# CS 305: Computer Networks
## Fall 2022

**Network Layer – The Control Plane**

**Ming Tang**

Department of Computer Science and Engineering
Southern University of Science and Technology (SUSTech)

# Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state

- distance vector

5.3 intra-AS routing in the Internet: OSPF
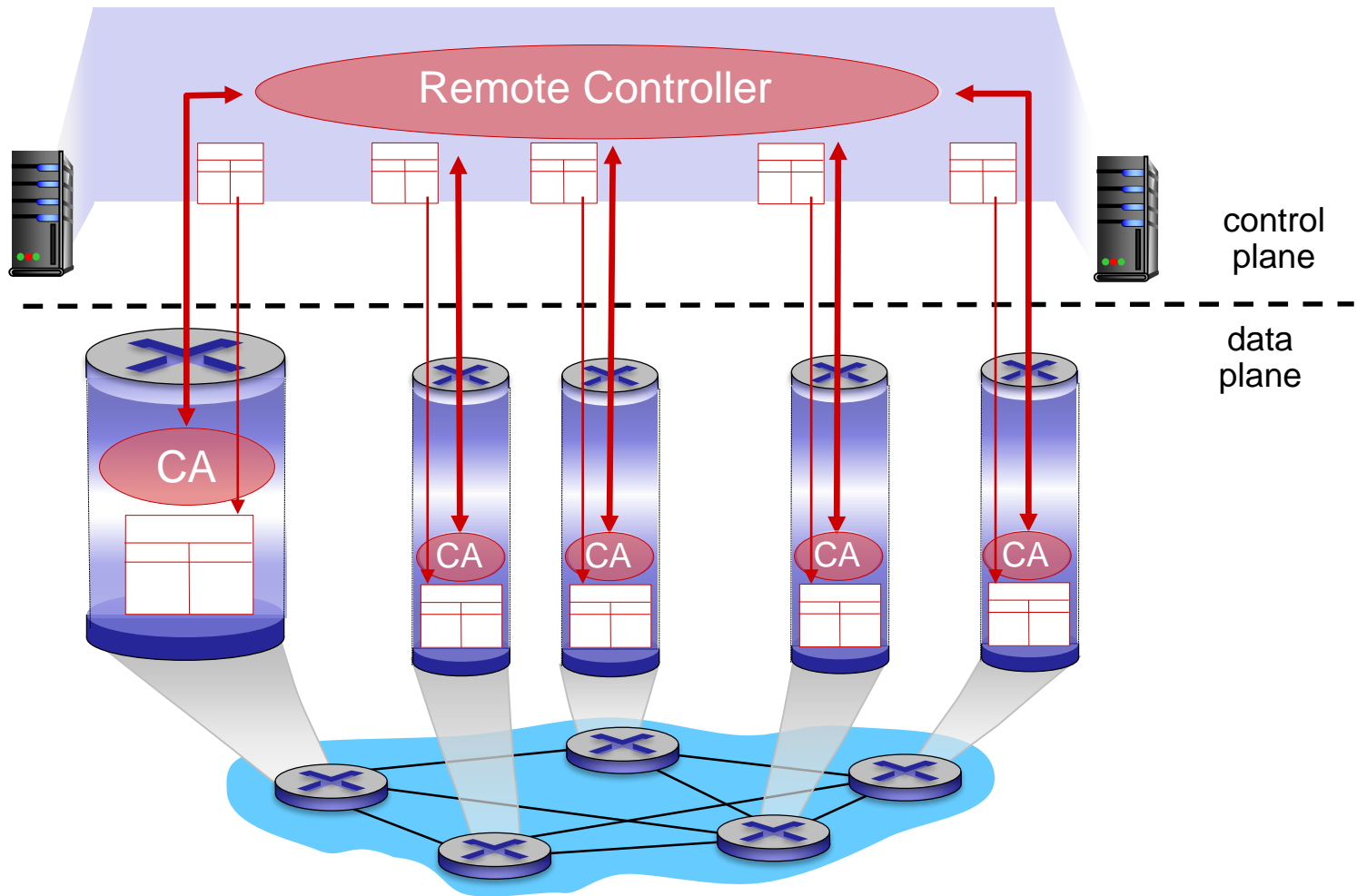
5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP

# Recall: SDN logically centralized control plane

A distinct (typically remote) controller interacts with local control agents (CAs) in routers to compute forwarding tables
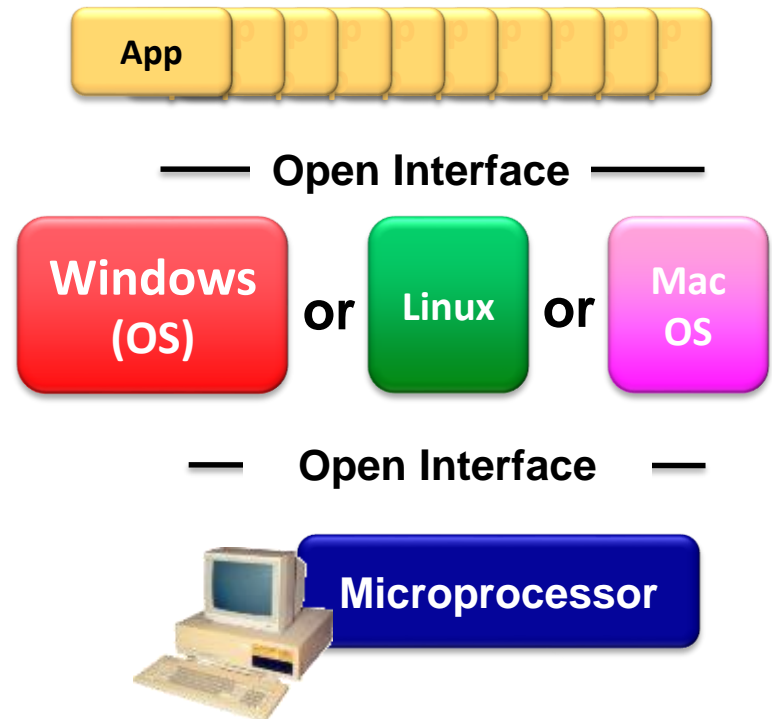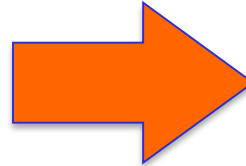
# Software defined networking (SDN)

*Why* a *logically centralized* control plane?

- ■ easier network management: avoid router misconfigurations, greater flexibility of traffic flows

- ■ table-based forwarding (recall OpenFlow API) allows "programming" routers

  - • centralized "programming" easier: compute tables centrally and distribute

  - • distributed "programming" more difficult: compute tables as result of distributed algorithm (protocol) implemented in each and every router

