

Lecture 12: Multicast Routing

CSE 123: Computer Networks
Stefan Savage



Last time

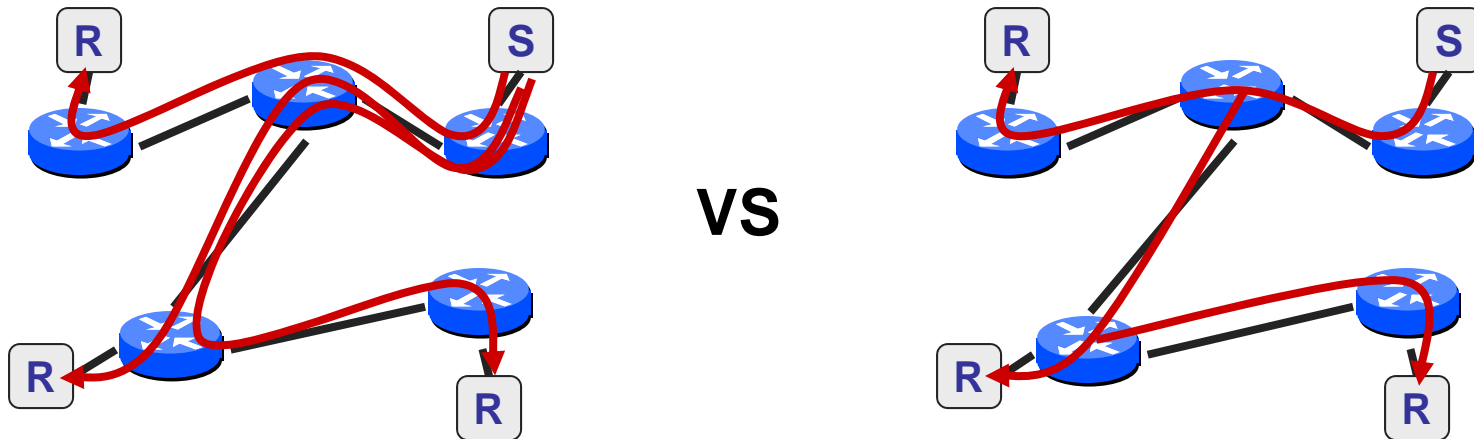
- Intra-domain routing won't work on the Internet
 - ◆ DV and LS do not scale
 - ◆ How do we pick a “global” metric?
- We talked about inter-domain routing
 - ◆ Mechanism and policy
 - ◆ BGP: Use path vector algorithm instead of DV or LS
 - ◆ Path vector looks at reachability info only
 - » Ignores metrics
 - » Route selection is based on local policy

Today: Multicast routing

- Multicast service model
- Host interface
- Host-router interactions (IGMP)
- Multicast Routing
 - ◆ Distance Vector
 - ◆ Link State
 - ◆ Shared tree
- Limiters
 - ◆ Deployment issues
 - ◆ Inter-domain routing
 - ◆ Operational/Economic issues

Motivation

- Efficient delivery to multiple destinations (e.g. video broadcast)



- Network-layer support for one-to-many addressing
 - ◆ Publish/subscribe communications model
 - ◆ Don't need to know destinations

IP Multicast service model

- **Communications based on groups**
 - ◆ Special IP addresses (Class D in IPv4) represent “multicast groups”
 - ◆ Anyone can join group to receive packets
 - ◆ Anyone can send to group
 - » Sender need not be part of group
 - ◆ Dynamic group membership – can join and leave at will
- **Unreliable datagram service**
 - ◆ Extension to unicast IP
 - ◆ Group membership not visible to hosts
 - ◆ No synchronization
- **Explicit *scoping* to limit spread of packets**

Elements of IP Multicast

- **(1) Host interface**
 - ◆ Application visible multicast API
 - ◆ Multicast addressing
 - ◆ Link-layer mapping
- **(2) Host-Router interface**
 - ◆ IGMP
- **(3) Router-Router interface**
 - ◆ Multicast routing protocols

(1) Host interface

- Senders (not much new)
 - ◆ Set TTL on multicast packets to limit “scope”
 - » Scope can be administratively limited on per-group basis
 - ◆ Send packets to *multicast address*, represents a group
 - ◆ Unreliable transport (no acknowledgements)
- Receivers (two new interfaces)
 - ◆ Join multicast group (group address)
 - ◆ Leave multicast group (group address)
 - ◆ Typically implemented as a socket option in most networking API

