# Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol
- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet
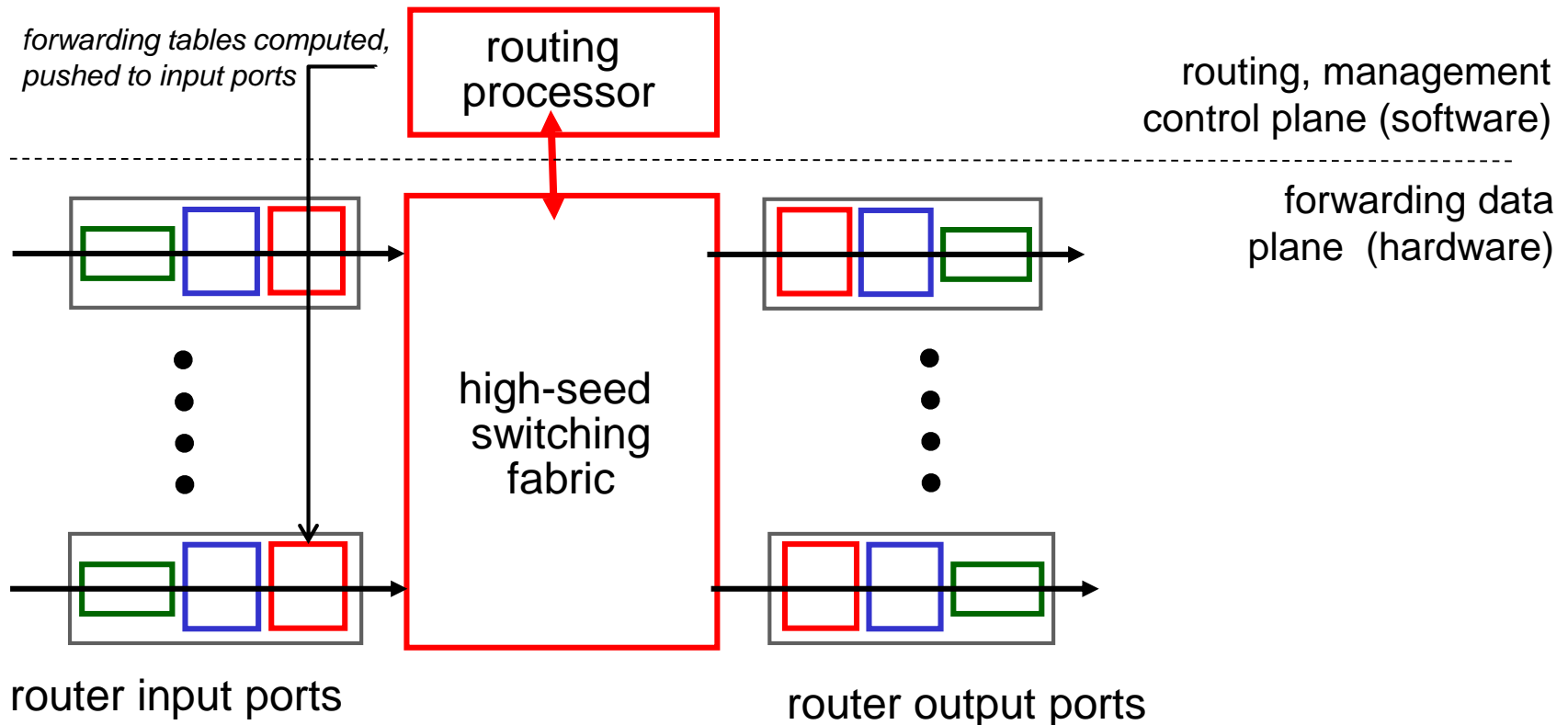- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

# What is a router?

❖ Modern routers have varying amounts of input ports and output ports.

❖ Home-grade, SMB "routers" typically have the following internal components:

- Switch
- Router
- Firewall
- Wireless Radio

❖ Advanced devices often have:

- VPN
- Port Forwarding
- QoS (Quality of Service)
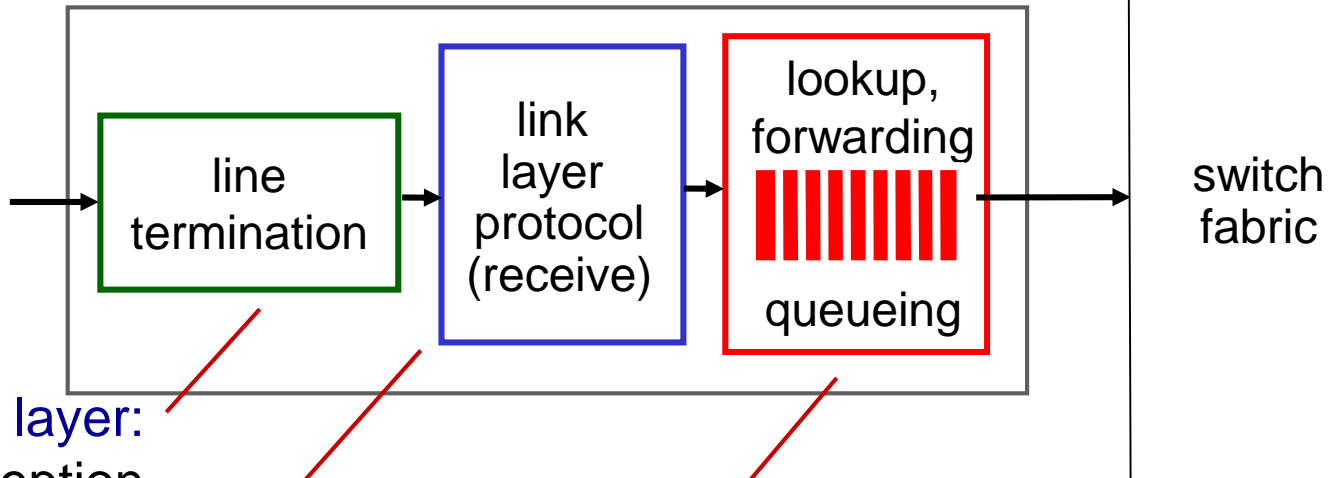- Separate control over each port

# Router (only) architecture overview

two key router functions:

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

*forwarding tables computed, pushed to input ports*

routing processor

routing, management control plane (software)

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

forwarding data plane (hardware)

high-seed switching fabric

router input ports

router output ports

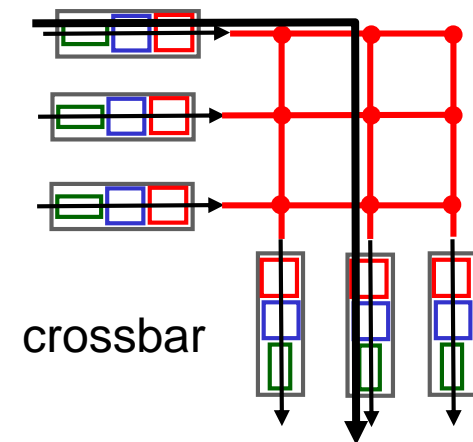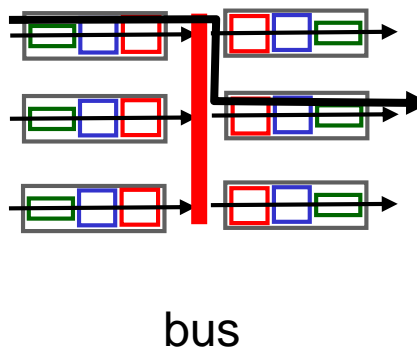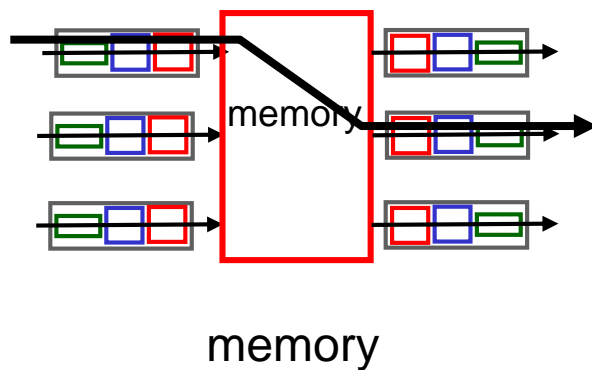# Input port functions



physical layer:
bit-level reception

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

decentralized switching:

❖ given datagram dest., lookup output port using forwarding table in input port memory (*"match plus action"*)

❖ goal: complete input port processing at 'line speed'

