

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

Chapter goals:

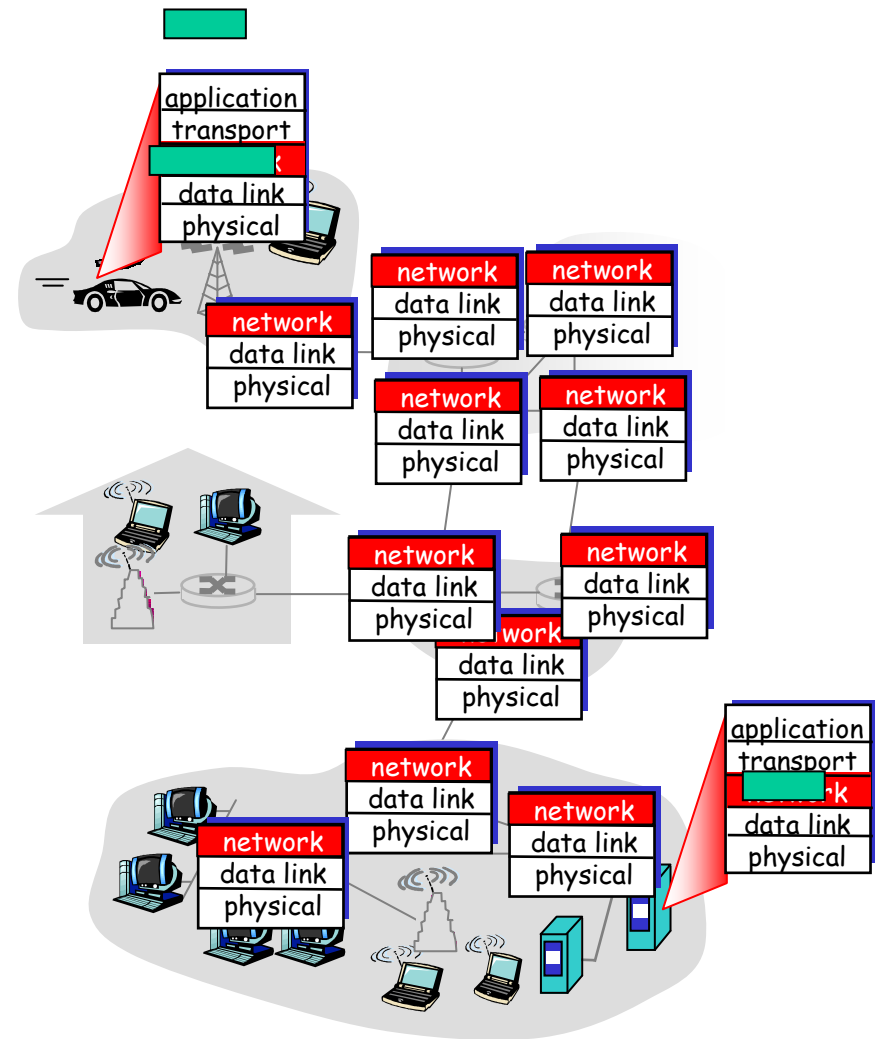
- r understand principles behind network layer services:
 - m network layer service models
 - m forwarding versus routing
 - m how a router works
 - m routing (path selection)
 - m dealing with scale
 - m advanced topics: IPv6, mobility
- r instantiation, implementation in the Internet

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

Network layer

- r transport segment from sending to receiving host
- r on sending side encapsulates segments into datagrams
- r on rcving side, delivers segments to transport layer
- r network layer protocols in *every* host, router
- r router examines header fields in all IP datagrams passing through it



Two Key Network-Layer Functions

r *forwarding*: move packets from router's input to appropriate router output

r *routing*: determine route taken by packets from source to dest.

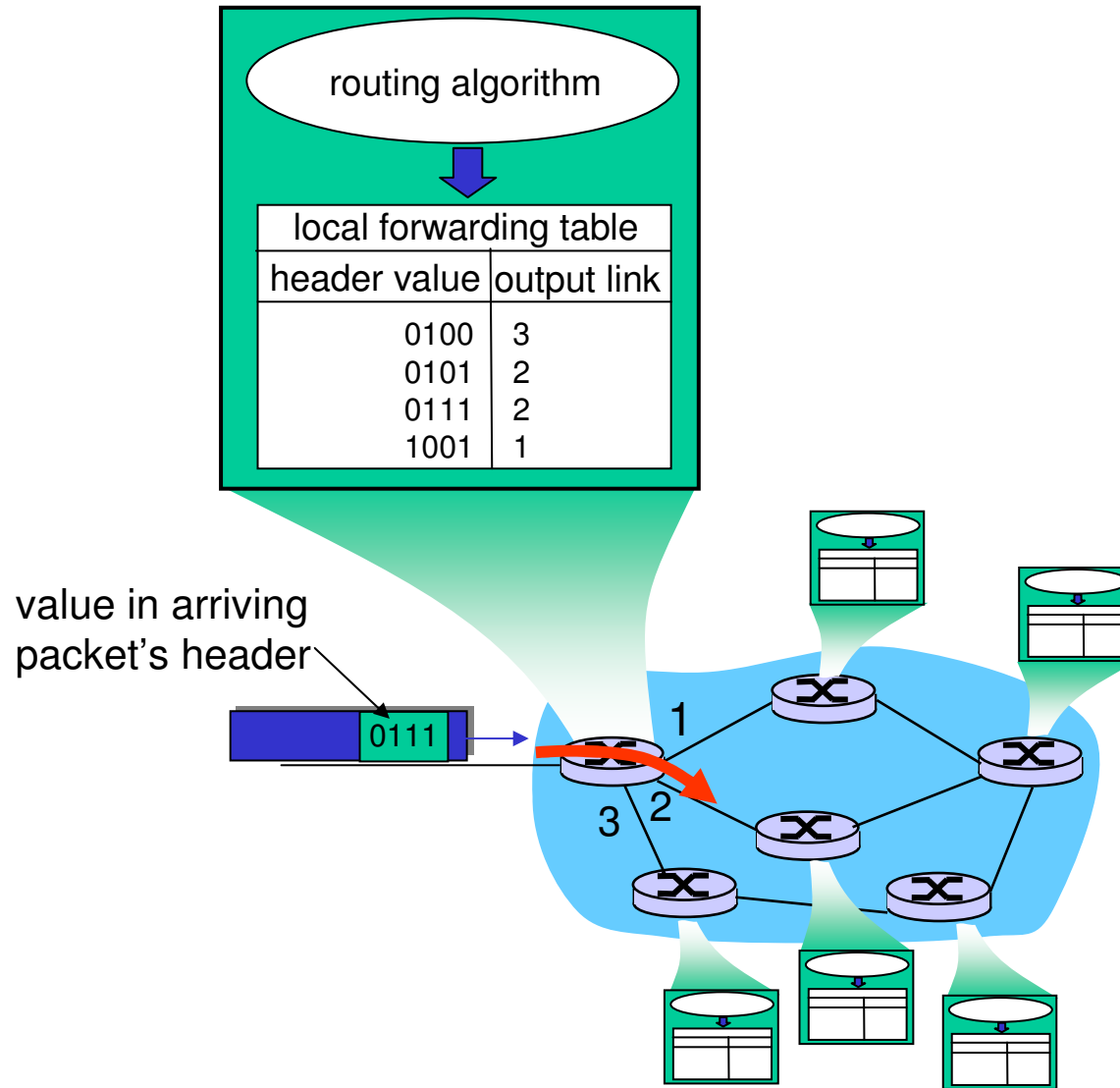
m *routing algorithms*

analogy:

r *routing*: process of planning trip from source to dest

r *forwarding*: process of getting through single interchange

Interplay between routing and forwarding



Network service model

Q: What *service model* for “channel” transporting datagrams from sender to receiver?

