

Copyright © William C. Cheng

1

# CS551

## Multicast Routing: IGMP

Bill Cheng

<http://merlot.usc.edu/cs551-f12>

---

Computer Communications - CSCI 551

Copyright © William C. Cheng

2

### Components of the IP Multicast Architecture

service model

- host-to-router protocol (IGMP)
- multicast routing protocols (various)

hosts

routers

---

Computer Communications - CSCI 551

Copyright © William C. Cheng

3

### Internet Group Management Protocol (IGMP)

- the protocol by which hosts report their multicast group memberships to neighboring routers
- version 1, the current Internet Standard, is specified in RFC-1112
- version 2: RFC 2236
- operates over broadcast LANs and point-to-point links
- occupies similar position and role as ICMP in the TCP/IP protocol stack

---

Computer Communications - CSCI 551

Copyright © William C. Cheng

4

### Link-layer Transmission/reception

**Transmission:**

- an IP multicast packet is transmitted as a link-layer multicast, on those links that support multicast
- the link-layer destination address is determined by an algorithm specific to the type of link (next slide)

**Reception:**

- the necessary steps are taken to receive desired multicasts on a particular link, such as modifying address reception filters on LAN interfaces
- multicast routers must be able to receive all IP multicasts on a link, without knowing in advance which groups will be sent to

---

Computer Communications - CSCI 551

Copyright © William C. Cheng

5

### Mapping to Link-layer Multicast Addresses

for Ethernet and other LANs using 802 addresses:

IP multicast address: 11 10

LAN multicast address: 11 10 0000000001011100

group bit

28 bits

23 bits

for point-to-point links: no mapping needed

---

Computer Communications - CSCI 551

Copyright © William C. Cheng

6

### IGMP Version 1 Message Format

Vers	Type	Reserved	Checksum
------	------	----------	----------

Group Address

Version : 1

Type : 1 = Membership Query

2 = Membership Report

Checksum : standard IP-style checksum of the IGMP Message

Group Address : group being reported (zero in Queries)

---

Computer Communications - CSCI 551

Computer Communications - CSC1 551

## How IGMP Works

hosts: routers:

- on each link, one router is elected the "querier"
- querier periodically sends a *Membership Query* message to the all-systems group (224.0.0.1), with TTL = 1
- on receipt, hosts start random timers (between 0 and 10 seconds) for each multicast group to which they belong

Copyright © William C. Cheng

Computer Communications - CSC1 551

## IGMP Implications

hosts: routers:

- In normal case, only one report message per group present is sent in response to a query (routers need not know who all the members are, only that members exist)
- Query interval is typically 60 – 90 seconds
- IGMPv2 adds explicit leave messages
- To reduce *join latency*, when a host first joins a group, it sends one or two immediate reports (unsolicited responses), instead of waiting for a query

Copyright © William C. Cheng

Computer Communications - CSC1 551

# CS551

## Multicast Routing

Bill Cheng

<http://merlot.usc.edu/cs551-f12>

Copyright © William C. Cheng

Computer Communications - CSC1 551

## IGMP Goal

hosts: routers:

- Determine what IP multicast groups have receivers present on the LAN
- just care about some vs. zero receivers, not how many
- Approach
  - designate one router as IGMP "querier"
  - it asks all hosts
  - get at least one response per active group
  - example of *soft state* (periodically query), so occasional losses are okay

Copyright © William C. Cheng

Computer Communications - CSC1 551

## How IGMP Works (Cont...)

hosts: routers:

- when a host's timer for group G expires, it sends a *Membership Report to group G*, with TTL = 1
- other members of G hear the report and stop their timers
- routers hear *all* reports, and time out nonresponding groups

Copyright © William C. Cheng

Computer Communications - CSC1 551

## IGMP Version 2

hosts: routers:

- changes from version 1:
  - new message and procedures to reduce "leave latency"
  - standard querier election method specified
  - version and type fields merged into a single field
  - backward-compatible with version 1
- soon to appear as a Proposed Standard RFC
- widely implemented already

