# Introduction to IPv4 Multicast

*SANOG 8 Tutorial*

Atif Khan
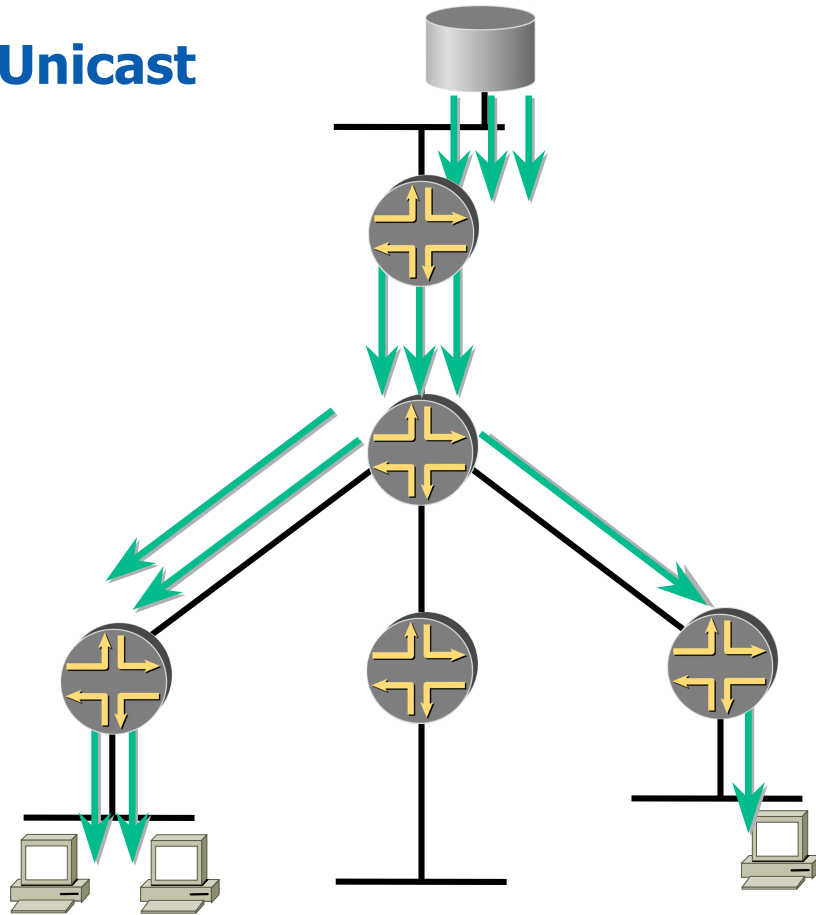
Juniper Networks

akhan@juniper.net

# Agenda

- Introduction

- Multicast addressing

- Group Membership Protocol

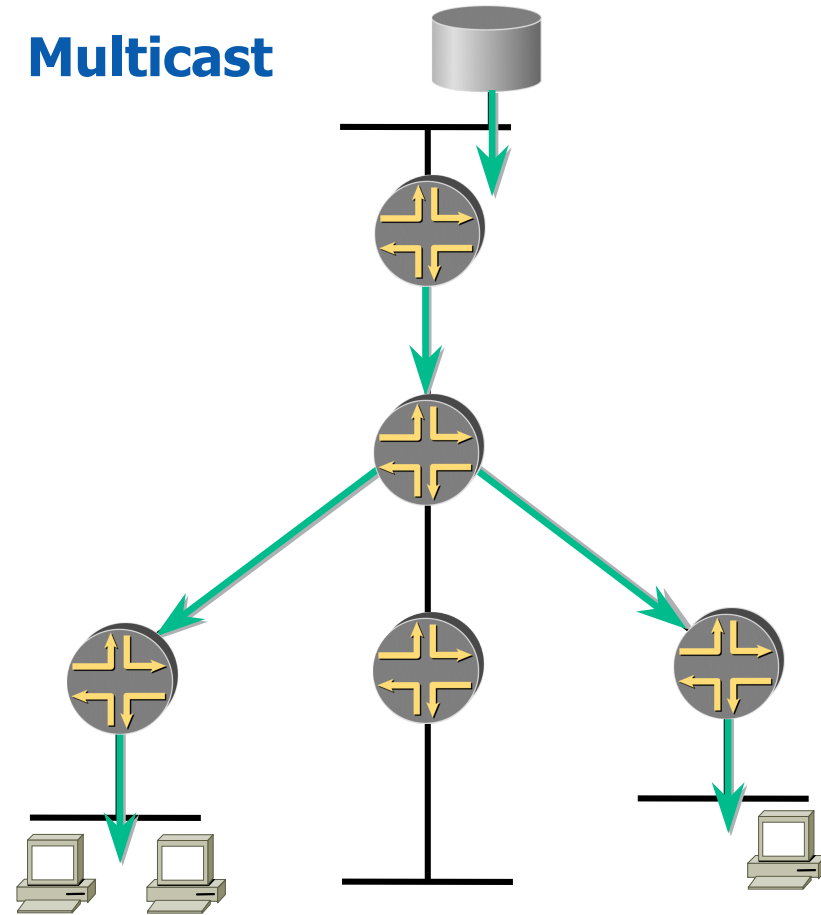- PIM-SM / SSM

- MSDP

- MBGP

- Summary

# Agenda

- **Introduction**

- Multicast addressing

- Group Membership Protocol

- PIM-SM / SSM

- MSDP

- MBGP
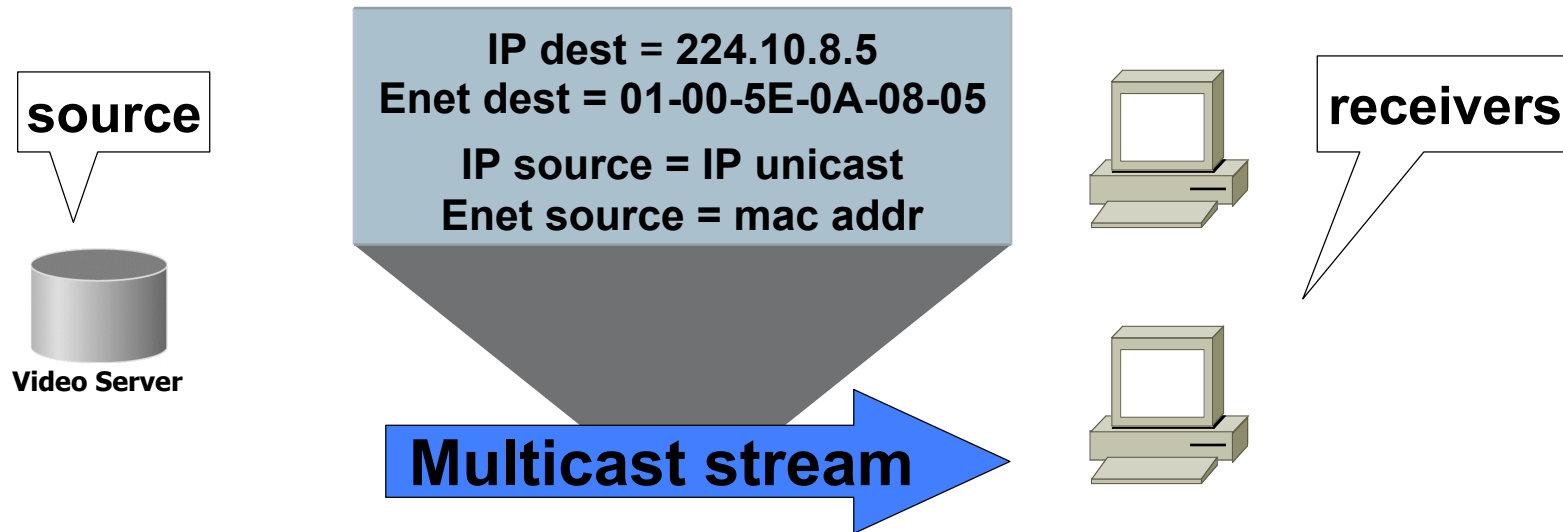
- Summary

# What is Multicasting?



**Unicast**

**Multicast**

# Multicast Uses

- Any Applications with multiple receivers
  - 1-to-many or many-to-many
- Live Video distribution
- Collaborative groupware
- Periodic Data Delivery - "Push" technology
  - stock quotes, sports scores, magazines, newspapers, ads
- Reducing Network/Resource Overhead
  - more efficient to establish multicast tree rather then multiple point-to-point links
- Distributed Interactive Simulation (DIS)
  - wargames
  - virtual reality

# Glossary of Terms: the basics



**source**

**Video Server**

IP dest = 224.10.8.5
Enet dest = 01-00-5E-0A-08-05

IP source = IP unicast
Enet source = mac addr

**receivers**

**Multicast stream**

- Source = source of multicast stream
- Multicast stream = IP packet with multicast address as IP destination address. a.k.a. multicast group.
    - s,g (unicast source, group) reference
    - UDP packets (TTL > 1 for routed nets)
- Receiver  = receiver (s) of multicast stream

# IP Multicast Building Blocks

- The SENDERS send
  - Multicast Addressing - rfc1700
  - class D (224.0.0.0 - 239.255.255.255)
- The RECEIVERS inform the routers what they want to receive
  - Internet Group Management Protocol (IGMP) - rfc2236 -> version 2
- The routers make sure the STREAMS make it to the correct receiving nets.
  - Multicast Routing Protocols (PIM-SM/SSM)
  - RPF  (reverse path forwarding) – against source address

# Multicast Forwarding

- 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.
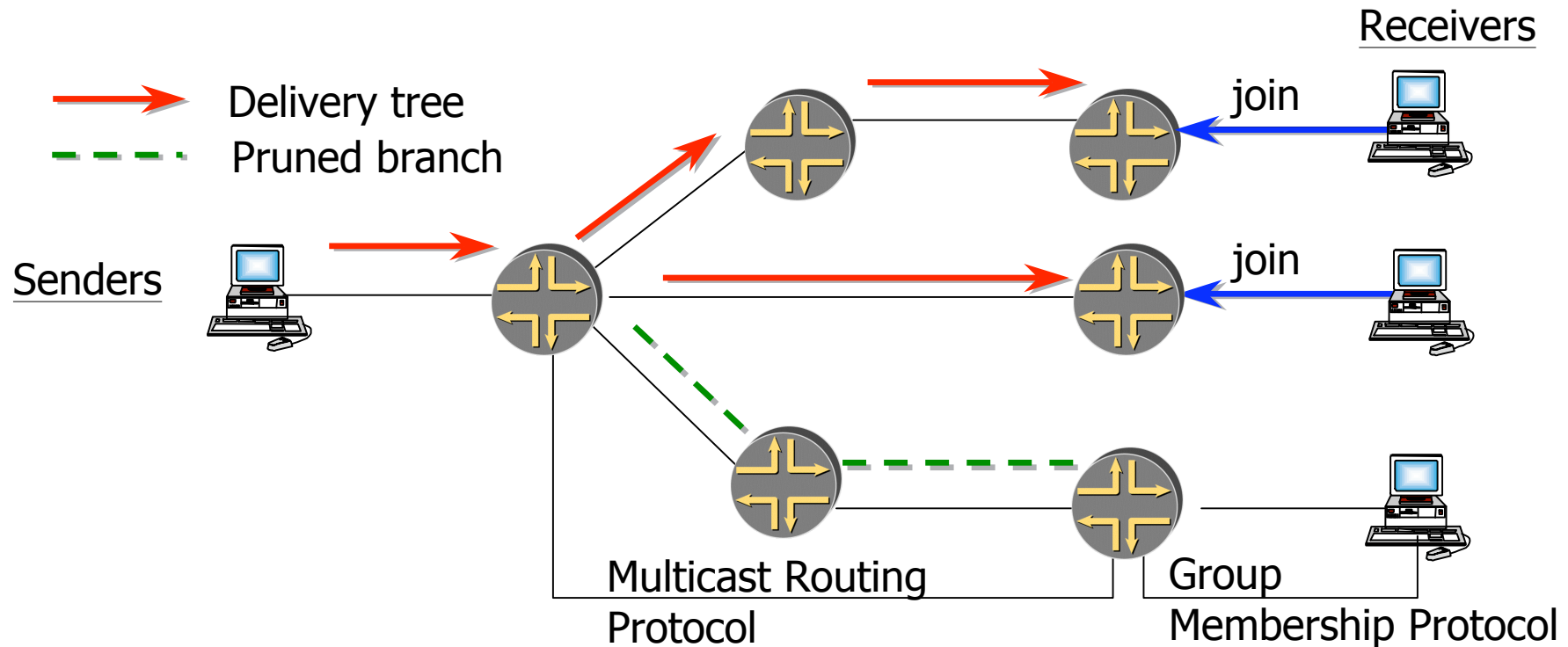
- Multicast Routing uses "Reverse Path Forwarding"

# Multicast Forwarding:
## Reverse Path Forwarding (RPF)

- ## 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 source IP address of incoming multicast packets are checked against a unicast routing table.
  - 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.

# Reverse Path Forwarding

- Multicast uses unicast routes to determine path back to source
- RPF checks ensure packets won't loop
- Routes contain incoming interface
  - Packets matching are forwarded
  - Packets mis-matching are dropped

# IP Multicast Components

Receivers

→ Delivery tree
--- Pruned branch

join

Senders

join

Multicast Routing
Protocol

Group
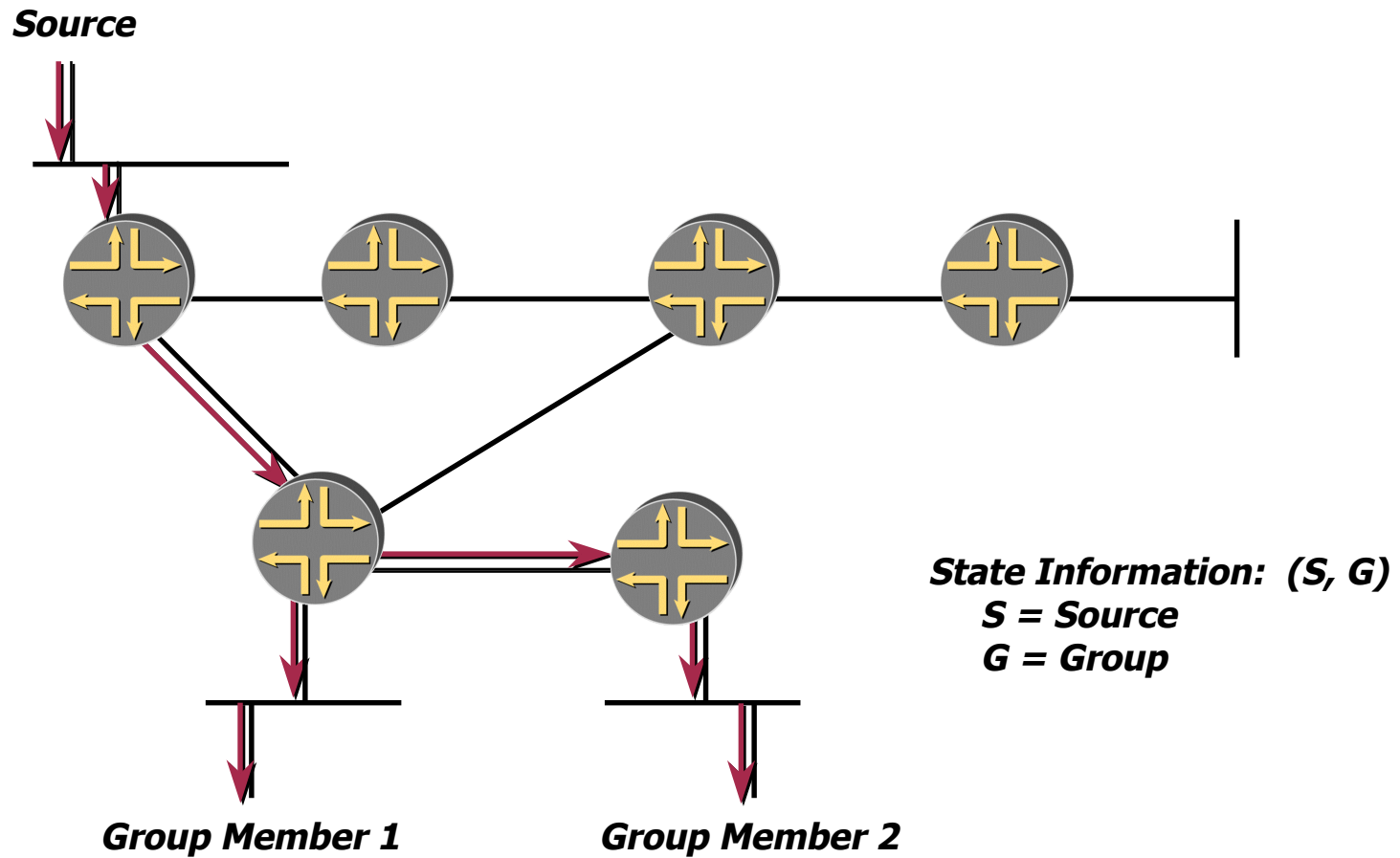Membership Protocol

- **Group Membership Protocol - enables hosts to dynamically join/leave multicast groups. Membership info is communicated to nearest router**