- ■ open (non-proprietary) implementation of control plane

# Analogy: mainframe to PC evolution*

**Specialized Applications**

**Specialized Operating System**

**Specialized Hardware**

App

— Open Interface —

**Windows (OS)** or **Linux** or **Mac OS**

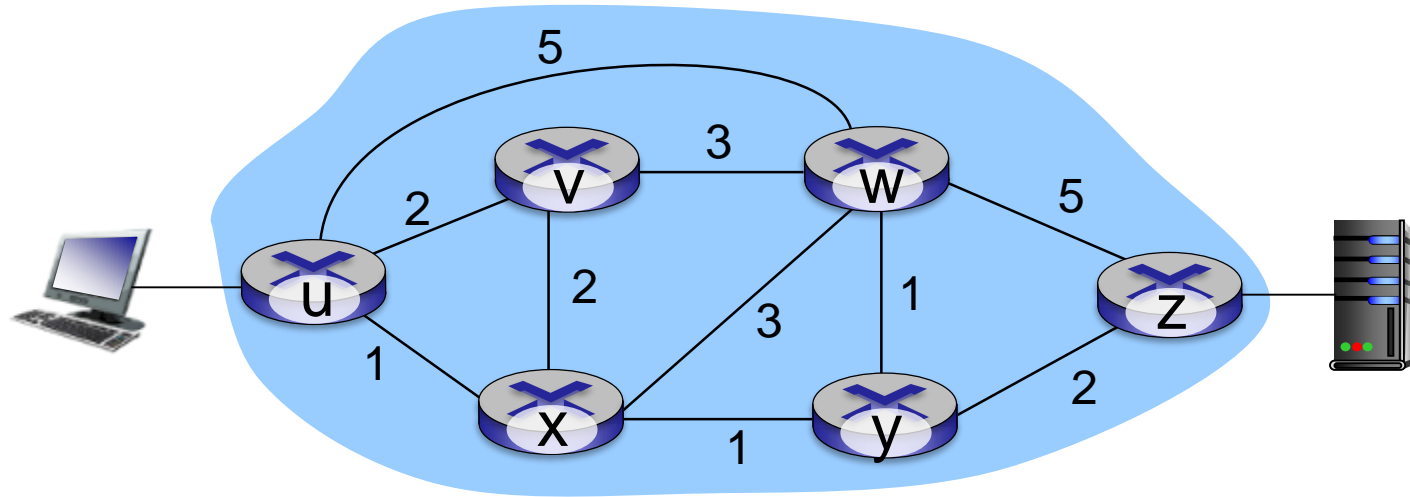— Open Interface —

**Microprocessor**

Vertically integrated
Closed, proprietary
Slow innovation
Small industry

Horizontal
Open interfaces
Rapid innovation
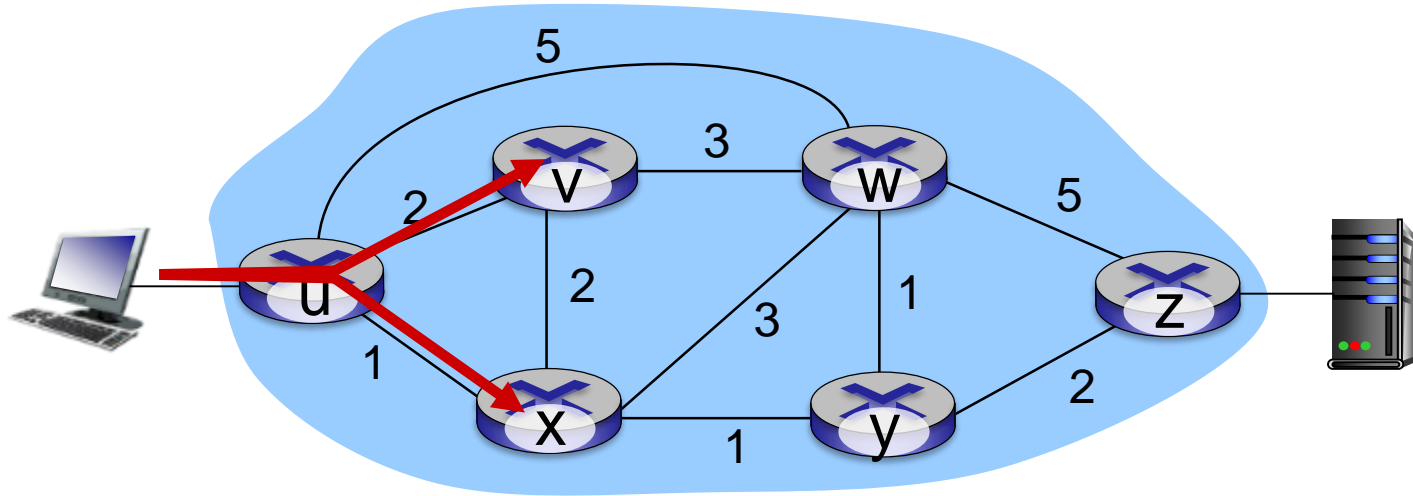Huge industry

# Traffic engineering: difficult traditional routing



*Q:* what if network operator wants u-to-z traffic to flow along *uvw*z, x-to-z traffic to flow *xwyz*?

*A:* need to define link weights so traffic routing algorithm computes routes accordingly (or need a new routing algorithm)!

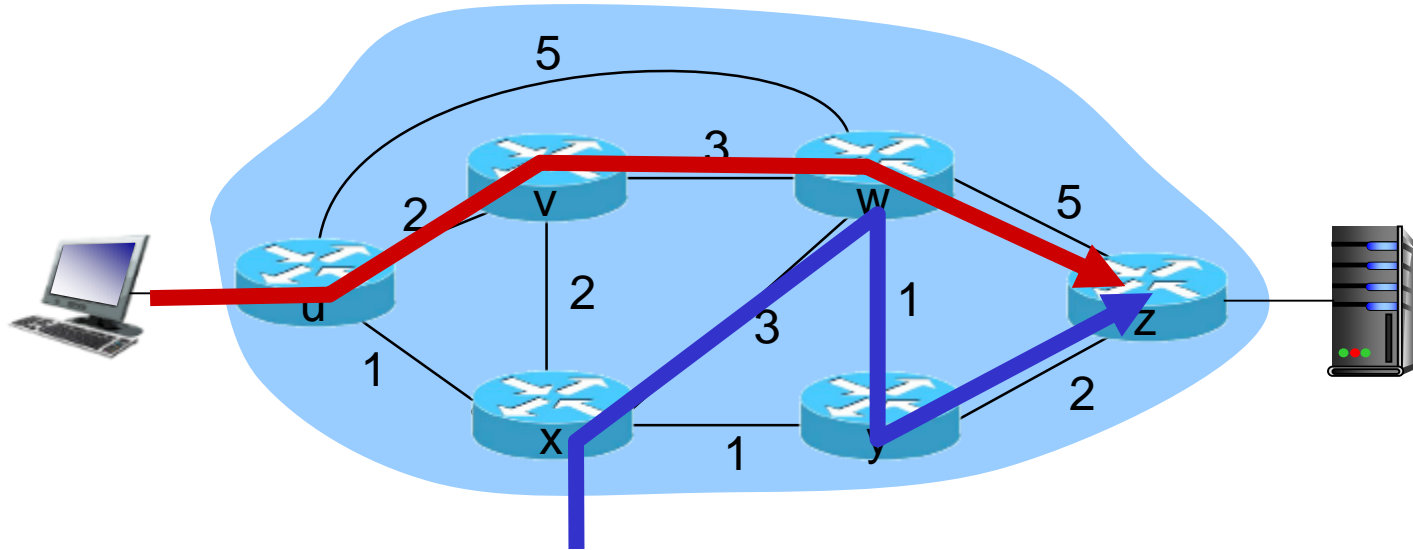But the link weights cannot be directly set to certain number

