

JANET Multicast Workshop



Stig Venaas - stig.venaas@uninett.no

Tim Chown - tjc@ecs.soton.ac.uk

David Mills - dgm@ecs.soton.ac.uk

University of Southampton

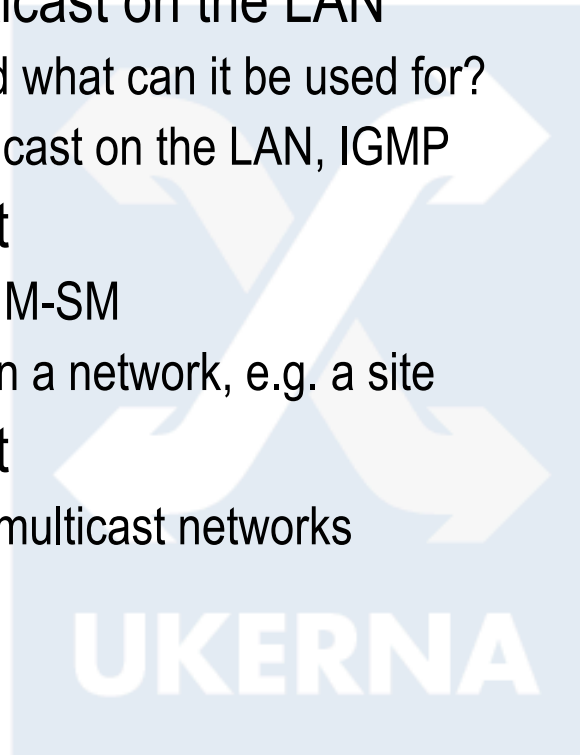
&

Gorry Fairhurst – gorry@erg.abdn.ac.uk

University of Aberdeen

Overview

- Introduction and multicast on the LAN
 - What is multicast and what can it be used for?
 - Addressing and multicast on the LAN, IGMP
- Intradomain multicast
 - Multicast routing – PIM-SM
 - Deploying multicast in a network, e.g. a site
- Interdomain multicast
 - How to interconnect multicast networks
 - MSDP and MBGP
 - Multicast in JANET
- IPv6 multicast
 - A taste of IPv6 multicast
 - IPv6 multicast addressing
 - MLD, SSM and embedded-RP





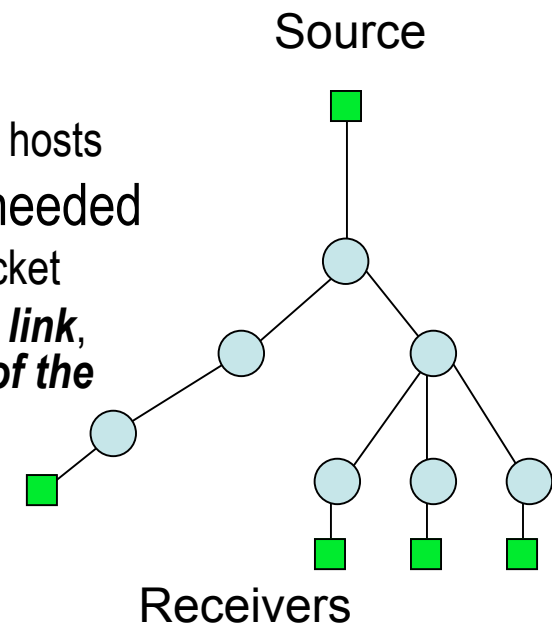
Session 1

Introduction and multicast on the LAN



What is IP multicast?

- Usually an IP packet is sent to one specific host
 - The IP destination address specifies which host
- With IP multicast, an IP packet is sent to a group of hosts
 - The IP destination address is a group and not a host address
 - IPv4 multicast addresses, class D. 224.0.0.0 – 239.255.255.255
 - The group can contain any number of hosts (0 to infinity)
 - The group members can be anywhere
 - Like IP subnet broadcast:
 - A single packet is received by all on the subnet.
 - Multicast is not restricted to the subnet, and not sent to all hosts
- Multicast packets will be replicated by routers where needed
 - Routers keep track of which interfaces should forward the packet
 - The same multicast packet is ***never sent twice on the same link***, hence the bandwidth used on a specific link is ***independent of the number of receivers***



Why is it useful?

- Imagine the BBC streaming TV on the Internet to every UK home
 - Multicast only needs a basic machine and typical home Internet connectivity
 - Remember, to send you don't need more bandwidth than a single receiver
- An ADSL user could send video to thousands of other users
 - The number of receivers is not an issue
- Useful for multi-party applications (conferencing or gaming)
 - Where each participant wants to send the same data to all others
- For financial and gaming applications
 - It may be important to deliver quickly and simultaneously to many recipients
- Multicast also useful for discovery
 - Imagine all printers on your network joining a specific multicast group
 - Can query all printers (and no other hosts) asking them to identify themselves

Service Model and Routing Challenges

- The basic multicast service model is:
 - Anyone can send to the multicast group
 - Senders don't need to know where the receivers are or how many (if any)
 - Hosts interested in the group join it
 - They don't need to know who is sending, where they are, or what other receivers there are
 - They just receive anything sent by anyone to the group while they are members
- Source-Specific Multicast (SSM, RFC 3569)
 - SSM is a new model where receivers specify the sources when joining
 - i.e. receivers need to know who the sources are
- The big challenge is routing
 - If anyone can be anywhere (only telling routers which group they are sending to or joining) how can routers learn from where and to where they should forward the data?
- In the beginning there was multicast only on Ethernet links, no routing
 - Trivial, especially with just a single coax cable or a hub
 - With switches it's more complicated
- Then one wanted to do routing across larger networks...

IPv4 Multicast Addressing

- IPv4 multicast addresses: 224.0.0.0 – 239.255.255.255 (224/4, class D)
- These are subdivided in rather complicated ways,
 - see <http://www.iana.org/assignments/multicast-addresses/> for details
- Examples:
 - 224.0.0.0 – 224.0.0.255 (224.0.0/24) – Local network control block, never forwarded
 - 224.0.0.1 - All local hosts
 - 224.0.0.2 - All local routers
 - 224.0.0.5 - OSPF
 - 224.0.0.13 - PIM
 - 224.0.0.22 - IGMP
 - 224.0.1.0 – 224.0.1.255 (224.0.1/24) – Internetwork control block, forwarded
 - 224.2/16 – SAP
 - 232/8 – SSM (only to be used for Source Specific Multicast)
 - 233/8 – GLOP
 - 234.0.0.0 – 238.255.255.255 – Reserved
 - 239/8 – Administrative scoping

Address Assignment, SAP and GLOP

- Knowing what addresses to use when creating a session seems rather complicated
- SAP (Session Announcement Protocol, RFC 2974)
 - Announces a session
 - SAP applications also help you pick what addresses to use
 - Uses dynamic groups in range 224.2.128.0 – 224.2.255.255 for global sessions
 - Global announcements sent to 224.2.127.254
 - sdr is the most common SAP application, but not used so much these days
- GLOP (not an acronym)
 - Assignment based on AS numbers, RFC 3180
 - 233.x.x/24 where x.x is an officially assigned AS number
 - For private AS space there is EGLOP (RFC 3138)
managed by registries, e.g. RIPE (still 233.x.x/24, but with private AS numbers)

Administrative Scoping – 239/8

- Addresses in the range 239/8 are used for administrative scoping
 - Private address space, not to be used globally
 - Different networks can use the same addresses
- 239.255/16 is the smallest administrative scope
 - Sometimes used for site-local
- 239.192/14 is organization-local scope
 - These addresses should work throughout JANET
 - All but 239.194/16 are restricted to JANET
 - 239.192/16, 239.193/16 and 239.195/16 used for sessions visible throughout JANET, but not outside
 - 239.194/16 is used for GÉANT
 - i.e. sessions using these groups are available throughout GÉANT (European academic networks), but not outside
- Multicast distribution can be restricted by specifying a small TTL value for packets
 - Limited use. With routing protocols like PIM-SM and MSDP, packets may travel very far even if TTL is small

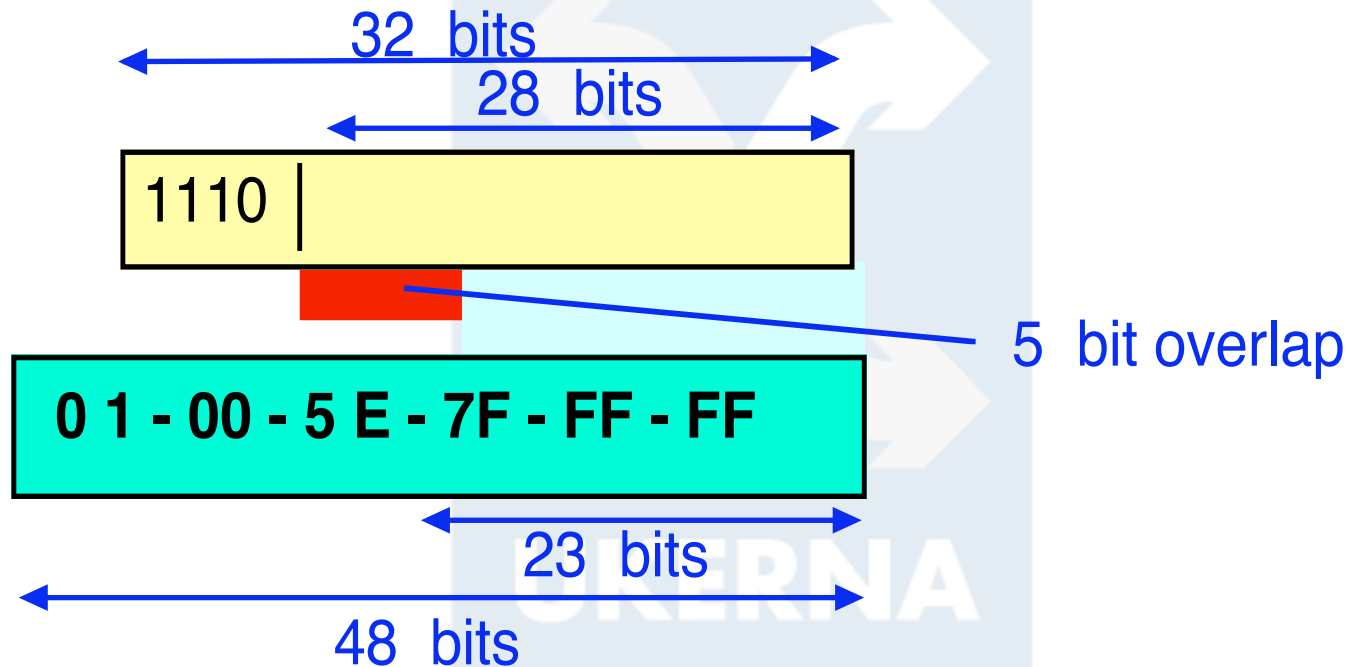
Multicast on the LAN

- Multicast is defined for Ethernet
 - Ethernet multicast is exactly like traditional IP multicast model
 - IP multicast service is based on the Ethernet service model extended from working on a LAN to the Internet
- Originally Ethernet multicast was very simple
 - Any host can send
 - All packets go everywhere (coax cables or hubs)
 - Any host on the LAN can choose to listen, only need to tell NIC what packets to pick up
 - But then came bridges and switches...

