

Network Layer – Part II

Lecturer: Prof. Xiaoming Fu

Assistant: Yachao Shao (MSc),

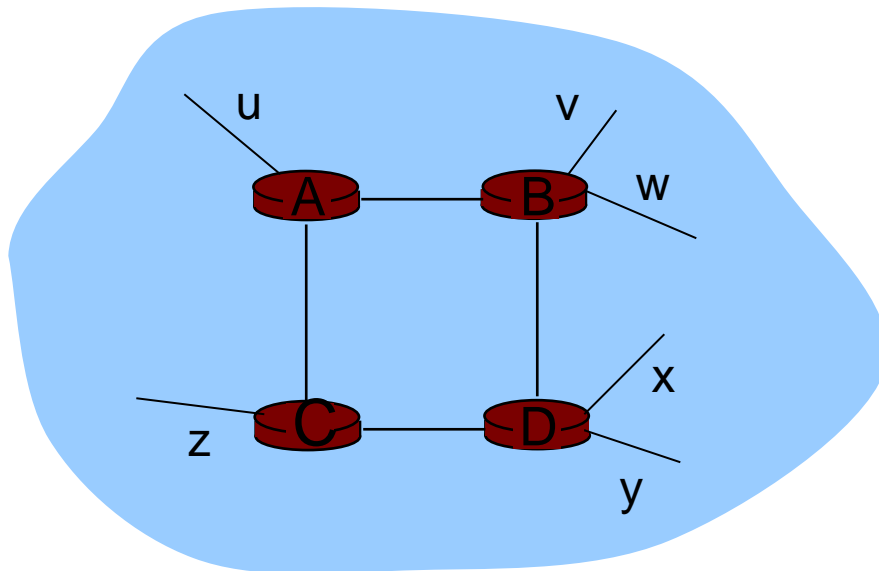
Fabian Wölk (MSc)

Network Layer II

- 4.5 Routing protocols
 - **Routing Information Protocol (RIP)**
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - Mobile IP

Routing Information Protocol (RIP)

- distance vector algorithm
- included in BSD-UNIX Distribution in 1982
- distance metric: # of hops (max = 15 hops)



From router A to subnets:

<u>destination</u>	<u>hops</u>
u	1
v	2
w	2
x	3
y	3
z	2

Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - **Open Shortest Path First (OSPF)**
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - Mobile IP

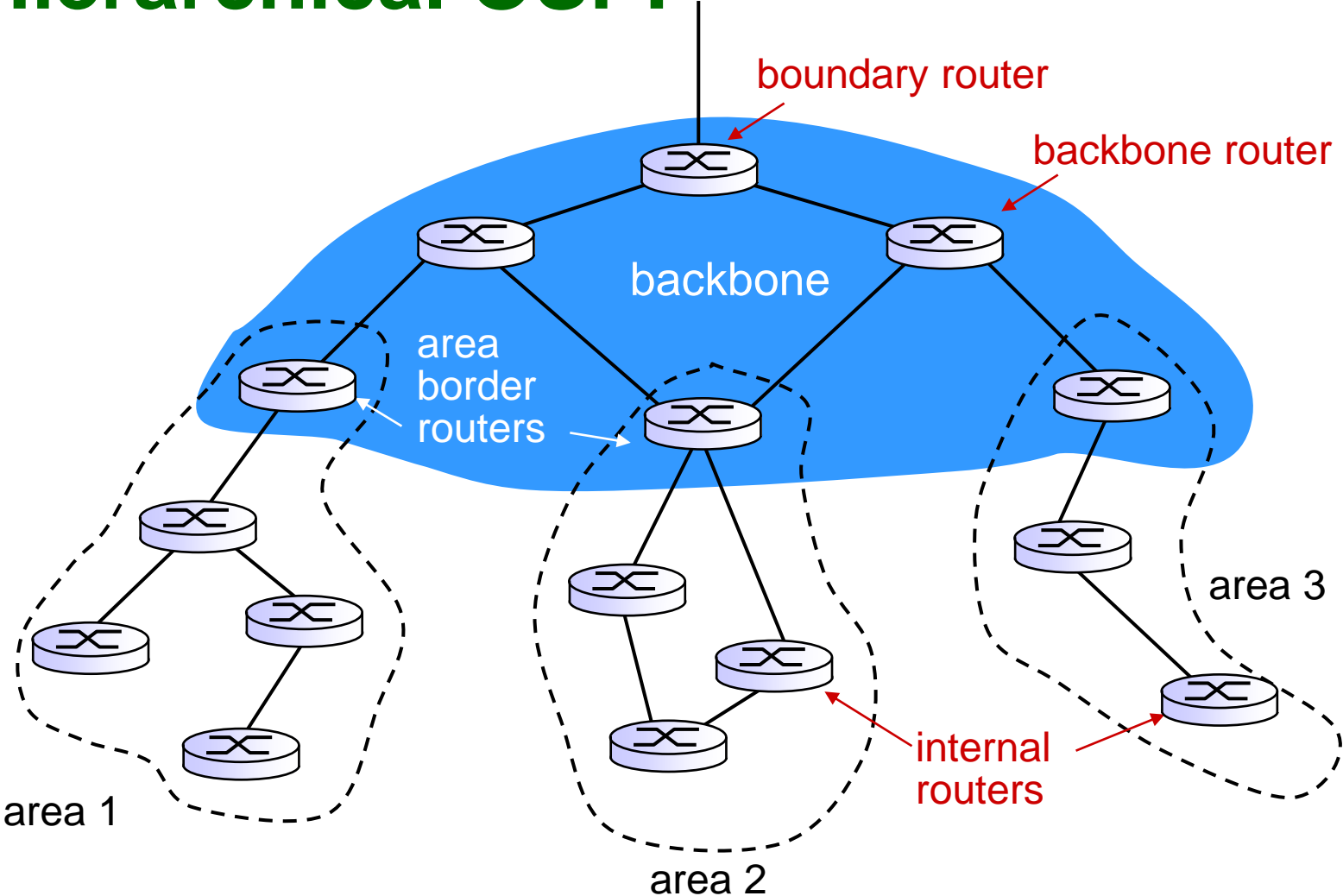
Open Shortest Path First (OSPF)

- “open”: publicly available
- uses Link State algorithm
 - LS packet dissemination
 - topology map at each node
 - route computation using Dijkstra’s algorithm
- OSPF advertisement carries one entry per neighbor router
- advertisements disseminated to entire AS (via flooding)
 - carried in OSPF messages directly over IP (rather than TCP or UDP)

OSPF “advanced” features (not in RIP)

- **Security**: all OSPF messages authenticated (to prevent malicious intrusion)
- **multiple** same-cost **paths** allowed (only one path in RIP)
- for each link, multiple cost metrics for different **TOS** (e.g., satellite link cost set low for best effort ToS; high for real-time ToS)
- Integrated uni- and **multicast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- **Hierarchical** OSPF in large domains.

Hierarchical OSPF



Hierarchical OSPF

- **two-level hierarchy:** local area, backbone.
 - Link-state advertisements only in area
 - each node has detailed area topology; only know direction (shortest path) to nets in other areas.
- **area border routers:** “summarize” distances to nets in own area, advertise to other Area Border routers.
- **backbone routers:** run OSPF routing limited to backbone.
- **boundary routers:** connect to other AS's.

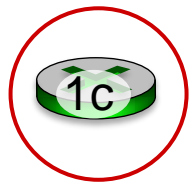
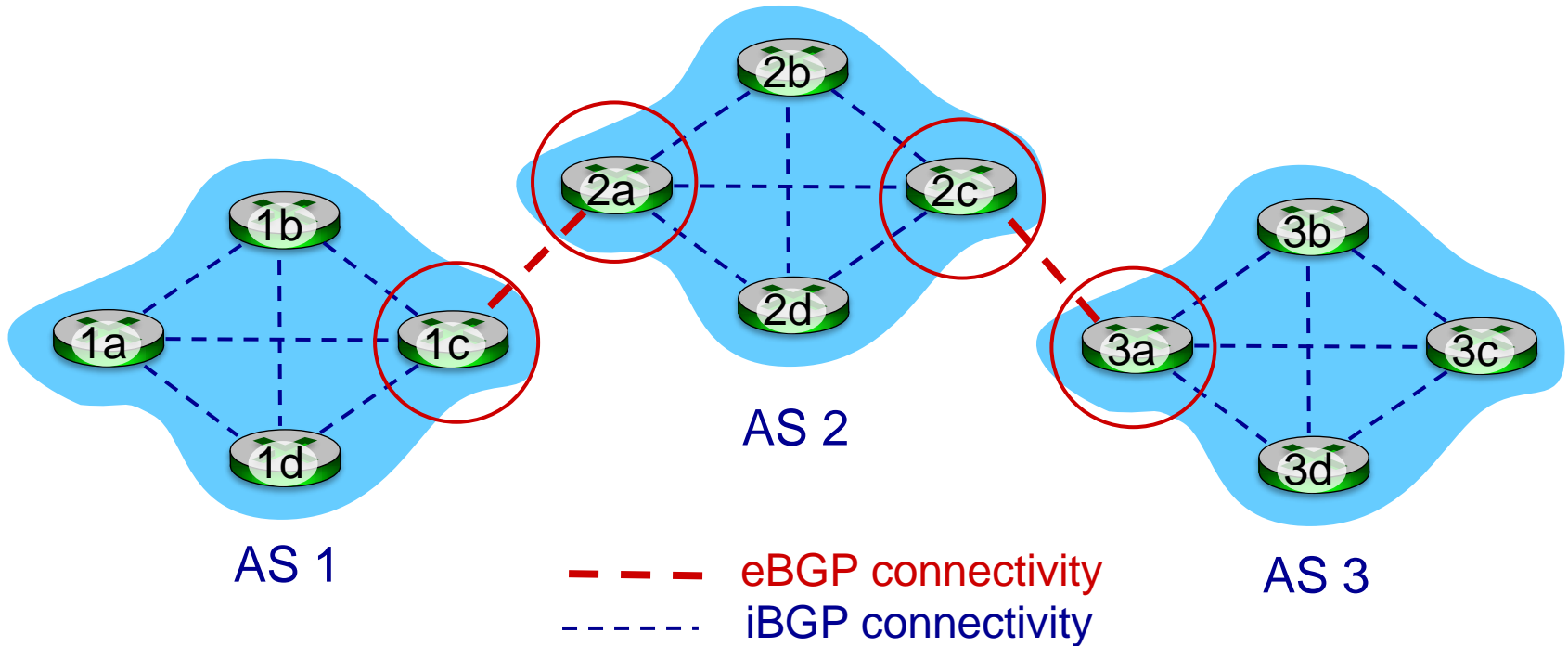
Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - **Border Gateway Protocol (BGP)**
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - Mobile IP

Inter-AS routing: BGP

- BGP (Border Gateway Protocol): the de facto standard
- BGP provides each AS means to:
 - **eBGP**: obtain subnet reachability information from neighboring ASes
 - **iBGP**: propagate reachability information to all AS-internal routers.
 - determine “good” routes to other networks based on reachability information and *policy*
- allows subnet to advertise its existence to rest of Internet: “I am here”

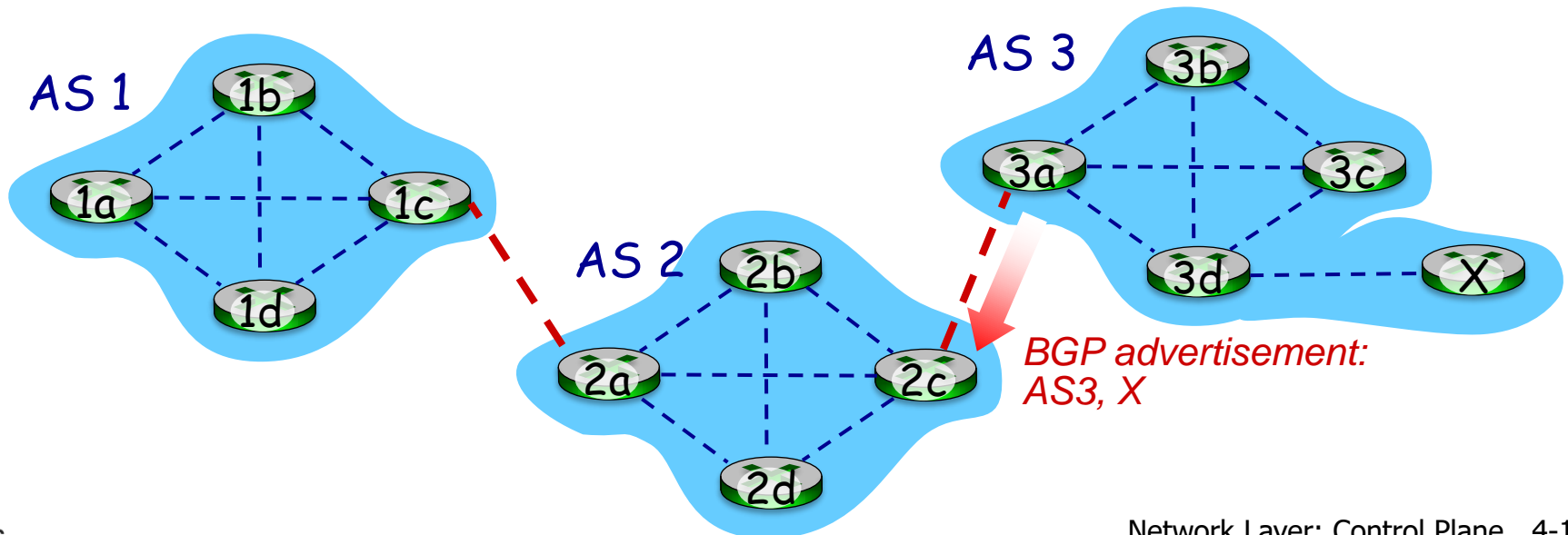
eBGP, iBGP connections



gateway routers run both eBGP and iBGP protocols

BGP basics

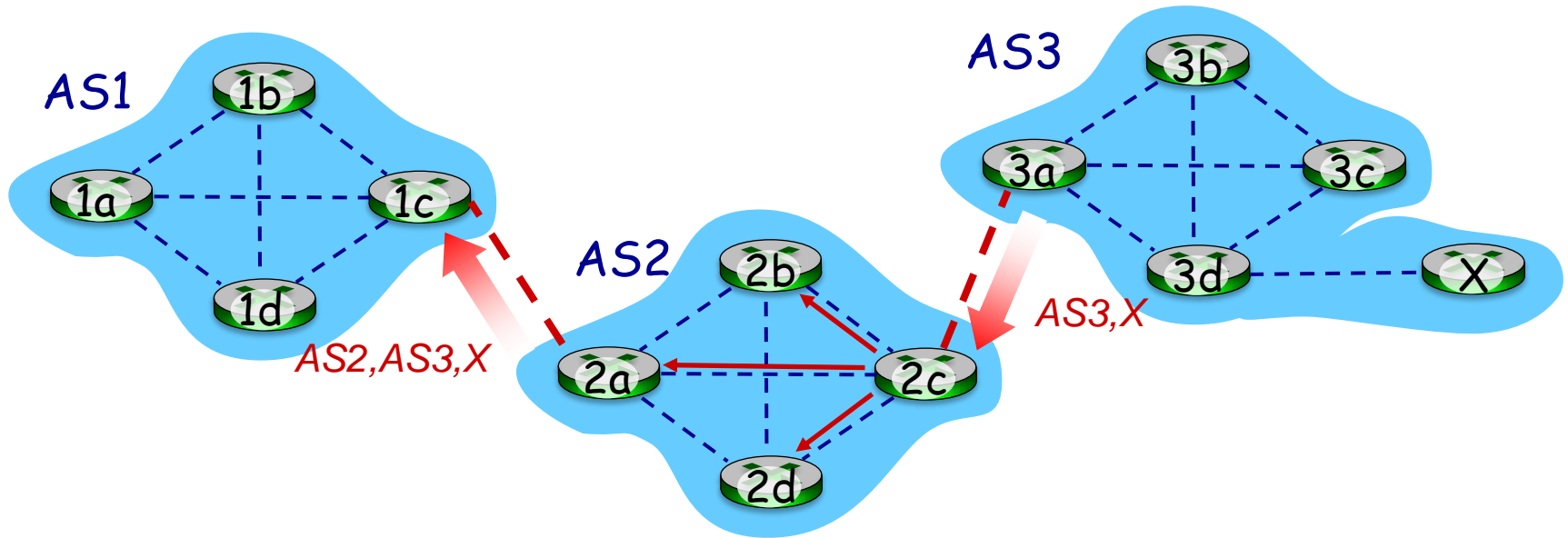
- **BGP session:** two BGP routers (“peers”) exchange BGP messages over semi-permanent TCP connection:
 - advertising *paths* to different destination network prefixes (BGP is a “path vector” protocol)
- when AS3 gateway router 3a advertises path **AS3,X** to AS2 gateway router 2c:
 - AS3 *promises* to AS2 it will forward datagrams towards X