# Traffic engineering: difficult



*Q:* what if network operator wants to split  u-to-z traffic along uvwz *and* uxyz (load balancing)?

*A:* can't do it (or need a new routing algorithm)

# Traffic engineering: difficult



*Q:* what if w wants to route blue and red traffic differently?

*A:* can't do it (with destination based forwarding, and LS, DV routing)

# Software defined networking (SDN)
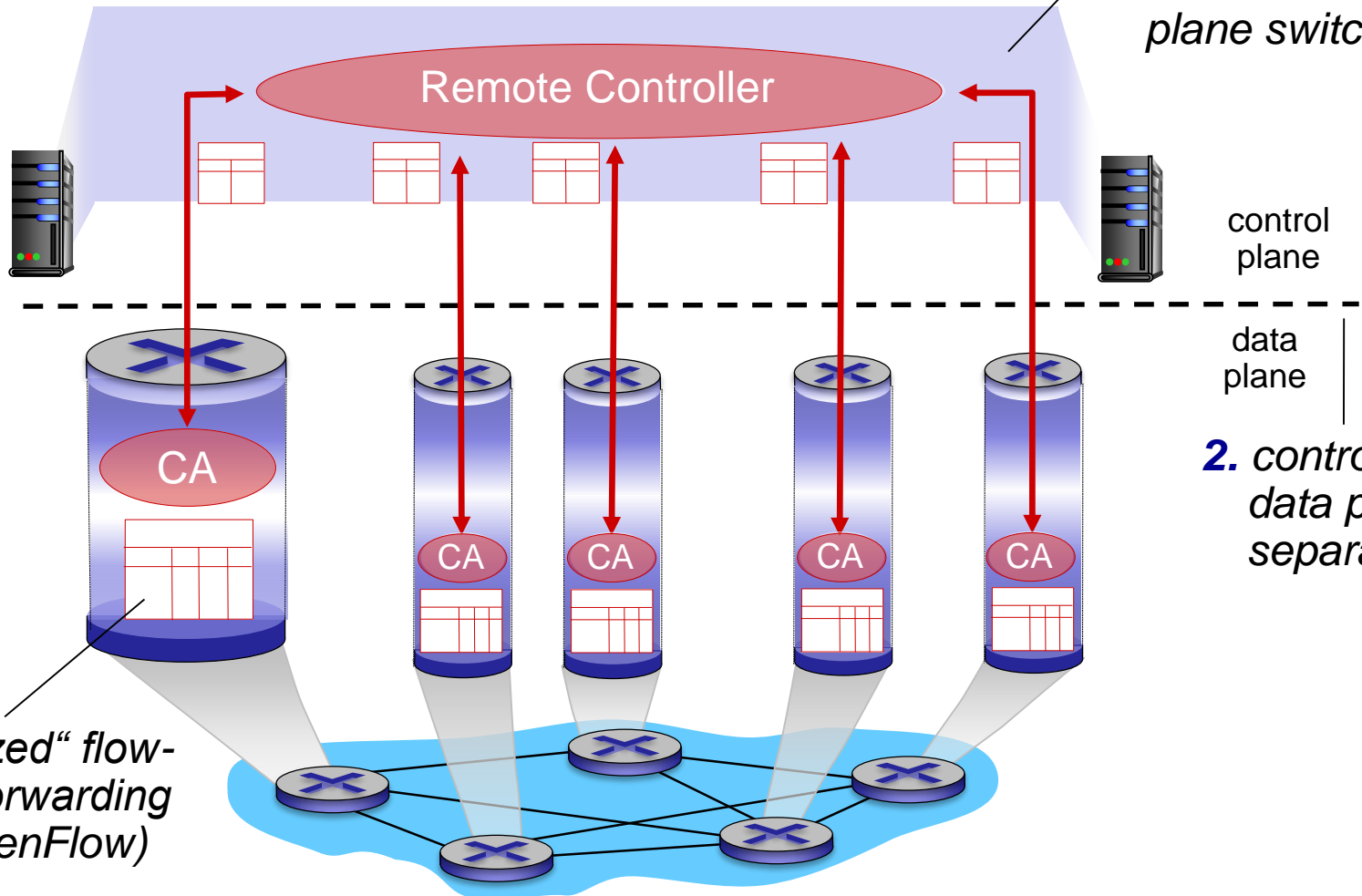
**4.** *programmable control applications*

routing

access control

. . .

