

NETWORK LAYER PART 6

Inter Routing Protocol: BGP

Broadcast Routing

Multicast Routing: DVMRP

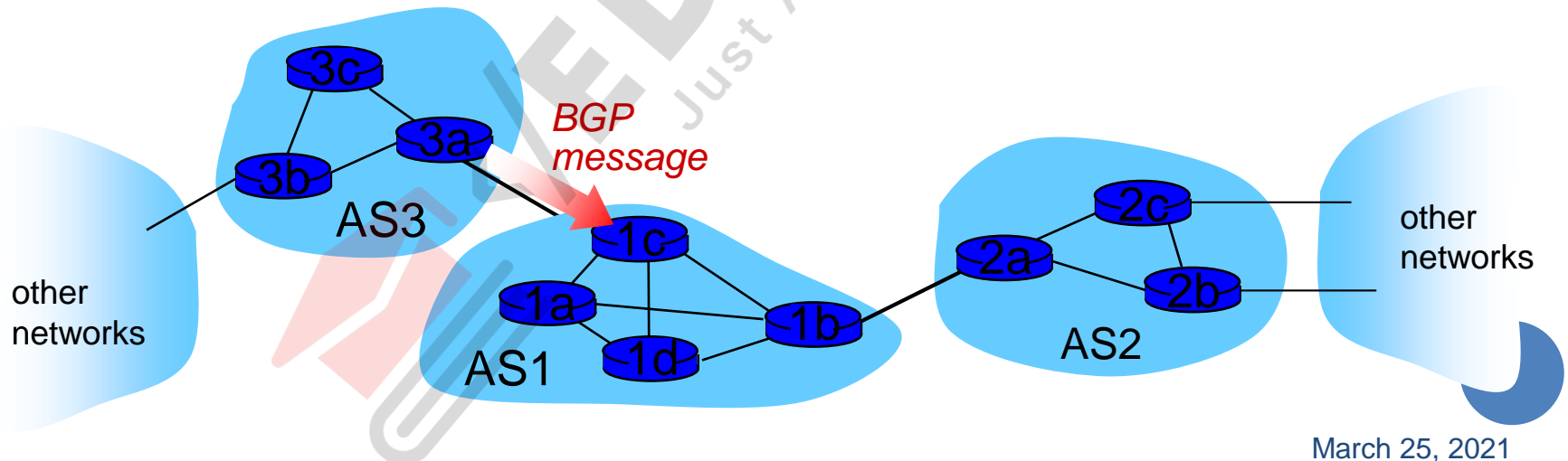
4-1

INTERNET INTER-AS ROUTING: BGP

- **BGP (Border Gateway Protocol):** *the* de facto inter-domain routing protocol
 - “glue that holds the Internet together”
- BGP provides each AS a means to:
 - **eBGP:** obtain subnet reachability information from neighboring ASs.
 - **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”*

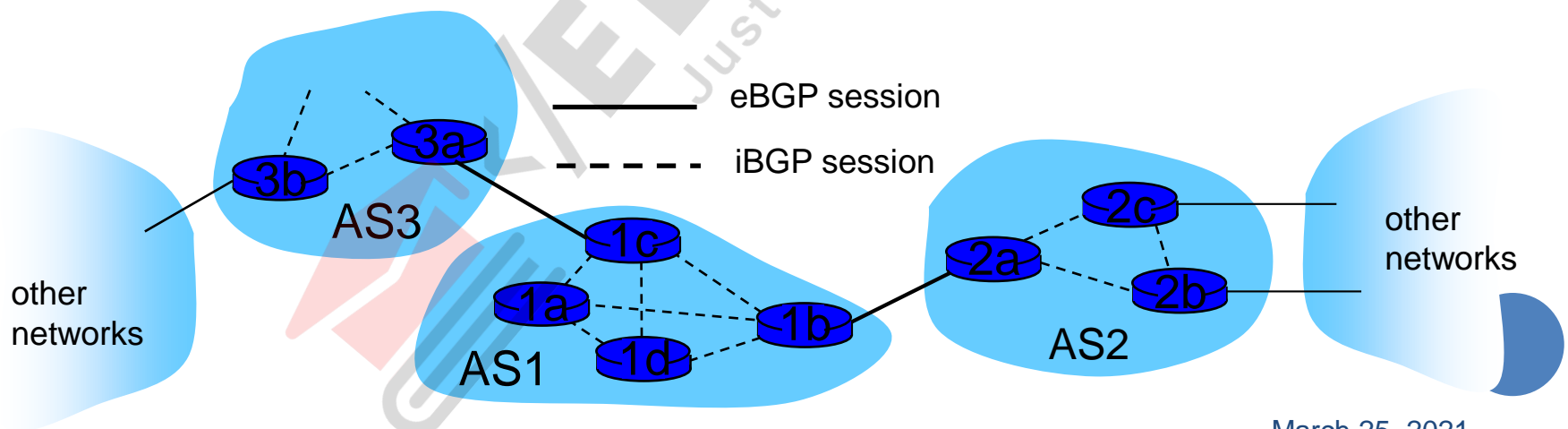
BGP BASICS

- ❖ **BGP session:** two BGP routers (“peers”) exchange BGP messages:
 - advertising *paths* to different destination network prefixes (“path vector” protocol)
 - exchanged over semi-permanent TCP connections
- when AS3 advertises a prefix to AS1:
 - AS3 *promises* it will forward datagrams towards that prefix
 - AS3 can aggregate prefixes in its advertisement



BGP BASICS: DISTRIBUTING PATH INFORMATION

- ❖ using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
 - 1c can then use iBGP to distribute new prefix info to all routers in AS1
 - 1b can then re-advertise new reachability info to AS2 over 1b-to-2a eBGP session
- ❖ when router learns of new prefix, it creates entry for prefix in its forwarding table.