Multicast addressing

- Special address range:
 - ◆ Class D (3 MSBs set to 1) 224.0.0.1- 239.255.255.255
 - ◆ Reserved by IANA for multicast
- Which address to use for a new group?
 - ◆ No standard
 - ◆ Global random selection
 - ◆ Per-domain addressing
- Which address to use to join an existing group?
 - ◆ No standard
 - ◆ Separate address distribution protocol (may use multicast)

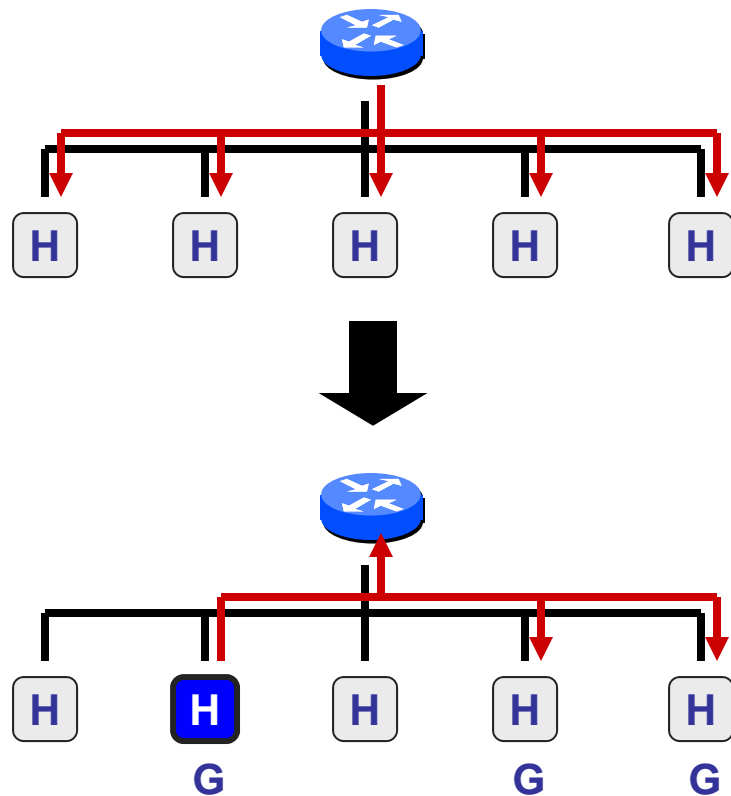
Link-layer multicast

- Many link-layers protocols have multicast capability
 - ◆ Ethernet, FDDI
- Translate IP Multicast address into LL address
 - ◆ E.g. Map 28 bits of IP Multicast address in 23bit Ethernet Multicast addresses
 - ◆ Senders send and receive on link-layer Multicast addresses
 - ◆ Routers must listen on all possible link-layer Multicast addresses
- Not an issue for point-to-point links

Internet Group Management Protocol (IGMP)

- **(2) Host-router interface**
- **Goal: Communicate group membership between hosts and routers**
- **Soft-state protocol**
 - ◆ Hosts explicitly inform their router about membership
 - ◆ Must periodically refresh membership report
 - ◆ Routers implicitly timeout groups that aren't refreshed
 - ◆ Why isn't explicit "leave group" message sufficient?
- Implemented in most of today's routers and switches

How IGMP works (roughly)



- Router broadcasts *membership query* to 224.0.0.1 (all-systems group) with TTL=1
- Hosts start random timer (0-10 sec) for each group they have joined
- When a host's timer expires for group G, send *membership report* to group G, with TTL=1
- When a member of G hears a report, they reset their timer for G
- Router times out groups that are not "refreshed" by some host's report