- **Multicast Routing Protocol - enables routers to build a delivery tree between the sender(s) and receivers of a multicast group**

# Multicast Distribution Trees

## Shortest Path or Source Based Distribution Tree



Source

Group Member 1        Group Member 2

State Information: (S, G)
S = Source
G = Group

# Multicast Distribution Trees

## Shared or Core Based Distribution Tree



Source 1

Core

Source 2

State Information:  (*, G)
   * = Any Source
   G = Group

Group Member 1          Group Member 2

# Multicast Distribution Trees

- Source or Shortest Path trees

  - More resource intensive; requires more state→ n(S x G)

  - You get optimal paths from source to all receivers, minimizes delay

  - Best for one-to-many distribution

- Shared or Core Based trees

  - Uses less resources; less  state →n(G)

  - You may get sub optimal paths from source to all receivers, depending on topology

  - The RP (core) itself and its location *may* affect performance

  - Best for many-to-many distribution

  - May be necessary for source discovery (PIM-SM)

# Agenda

# Multicast Addressing

- IP Multicast Group Addresses
  - 224.0.0.0–239.255.255.255
  - Class "D" Address Space
    - High order bits of 1st Octet = "1110"

# Multicast Addressing

- Contolled by Internet Assigned Numbers Authority - IANA
  - http://www.iana.org/assignments/multicast-addresses
    - 224.0.0.0/24: link local multicast range
    - 224.2.0.0/16: SAP/SDP range
    - 232.0.0.0/8:   SSM range
    - 233.0.0.0/8: AS-encoded statically assigned GLOP range
    - 239.0.0.0/8: administratively scoped multicast range

# Multicast Addresses - Layer 2

- **RFC 1700 - ethernet**

```
                              224.       10.       8.        5        <--  Class D IP address
  0000 0001   0000 0000   0101 1110   0xxx xxxx   xxxx xxxx   xxxx xxxx   <--  IANA's reserved block 01-00-5E
           |                                   |
       Multicast Bit                       0 = Internet Multicast
                                           1 = IANA reserved

  0000 0001   0000 0000   0101 1110   0000 1010   0000 1000  0000 0101   <--  MAC address 01-00-5E-0A-08-05
```

224.10.8.5 multicast stream maps to 01-00-5E-0A-08-05 ethernet layer 2 address.

- **rfc 1469 TR**
- **rfc 1390 FDDI**
- **rfc 2226 & 2022 - ATM**
- **rfc 1209 SMDS (broadcast)**

# Ethernet Multicast Addressing

- IANA Owns 01-00-5E Vendor Address Block
- Half of It Assigned for IP Multicast

Class D Address

0    8                    31

**32-Bit IP Address**

| 1110 | ignored | |
|------|---------|--|

**48-Bit Ethernet Address**

23 Bits

IEEE Ethernet Multicast bit    24    47

0

00000001000000000001011110 0

01-00-5E-

0 = Internet Multicast
1 = Reserved for Other Use

00-00-00
through
7F-FF-FF

# Agenda

- Introduction

- Multicast addressing

- **Group Membership Protocol**

- PIM-SM / SSM

- M-BGP

- MSDP

- Summary

# Internet Group Management Protocol (IGMP)

- How hosts tell routers about group membership

- Routers solicit group membership from directly connected hosts

- RFC 1112 specifies version 1 of IGMP

  - Supported on Windows 95

- RFC 2236 specifies version 2 of IGMP

  - Supported on latest service pack for Windows, newer Windows releases, and most UNIX systems

- IGMP version 3 is specified in RFC 3376

  - Provides source include-list capabilities

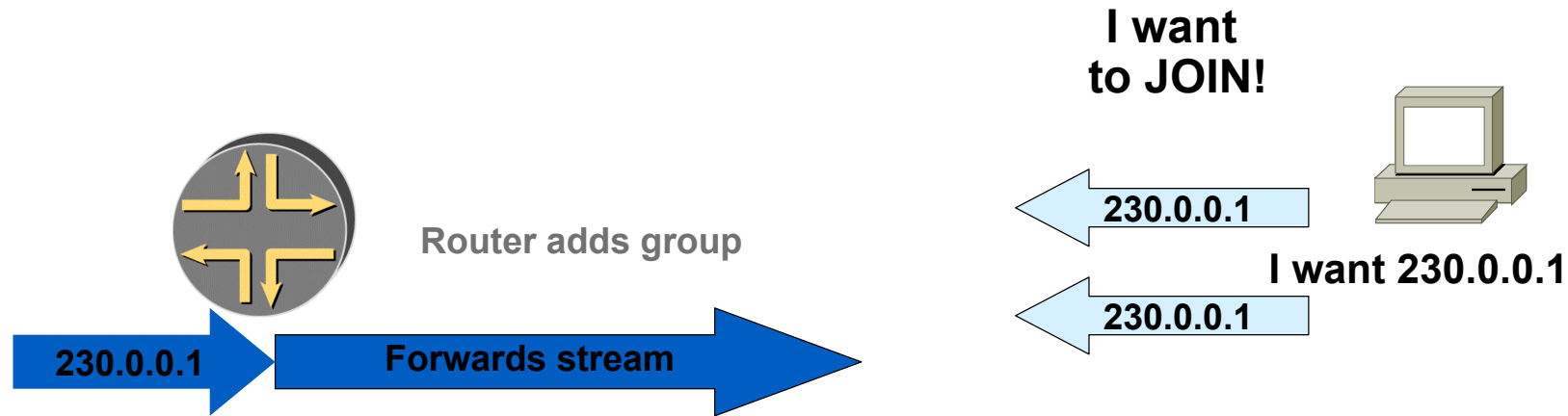  - See www.ietf.org for more information

# IGMP Details

- **Router:**
  - **sends Membership Query messages to All Hosts (224.0.0.1)**
    - query-interval = 125 secs default
  - **router with lowest IP address is Querier (rest non-queriers)**
  - **If lower-IP address query heard, backoff to non-querier state**
  - **listens for reports (whether querier or not) and adds group to membership list for that interface**
    - query-response-interval = 10 secs default
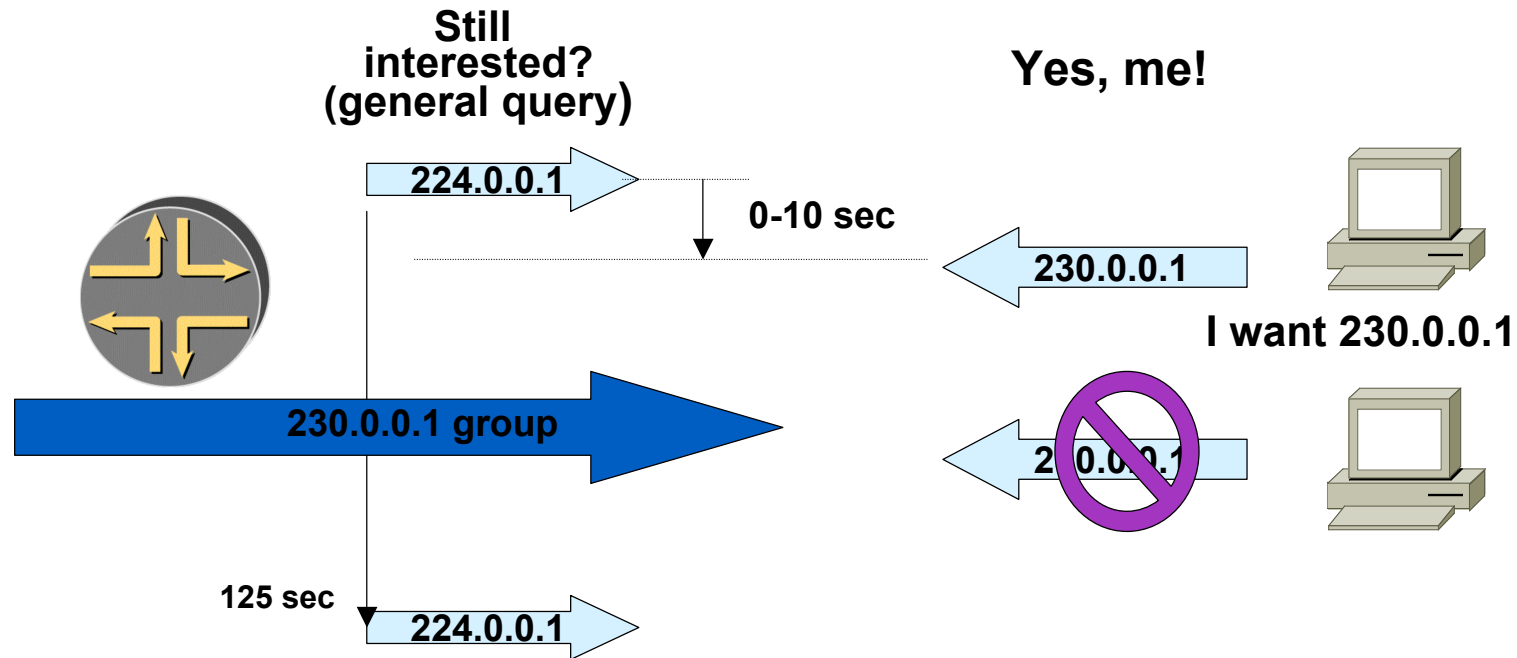
# IGMP Details

- Host:
  - sends Membership Report messages, if joined
    - waits 0-10 sec (def).
    - Hosts listen to other host reports
    - Only 1 host responds
  - Join messages (unsolicited Membership Report) to group address (e.g. 224.10.8.5)
  - Leave messages to All Routers (224.0.0.2)
  - IGMPv1/2 reports group membership ONLY – No sources

# IGMP Protocol Flow - Join a Group

**I want to JOIN!**

**Router adds group**

230.0.0.1

230.0.0.1

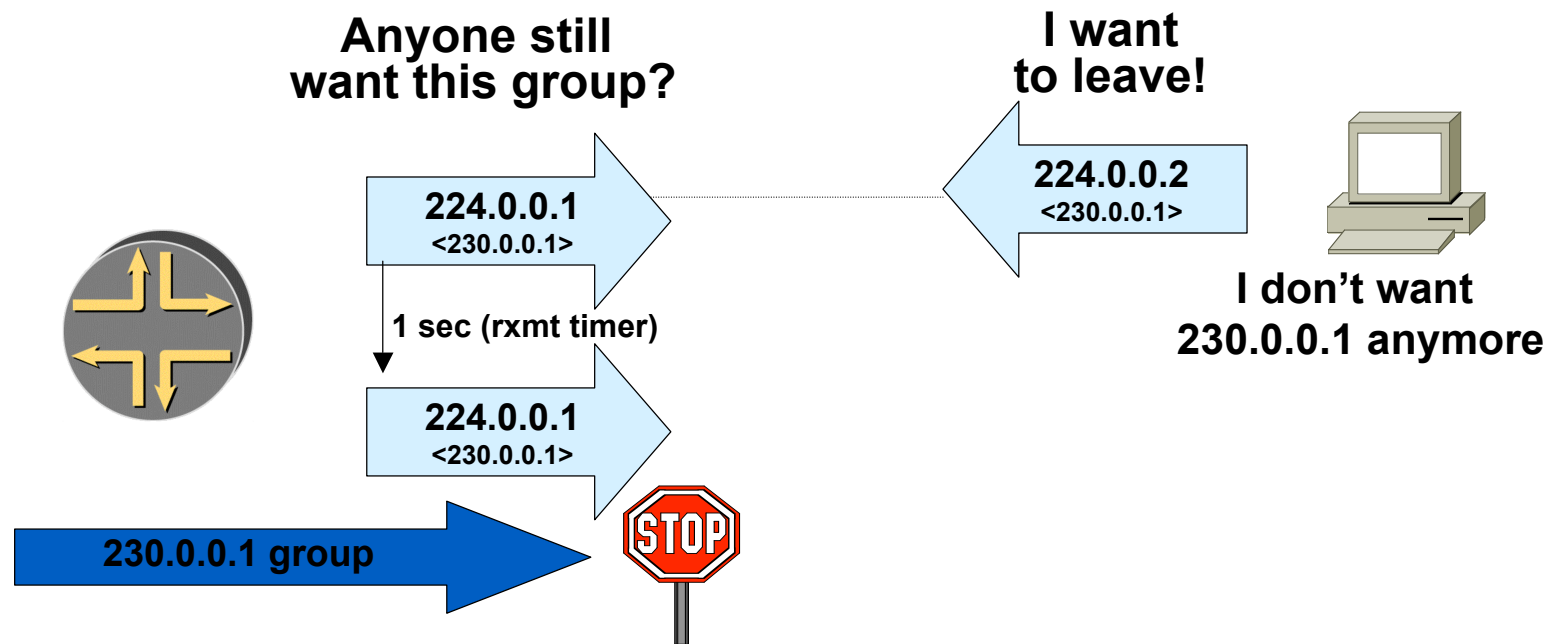**I want 230.0.0.1**

230.0.0.1

**Forwards stream**

230.0.0.1

- Router triggers group membership request to PIM.
- Hosts can send unsolicited *join* membership messages – called reports in the RFC (usually more than 1)
- Or hosts can join by responding to periodic query from router

