

Chapter 4

Network Layer

Reti di Elaboratori

Corso di Laurea in Informatica

Università degli Studi di Roma "La Sapienza"

Canale A-L

Prof.ssa Chiara Petrioli

Parte di queste slide sono state prese dal materiale associato al libro
Computer Networking: A Top Down Approach, 5th edition.

All material copyright 1996-2009

J.F Kurose and K.W. Ross, All Rights Reserved

Thanks also to Antonio Capone, Politecnico di Milano, Giuseppe Bianchi and
Francesco LoPresti, Un. di Roma Tor Vergata

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6
- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing

OSPF (Open Shortest Path First)

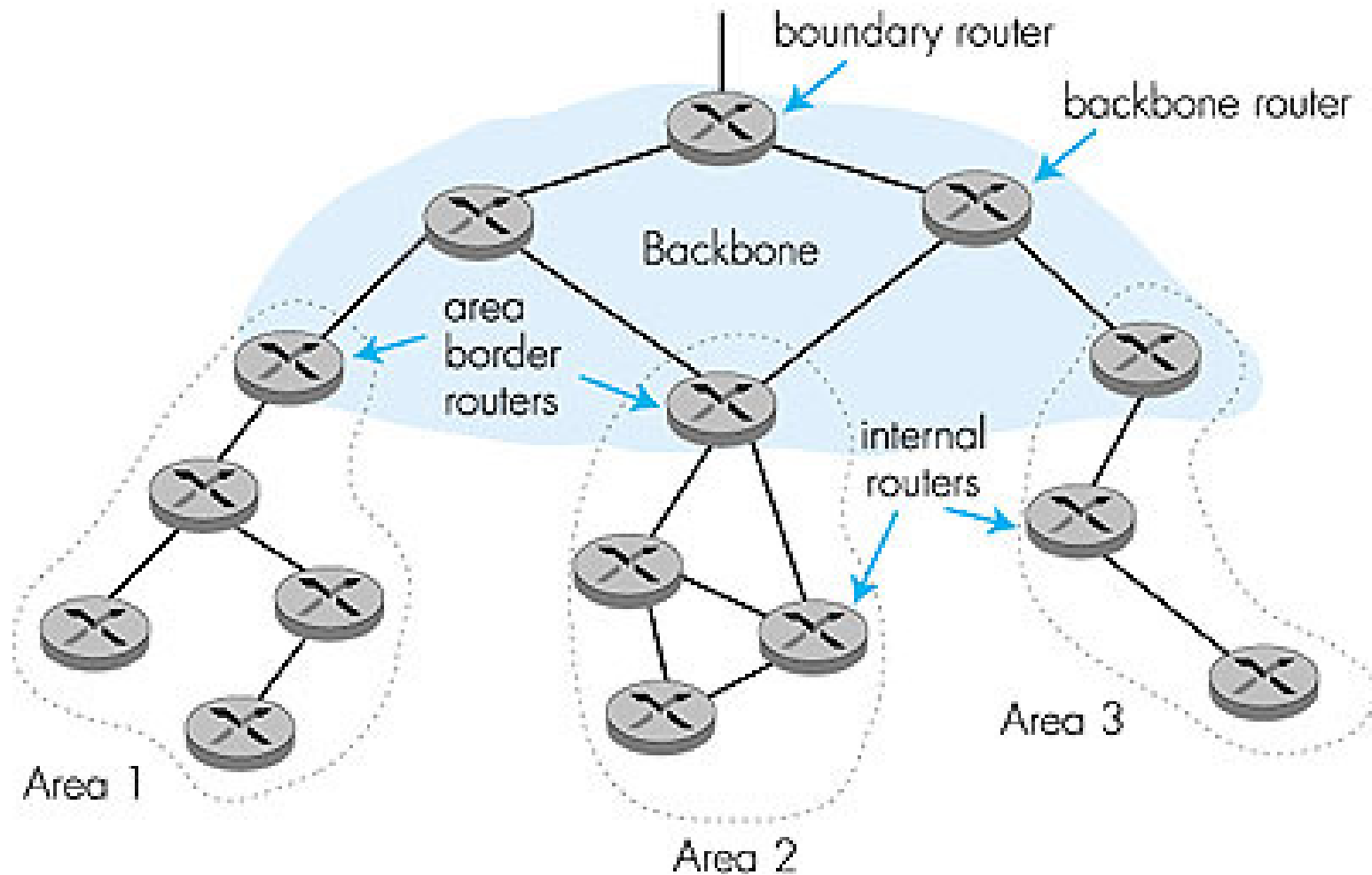
- r “open”: publicly available
- r uses Link State algorithm
 - m LS packet dissemination
 - m topology map at each node
 - m route computation using Dijkstra's algorithm
- r OSPF advertisement carries one entry per neighbor router
 - m Each node disseminates its local view of the topology
 - m i.e., the router usable interfaces and reachable neighbors
- r advertisements disseminated to **entire** AS (via flooding)
 - m carried in OSPF messages directly over IP (using protocol number 89)

OSPF "advanced" features (not in RIP)

- r **security**: all OSPF messages authenticated (to prevent malicious intrusion)
- r **multiple** same-cost **paths** allowed (only one path in RIP)
- r For each link, multiple cost metrics for different **TOS** (e.g., satellite link cost set "low" for best effort; high for real time)
 - m Periodic updates (30 min) or event driven (link cost change)
- r integrated uni- and **multicast** support:
 - m Multicast OSPF (MOSPF) uses same topology data base as OSPF
- r **hierarchical** OSPF in large domains.

Externally derived routing data is advertised throughout the Autonomous System unaltered

Hierarchical OSPF

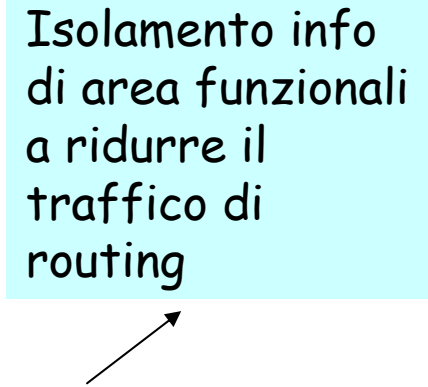


Splitting the AD into areas

- r OSPF allows collections of contiguous networks and hosts to be grouped together
 - m Such a group together with the routers with interfaces to any of the included networks is called an area
 - m Each area runs a separate copy of the basic link-state routing algorithm
 - has its own link state database
 - m The topology of an area is invisible from the outside
 - m Routers internal to a given area know nothing of the detailed topology external to the area