load balance

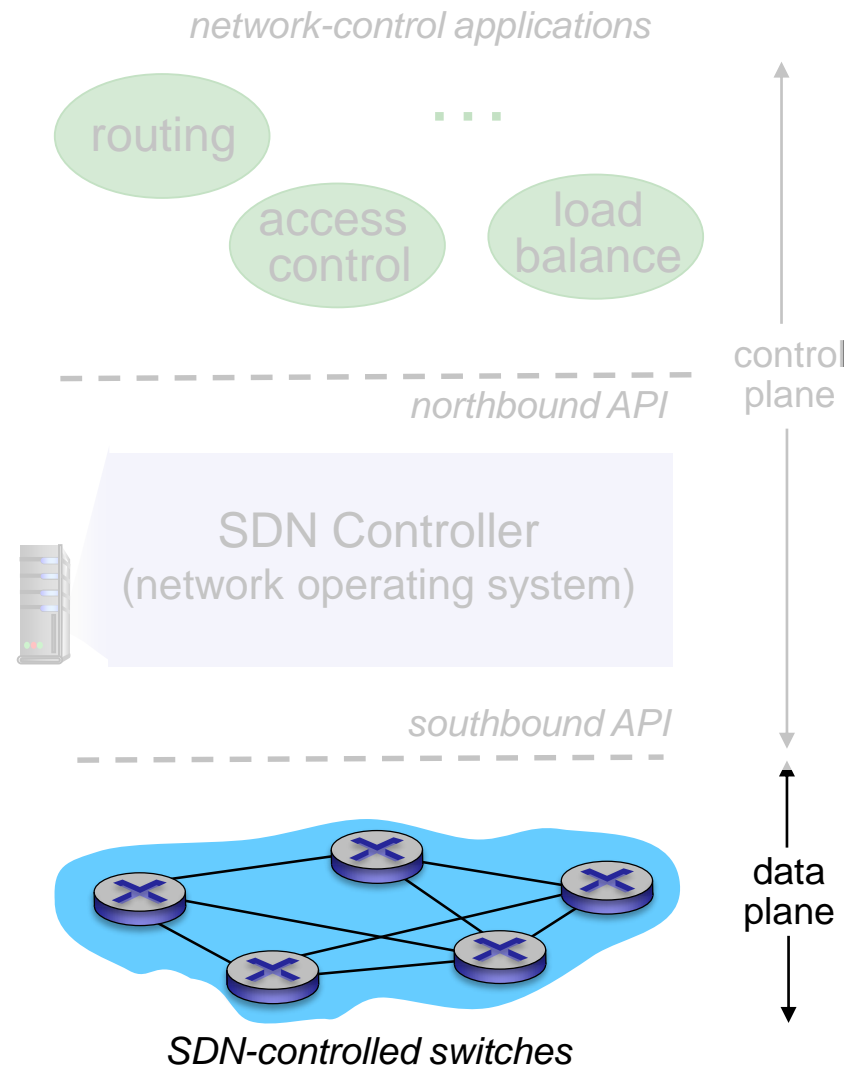**3.** *control plane functions external to data-plane switches*