Path attributes and BGP routes

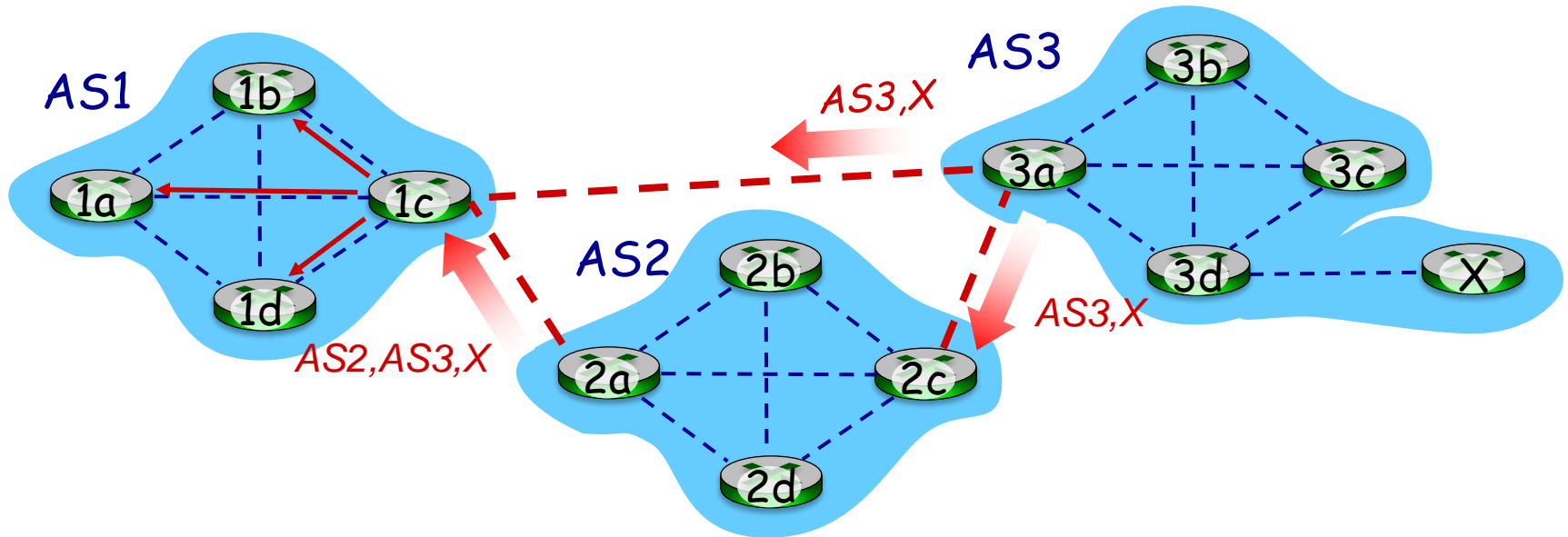
- advertised prefix includes BGP attributes
 - prefix + attributes = “route”
- two important attributes:
 - **AS-PATH**: list of ASes through which prefix advertisement has passed
 - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS
- *Policy-based routing*:
 - gateway receiving route advertisement uses *import policy* to accept/decline path (e.g., never route through AS Y).
 - AS policy also determines whether to *advertise* path to other neighboring ASes

BGP path advertisement



- AS2 router 2c receives path advertisement **AS3,X** (via eBGP) from AS3 router 3a
- Based on AS2 policy, AS2 router 2c accepts path **AS3,X**, propagates (via iBGP) to all AS2 routers
- Based on AS2 policy, AS2 router 2a advertises (via eBGP) path **AS2, AS3,X** to AS1 router 1c

BGP path advertisement



gateway router may learn about **multiple** paths to destination:

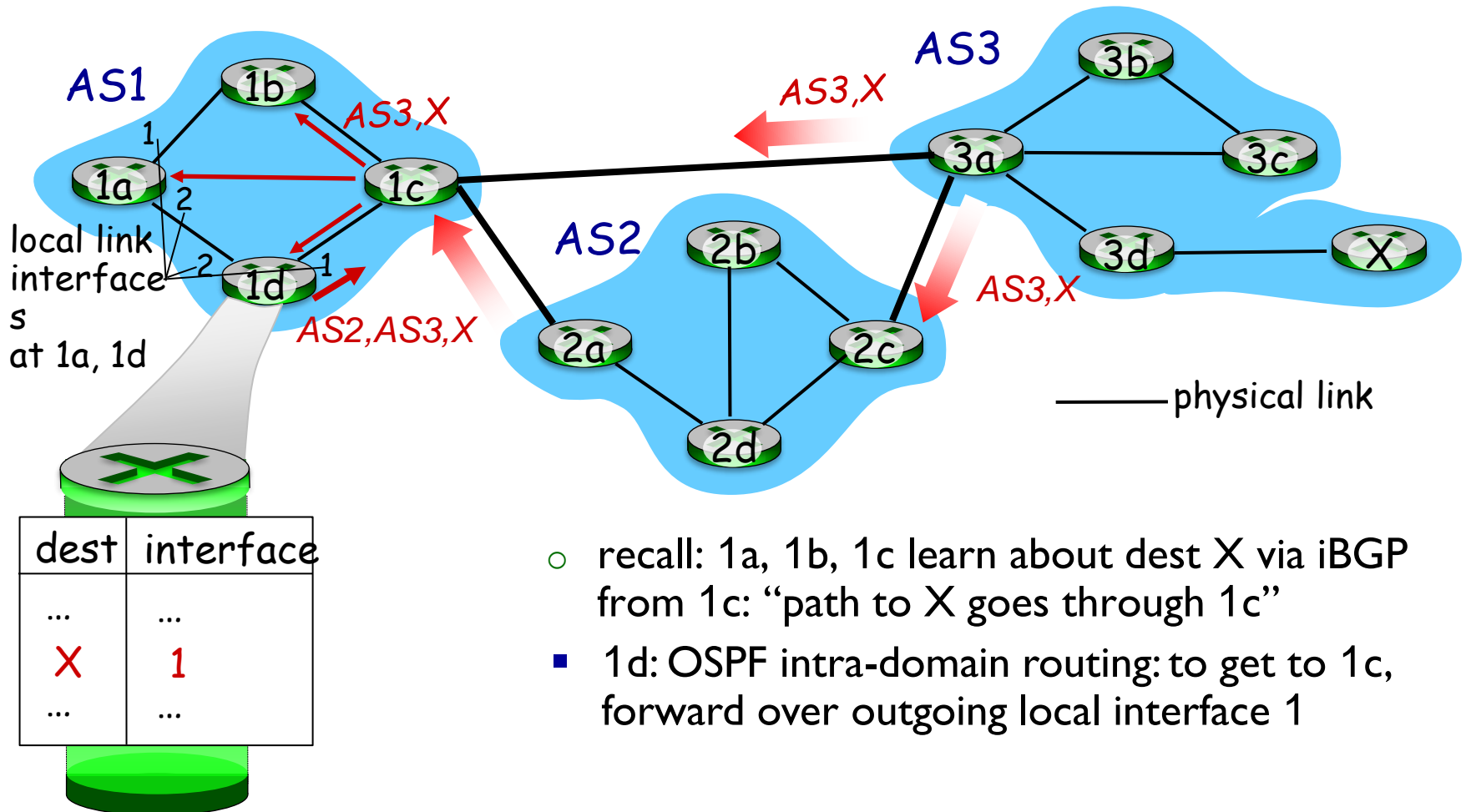
- AS1 gateway router 1c learns path **AS2,AS3,X** from 2a
- AS1 gateway router 1c learns path **AS3,X** from 3a
- Based on policy, AS1 gateway router 1c chooses path **AS3,X**, and **advertises path within AS1 via iBGP**

BGP messages

- BGP messages exchanged between peers over TCP connection
- BGP messages:
 - **OPEN:** opens TCP connection to remote BGP peer and authenticates sending BGP peer
 - **UPDATE:** advertises new path (or withdraws old)
 - **KEEPALIVE:** keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - **NOTIFICATION:** reports errors in previous msg; also used to close connection

BGP, OSPF, forwarding table entries

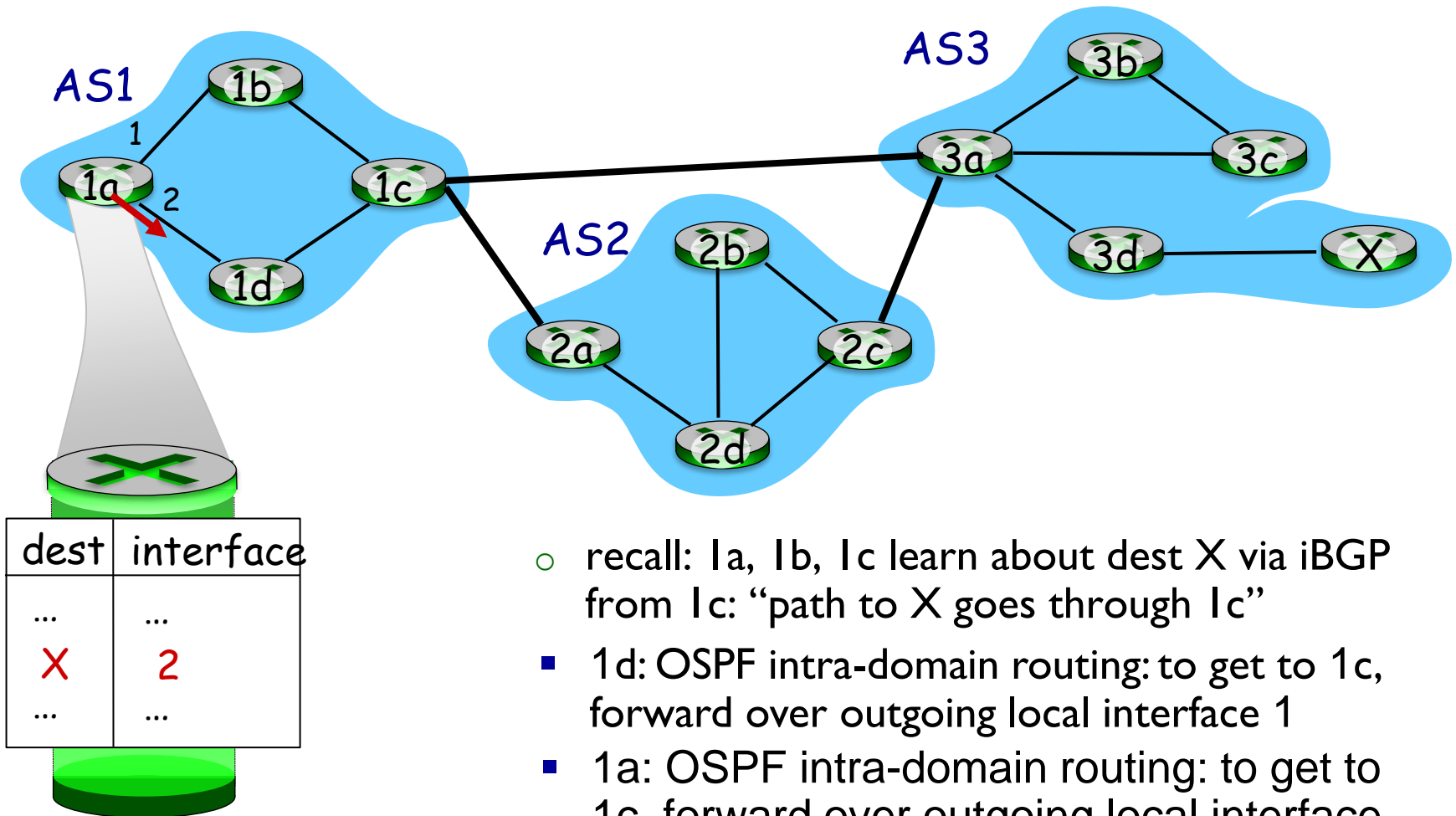
Q: how does router set forwarding table entry to distant prefix?



- recall: 1a, 1b, 1c learn about dest X via iBGP from 1c: “path to X goes through 1c”
- 1d: OSPF intra-domain routing: to get to 1c, forward over outgoing local interface 1

BGP, OSPF, forwarding table entries

Q: how does router set forwarding table entry to distant prefix?