Multicast routing

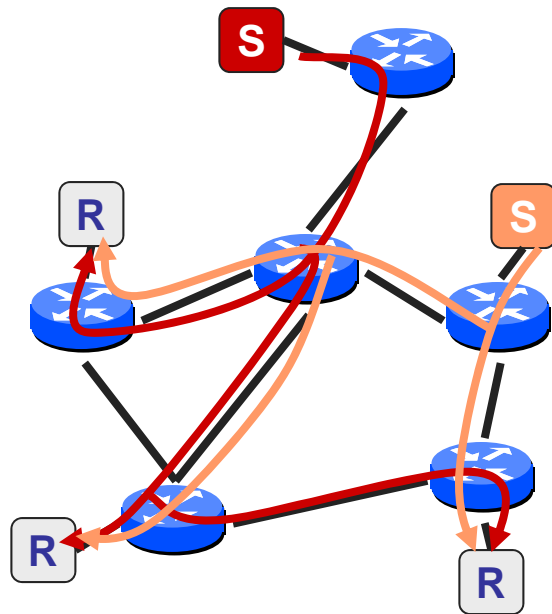
- **(3) Router-router interface**
- **Goal: Build distribution tree for multicast packets**
 - ◆ Efficient tree (ideally, shortest path)
 - ◆ Low join/leave latency
- **Several approaches**
 - ◆ Distance Vector/Link State
 - » Leverage existing unicast routing protocols
 - ◆ Shared tree
 - » Unicast/multicast hybrids

Multicast routing taxonomy

- **Source-based tree**
 - ◆ *Separate shortest path tree for each source*
 - ◆ **Flood and prune** (DVMRP, PIM-DM)
 - » Send multicast traffic everywhere
 - » Prune edges that are not actively subscribed to group
 - ◆ **Link-state** (MOSPF)
 - » Routers flood groups they would like to receive
 - » Compute shortest-path trees on demand
- **Shared tree** (PIM-SM)
 - ◆ *Single distributed tree shared among all sources*
 - ◆ Specify rendezvous point (RP) for group
 - ◆ Senders send packets to RP, receivers join at RP
 - ◆ RP multicasts to receivers; Fix-up tree for optimization

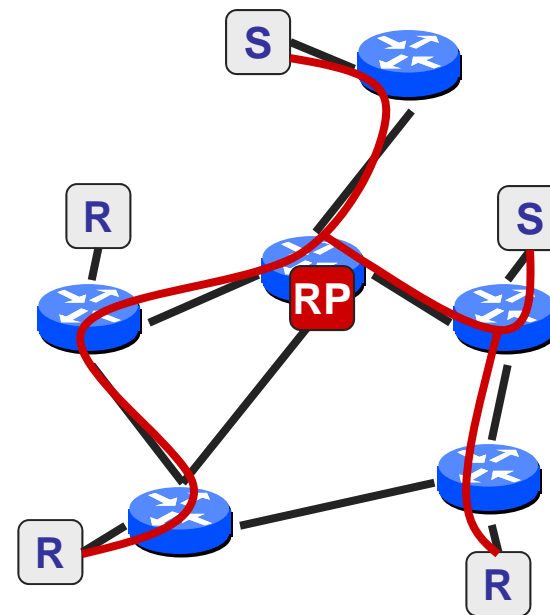
Source-based vs Shared

Source-based tree



- Efficient trees; low delay, even load
- Per-source state in routers (S,G)

Shared-tree



- Higher delay, skewed load
- Per-group state only (G)