PATH ATTRIBUTES AND BGP ROUTES

- advertised prefix includes BGP attributes
 - prefix + attributes = “route”
- two important attributes:
 - **AS-PATH**: contains ASs through which prefix advertisement has passed: e.g., AS 67, AS 17
 - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- gateway router receiving route advertisement uses **import policy** to accept/decline
 - e.g., never route through AS x
 - *policy-based* routing

BGP ROUTE SELECTION

- ❖ router may learn about more than 1 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

BGP MESSAGES

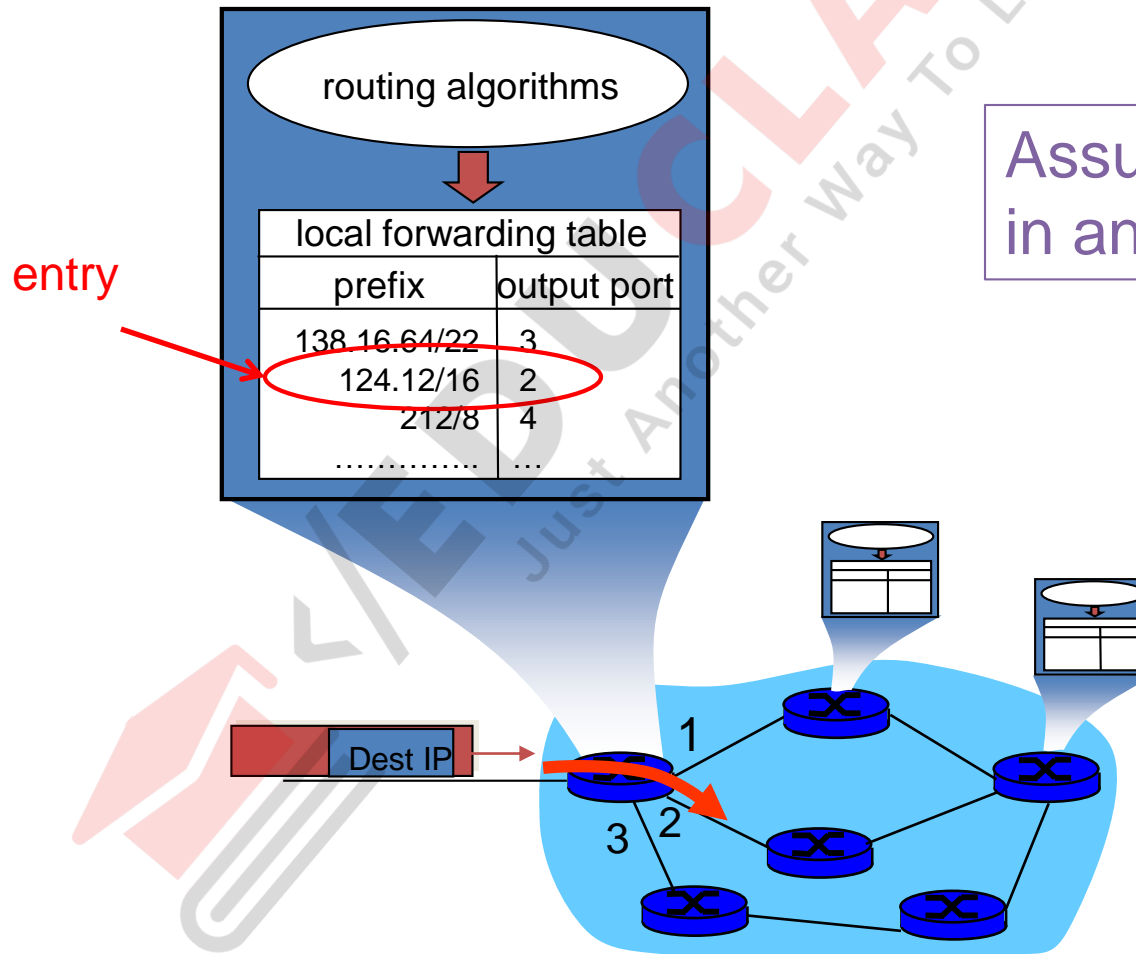
- BGP messages exchanged between peers over TCP connection
- BGP messages:
 - **OPEN:** opens TCP connection to peer and authenticates sender
 - **UPDATE:** advertises new path (or withdraws old)
 - **KEEPALIVE:** keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - **NOTIFICATION:**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

PUTTING IT ALTOGETHER: *HOW DOES AN ENTRY GET INTO A ROUTER'S FORWARDING TABLE?*

- Answer is complicated!
- Ties together hierarchical routing (Section 4.5.3) with BGP (4.6.3) and OSPF (4.6.2).
- Provides nice overview of BGP!

HOW DOES ENTRY GET IN FORWARDING TABLE?

March 25, 2021



How does entry get in forwarding table?

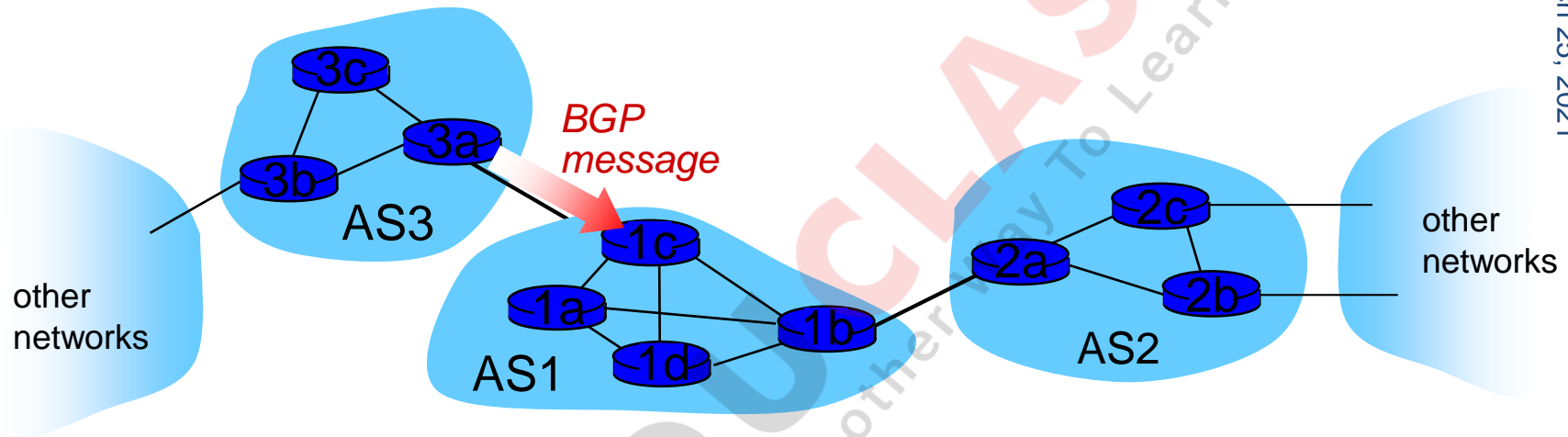
March 25, 2021

High-level overview

1. Router becomes aware of prefix
2. Router determines output port for prefix
3. Router enters prefix-port in forwarding table