Remote Controller

control plane

data plane

**2.** *control, data plane separation*

CA

CA

CA

CA

CA

**1:** *generalized" flow-based" forwarding (e.g., OpenFlow)*

# SDN perspective: data plane switches
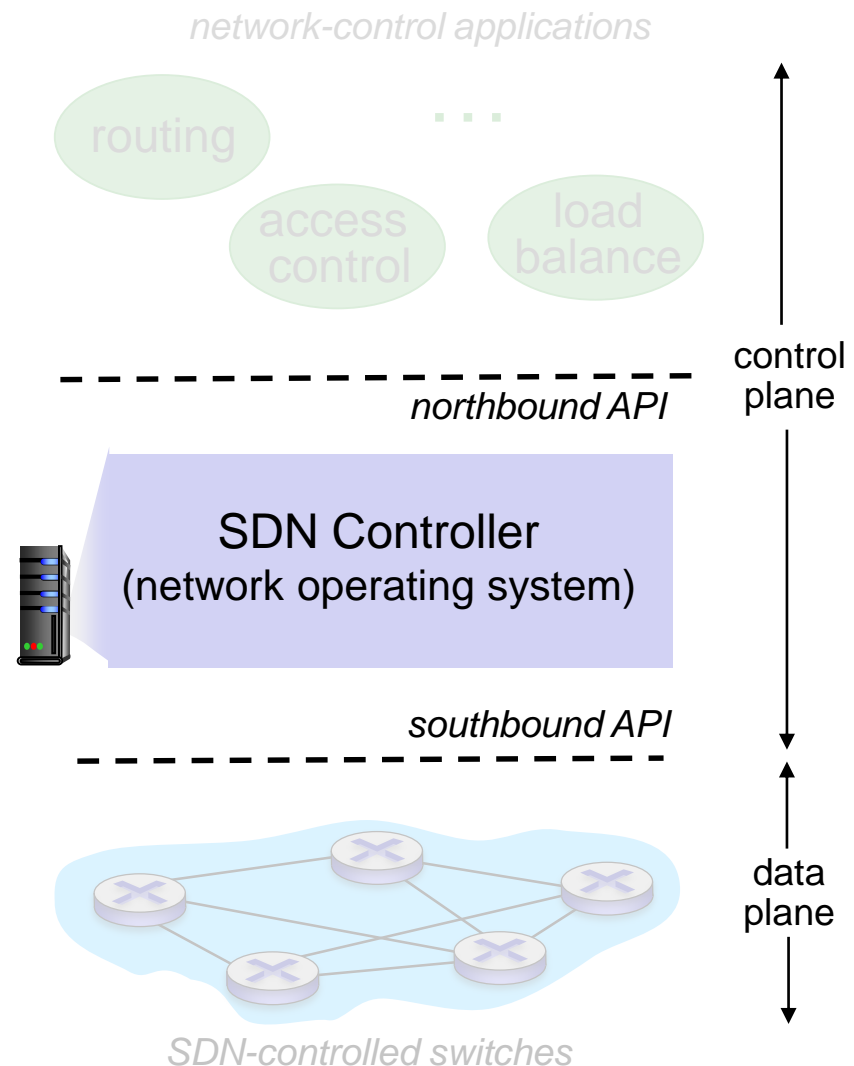
## *Data plane switches*

- fast, simple, commodity switches implementing generalized data-plane forwarding (Section 4.4) in hardware

- switch flow table computed, installed by controller

- API for table-based switch control (e.g., OpenFlow)

  - defines what is controllable and what is not

- protocol for communicating with controller (e.g., OpenFlow)

network-control applications

routing

. . .

access control

load balance

control plane

*northbound API*

SDN Controller
(network operating system)