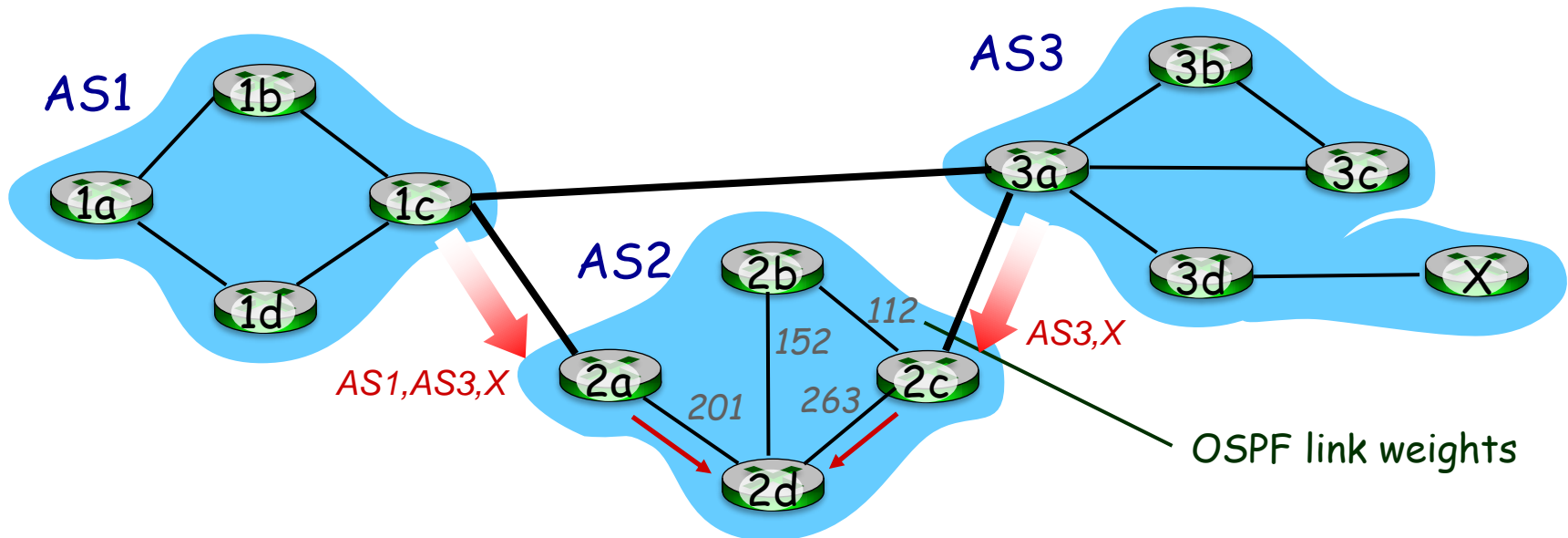


- recall: 1a, 1b, 1c learn about dest X via iBGP from 1c: “path to X goes through 1c”
- 1d: OSPF intra-domain routing: to get to 1c, forward over outgoing local interface 1
- 1a: OSPF intra-domain routing: to get to 1c, forward over outgoing local interface 2

BGP route selection

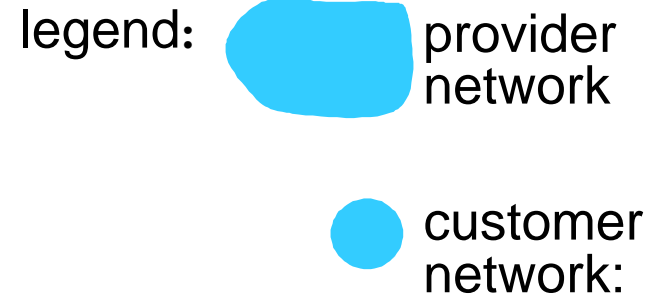
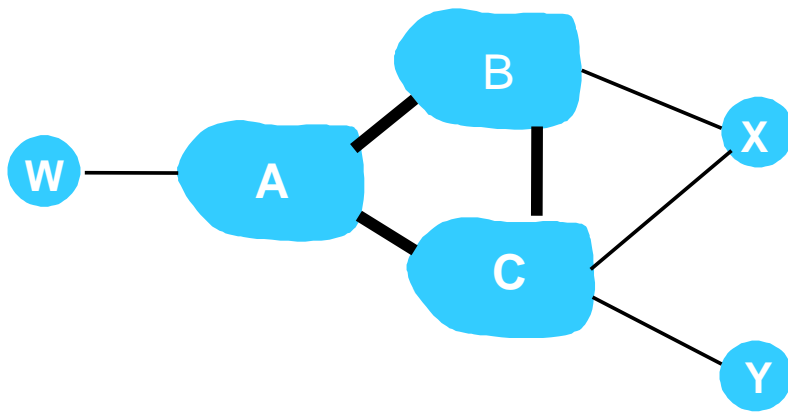
- router may learn about more than one route to destination AS, selects route based on:
 1. local preference value attribute: policy decision
 2. shortest AS-PATH
 3. closest NEXT-HOP router: hot potato routing
 4. additional criteria

Hot Potato Routing



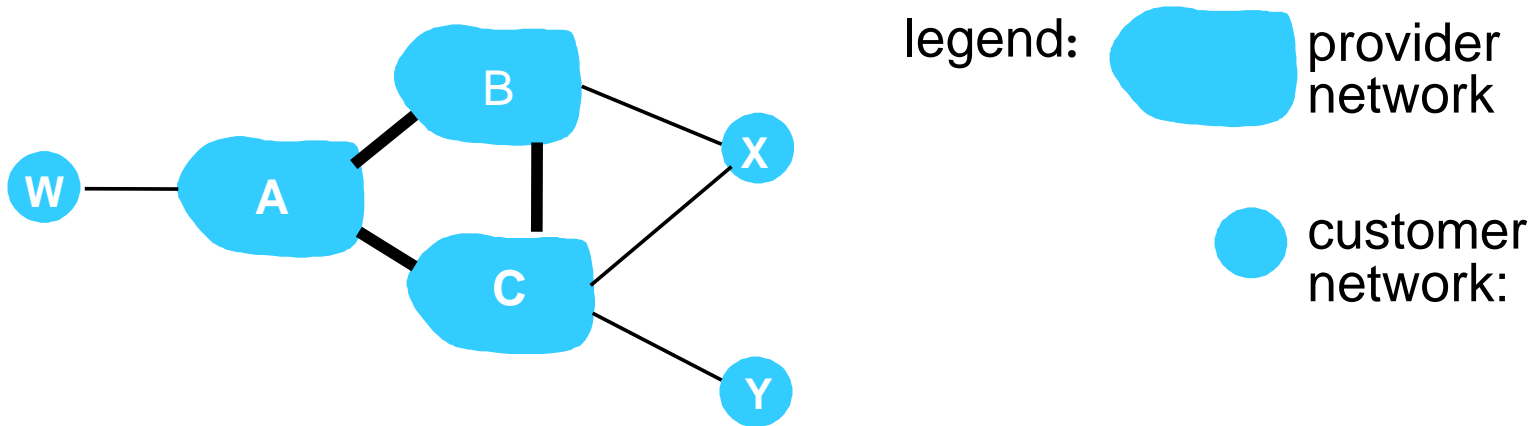
- 2d learns (via iBGP) it can route to X via 2a or 2c
- *hot potato routing*: choose local gateway that has least intra-domain cost (e.g., 2d chooses 2a, even though more AS hops to X): don't worry about inter-domain cost!

BGP routing policy



- A,B,C are provider networks
- X,W,Y are customer networks
- X is dual-homed: attached to two networks
 - X does not want to route from B via X to C
 - .. so X will not advertise to B a route to C

BGP routing policy (2)

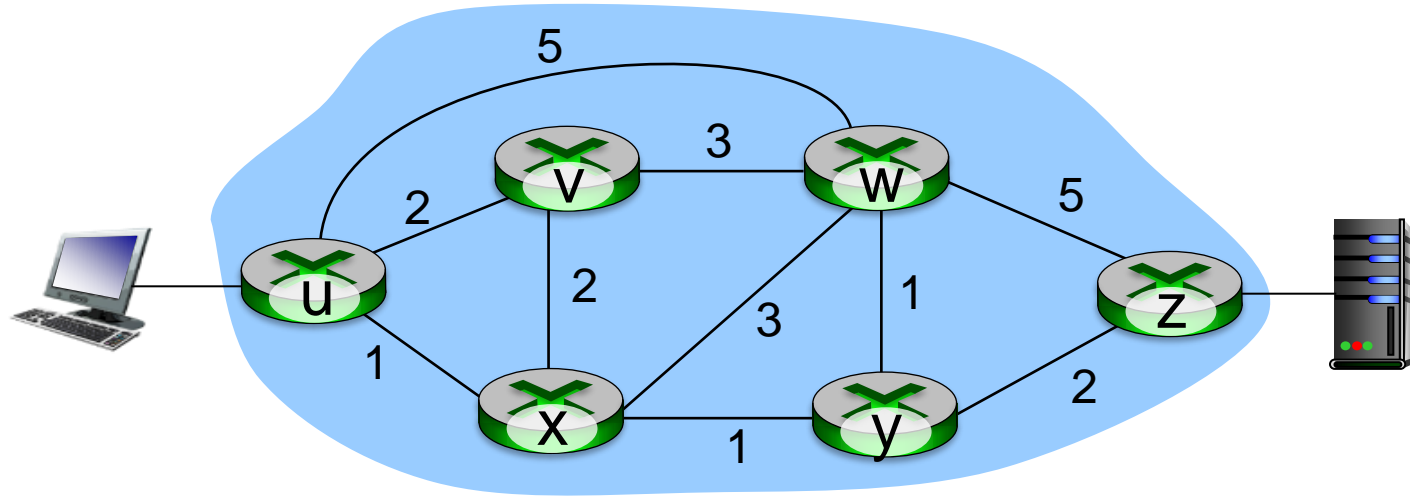


- A advertises path AW to B
- B advertises path BAW to X
- Should B advertise path BAW to C?
 - No way! B gets no “revenue” for routing CBAW since neither W nor C are B’s customers
 - B wants to force C to route to w via A
 - B wants to route *only* to/from its customers!

Why different Intra- and Inter-AS routing?

- Policy
 - Inter-AS: admin wants control over how its traffic routed, who routes through its net.
 - Intra-AS: single admin, so no policy decisions needed
- Scale
 - hierarchical routing saves table size, reduced update traffic
- Performance
 - Intra-AS: can focus on performance
 - Inter-AS: policy may dominate over performance

Traffic engineering: difficult traditional routing

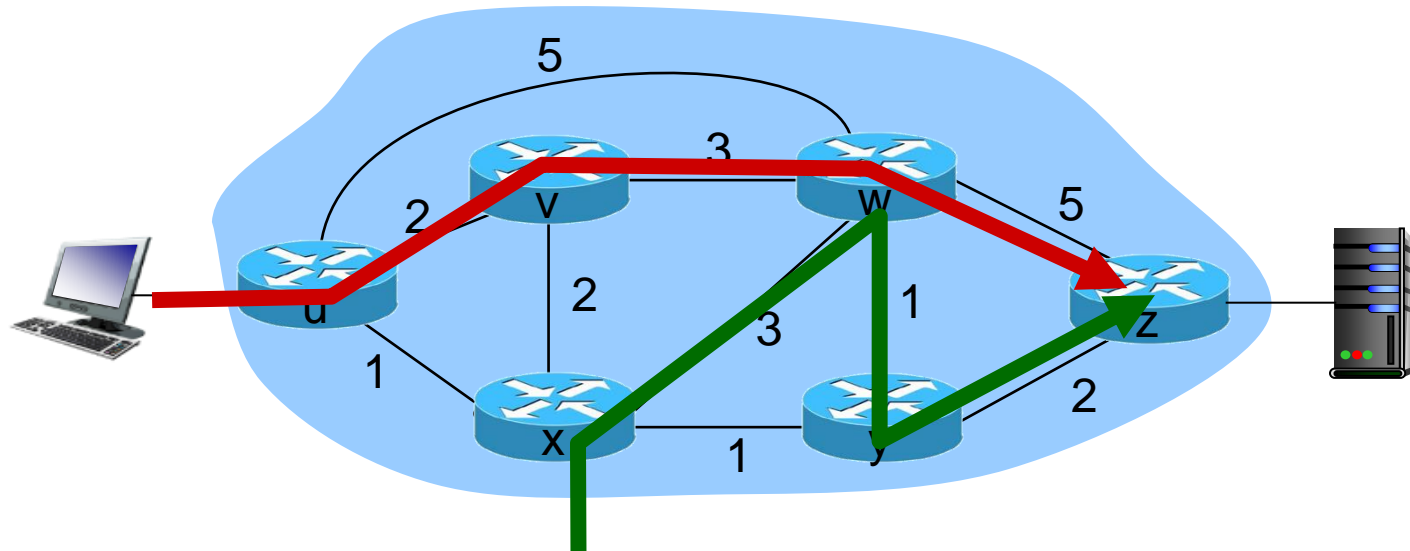


Q: what if network operator wants u-to-z traffic to flow along $uvwz$, x-to-z traffic to flow $xwyz$?

A: need to define link weights so traffic routing algorithm computes routes accordingly (or need a new routing algorithm)!

Link weights are only control “knobs”: wrong!

Traffic engineering: difficult



Q: what if w wants to route blue and red traffic differently?

A: can't do it (with destination based forwarding, and LS, DV routing)

Software defined networking (SDN)

4. programmable control applications

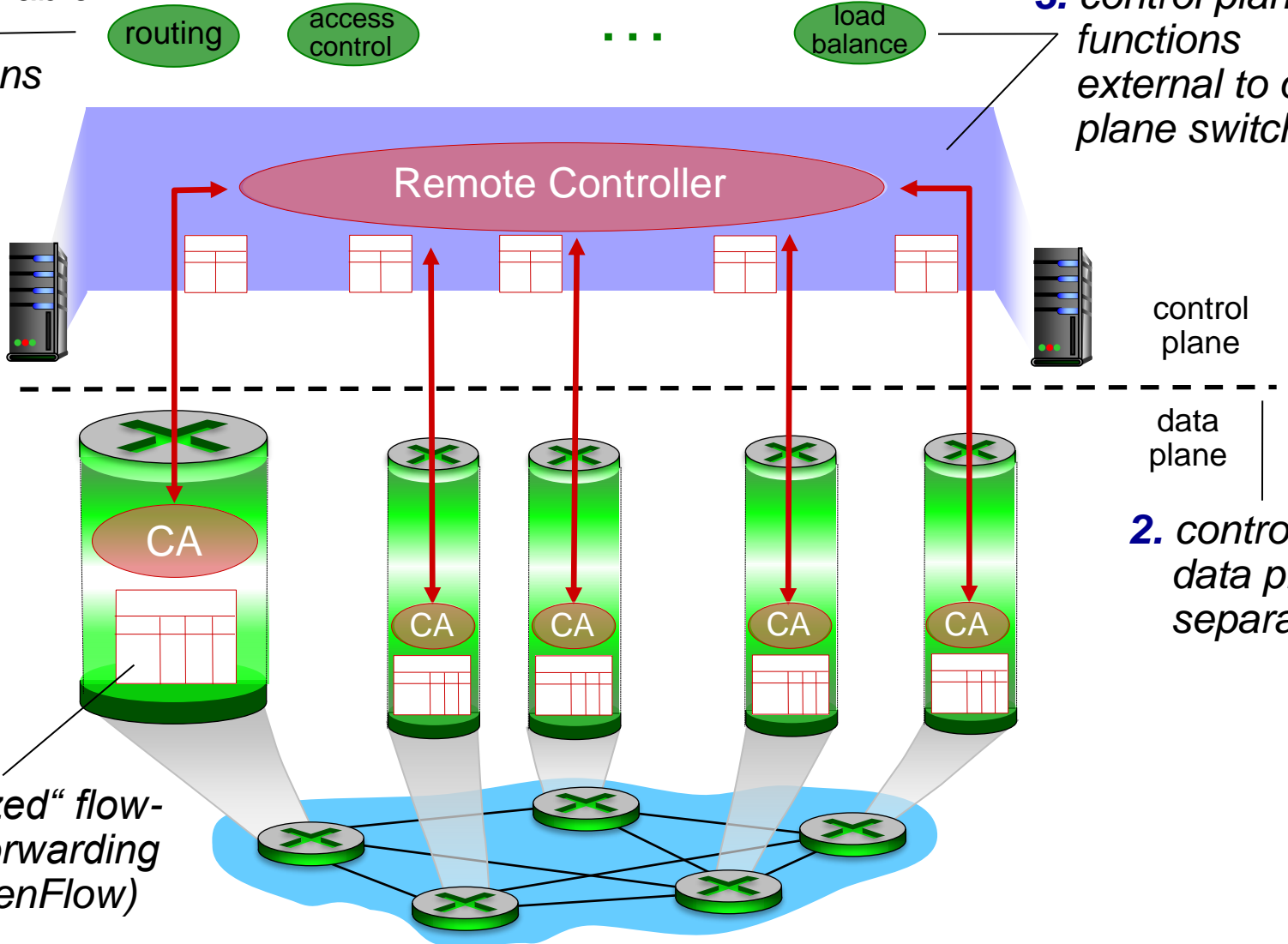
routing

access control

...