Example services for individual datagrams:

- r guaranteed delivery
- r guaranteed delivery with less than 40 msec delay

Example services for a flow of datagrams:

- r in-order datagram delivery
- r guaranteed minimum bandwidth to flow
- r restrictions on changes in inter-packet spacing

Network layer service models:

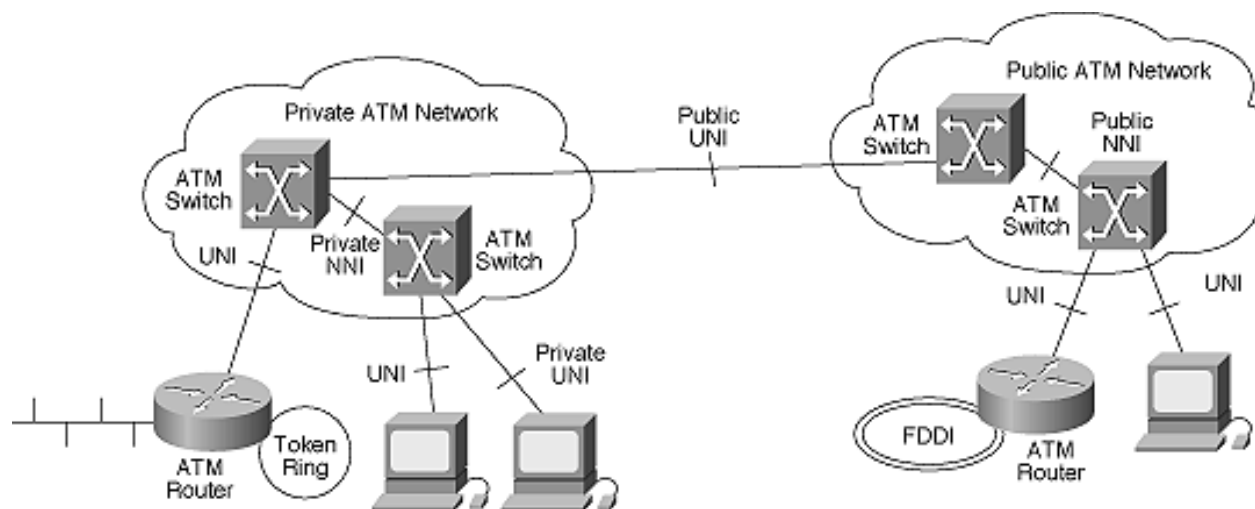
Network Architecture	Service Model	Guarantees ?			Congestion feedback	
		Bandwidth	Loss	Order Timing		
Internet	best effort	none	no	no	no (inferred via loss)	
ATM	CBR	constant rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

Connection setup

- r 3rd important function in *some* network architectures:
 - m ATM, frame relay, X.25
- r before datagrams flow, two end hosts *and* intervening routers establish virtual connection
 - m routers get involved
- r network vs transport layer connection service:
 - m **network**: between two hosts (may also involve intervening routers in case of VCs)
 - m **transport**: between two processes

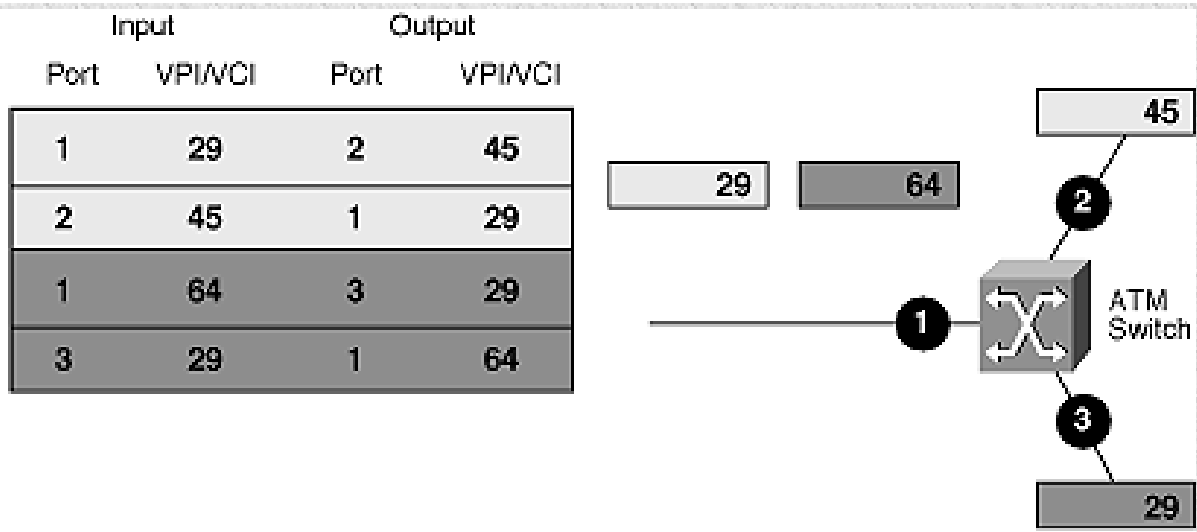
ATM networks

- r **UNI User-Network Interface:** connette un host con lo switch a cui é collegato
- r **NNI Network-Network Interface:** connette due switch.

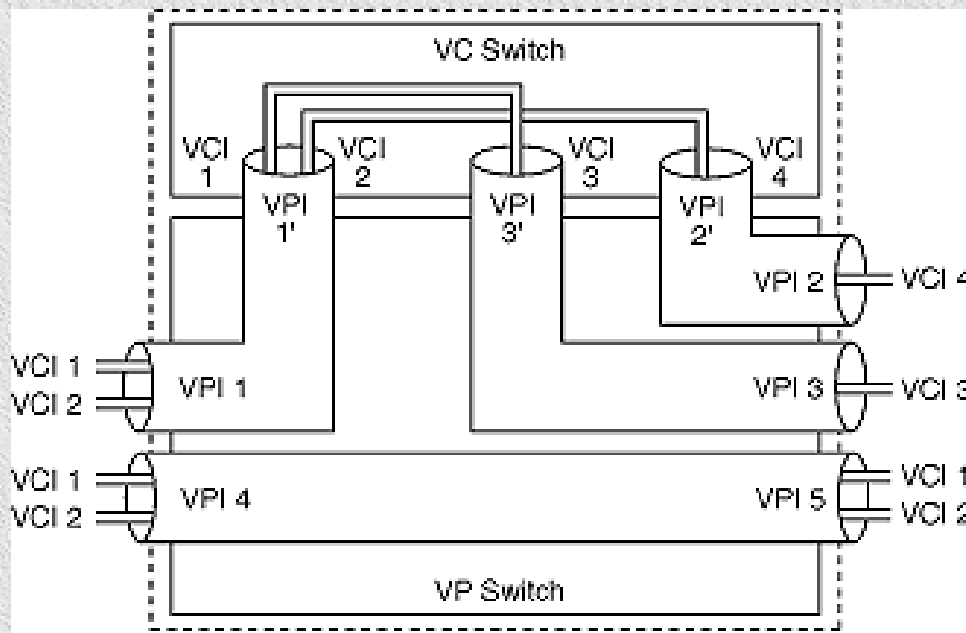


- r Virtual Path (insieme di virtual channel): identificate dal valore VPI
- r Virtual Channel : identificate dalla coppia VPI/VCI