Copyright © William C. Cheng

Copyright © William C. Cheng

14

### Components of the IP Multicast Architecture

service model

- host-to-router protocol (IGMP)
- multicast routing protocols (various)

routers

hosts

- ↳ Multicast service model makes it hard to locate receivers
- ↳ anonymity
- ↳ dynamic join/leave
- ↳ Options so far (not very efficient)
  - ↳ flood data packets to entire network, or
  - ↳ tell routers about all possible groups and receivers so they can create routes (trees)

Computer Communications - CSC1 551

Copyright © William C. Cheng

15

### Early Routing Techniques

- ↳ **Flood and prune**
  - begin by flooding traffic to entire network
  - prune branches with no receivers
  - unwanted state where there are no receivers
  - examples: DVMRP, PIM-DM
- ↳ Link-state multicast protocols
  - routers advertise groups for which they have receivers to entire network
  - compute trees on demand
  - unwanted state where there are no senders
  - examples: MOSPF

Computer Communications - CSC1 551

Copyright © William C. Cheng

16

### Source-based Trees

output link determined from input link, multicast address, and source address

### Rendezvous Options

- ↳ Specify **rendezvous** (or meeting place) to which sources send initial packets, and receivers join; requires mapping between multicast group address and meeting place
- ↳ examples: CBT, PIM-SM

Computer Communications - CSC1 551

Copyright © William C. Cheng

17

### Multicast Tree Taxonomy

- ↳ Multicast routing can build different types of distribution trees
- ↳ **Source-based trees**
  - separate shortest path tree (SPT) for each sender
  - can have multiple senders per group
  - examples: DVMRP, MOSPF, PIM-DM, PIM-SM
- ↳ **Shared trees**
  - single tree shared by all members
  - shared tree rooted at group core/rendezvous point
  - examples: CBT, PIM-SM

Computer Communications - CSC1 551

Copyright © William C. Cheng

18

### Source-based Trees

output link determined from input link, multicast address, and source address

Computer Communications - CSC1 551

Copyright © William C. Cheng

24

<http://merlot.usc.edu/cs551-f12>

Bill Cheng

[Deering88b]

**DVMRP & MOSPF**

**CS551**

Computer Communications - CSCI 551

Copyright © William C. Cheng

22

Who Can Send?

- Anyone (Deering's service model)
  - model used by most multicast applications
- Single-source
  - only one node can send (others must make their own group)
- EXPRESS [Holbrook99a]

Computer Communications - CSCI 551

Copyright © William C. Cheng

20

Shared v.s. Source-Based Trees

- Source-based trees
  - shortest path trees - low delay, better load distribution
  - more state at routers (per-source state)
  - efficient for dense-area multicast
- Shared trees
  - higher delay (bounded by factor of 2), traffic concentration
  - per-group state at routers
  - efficient for sparse-area multicast

Computer Communications - CSCI 551

Copyright © William C. Cheng

23

Multicast Status

- MBone exists
  - moderately widely used in research
  - but not always stable
  - multi-domain routing is hard, need to coordinate people and often people don't talk about experimental services
- Some commercial use (applications)
  - but very little ISP support
  - concerned about how to charge, and potential over-use
- Multicast widely used on LANs
  - e.g., Google, Inktomi use it for load balancing

Computer Communications - CSCI 551

Copyright © William C. Cheng

21

Protocol Taxonomy

- DVMRP - source-based trees
- MOSPF - source-based trees
- PIM - shared and source-based trees

Computer Communications - CSCI 551

Copyright © William C. Cheng

19

A Shared Tree

output link determined from input link & multicast address

Computer Communications - CSCI 551

Copyright © William C. Cheng

25

### Key Ideas

- ↳ Lays foundation for IP multicast
- ↳ defines IP multicast service model
  - e.g., best effort, packet based, anonymous groups
  - compare to ISIS with explicit group membership, guaranteed ordering (partial or total ordering)
- ↳ Several algorithms
  - = extended/bridged LANs
  - = distance-vector extensions (DVMRP)
  - = link-state extensions (MOSP)
- ↳ Cost analysis

### Characterizing Groups

- ↳ Pervasive or dense
  - = most LANs have a receiver
- ↳ Sparse
  - = few LANs have receivers
- ↳ Local
  - = inside a single administrative domain