Mapping to Ethernet MAC

Class D IPv4 destination address

224.0.0.0-239.255.255.255



MAC hardware destination address

*One L2 (MAC) address
may carry multiple L3 (IPv4) addresses*

Multicast IP Packet Decode

ETHER

Packet size = 218 bytes

Destination = 1:0:5e:2:dc:3e, (multicast) (01-00-5e-02-dc-3e)

Source = 0:d0:bb:f7:c6:c0,

Ethertype = 0800 (IPv4)

IP

Version = 4, Header length = 20 bytes

Type of service = 0x00

Total length = 204 bytes (00cc)

ID = 57862, Flags = 0x00, Frags = 0

Time To Live = 113 seconds/hops

Protocol = 17 (UDP)

Header checksum = a1a9

Source address = 132.185.132.118

Destination address = 224.2.220.62

No options

UDP

Source port = 31106 (7982)

Destination port = 31106 (7982)

Length = 184 (00b8)

Checksum = 08a0

RTP

180B of Data

0:	0100	5e02	dc3e	00d0	bbf7	c6c0	0800	4500
16:	00cc	e206	0000	7111	a1a9	84b9	8476	e002
32:	dc3e	7982	7982	00b8	08a0	8005	dbc6	d721
48:	69c0	0752	bb5f	fe39	3600	8808	b120	8933
64:	6219	9118	5128	ffc8	1321	bc10	933e	aa23
80:	3233	ba00	e892	a00c	1a3c	0a28	37ab	012d
96:	aca5	4819	9088	0b39	64ba	43a0	b9a8	04b3
112:	88b8	4bf8	3940	d024	0a98	8b0b	1703	0a3a
128:	8820	a381	a21f	3bc0	9298	e893	90bd	042a
144:	0a88	3287	59ab	e980	1211	4002	2208	98b1
160:	7039	0b26	e898	99ab	b118	a1aa	a702	9ac4
176:	9128	ca21	7822	2971	090a	2194	98d0	27bb
192:	0958	8092	993f	b3b0	2922	337a	0f88	8810
208:	8a29	0183	fb15	b888	0d4c			

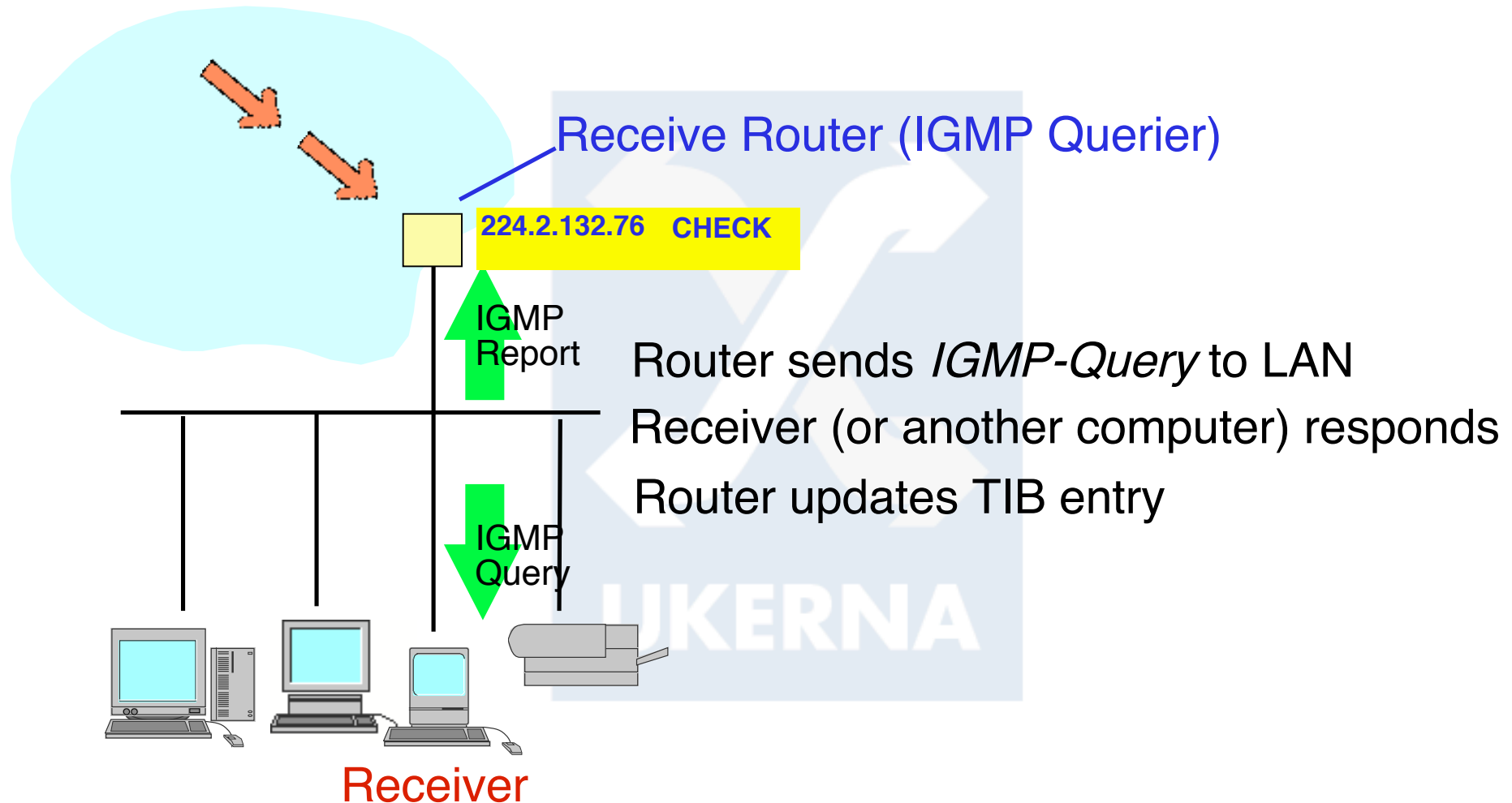
Internet Group Management Protocol – IGMP

- IGMP is a protocol used on the LAN for hosts to tell routers which groups they are interested in
 - May not be necessary to receive from a source on the same LAN
 - However, IGMP is often used by switches to restrict multicast flow on the LAN
- Three versions of IGMP
- IGMPv1 (RFC 1112)
 - Hardly no one uses this anymore, used by e.g. Windows 95
- IGMPv2 (RFC 2236)
 - Maybe the most commonly used version today
- IGMPv3 (RFC 3376)
 - This is the recommended version, backwards compatible with IGMPv2
 - This is needed for SSM (to specify sources to join)
 - Supported by Windows XP, recent Linux and some UNIX systems

IGMP Overview

- Multicast router with the lowest address is elected as *IGMP querier*
- The querier sends periodic queries (default 125s intervals)
- Hosts respond with which groups they want to receive
- Hosts also immediately send a report when host initially joins
 - Do not have to wait for the periodic query
- Hosts immediately send a message when they leave
 - This was not the case for IGMPv1
- State times out if there are no responses to the queries
 - This is important if a host crashes and never says stop

IGMP Query



TCPdump: Receiver Joining

Receiver enables Group (G) at Ethernet interface
Receiver sends *IGMP Membership Report* for (G)

↑
IGMP
Report

```
16:00:35.401923 churchward.erg.abdn.ac.uk > 224.2.132.76: igmp v2 report 224.2.132.76 [ttl 1]
```



Data flows from Source to requested Group

```
16:00:35.594515 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 986 (ttl 127)
16:00:35.700825 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 969 (ttl 127)
16:00:35.706132 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 986 (ttl 127)
....
16:00:44.363832 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 1011 (ttl 127)
16:00:44.369947 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 986 (ttl 127)
```

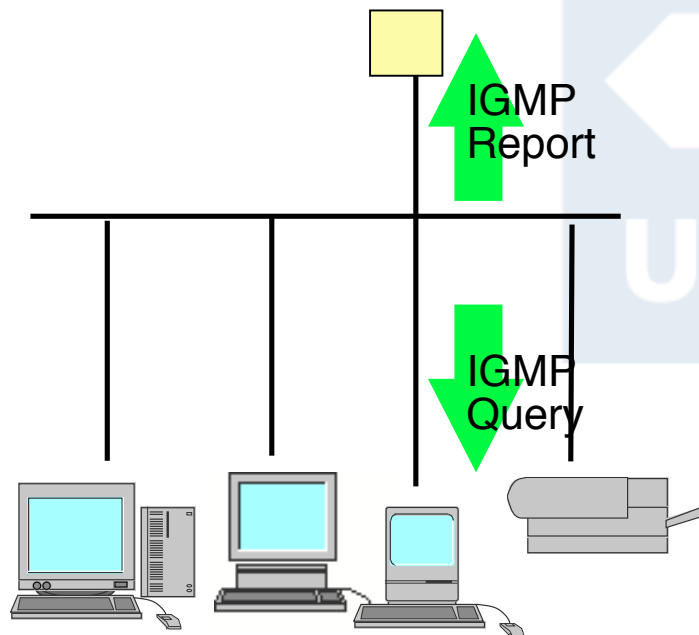
↑
IGMP
Report

Receiver resends *IGMP Membership Report* for (G)- in case first report was lost

```
16:00:44.370493 churchward.erg.abdn.ac.uk > 224.2.132.76: igmp v2 report 224.2.132.76 [ttl 1]
16:00:44.711456 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 998 (ttl 127)
16:00:44.801106 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 970 (ttl 127)
16:00:44.806525 source.x.ac.uk.64276 > 224.2.132.76.49280: udp 981 (ttl 127)
```

