

CS 268: IP Multicast Routing

Ion Stoica
April 8, 2003

Motivation

- Many applications requires one-to-many communication
 - E.g., video/audio conferencing, news dissemination, file updates, etc.
- Using unicast to replicate packets not efficient → thus, IP multicast needed
 - What about the e2e arguments?

istoica@cs.berkeley.edu

2

Semantic

- Open group semantic
 - A group is identified by a location-independent address
 - Any source (not necessary in the group) can multicast to all members in a group
- Advantages:
 - Query an object/service when its location is not known
- Disadvantage
 - Difficult to protect against unauthorized listeners

istoica@cs.berkeley.edu

3

Problem

- Multicast delivery widely available on individual LANs
 - Example: Ethernet multicast
- But not across interconnection of LANs
 - I.e., can't do Internet multicast

istoica@cs.berkeley.edu

4

Three Approaches [Deering & Cheriton '89]

- Single spanning-tree (SST)
- Distance-vector multicast (DVM)
- Link-state multicast (LSM)
- Also: Sketches hierarchical multicast

istoica@cs.berkeley.edu

5

Multicast Service Model

- Built around the notion of group of hosts:
 - Senders and receivers need not know about each other
- Sender simply sends packets to “logical” group address
- No restriction on number or location of receivers
 - Applications may impose limits
- Normal, best-effort delivery semantics of IP
 - Same recovery mechanisms as unicast

istoica@cs.berkeley.edu

6

Multicast Service Model (cont'd)

- Dynamic membership
 - Hosts can join/leave at will
- No synchronization or negotiation
 - Can be implemented a higher layer if desired

istoica@cs.berkeley.edu

7

Key Design Goals

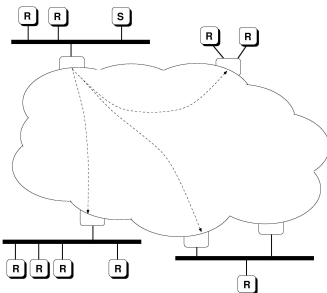
1. Delivery efficiency as good as unicast
2. Low join latency
3. Low leave latency

istoica@cs.berkeley.edu

8

Network Model

- Interconnected LANs
- LANs support link-level multicast
- Map globally unique multicast address to LAN-based multicast address (LAN-specific algorithm)



istoica@cs.berkeley.edu

9

Distance Vector Multicast Routing

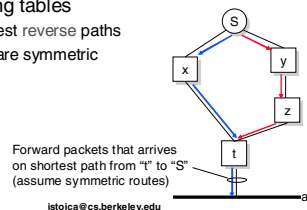
- An elegant extension to DV routing
- Use shortest path DV routes to determine if link is on the source-rooted spanning tree

istoica@cs.berkeley.edu

10

Reverse Path Flooding (RPF)

- A router forwards a broadcast packet from source (S) iff it arrives via the shortest path from the router back to S
- Packet is replicated out all but the incoming interface
- Reverse shortest paths easy to compute → just use info in DV routing tables
 - DV gives shortest reverse paths
 - Works if costs are symmetric

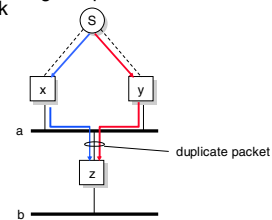


istoica@cs.berkeley.edu

11

Problem

- Flooding can cause a given packet to be sent multiple times over the same link



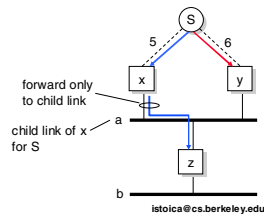
istoica@cs.berkeley.edu

12

- Solution: Reverse Path Broadcasting

Reverse Path Broadcasting (RPB)

- Basic idea: forward a packet from S only on child links for S
- Child link of router R for source S: link that has R as parent on the shortest path from the link to S