*southbound API*

data plane

*SDN-controlled switches*

# SDN perspective: SDN controller
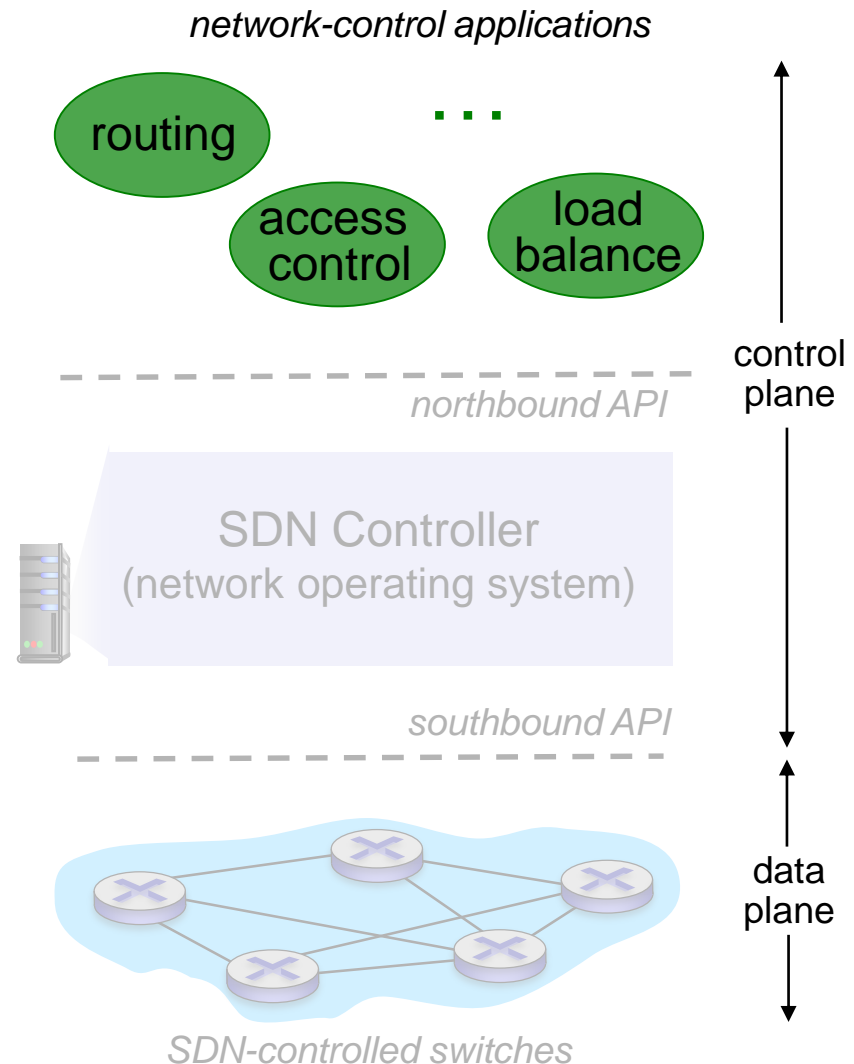
*SDN controller (network OS):*

- maintain network state information

- interacts with network control applications "above" via northbound API

- interacts with network switches "below" via southbound API

- implemented as distributed system for performance, scalability, fault-tolerance, robustness

*network-control applications*

routing

access control

load balance

- - -

*northbound API*

SDN Controller
(network operating system)

*southbound API*

*SDN-controlled switches*

control plane

data plane

# SDN perspective: control applications

*network-control apps:*

- "brains" of control: implement control functions using lower-level services, API provided by SND controller
- *unbundled:* can be provided by 3rd party: distinct from routing vendor, or SDN controller
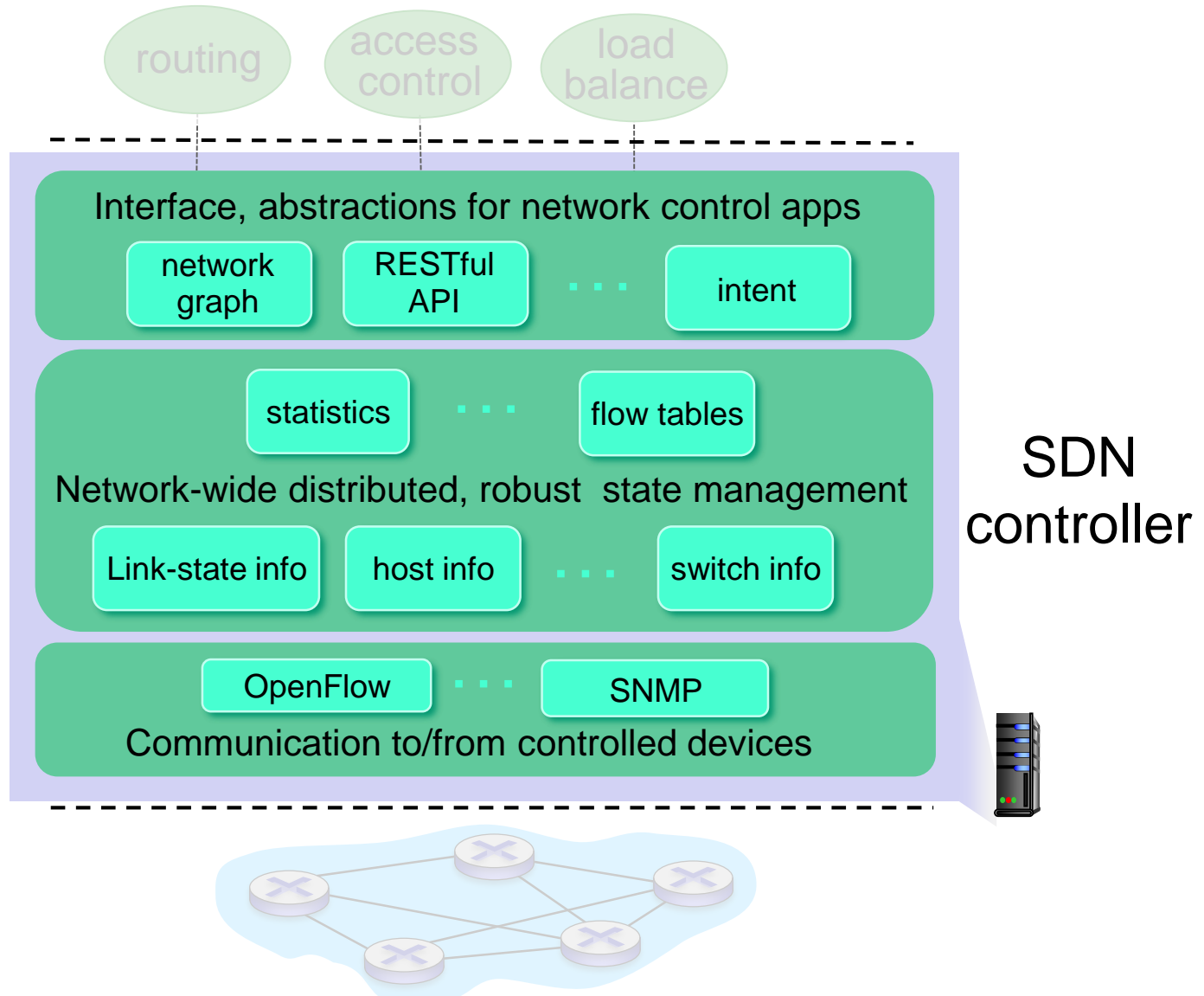
*network-control applications*