# IGMP Protocol Flow - Querier

**Still interested? (general query)**

**Yes, me!**

224.0.0.1

0-10 sec

230.0.0.1

I want 230.0.0.1

230.0.0.1 group

230.0.0.1

125 sec

224.0.0.1

- Hosts respond to *query* to indicate (new or continued) interest in group(s)
  - only 1 host should respond per group
    - Hosts fall into idle-member state when same-group report heard.
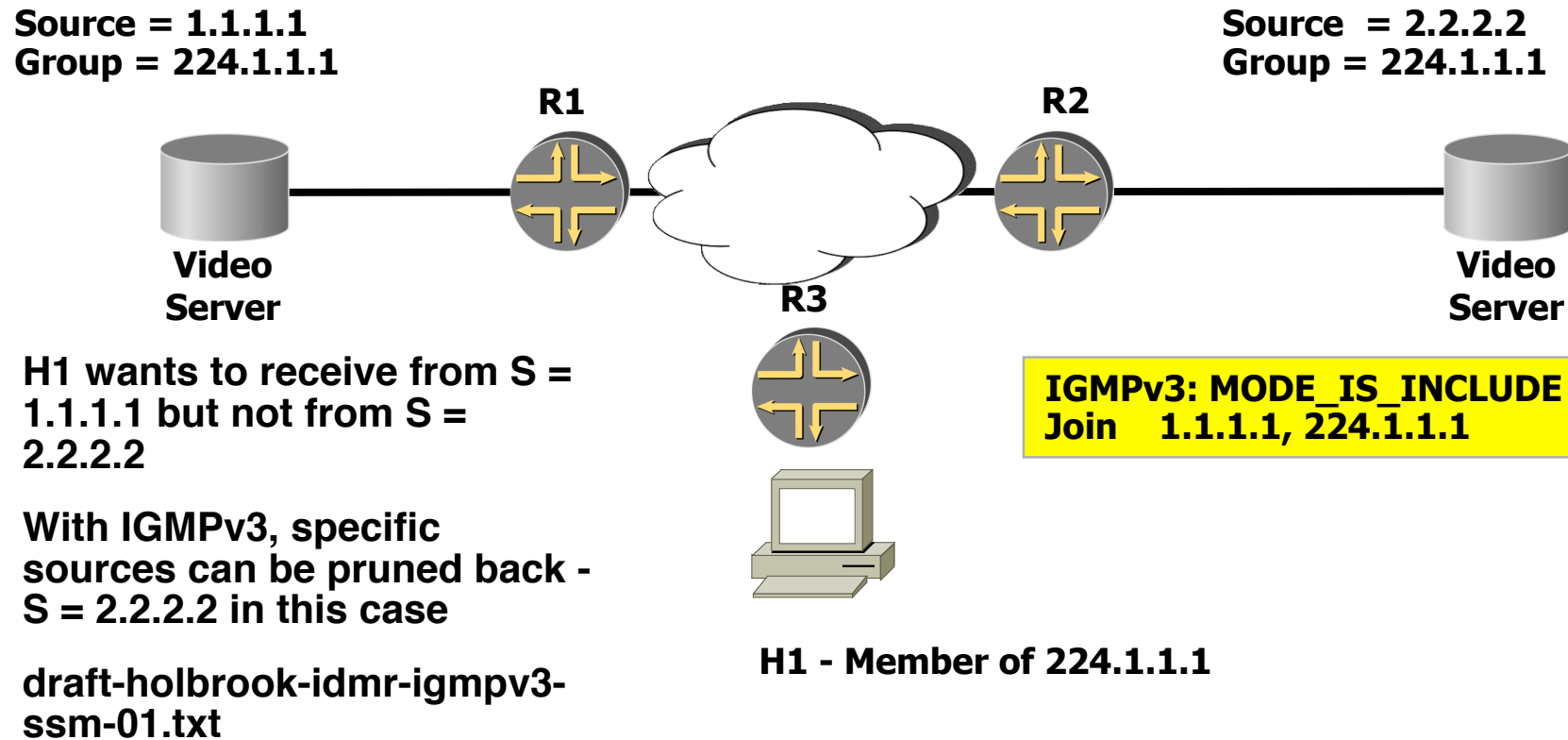- After 260 sec with no response, router times out group

# IGMP Protocol Flow - Leave a Group

**Anyone still
want this group?**

**I want
to leave!**

**224.0.0.1**
<230.0.0.1>

**224.0.0.2**
<230.0.0.1>

**1 sec (rxmt timer)**

**224.0.0.1**
<230.0.0.1>

**I don't want
230.0.0.1 anymore**

**230.0.0.1 group**

**STOP**

- Hosts that support IGMP v2 send *leave* messages to all routers group indicating group they're leaving.
  - Router follows up with 2 *group-specific queries* messages
- IGMP v1 hosts leave by not responding to *queries* (260 sec timeout)

# IGMPv3

**Enables hosts to listen only to a specified subset of the hosts sending to the group**

**Source = 1.1.1.1**
**Group = 224.1.1.1**

**Source = 2.2.2.2**
**Group = 224.1.1.1**

**R1**

**R2**

**Video Server**

**Video Server**

**R3**

- H1 wants to receive from S = 1.1.1.1 but not from S = 2.2.2.2

**IGMPv3: MODE_IS_INCLUDE**
**Join      1.1.1.1, 224.1.1.1**

- With IGMPv3, specific sources can be pruned back - S = 2.2.2.2 in this case

- draft-holbrook-idmr-igmpv3-ssm-01.txt

**H1 - Member of 224.1.1.1**

# IGMP Details

- IGMP Version 2
  - Multicast router with lowest IP address is elected querier
    - IGMPv1 was mcast protocol specific and potentially conflicted.
  - Group-Specific Query message is defined. Enables router to transmit query to specific multicast address rather than to the "all-hosts" address of 224.0.0.1
  - Leave Group message is defined. Last host in group wishes to leave, it sends Leave Group message to the "all-routers" address of 224.0.0.2. Router then transmits Group-Specific query and if no reports come in, then the router removes that group from the list of group memberships for that interface
- IGMP Version 3
  - Group-Source Report message is defined. Enables hosts to specify which senders it can receive or not receive data from.
  - Group-Source Leave message is defined. Enables host to specify the specific IP addresses of a (source,group) that it wishes to leave.

# Agenda

- Introduction

- Multicast addressing

- Group Membership Protocol

- **PIM-SM / SSM**

- MSDP

- MBGP

- Summary

# PIM-SM

- Protocol Independent Multicast - sparse mode
  - explicit join: assumes everyone does not want the data
  - uses unicast routing table for RPF checking
  - data and joins are forwarded to RP for initial rendezvous
  - all routers in a PIM domain must have RP mapping
  - when load exceeds threshold forwarding swaps to shortest path tree (default is first packet)
  - state increases (not everywhere) as number of sources and number of groups increase
  - source-tree state is refreshed when data is forwarded and with Join/Prune control messages
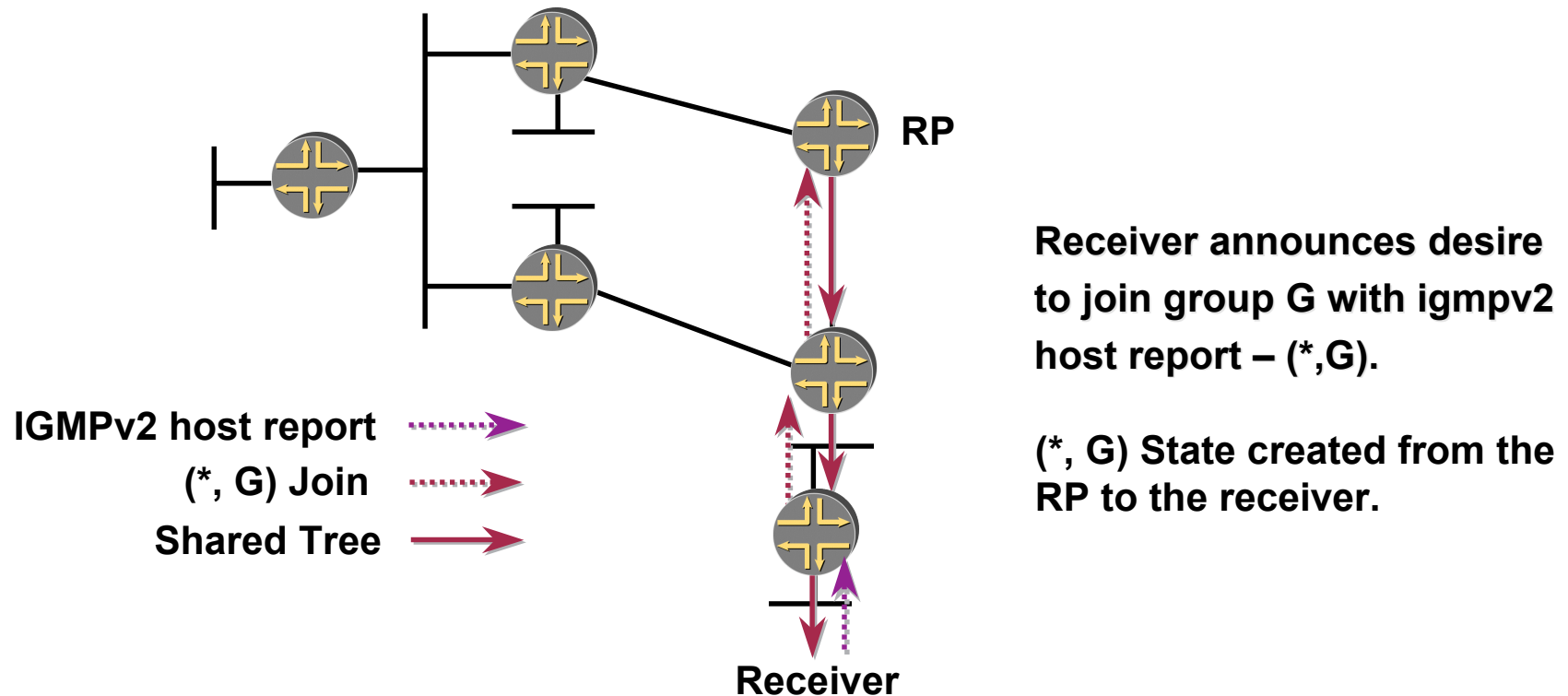
# PIM-SM Operation
# Designated Router (DR)

- Neighboring PIM-SM routers multicast periodic "Hello" messages to each other - default 30 secs.
    - Hello-interval tunable for faster convergence
- On receipt of a Hello message
    - a router stores the IP address and priority for that neighbor
- Router with highest IP address is selected as the DR, if the priorities match
- When DR goes down:
    - new one selected by scanning all neighbors on the interface and choosing the one with the highest IP address
- DR sends
    - "Join/Prune" messages toward the RP from receiver network
    - "Register" messages toward the RP from source network
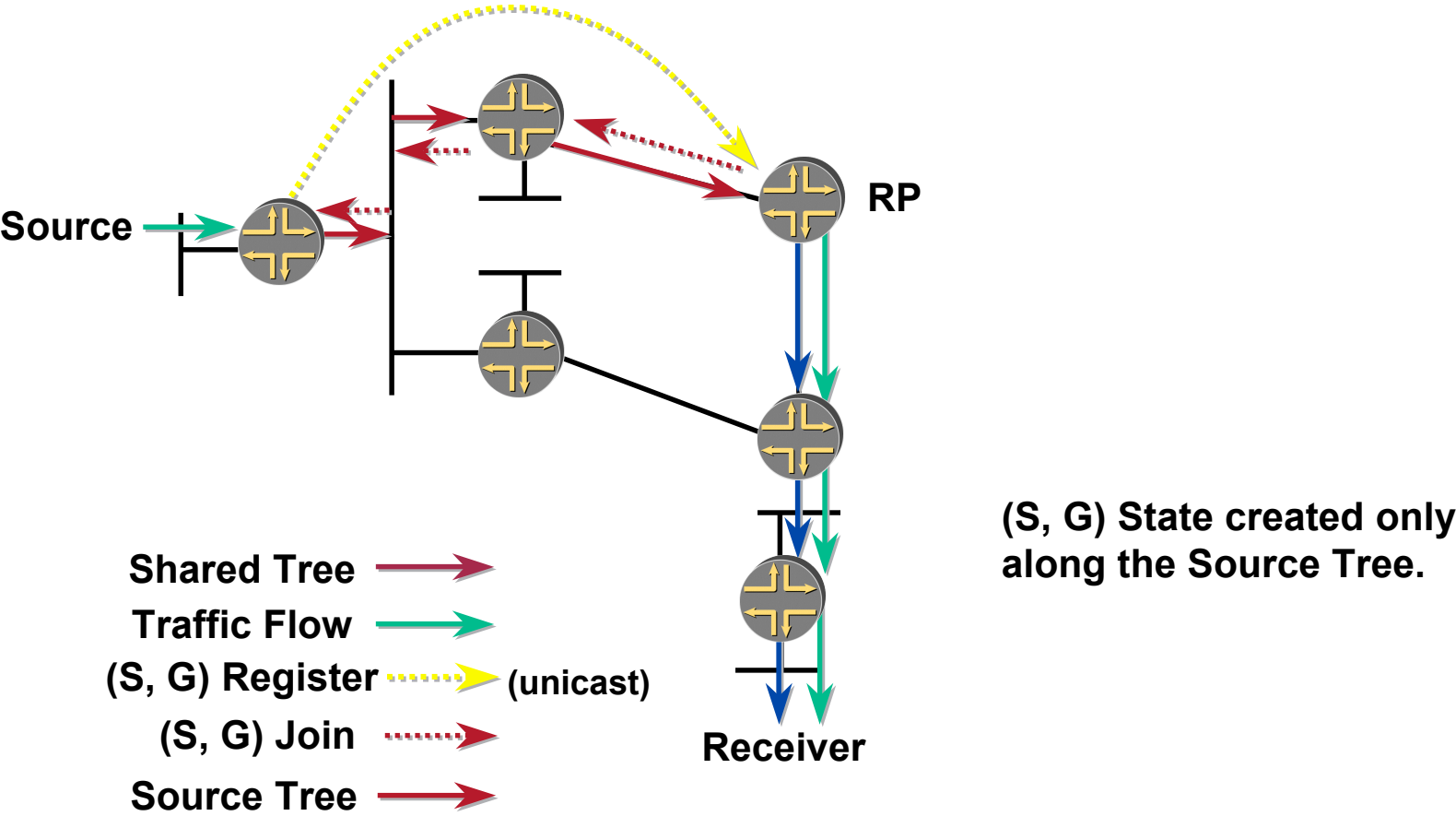
# PIM Sparse-Mode :RP

- Allows Source Trees or Shared Trees
- Rendezvous Point (RP)
  - Matches senders with receivers
  - Provides network source discovery
  - Root of shared tree
- Typically use shared tree to bootstrap source tree
- RP's can be learned via:
  - Static configuration – RECOMMENDED
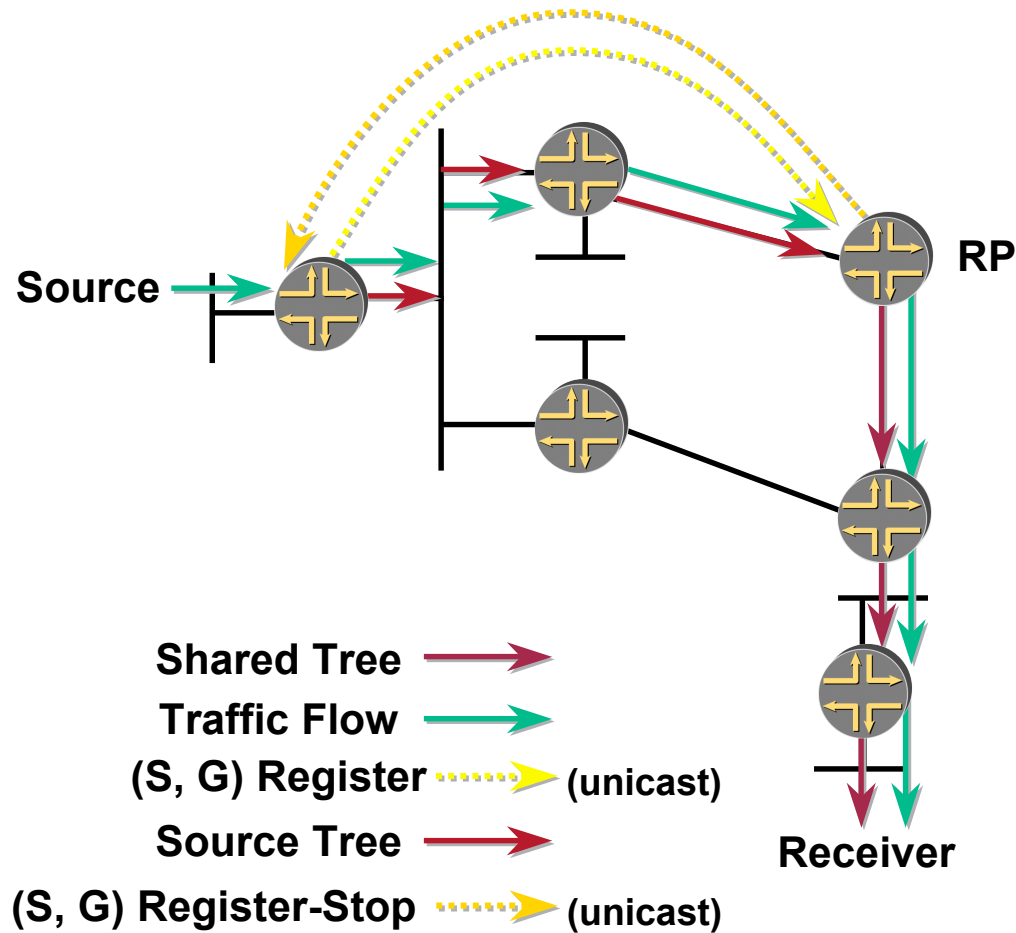  - Auto-RP (V1 & V2)
  - Bootstrap Router (V2)

# PIM-SM Shared Tree Join

**IGMPv2 host report** ┄┄┄►
**(*, G) Join** ┄┄┄►
**Shared Tree** ───►

**RP**

**Receiver**

Receiver announces desire to join group G with igmpv2 host report – (*,G).