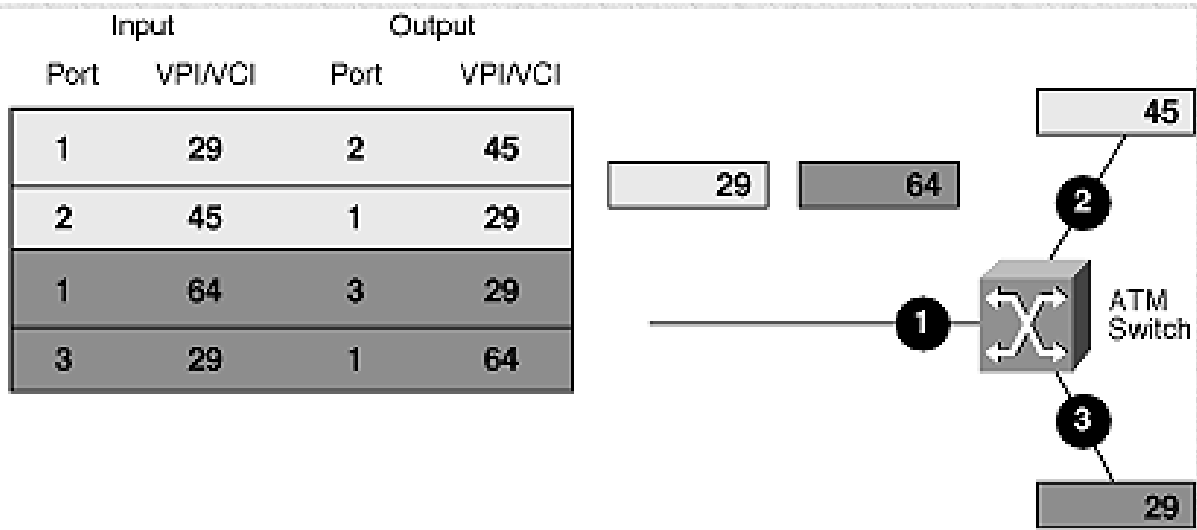
ATM switching



La coppia VPI/VCI ha solo significato *locale* nel senso che vale esclusivamente all'interno di un link, ogni volta che la connessione attraversa uno switch la coppia VPI/VCI in genere viene cambiata ed in questo consiste il meccanismo di switching ovvero nel determinare in base ad una tabella di instradamenti come si mappano le connessioni in base alle porte di provenienza ed al VPI/VCI delle celle.



ATM switching



Rete a commutazione di Pacchetto, a circuito Virtuale

Orientata alla connessione

→ I pacchetti arrivano in ordine

Tipi di trasmissione

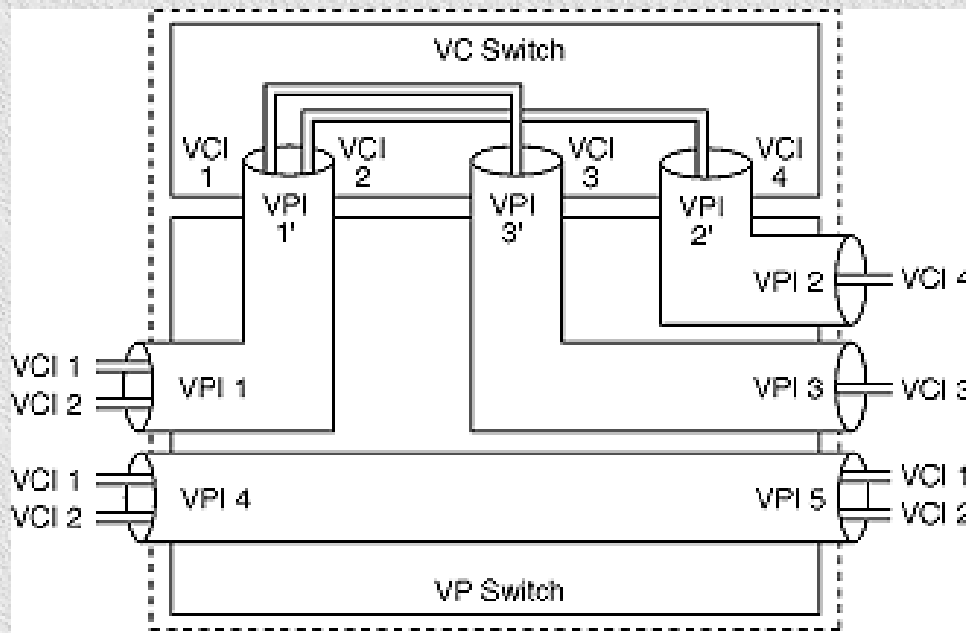
→ CBR

→ VBR

→ ABR (Available Bit Rate)

→ UBR (Unspecified Bit Rate)

Call Admission



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

Network layer connection and connection-less service

- r datagram network provides network-layer connectionless service
- r VC network provides network-layer connection service
- r analogous to the transport-layer services, but:
 - m **service:** host-to-host
 - m **no choice:** network provides one or the other
 - m **implementation:** in network core

Virtual circuits

“source-to-dest path behaves much like telephone circuit”

- m performance-wise

- m network actions along source-to-dest path

- r call setup, teardown for each call *before* data can flow
- r each packet carries VC identifier (not destination host address)
- r *every* router on source-dest path maintains “state” for each passing connection
- r link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)

VC implementation

a VC consists of:

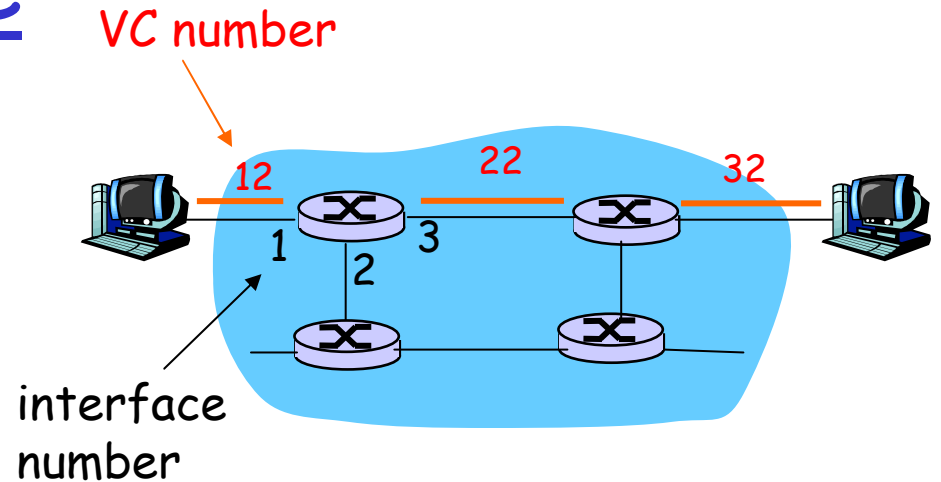
1. path from source to destination
2. VC numbers, one number for each link along path
3. entries in forwarding tables in routers along path

r packet belonging to VC carries VC number (rather than dest address)

r VC number can be changed on each link.