Router IGMP Table

```
•gate#show ip igmp group
•IGMP Connected Group Membership
•Group Address      Interface      Uptime      Expires      Last Reporter
•239.255.255.255    FastEthernet0  2d02h       00:02:38     139.133.204.131
•239.255.255.250    FastEthernet0  2d02h       00:02:40     139.133.204.220
•224.2.132.76       FastEthernet0  2d02h       00:02:58     139.133.204.110
•224.255.222.239    FastEthernet0  2d02h       00:02:59     139.133.204.110
•233.2.171.1        FastEthernet0  2d02h       00:02:46     139.133.204.118
•224.2.127.254      FastEthernet0  2d02h       00:02:38     139.133.204.131
•224.0.1.75         FastEthernet0  2d02h       00:02:38     139.133.204.131
•224.0.1.40         FastEthernet0  2d02h       never        139.133.204.210
```



IGMPv2 Details

- **General periodic queries**
 - Sent to 224.0.0.1 (all nodes)
 - All multicast hosts must listen to this group
- **Host reports:** Hosts respond when they receive a General Query
 - Sets a timer for each group with a random value of 0-10s (10s is the default).
 - When the timer expires send a report for each group,
 - Send a report for the group unless another host responded for the same group
 - Hosts immediately send unsolicited reports when joining a new group
 - Reports usually repeated 1-2 times (in case they are lost)
 - All reports are sent to the multicast group itself
- **Host Leave messages:** Immediately send a Leave Message when leaving a group
 - Leave Messages are sent to 224.0.0.2 (all routers)
- **Group-Specific Query:** Sent when a Querier receives a Leave message
 - Sends a group specific query to the group address
 - Checks if there still are listeners, resent 1-2 times every 10s
 - Interested hosts respond, similar to general queries
 - Router assumes no more interest, if no reports are seen

IGMPv3 Details 1/2

- The general periodic queries are sent to 224.0.0.1 (all nodes)
 - All multicast hosts must listen to this group
- Host reports
 - All reports are sent to 224.0.0.22
 - No idle member state, hosts always report all groups they are members of
 - Hosts immediately sends unsolicited reports when joining a new group
 - Reports are usually repeated 1-2 times (in case they are lost)
 - Reports can contain several joins (and leaves) for multiple groups
- No host Leave message
 - Hosts immediately send a report when leaving a group, but not a specific message type
- Group-Specific queries
 - When a host leaves, the querier sends a group specific query to the group address to check if there still are listeners, resent 1-2 times every 10s
 - Interested hosts respond in a similar manner to general queries
 - Router assumes no more interest, if no reports are seen

IGMPv3 Details 2/2

- IGMPv3 allows sources to be specified
- IGMPv3 allows *include mode* and *exclude mode*
 - Can do include mode for some groups while exclude for others
- In *include mode* one specifies which sources to receive
 - This is used for SSM, but one may also specify sources for non-SSM groups
- In *exclude mode* one specifies sources to block
 - Receive from all but the listed set of sources
- IGMPv2 compatibility mode
 - Routers fall back to v2 mode if they see IGMPv2 queries
 - Hosts and routers fall back to v2 mode for a specific group if they see v2 reports
 - All multicast routers must support IGMPv3 for it to be used
 - IGMPv2 and IGMPv3 hosts can co-exist

Configuring IGMP on Cisco IOS

For IOS to do IGMP we need to enable multicast routing and enable PIM on the interfaces

```
ip multicast-routing
!  
interface ...  
  ip pim sparse-mode
```

IOS uses IGMPv2 by default, to use IGMPv3:

```
interface ...  
  ip igmp version 3
```

Checking IGMP state on Cisco IOS 1/2

To see group memberships we can do:

```
cisco> show ip igmp groups
IGMP Connected Group Membership
```

Group Address	Interface	Uptime	Expires	Last Reporter
239.255.255.255	Vlan26	7w0d	00:02:32	158.38.63.1
239.255.255.255	Vlan2	7w0d	00:02:29	128.39.47.90
239.255.255.253	Vlan20	4d10h	00:02:33	158.38.60.95
239.255.255.253	Vlan10	1w4d	00:02:33	158.38.60.11
239.255.255.253	Vlan95	7w0d	00:02:12	158.38.152.196
239.255.255.250	Vlan80	1d00h	00:02:56	158.38.61.149
239.255.255.250	Vlan10	5d12h	00:02:26	158.38.60.44
239.255.255.250	Vlan20	7w0d	00:02:35	158.38.60.95
224.2.127.254	Vlan26	7w0d	00:02:35	158.38.63.1
224.2.127.254	Vlan2	7w0d	00:02:32	128.39.47.90
239.255.67.250	Vlan20	4d10h	00:02:38	158.38.62.97
232.26.17.81	Vlan26	4d07h	stopped	158.38.63.22

UKERNA

Checking IGMP state on Cisco IOS 2/2

To show sources and not just groups for IGMPv3 we can do

```
cisco> show ip igmp groups detail
Flags: L - Local, U - User, SG - Static Group, VG - Virtual Group,
      SS - Static Source, VS - Virtual Source

Interface:      Vlan26
Group:          239.255.255.255
Flags:          L U
Uptime:         7w0d
Group mode:     EXCLUDE (Expires: 00:02:24)
Last reporter:  158.38.63.1
Source list is empty

Interface:      Vlan26
Group:          232.26.17.81
Flags:          SSM
Uptime:         4d07h
Group mode:     INCLUDE
Last reporter:  158.38.63.22
Group source list: (C - Cisco Src Report, U - URD, R - Remote, S - Static,
                  V - Virtual, Ac - Accounted towards access control limit,
                  M - SSM Mapping)

Source Address  Uptime    v3 Exp   CSR Exp   Fwd  Flags
129.177.30.248  4d07h     00:02:32 stopped   Yes   R
129.242.2.140   4d07h     00:02:32 stopped   Yes   R
152.94.26.6     4d07h     00:02:32 stopped   Yes   R
158.36.22.14    2d07h     00:02:32 stopped   Yes   R
```

Multicast and Ethernet Switches

- Many switches have features that can restrict multicast flow to only ports where there are members
 - This is usually good, but sometimes switches misbehave
- There are generally three possible methods
 - GARP/GMRP – hosts use I2 protocol to tell switches
 - Supported by very few systems
 - CGMP – routers tell switches what to do
 - A Cisco proprietary protocol and Cisco are dropping support
 - IGMPv3 “leaves” not supported
 - IGMP snooping/proxy
 - The switches snoop the IGMP messages going between hosts and routers

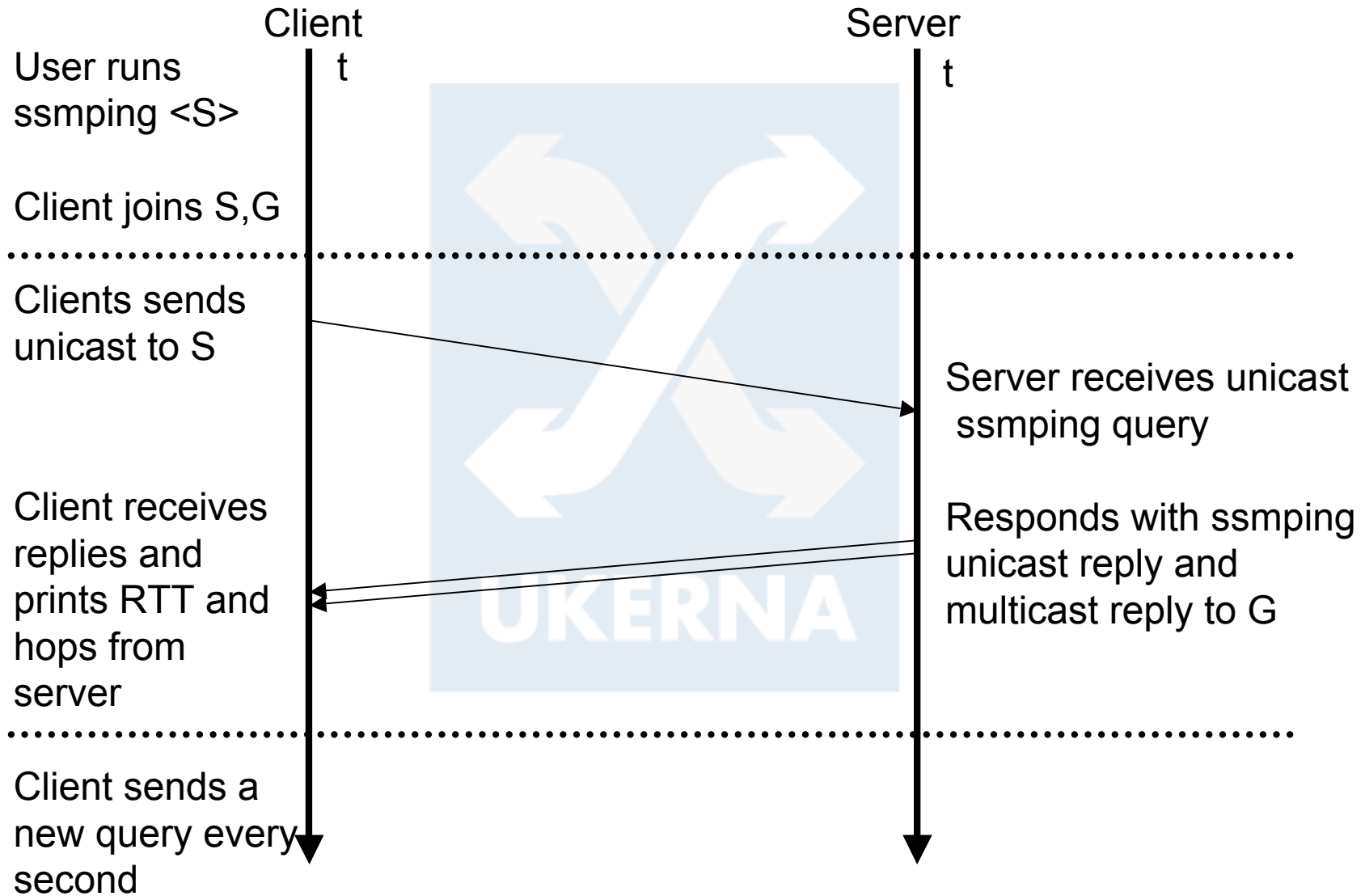
IGMP Snooping

- Switches must know where the router is
 - Switches must always forward multicast on router port (may snoop for IGMP queries)
- IGMPv1/v2 reports from one host must not be seen by other hosts
 - Due to the report suppression (idle member state)
 - Switches must send reports upstream (router port)
 - Some forward one of the host reports, others fake a report using a host's IP address, switch's management address or 0.0.0.0 as source address
 - Similarly switches may spoof group specific queries when they receive leave messages from hosts
- IGMPv1/v2 reports are sent to the group
 - Switches may have difficulties distinguishing data and reports (look for IP Router Alert option)
- Switches may only use MAC addresses to distinguish packets
 - So data may still go to more ports than necessary if different groups have same MAC address. Some layer-3 switches check IP group addresses

ssmping

- A tool for testing multicast connectivity
- Behavior is a bit like normal ping
- A server must run ssmpingd
- A client can ping a server by sending unicast ssmping query
- Server replies with both unicast and multicast ssmping replies
- In this way a client can check that it receives SSM from the server
 - also parameters like delay, number of router hops etc.
- JANET is running a server at ssmping.beacon.ja.net
- There is a similar tool called asmping for checking ASM connectivity
- See <http://www.venaas.no/multicast/ssmping/>