(*, G) State created from the RP to the receiver.

# PIM-SM Sender Registration



Source

RP

(S, G) State created only
along the Source Tree.

Shared Tree

Traffic Flow

(S, G) Register (unicast)

(S, G) Join

Source Tree

Receiver

# PIM-SM Sender Registration



**Source**

**RP**

**Receiver**

Shared Tree ➝

Traffic Flow ➝

(S, G) Register ┈┈➤ (unicast)

Source Tree ➝

(S, G) Register-Stop ┈┈➤ (unicast)
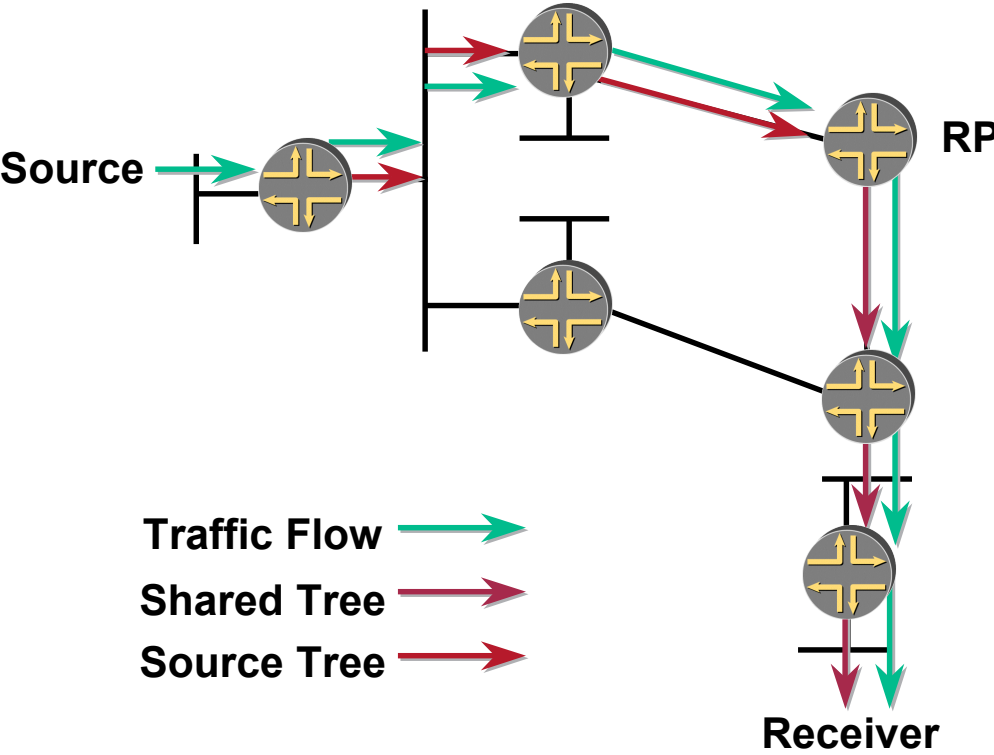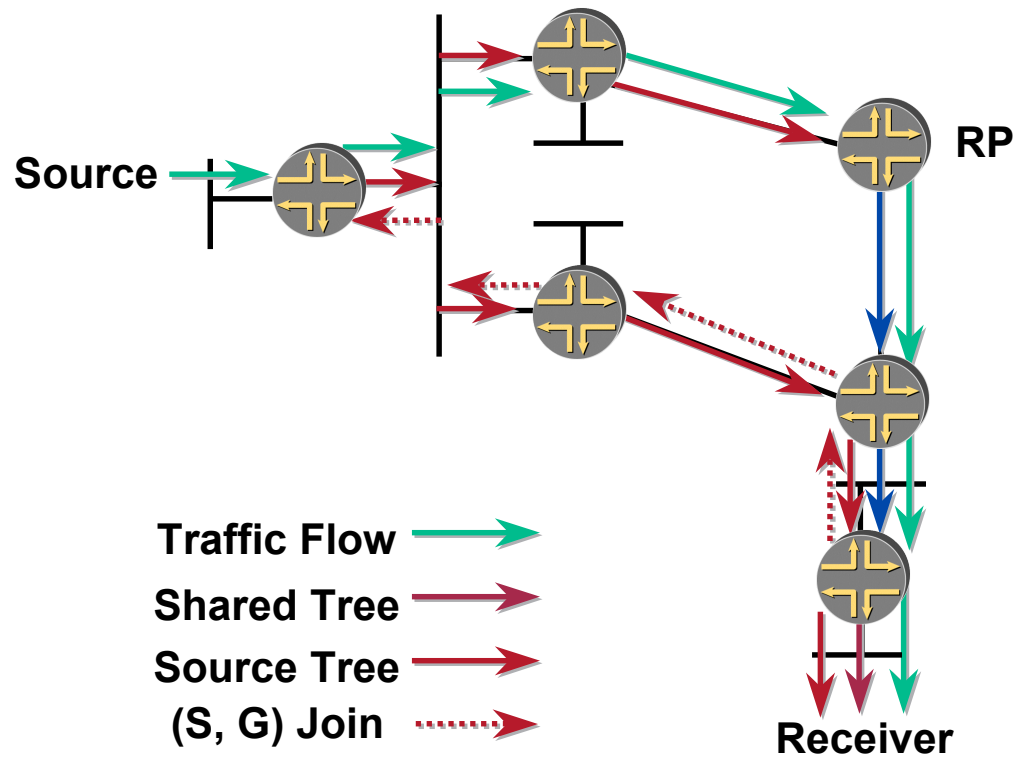
(S, G) traffic begins arriving at the RP via the Source tree.

RP sends a Register-Stop back to the first-hop router to stop the Register process.

# PIM-SM Sender Registration

**Source**

**RP**

**Receiver**

Traffic Flow →

Shared Tree →

Source Tree →

Source traffic flows natively along SPT to RP.

From RP, traffic flows down the Shared Tree to Receivers.

# PIM-SM SPT Cutover



**Source**

**RP**

Traffic Flow

Shared Tree

Source Tree

(S, G) Join

**Receiver**
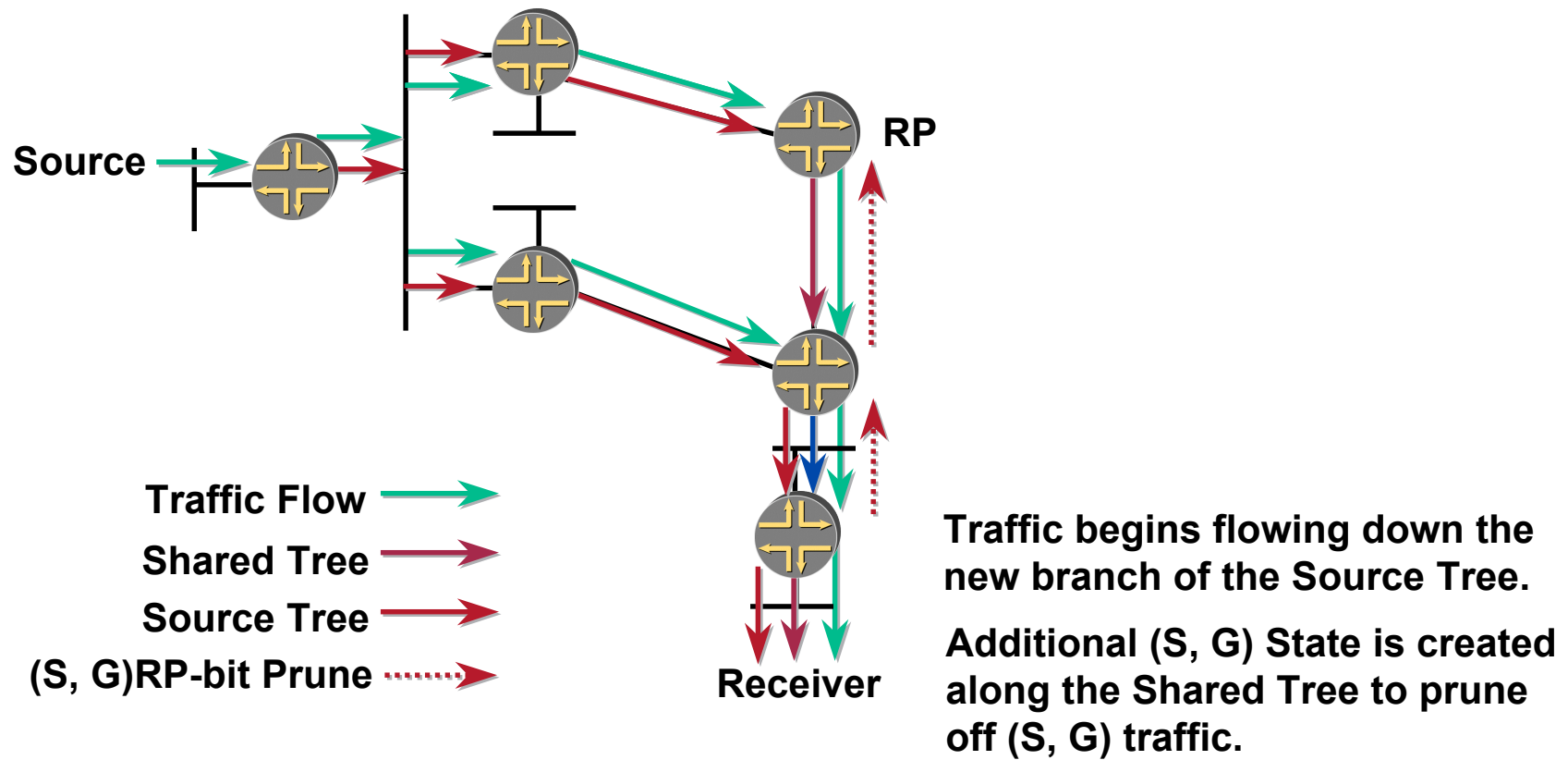
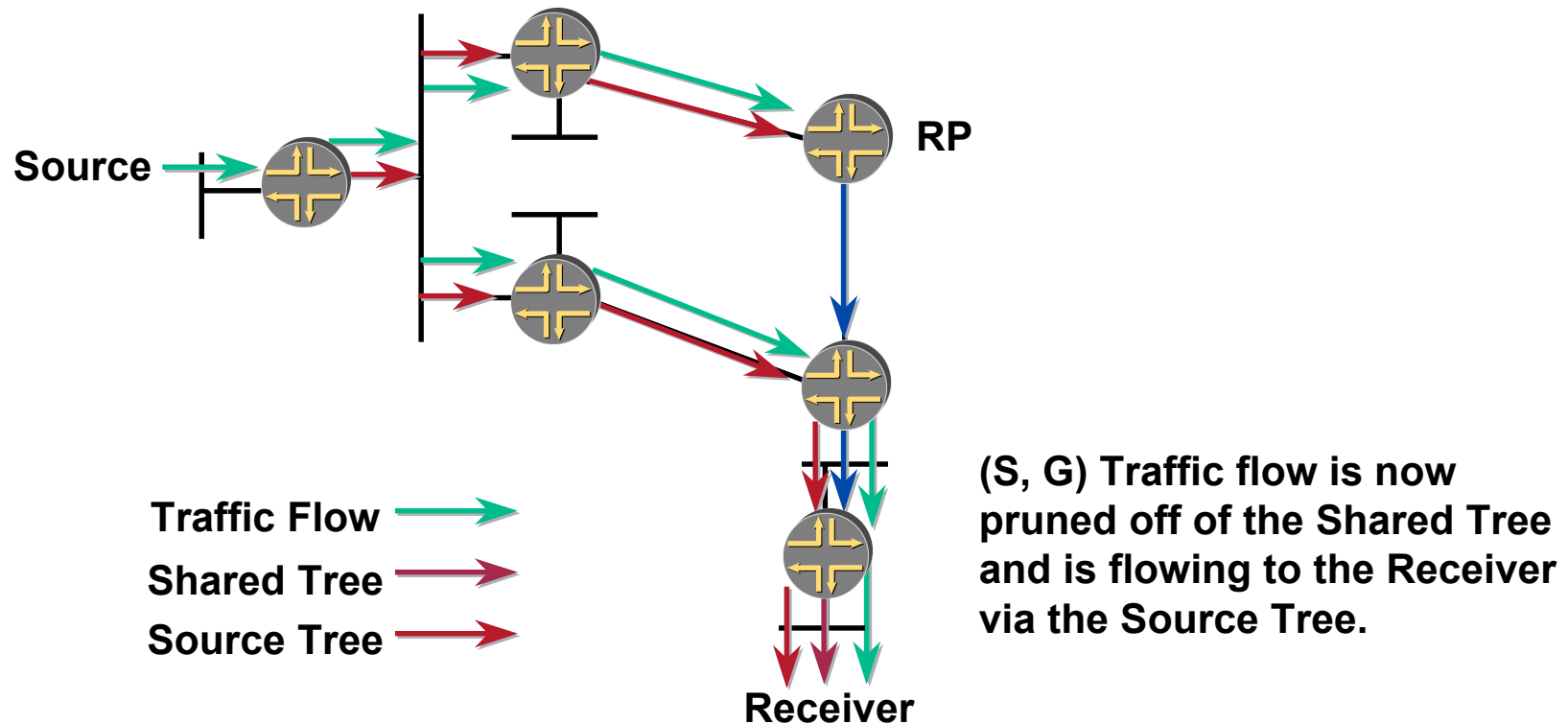Last-hop router joins the Source Tree.

Additional (S, G) State is created along new part of the Source Tree.