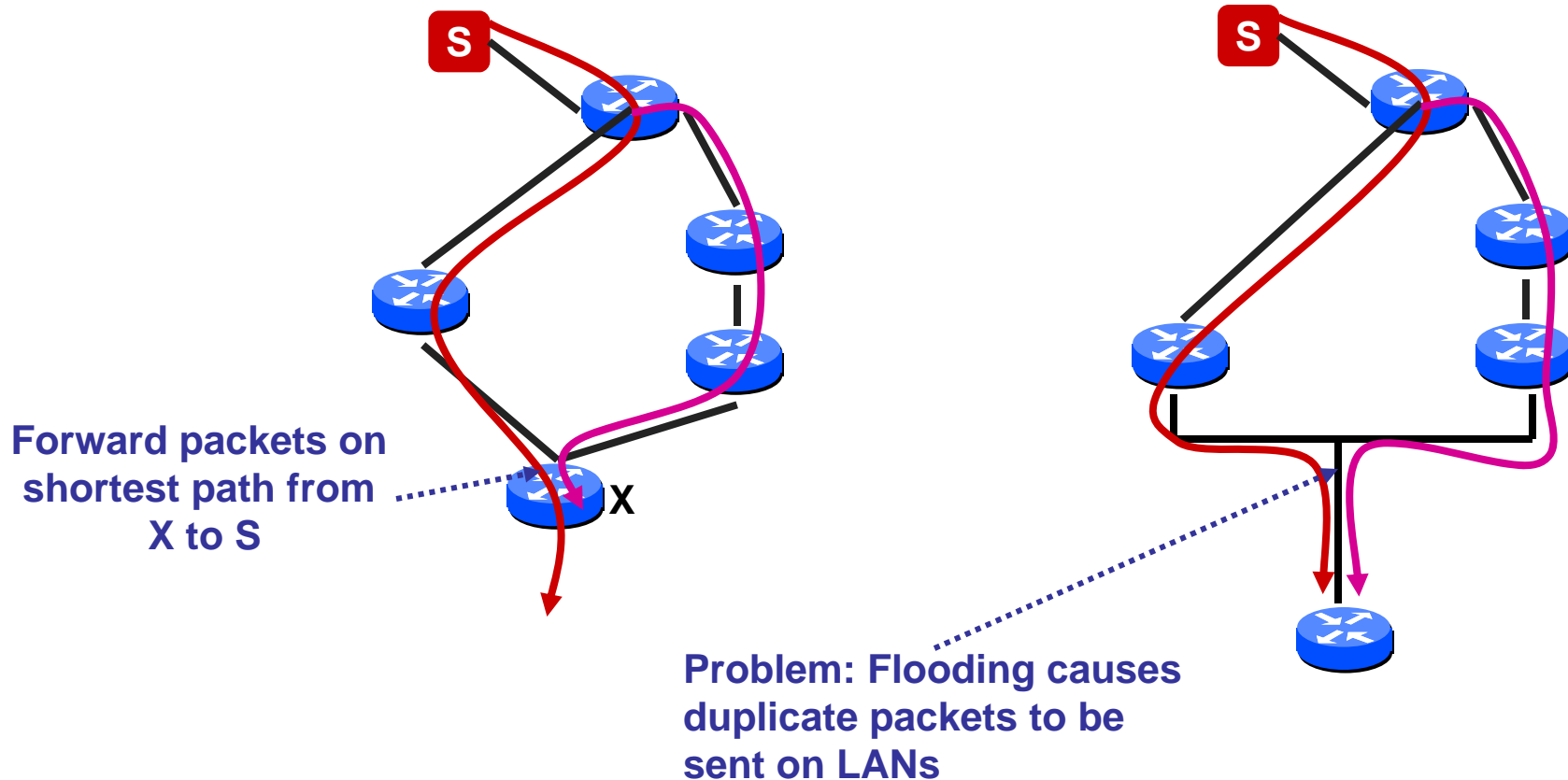
Source-based: Flood and Prune (DV)

- Extensions to unicast distance vector algorithm
- Goal
 - ◆ Multicast packets delivered along shortest-path tree from sender to members of the multicast group
 - ◆ Likely have different tree for different senders
- Distance Vector Multicast Routing Protocol (DVMRP) developed as a progression of algorithms
 - ◆ Reverse Path Flooding (RPF)
 - ◆ Reverse Path Broadcast (RPB)
 - ◆ Reverse Path Multicast (RPM)

Source-based: Reverse Path Flooding (RPF)

- Observation: Shortest-path multicast tree is subtree of shortest-path broadcast tree
- Approach: Use shortest-path broadcast tree
- Use **reverse path** to determine shortest path
 - ◆ Router forwards a packet **from S** iff received from the shortest-path link **to S**
 - ◆ Exactly what is in entry in forwarding table?
 - » To reach S along shortest path, use link L
 - » If received packet from S on L, it came along shortest path
- How are packets forwarded?
 - ◆ **Flooding** – forward packets to multicast address out to all links except incoming link (hence reverse path flooding)