How ssmpling works



Example IPv4 ssm ping output (v6 supported)

```
$ ssm ping -4 -c 5 ssm ping.beacon.ja.net
ssm ping joined (S,G) = (193.60.199.162,232.43.211.234)
pinging S from 158.38.63.22
  unicast from 193.60.199.162, seq=1 dist=16 time=39.331 ms
  unicast from 193.60.199.162, seq=2 dist=16 time=39.394 ms
multicast from 193.60.199.162, seq=2 dist=16 time=43.905 ms
  unicast from 193.60.199.162, seq=3 dist=16 time=39.542 ms
multicast from 193.60.199.162, seq=3 dist=16 time=39.547 ms
  unicast from 193.60.199.162, seq=4 dist=16 time=39.137 ms
multicast from 193.60.199.162, seq=4 dist=16 time=39.142 ms
  unicast from 193.60.199.162, seq=5 dist=16 time=39.535 ms
multicast from 193.60.199.162, seq=5 dist=16 time=39.539 ms

--- 193.60.199.162 ssm ping statistics ---
5 packets transmitted, time 5000 ms
unicast:
  5 packets received, 0% packet loss
  rtt min/avg/max/std-dev = 39.137/39.387/39.542/0.292 ms
multicast:
  4 packets received, 0% packet loss since first mc packet (seq 2) recvd
  rtt min/avg/max/std-dev = 39.142/40.533/43.905/1.958 ms
$
```

What does ssmping output tell us?

- 16 unicast hops from source, also 16 for multicast, might indicate that unicast and multicast follow the same path
- Multicast RTTs are about same for unicast and multicast
 - However, the delay for the first multicast packet is large, would need to send more queries for a proper test
 - Note the difference in unicast and multicast RTT shows one way difference for unicast and multicast replies, since they are replies to the same request packet
- Multicast tree not ready for first multicast reply, ok for 2nd so the tree was in place after about one second when the second packet was sent
- No unicast loss, no multicast loss after tree established



Session 2

Intradomain multicast routing

Intradomain Multicast Routing

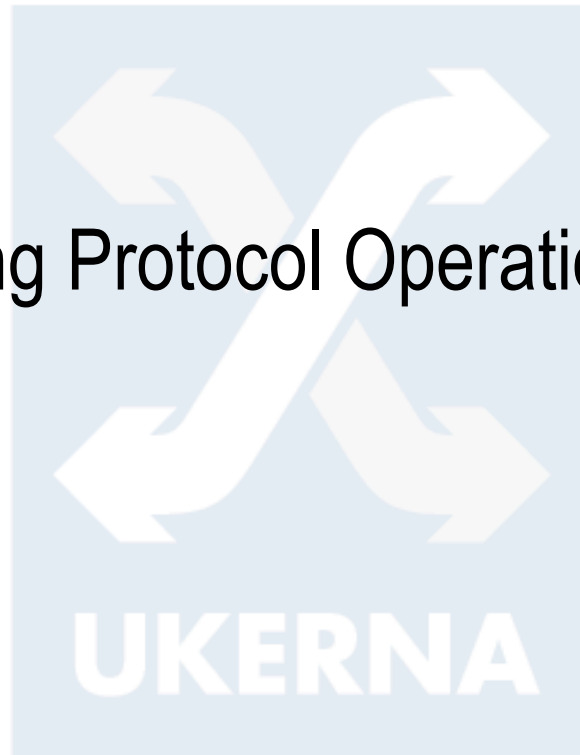
- We will look at how multicast can be deployed in a site network
 - Or any other reasonably small network within one management domain
- We want the traditional service model where anyone in the network can join without knowing where sources are, and sources can send without knowing the receivers
- Multicast routing is all about efficiently creating multicast distribution trees, where each tree is rooted at the source and spreads out to the receivers
- May be difficult since neither sources nor receivers know where everyone else is
- Some multicast routing protocols are based on “flood and prune”
 - Data is initially flooded everywhere
 - Trees are then pruned to be only where needed
 - However does not scale well, unless receivers are almost everywhere
- The most common routing protocols today are PIM-DM and PIM-SM
 - PIM is Protocol Independent Multicast in that it can make use of the unicast routing protocol, and is independent of which Internet Protocol is used

PIM-SM/PIM-DM

- PIM-DM (PIM Dense Mode, RFC 3973)
 - A flood and prune protocol
 - May be okay for a dense population of receivers
- PIM-SM (PIM Sparse Mode, RFC 2362)
 - Does not flood
 - Works with a sparse population of receivers, scales much better
- PIM-SM is by far the most commonly used protocol today
- PIM-SM makes use of a so-called Rendezvous Point where sources and receivers meet
- All routers in the network agree where the RP is for a group
 - Hosts and receivers do not need to know where the others are.
 - Trees, at least initially, pass through the RP

PIM-SM

- Multicast Routing Protocol Operation



Configuring PIM-SM 1/2

- Routers must have:
 - Multicast routing enabled
 - PIM on interfaces where they face one another
 - IGMP on host interfaces
 - Note that on some routers, incl IOS, you need to enable PIM on host interfaces to get IGMP
- Also essential to configure RPs
 - All routers in the domain must agree which RP to use for a group. May have just one RP for all, or different RPs for different ranges
 - RP routers must be configured to be RPs
 - Other routers must know the addresses of the RPs for the different group ranges

Configuring PIM-SM 2/2

- One common way, and perhaps the best, is to just statically configure the RP address(es) on each router
 - RP addresses can be configured as additional loopback interfaces on routers and announced as host routes into the routing table
 - You can then move the RP without configuring all the routers
 - Anycast-RP (RFC 3446) allows failover between multiple RPs
 - Today done with MSDP (see later)
- Another option is BSR (bootstrap router protocol)
 - Allows routers to dynamically learn which RPs to use
 - One router is elected as so called BSR
 - Potential RPs announce themselves and the BSR forms announces an RP-set throughout the domain
 - Allows failover as well, but not as well as the anycast with MSDP technique
- A Cisco specific protocol similar to BSR is Auto-RP
 - Some advantages over BSR, but suggest using BSR in most cases
 - Auto-RP makes use of PIM-DM for its operation, PIM-SM used for other groups

Configuring RPs on Cisco IOS

Recommend static RP config, on RPs and other routers you may simply do e.g.

```
ip pim rp-address 192.0.2.1
```

This would specify 192.0.2.1 to be the RP for all multicast groups. In some cases you may prefer to have your own RP for groups used internally, while using your provider's RP for global groups. Below is an example of that

```
ip pim rp-address 192.0.2.1 21
ip pim rp-address 192.0.2.129 20
!
access-list 20 permit 239.255.0.0 0.0.255.255
access-list 20 permit 229.55.150.208
access-list 20 permit 224.0.1.0 0.0.0.255
access-list 20 deny any
!
access-list 21 deny 239.255.0.0 0.0.255.255
access-list 21 deny 229.55.150.208
access-list 21 deny 224.0.1.0 0.0.0.255
access-list 21 permit any
```

SSM – Source Specific Multicast

- SSM is a new multicast service model
 - Receivers specify the source address(es) in addition to the group
- This avoids rogue sources sending to the group
 - Imagine watching one video stream and someone else also sends their video, or just random data at high rates
- The main benefit is hugely simplified routing
- PIM-SM works very well with SSM
 - Last-hop routers know the sources, can immediately join Shortest-Path Trees.
 - No need for an RP and a shared tree to do source discovery
 - Multicast much easier to deploy and manage
- SSM requires source discovery at the application layer
 - Very easy for streaming with one fixed source (maybe most important use of multicast)
 - May be difficult for some multi-party or discovery applications

Configuring SSM on Cisco IOS

- For SSM to work you need to enable IGMPv3 at the edges (see earlier slide)
- Must say the standard (default) SSM range 232/8 to be treated as SSM groups
 - Do not allow (*,G)-join for SSM groups
 - Never send PIM registers for 232/8
 - RPs ignore (*,G)-joins and registers for such groups