# PIM-SM SPT Cutover



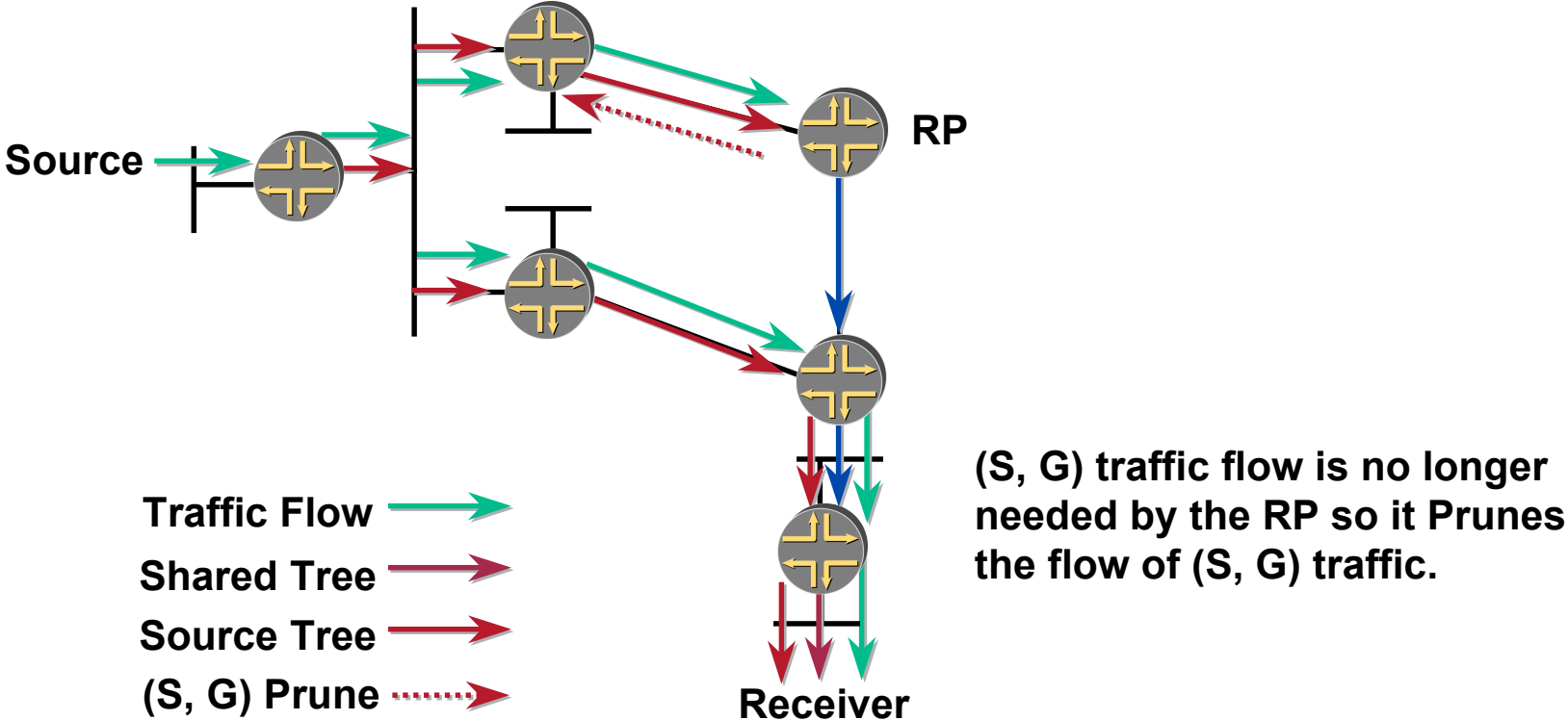**Source**

**RP**

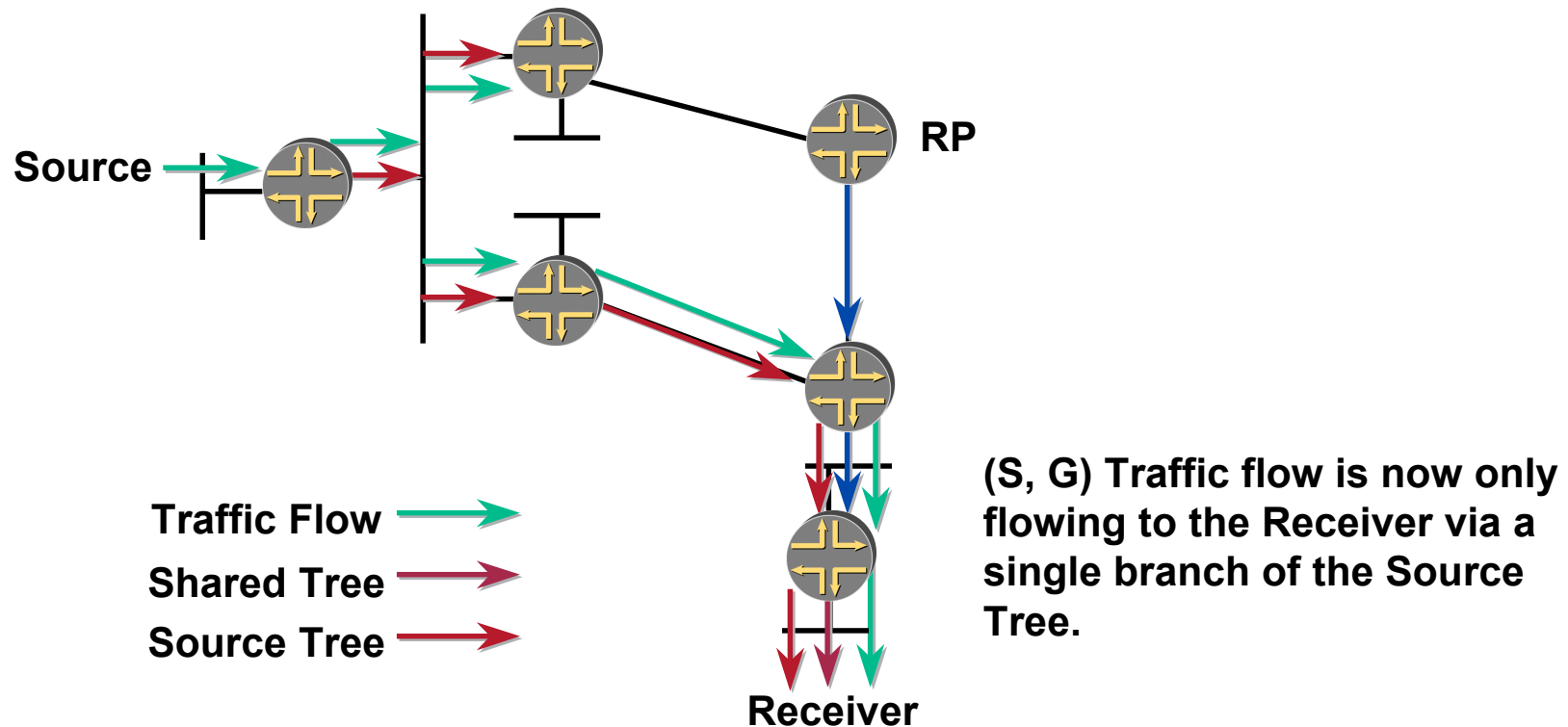**Receiver**

Traffic Flow ———→

Shared Tree ———→

Source Tree ———→

(S, G)RP-bit Prune ┈┈┈→

**Traffic begins flowing down the new branch of the Source Tree.**

**Additional (S, G) State is created along the Shared Tree to prune off (S, G) traffic.**

# PIM-SM SPT Cutover



Source

RP

Traffic Flow

Shared Tree

Source Tree

Receiver

(S, G) Traffic flow is now pruned off of the Shared Tree and is flowing to the Receiver via the Source Tree.

# PIM-SM SPT Cutover

**Source**

**RP**

(S, G) traffic flow is no longer needed by the RP so it Prunes the flow of (S, G) traffic.

**Traffic Flow** ⟶

**Shared Tree** ⟶

**Source Tree** ⟶

**(S, G) Prune** ┈┈▶

**Receiver**

# PIM-SM SPT Cutover

Source

RP

Receiver

Traffic Flow
Shared Tree
Source Tree

(S, G) Traffic flow is now only flowing to the Receiver via a single branch of the Source Tree.
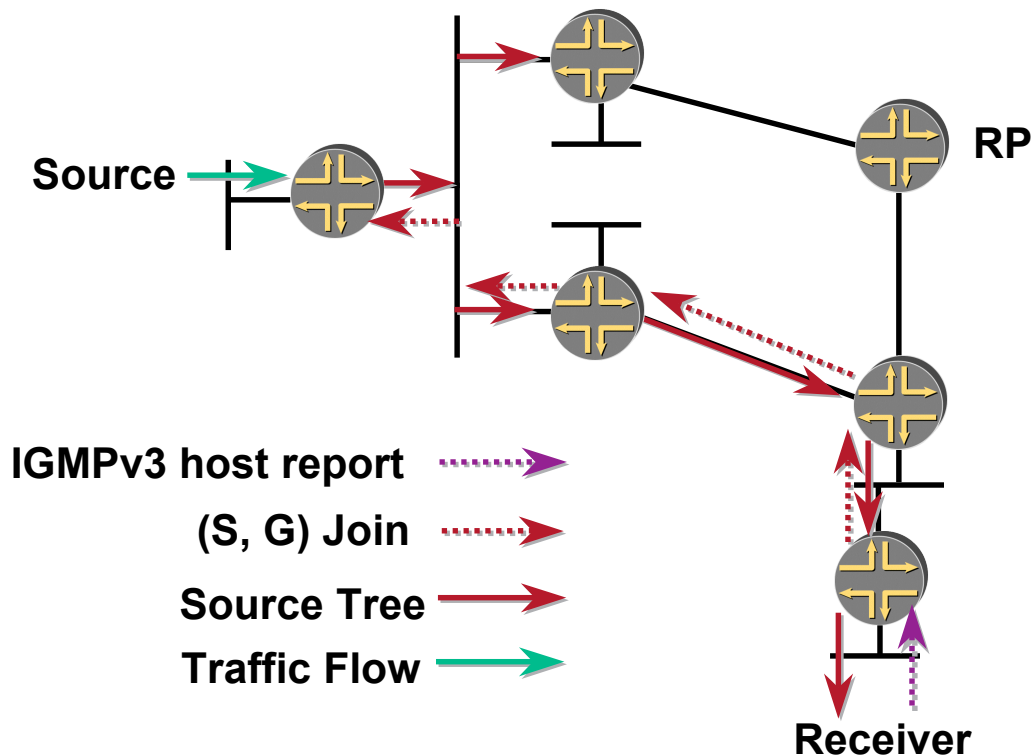
# PIM-SSM

- No shared trees

- No register packets

- No RP mapping required (no RP required!)

- No RP-to-RP source discovery (MSDP)

- Requires IGMP include-source list – IGMPv3

- Hard-coded behavior in 232/8

  - Configurable to expand range

# PIM-SSM



**Source**

**RP**

IGMPv3 host report ········▶

(S, G) Join ········▶

Source Tree ────▶

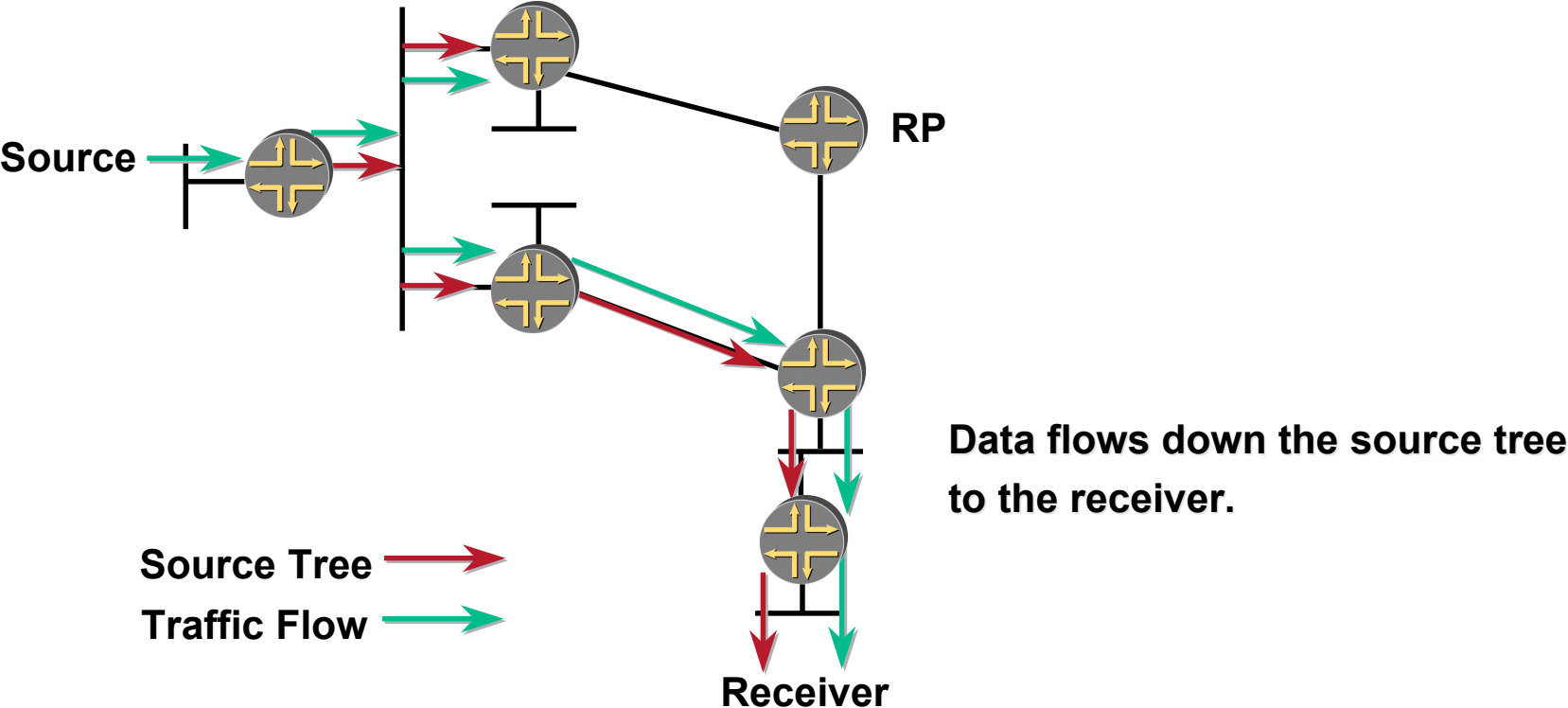Traffic Flow ────▶

**Receiver**

Receiver announces desire
to join group G AND source
S with an IGMPv3 include-list.

Last-hop router joins the Source
Tree.

(S,G) state is built between the
source and the receiver.

# PIM-SSM

Source

RP

Data flows down the source tree
to the receiver.

Source Tree →

Traffic Flow →

Receiver

# Agenda

- Introduction

- Multicast addressing

- Group Membership Protocol

- PIM-SM / SSM

- **MSDP**

- MBGP

- Summary

# MSDP

- Multicast Source Discovery Protocol
  - Allows each domain to control its own RP(s)
  - Interconnect RPs between domains with TCP connections to pass source active messages (SAs)
  - Can also be used within a domain to provide RP redundancy (Anycast-RP)
  - RPs send SA messages for internal sources to MSDP peers
  - SAs are Peer-RPF checked before accepting or forwarding
  - RPs learn about external sources via SA messages and may trigger (S,G)joins on behalf of local receivers
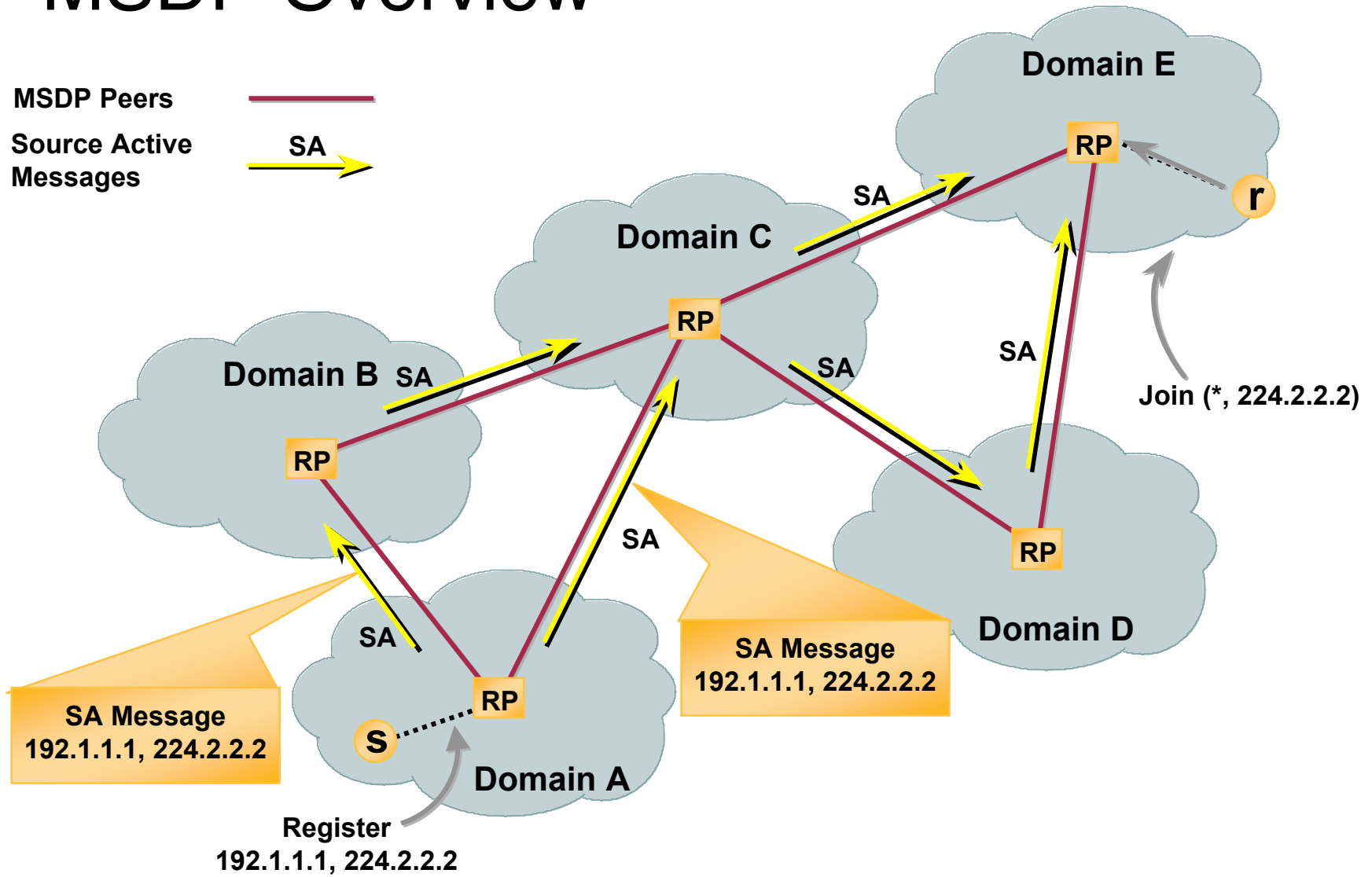  - MSDP connections typically parallel MBGP connections