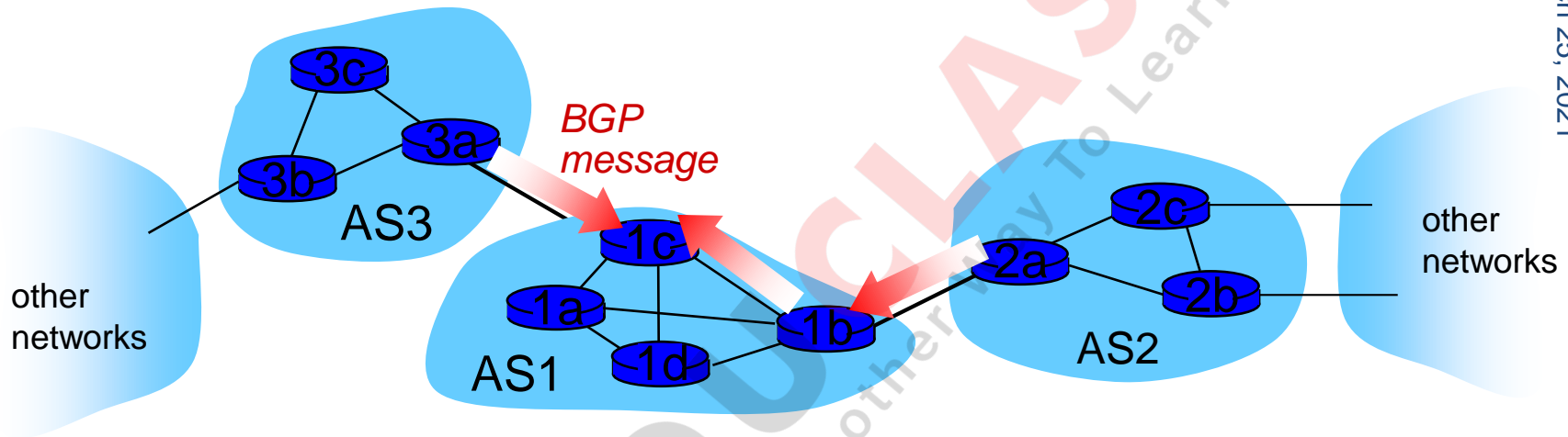
ROUTER BECOMES AWARE OF PREFIX



- ❖ BGP message contains “routes”
- ❖ “route” is a prefix and attributes: AS-PATH, NEXT-HOP,...
- ❖ Example: route:
 - ❖ Prefix: 138.16.64/22 ; AS-PATH: AS3 AS131 ; NEXT-HOP: 201.44.13.125

ROUTER MAY RECEIVE MULTIPLE ROUTES

March 25, 2021



- ❖ Router may receive multiple routes for same prefix
- ❖ Has to select one route

SELECT BEST BGP ROUTE TO PREFIX

March 25, 2021

- Router selects route based on shortest AS-PATH

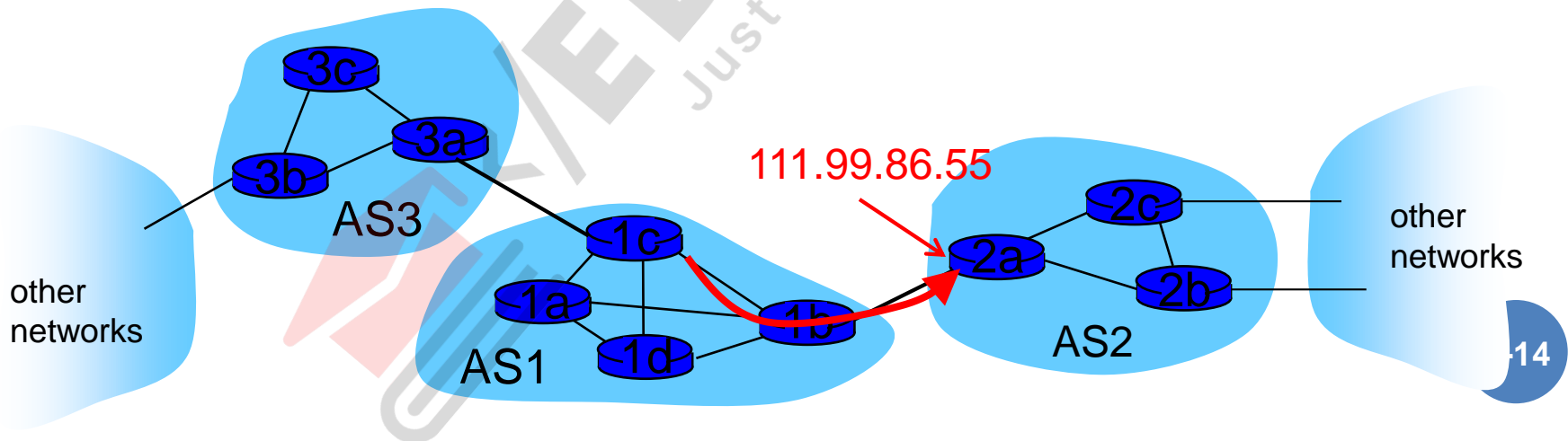
❖ Example:

- ❖ AS2 AS17 to 138.16.64/22 select
- ❖ AS3 AS131 AS201 to 138.16.64/22
- ❖ What if there is a tie? We'll come back to that!

FIND BEST INTRA-ROUTE TO BGP ROUTE

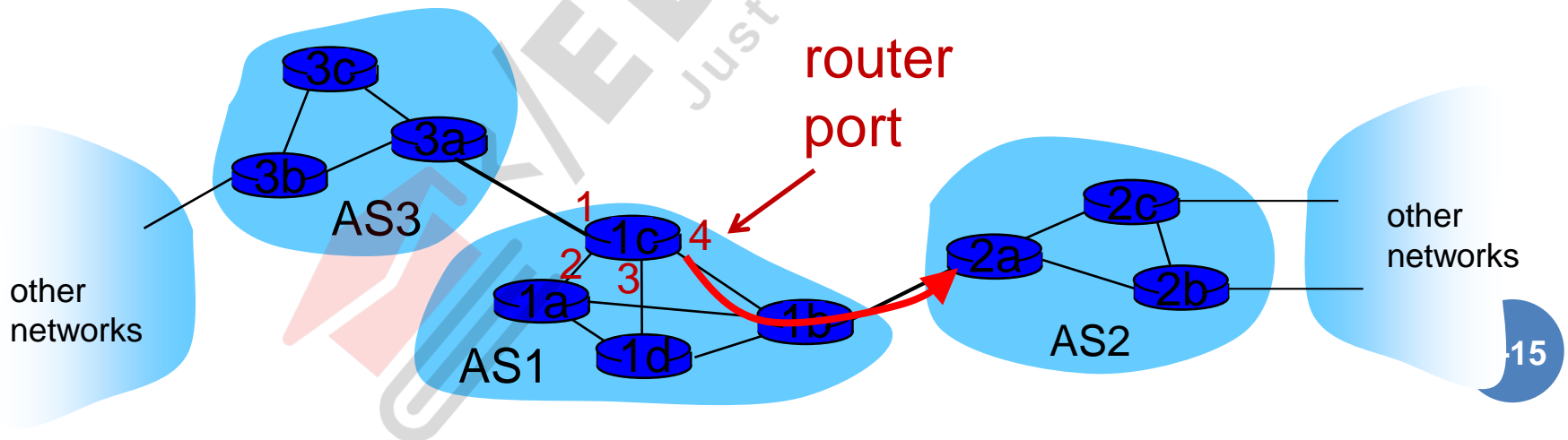
March 25, 2021

- Use selected route's NEXT-HOP attribute
 - Route's NEXT-HOP attribute is the IP address of the router interface that begins the AS PATH.
- Example:
 - ❖ AS-PATH: AS2 AS17 ; NEXT-HOP: 111.99.86.55
- Router uses OSPF to find shortest path from 1c to 111.99.86.55