m New VC number comes from forwarding table

Forwarding table



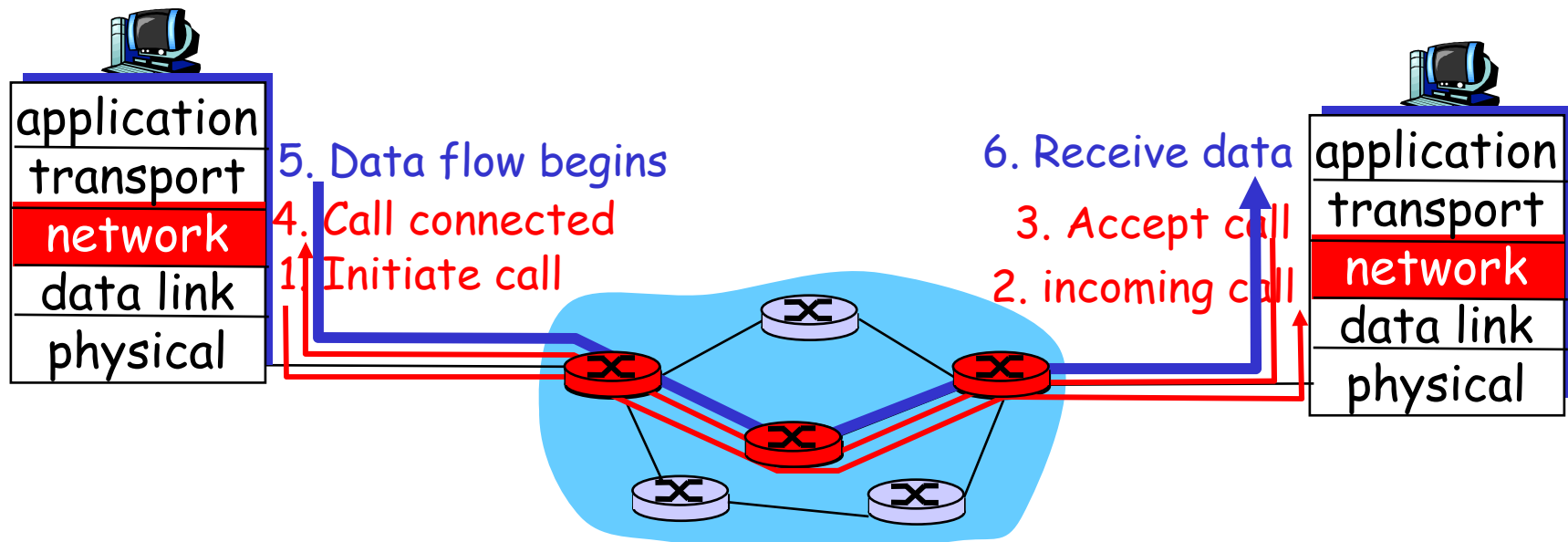
Forwarding table in northwest router:

Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...

Routers maintain connection state information!

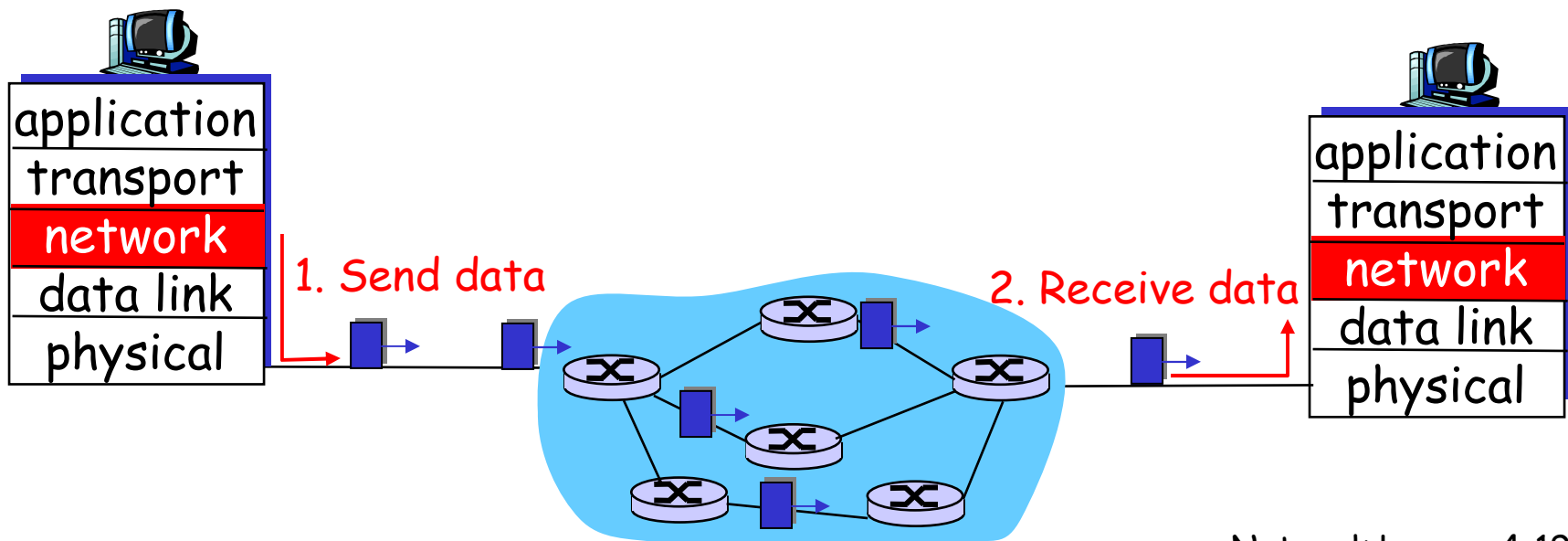
Virtual circuits: signaling protocols

- r used to setup, maintain teardown VC
- r used in ATM, frame-relay, X.25
- r not used in today's Internet



Datagram networks

- r no call setup at network layer
- r routers: no state about end-to-end connections
 - m no network-level concept of "connection"
- r packets forwarded using destination host address
 - m packets between same source-dest pair may take different paths



Forwarding table

4 billion
possible entries

<u>Destination Address Range</u>	<u>Link Interface</u>
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2
otherwise	3

Longest prefix matching

<u>Prefix Match</u>	<u>Link Interface</u>
11001000 00010111 00010	0
11001000 00010111 00011000	1
11001000 00010111 00011	2
otherwise	3

Examples

DA: 11001000 00010111 00010110 10100001

Which interface?

DA: 11001000 00010111 00011000 10101010

Which interface?

Datagram or VC network: why?

Internet (datagram)

- r data exchange among computers
 - m "elastic" service, no strict timing req.
- r "smart" end systems (computers)
 - m can adapt, perform control, error recovery
 - m simple inside network, complexity at "edge"
- r many link types
 - m different characteristics
 - m uniform service difficult

ATM (VC)

- r evolved from telephony
- r human conversation:
 - m strict timing, reliability requirements
 - m need for guaranteed service
- r "dumb" end systems
 - m telephones
 - m complexity inside network

