

# Chapter -

## Network Layer

Introduction: 1-2

2

### Network layer: our goals

- understand principles behind network layer services, focusing on data plane:
  - network layer service models
  - forwarding versus routing
  - how a router works
  - addressing
  - generalized forwarding
  - Internet architecture
- instantiation, implementation in the Internet
  - IP protocol
  - NAT, middleboxes (rfc 3234)

Network Layer: 4-3

3

## Network layer: “data plane” roadmap

- **Network layer: overview**
  - data plane
  - control plane
- **What’s inside a router**
  - input ports, switching, output ports
  - buffer management, scheduling
- **IP: the Internet Protocol**
  - datagram format
  - addressing
  - network address translation
  - IPv6
- Generalized Forwarding, SDN
- Middleboxes

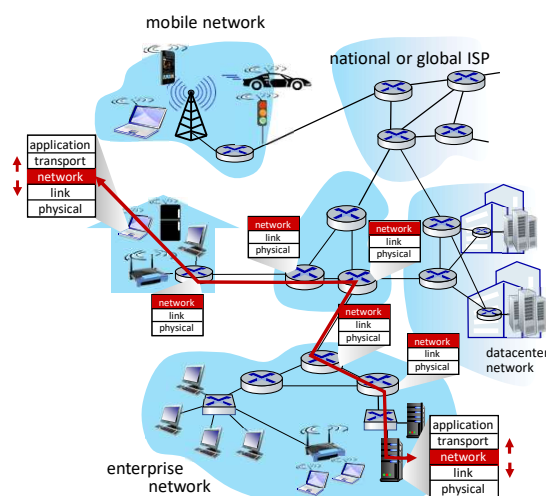


Network Layer: 4-4

4

## Network-layer services and protocols

- transport segment from sending to receiving host
  - sender: encapsulates segments into **datagrams**, passes to link layer
  - receiver: delivers segments to transport layer protocol
- network layer protocols in *every Internet device*: **hosts**, **routers**
- **routers**:
  - examines header fields in all IP datagrams passing through it
  - moves datagrams from input ports to output ports to transfer datagrams along end-end path



Network Layer: 4-6

6

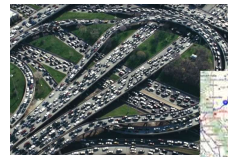
## Two key network-layer functions

### network-layer functions:

- *forwarding*: move packets from a router's input link to appropriate router output link
- *routing*: determine route taken by packets from source to destination
  - *routing algorithms*

### analogy: taking a trip

- *forwarding*: process of getting through single interchange
- *routing*: process of planning trip from source to destination



forwarding



routing

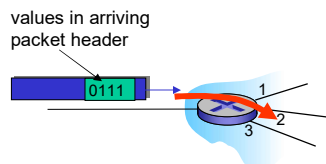
Network Layer: 4-7

7

## Network layer: data plane, control plane

### Data plane:

- *local*, per-router function
- determines how datagram arriving on router input port is forwarded to router output port



### Control plane

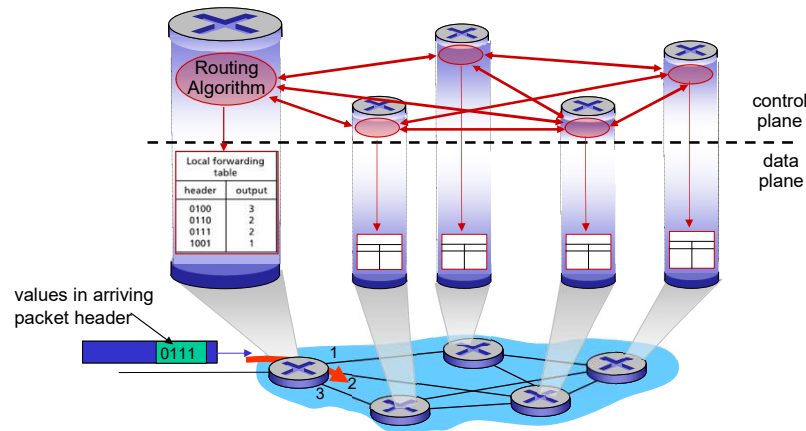
- *network-wide* logic
- determines how datagram is routed among routers along end-end path from source host to destination host
- two control-plane approaches:
  - *traditional routing algorithms*: implemented in routers
  - *software-defined networking (SDN)*: implemented in (remote) servers