load balance

3. control plane functions external to data-plane switches



1. generalized "flow-based" forwarding (e.g., OpenFlow)

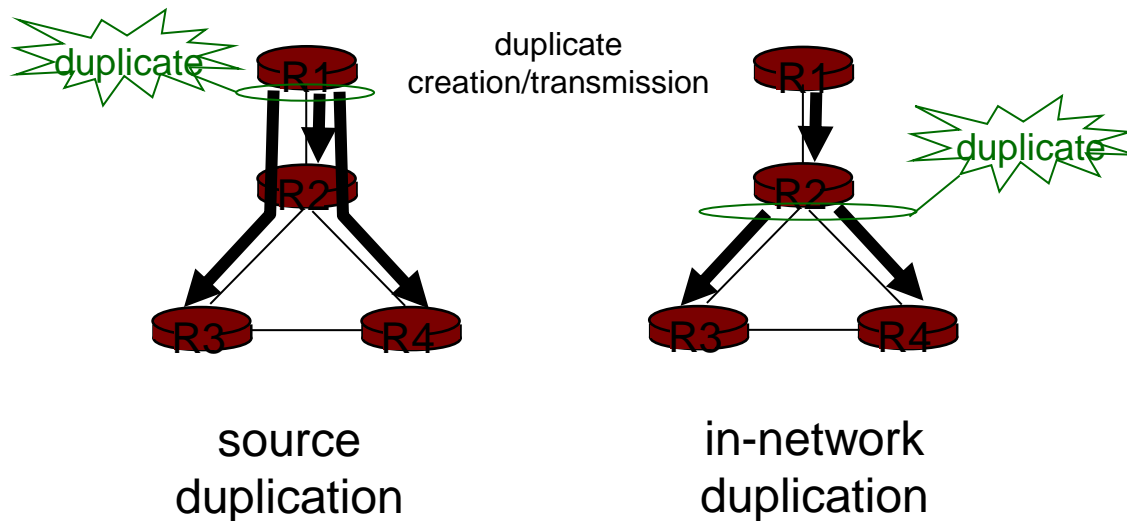
2. control, data plane separation

Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - **Broadcast routing**
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - Mobile IP

Broadcast Routing

- Deliver packets from source to all other nodes
- Source duplication is inefficient:



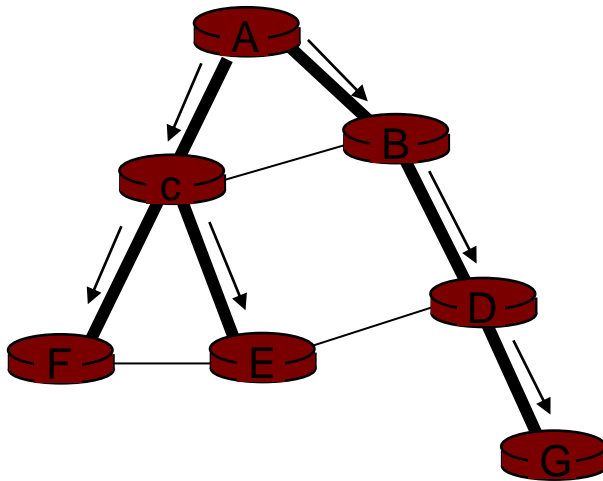
- Where does info come from? How to use in link state?

In-network duplication

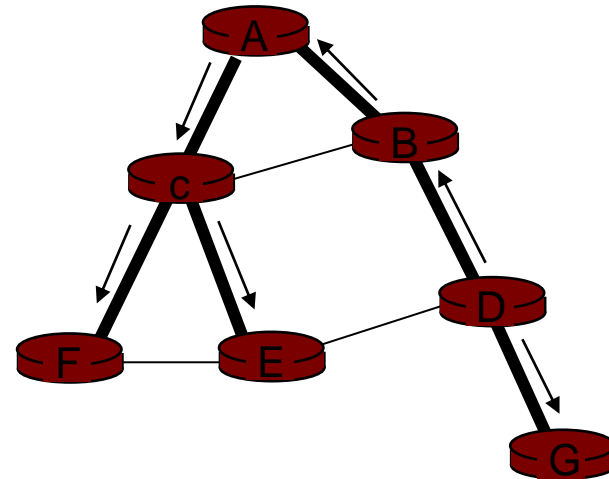
- **Flooding:** when node receives broadcast packets, sends copy to all neighbors
 - Problems: cycles & broadcast storm
- **Controlled flooding:** node only broadcast pkt if it hasn't broadcasted same pkt before
 - Node keeps track of pkt ids already broadcasted
 - Reverse path forwarding (RPF): only forward pkt if it arrived on shortest path between node and source
- **Spanning tree**
 - No redundant packets received by any node

Spanning Tree

- First construct a spanning tree
- Nodes forward copies only along spanning tree



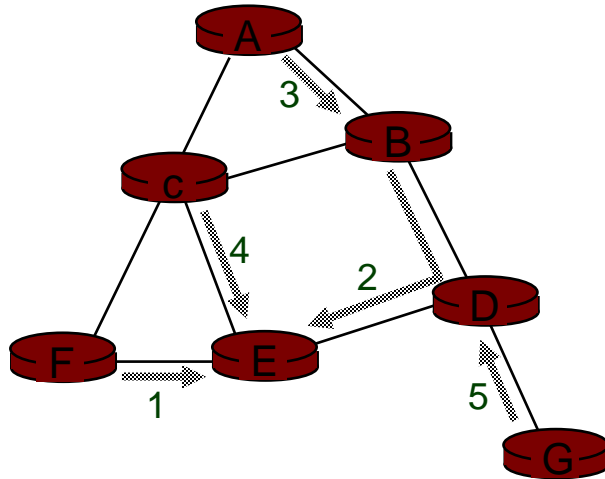
(a) Broadcast initiated at A



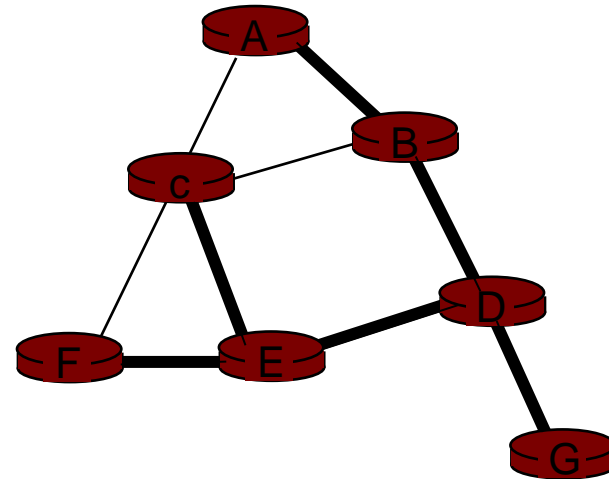
(b) Broadcast initiated at D

Spanning Tree: Creation

- Center node
- Each node sends unicast join message to center node 'E'
 - Message forwarded until it arrives at a node already belonging to spanning tree



(a) Stepwise construction of spanning tree



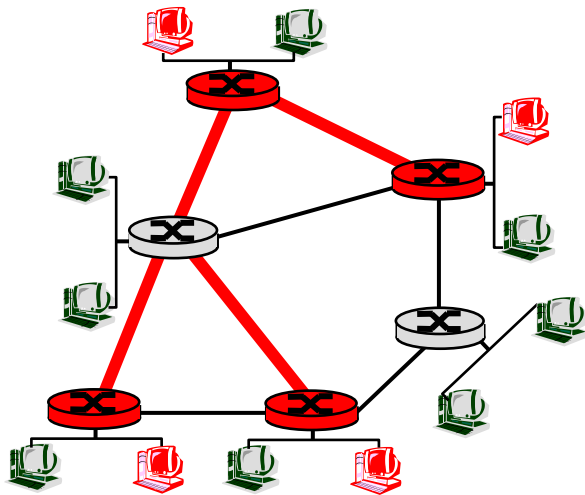
(b) Constructed spanning tree

Network Layer II

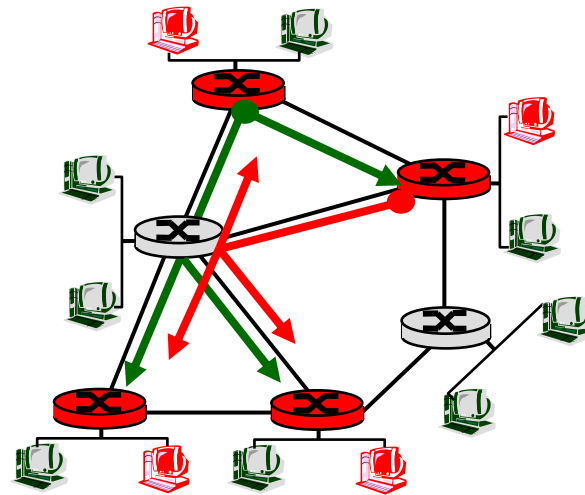
- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - **Multicast routing**
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - Mobile IP

Multicast Routing: Problem Statement

- **Goal:** find a tree (or trees) connecting routers that have local multicast group members
 - Tree: not all paths between routers used
 - Source-based: different tree from each sender to receiver
 - Shared-tree: same tree used by all group members



Shared tree



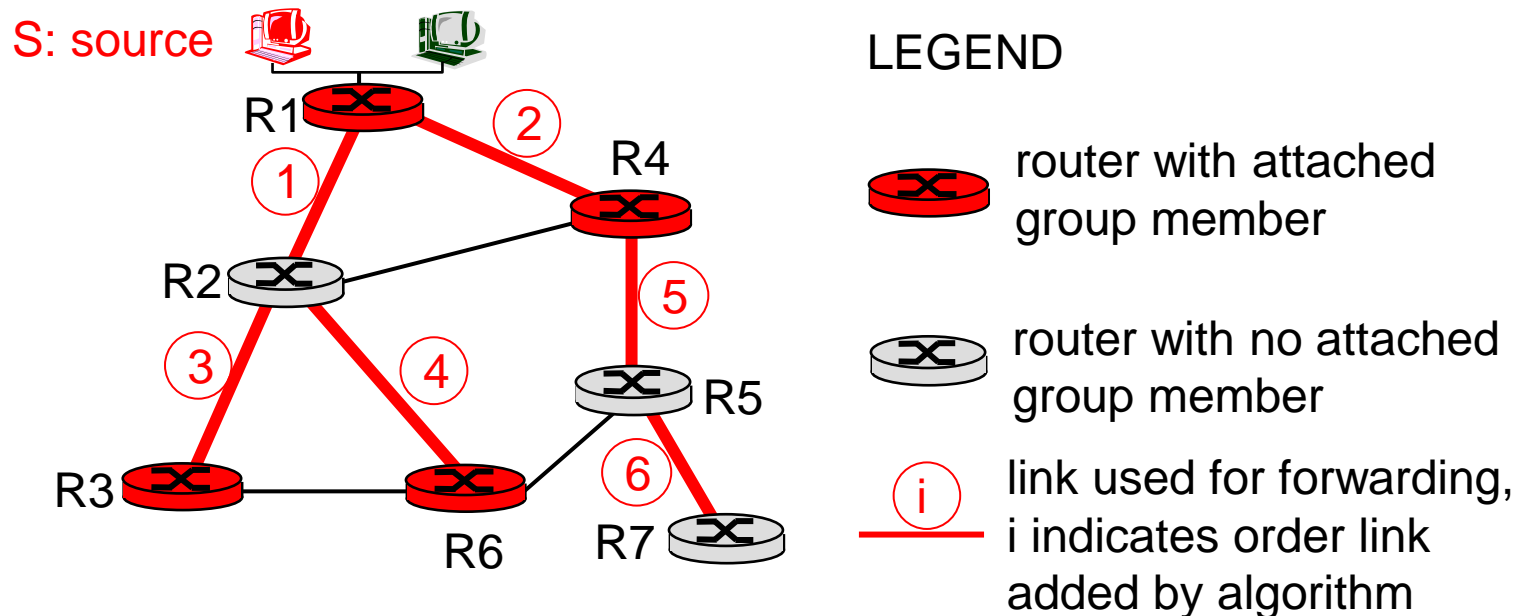
Source-based trees

Approaches for building mcast trees

- Source-based tree: one tree per source
 - shortest path trees
 - reverse path forwarding
- Group-shared tree: group uses one tree
 - minimal spanning (Steiner)
 - center-based trees

Shortest Path Tree

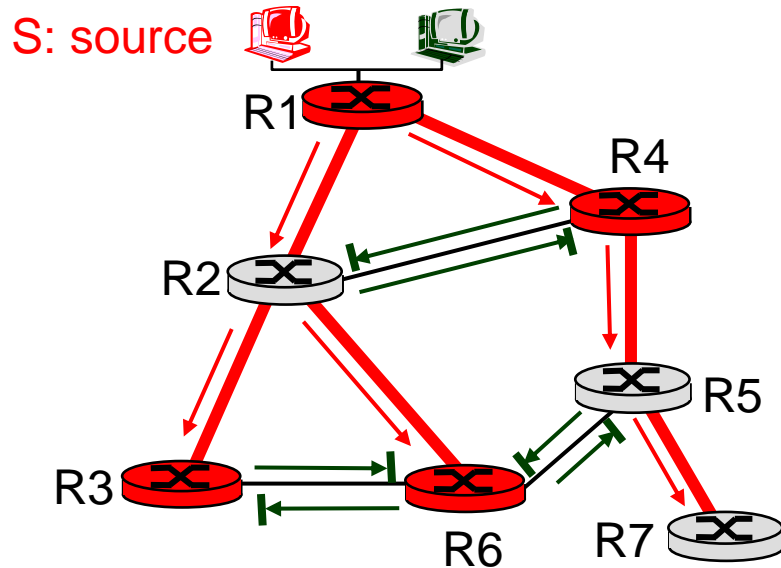
- Multicast forwarding tree: tree of shortest path routes from source to all receivers
 - Dijkstra's algorithm



Reverse Path Forwarding

- Relies on router's knowledge of unicast shortest path from it to sender
- Each router has simple forwarding behavior:
 - if (multicast datagram received on incoming link on shortest path back to center)
 - then flood datagram onto all outgoing links
 - else ignore datagram