13

Identify Child Links

- Routing updates identify parent
- Since distances are known, each router can easily figure out if it's the parent for a given link
- In case of tie, lower address wins

istoica@cs.berkeley.edu

14

Problem

- This is still a broadcast algorithm – the traffic goes everywhere
- First order solution: Truncated RPB

istoica@cs.berkeley.edu

15

Truncated RPB

- Don't forward traffic onto network with no receivers
 - Identify leaves
 - Detect group membership in leaf

istoica@cs.berkeley.edu

16

Reverse Path Multicast (RPM)

- Prune back transmission so that only absolutely necessary links carry traffic
- Use on-demand pruning so that router group state scales with number of active groups (not all groups)

istoica@cs.berkeley.edu

17

Basic RPM Idea

- Prune (Source,Group) at leaf if no members
 - Send Non-Membership Report (NMR) up tree
- If all children of router R prune (S,G)
 - Propagate prune for (S,G) to parent R
- On timeout:
 - Prune dropped
 - Flow is reinstated
 - Down stream routers re-prune
- Note: again a soft-state approach

istoica@cs.berkeley.edu

18

Details

- How to pick prune timers?
 - Too long → large join time
 - Too short → high control overhead
- What do you do when a member of a group (re)joins?
 - Issue prune-cancellation message (grafts)
- Both NRM and graft messages are positively acknowledged (why?)

istoica@cs.berkeley.edu

19

RMP Scaling

- State requirements:
 - $O(\text{Sources} \times \text{Groups})$ active state
- How to get better scaling?
 - Hierarchical Multicast
 - Core-based Trees

istoica@cs.berkeley.edu

20

Core Based Trees (CBT)

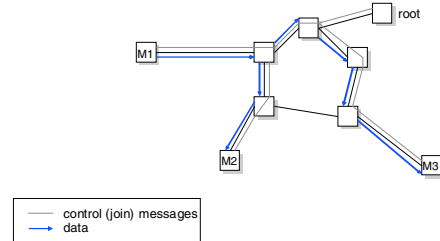
- Ballardie, Francis, and Crowcroft,
 - "Core Based Trees (CBT): An Architecture for Scalable Inter-Domain Multicast Routing", SIGCOMM 93
- Similar to Deering's Single-Spanning Tree
- Unicast packet to core and bounce it back to multicast group
- Tree construction is receiver-based
 - One tree per group
 - Only nodes on tree involved
- Reduce routing table state from $O(S \times G)$ to $O(G)$

istoica@cs.berkeley.edu

21

Example

- Group members: M1, M2, M3
- M1 sends data



istoica@cs.berkeley.edu

22

Disadvantages

- Sub-optimal delay
- Single point of failure
 - Core goes out and everything lost until error recovery elects a new core
- Small, local groups with non-local core
 - Need good core selection
 - Optimal choice (computing topological center) is NP complete

istoica@cs.berkeley.edu

23

IP Multicast Revisited

- Despite many years of research and many compelling applications, and despite the fact that the many of today routers implement IP multicast, this is still not widely deployed
- Why?

istoica@cs.berkeley.edu

24

Possible Explanations [Holbrook & Cheriton '99]

- Violates ISP input-rate-based billing model
 - No incentive for ISPs to enable multicast!
- No indication of group size (again needed for billing)
- Hard to implement sender control → any node can send to the group (remember open group semantic?)
- Multicast address scarcity

istoica@cs.berkeley.edu

25

Solution: EXPRESS

- Limit to single source group
- Use both source and destination IP fields to define a group
 - Each source can allocate 16 millions channels (i.e., multicast groups)
- Use RPM algorithm
- Add a counting mechanism
 - Use a recursive CountQuery message
- Use a session rely approach to implement multiple source multicast trees

istoica@cs.berkeley.edu

26

Summary

- Deering's DV-RMP an elegant extension of DV routing
- CBT addresses some of the DV-RMP scalability concerns but is sub-optimal and less robust
- Protocol Independent Multicast (PIM)
 - Sparse mode similar to CBT
 - Dense mode similar to DV-RMP
- Lesson: economic incentives plays a major role in deploying a technical solution
 - See EXPRESS work

istoica@cs.berkeley.edu

27