution of multicast receivers
 - Interoperates with DVMRP
- Potential issue
 - Has no capability to switch to SPT
 - Can suffer from latency problems since traffic must flow through the Core router.
 - Core routers can become bottlenecks if not selected with great care, especially when senders and receivers are located very far from each other

Protocol Summary

Cisco.com

CONCLUSION

**“Virtually all production networks
should be configured to run PIM in
Sparse mode!”**

Module1.ppt

©1998 – 2001, Cisco Systems, Inc. All rights reserved.

8/10/2001 3:45 PM

90

- **Protocol Summary**

- Given the pros and cons of all the multicast routing protocols available, virtually all production networks should be configured to run PIM SM.

CISCO SYSTEMS



EMPOWERING THE
INTERNET GENERATIONSM