Reverse Path Forwarding: example



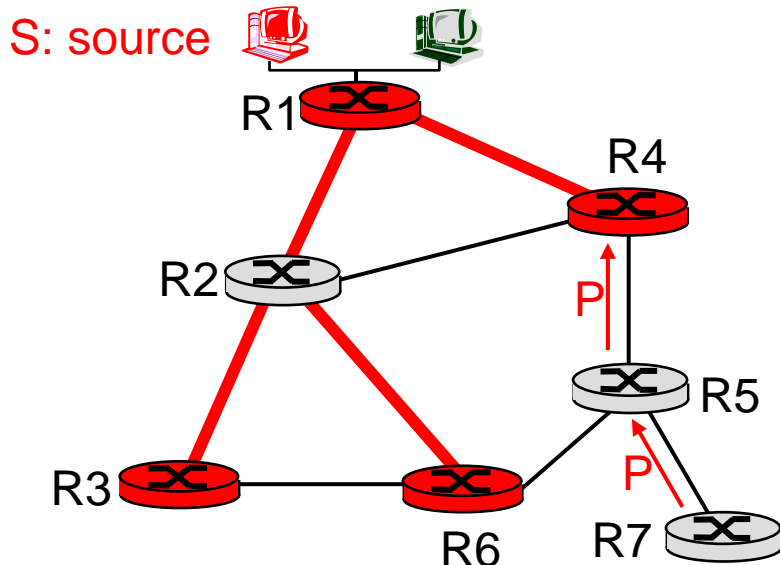
LEGEND

- router with attached group member
- router with no attached group member
- datagram will be forwarded
- datagram will not be forwarded





- result is a source-specific *reverse* SPT
 - may be a bad choice with asymmetric links

Reverse Path Forwarding: pruning

- forwarding tree contains subtrees with no multicast group members
 - no need to forward datagrams down subtree
 - “prune” msgs sent upstream by router with no downstream group members



LEGEND

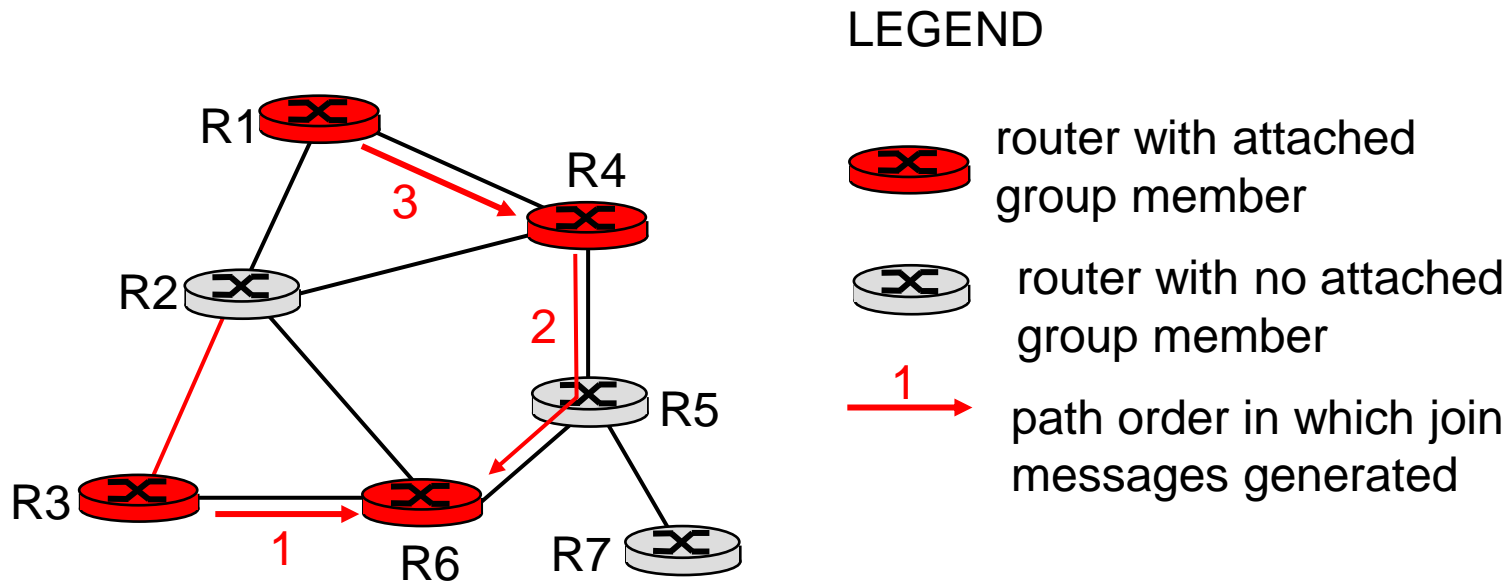
-  router with attached group member
-  router with no attached group member
-  prune message
-  links with multicast forwarding

Center-based trees

- Single delivery tree shared by all
- One router identified as “center” of tree
- To join:
 - edge router sends unicast join-msg addressed to center router
 - join-msg “processed” by intermediate routers and forwarded towards center
 - join-msg either hits existing tree branch for this center, or arrives at center
 - path taken by join-msg becomes new branch of tree for this router

Center-based trees: an example

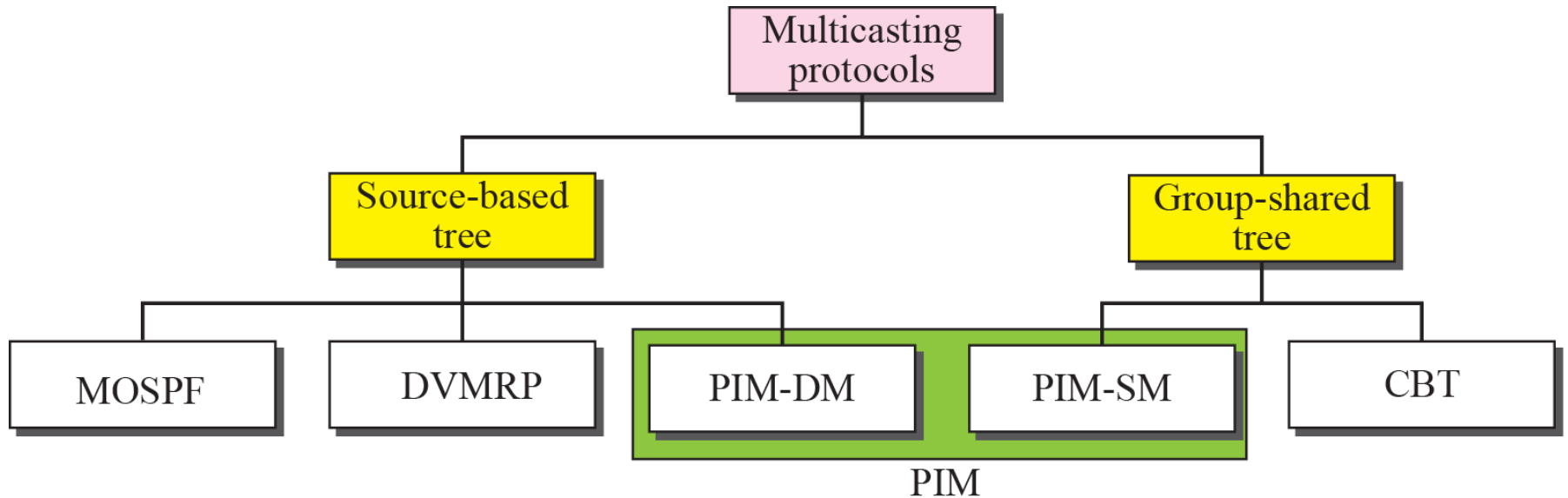
Suppose R6 chosen as center:



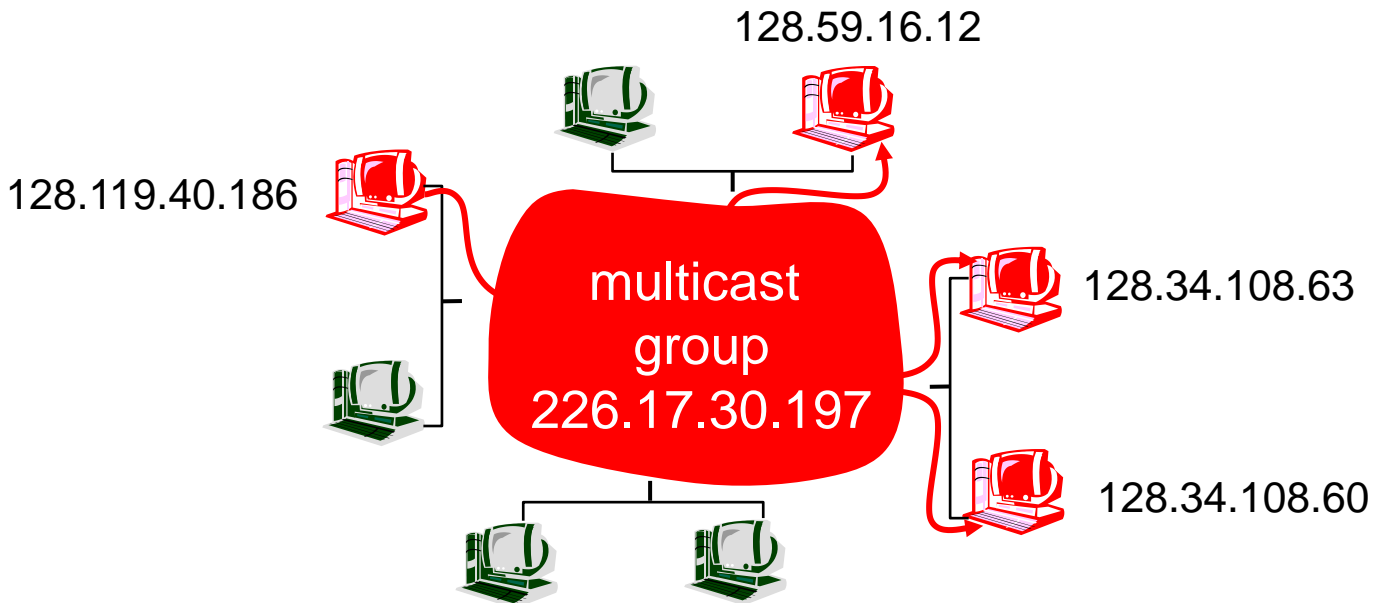
Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - **Multicast routing protocols**
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - Mobile IP

Multicast routing protocols



Internet Multicast Service Model

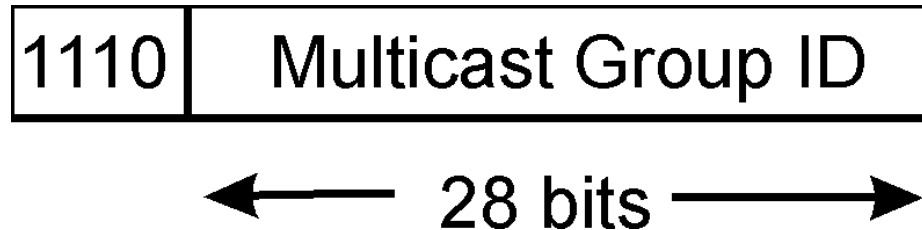


Multicast group concept: use of **indirection**

- hosts addresses IP datagram to multicast group
- routers forward multicast datagrams to hosts that have “joined” that multicast group

Multicast Groups

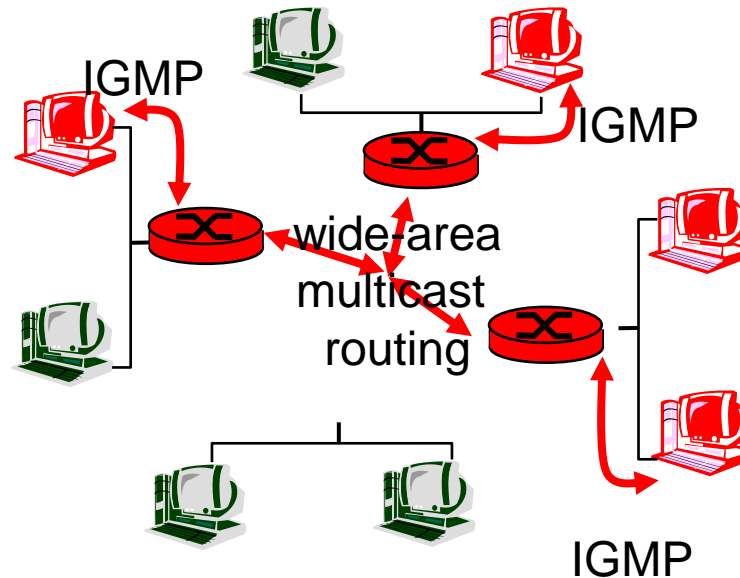
- Class D Internet addresses reserved for multicast:



- Host group semantics:
 - anyone can “join” (receive pkts) multicast group
 - anyone can send pkts to multicast group
 - no network-layer identification to hosts of the members
- Needed: infrastructure to deliver mcast-addressed datagrams to all hosts that have joined that multicast group

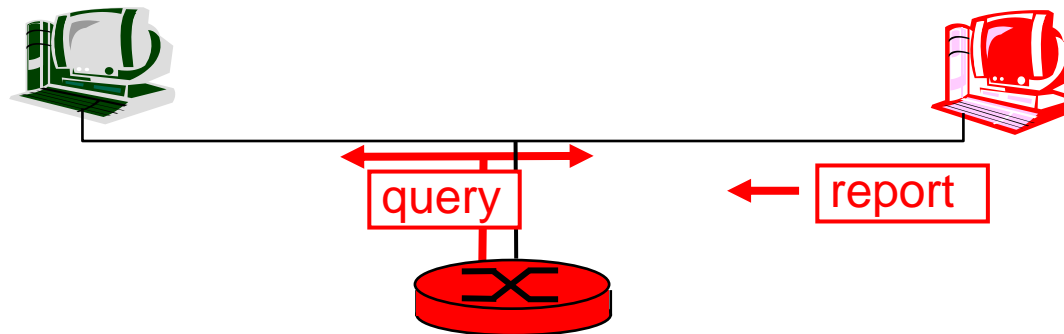
Joining a mcast group: two-step process

- **Local:** host informs local mcast router of a desire to join group:
 - IGMP (Internet Group Management Protocol)
- **Wide area:** local router interacts with other routers to receive mcast datagram flow
 - many protocols (e.g., DVMRP, MOSPF, PIM)



IGMP: Internet Group Management Protocol

- **Host:** sends IGMP report when application joins mcast group
 - IP_ADD_MEMBERSHIP socket option
 - host needs not explicitly “disjoin” group when leaving
- **Router:** sends IGMP query at regular intervals
 - host belonging to a mcast group must reply to query



Internet Multicasting Routing: DVMRP

- DVMRP: distance vector multicast routing protocol, RFC1075
- flood and prune: reverse path forwarding, source-based tree
 - RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers
 - no assumptions about underlying unicast
 - initial datagram to mcast group flooded everywhere via RPF
 - routers not wanting group: send upstream prune msgs

DVMRP: continued...