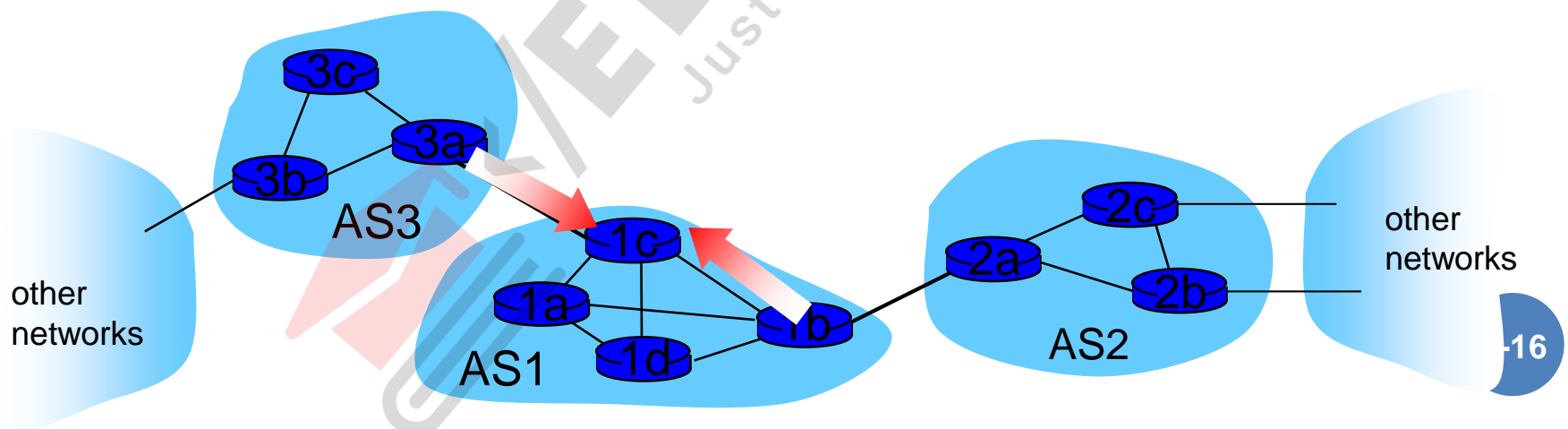
ROUTER IDENTIFIES PORT FOR ROUTE

- Identifies port along the OSPF shortest path
- Adds prefix-port entry to its forwarding table:
 - (138.16.64/22 , port 4)



HOT POTATO ROUTING

- Suppose there two or more best inter-routes.
- Then choose route with closest NEXT-HOP
 - Use OSPF to determine which gateway is closest
 - Q: From 1c, chose AS3 AS131 or AS2 AS17?
 - A: route AS3 AS201 since it is closer



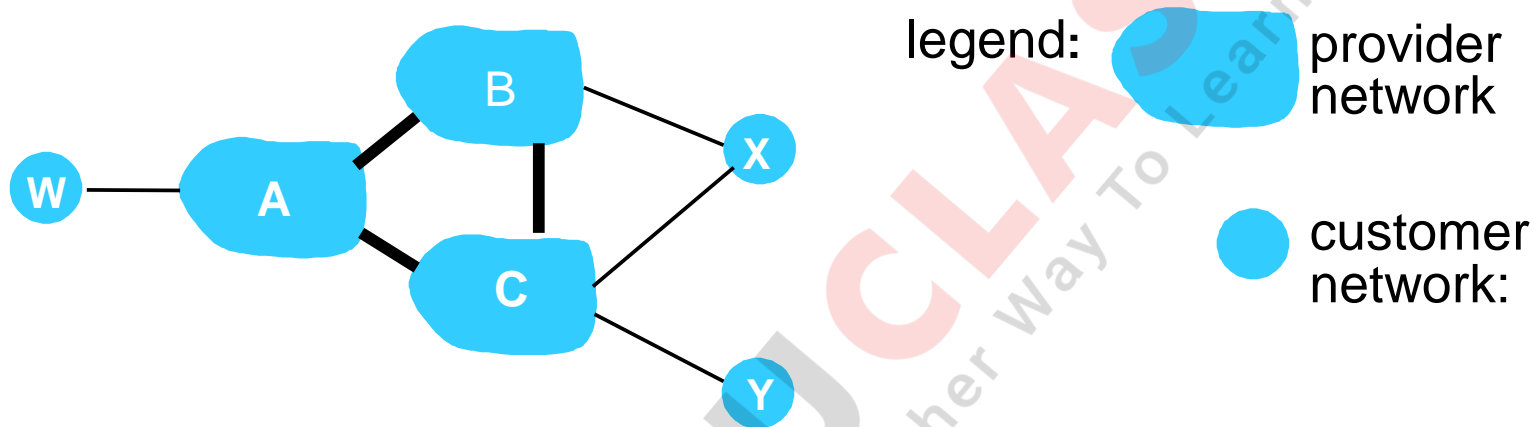
HOW DOES ENTRY GET IN FORWARDING TABLE?