Example: Reverse Path Flooding

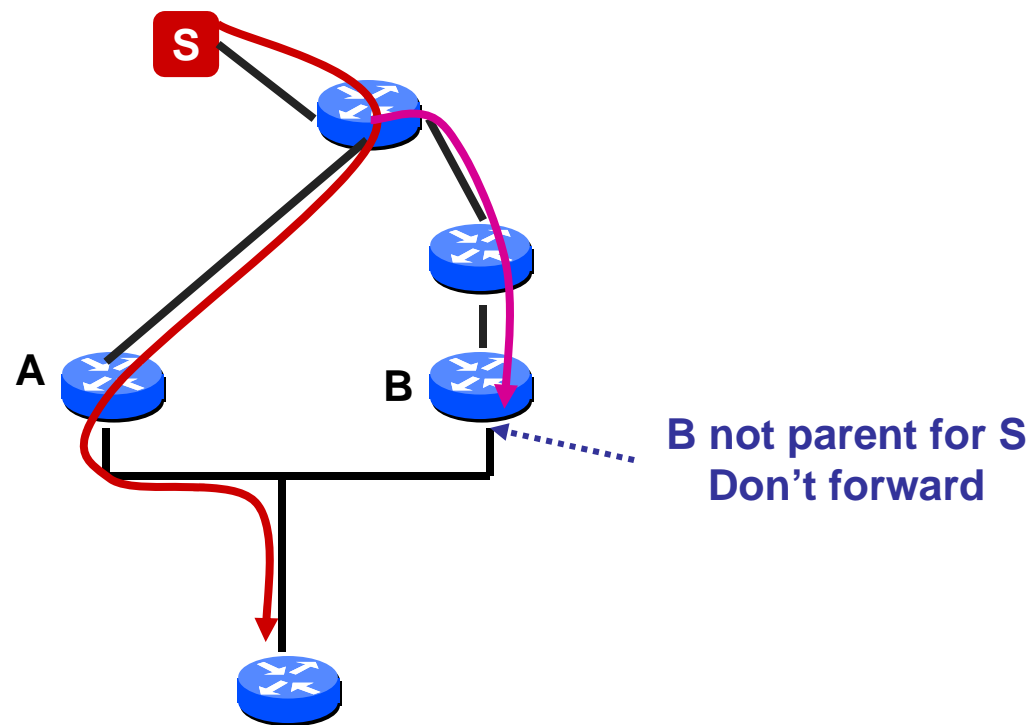


Solution:

Reverse Path Broadcast (RPB)

- Flooding vs. broadcast
 - ◆ With flooding, a single packet can be sent along an individual link multiple times
 - » Each router attached to link can potentially forward same packet
 - ◆ RPB sends a packet along a link at most once
- Approach: Define **parent** and **child** routers for each link
 - ◆ Relative to each link and each source S
 - ◆ Router is a parent for link if it has minimum path to S
 - ◆ All other routers on the link are children
 - ◆ Only parent router is allowed to forward multicast packets on link
- How to decide parent and children routers for link?
 - ◆ In routing updates; router determines if is parent

Example: Reverse Path Broadcasting



Source-based: Reverse Path Multicast (RPM)

- Problem: Still **broadcasting** up to leaf networks
- Idea: Instead of **actively building tree**, use reports to **actively prune tree**
- Start with a full broadcast tree to all links (RPB),
- Prune (S,G) at leaf if it has no members
 - ◆ Send **Non-Membership Report** (NMR) to prev-hop for S
- If all children of router R prune (S,G)
 - ◆ Send NMR for (S,G) to parent of R
- Soft-state management (must refresh NMR or rejoin)
- New group member sends graft (anti-prune) message

Source-based: Link State

- Use existing link-state routing algorithm (e.g. OSPF)
- Idea: Include active groups in LSPs
 - ◆ Each router can compute shortest path tree from source to all destinations for any group
 - ◆ Trigger new flood of LSPs on group membership change
- Performance issues
 - ◆ Expensive to precompute all (S,G) trees
 - ◆ Keep cache of trees and compute new trees on demand when new (S,G) packet arrives
 - ◆ Workload/topology dependent
- Best known example: MOSPF

Shared tree approaches

- Unicast packets to Rendezvous Point (RP), which multicasts packet on shared tree
- Tree construction
 - ◆ Receivers send join messages to RP for each group
 - ◆ Intermediate routers install state to create per-group tree
 - ◆ Key advantage is routers only store $O(G)$ state
 - ◆ Potential optimizations: reroute to source-specific trees for local group members or high data-rate sources
 - ◆ Example: PIM-SM
- Issues
 - ◆ Delay, fault tolerance, RP selection