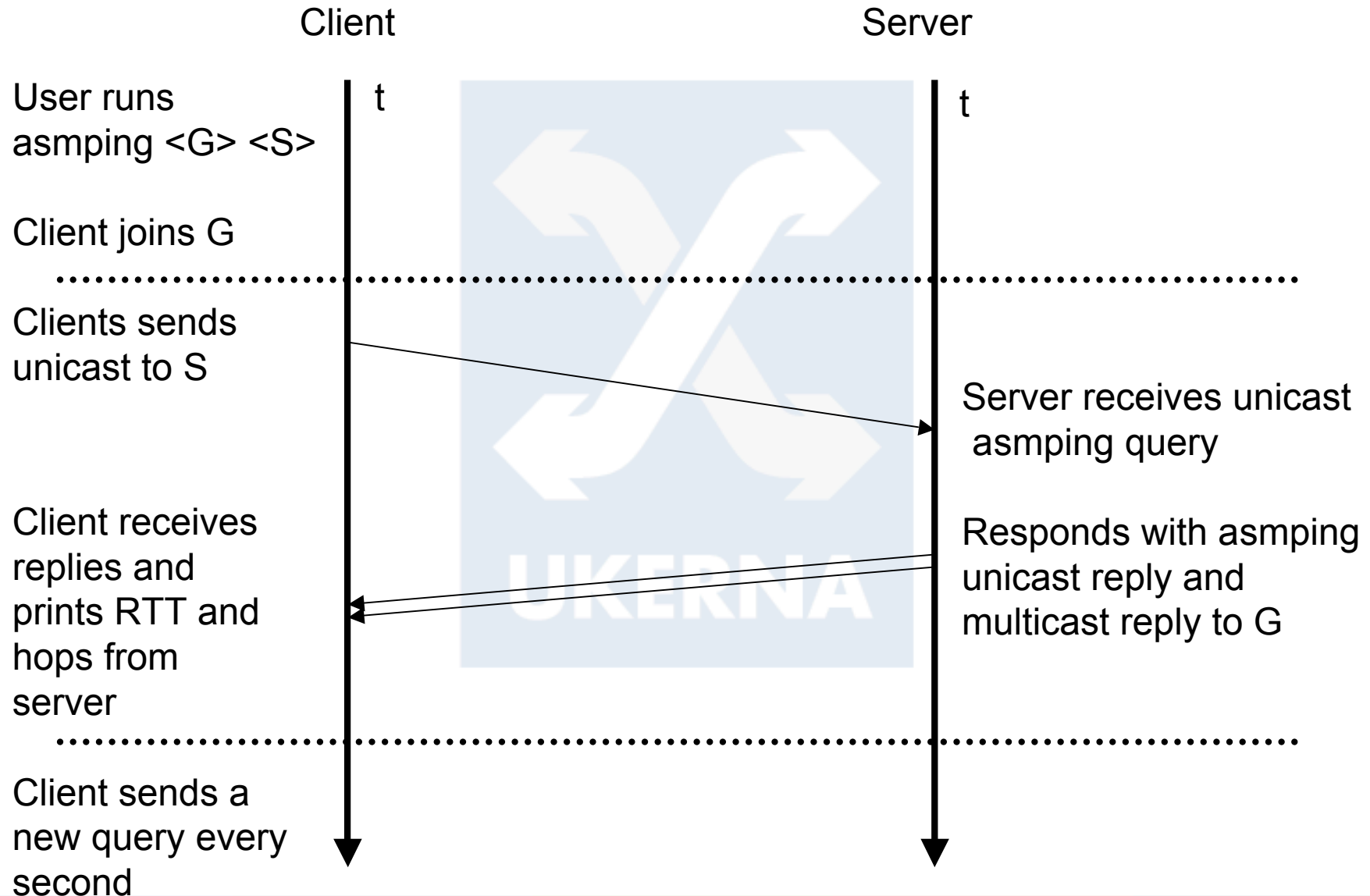
```
ip pim ssm default
```

UKERNA

asmping

- A tool for testing multicast connectivity
- Behavior is a bit like normal ping
- A server must run ssm pingd (latest version supports asmping)
- asmping is ASM version of ssm ping
- A client can ping a server by sending unicast asmping query
- Server replies with both unicast and multicast asmping replies
- In this way a client can check that it receives ASM from the server
 - And also parameters like delay, number of router hops etc
- JANET is running a server at ssmping.beacon.ja.net
- With asmping one must specify both group and server to ping as arguments
 - Note that the last byte of the IPv4 group address is always set to 234 to prevent servers being made to send to arbitrary groups

How it Works



Example Output

```
$ asmping -c 5 224.3.4.234 ssmping.beacon.ja.net
ssmping joined (S,G) = (193.60.199.162,224.3.4.234)
pinging S from 158.38.63.22
  unicast from 193.60.199.162, seq=1 dist=16 time=39.275 ms
multicast from 193.60.199.162, seq=1 dist=4 time=238.958 ms
  unicast from 193.60.199.162, seq=2 dist=16 time=39.502 ms
multicast from 193.60.199.162, seq=2 dist=16 time=42.980 ms
  unicast from 193.60.199.162, seq=3 dist=16 time=39.553 ms
multicast from 193.60.199.162, seq=3 dist=16 time=39.558 ms
  unicast from 193.60.199.162, seq=4 dist=16 time=39.591 ms
multicast from 193.60.199.162, seq=4 dist=16 time=39.595 ms
  unicast from 193.60.199.162, seq=5 dist=16 time=39.205 ms
multicast from 193.60.199.162, seq=5 dist=16 time=39.210 ms

--- 193.60.199.162 ssmping statistics ---
5 packets transmitted, time 5000 ms
unicast:
  5 packets received, 0% packet loss
  rtt min/avg/max/std-dev = 39.205/39.425/39.591/0.199 ms
multicast:
  5 packets received, 0% packet loss since first mc packet (seq 1) recvd
  rtt min/avg/max/std-dev = 39.210/80.060/238.958/79.460 ms
$
```

What Does the Output tell us?

- 16 unicast hops from source
- For multicast we first have 4, later stays at 16
- We also get some info about loss and RTTs
- Number of hops is perhaps the most interesting
- Initially only 4 hops, probably PIM registers involved (and MSDP)
- In the stable situation, with 16, there is probably native forwarding all the way
- Forwarding is probably on shortest path tree from source to receiver
- In theory may also detect switch from RPT to SPT if number of hops differ
 - Number of hops on RPT are then probably larger than number of unicast hops
- We got the first multicast packet (not with SSM), but with a very long delay
 - This is probably due to MSDP caching
 - MSDP was involved because we pinged from another domain
 - More about MSDP later



Session 3

Interdomain multicast routing

Interdomain Multicast Routing

- We will look at how multicast domains can be interconnected to get multicast connectivity throughout the Internet
- Many networks don't support multicast, so multicast often needs to be routed differently from unicast
 - BGP is used for setting up peerings to route unicast between networks
 - If you use BGP, you may need MBGP (Multiprotocol BGP)
 - If you use only static unicast routes, you are fine with just static multicast routes
- Each domain typically uses their own RP for all groups
 - Each RP only knows about sources in its own domain
 - To get connectivity between domains, MSDP is used so that RPs learn about sources in other domains
 - This is not needed for SSM to work between domains (no RPs)
- You may wish to configure boundaries to restrict some multicast groups to be local
 - Filtering what can flow across the domain's boundary

Multiprotocol BGP (MBGP)

- Originally BGP supported just IPv4 unicast
- Multiprotocol BGP (RFC 2858) may have AFI (Address Family Identifier) IPv4/IPv6 and for those, SAFI (Subsequent AFI)
 - SAFI = 1 for unicast, 2 for multicast, 3 for both
 - Seems IETF wants to deprecate 3
- For a peering one may configure IPv4/IPv6 unicast/multicast separately with different policies
- Multicast routes used for RPF
 - Sometimes in addition to unicast routes
 - One may also sometimes translate unicast routes into multicast
- Recommend only using multiprotocol BGP
 - Whenever unicast BGP between multicast networks exists; enable multicast peering if multicast connectivity is desired
 - Try to avoid tricks like translation or unicast routes for RPF
- Note that you only need to worry about this if you have BGP peerings to networks you want to have multicast connectivity with
- We sometimes talk about multicast BGP, meaning multiprotocol BGP for exchanging multicast routes

Pv4 Multiprotocol BGP on IOS

```
router bgp 224
  no bgp default ipv4-unicast
  neighbor 192.0.2.1 remote-as 64001
  neighbor 192.0.2.129 remote-as 64002
  !
  address-family ipv4
  ...
  !
  address-family ipv4 multicast
  neighbor 192.0.2.1 activate
  neighbor 192.0.2.65 activate
  network 192.0.2.128 mask 255.255.255.128
  exit-address-family
```

- Basically you configure multiprotocol BGP like you do for unicast
- If the policies are the same you can, for each multicast peer, copy from that peer's unicast config

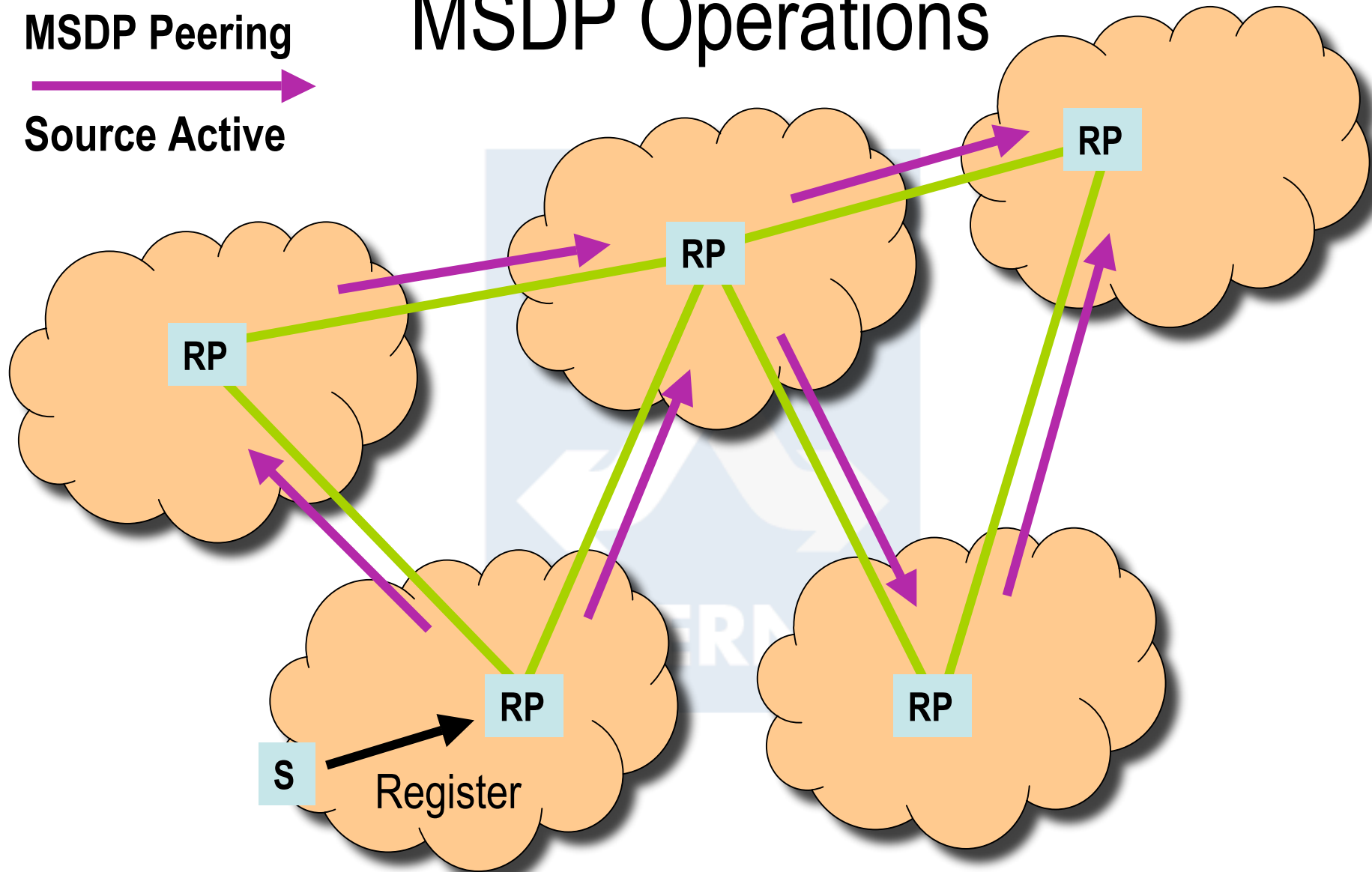
MSDP – Multicast Source Discovery Protocol