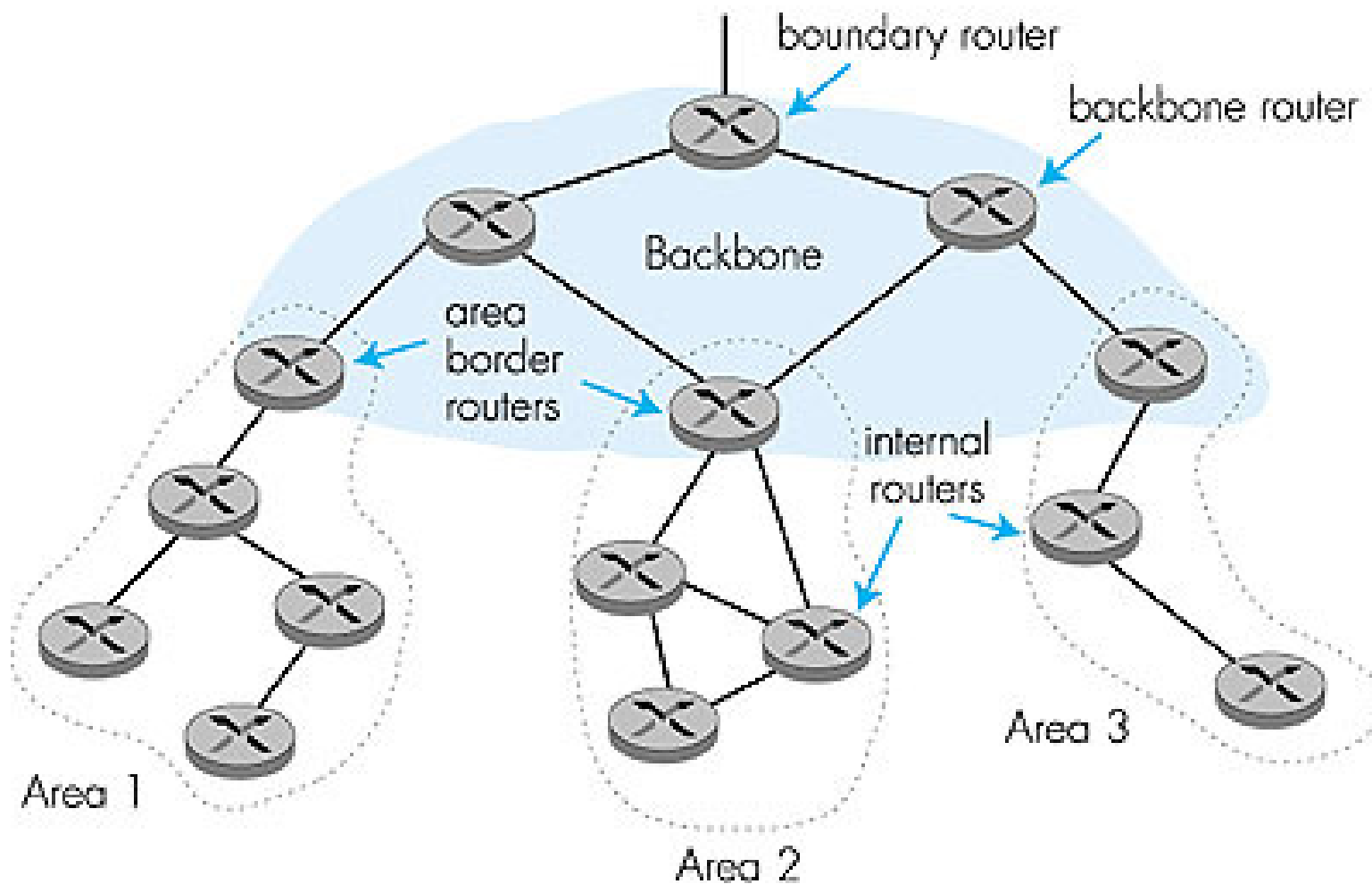
Hierarchical OSPF

Isolamento info
di area funzionali
a ridurre il
traffico di
routing



- r **two-level hierarchy:** local area, backbone.
 - m Link-state advertisements only in area
 - m each node has detailed area topology; only known direction (shortest path) to nets in other areas.
- r **area border routers:** "summarize" distances to nets in own area, advertise to other Area Border routers.

Hierarchical OSPF



Hierarchical OSPF

- r The OSPF backbone is the special OSPF Area 0
- r The OSPF backbone always contains backbone routers
- r Backbone routers: run OSPF routing limited to backbone.
- r The backbone is responsible for distributing routing information between non backbone areas
 - m Every area border router hears the area summaries from all other area border routers
 - m adding backbone distance+distance in summaries each router knows distance to different destinations
 - m These distances are then advertised internally to their areas
- r The backbone must be contiguous but not physically contiguous
 - m Backbone connectivity can be established/maintained through the configuration of virtual links (part of the backbone with actual way to route between end point of the virtual link based on intra_AS routing)
- r Boundary routers: connect to other AS's.
- r AS external LSAs are advertised in the AS WITH THE EXCEPTION OF stub areas
 - m Stub areas use a default routing

Types of networks

- r Transit networks are capable of carrying data traffic which is neither locally originated nor locally destined
- r A stub network only carries traffic it either generates or addressed to it

LSA (Link State Advertisement)

- r Periodic advertisement
- r Link state is also advertised when a router state changes
 - m Hello packets used to discover and maintain neighbor relationships
- r Disseminated via flooding
- r Flooding algorithm is reliable ensuring that all routers in the area have the same link state database

Chapter 4: Network Layer

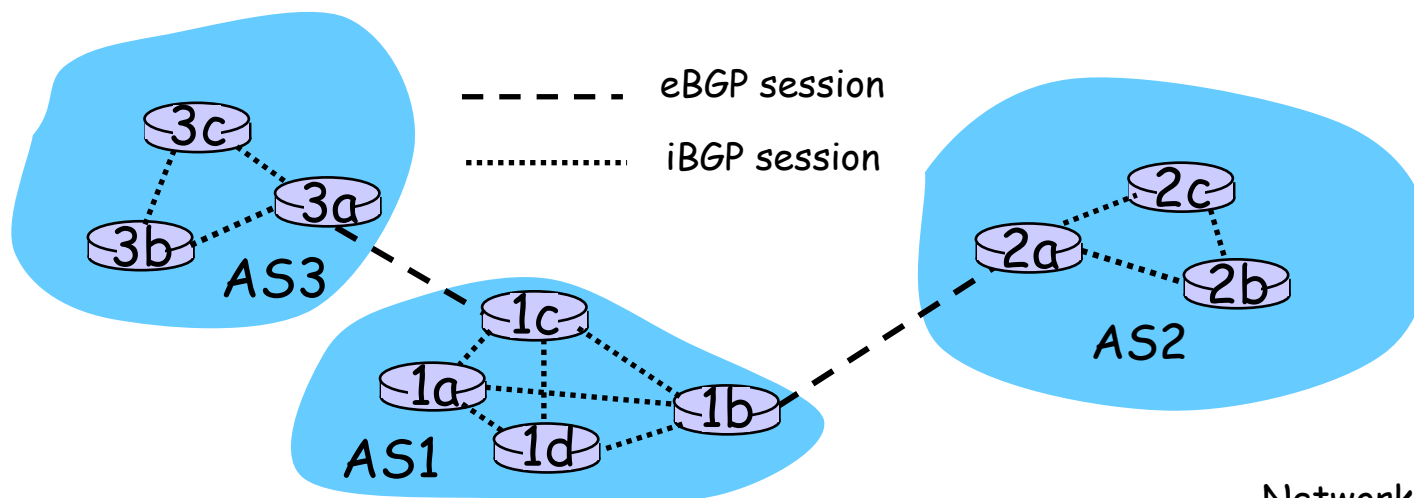
- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6
- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing

Internet inter-AS routing: BGP

- r **BGP (Border Gateway Protocol):** *the de facto standard*
- r BGP provides each AS a means to:
 1. Obtain subnet reachability information from neighboring ASs.
 2. Propagate reachability information to all AS-internal routers.
 3. Determine "good" routes to subnets based on reachability information and policy.
- r allows subnet to advertise its existence to rest of Internet: *"I am here"*

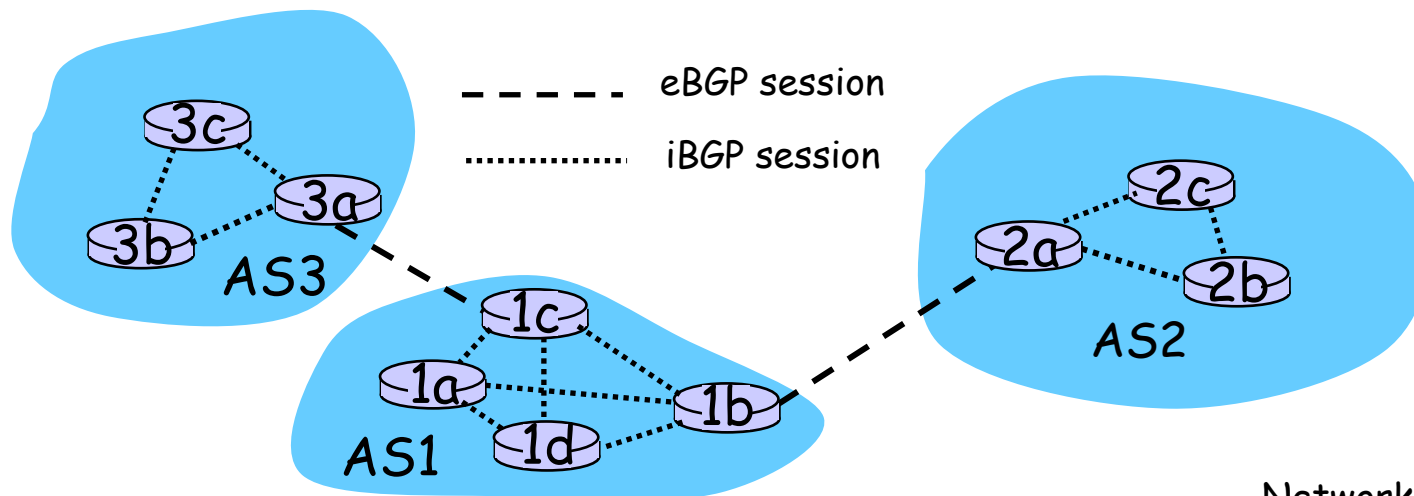
BGP basics

- r pairs of routers (BGP peers) exchange routing info over semi-permanent TCP connections: **BGP sessions**
 - m BGP sessions need not correspond to physical links.
- r when AS2 advertises a prefix to AS1:
 - m AS2 **promises** it will forward datagrams towards that prefix.
 - m AS2 can aggregate prefixes in its advertisement



Distributing reachability info

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



Path attributes & BGP routes

- r advertised prefix includes BGP attributes.
 - m prefix + attributes = "route"
- r two important attributes:
 - m **AS-PATH**: contains ASs through which prefix advertisement has passed: e.g, AS 67, AS 17
 - m **NEXT-HOP**: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- r when gateway router receives route advertisement, uses **import policy** to accept/decline.

BGP route selection

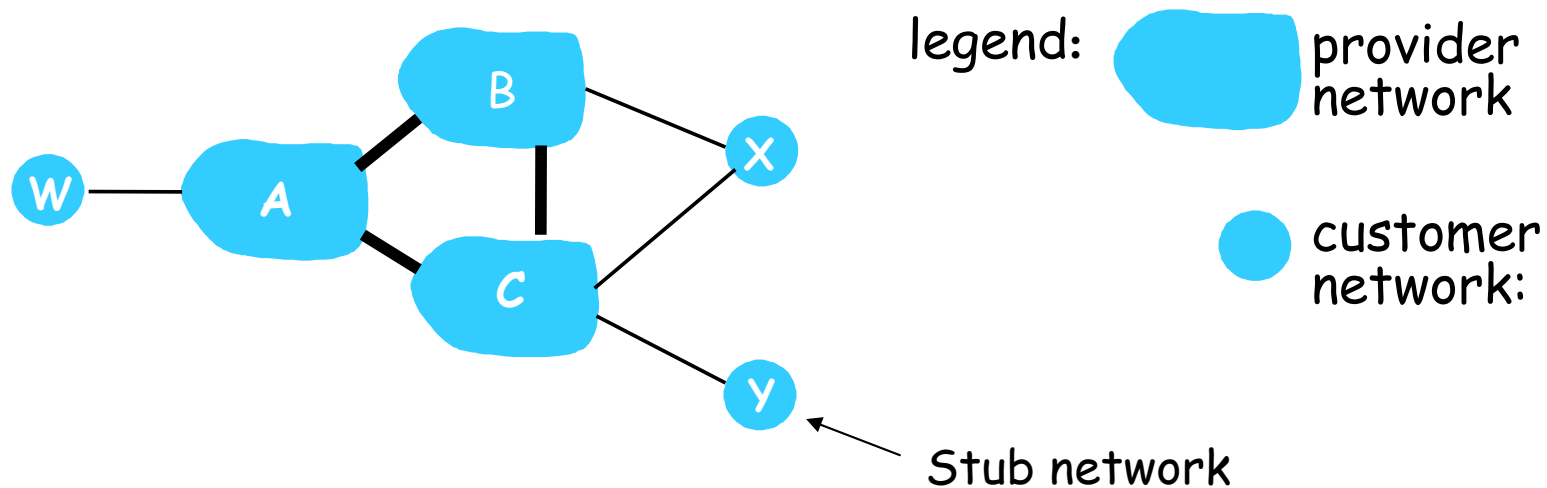
- r router may learn about more than 1 route to some prefix. Router must select route.
- r elimination rules (in priority order):
 1. local preference value attribute: policy decision
 2. (in case of same preference) shortest AS-PATH
 3. (in case of same preference and AS-PATH length) closest NEXT-HOP router: hot potato routing
 4. additional criteria to break the tie

BGP messages

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

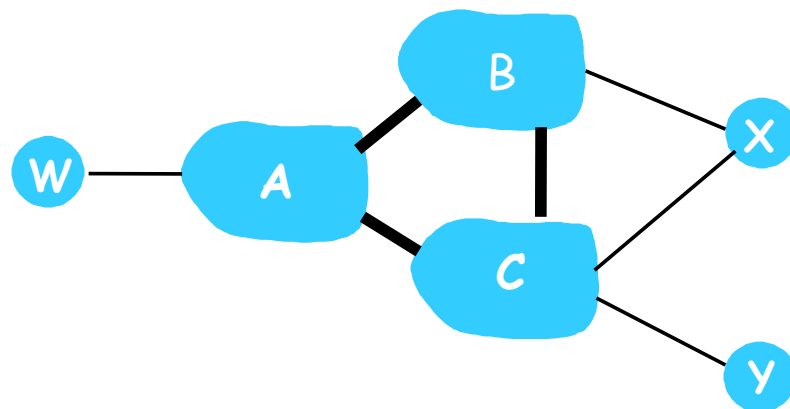
Sent periodically or in case of selected routes changes



BGP routing policy



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

BGP routing policy (2)



legend:  provider network
 customer network:

- r A advertises path AW to B
- r B advertises path BAW to X (who is its client)
- r Should B advertise path BAW to C?
 - m No way! B gets no "revenue" for routing $CBAW$ since neither W nor C are B's customers
 - m B wants to force C to route to w via A
 - m B wants to route *only* to/from its customers!
 - m Peering agreements amongs pairs of ISP possible to solve this problem

Decision process

- r The decision process selects routes for subsequent advertisement applying the policies in the local Policy Information Base (PIB) to the routes stored in its Adj-RIBs_In (Incoming Routing Information Base)
- r A function takes as argument the attributes of a give route and returns a) either a non negative integer identifying the degree of preference for the route or b) a value indicating the route is inelegible

Why different Intra- and Inter-AS routing ?