❖ queuing: if datagrams arrive faster than forwarding rate into switch fabric
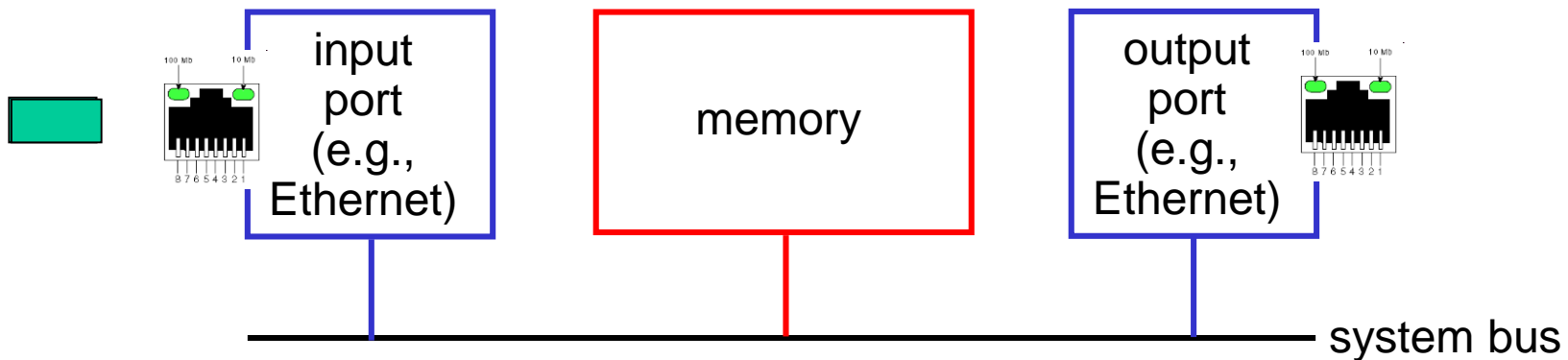
# Switching fabrics

❖ transfer packet from input buffer to appropriate output buffer

❖ 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 types of switching fabrics

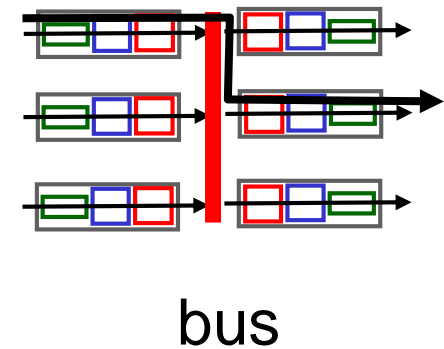memory                          bus                          crossbar

# 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)

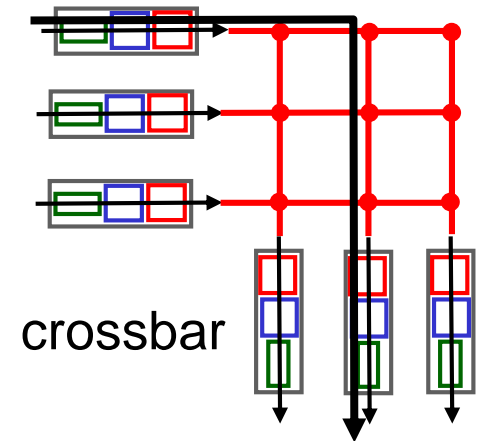# Switching via a bus

❖ datagram from input port memory copied 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 and enterprise routers
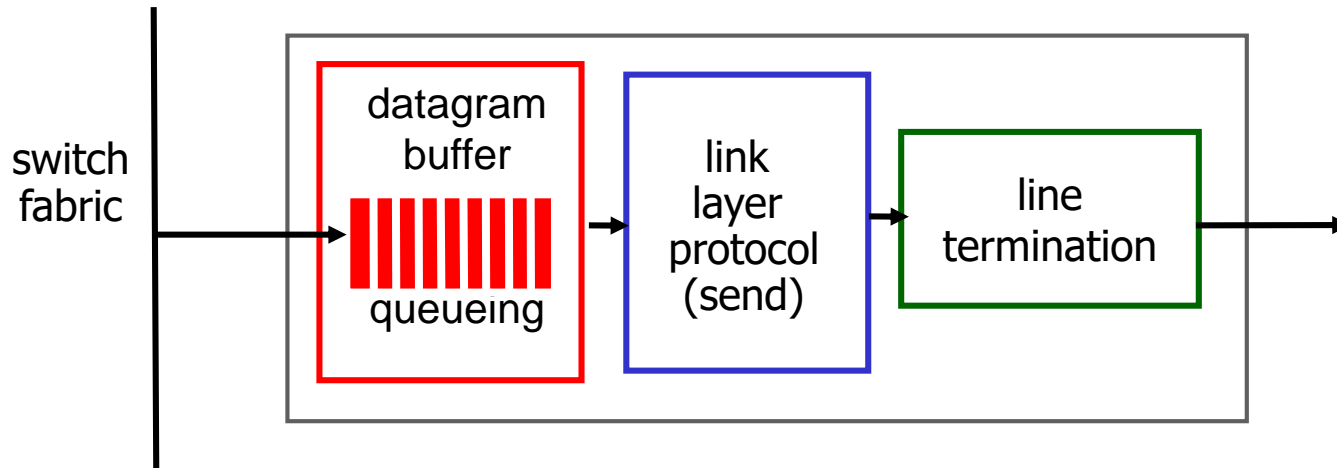
bus

# Switching via interconnection network

❖ overcome bus bandwidth limitations

❖ banyan networks, crossbar, other interconnection nets initially developed to connect processors in multiprocessor



crossbar

❖ advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.