# MSDP Operation

- MSDP peers (inter or intra domain)
  - (TCP port 639 w/ higher IP addr LISTENS)
- "FLOOD & join"
  - SA (source active) packets periodically sent to MSDP peers indicating:
    - source address of active streams
    - group address of active streams
    - IP address of RP originating the SA
    - only originate SA's for your sources w/in your domain
- "flood & JOIN"
  - interested parties can send PIM JOIN's towards source (creates inter-domain source trees)
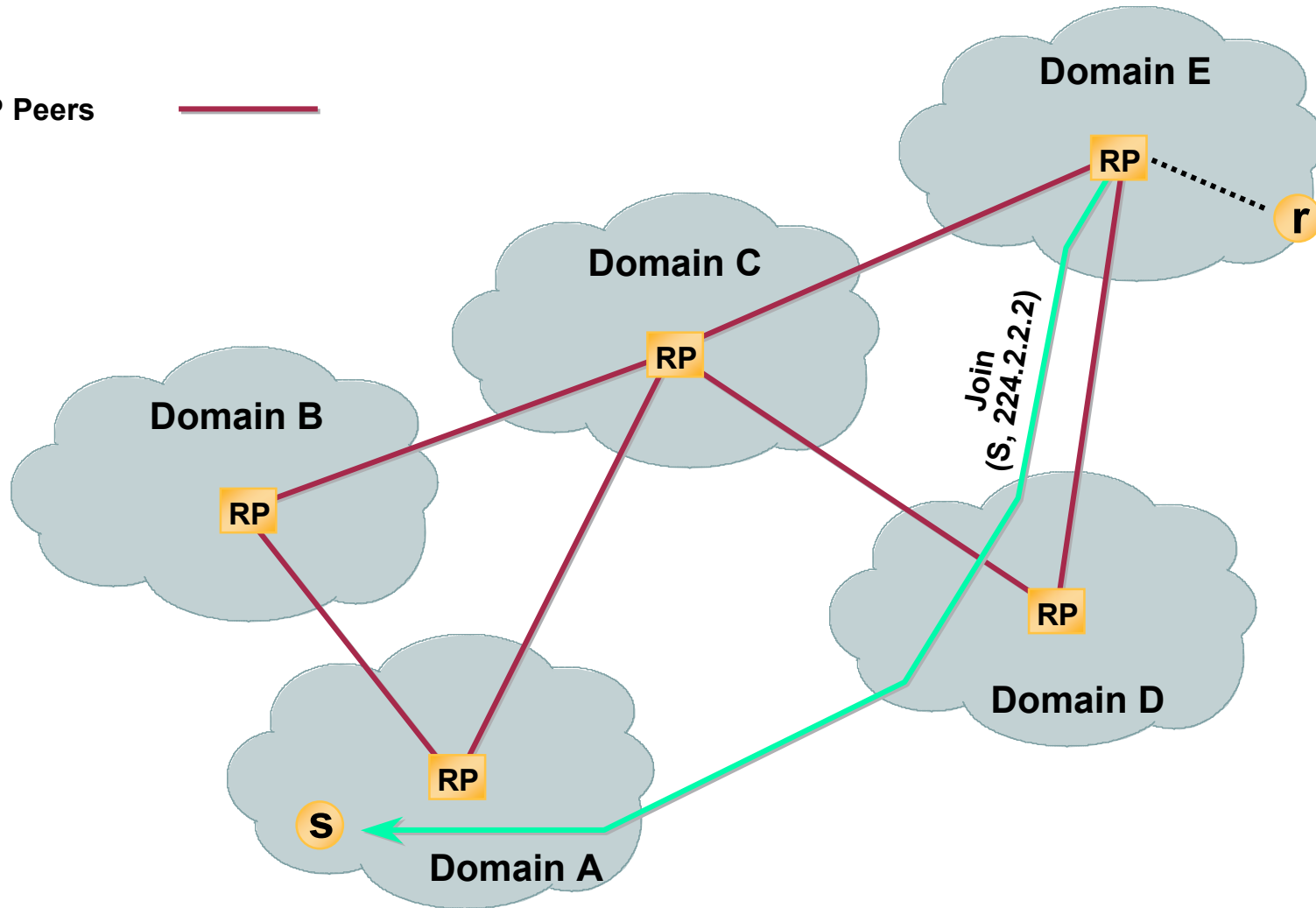
# MSDP Source Active Msgs

- Initial SA message sent when source first registers
  - May optionally encapsulate first data packet
- Subsequent SA messages periodically refreshed every 60 seconds as long as source still active by originating RP
- Other MSDP peers don't originate this SA but only forward it if received
- SA messages cached on router for new group members that may join
  - Reduced join latency
  - Prevent SA storm propagation

# MSDP Overview

**MSDP Peers** ———————

**Source Active Messages** SA ——→

**Domain E**

RP

r

SA

**Domain C**

RP

SA

**Domain B** SA

RP

SA

SA

**Join (*, 224.2.2.2)**

SA

RP

**Domain D**

SA

SA

**SA Message
192.1.1.1, 224.2.2.2**

RP

S

**SA Message
192.1.1.1, 224.2.2.2**

**Domain A**

**Register
192.1.1.1, 224.2.2.2**

# MSDP Overview



**MSDP Peers** ——————

Domain E

RP ......... r

Domain C

RP

Join (S, 224.2.2.2)

Domain B

RP

Domain D

RP

Domain A

RP

S

# MSDP Overview



**MSDP Peers** ————

**Multicast Traffic** ————

Domain E

Domain C

Domain B

Domain D

Domain A

RP

RP

RP

RP

RP

S

r

# MSDP Overview

**MSDP Peers** ————

**Multicast Traffic** ————

Domain E

Domain C

Domain B

Domain D

Domain A

RP

RP

RP

RP

RP

S

r

Join
(S, 224.2.2.2)

# MSDP Overview

MSDP Peers ────────

Multicast Traffic ────────

Domain E

Domain C

Domain B

RP

RP

RP

RP

RP

S

r

Domain A

Domain D

# MSDP Peers

- MSDP establishes a neighbor relationship between MSDP peers
  - Peers connect using TCP port 639
  - Peers send keepalives every 60 secs (fixed)
  - Peer connection reset after 75 seconds if no MSDP packets or keepalives are received
- MSDP peers must run mBGP!
  - May be an MBGP peer, a BGP peer or both
  - Required for peer-RPF checking of the RP address in the SA to prevent SA looping
  - Exception: BGP is unnecessary when peering with only a single MSDP peer (default-peer)
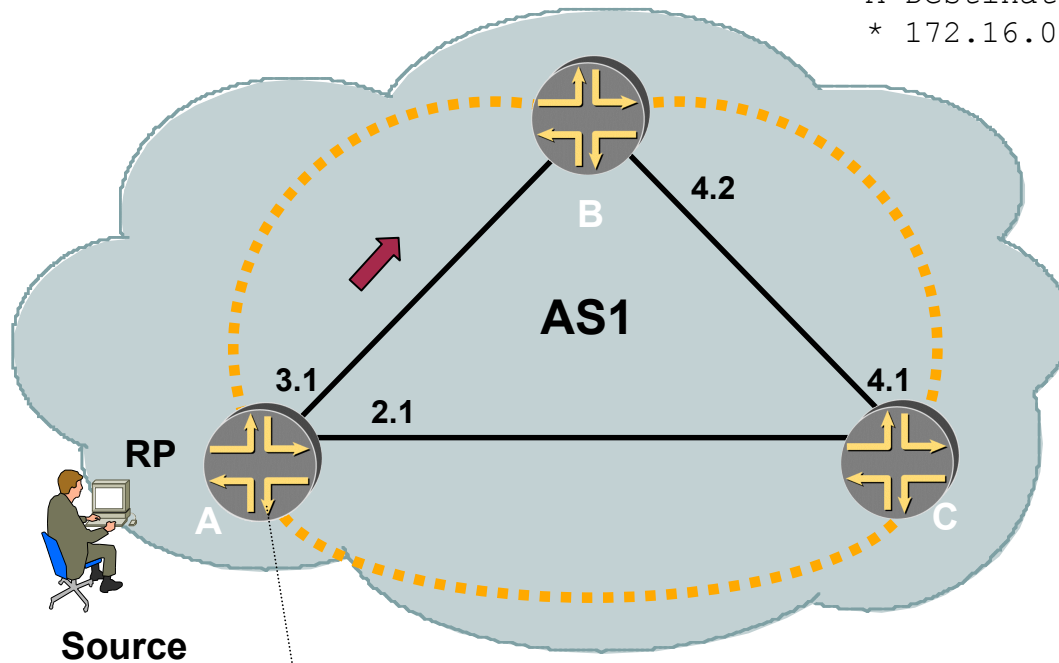
# Receiving SA Messages

- Skip RPF Check and accept SA if:
  - Sending MSDP peer is default-peer
  - Sending MSDP peer = Mesh-Group peer
- RPF Check the received SA message
  - If the MSDP peer IS THE originating RP – then accept.
  - Lookup best MBGP path to RP in SA message
  - Is the sending MSDP Peer also an MBGP peer?
    - Yes: Is best path to RP via this MBGP peer?
      - If yes, RPF Check Succeeds; process SA message
    - No: Is the first AS in the best path to RP = the first AS in the best path to MSDP peer?
      - If yes, RPF Check Succeeds; process SA message

# Receiving SA Messages

- RPF Check rule example cases

  - Case 1: Sending MSDP Peer = iMBGP peer

    - Is best path to RP via this MBGP peer?

  - Case 2: Sending MSDP Peer = eMBGP peer

    - Is best path to RP via this MBGP peer?

  - Case 3: Sending MSDP Peer != BGP peer

    - Is the next AS in best path to RP = AS of the sending MSDP peer?

# RPF Check Example

**RPF rule when MSDP peer == iMBGP peer**



```
                                        MBGP Table router B
A Destination              Next hop              AS path
* 172.16.0.2/32           >172.16.3.1           i
```

Who is the iMBGP peer adverting this route?
in our example 172.16.3.1

```
             MSDP Peers router B

Peer Address              State
172.16.3.1                Established
172.16.4.1                Established
```

Is the MSDP == MBGP peer?

## RPF Success!

B

4.2

**AS1**

3.1

2.1

4.1

**RP**

A

C

**Source**

```
rp {
    local {
        address 172.16.0.2;
```

➡ **MSDP SA**

••••• **MSDP/iMBGP mesh-peering**

# RPF Check Example

**RPF rule when MSDP peer == iMBGP peer**

```
                                    MBGP Table router B
                           A Destination        Next hop        AS path
                           * 172.16.0.2/32      >172.16.3.1     i
```



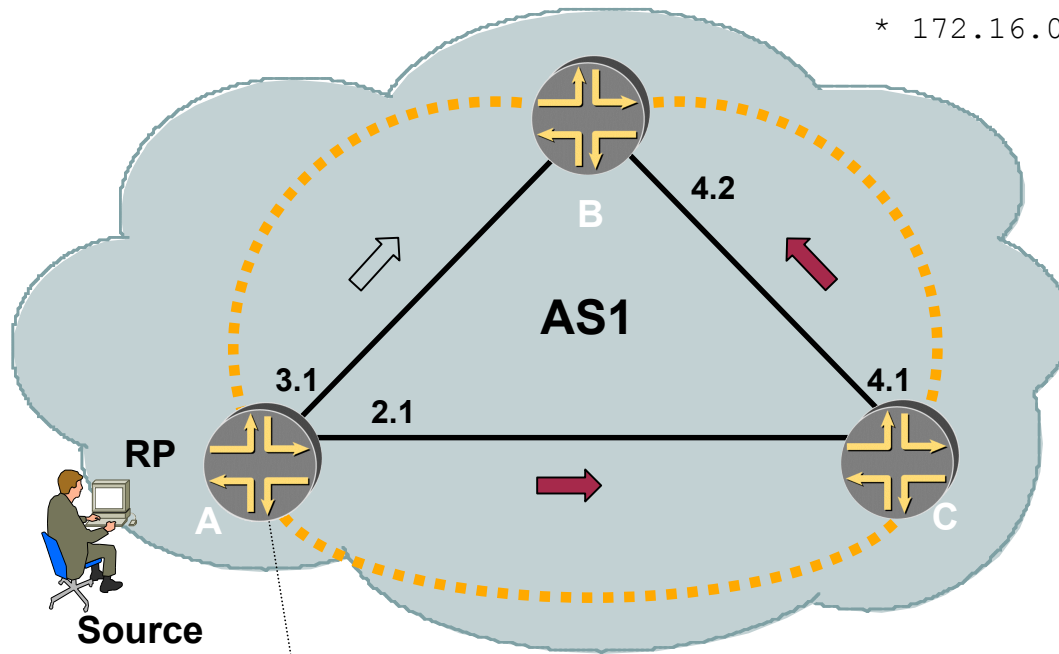Who is the iMBGP peer adverting this route?
In our example 172.16.3.1

```
                MSDP Peers router B

Peer Address          State
172.16.3.1            Established
172.16.4.1            Established
```
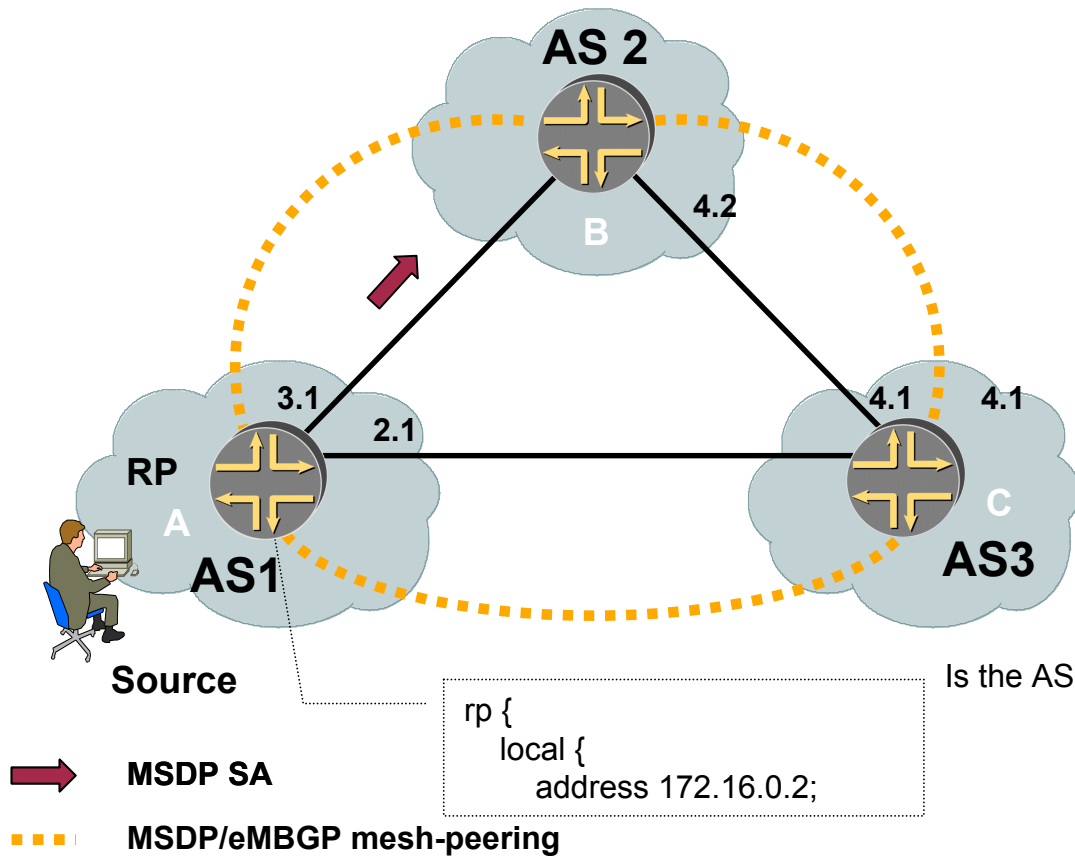
Is the MSDP == MBGP peer?