- soft state: DVMRP router periodically (1 min.) “forgets” branches are pruned:
 - mcast data again flows down unpruned branch
 - downstream router: re prune or else continue to receive data
- routers can quickly regraft to tree
 - following IGMP join at leaf
- odds and ends
 - commonly implemented in commercial routers
 - Mbone routing done using DVMRP

PIM: Protocol Independent Multicast

- not dependent on any specific underlying unicast routing algorithm (works with all)
- two different multicast distribution scenarios :
 - Dense:
 - group members densely packed, in “close” proximity.
 - bandwidth more plentiful
 - Sparse:
 - # networks with group members small wrt # interconnected networks
 - group members “widely dispersed”
 - bandwidth not plentiful

Consequences of Sparse-Dense Dichotomy

○ Dense

- group membership by routers assumed until routers explicitly prune
- data-driven construction on mcast tree (e.g., RPF)
- bandwidth and non-group-router processing profligate

○ Sparse

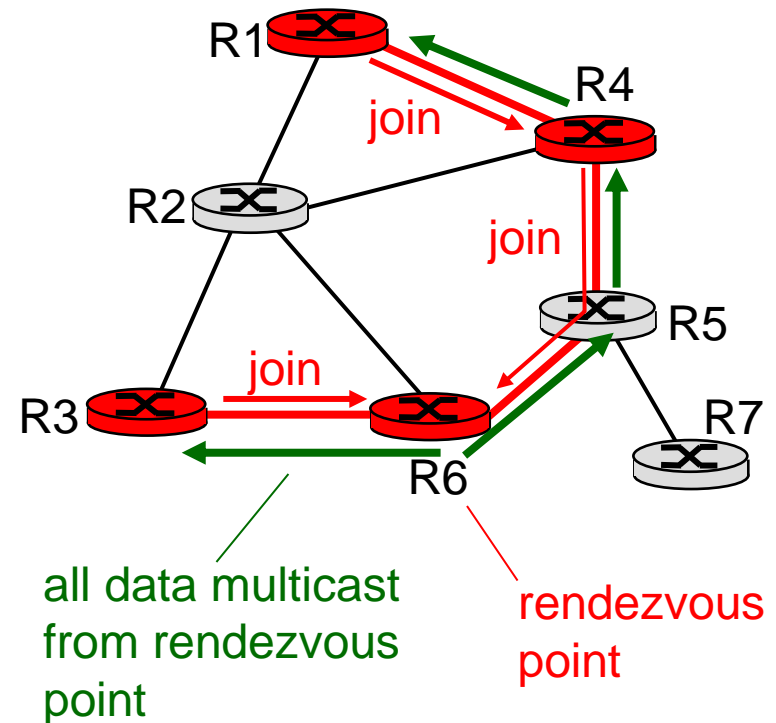
- no membership until routers explicitly join
- receiver-driven construction of mcast tree (e.g., center-based)
- bandwidth and non-group-router processing conservative

PIM- Dense Mode

- Flood-and-prune RPF, similar to DVMRP but
 - underlying unicast protocol provides RPF info for incoming datagram
 - less complicated (less efficient) downstream flood than DVMRP reduces reliance on underlying routing algorithm
 - has protocol mechanism for router to detect it is a leaf-node router

PIM - Sparse Mode

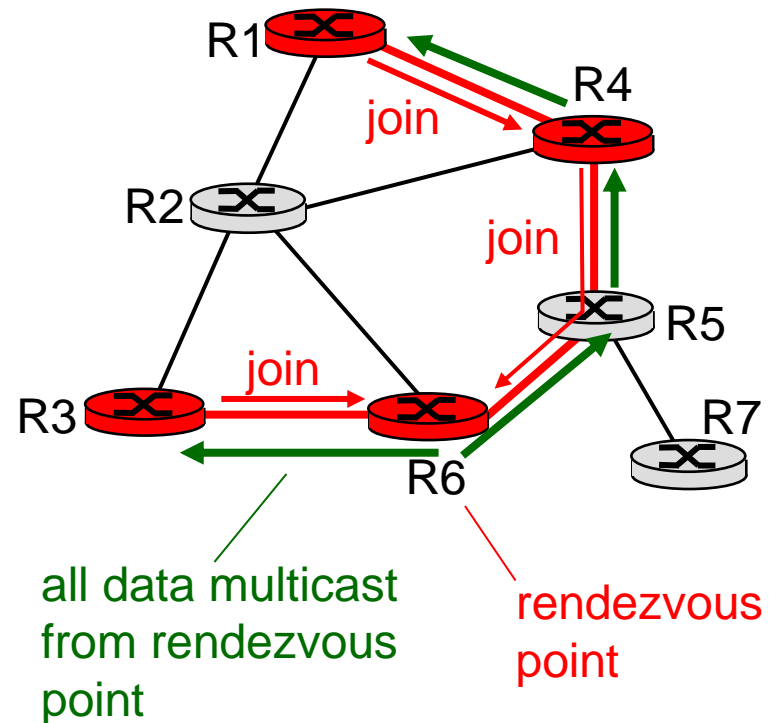
- center-based approach
- router sends *join* msg to rendezvous point (RP)
 - intermediate routers update state and forward *join*
- after joining via RP, router can switch to source-specific tree
 - increased performance: less concentration, shorter paths



PIM - Sparse Mode

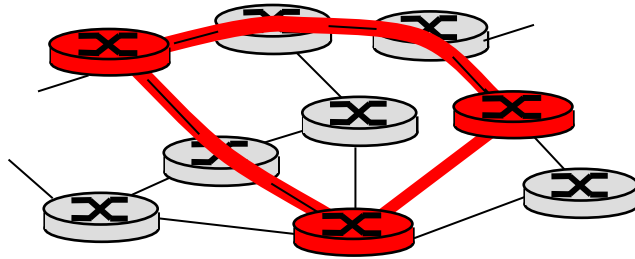
sender(s):

- unicast data to RP, which distributes down RP-rooted tree
- RP can extend mcast tree upstream to source
- RP can send *stop* msg if no attached receivers
 - “no one is listening!”

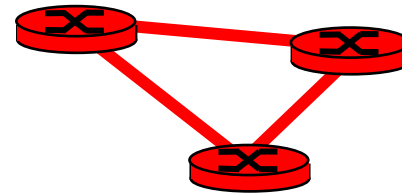


Tunneling

- Q: How to connect “islands” of multicast routers in a “sea” of unicast routers?



physical topology



logical topology

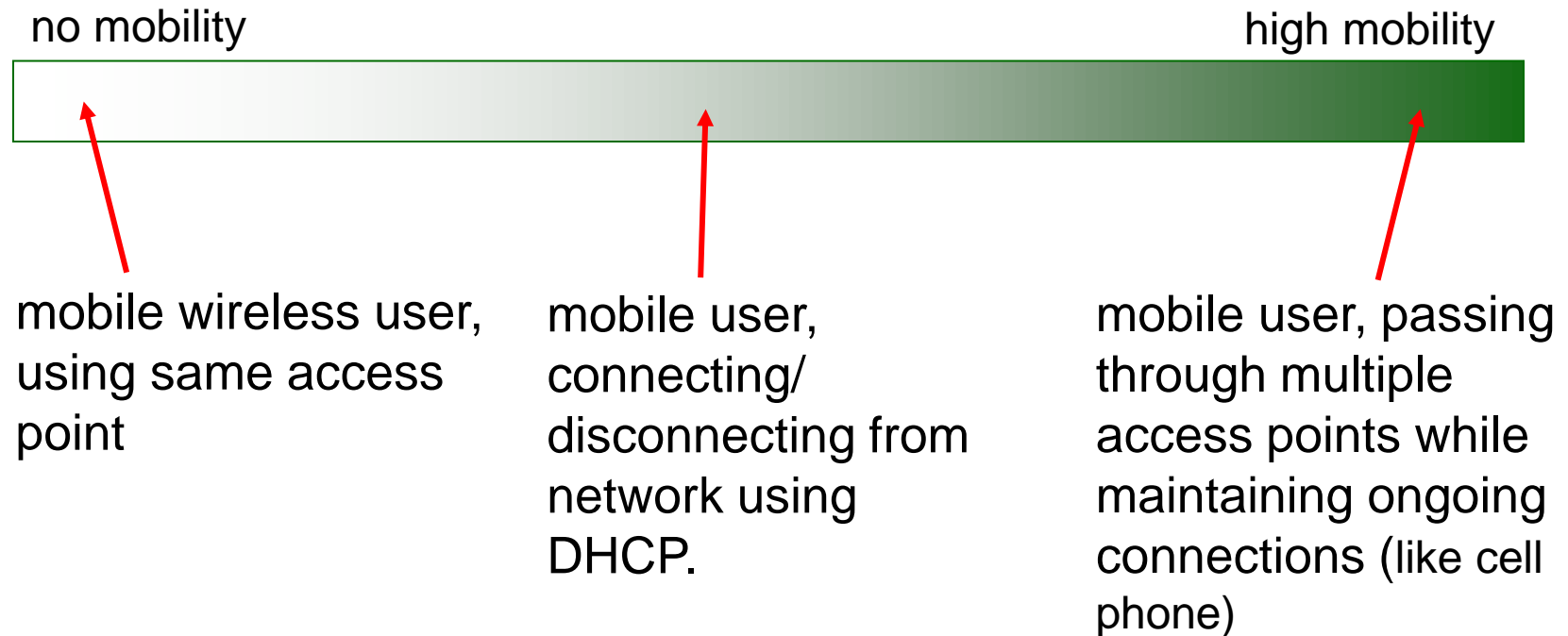
- mcast datagram encapsulated inside “normal” (non-multicast-addressed) datagram
- normal IP datagram sent through “tunnel” via regular IP unicast to receiving mcast router
- receiving mcast router de-capsulates pkt to get mcast datagram

Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - **What is Mobility?**
 - Network layer mobility concepts and principles
 - Mobile IP

What is mobility?

- spectrum of mobility, from the *network* perspective:

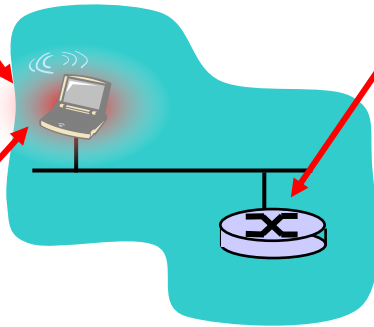


Mobility: Vocabulary

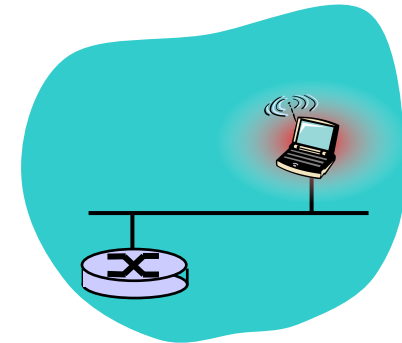
home network: permanent
“home” of mobile
(e.g., 128.119.40/24)

home agent: entity that will
perform mobility functions on
behalf of mobile, when mobile is
remote

Permanent address:
address in home
network, *can always* be
used to reach mobile
e.g., 128.119.40.186



wide area
network



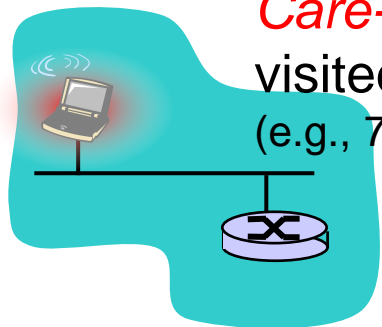
correspondent

Mobility: more vocabulary

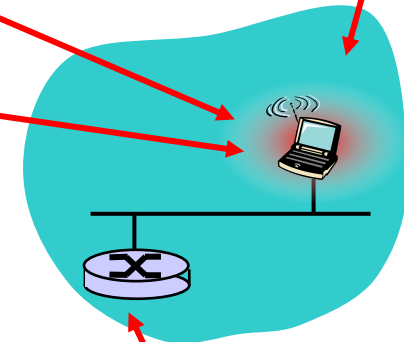
Permanent address: remains constant (e.g., 128.119.40.186)

visited network: network in which mobile currently resides (e.g., 79.129.13/24)

Care-of-address: address in visited network. (e.g., 79.129.13.2)

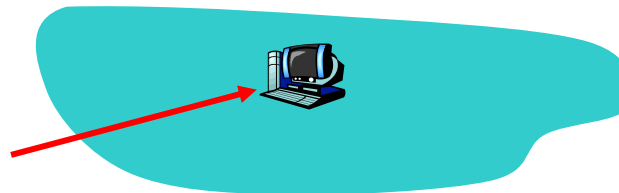


wide area network



foreign agent: entity in visited network that performs mobility functions on behalf of mobile.

correspondent: wants to communicate with

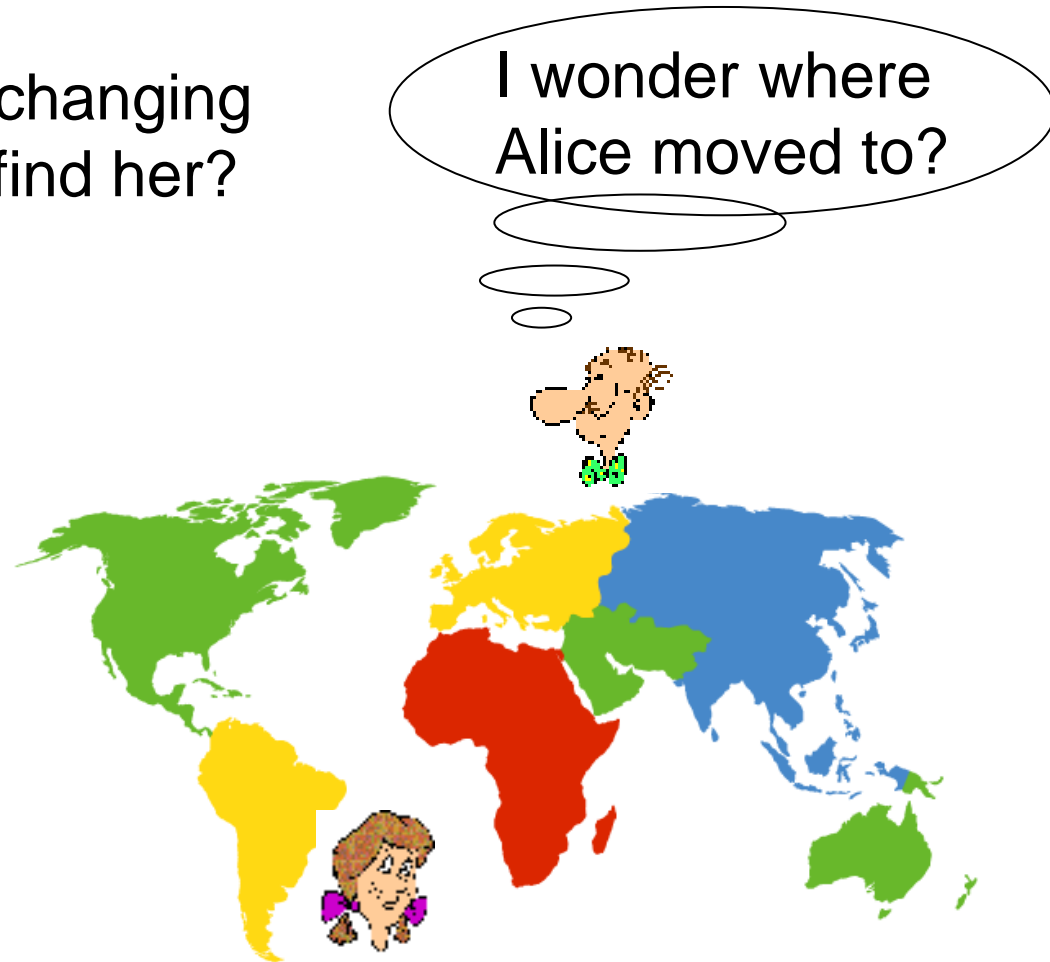


mobile

How do *you* contact a mobile friend:

Consider friend frequently changing addresses, how do you find her?