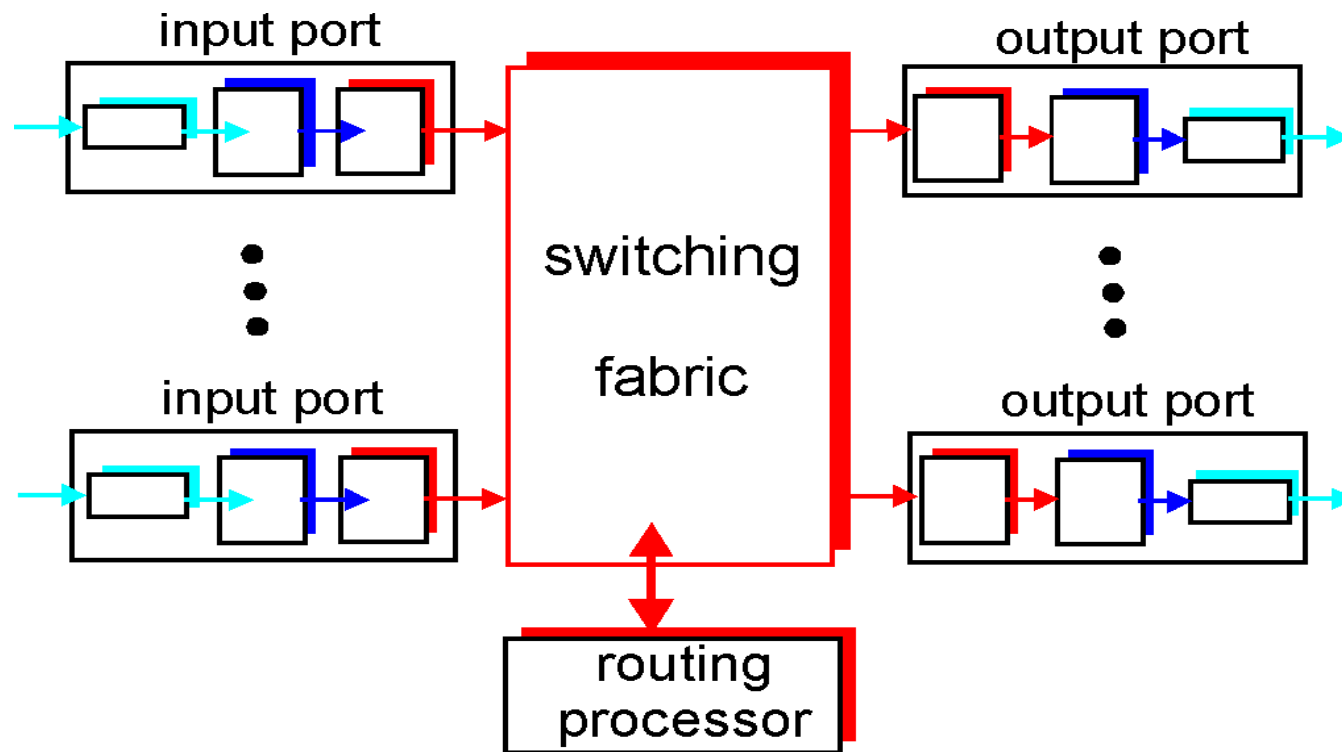
Chapter 4: Network Layer

- 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

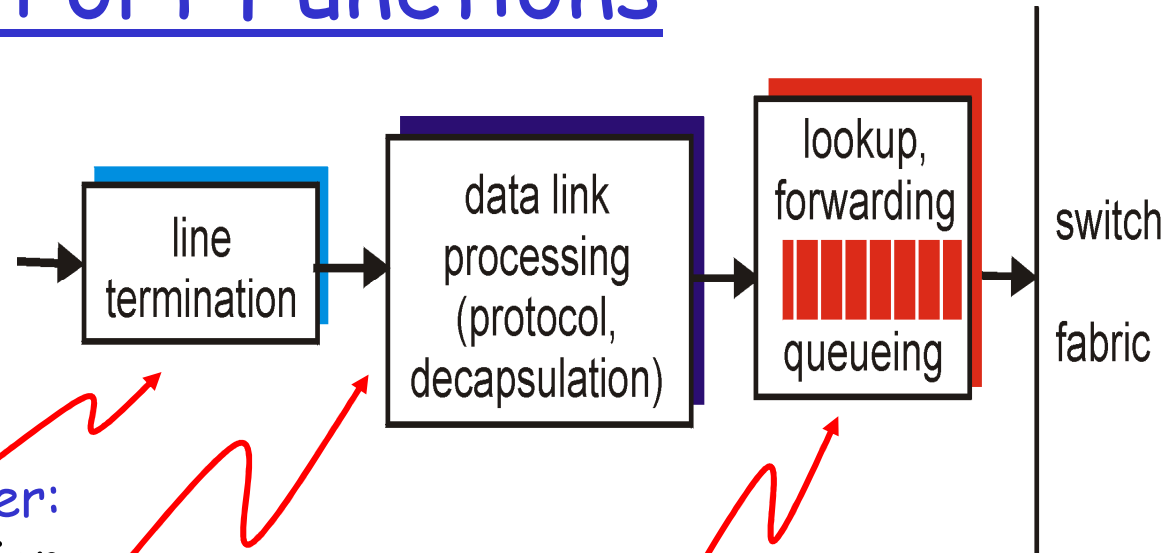
Router Architecture Overview

Two key router functions:

- r run routing algorithms/protocol (RIP, OSPF, BGP)
- r *forwarding* datagrams from incoming to outgoing link



Input Port Functions



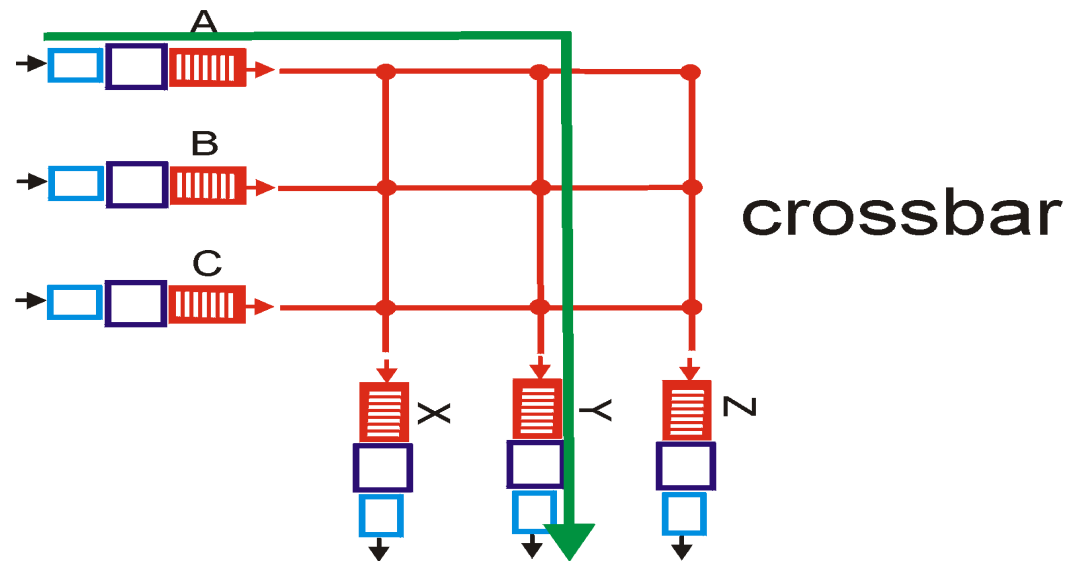
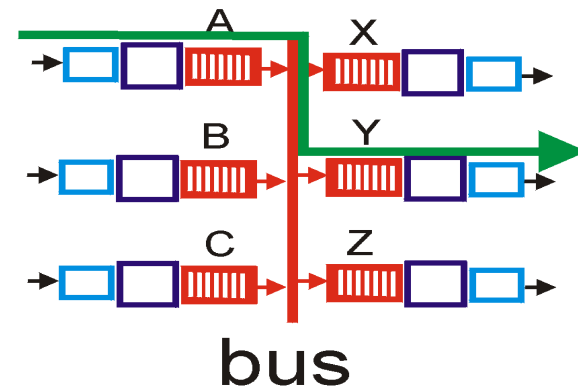
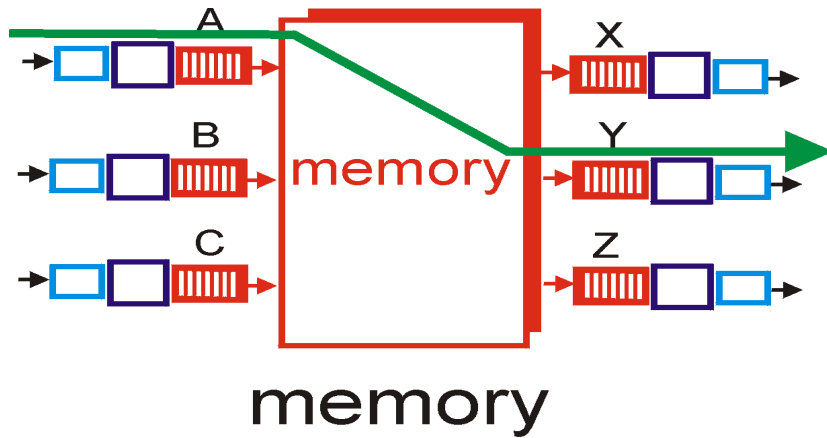
Physical layer:
bit-level reception