March 25, 2021

Summary

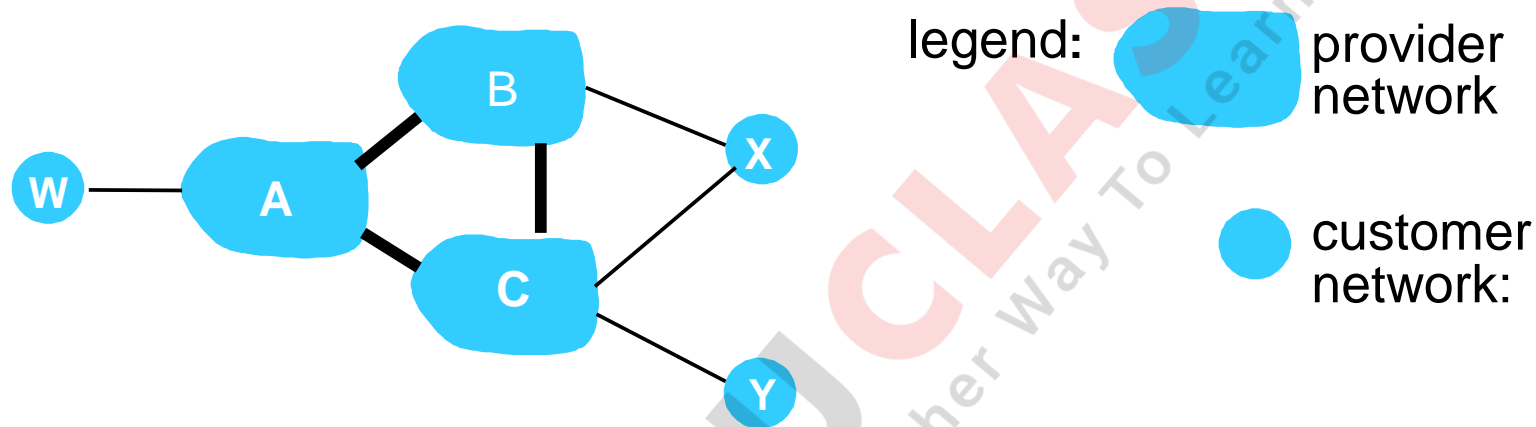
1. Router becomes aware of prefix
 - via BGP route advertisements from other routers
2. Determine router output port for prefix
 - Use BGP route selection to find best inter-AS route
 - Use OSPF to find best intra-AS route leading to best inter-AS route
 - Router identifies router port for that best route
3. Enter prefix-port entry in forwarding table

BGP ROUTING POLICY



- ❖ A,B,C are *provider networks*
- ❖ X,W,Y are customer (of provider 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-, 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:

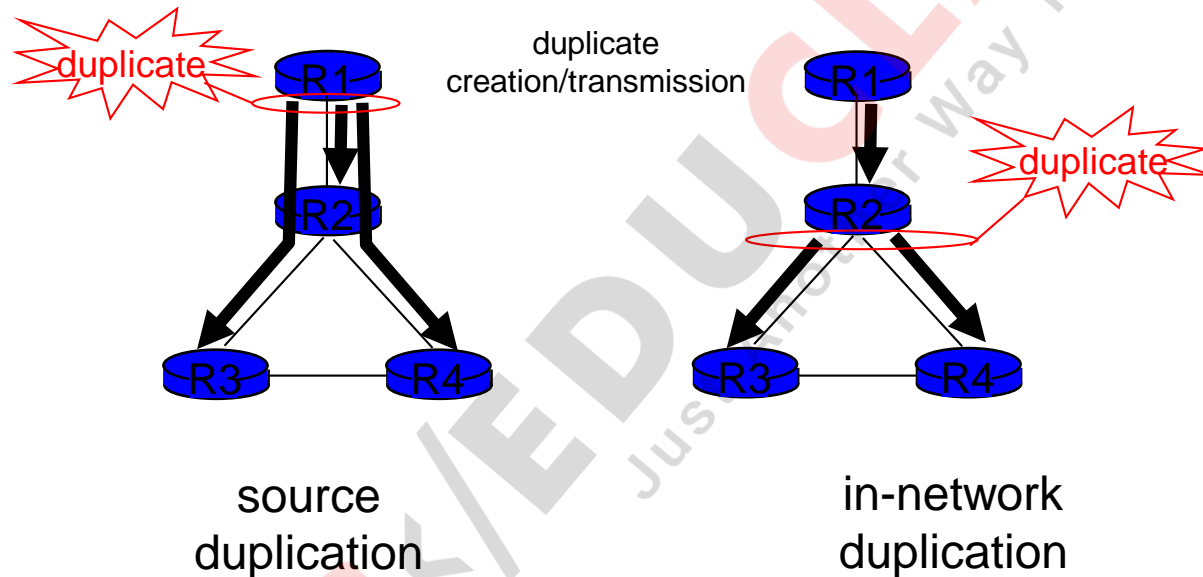
- hierarchical routing saves table size, reduced update traffic

performance:

- intra-AS: can focus on performance
- inter-AS: policy may dominate over performance

BROADCAST ROUTING

- ❖ deliver packets from source to all other nodes
- ❖ source duplication is inefficient:



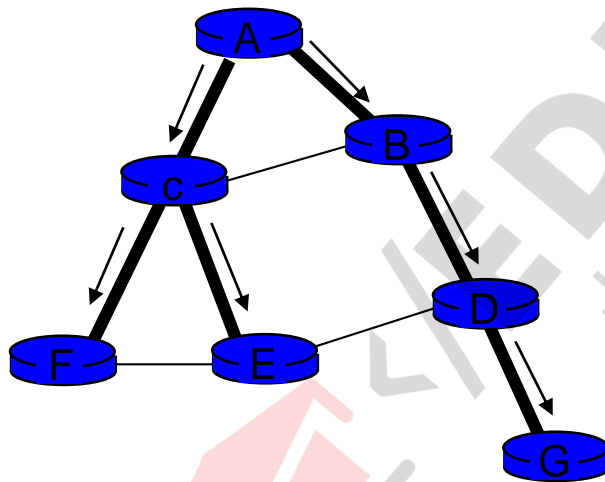
- ❖ source duplication: how does source determine recipient addresses?

IN-NETWORK DUPLICATION

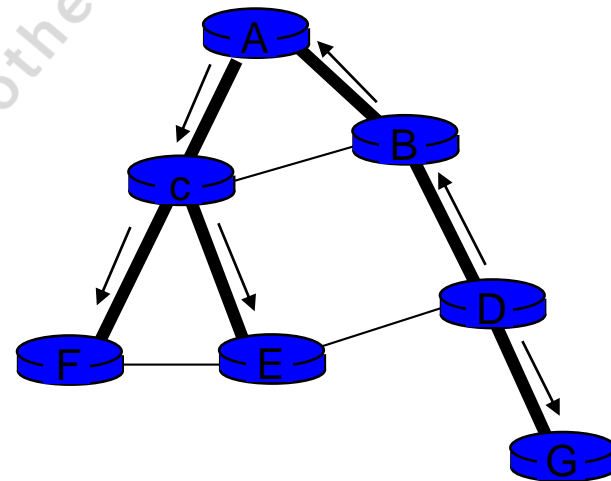
- *flooding*: when node receives broadcast packet, sends copy to all neighbors
 - problems: cycles & broadcast storm
- *controlled flooding*: node only broadcasts pkt if it hasn't broadcast same packet before
 - *Sequence number controlled flooding* : node keeps track of packet ids already broadcasted
 - *Reverse path forwarding (RPF)*: only forward packet 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 then forward/make copies only along spanning tree