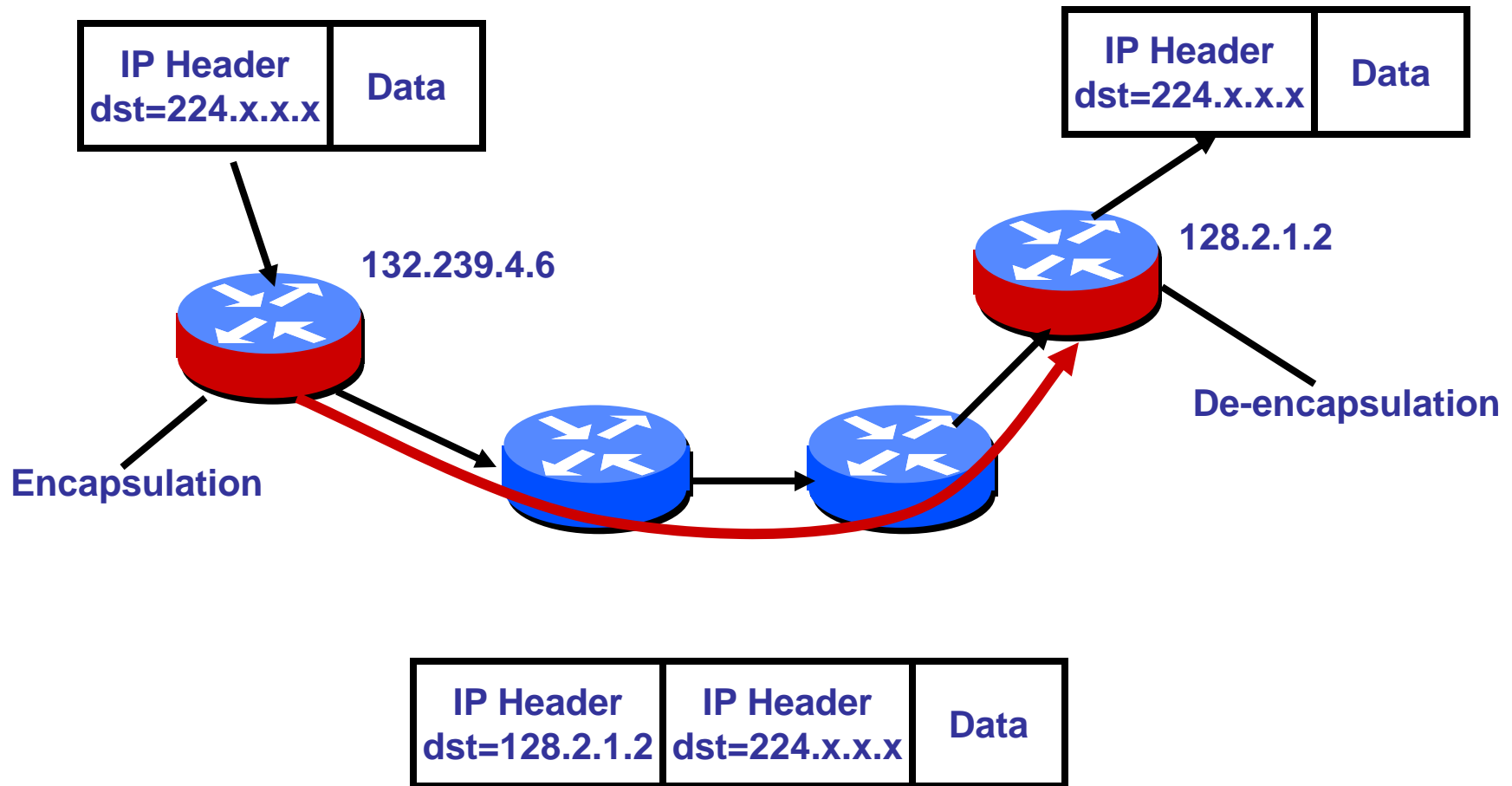
IP Multicast today

- IP Multicast has generated 1000s of papers, but has not been widely deployed in the Internet...
- Why?
 - ◆ General deployment difficulties
 - » Most routers don't support multicast
 - » Mbone (multicast backbone)
 - ◆ Inter-domain multicast complexity
 - » Policy issues again
 - ◆ Economics of multi-source multicast
 - » What/who do you charge for multicast traffic?