Data link layer:
e.g., Ethernet
see chapter 5

Decentralized switching:

- r given datagram dest., lookup output port using forwarding table in input port memory
- r goal: complete input port processing at 'line speed'
- r queuing: if datagrams arrive faster than forwarding rate into switch fabric

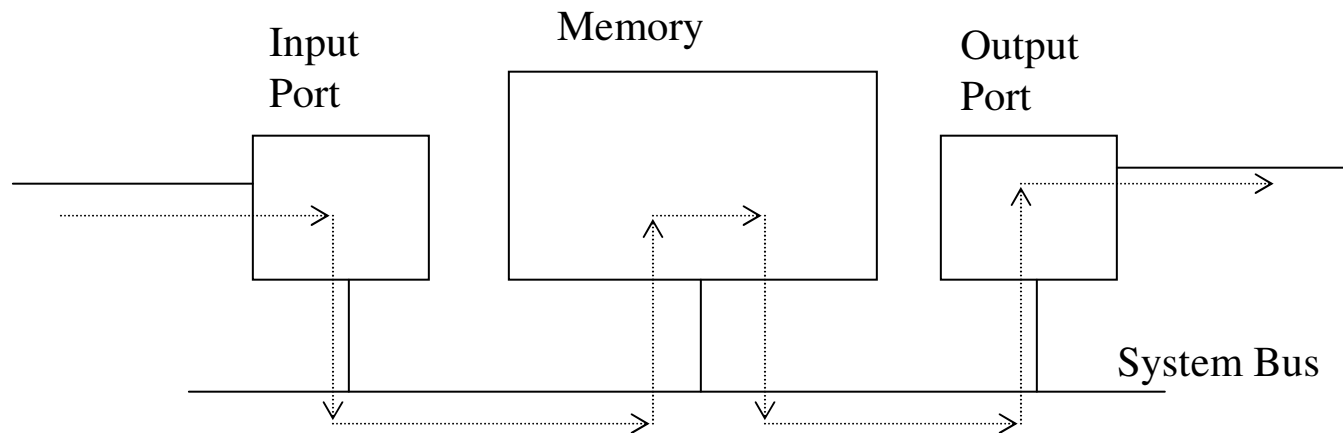
Three types of switching fabrics



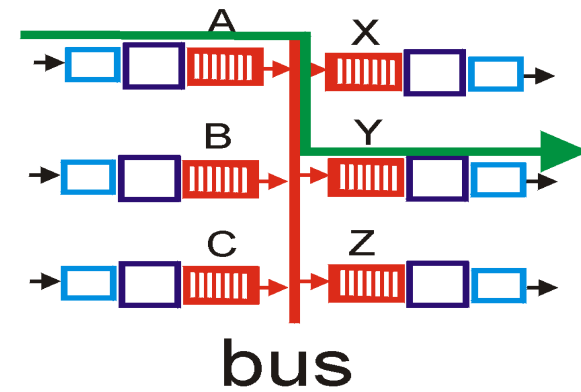
Switching Via Memory

First generation routers:

- r traditional computers with switching under direct control of CPU
- r packet copied to system's memory
- r speed limited by memory bandwidth (2 bus crossings per datagram)



Switching Via a Bus

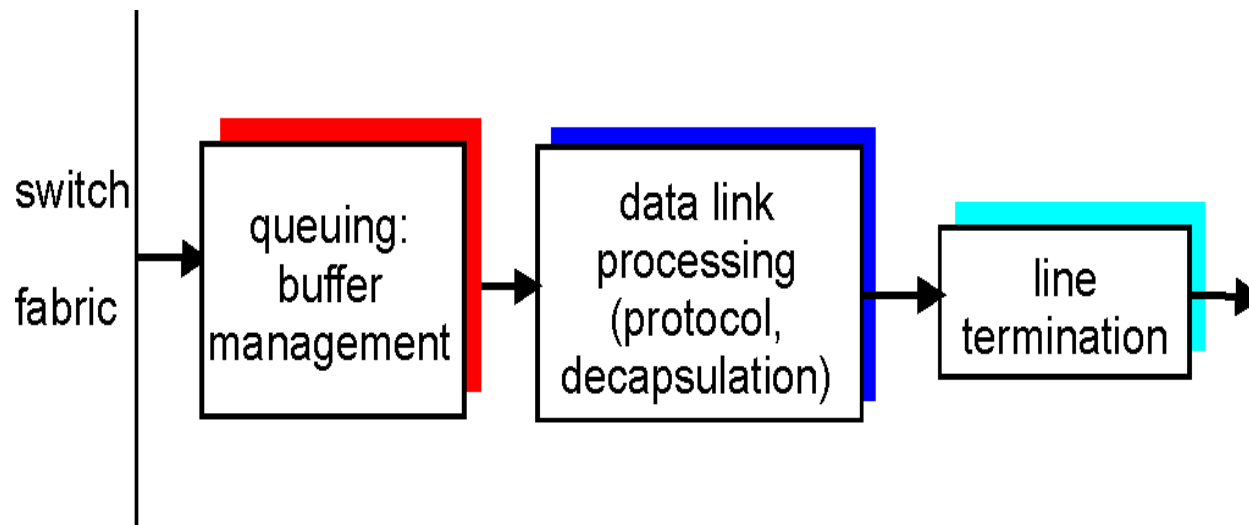


- r datagram from input port memory to output port memory via a shared bus
- r **bus contention:** switching speed limited by bus bandwidth
- r 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers

Switching Via An Interconnection Network

- r overcome bus bandwidth limitations
- r Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor
- r advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- r Cisco 12000: switches 60 Gbps through the interconnection network