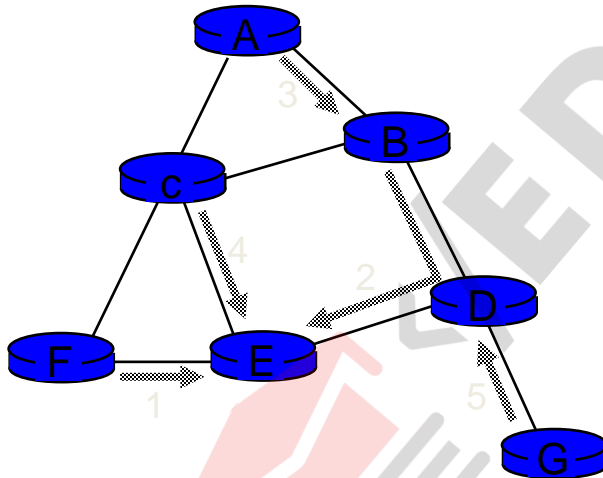
(a) broadcast initiated at A



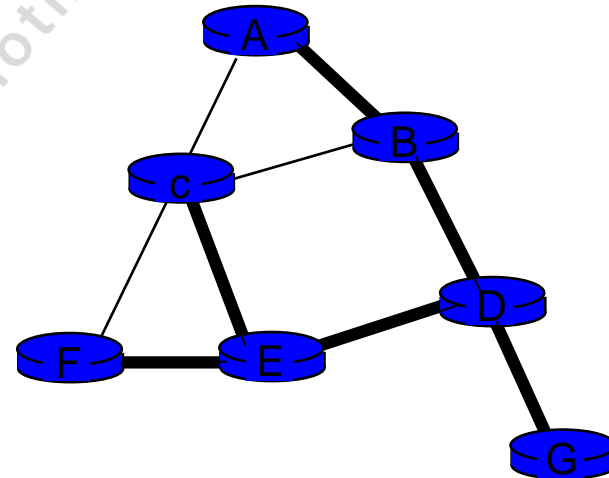
(b) broadcast initiated at D

SPANNING TREE: CREATION

- ❖ center node
- ❖ each node sends unicast join message to center node
 - message forwarded until it arrives at a node already belonging to spanning tree



(a) stepwise construction of spanning tree (center: E)



(b) constructed spanning tree

MULTICAST ROUTING: PROBLEM STATEMENT

goal: find a tree (or trees) connecting routers having local mcast group members

- ❖ *tree:* not all paths between routers used
- ❖ *shared-tree:* same tree used by all group members
- ❖ *source-based:* different tree from each sender to rcvrs

legend



group member



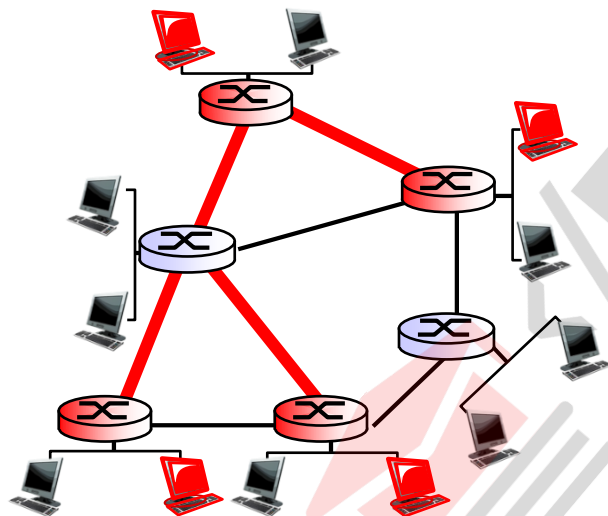
not group member



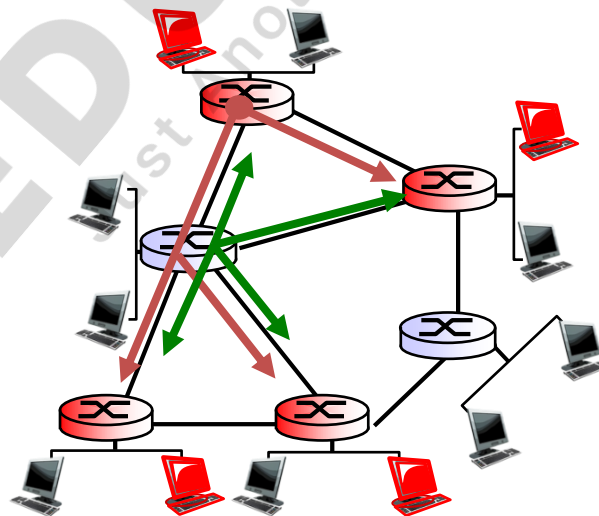
router with a group member



router without group member



shared tree



source-based trees

APPROACHES FOR BUILDING MCAST TREES

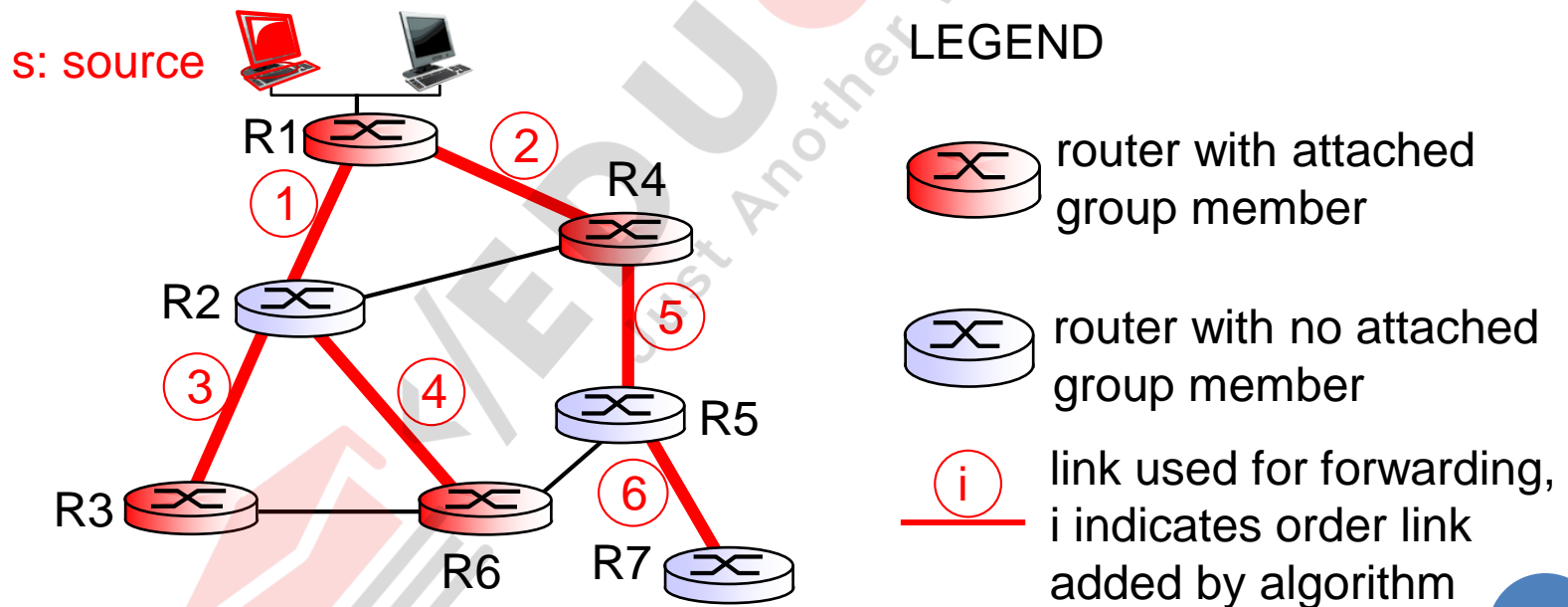
approaches:

- ❖ *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

...we first look at basic approaches, then specific protocols adopting these approaches

SHORTEST PATH TREE

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

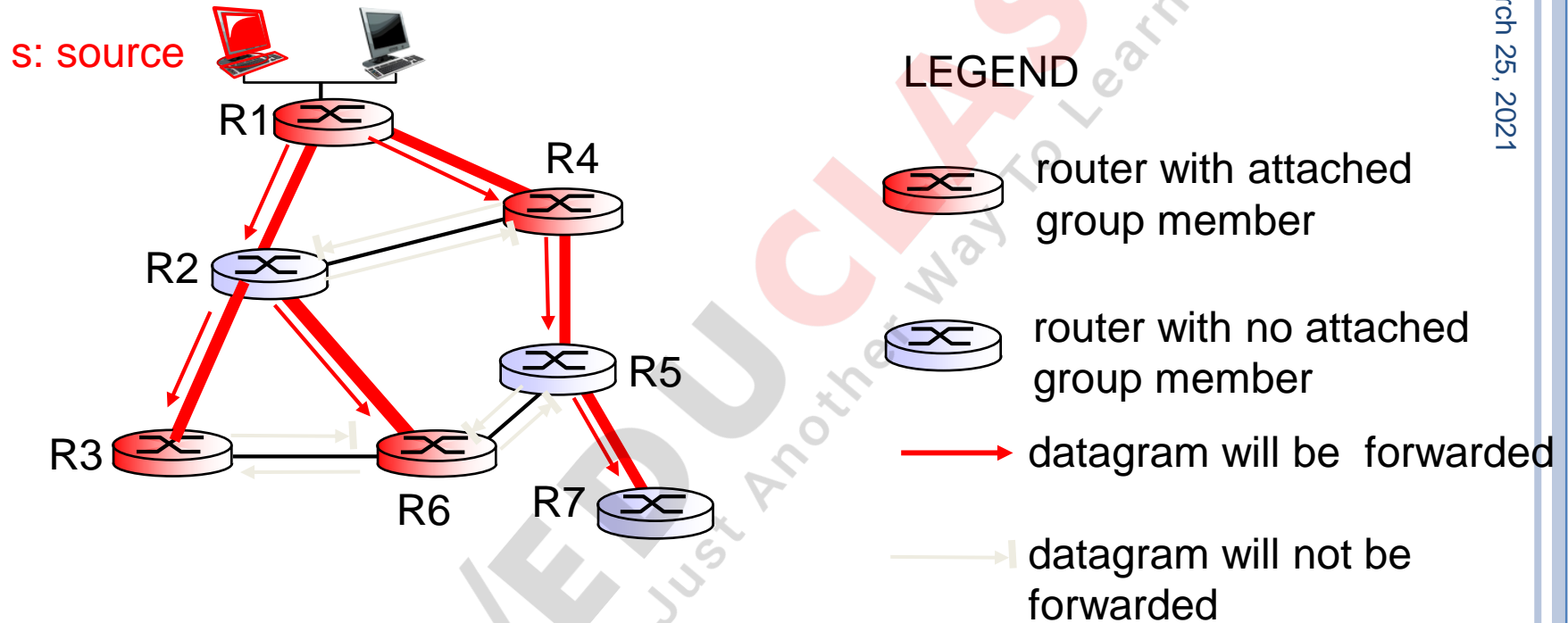


REVERSE PATH FORWARDING

- ❖ rely on router's knowledge of unicast shortest path from it to sender
- ❖ each router has simple forwarding behavior:

if (mcast datagram received on incoming link
on shortest path back to center)
then flood datagram onto all outgoing links
else ignore datagram