Network Layer: 4-8

8

## Per-router control plane

Individual routing algorithm components *in each and every router* interact in the control plane

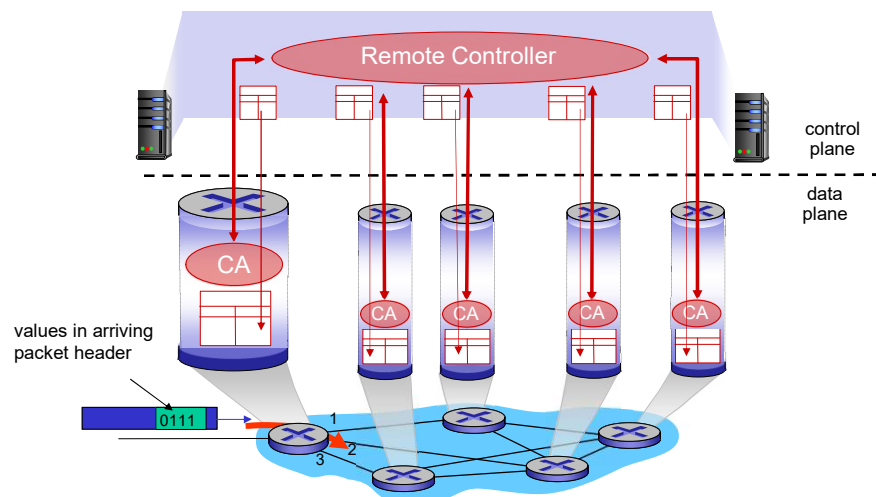


Network Layer: 4-9

9

## Software-Defined Networking (SDN) control plane

**Remote controller** computes, installs forwarding tables in routers



Network Layer: 4-10

10

## Network service model

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

example services for *individual* datagrams:

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

example services for a *flow* of datagrams:

- in-order datagram delivery
- guaranteed minimum bandwidth to flow

Network Layer: 4-11

11

## Network-layer service model

Network Architecture	Service Model	Quality of Service (QoS) Guarantees ?			
		Bandwidth	Loss	Order	Timing
Internet	best effort	none	no	no	no

Internet “best effort” service model

**No** guarantees on:

- i. successful datagram delivery to destination
- ii. timing or order of delivery
- iii. bandwidth available to end-end flow

Network Layer: 4-13

13

## Network-layer service model

Network Architecture	Service Model	Quality of Service (QoS) Guarantees ?			
		Bandwidth	Loss	Order	Timing
Internet	best effort	none	no	no	no
ATM	Constant Bit Rate	Constant rate	yes	yes	yes
ATM	Available Bit Rate	Guaranteed min	no	yes	no
Internet	Intserv Guaranteed (RFC 1633)	yes	yes	yes	yes
Internet	Diffserv (RFC 2475)	possible	possibly	possibly	no

Network Layer: 4-14

14

## Network layer: “data plane” roadmap

- Network layer: overview
  - data plane
  - control plane
- What’s inside a router
  - input ports, switching, output ports
  - buffer management, scheduling
- IP: the Internet Protocol
  - datagram format
  - addressing
  - network address translation
  - IPv6