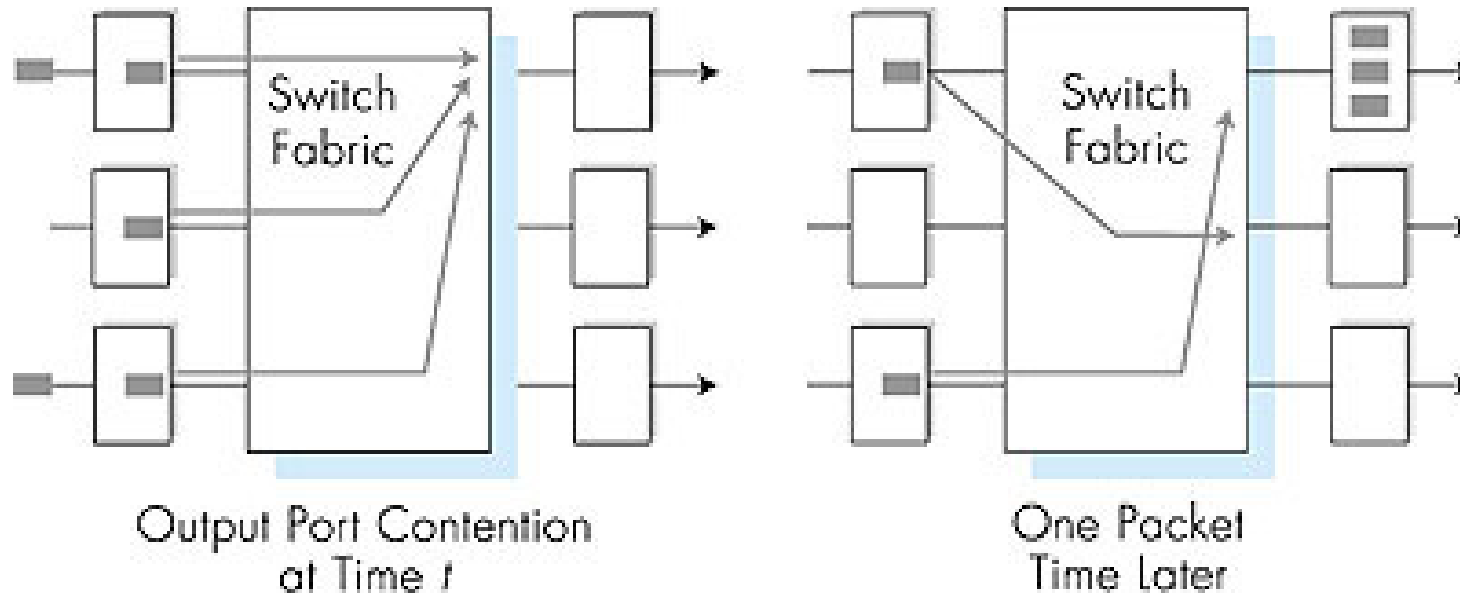
Output Ports



- r *Buffering* required when datagrams arrive from fabric faster than the transmission rate
 - m What if the queue builds up?
 - Drop Tail
 - Random Early Discard
- r *Scheduling discipline* chooses among queued datagrams for transmission
 - m First Come First Served
 - m Weighted Fair Queueing

$$\frac{Rw_i}{(w_1 + w_2 + \dots + w_N)}$$

Output port queueing



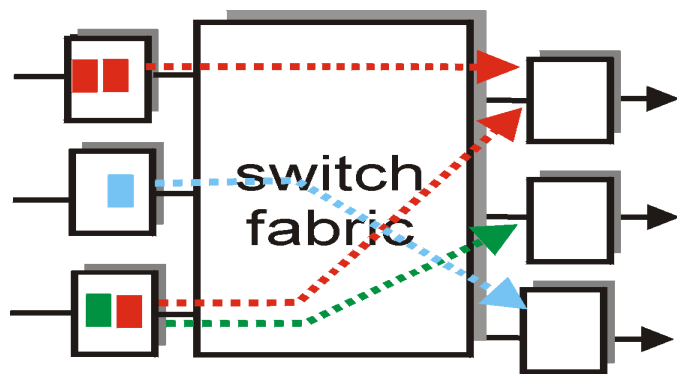
- r buffering when arrival rate via switch exceeds output line speed
- r *queueing (delay) and loss due to output port buffer overflow!*

How much buffering?

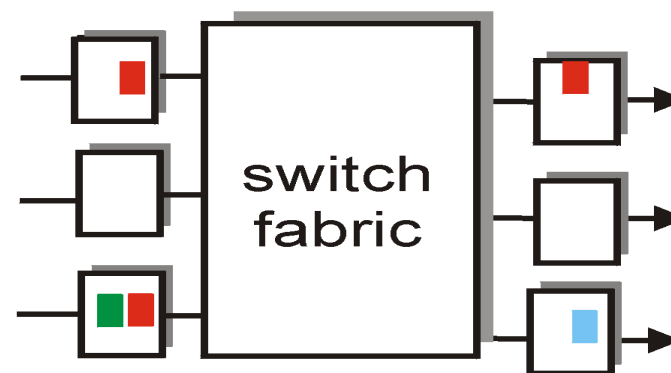
- r RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity C
 - m e.g., $C = 10$ Gps link: 2.5 Gbit buffer
- r Recent recommendation: with N flows, buffering equal to $\frac{RTT \cdot C}{\sqrt{N}}$

Input Port Queuing

- r Fabric slower than input ports combined -> queueing may occur at input queues
- r **Head-of-the-Line (HOL) blocking:** queued datagram at front of queue prevents others in queue from moving forward
- r *queueing delay and loss due to input buffer overflow!*



output port contention
at time t - only one red
packet can be transferred



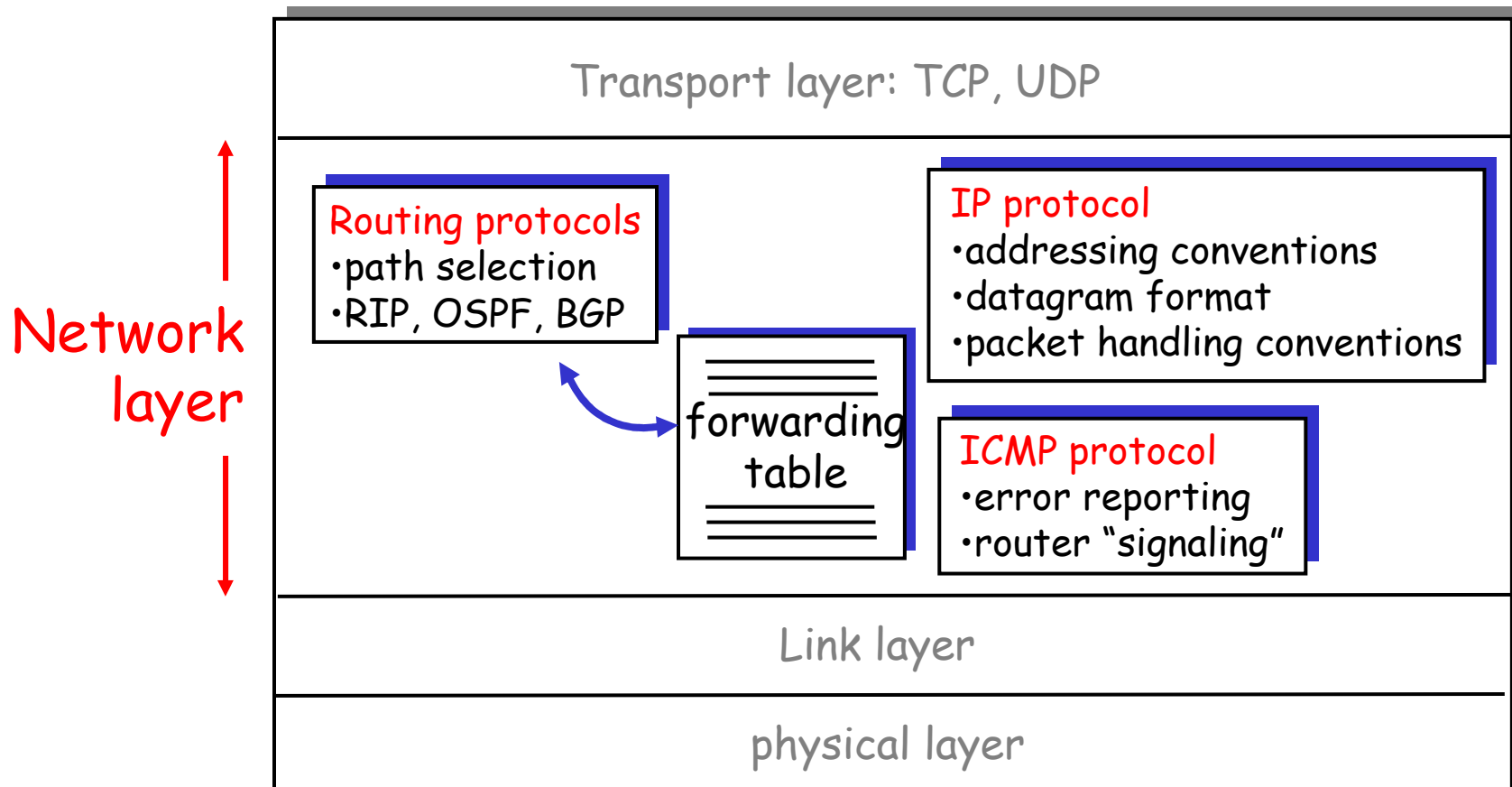
green packet
experiences HOL blocking

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

The Internet Network layer

Host, router network layer functions:



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

IP datagram format

IP protocol version number

header length (bytes)

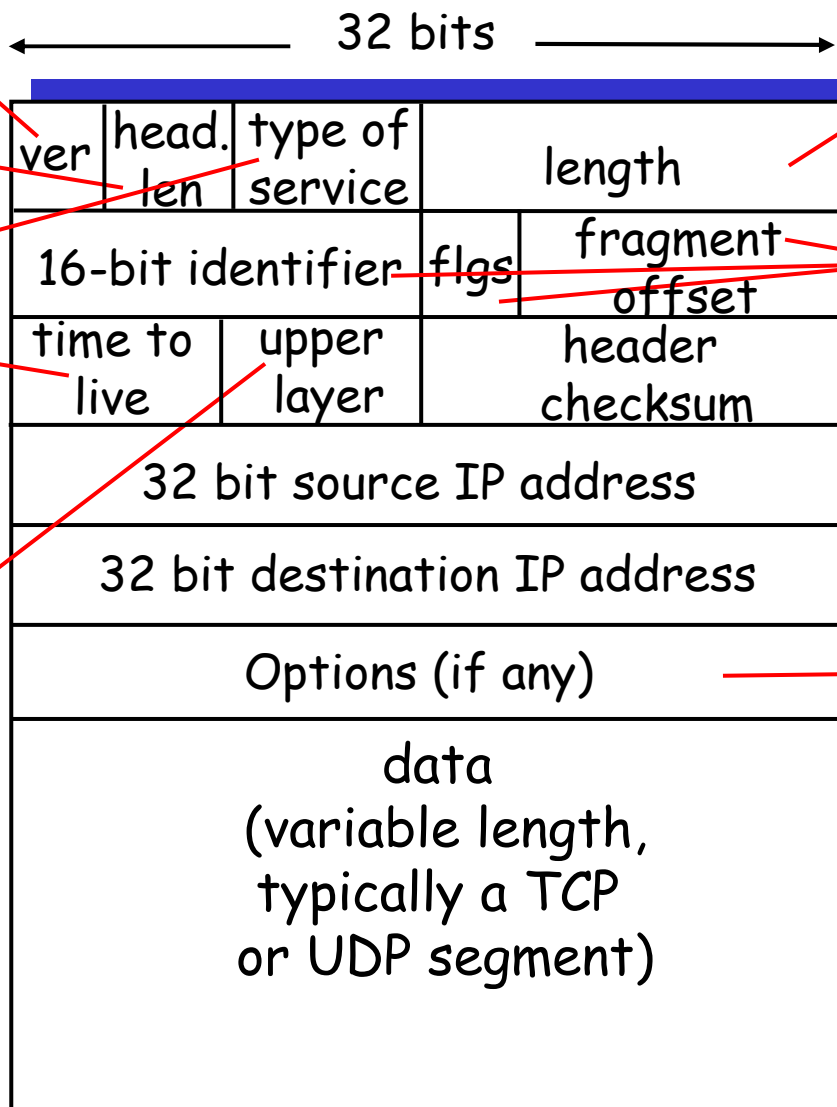
"type" of data

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

how much overhead with TCP?

- r 20 bytes of TCP
- r 20 bytes of IP
- r = 40 bytes + app layer overhead



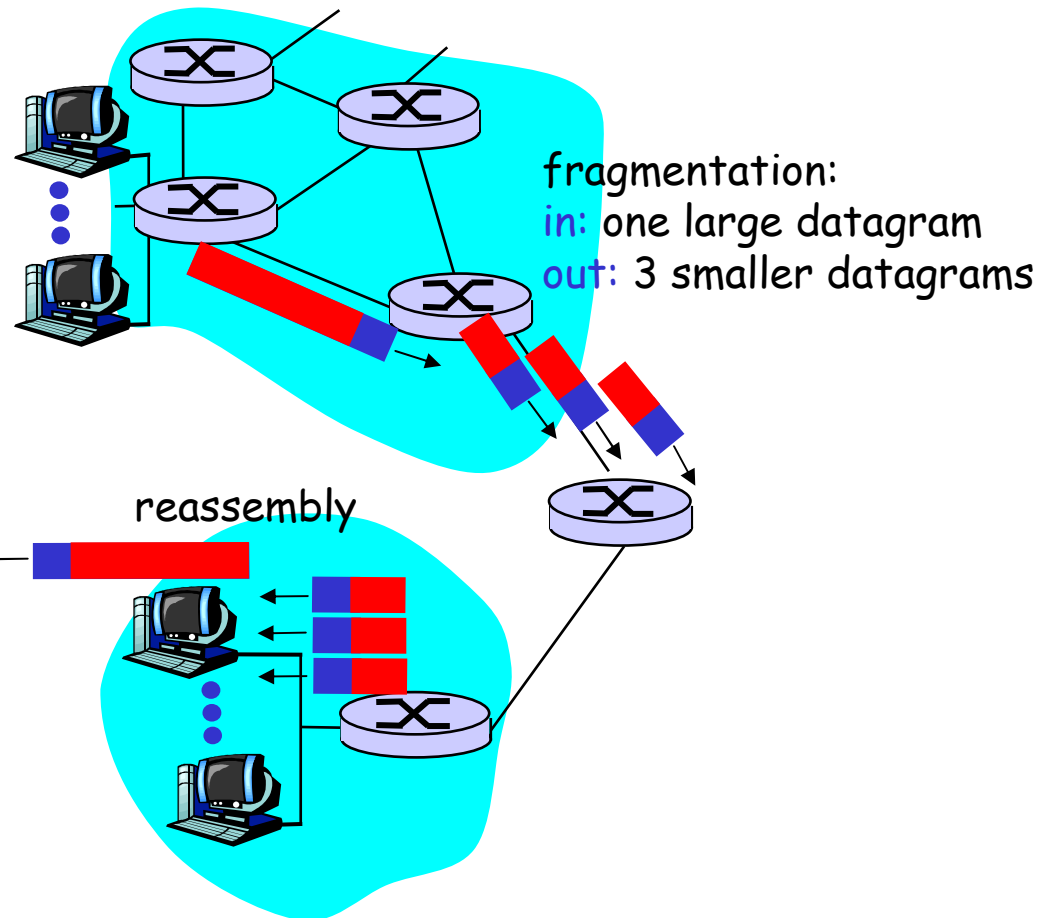
total datagram length (bytes)

for fragmentation/reassembly

E.g. timestamp, record route taken, specify list of routers to visit.

IP Fragmentation & Reassembly

- r network links have MTU (max.transfer size) - largest possible link-level frame.
 - m different link types, different MTUs
- r large IP datagram divided ("fragmented") within net
 - m one datagram becomes several datagrams
 - m "reassembled" only at final destination
 - m IP header bits used to identify, order related fragments



IP Fragmentation and Reassembly

Example

- r 4000 byte datagram
- r MTU = 1500 bytes

length	ID	fragflag	offset
=4000	=x	=0	=0

One large datagram becomes several smaller datagrams

1480 bytes in data field

offset =
1480/8

length	ID	fragflag	offset
=1500	=x	=1	=0

length	ID	fragflag	offset
=1500	=x	=1	=185

length	ID	fragflag	offset
=1040	=x	=0	=370

Offset: i dati saranno inseriti a partire Dal byte YY