- search all phone books?
- call her parents?
- expect her to let you know where he/she is?



Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - **Network layer mobility concepts and principles**
 - Mobile IP

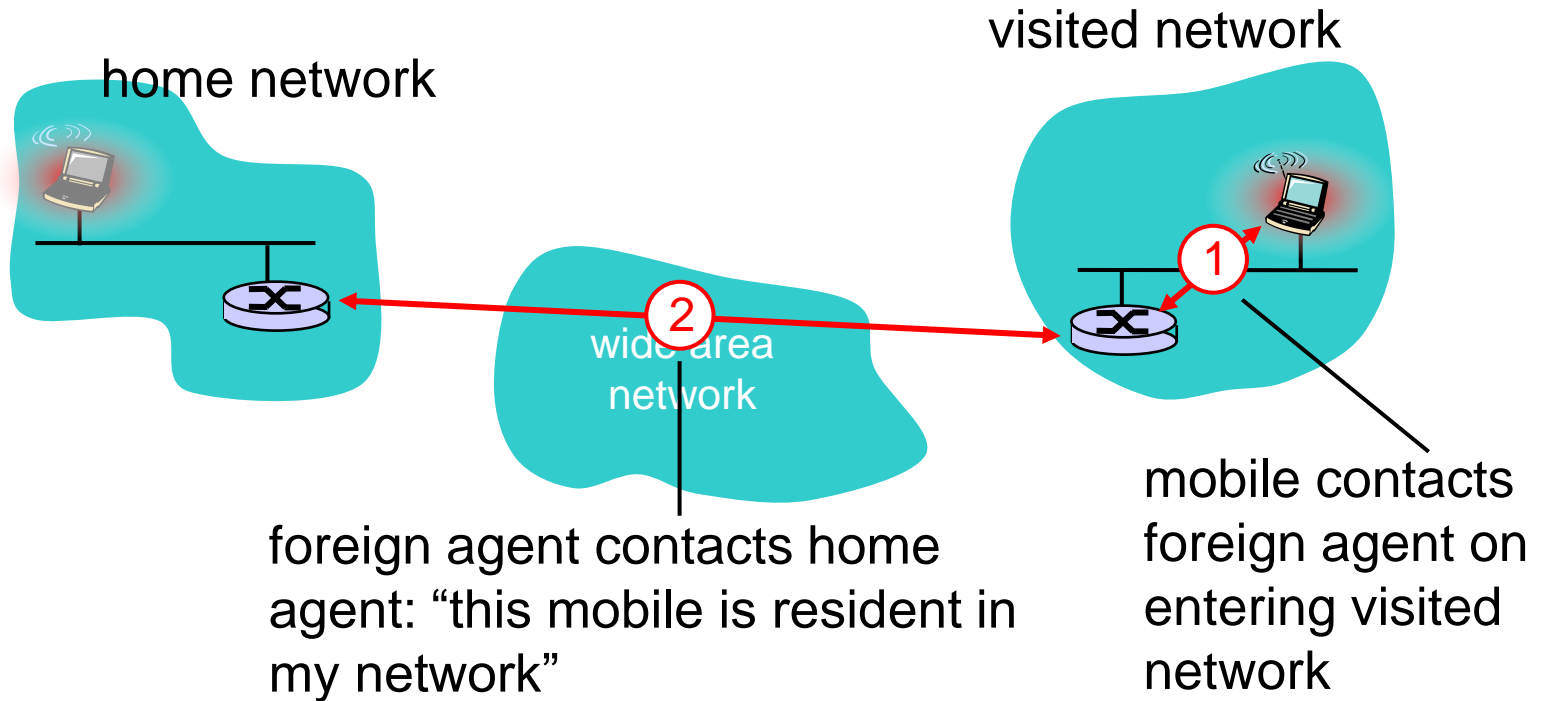
Mobility: approaches

- Let routing handle it
 - routers advertise permanent address of mobile-nodes via usual routing table exchange.
 - routing tables indicate where each mobile located
 - no changes to end-systems
 - **does not scale well!**

Mobility: approaches

- Let end-systems handle it
 - **Indirect routing:** communication from correspondent to mobile goes through home agent, then forwarded to remote
 - **Direct routing:** correspondent gets foreign address of mobile, sends directly to mobile

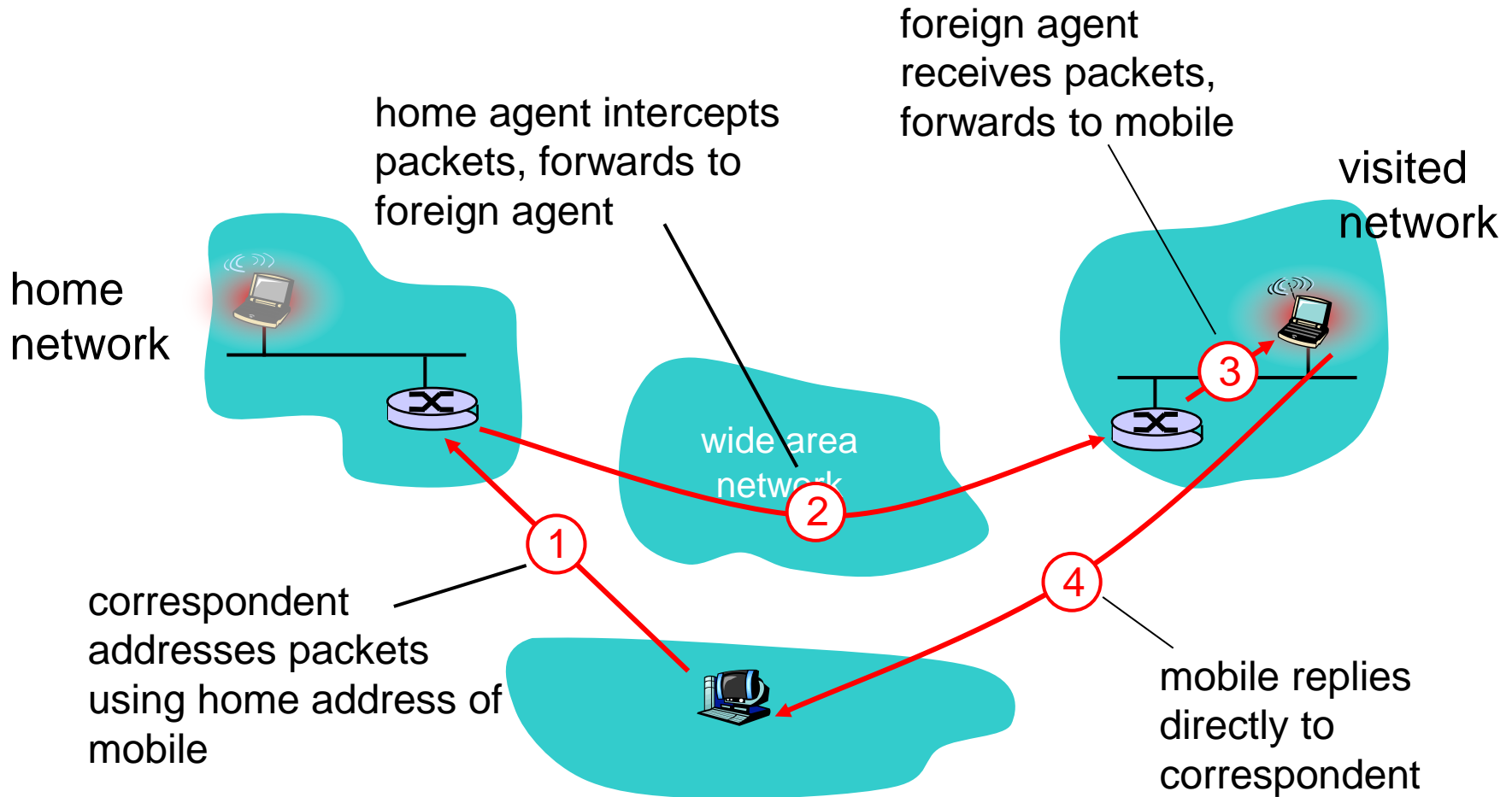
Mobility: registration



End result:

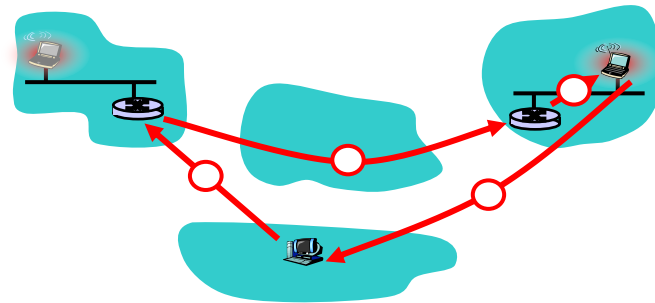
- Foreign agent knows about mobile
- Home agent knows location of mobile

Mobility via Indirect Routing



Indirect Routing: comments

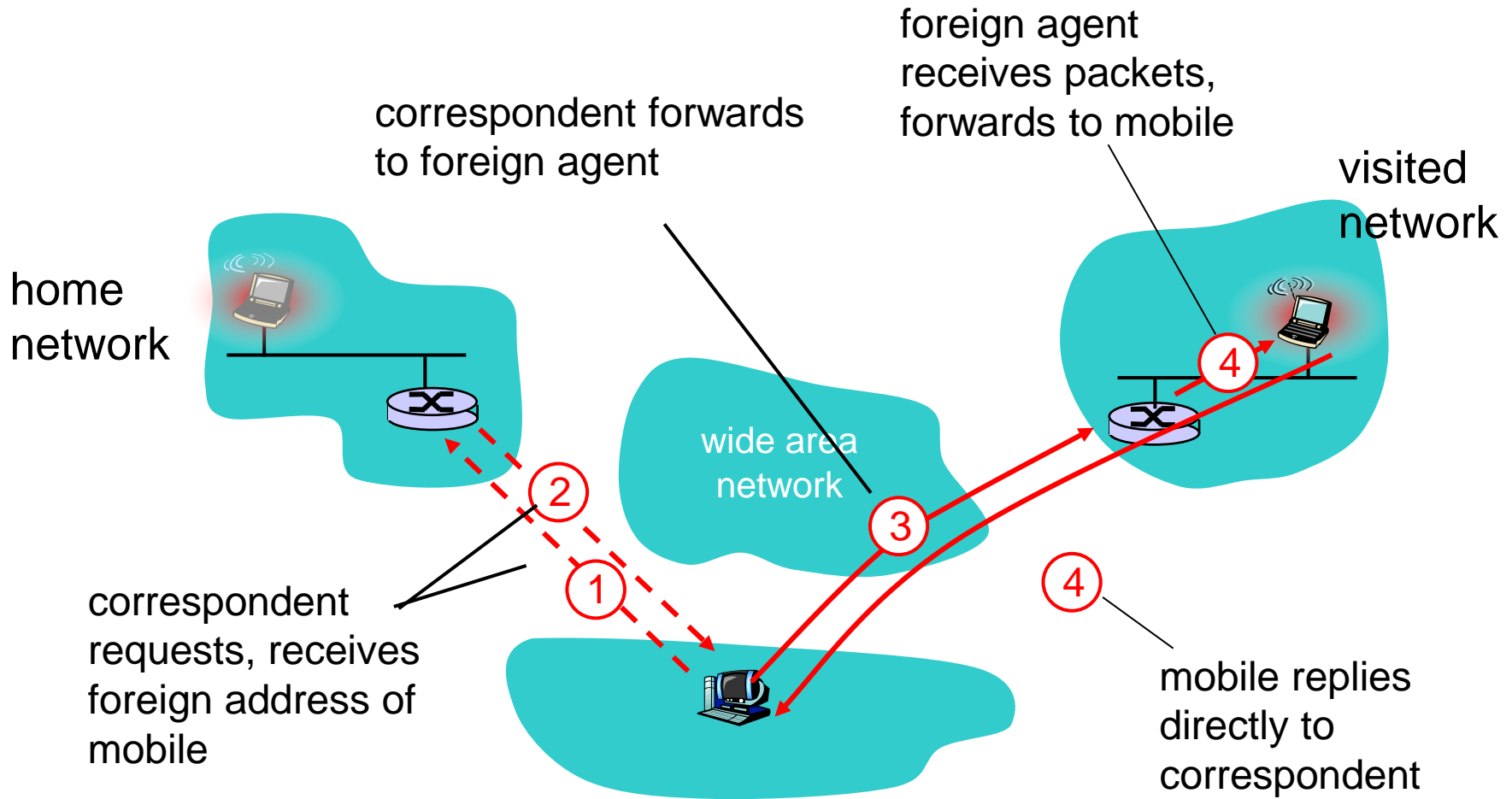
- Mobile uses two addresses:
 - **permanent address**: used by correspondent (hence mobile location is *transparent* to correspondent)
 - **care-of-address**: used by home agent to forward datagrams to mobile
- foreign agent functions may be done by mobile itself
- **triangle routing**: correspondent-home-network-mobile
 - inefficient when correspondent, mobile are in same network



Indirect Routing: moving between networks

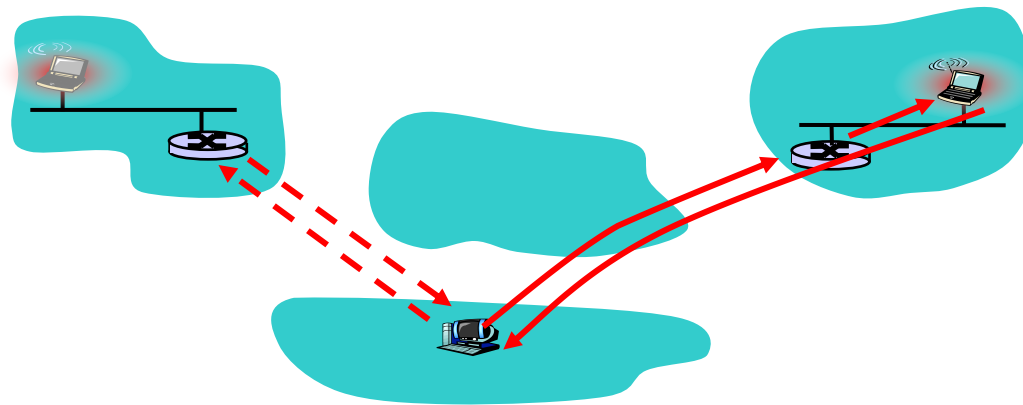
- suppose mobile user moves to another network
 - registers with new foreign agent
 - new foreign agent registers with home agent
 - home agent update care-of-address for mobile
 - packets continue to be forwarded to mobile (but with new care-of-address)
- mobility, changing foreign networks
transparent: *ongoing connections can be maintained!*