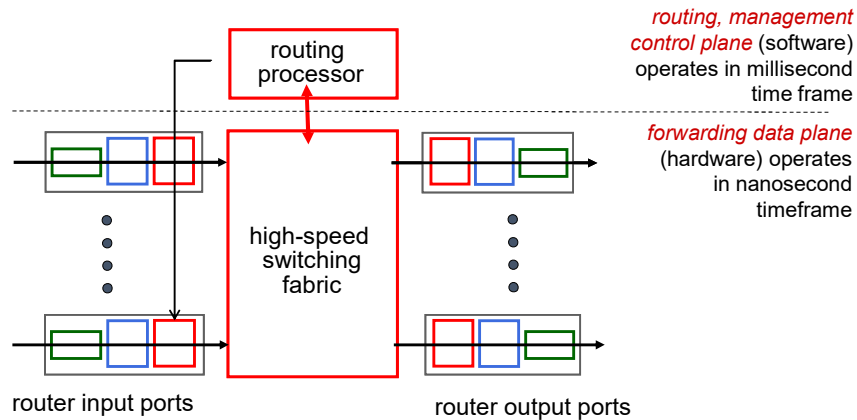
- Generalized Forwarding, SDN
- Middleboxes

Network Layer: 4-16

16

## Router architecture overview

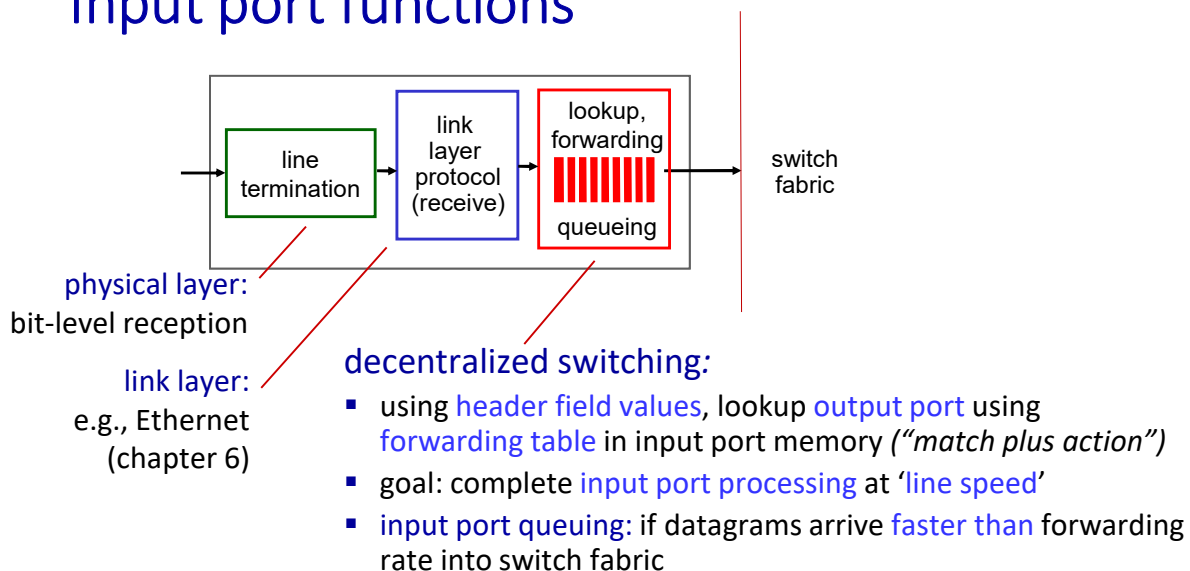
high-level view of generic router architecture:



Network Layer: 4-18

18

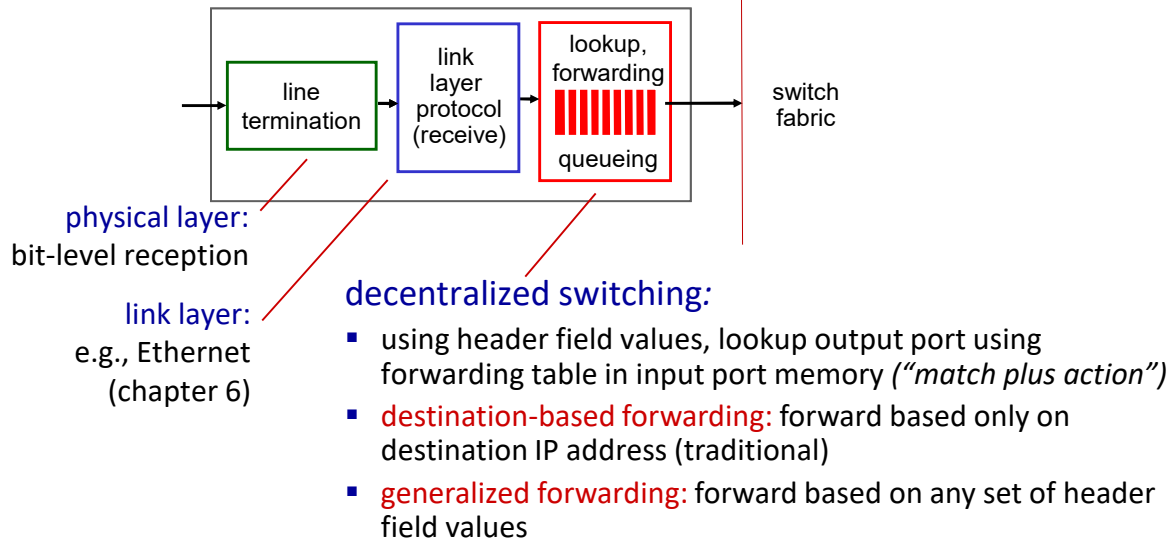
## Input port functions



Network Layer: 4-19

19

## Input port functions



Network Layer: 4-20

20

## Destination-based forwarding

*forwarding table*

Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through 11001000 00010111 00010000 00000100	n
11001000 00010111 00010000 00000100 through 11001000 00010111 00010000 00000111	3
11001000 00010111 00011000 11111111 through 11001000 00010111 00011001 00000000	2
11001000 00010111 00011111 11111111 otherwise	3

**Q:** but what happens if ranges don't divide up so nicely?

Network Layer: 4-21

21



## Longest prefix matching

### longest prefix match

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

11001000 00010111 00010110 10100001    which interface?  
 11001000 00010111 00011000 10101010    which interface?

Network Layer: 4-22

22

## Longest prefix matching

### longest prefix match

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

match!

examples:

11001000 00010111 00010110 10100001    which interface?  
 11001000 00010111 00011000 10101010    which interface?

Network Layer: 4-23

23

## Longest prefix matching

### longest prefix match

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

11001000 00010111 00010110 10100001	which interface?
11001000 00010111 00011000 10101010	which interface?

Network Layer: 4-24

24

## Longest prefix matching

### longest prefix match

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

11001000 00010111 00010110 10100001	which interface?
11001000 00010111 00011000 10101010	which interface?

Network Layer: 4-25

25

## Longest prefix matching

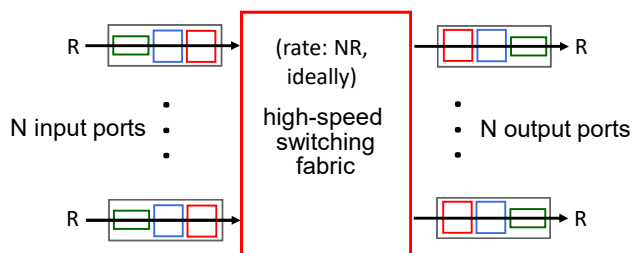
- we'll see *why* longest prefix matching is used shortly, when we study addressing
- longest prefix matching: often performed using **ternary content addressable memories (TCAMs)**
  - *content addressable*: present address to TCAM: retrieve address in one clock cycle, regardless of table size
  - Cisco Catalyst: ~1M routing table entries in TCAM

Network Layer: 4-26

26

## Switching fabrics

- transfer packet from input link to appropriate output link
- **switching rate**: rate at which packets can be transfer from inputs to outputs
  - often measured as multiple of input/output line rate
  - N inputs: **switching rate** N times **line rate** desirable