Policy:

- r Inter-AS: admin wants control over how its traffic routed, who routes through its net.
- r Intra-AS: single admin, so no policy decisions needed

Scale:

- r hierarchical routing saves table size, reduced update traffic

Performance:

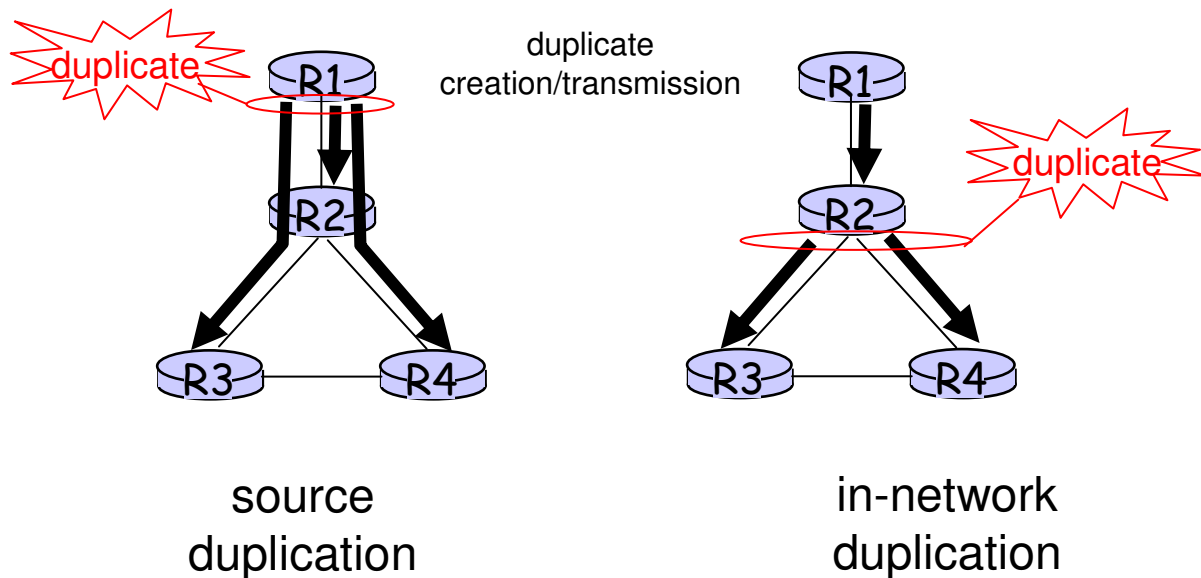
- r Intra-AS: can focus on performance
- r Inter-AS: policy may dominate over performance

Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6
- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r **4.7 Broadcast and multicast routing**

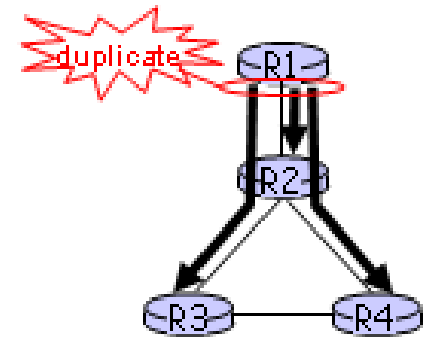
Broadcast Routing

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



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

Unicast ad N vie

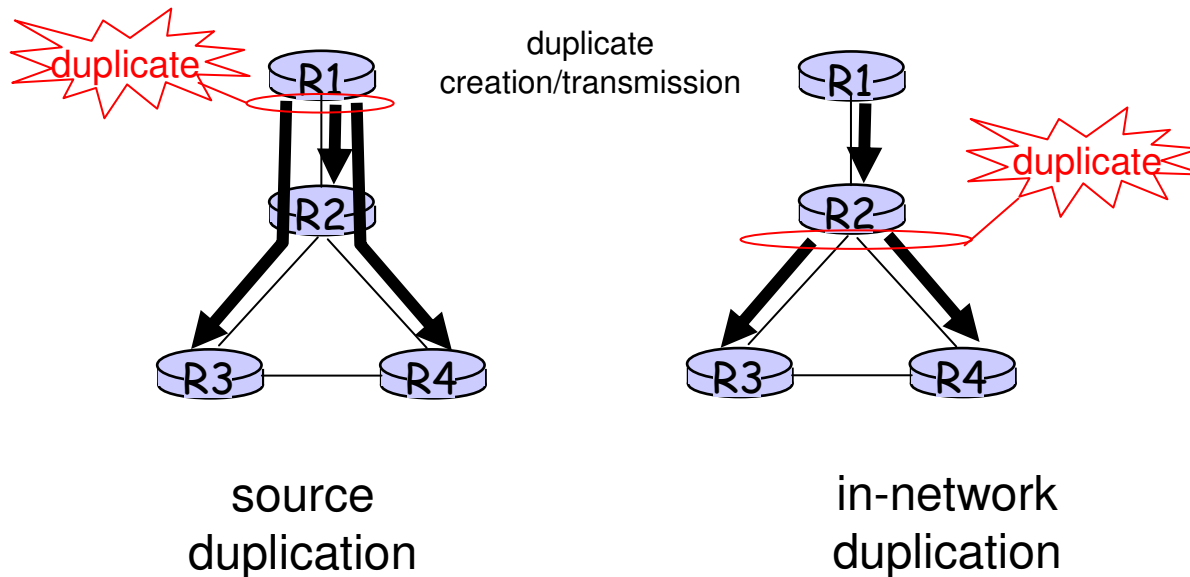


source
duplication

- r Inefficiente
 - m Un singolo collegamento attraversato da N copie del messaggio se il nodo origine è connesso al resto della rete tramite un unico collegamento
- r Indirizzi di tutte le destinazioni devono essere noti al mittente
 - m altri meccanismi protocollari sono richiesti
- r Broadcast può essere usato per inoltrare informazioni di topologia in una situazione in cui le rotte non sono ancora note
 - m es. OSPF

Broadcast Routing

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

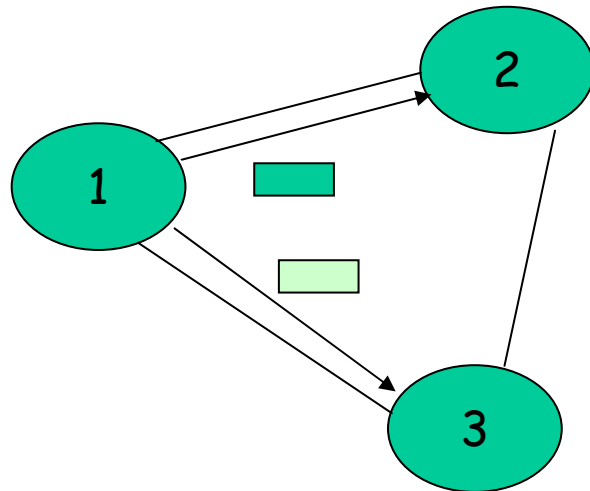


- r source duplication: how does source determine recipient addresses?