Computer Communications - CSC1 551

Copyright © William C. Cheng

28

### Multicast Forwarding

- ↳ A DVMRP router forwards a packet if
  - = **Reverse Path Forwarding (RPF)**
  - the packet arrived from the link used to reach the source
  - take advantage of what is available from unicast
  - = similar (but not quite the same) to flooding each packet once
  - if downstream links have not pruned the tree

Computer Communications - CSC1 551

Copyright © William C. Cheng

30

### Phase 1: Flood Using Truncated Broadcast

This router knows it has no group members on its LAN, so it does not broadcast over its LAN

Computer Communications - CSC1 551

Copyright © William C. Cheng

27

### Distance-vector Multicast Routing Protocol (DVMRP)

- ↳ Basic idea: **flood and prune**
  - = flood: send information about new **sources** everywhere
  - = prune: routers will tell us if they don't have receivers
- ↳ Routing information is soft state; periodically re-flood (and prune) to refresh this information
  - = if no refresh, then the information goes away
  - = easy fault recovery
- ↳ DVMRP consists of two major components:
  - = a conventional distance-vector routing protocol (like RIP)
  - = a protocol for determining how to forward multicast packets, based on the routing table

Computer Communications - CSC1 551

Copyright © William C. Cheng

29

### Example Topology

Computer Communications - CSC1 551

Copyright © William C. Cheng

36

### Link-State Multicast Routing

- Basic idea: treat group members (receivers) as new links
- flood information about them to everyone in LSA message (just like LSA routing)
- Realized as MOSPF (Multicast Open Shortest-Path First)
- add-on to OSPF
- each router indicates groups for which there are directly-connected members
- link-state advertisements augmented with multicast group addresses to which local members have joined
- link-state routing algorithm augmented to compute shortest-path distribution tree from any source to any set of destinations

Computer Communications - CSC1 551

Copyright © William C. Cheng

34

### Sending Data in DVMRP

- Data packets are sent on all branches of the tree
- send on all interfaces except the one they came in on
- RPF (Reverse Path Forwarding) check:
  - drop packets that arrive on incorrect interfaces (i.e., not from the unicast direction to the sending host)
  - why? suppress errant packets

Computer Communications - CSC1 551

Copyright © William C. Cheng

32

### Phase 3: Graft

Computer Communications - CSC1 551

Copyright © William C. Cheng

35

### DVMRP Pros and Cons

- Pros
  - simple
  - works well with many receivers
  - overhead is per-sender; receivers are passive
- Cons
  - works poorly with many groups
  - every send in every group floods the nets
  - works poorly with sparse groups
  - flood data everywhere and then prune back, expensive
  - if only needed at some places

Computer Communications - CSC1 551

Copyright © William C. Cheng

33

### Phase 4: Steady State

Computer Communications - CSC1 551

Copyright © William C. Cheng

31

### Phase 2: Prune

Computer Communications - CSC1 551

Copyright © William C. Cheng

42

<http://merit.utor.usc.edu/cs51-f12>

Bill Cheng

[Deering96a]

# Multicast Routing: PIM

## CS51

Computer Communications - CSCI 551

Copyright © William C. Cheng

40

Link state advertisement (T) with new membership (R3) may require incremental computation and addition of interface to outgoing interface list (Z)

Overhead: all these inactive nodes must keep multicast states

Computer Communications - CSCI 551

Copyright © William C. Cheng

38

Z has network map, including membership at X and Y

Z computes shortest path tree from S1 to X and Y (lazily, when it gets a data packet for group) W, Q, R, each do same thing as data arrives at them

Computer Communications - CSCI 551

Copyright © William C. Cheng

41

## MOSPF Pros and Cons

Computer Communications - CSCI 551

- ↳ Pros
  - = simple add on to OSPF
  - = works well with many senders
  - = no per-sender state
- ↳ Cons
  - = works poorly with many receivers
  - = per-receiver costs
  - = works poorly with sparse groups
  - = lots of information goes places that don't want it
  - = works poorly with large domains
  - = link-state scales with respect to number of links
  - = many links causes frequent changes

Copyright © William C. Cheng

39

Link state advertisement with new topology may require re-computation of tree and forwarding entry (only Z and W send new LSA messages, but all on path recompute)

Computer Communications - CSCI 551

Copyright © William C. Cheng

37

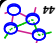
Link state: Each router floods membership

Multicast: add membership information to "link state"

Each router computes multicast tree for each active source, builds forwarding entry with outgoing interface list

Computer Communications - CSCI 551

Copyright © William C. Cheng



## Key Ideas

- Want a multicast routing protocol that works well with sparse users
- Use a single shared tree; fix one host as rendezvous point

How do we solve the problem?