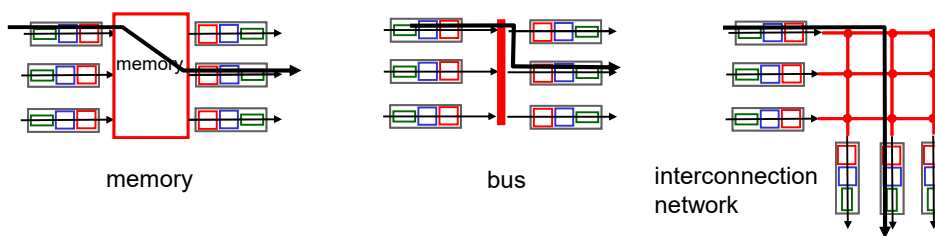


Network Layer: 4-27

27

## Switching fabrics

- transfer packet from input link to appropriate output link
- **switching rate**: rate at which packets can be transfer from inputs to outputs
  - often measured as multiple of input/output line rate
  - N inputs: switching rate N times line rate desirable
- three major types of switching fabrics:



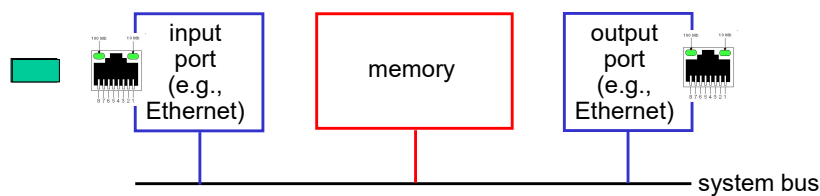
Network Layer: 4-28

28

## Switching via memory

### first generation routers:

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

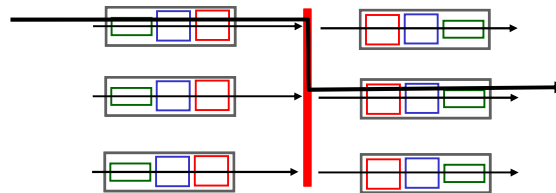


Network Layer: 4-29

29

## Switching via a bus

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

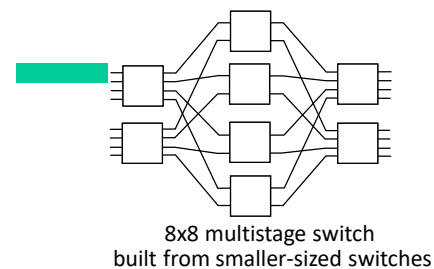
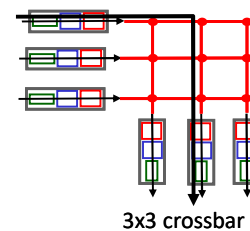


Network Layer: 4-30

30

## Switching via interconnection network

- **Crossbar**, **Clos networks**, other interconnection nets initially developed to connect processors in multiprocessor
- **multistage switch**:  **$n \times n$  switch** from multiple stages of smaller switches
- **exploiting parallelism**:
  - fragment datagram into fixed length cells on entry
  - switch cells through the fabric, reassemble datagram at exit

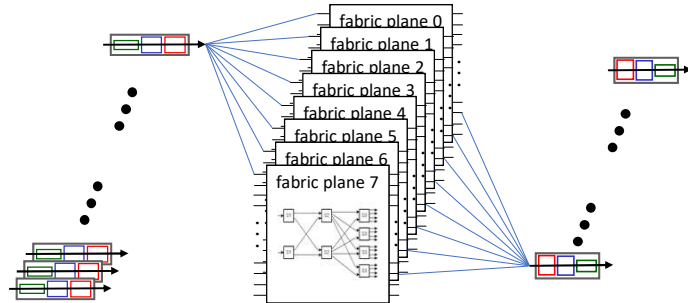


Network Layer: 4-31

31

## Switching via interconnection network

- scaling, using multiple switching “planes” in parallel:
  - speedup, scaleup via parallelism
- Cisco CRS router:
  - basic unit: 8 switching planes
  - each plane: 3-stage interconnection network
  - up to 100's Tbps switching capacity