- MSDP (RFC 3618) sets up a mesh of MSDP peerings between domains
 - Usually between MBGP pairs using the same addresses for the end-points
 - Internal MSDP peerings can be set up if the RP is not where the external peerings are terminated
 - Internal peerings also allow multiple RPs for the same group
 - Used for anycast-RP
- When an RP learns a new local source:
 - Sends a Source Active (SA) advertisement
 - SAs are flooded through the mesh of MSDP peerings to other RPs
 - SAs may optionally contain data packet
 - SAs are cached
 - RPF is used to prevent loops
- If an RP knows of sources within other domains, it can construct Shortest Path Trees (SPTs) to them when a local host joins the group

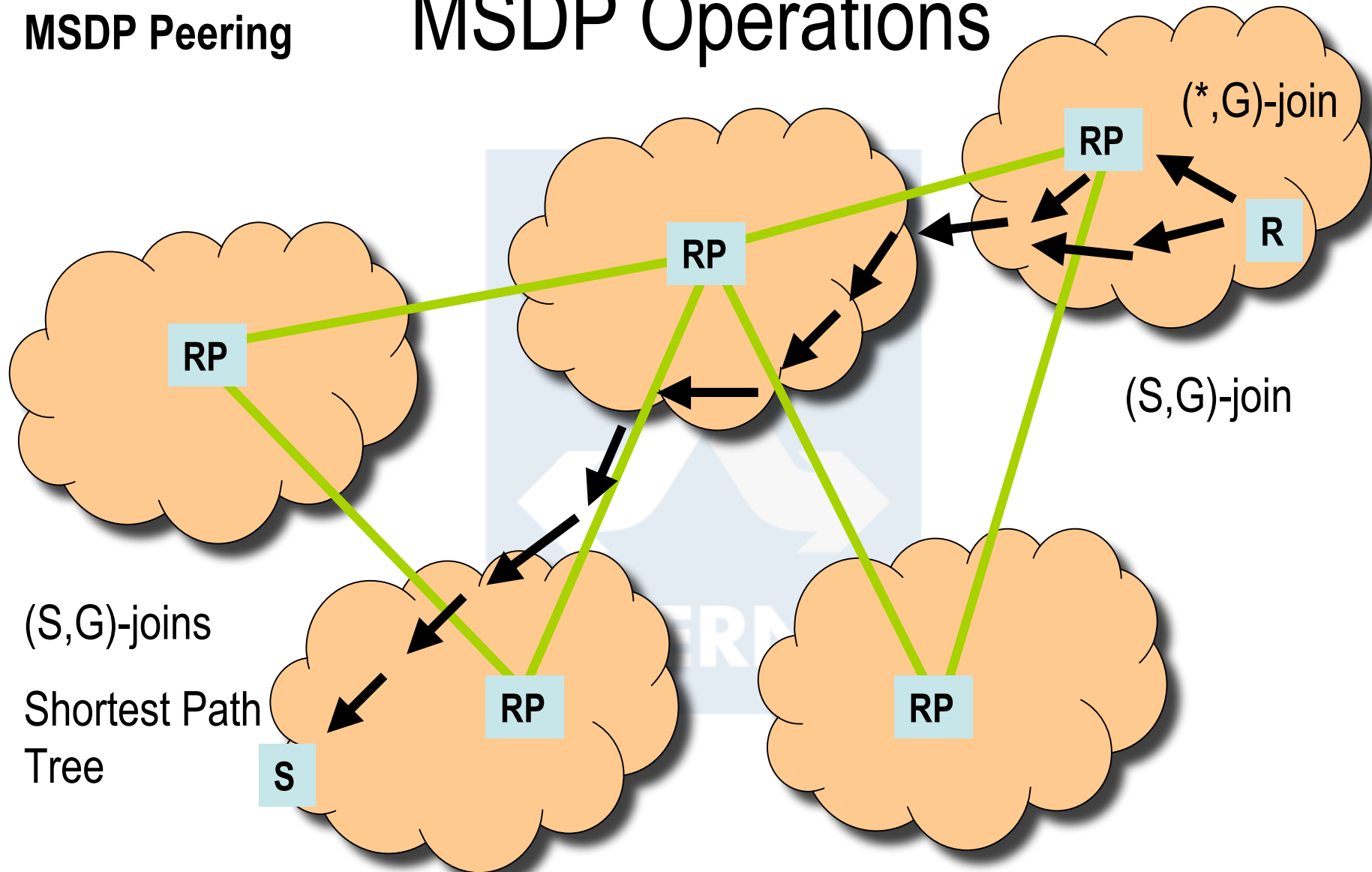
MSDP Peering

Source Active

MSDP Operations



MSDP Operations



MSDP configuration on IOS

```
ip msdp peer 192.0.2.1 connect-source Loopback0
ip msdp sa-filter in 192.0.2.1 list 21
ip msdp sa-filter out 192.0.2.1 list 21
ip msdp peer 192.0.2.65 connect-source Loopback0
ip msdp cache-sa-state
ip msdp originator-id Loopback0
ip msdp mesh-group example 192.0.2.65
ip msdp mesh-group example 192.0.2.66
!
access-list 21 deny 239.255.0.0 0.0.255.255
access-list 21 deny 229.55.150.208
access-list 21 deny 224.0.1.0 0.0.0.255
access-list 21 permit any
```

- Here we have one external MSDP peer
- We also do anycast RP between 3 RPs, consisting of this router plus 2 others in mesh-group example

dbeacon

- dbeacon is a new multicast beacon
 - <http://dbeacon.innerghost.net/>
 - Alternative to NLNR beacon
 - IPv4 and IPv6, ASM and SSM
 - Written in C, light and easy to install
 - No central server, ASM used for signalling
 - Any beacon client can be configured to provide a matrix
 - Beacon options, e.g.
 - `dbeacon -b ff7e:a30:2001:db8:10::beac -S -a admin@email`
 - With Apache:
 - `ScriptAlias /matrix/ /usr/share/dbeacon-matrix/matrix.pl`
 - Edit script for path to matrix `$dumpfile = '/var/lib/dbeacon/dump.xml';`

Example dbeacon Matrix

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	S13	S14	S15	S16	S17	S18	S19	S20	S21	S23	S24	S25
HS-MRD6 R1		12	12	9	6	13	4	7	8	8	7	5	8	7	5	6	5	5	9	7	7			
New York University R2	12			4		10	12	12		13	12	10	13	12	10	11	10		10	12		12		10
CNR Pisa R3	12	14		10			12	12		13	12	10	13	12	11	11	10	10	10	12		12	10	10
6pack.org R5	6	10		7		11	9	9		10	9	4	10	9	7	8	7	9	5	9		9	9	7
Internet2-Ann Arbor R6	13	10		7			13	13		14	13	11	14	13	11	12	11	11	11	13		13	11	11
ITIN-IABG R7	4	12	12	9	6	13		7	8	8	7	5	8	7	5	6	5	5	9	7	7			
RENATER R8	7	12		5			7			4	3	5	4	3	5	6	5	5	5	3		3	5	5
cemp1.switch.ch R9	8	10	9	7		11	8	8		9	8	6	9	8	9	7	6	8	8	8	9	8	8	6
Phocean R10	8	13		6			8	4			4	6	5	4	6	7	6	6	6	4		4	6	6
canet.u-strasbg.fr R11	7	12	12	9	9	13	7	3	8	4		12	4	3	10	6	5	9	9	3	10	3	9	5
ssmping.uninett.no R12	5	13		10	4	14	5	12		13	12		13	12	10	11	10	12	10	12		12	12	10
Universite-Paris13 R13	8	13		6			8	4		5	4	13		4	6	7	6	10	6	4		4	10	6
ITIN-Renater R14	7	12	12	9	9	13	7	3	8	4	3	12	4		10	6	5	9	9	3	10	3	9	5
CESGA R15	10	13	11	10	10	14	10	10	9	11	10	13	11	10		9	8	6	10	10	11	10	6	8
CESNET2 R17		11	12	8		12			8									11					11	
UC3M R18	9	12		9	9	13	9	9	8	10	9	7	10	9	6	8	7		9	9		9	3	7
UofA-ERG R19	9	10	12	7	5	11	9	9	8	10	9	10	10	9	7	8	7	9		9	6	9	9	7
IUT_Colmar R20	7	12	12	9	9	13	7	3	8	4	3	12	4	3	5	6	5	9	9		10	3	9	5
ECS Southampton R21	7	11	13	8	2	12	7	10	9	11	10	5	11	10	8	9	8	10	6	10		10	10	8
hadron.switch.ch R22	9	11	10	8	9	12	9	9	2	10	9	12	10	9	10	8	7	9	9	9	10	9	9	7

Multicast Applications

- Mbone tools, vic/rat etc
 - IPv6 multicast conferencing applications
 - <http://www-mice.cs.ucl.ac.uk/multimedia/software/>
- AccessGrid
 - Uses vic and rat for high quality room-based conferencing
 - <http://www.accessgrid.org>
- VideoLAN (vlc)
 - Video streaming, also IPv6 multicast. Server and client
 - <http://www.videolan.org/>
- DVTS
 - Streaming DV over RTP over IPv4/IPv6
 - DV devices using Firewire can be connected to two different machines and you can stream video between them over the Internet
 - <http://www.sfc.wide.ad.jp/DVTS/>
- Mad flute
 - Streaming of files using multicast (IPv4/IPv6 ASM/SSM)
 - <http://www.atm.tut.fi/mad/>