In-network duplication

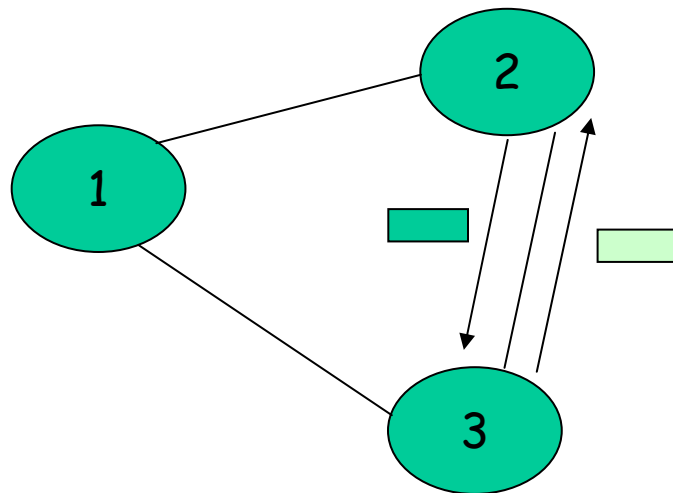
r flooding: when node receives brdcst pckt, sends copy to all neighbors EXCEPT the one from which the pckt was received

m Problems: cycles & broadcast storm



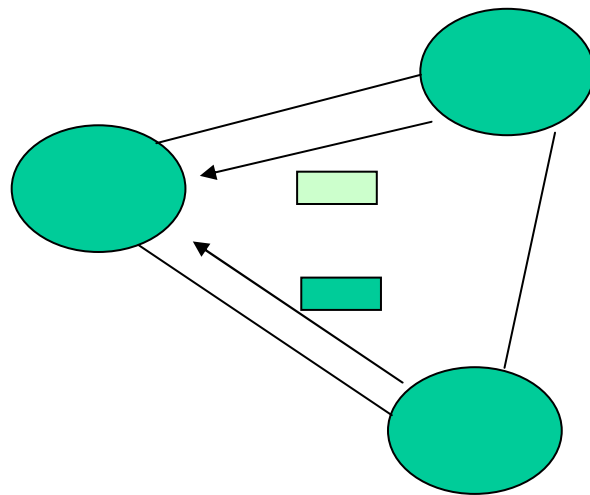
In-network duplication

- r flooding: when node receives brdcst pckt, sends copy to all neighbors
 - m Problems: cycles & broadcast storm



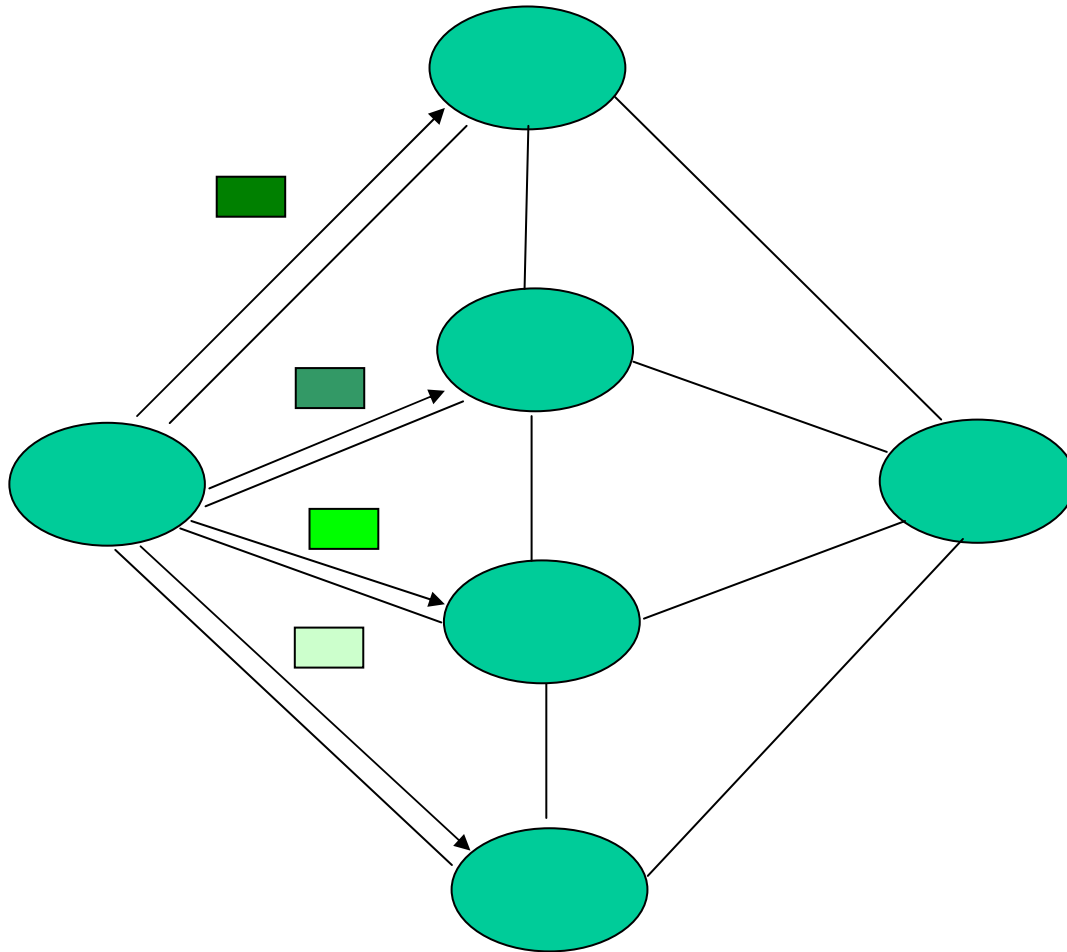
In-network duplication

- r flooding: when node receives brdcst pkt, sends copy to all neighbors
 - m Problems: cycles & broadcast storm