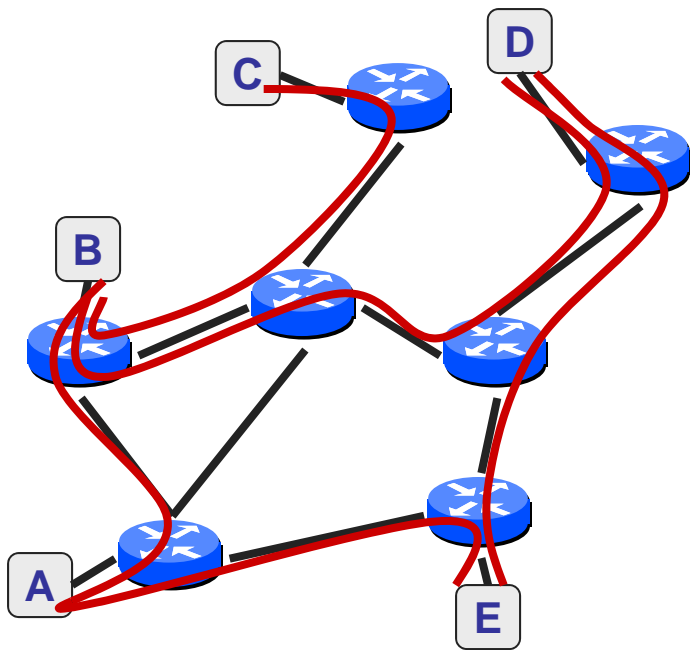
Multicast evolution

- How to deploy a new network-layer service?
 - ◆ Difficult to change router software
 - ◆ Difficult to change all routers
- Mbone (tunneling)
 - ◆ Special multicast routers (built from PCs/Workstations)
 - ◆ Construct virtual topology between them (overlay)
 - ◆ Run routing protocol over virtual topology
 - ◆ Virtual point-to-point links called **tunnels**
 - » Multicast traffic encapsulated in IP datagrams
 - » Multicast routers forward over tunnels according to computed virtual next-hop

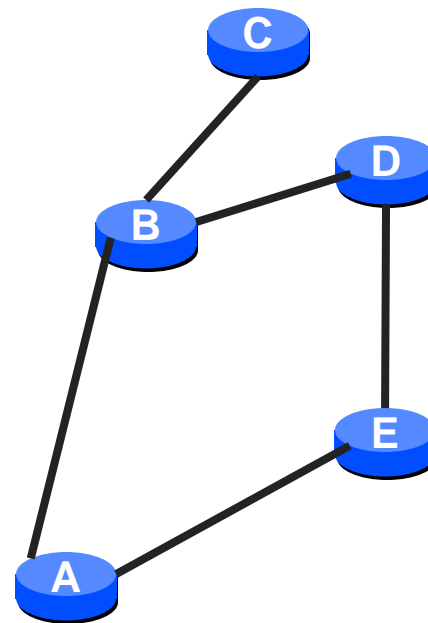
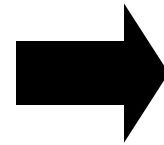
Tunneling



Virtual overlay network



Real topology with tunnels



Virtual overlay topology

Mbone Pro/Con

- Success story
 - ◆ Multicast video to 20 sites in 1992
 - ◆ Easy to deploy, no explicit router support
 - ◆ Ran DVMRP and had 100s of routers
- Drawbacks
 - ◆ Manual tunnel creation/maintenance
 - ◆ Inefficient
 - ◆ No routing policy (single tree)
 - ◆ Why would an ISP deploy a new Mbone node?

Inter-domain Multicast routing

- Technical issues
 - ◆ How to exchange reachability information?
 - ◆ How to construct trees?
 - ◆ Who controls RP in shared tree?
- MBGP (Multicast BGP): reachability to multicast sources per prefix
- PIM-SM: shared tree multicast protocol
- MSDP (Multicast source discovery protocol): RP per group per AS, communication presence of group sources between RPs
- BGMP (Border gateway multicast protocol): alternative proposal, single shared tree with group addresses owned by individual ASs

Economic issues

- ISP router migration cycle
 - ◆ Can't afford new routers on edge
- Domain independence
 - ◆ Do I want my customers multicast controlled by an RP in a competitors domain?
 - ◆ Why run an RP for which I have no senders or receivers?
- Billing model
 - ◆ Inconsistent with input-rate-based billing
 - ◆ No group management (how big is group?)
- Group management
 - ◆ Who is in the group? Who can send? Security...
- Limited Multicast addresses

Summary

- Multicast service model
 - ◆ One-to-many, anonymous communication
 - ◆ Simple host interface
- Per-source tree routing
 - ◆ Efficient trees
 - ◆ $O(S \cdot G)$ state explosion for large networks/groups
- Shared tree
 - ◆ More complex, fragile, hard to manage
 - ◆ Inefficient trees
 - ◆ Only requires $O(G)$ state on routers
- Operational and Economic issues matter in deployment
- Killer app not found

Next time

- Router design
 - ◆ How routers are actually built