Workshop 34 - Multicast Workshop
Status of Multicast on JANET

Duncan Rogerson
<d.rogerson@ukerna.ac.uk>

April, 2006



Status of IP Multicast on JANET

- No changes
- Architecture based on
 - PIM Sparse Mode
 - MSDP and MBGP in the backbone
 - MSDP between backbone and Regional Networks for distribution of active source information
 - MBGP between backbone and Regional Networks for populating the multicast forwarding table
 - Direct backbone connects – PIM SP
 - Regional Network to site – models may vary

Multicast Monitoring

- Still problematic
 - Esp troubleshooting tools
 - Incidents still usually not reported until after the event
- New version of multicast beacon deployed
 - Retains global, Regional and AccessGrid views

SuperJANET 4 to SuperJANET 5

- SuperJANET 5 Backbone Routers (SBR)
 - Will not be configured as RPs
 - Router hardware specification does not support RP functionality
- SuperJANET 4 Backbone Access Routers (BAR)
 - Were configured as RPs
 - Use as RPs was never recommended or supported

SuperJANET 4 to SuperJANET 5

- Dual links to Regional Networks
 - More complex routing configuration
 - However, no more so than tuning the unicast routing config for dual links
- IPv6 multicast?
 - Most likely implemented at start
 - Mature JunOS and IOS code available



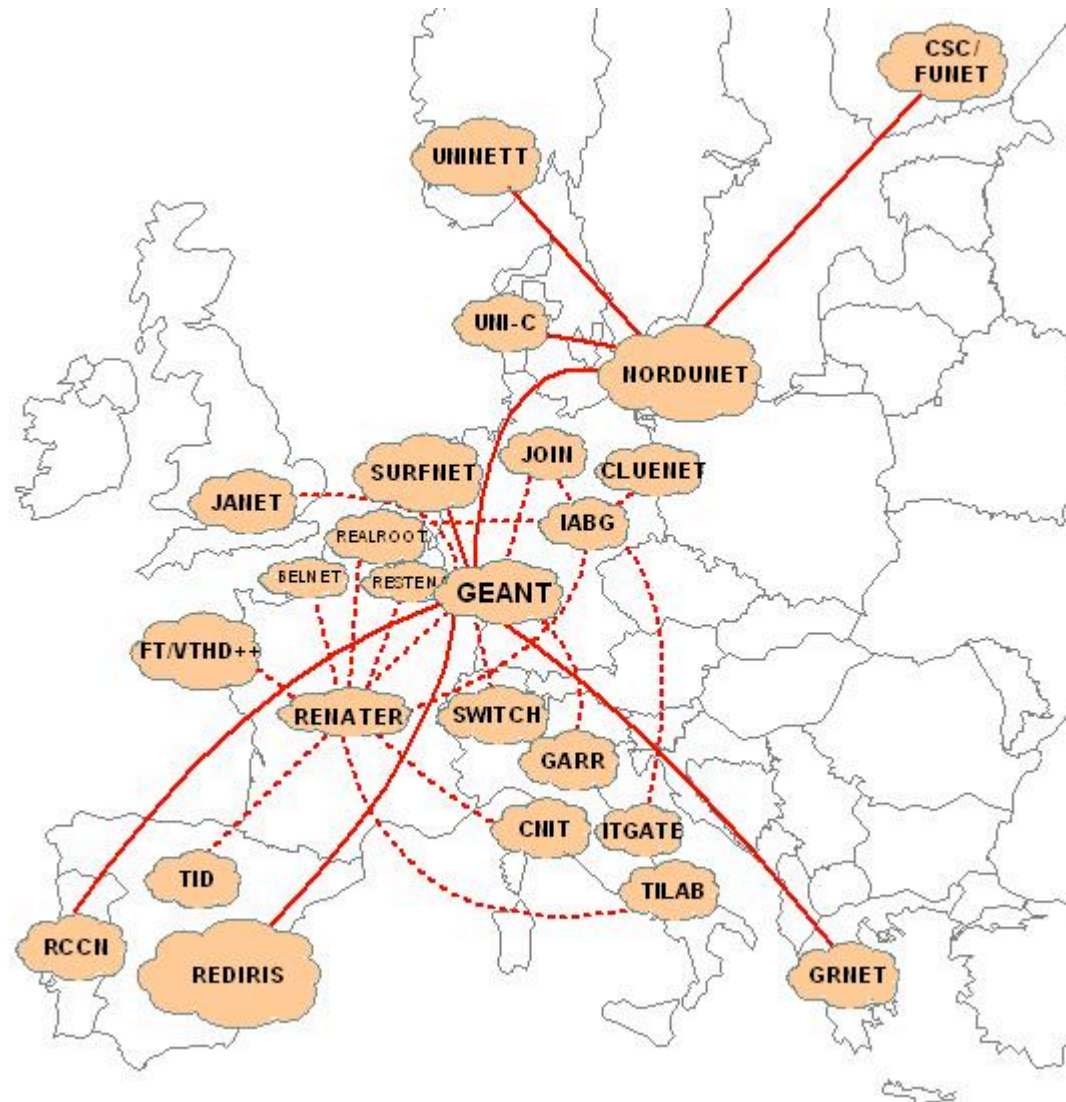
Session 4

IPv6 multicast

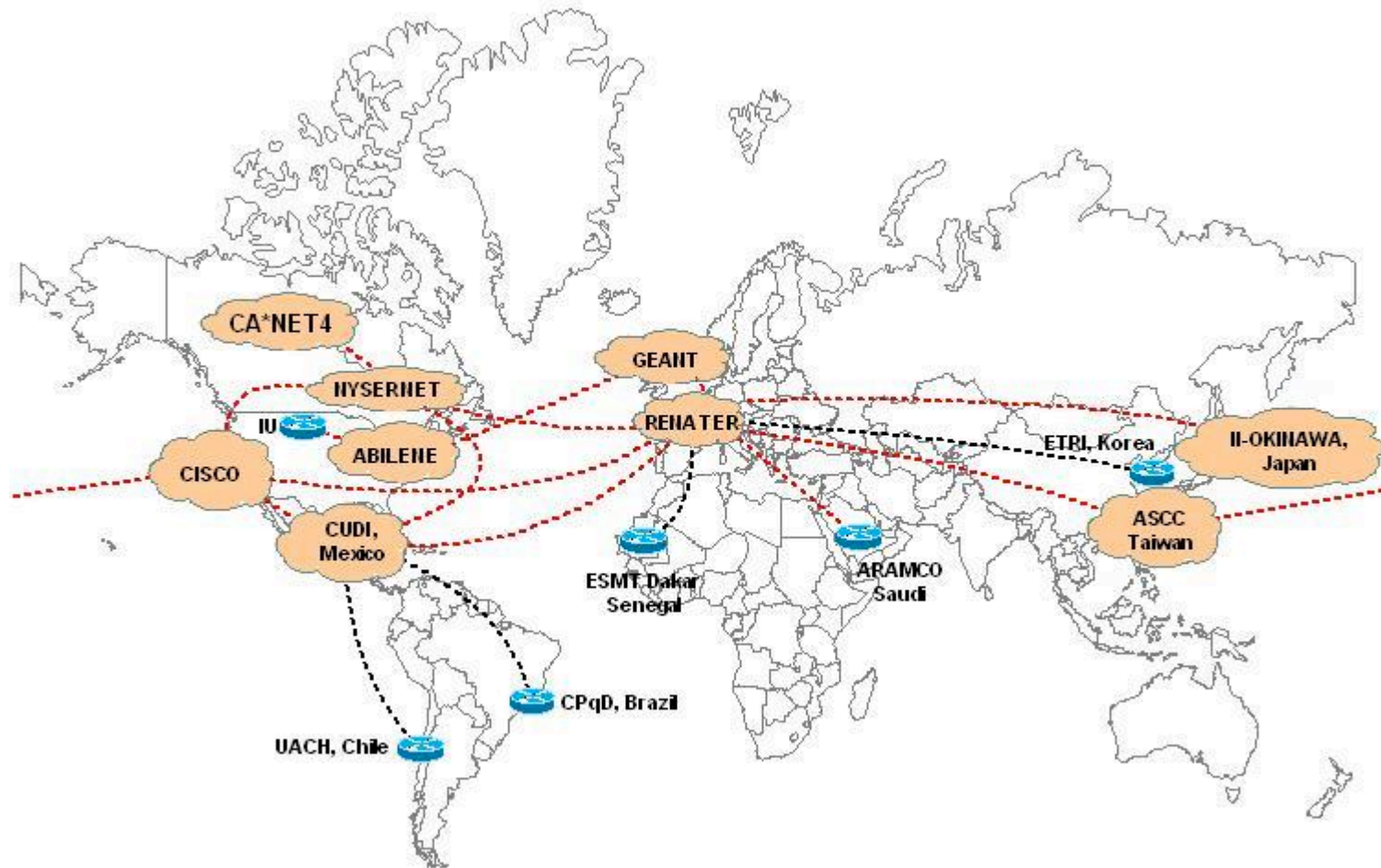
Current IPv6 Multicast Deployment

- A tunneled multicast network called M6Bone was established in France in 2001
- Today more than 50 sites from four continents are connected
- Networks like Abilene, GEANT and NORDUnet have native IPv6 multicast in their production networks and native peerings
- JANET may provide IPv6 tunnels for experimenting with IPv6 multicast
- The M6Bone has both a web site and a mailing list for information and discussion on IPv6 multicast
 - Web site at <http://www.m6bone.net/>
 - Mailing list m6bone@ml.renater.fr
- The largest router vendors support IPv6 multicast on most routers. Many other vendors have some support, or plan to offer support soon

M6Bone – Europe



M6Bone – The World



IPv6 Multicast Addresses

11111111	Flags	Scope	Group ID
8 bits	4 bits	4 bits	112 bits

Scopes:

- 0 Reserved
- 1 Interface local
- 2 Link local
- 3 Reserved
- 4 Admin local
- 5 Site local
- 6-7 Unassigned (admin)
- 8 Organization local
- 9-D Unassigned (admin)
- E Global
- F Reserved

IPv6 multicast addresses are in the range FF00::/8

Flag bits: 0 R P T

T = 0 – permanent IANA assigned address

T = 1 – transient address (not IANA)

P = 1 – derived from unicast prefix (RFC 3306)

R = 1 – embedded-RP address (RFC 3956)