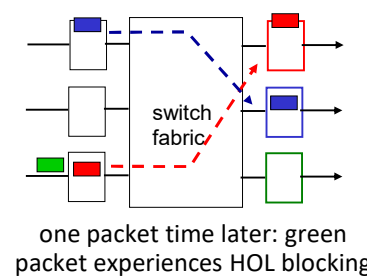
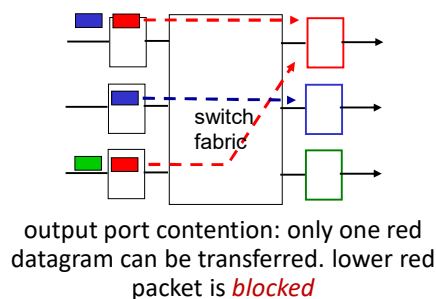


Network Layer: 4-32

32

## Input port queuing

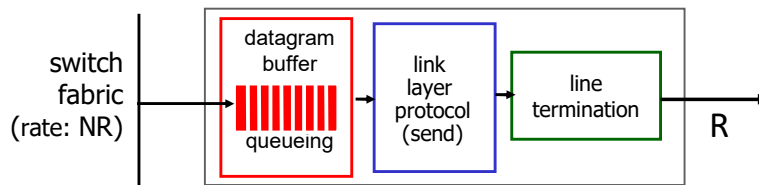
- If switch fabric slower than input ports combined -> queueing may occur at input queues
  - queueing delay and loss due to input buffer overflow!
- **Head-of-the-Line (HOL) blocking:** queued datagram at front of queue prevents others in queue from moving forward



Network Layer: 4-33

33

## Output port queuing



This is a really important slide

- **Buffering** required when datagrams arrive from fabric faster than link transmission rate. **Drop policy**: which datagrams to drop if no free buffers?
- **Scheduling discipline** chooses among queued datagrams for transmission



Datagrams can be lost due to congestion, lack of buffers

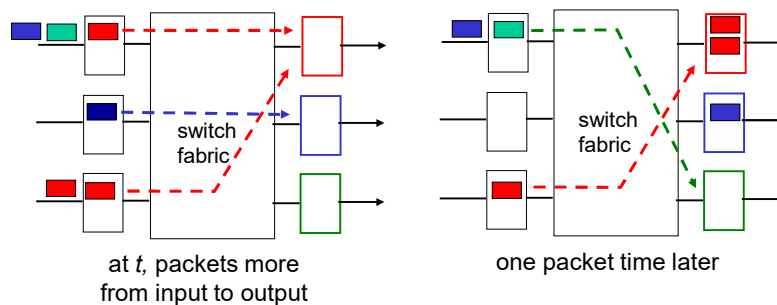


Priority scheduling – who gets best performance, network neutrality

Network Layer: 4-34

34

## Output port queuing



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

Network Layer: 4-35

35

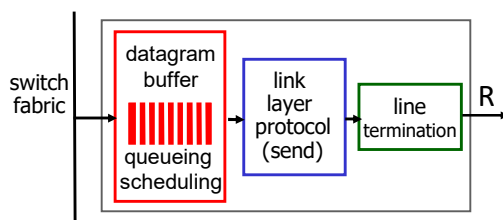
## How much buffering?

- **RFC 3439** rule of thumb: **average buffering** equal to “typical” **RTT** (say 250 msec) times **link capacity C**
  - e.g., C = 10 Gbps link: 2.5 Gbit buffer
- more recent recommendation: with  $N$  flows, **buffering** equal to
 
$$\frac{RTT \cdot C}{\sqrt{N}}$$
- but **too much buffering** can increase **delays** (particularly in home routers)
  - **long RTTs**: poor performance for realtime apps, sluggish TCP response
  - recall **delay-based congestion control**: “keep **bottleneck link** just **full enough** (busy) but no fuller”