## RPF Failure!

```
rp {
    local {
        address 172.16.0.2;
```

→ **MSDP SA**

⋯ **MSDP/iMBGP mesh-peering**

# RPF Check Example

**RPF rule when MSDP == MBGP peer**



AS 2

B

4.2

3.1

2.1

RP

A

AS1

C

4.1    4.1

AS3

Source

rp {
    local {
        address 172.16.0.2;

➡ **MSDP SA**

▪▪▪ **MSDP/eMBGP mesh-peering**

```
         MSDP Peers router B

MSDP Peer          State
172.16.3.1         Established
172.16.4.1         Established

    BGP Neighbours router B

Peer               AS
172.16.3.1         1
172.16.4.1         3
```

MBGP Table

```
Destination        Next Hop        Path
* 172.16.0.2/32    >172.16.3.1      1 i
  172.16.0.2/32     172.16.4.1      3 1 i
```
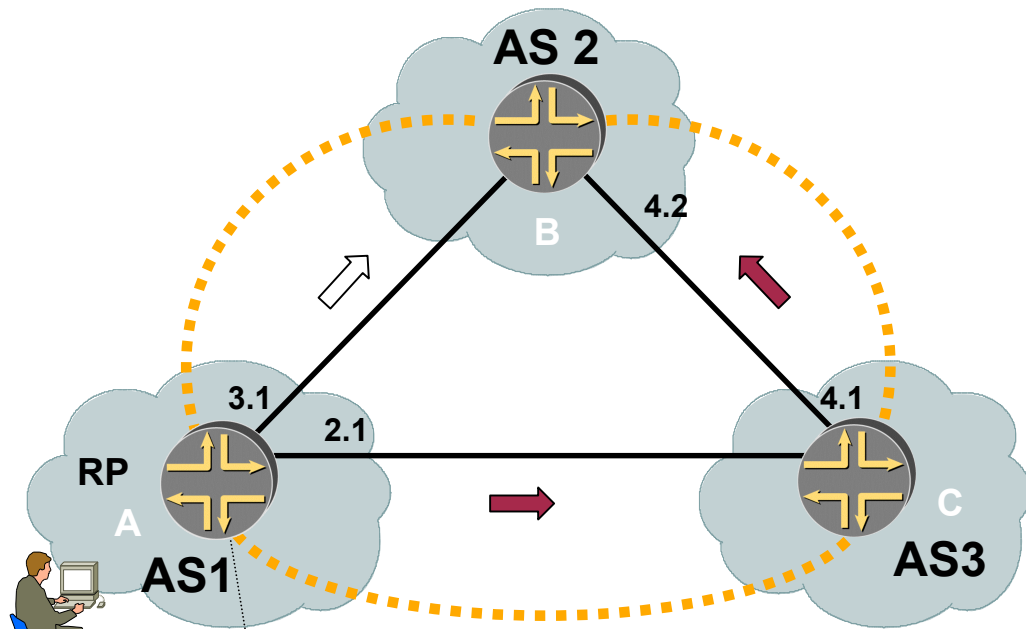
Is the AS of the sending MSDP peer  == First AS in RP route?

## RPF Success!

Who is the BGP peer adverting this route

# RPF Check Example

**RPF rule when MSDP == MBGP peer**

AS 2

B

4.2

3.1

2.1

RP

A

AS1

4.1

C

AS3

```
rp {
    local {
        address 172.16.0.2;
```

**Source**

MSDP SA

MSDP/eMBGP mesh-peering

```
MSDP Peers router B

MSDP Peer       State
172.16.3.1      Established
172.16.4.1      Established


    BGP Neighbours router B

Neighbor        AS
172.16.3.1      1
172.16.4.1      3
```

BGP Table

```
Destination          Next Hop          Path
* 172.16.0.2/32      >172.16.3.1       1 i
  172.16.0.2/32       172.16.4.1        3 1 i
```
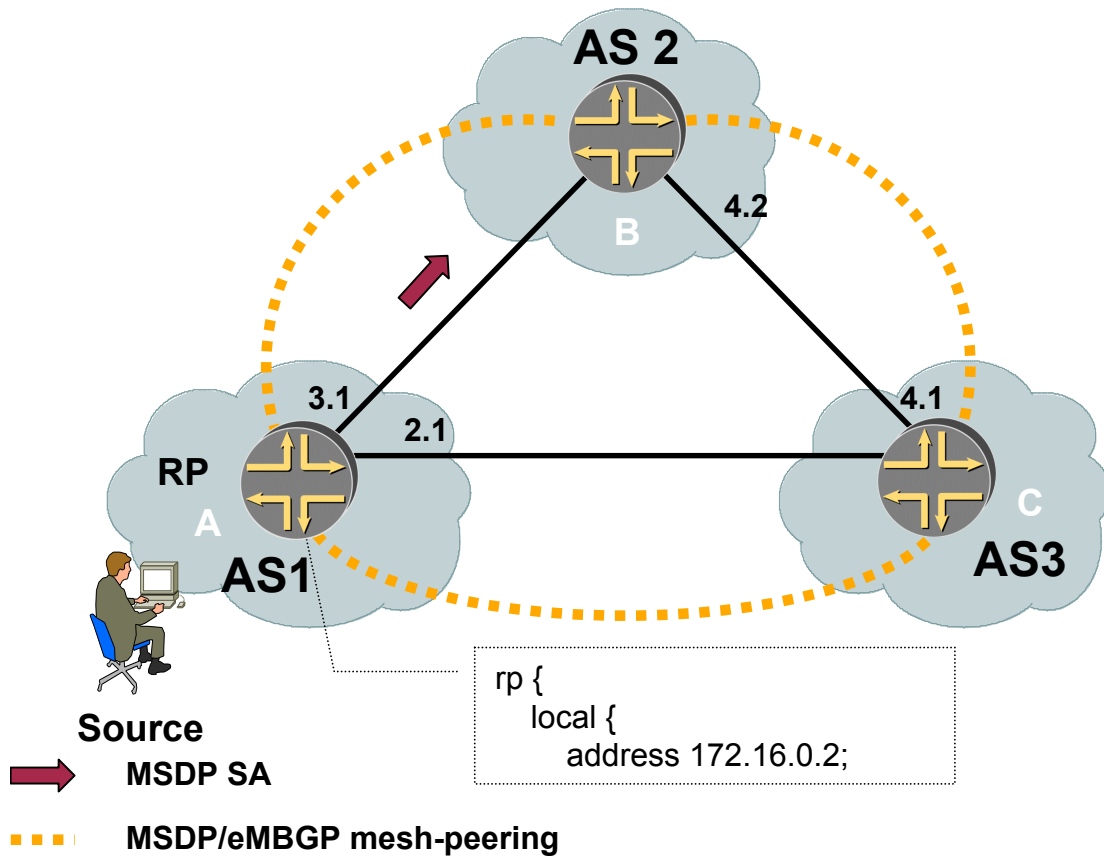
Is the AS of the sending MSDP peer  == First AS in RP route?

## RPF Failure!

Who is the BGP peer adverting this route

# RPF Check Example

**RPF rule when MSDP != MBGP peer**



AS 2

B

4.2

3.1

2.1

RP

A

AS1

4.1

C

AS3

```
rp {
    local {
        address 172.16.0.2;
```

**Source**

→ **MSDP SA**

▪ ▪ ▪ **MSDP/eMBGP mesh-peering**

```
              MSDP Peers router B

    MSDP Peer         State
    172.16.3.1        Established
    172.16.4.1        Established


            BGP Table router B

    Destination          Next Hop   Path
    *172.16.0.2/32    >172.16.3.1    1 i
     172.16.0.2/32    >172.16.4.1    3 1 I
    *172.16.4.0/24    >172.16.4.1    3 i
    *172.16.3.0/24    >172.16.3.1    1 i
```
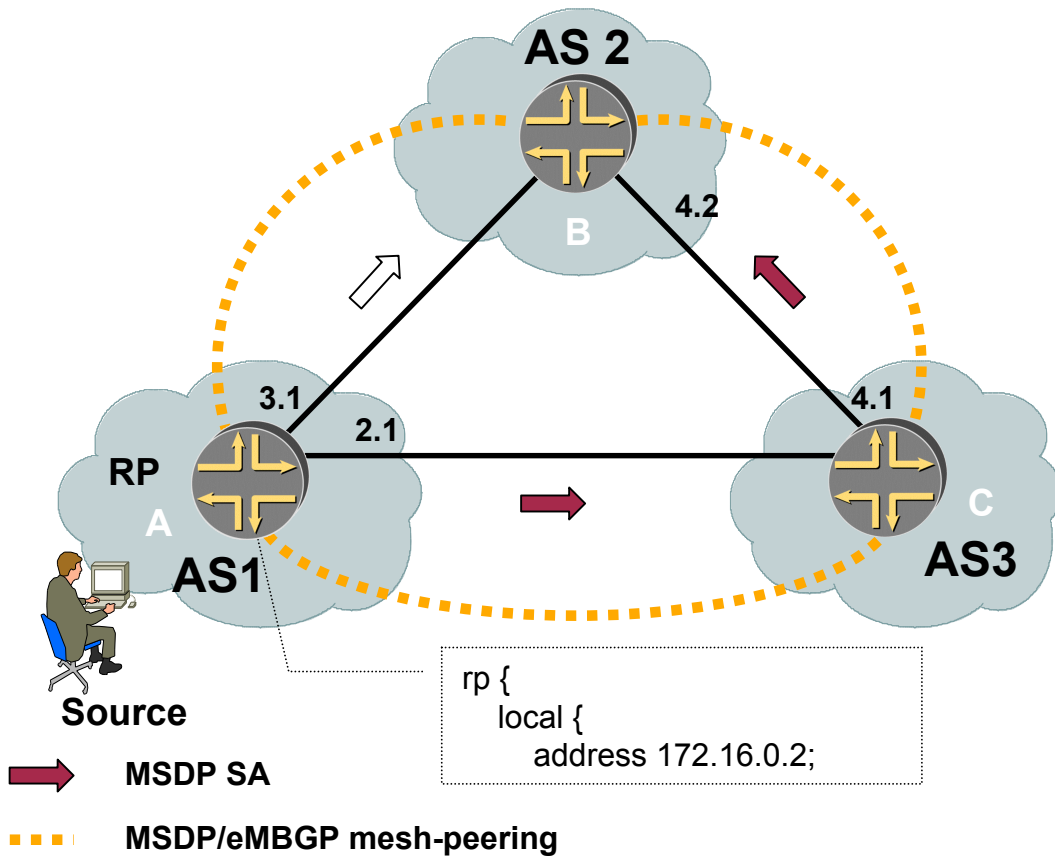
Is the first AS in the path to the MSDP peer
== First AS in best path to the RP?

## RPF Success!

# RPF Check Example

**RPF rule when MSDP != MBGP peer**



**AS 2**

B

4.2

3.1

2.1

RP

A

**AS1**

4.1

C

**AS3**

**Source**

rp {
    local {
        address 172.16.0.2;

➡ **MSDP SA**

···· **MSDP/eMBGP mesh-peering**

```
          MSDP Peers router B

   MSDP Peer        State
   172.16.3.1       Established
   172.16.4.1       Established


          BGP Table router B

 Destination      Next Hop       Path
 *172.16.0.2/32   >172.16.3.1    1 i
  172.16.0.2/32   >172.16.4.1    3 1 I
 *172.16.4.0/24   >172.16.4.1    3 i
 *172.16.3.0/24   >172.16.3.1    1 i
```
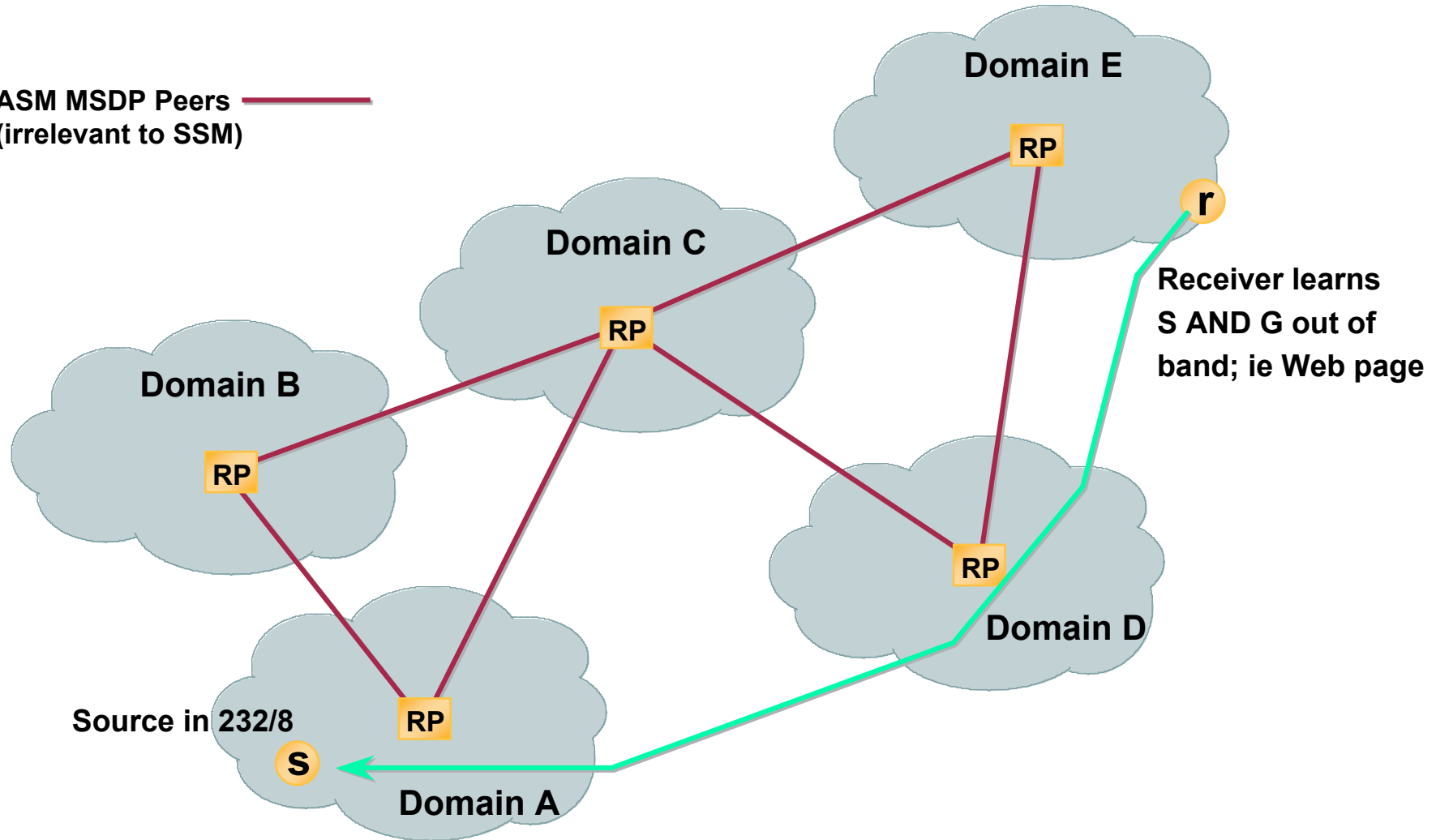
Is the first AS in the path to the MSDP peer
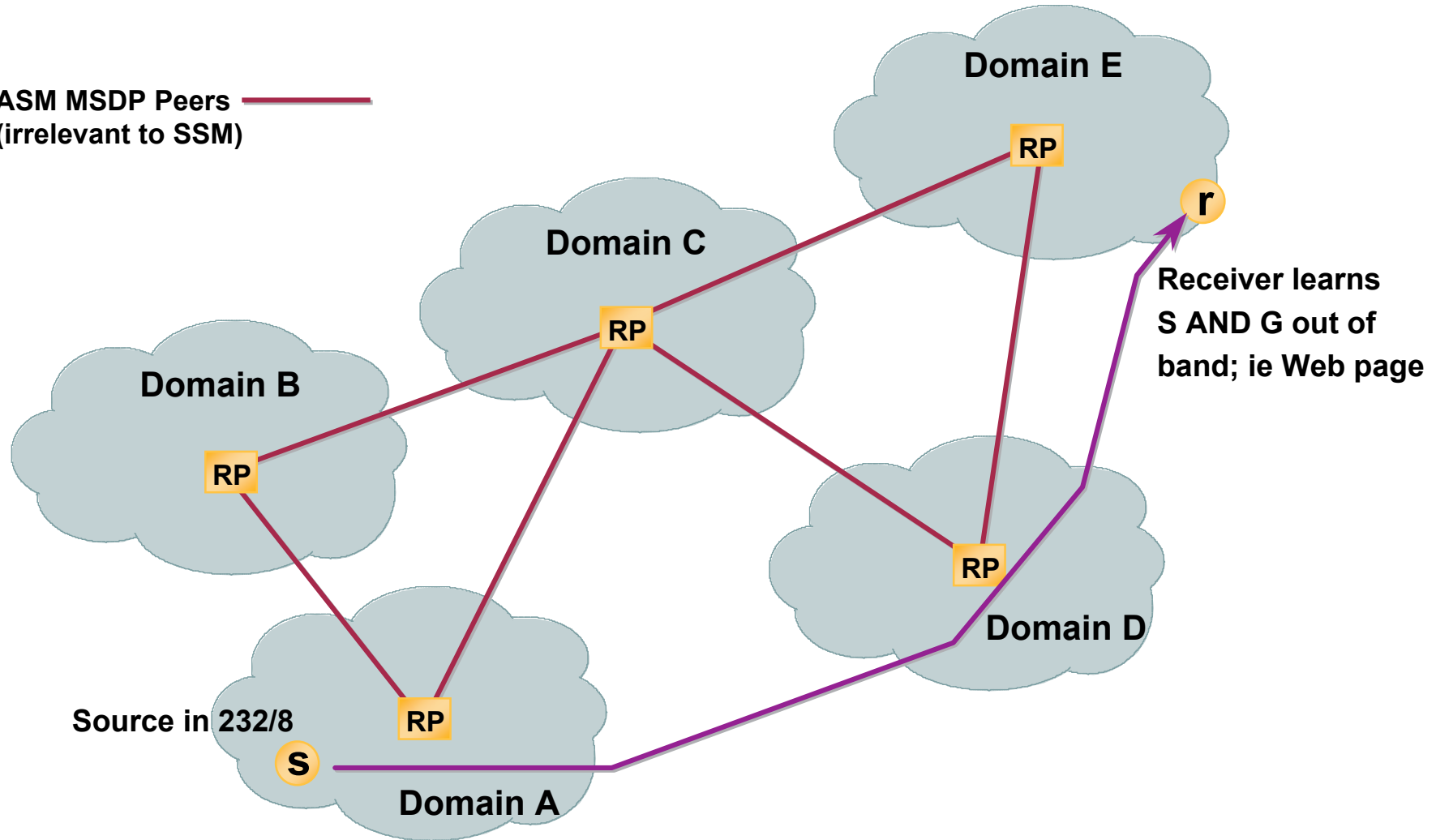== First AS in best path to the RP?