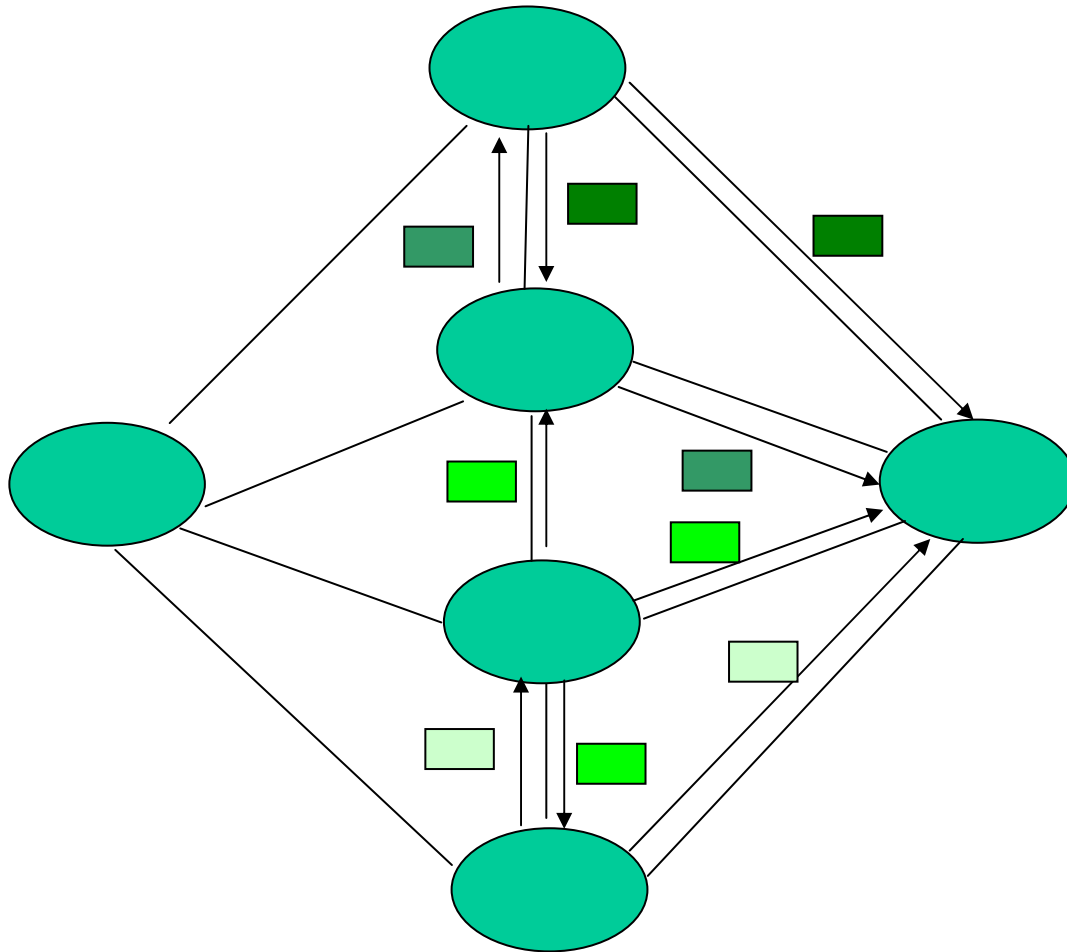


E ricominciamo come nella prima situazione
Bisogna saper distinguere tra quando
mandiamo un nuovo messaggio e quando
stiamo ritrasmettendo qualcosa che
abbiamo già visto
→ Sequence numbers!

Broadcast storm

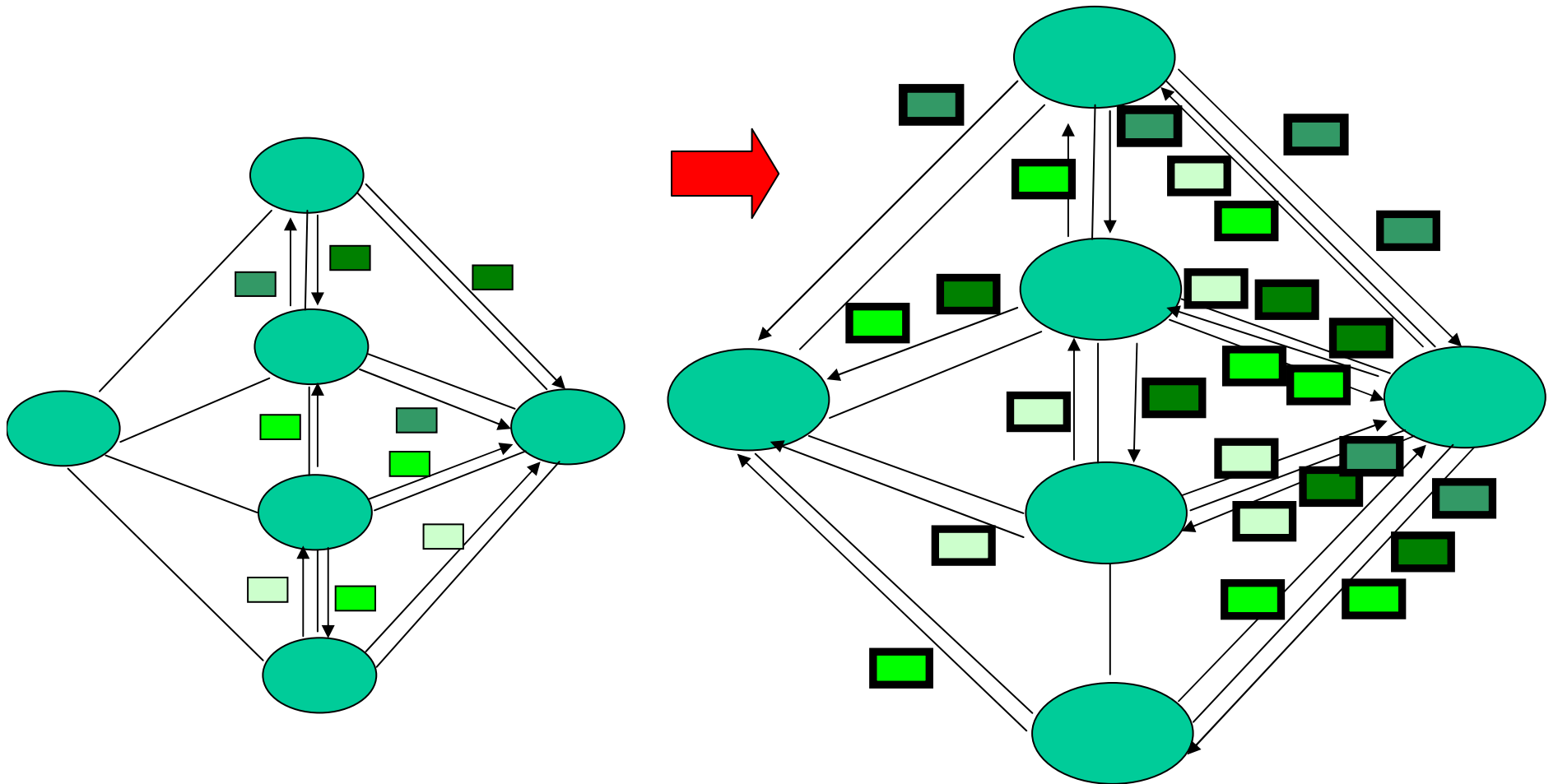


Broadcast storm



Broadcast storm

Il numero di pacchetti in rete cresce significativamente!!