❖ Cisco 12000: switches 60 Gbps through the interconnection network

# Output ports

*This slide is HUGELY important!*



switch fabric → datagram buffer (queueing) → link layer protocol (send) → line termination →

❖ *buffering* required from fabric faster rate

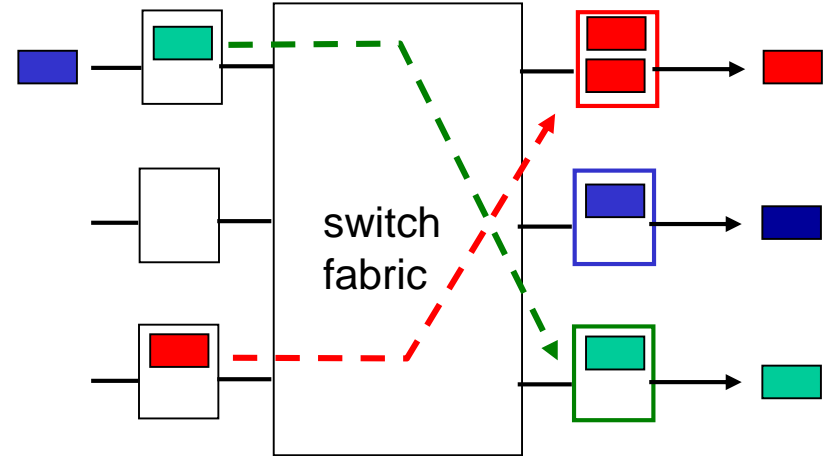Datagram (packets) can be lost due to congestion, lack of buffers

❖ *scheduling* datagrams

Priority scheduling – who gets best performance, network neutrality

# Output port queueing



at time *t,* packets move
from input to output

one packet time later

❖ **buffering when arrival rate via switch exceeds output line speed**

❖ *queueing (delay) and loss due to output port buffer overflow!*

# How much buffering?
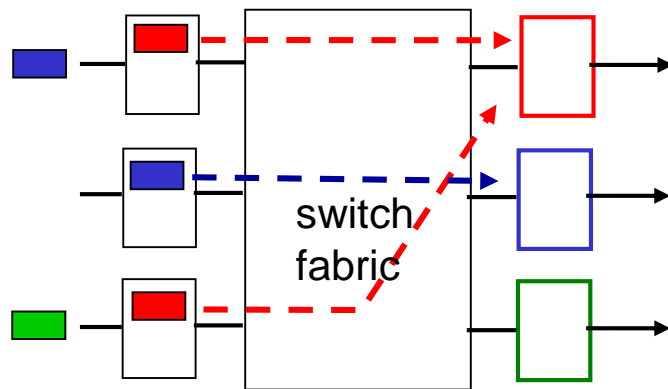
❖ RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity C
  ▪ e.g., 250ms * 10 Gbps => 2.5 Gbit buffer

❖ recent recommendation: with *N* TCP "flows", buffering equal to $\dfrac{\text{RTT} \cdot \text{C}}{\sqrt{N}}$

thus, with RTT = 0.25s, C = 10 Gbps, and 10 flows:
  ▪ e.g., (0.25s * 10 Gbps) / sqrt(10) => 0.79 Gbit buffer

# Input port queuing

❖ 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



output port contention:
Assume only one **red** datagram
can be transferred per time *t*.
*lower red packet is blocked*

one packet time later: **green**
packet experiences HOL
blocking - can't be sent, even
though it's desired output is
not busy!

# Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol
- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet
- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

# The Internet network layer

host, router network layer functions:

transport layer: TCP, UDP

**network layer**

*routing protocols*
• path selection
• RIP, OSPF, BGP

forwarding table

*IP protocol*
• addressing conventions
• datagram format
• packet handling conventions

*ICMP protocol*
• error reporting
• router "signaling"

link layer

physical layer

# 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

32 bits

total datagram length (bytes)

for fragmentation/ reassembly

| ver | head. len | type of service | length | | |
|------|------|------|------|------|------|
| 16-bit identifier | | | flgs | fragment offset | |
| time to live | | upper layer | | header checksum | |
| 32 bit source IP address | | | | | |
| 32 bit destination IP address | | | | | |
| options (if any) | | | | | |
| data (variable length, typically a TCP or UDP segment) | | | | | |

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

*how much overhead?*
- 20 bytes of TCP
- 20 bytes of IP
- = 40 bytes + app layer overhead
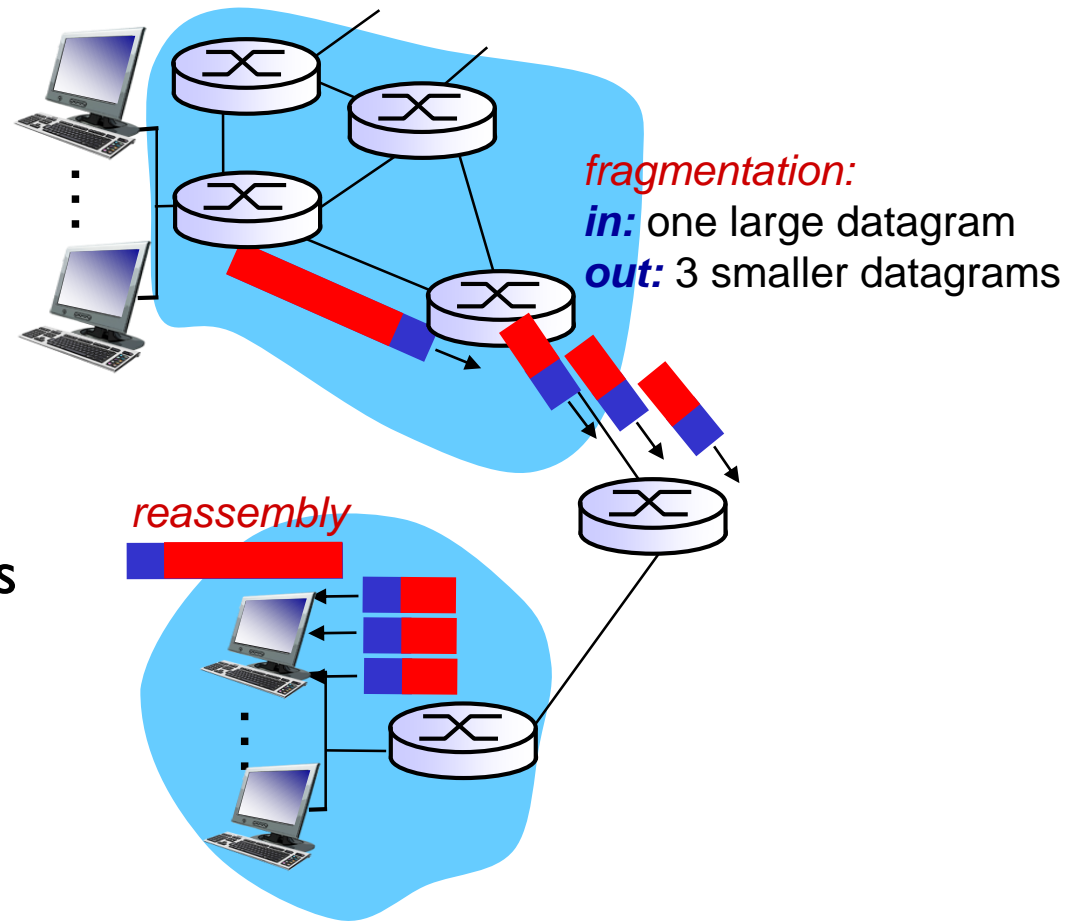
# IP fragmentation, reassembly

❖ network links have MTU (max.transfer size) - largest possible link-level frame
  ▪ different link types, different MTUs

❖ large IP datagram divided ("fragmented") within net
  ▪ one datagram becomes several datagrams
  ▪ "reassembled" only at final destination
  ▪ IP header bits used to identify, order related fragments

*fragmentation:*
*in:* one large datagram
*out:* 3 smaller datagrams

*reassembly*

# IP fragmentation, reassembly

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

*example:*

❖ 4000 byte datagram

❖ MTU = 1500 bytes

*one large datagram becomes several smaller datagrams*

1480 bytes in data field

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

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

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

offset = 1480/8

# Chapter 4: outline

4.1 introduction

4.2 virtual circuit and
   datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol
- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet
- RIP
- OSPF
- BGP

4.7 broadcast and multicast
   routing