Unicast Prefix-Based Addresses

11111111	Flags	Scope	Resrvd	Plen	Prefix	Group ID
8 bits	4 bits	4 bits	8 bits	8 bits	64 bits	32 bits

- Unicast prefix based addresses are defined in RFC 3306
- **Flags** are set to **3** (P=T=1)
- Reserved must be 0
- **Plen** is the length of the prefix used
- **Prefix** is a unicast prefix
- **Group ID** consists of any 32 bits
- In this way every site will have many multicast addresses unique to them
- One can also use /64 prefixes to get addresses unique to a link
- Example:
 - Site with unicast prefix **2001:db8:1::/48**
 - Example multicast address **ff3e:30:2001:db8:1:0:dead:beef**

Embedded-RP Addresses

11111111	Flgs	Scp	Res	RIID	Plen	Prefix	Group ID
8 bits	4	4	4	4	8	64 bits	32 bits

- Embedded-RP addresses are defined in RFC 3956
- Variation of the unicast prefix based addresses
- **Flags** are set to 7 (R=P=T=1). R is a new flag for embedded-RP
- Reserved must be 0
- **RIID** contains the last 4 bits of the RP address
- **Plen** is the length of the prefix used
- **Prefix** is a unicast prefix
- **Group ID** consists of any 32 bits
- One may use /64 prefixes as well to have an RP per link
- Example:
 - Site with unicast prefix 2001:db8:1::/48 can pick RP address 2001:db8:1::a
 - Example multicast address ff3e:0a30:2001:db8:1:0:dead:beef

SSM Addresses

11111111	Flags	Scope	Resrvd	Plen	Prefix	Group ID
----------	-------	-------	--------	------	--------	----------

8 bits 4 bits 4 bits 8 bits 8 bits 64 bits 32 bits

- SSM addresses are a special subset of unicast prefix based addresses
- **Flags** are set to 3 (P=T=1)
- Reserved must be 0
- **Plen** is set to 0
- **Prefix** is set to all zeroes
- **Group ID** consists of any 32 bits
- This results in ff3x::/96 for SSM (where x is any scope)

11111111	3	Scope	0	0	0	Group ID
----------	---	-------	---	---	---	----------

8 bits 4 bits 4 bits 8 bits 8 bits 64 bits 32 bits

MLD – Multicast Listener Discovery

- MLD is used between hosts and routers to signal which groups (and sources) a host is interested in
- Similar to IGMP for IPv4, but uses ICMP
- MLDv1 – RFC 2710 supports only ASM, similar to IGMPv2
- MLDv2 – RFC 3810 also supports SSM, similar to IGMPv3
- Querier periodically (125s default) sends queries to **ff02::1** to see whether there still is interest for current groups
- MLD is used for all groups of scope 2 or larger
 - Except for ff02::1
- For robustness messages are retransmitted
- Both the robustness and the timers can be tuned by querier

Multicast and Ethernet Switches

- Some switches don't correctly forward multicast
 - One witnessed problem is when hosts receive multicast for a couple of minutes, and then it stops
 - The likely reason for this is that the initial unsolicited report from the host reaches the router, but the host is not receiving MLD queries from the router. Try upgrading your switch software if possible
- To restrict the flow of multicast, one may want to use MLD snooping
 - Switches supporting MLDv2 are now available
 - Not aware of any MLDv1 snooping switches
- Due to snooping switches, hosts should also use MLD for link-local groups like the solicited node address
 - Only exception is the all-nodes address ff02::1
 - One drawback is that all IPv6 hosts join solicited node address (usually unique), so the switch will get a lot of state. Some switches may choose to flood these groups or possibly all link-local multicast

PIM-SM – RP-set Configuration

- Routers can learn their RP-set via static configuration or BSR
- Routers may also learn RPs from embedded-RP
 - Router learns the RP when it sees a multicast packet or a join/prune with embedded-RP address
 - Embedded-RP results in the RP-set being highly dynamic
 - IOS has started to use one fixed tunnel interface for embedded-RP to reduce amount of tunnel interfaces going up and down

RP-set configuration on IOS

- Here we have statically configured one RP
 - We use acl, else for ff00::/8
 - Embedded-RP on by default

```
ipv6 multicast-routing
ipv6 pim rp-address 2001:660:3007:300:1:: rpm6bone
!
ipv6 access-list rpm6bone
 permit ipv6 any FF0E::/16
 permit ipv6 any FF1E::/16
 permit ipv6 any FF3E::/16
```

- This router has the above config, but is also an embedded-RP

```
ipv6 pim rp-address 2001:700:0:F000::1 rpemblo1
!
ipv6 access-list rpemblo1
 permit ipv6 any FF7E:140:2001:700:0:F000::/96
```

Site Router Configuration

- In general no configuration apart from RP-set is needed
- You may want to have a static RP for site scope
 - Some applications make use of FF05::/16
 - E.g. the IANA assigned group FF05::1:3 for all site DHCPv6 servers
 - Recommend configuring a loopback interface with an address that can be moved around, and static config on all site routers
 - Might use BSR for this
- You may want to use embedded-RP
 - Only need to configure on the RP itself
 - For example on Juniper, embedded-RP must be enabled on each router
 - Recommend one for global scope if you will create ASM multicast sessions that should be received externally
 - Might also use with internal scopes
- Edge routers need to use MLDv2 for SSM, normally the default

Interdomain IPv6 Multicast

- For IPv4, each site typically has their own RP for all global groups. RPs in different sites use MSDP to learn of remote sources
- This avoids relying on a 3rd party to host some central RP
- MSDP does not scale, and there is no MSDP for IPv6
- The lack of MSDP means that for a given global group there can be only one single RP on the Internet
 - It is not possible to have scalable services with global well known addresses
 - For example, there is SAP (sdr) with group ff0e::2:7ffe, which would require a central common RP for the service, so SAP is not really useful as a global announcement service

Embedded-RP 1/2

- Despite having a single RP per group, there is something we can do
- With embedded-RP the site hosting some multicast content or session can pick a group using their RP
 - Everyone on the Internet will then use their RP for that group
 - There is no 3rd party responsible
- Ideally embedded-RPs should be located near the edge
- A backbone network may not need to do any RP configuration
 - On some routers embedded-RP must be enabled, but just a simple toggle
 - Embedded-RP on IOS is on by default. Needs to be enabled on JunOS
- Some providers may wish to configure an RP and offer that as a service to customers. Several customers can use the same RP, but provider should then assign them different group ranges.
- Since some routers don't support embedded-RP yet, we have in the M6Bone (incl. Abilene) agreed on one central RP that may be used for ff0e::/16 and ff1e::/16. This does obviously not scale and is expected to be phased out soon.

Embedded-RP 2/2

- One difference from MSDP is that we get shared trees crossing the internet from receiver to the remote RP
- Embedded-RP addresses are based on your unicast addresses
 - Does not work with IANA-assigned addresses
 - or for application hard-coded groups
- Users and/or applications need to configure correct embedded-RP group range
 - Users should not need to understand embedded-RP
- With embedded-RP one should try to use an RP address that never changes
 - Good idea to use a loopback address that can be moved as needed

Router Configuration

- Most routers need minimal configuration
- As for IPv4 multicast, one may need to run multicast BGP peerings
 - Usually needed where unicast BGP is used and peers wish to exchange multicast
- Scope boundaries may need to be configured on border routers
- Edge routers need to use MLDv2 for SSM
 - May need to be enabled (if not default)
- For ASM, one may need to enable embedded-RP (if not default)
 - RP routers need to be statically configured
 - RPs only needed near the edge
- IOS will by default have IPv6 PIM on all IPv6 interfaces, MLDv2 and embedded-RP
- JUNOS will need explicit config for PIM, MLD, embedded-RP (and a Tunnel PIC in each router acting as an RP)

Configuring Scope Borders on IOS

- IOS can filter PIM messages on scope borders
 - If we use scope 8 (and smaller) for our site, we don't want any join, prunes, registers, BSR messages etc. for scope ≤ 8 to cross the border
 - The filter below filters all but register messages (in decimal, e.g. 10, not A)

```
interface ...  
  ipv6 multicast boundary scope 8
```

- Note that the command has changed, old:

```
interface ...  
  ipv6 zone boundary 8
```

Multiprotocol BGP

- As for IPv4 multicast, one may need to run multicast BGP peerings for RPF
- Multiprotocol BGP (MBGP) may have AFI (Address Family Identifier) IPv4/IPv6 and for those, SAFI (Subsequent AFI)
 - SAFI = 1 for unicast, 2 for multicast, 3 for both
 - Seems IETF wants to deprecate 3
- A peering may configure different policies for IPv4/IPv6 unicast/multicast
- Multicast routes used for RPF
 - Sometimes in addition to unicast routes
 - One may also sometimes translate unicast routes into multicast
- Recommend only using multiprotocol BGP
 - Whenever unicast BGP between multicast networks exists, enable multicast peering if multicast connectivity is desired
 - Try to avoid tricks like translation or unicast routes for RPF

IPv6 Multiprotocol BGP on IOS

```
router bgp 224
no bgp default ipv4-unicast
neighbor 2001:700:0:F019::2 remote-as 8933
neighbor 2001:700:0:F020::2 remote-as 25689
!
address-family ipv6 unicast
!
address-family ipv6 multicast
neighbor 2001:700:0:F019::2 activate
neighbor 2001:700:0:F019::2 prefix-list emb-rp-tu-prefixes in
neighbor 2001:700:0:F019::2 prefix-list le48 out
neighbor 2001:700:0:F020::2 activate
network 2001:700::/32
exit-address-family
```

- For multicast RPF, routes from static config, internal routing protocols and MBGP are used
 - Unicast BGP is usually not used
 - **ipv6 rpf use-bgp** enables use of unicast BGP

IPv6 Multicast Applications

- Mbone tools, vic/rat etc
 - IPv6 multicast conferencing applications
 - <http://www-mice.cs.ucl.ac.uk/multimedia/software/>
- VideoLAN (vlc)
 - Video streaming, also IPv6 multicast. Server and client
 - Many operating systems, both Windows and UNIX
 - <http://www.videolan.org/>
- DVTS <http://www.sfc.wide.ad.jp/DVTS/>
 - Streaming DV over RTP over IPv4/IPv6
 - DV devices using Firewire can be connected to two different machines and you can stream video between them over the Internet
 - Also pure software based playback
- Mad flute
 - Streaming of files using multicast (IPv4/IPv6 ASM/SSM)
 - Linux and Windows (not totally sure about *BSD status)
 - <http://www.atm.tut.fi/mad/>
- Tools for debugging/management
 - dbeacon, ssmping, asmping