Controlled flooding

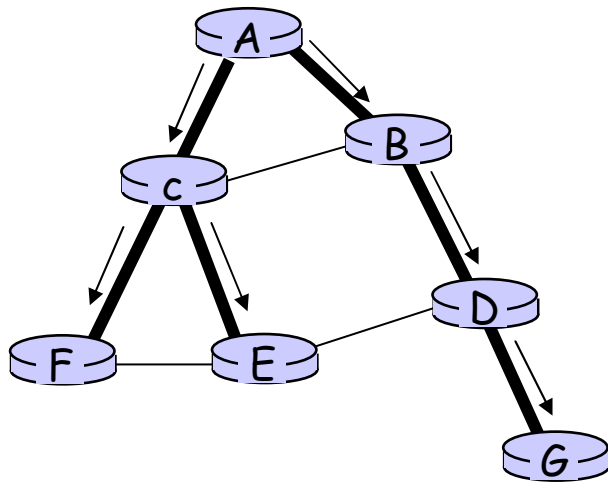
- r Il nodo origine pone il proprio indirizzo ed il numero di sequenza nei pacchetti che invia in broadcast
- r Ciascun nodo mantiene una lista di ID origine, SEQN per i broadcast ricevuti, trasmesso o inoltrato
- r Se riceve un pacchetto broadcast per prima cosa verifica se $\langle \text{ID}, \text{SEQN} \rangle$ compare nella lista dei pacchetti già gestiti
 - m Se si scarta
 - m Altrimenti riinvia su tutte le interfacce tranne quella da cui ha ricevuto

Controlled flooding, altre opzioni

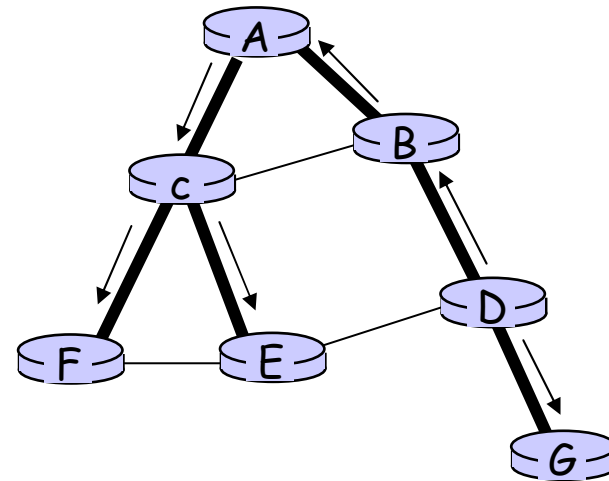
- r Reverse path forwarding (RPF): only forward pckt (on all links but the one from which the packet was received) if it arrived on shortest path between node and source

Spanning Tree

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



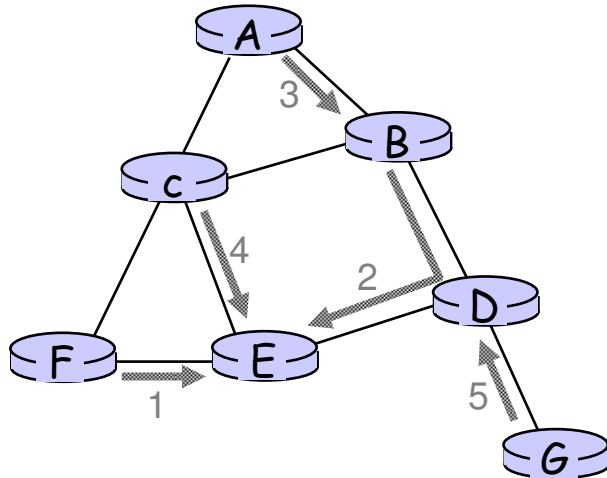
(a) Broadcast initiated at A



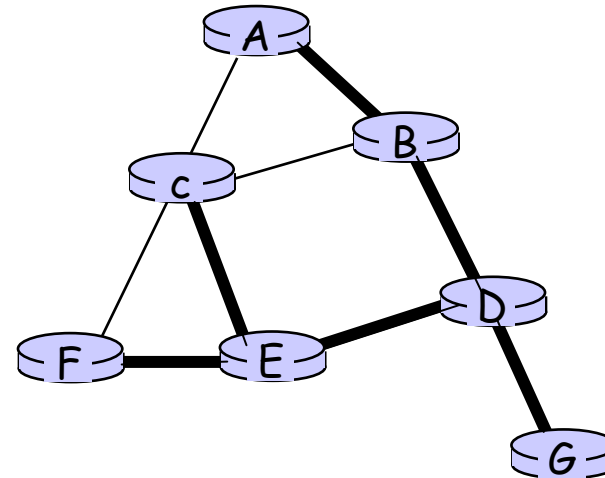
(b) Broadcast initiated at D

Spanning Tree: Creation

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



(a) Stepwise construction of spanning tree



(b) Constructed spanning tree

Multicasting

- r Molte applicazioni richiedono il trasferimento di pacchetti da uno o più mittenti ad un gruppo di destinatari
 - m trasferimento di un aggiornamento SW su un gruppo di macchine
 - m streaming (audio/video) ad un gruppo di utenti o studenti
 - m applicazioni con dati condivisi (lavagna elettronica condivisa da più utenti)
 - m aggiornamento di dati (adnamento di borsa)
 - m giochi multi-player interattivi
 - m ...

Indirizzamento Multicast

- r L'identificatore che rappresenta un gruppo multicast è un indirizzo IP multicast di classe D
- r Come ci si affilia ad un indirizzo multicast?
Come vengono gestiti i cambiamenti dinamici (join/remove) nel gruppo?
 - m Gestione dinamica del gruppo OLTRE a
 - m Algoritmi per la consegna delle informazioni ad un gruppo multicast

IGMP Internet Group Management Protocol

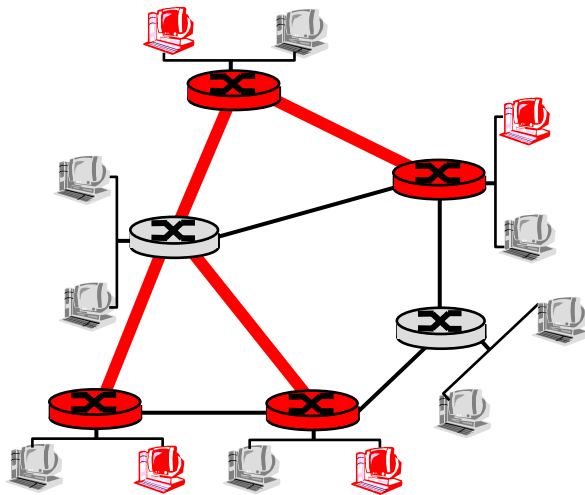
- r Messaggi incapsulati in datagrammi IP, con IP protocol number 2
- r Mandati con TTL a 1
- r Messaggi IGMP
 - m Type (8bit) Query (richiesta dal router)/ Membership Report (risposta dagli host)/ Leave group (ma anche possibile timeout + mancata risposta alla richiesta del router → soft state)
- r Max Response Time (per rispondere ad una query)
- r Checksum
- r Group Address (0 se si manda una general query, indirizzo IP del gruppo nel caso di una group specific query con cui si richiede chi sia affiliato a quel gruppo)