Mobility via Direct Routing



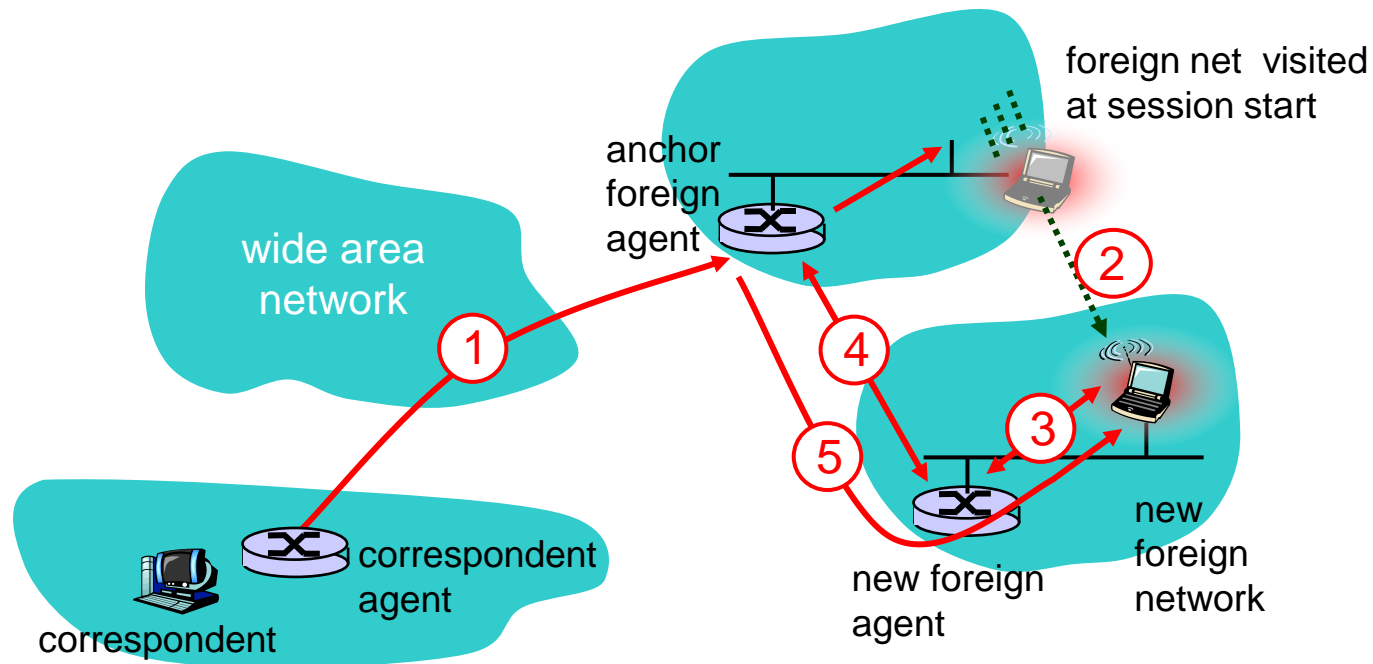
Mobility via Direct Routing: comments

- overcome triangle routing problem
- **non-transparent to correspondent:**
correspondent must get care-of-address from home agent
 - what if mobile changes visited network?



Accommodating mobility with direct routing

- anchor foreign agent: FA in first visited network
- data always routed first to anchor FA
- when mobile moves: new FA arranges to have data forwarded from old FA (chaining)



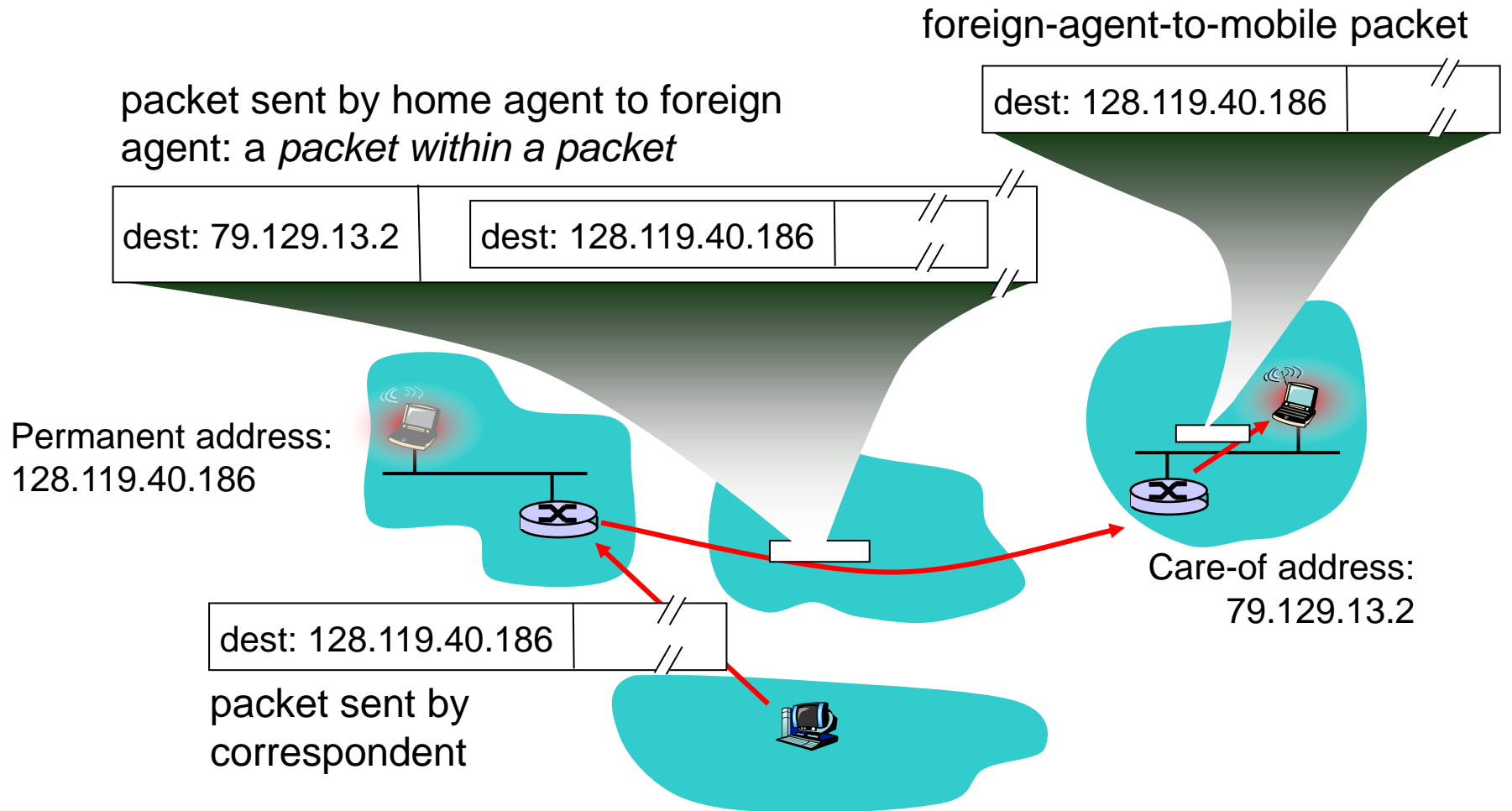
Network Layer II

- 4.5 Routing protocols
 - Routing Information Protocol (RIP)
 - Open Shortest Path First (OSPF)
 - Border Gateway Protocol (BGP)
- 4.6 Multicast
 - Broadcast routing
 - Multicast routing
 - Multicast routing protocols
- 4.7 Mobility
 - What is Mobility?
 - Network layer mobility concepts and principles
 - **Mobile IP**

Mobile IP

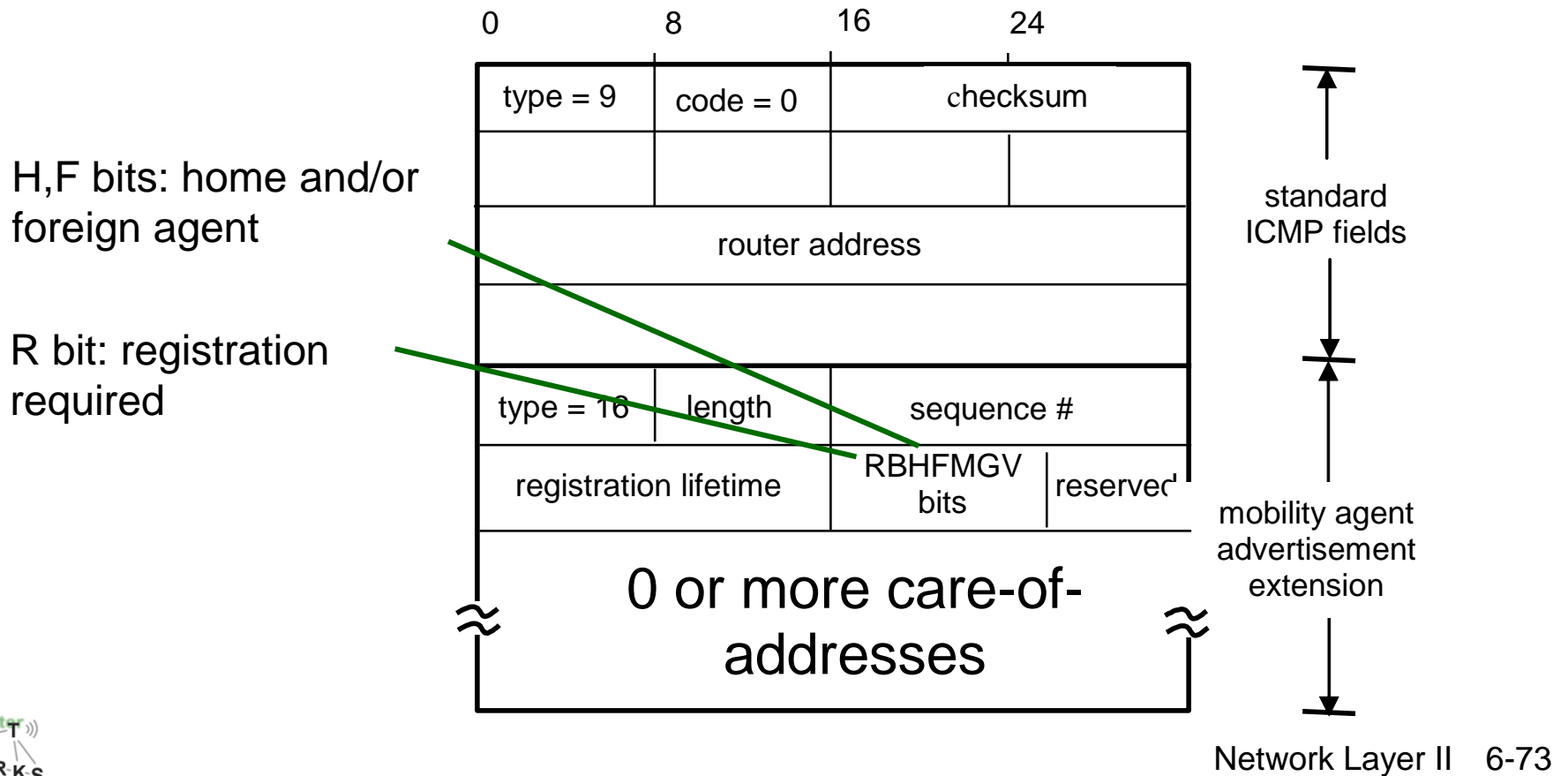
- RFC 3344
- has many features we've seen:
 - home agents, foreign agents, foreign-agent registration, care-of-addresses, encapsulation (packet-within-a-packet)
- three components to standard:
 - indirect routing of datagrams
 - agent discovery
 - registration with home agent

Mobile IP: indirect routing

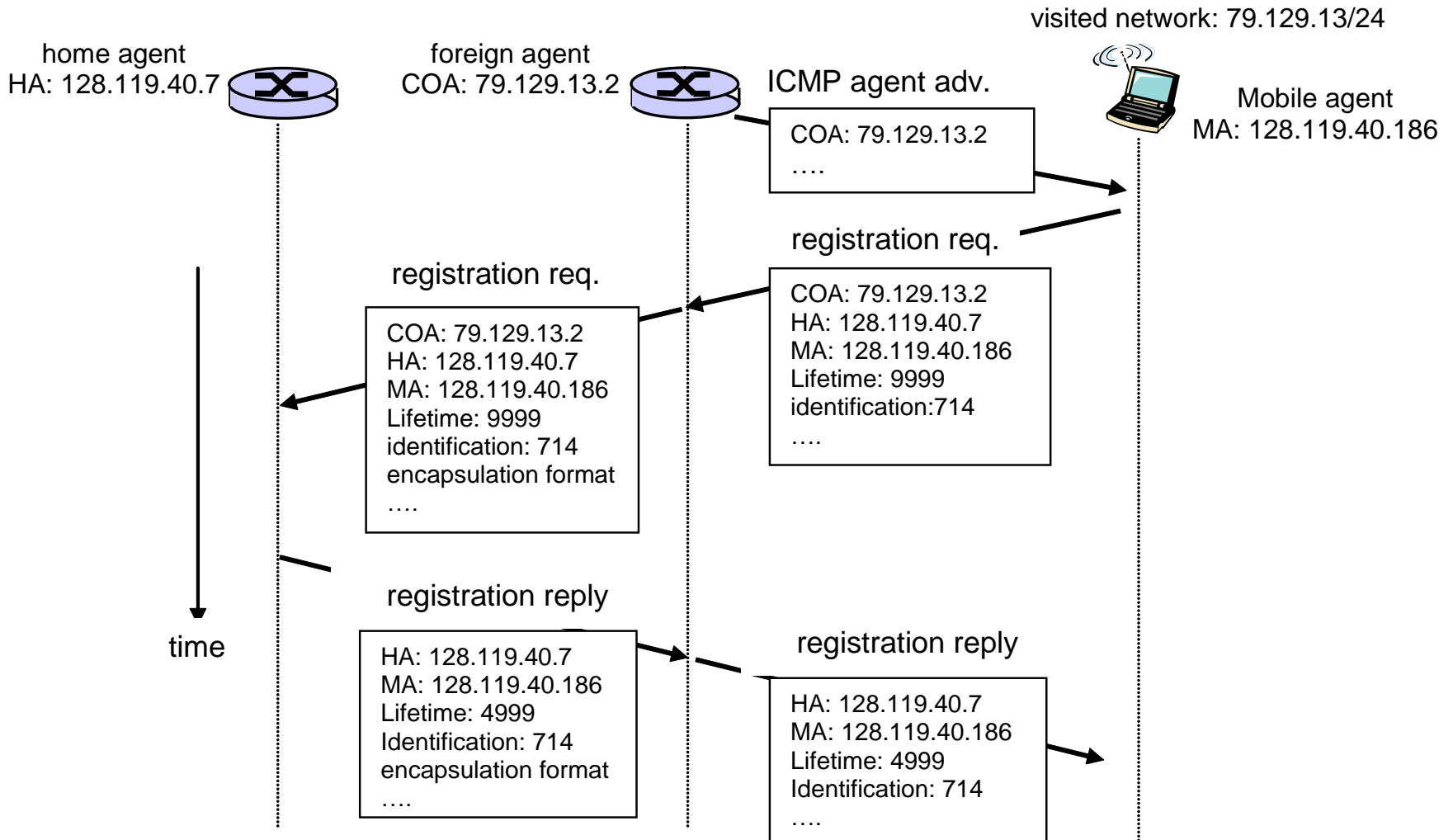


Mobile IP: agent discovery

- **agent advertisement:** foreign/home agents advertise service by broadcasting ICMP messages (typefield = 9)



Mobile IP: registration example



Thank you

Any questions?