- = *shared trees*
- = establish a meeting place: center, core or rendezvous point
- trade-off: shared trees can be inefficient

With source-based trees senders and receivers meet by:

- = flooding and pruning
- = LS distribution of group and receiver state

Computer Communications - CSC1 551

Copyright © William C. Cheng

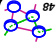


## PIM Terminology

- *incoming interface (iii)*: interface from which multicast packet is accepted and forwarded
- *outgoing interface list (oif list)*: interfaces out of which multicast packets are forwarded
- *Rendezvous Point (RP)*: used in PIM as alternative to broadcast
- *Designated Router (DR)*: one router per multi-access LAN elected to track group membership, and then join/prune accordingly

Computer Communications - CSC1 551

Copyright © William C. Cheng




## How to Build A Shared Tree

- Quite easy if you have a RP
- simply send a message towards the RP
- use the *unicast* routing table to get there
- add links to the tree as you go
- stop if you get to a router that's already in the tree
- get *reverse* shortest path to RP

Computer Communications - CSC1 551

Copyright © William C. Cheng




## PIM Protocol Overview

- Basic protocol steps
  - routers with local members *Join* toward *Rendezvous Point (RP)* to join *Shared Tree*
  - routers with local sources encapsulate data in *Register* messages to RP
  - routers with local members may initiate data-driven switch to *source-specific shortest path trees*
- Soft state: periodic state-driven refreshes, time-out idle state
- See PIM v.2 Specification (RFC2362)

Computer Communications - CSC1 551

Copyright © William C. Cheng




## PIM Terminology (Cont...)

- *Shared tree*: reverse-shortest-path tree rooted at RP
- *Source-specific tree*: reverse-shortest-path tree rooted at source. Also referred to as *Shortest Path Tree (SPT)*
- *Entry*: Multicast forwarding state for a particular source-specific or Shared tree
- *Reverse-path forwarding (RPF) check*: checks if a packet arrived on the interface used to reach the source of the packet

Computer Communications - CSC1 551

Copyright © William C. Cheng



## Key Ideas

- Want a multicast routing protocol that works well with sparse users
- Use a single shared tree; fix one host as rendezvous point

How do we solve the problem?

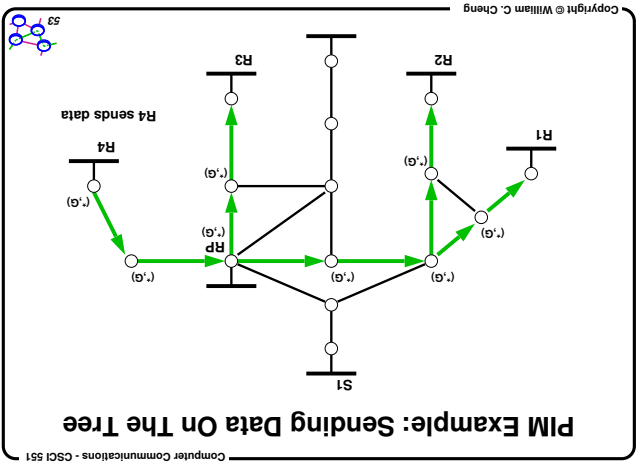
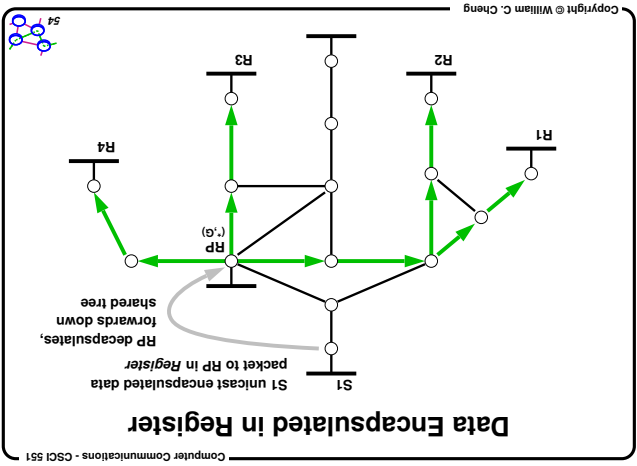
- = LS distribution of group and receiver state
- = flooding and pruning
- trade-off: shared trees can be inefficient