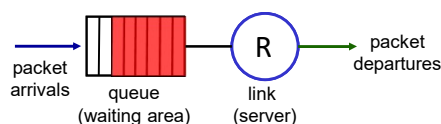
Network Layer: 4-36

36

## Buffer Management



### Abstraction: queue



### buffer management:

- **drop**: which packet to add, drop when **buffers are full**
  - **tail drop**: drop arriving packet
  - **priority drop**: drop/remove on priority basis
- **marking**: which packets to mark to signal congestion (ECN, RED)

Network Layer: 4-37

37



## Packet Scheduling: FCFS

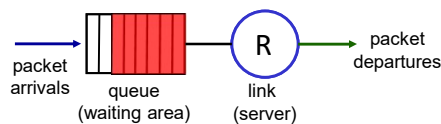
**packet scheduling:** deciding which packet to send next on link

- first come, first served
- priority
- round robin
- weighted fair queueing

**FCFS:** packets transmitted in order of arrival to output port

- also known as: First-in-first-out (FIFO)
- real world examples?

Abstraction: queue



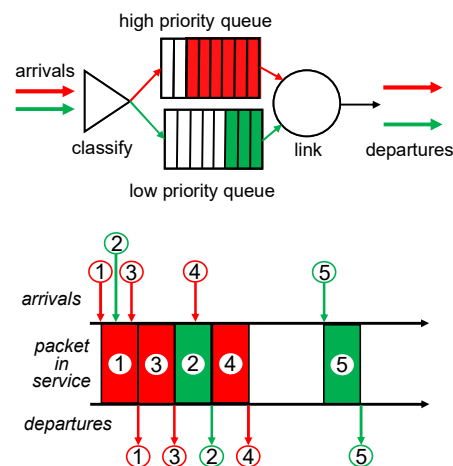
Network Layer: 4-38

38

## Scheduling policies: priority

**Priority scheduling:**

- arriving traffic classified, queued by class
  - any header fields can be used for classification
- send packet from highest priority queue that has buffered packets
  - FCFS within priority class



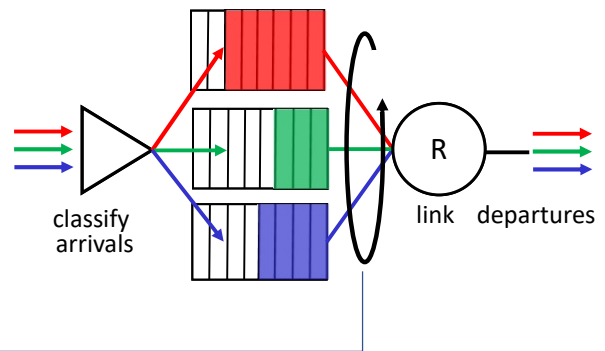
Network Layer: 4-39

39

## Scheduling policies: round robin

### Round Robin (RR) scheduling:

- arriving traffic classified, queued by class
  - any header fields can be used for classification
- server cyclically, repeatedly scans class queues, sending one complete packet from each class (if available) in turn



Network Layer: 4-40

40

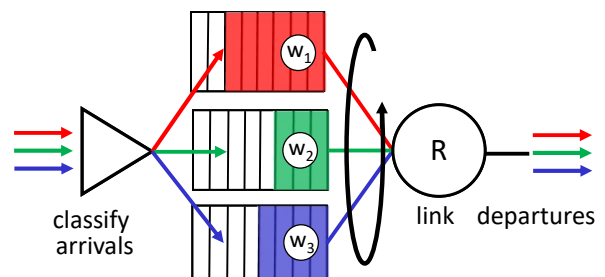
## Scheduling policies: weighted fair queueing

### Weighted Fair Queuing (WFQ):

- generalized Round Robin
- each class,  $i$ , has weight,  $w_i$ , and gets weighted amount of service in each cycle:

$$\frac{w_i}{\sum_j w_j}$$

- minimum bandwidth guarantee (per-traffic-class)



Network Layer: 4-41

41