## RPF Failure!

# MSDP wrt SSM – Unnecessary!



**ASM MSDP Peers**
**(irrelevant to SSM)**

**Domain E**

**RP**

**r**

**Receiver learns**
**S AND G out of**
**band; ie Web page**

**Domain C**

**RP**

**Domain B**

**RP**

**RP**

**Domain D**

**Source in 232/8**

**RP**

**S**

**Domain A**

# MSDP wrt SSM – Unnecessary!

**ASM MSDP Peers (irrelevant to SSM)**

**Domain E**

**Domain C**

**Domain B**

RP

RP

RP

RP

RP

**r**

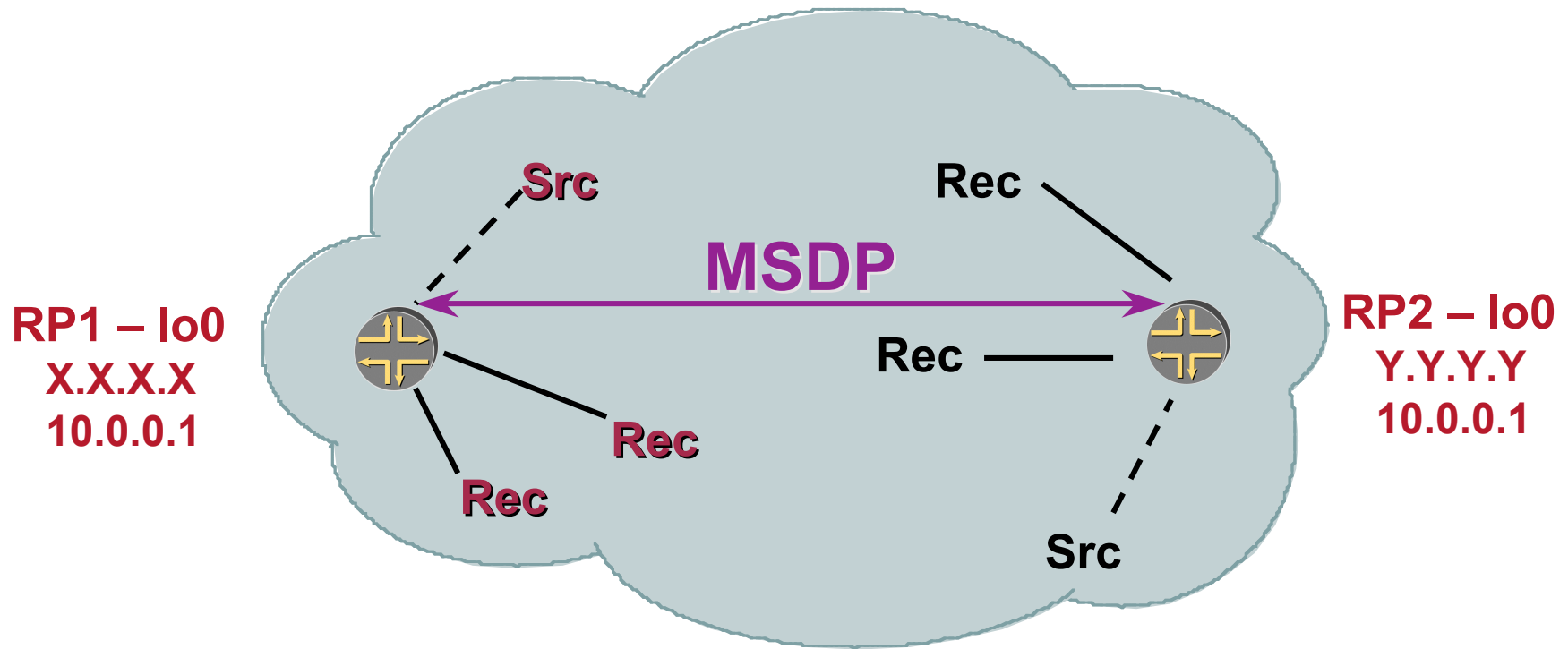Receiver learns S AND G out of band; ie Web page

**Source in 232/8**
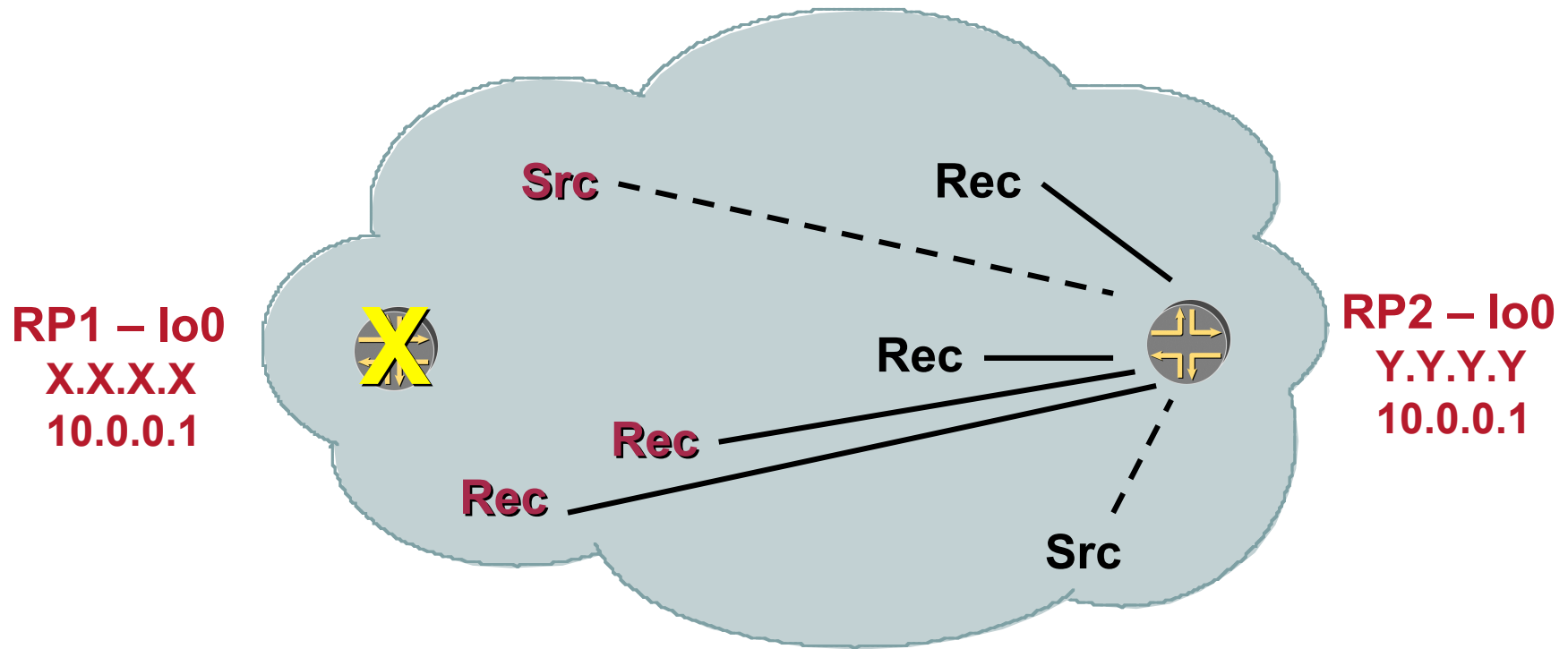
**S**

**Domain A**

**Domain D**

# MSDP Application: Anycast-RP

- Within a domain, deploy more than one RP for the same group range
- Sources from one RP are known to other RPs using MSDP
- Give each RP the same /32 IP address
- Sources and receivers use closest RP, as determined by the IGP
- Used intra-domain to provide redundancy and RP load sharing, when an RP goes down, sources and receivers are taken to new RP via unicast routing
  - Fast convergence!

# Anycast-RP

# Anycast-RP



**RP1 – lo0**
**X.X.X.X**
**10.0.0.1**

**RP2 – lo0**
**Y.Y.Y.Y**
**10.0.0.1**

Src

Rec

Rec

Rec

Rec

Src

# Agenda

- Introduction

- Multicast addressing

- Group Membership Protocol

- PIM-SM / SSM

- MSDP

- MBGP

- Summary

# MBGP—Multiprotocol BGP

- MBGP overview

- MBGP capability negotiation

- MBGP NLRI exchange

- Configuration guidelines

# MBGP

- Multiprotocol Extensions to BGP (RFC 2283).
- Tag unicast prefixes as multicast source prefixes for intra-domain mcast routing protocols to do RPF checks.
- WHY?  Allows for interdomain RPF checking where unicast and multicast paths are non-congruent.
- DO I REALLY NEED IT?
  - YES, if:
    - ISP to ISP peering
    - Multiple-homed networks
  - NO, if:
    - You are single-homed

# MBGP Overview

- MBGP: Multiprotocol BGP
  - Defined in RFC 2283 (extensions to BGP)
  - Can carry different route types for different purposes
    - Unicast
    - Multicast
  - Both route types carried in same BGP session
  - Does not propagate multicast state information
  - Same path selection and validation rules
    - AS-Path, LocalPref, MED, …

# MBGP Overview

- New multiprotocol attributes
  - MP_REACH_NLRI
  - MP_UNREACH_NLRI
- MP_REACH_NLRI and MP_UNREACH_NLRI
  - Address Family Information (AFI) = 1 (IPv4)
    - Sub-AFI = 1 (NLRI is used for unicast)
    - Sub-AFI = 2 (NLRI is used for multicast RPF check)
    - Sub-AFI = 3 (NLRI is used for both unicast and multicast RPF check)
- SAFI 1 -> RIB inet.0
- SAFI 2 -> RIB inet.2
- Multicast uses SAFI 2 routes for RPF
- Allows for different policies between multicast and unicast

# Ribs & Rib groups

- Routing Information Base (RIB)
  - Simply a routing table with a purpose
- RIB Group
  - Primary import RIB
  - Optional list of secondary import RIBs
  - Export RIB
- Well known JUNOS ribs
  - Inet.0 – Primary unicast rib
  - Inet.1 – Multicast forwarding rib
  - Inet.2 – Multicast source rib (RPF)
  - Inet.3 – MPLS rib
  - Inet.4 – MSDP SA rib

# MBGP—Capability Negotiation

- BGP routers establish BGP sessions through the OPEN message
- OPEN message contains optional parameters
- BGP session is terminated if OPEN parameters are not recognised
- New parameter: CAPABILITIES
  - Multiprotocol extension
  - Multiple routes for same destination
- Configures router to negotiate either or both NLRI
  - If neighbor configures both or subset, common NRLI is used in both directions
  - If there is no match, notification is sent and peering doesn't come up
  - If neighbor doesn't include the capability parameters in open, session backs off and reopens with no capability parameters
    - Peering comes up in unicast-only mode

# MBGP—Summary

- Solves part of inter-domain problem

  - Can exchange unicast prefixes for multicast RPF checks

  - Uses standard BGP configuration knobs

  - Permits separate unicast and multicast
    topologies if desired

- Still must use PIM to:

  - Build distribution trees

  - Actually forward multicast traffic

  - PIM-SM recommended

# Agenda

- Introduction
- Multicast addressing
- Group Membership Protocol
- Multicast Forwarding Algorithm
- PIM-SM / SSM
- MSDP
- MBGP
- Summary

# The Soup

- **IGMP -** Internet Group Management Protocol is used by hosts and routers to tell each other about group membership.

- **PIM-SM -** Protocol Independent Multicast-Sparse Mode is used to propagate forwarding state between routers.

- **SSM -** Source Specific Multicast utilizes a subset of PIM's functionality to guaranty source-only trees in the 232/8 range.

- **MBGP -** Multiprotocol Border Gateway Protocol is used to exchange routing information for interdomain RPF checking.

- **MSDP -** Multicast Source Discovery Protocol is used to exchange ASM active source information between RPs.

# Multicast Transit Design Objectives

- PIM Border Constraints
  - Confine registers within domain
  - Confine local groups
  - Confine RP announcements
  - Control SA advertisements via MSDP
- Border RPF check
  - RPF check against unicast routes to multicast sources
- MSDP RPF check
  - RPF check toward RP in received SAs

# ISP Requirements at the MIX

- Current solution: MBGP + PIM-SM + MSDP

  - Environment

    - ISPs run iMBGP and PIM-SM (internally)

    - ISPs multicast peer at a public interconnect

  - Deployment

    - Border routers run eMBGP

    - The interfaces on interconnect run PIM-SM

    - RPs' MSDP peering must be consistant with eMBGP peering

    - All peers set a common distance for eMBGP

# Thank you!

# More Information

- For more information on Multicast, please refer to the following intranet sites:
  - http://www-int.juniper.net/sales/sales_training/technology_detail.html#14