With source-based trees senders and receivers meet by:

- = *shared trees*
- = establish a meeting place: center, core or rendezvous point
- trade-off: shared trees can be inefficient

Computer Communications - CSC1 551



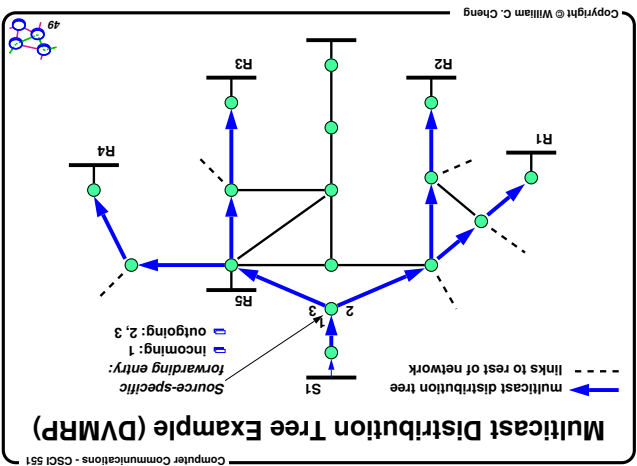
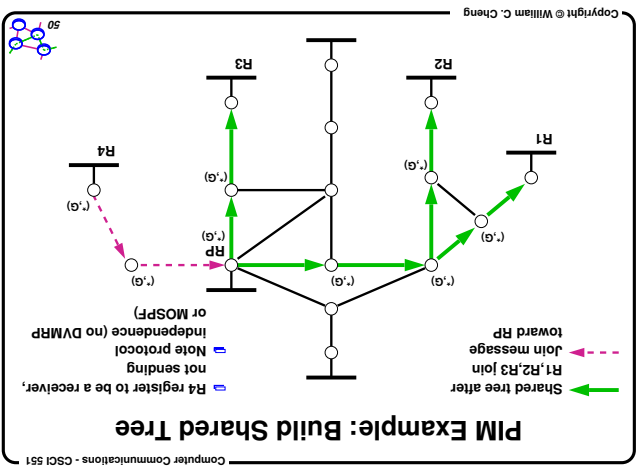


**PIM: Sending Data**

- If you are on the tree, you just send it as with other multicast protocols
- it follows the multicast tree
- If you are not on the tree (say, you are a sender but not a group member), the packet is tunneled to the RP that sends it
- this makes central placement of RP important

**How Do Routers Know RPs?**

- RP information is flooded through the network
- cannot avoid flooding something!
- but flooding control information is OK
- If there are multiple RPs, each router uses the same hash function to pick a unique RP for the group
- hash based on group address



Discussion

- Context
  - Interest in multicast motivated by audio and video apps
  - PIM was part of a large body of work in multicast routing
- Impact
  - Improved scalability compared to DVMRP and MOSPF
  - standardize and implemented
- Multicast status
  - PIM is an intra-domain routing protocol
  - RP flooding limits scalability
  - subsequent work developed inter-domain multicast protocols
  - BGMP & MSDP
- management of multicast is hard
  - multicast deployment deadlock

Computer Communications - CSC1 551

Forward Packets on "Longest Match" Entry

- Source-specific entry is "longer match" for source S1 than is Shared tree entry
- Shared tree distribution tree

Computer Communications - CSC1 551

RP May Ask High-rate Src to Join (Cont..)

Computer Communications - CSC1 551

Prune S1 off Shared Tree to Avoid Duplicates

- S1 distribution tree
- Shared tree
- Prune S1 off shared tree where IIF of S1 and RP entries differ

Computer Communications - CSC1 551

Build Source-specific Distribution Tree

- RP distribution tree
- Join messages toward S1
- Build source-specific tree for high data rate source

Computer Communications - CSC1 551

RP May Ask High-rate Src to Join

Computer Communications - CSC1 551