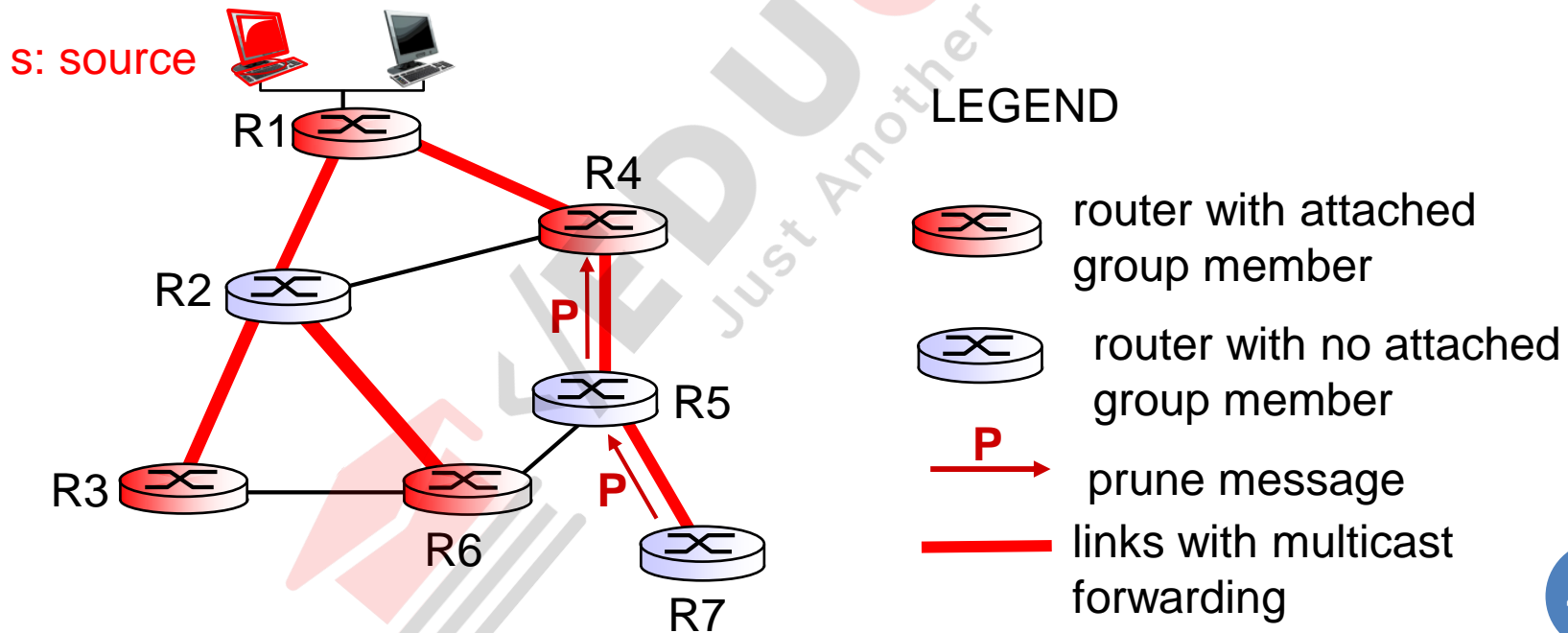
REVERSE PATH FORWARDING: EXAMPLE



- ❖ 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 mcast group members
 - no need to forward datagrams down subtree
 - “prune” msgs sent upstream by router with no downstream group members



SHARED-TREE: STEINER TREE

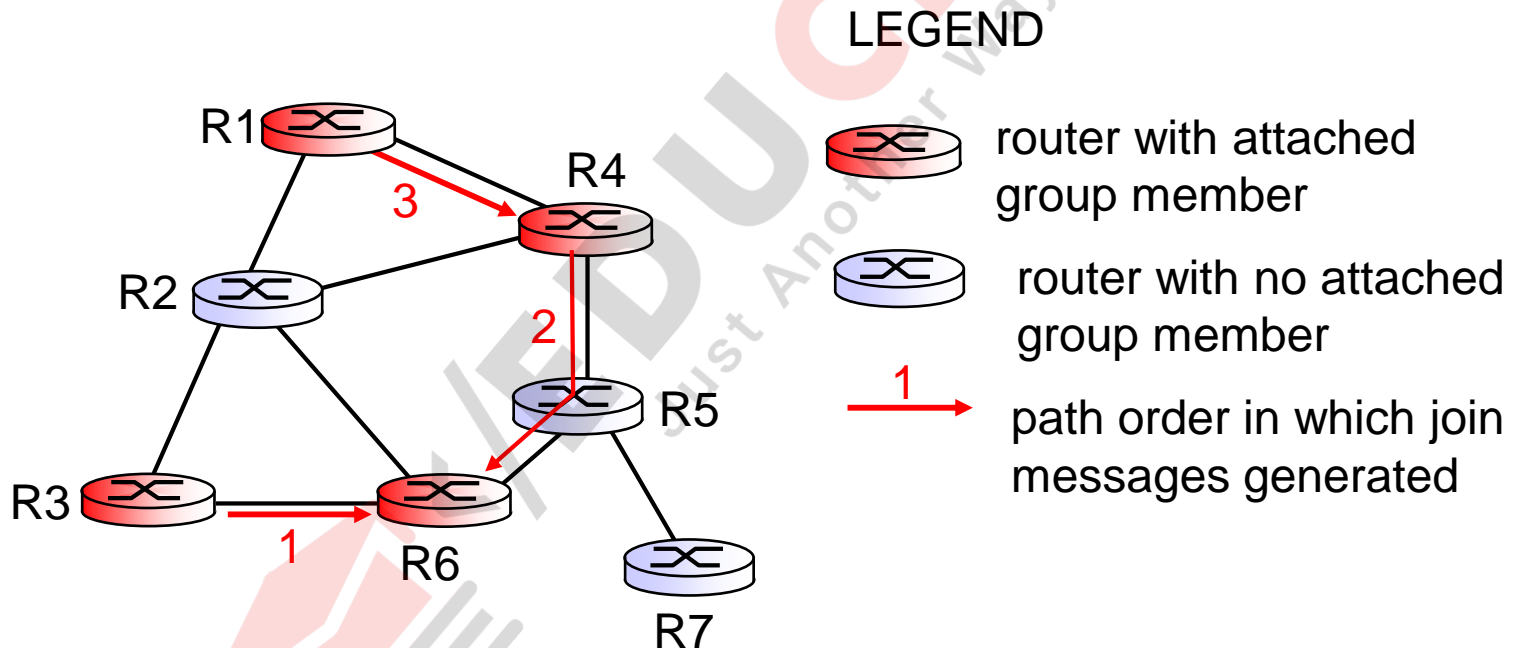
- ❖ *steiner tree*: minimum cost tree connecting all routers with attached group members
- ❖ problem is NP-complete
- ❖ excellent heuristics exists
- ❖ not used in practice:
 - computational complexity
 - information about entire network needed
 - monolithic: rerun whenever a router needs to join/leave

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: EXAMPLE

suppose R6 chosen as center:



INTERNET MULTICASTING ROUTING: DVMRP

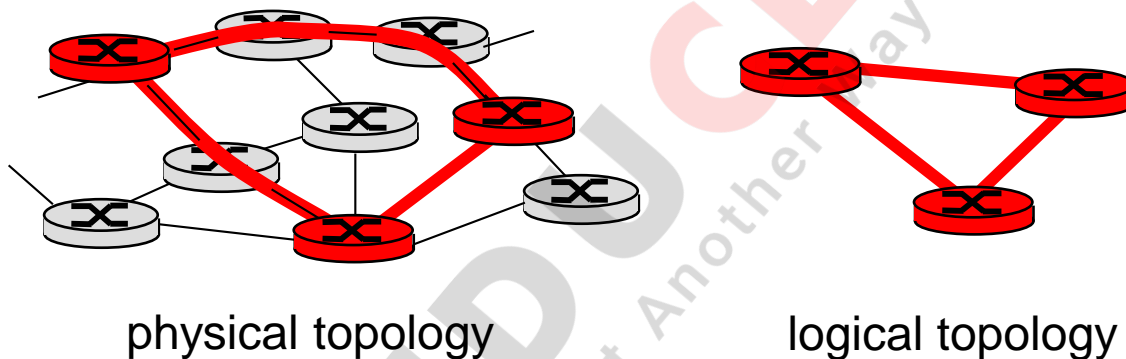
- **DVMRP**: distance vector multicast routing protocol, RFC1075
- *flood and prune*: reverse path forwarding, source-based tree
 - RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers
 - no assumptions about underlying unicast
 - initial datagram to mcast group flooded 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 router

TUNNELING

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



- ❖ mcast datagram encapsulated inside “normal” (non-multicast-addressed) datagram
- ❖ normal IP datagram sent thru “tunnel” via regular IP unicast to receiving mcast router (recall IPv6 inside IPv4 tunneling)
- ❖ receiving mcast router unencapsulates to get mcast datagram