IGMP Internet Group Management Protocol

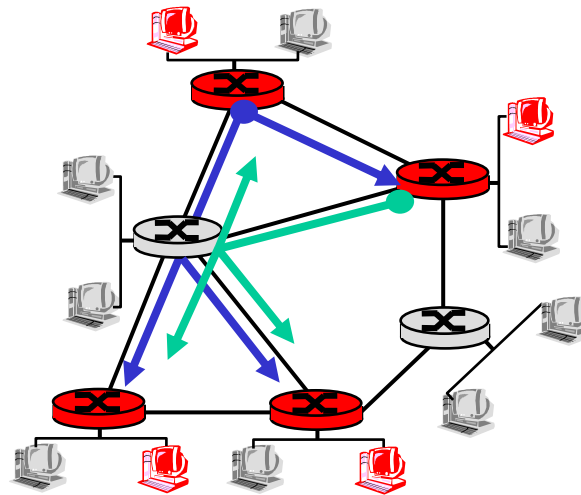
- r IGMP consente ad un router di imparare quali gruppi multicast hanno affiliati sulle sottoreti connesse a ciascuna delle loro interfacce
- r Un router multicast tiene una lista per ciascuna sottorete dei multicast group (multicast group membership → almeno un elemento del gruppo fa parte della sottorete) con un timer per membership
 - m la membership deve essere aggiornata da report inviati prima della scadenza del timer
 - m può essere anche aggiornata tramite messaggi di leave espliciti

Multicast Routing: Problem Statement

- r **Goal:** find a tree (or trees) connecting routers having local mcast group members
- m **tree:** not all paths between routers used
- m **source-based:** different tree from each sender to rcvrs
- m **shared-tree:** same tree used by all group members



Shared tree



Source-based trees

Approaches for building mcast trees

Approaches:

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

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

Instradamento multicast con albero condiviso dal gruppo

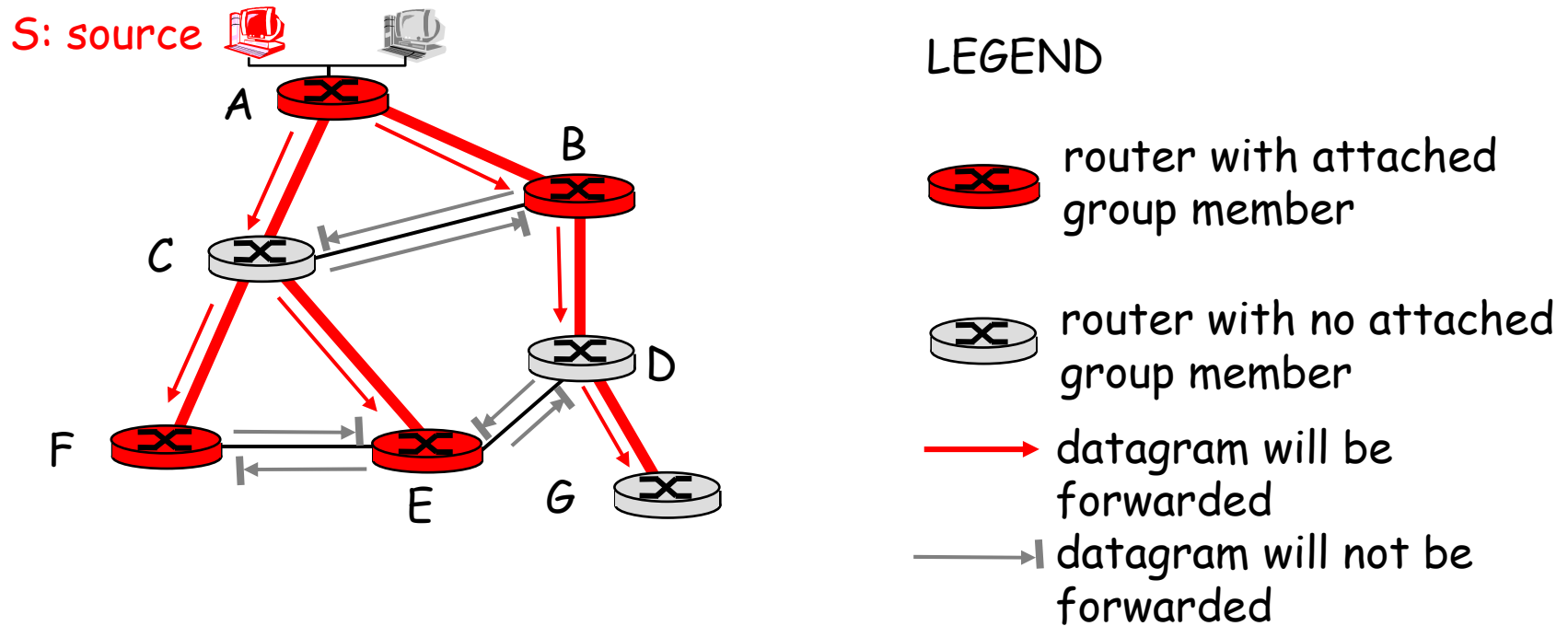
- r Stesso approccio visto per calcolare un albero di ricoprimento per il flooding
- r **I router multicast connessi a sottoreti con elementi del gruppo si affiliano tramite un unicast message mandato all'iniziatore dell'albero**
- r **Tutti i nodi attraversati prima di unirsi ad un elemento già parte dell'albero di multicast fanno parte dell'albero**

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 for source based multicasting: example

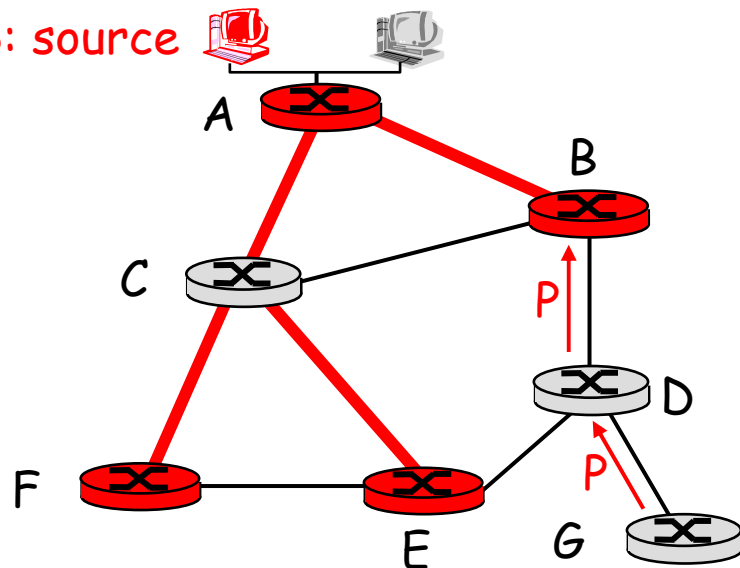


- result is a source-specific *reverse* SPT
 - may be a bad choice with asymmetric links
 - unneeded forwarding e.g., from D to G





Reverse Path Forwarding: pruning

- r forwarding tree contains subtrees with no mcast group members
- m no need to forward datagrams down subtree
- m "prune" msgs sent upstream by router with no downstream group members

S: source



LEGEND

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

Shared-Tree: Steiner Tree

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

Center-based trees

- r single delivery tree shared by all
- r one router identified as "*center*" of tree
- r to join:
 - m edge router sends unicast *join-msg* addressed to center router
 - m *join-msg* "processed" by intermediate routers and forwarded towards center
 - m *join-msg* either hits existing tree branch for this center, or arrives at center
 - m path taken by *join-msg* becomes new branch of tree for this router

Chapter 4: summary

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
 - m Datagram format
 - m IPv4 addressing
 - m ICMP
 - m IPv6
- r 4.5 Routing algorithms
 - m Link state
 - m Distance Vector
 - m Hierarchical routing
- r 4.6 Routing in the Internet
 - m RIP
 - m OSPF
 - m BGP
- r 4.7 Broadcast and multicast routing