routing · · ·

access control

load balance

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

*northbound API*

control plane

SDN Controller
(network operating system)

*southbound API*

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

data plane

*SDN-controlled switches*

# Components of SDN controller

routing   access control   load balance

Interface layer to network control apps: abstractions API

**Interface, abstractions for network control apps**

network graph   RESTful API   · · ·   intent

Network-wide state management layer: state of networks links, switches, services: a *distributed database*

statistics   · · ·   flow tables

**Network-wide distributed, robust state management**

Link-state info   host info   · · ·   switch info

*communication layer*: communicate between SDN controller and controlled switches

OpenFlow   · · ·   SNMP

**Communication to/from controlled devices**

SDN controller

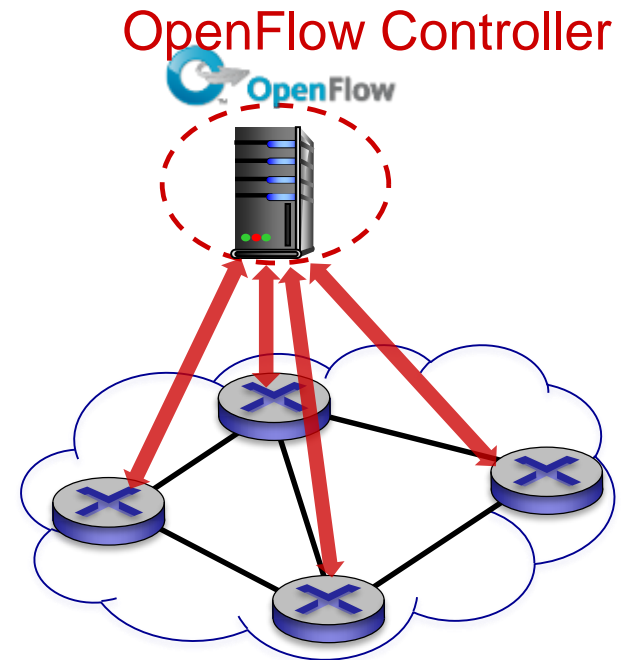# OpenFlow protocol

OpenFlow Controller



- operates between controller, switch
- TCP used to exchange messages
- OpenFlow messages:
  - controller-to-switch
  - switch to controller

# OpenFlow: controller-to-switch messages

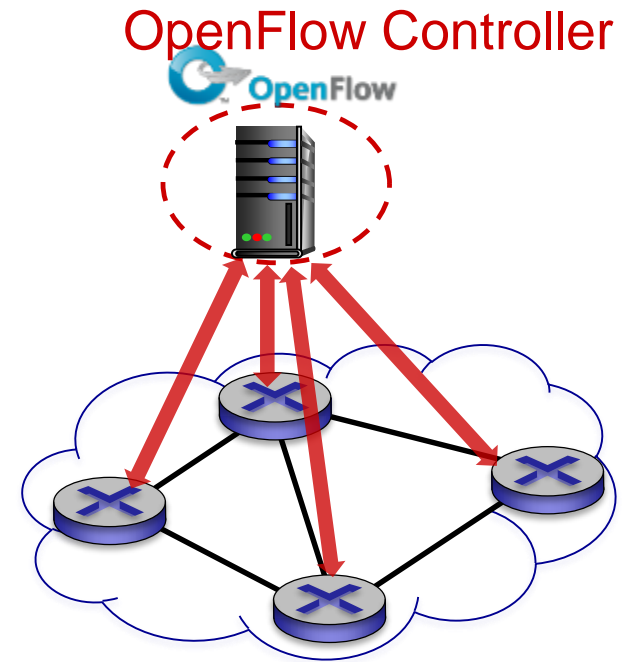*Key controller-to-switch messages*

- *configure:* controller queries/sets switch configuration parameters

- *modify-state:* add, delete, modify flow entries in the OpenFlow tables

- *Read-state:* collect statistics and counter values from the switch's flow table and ports

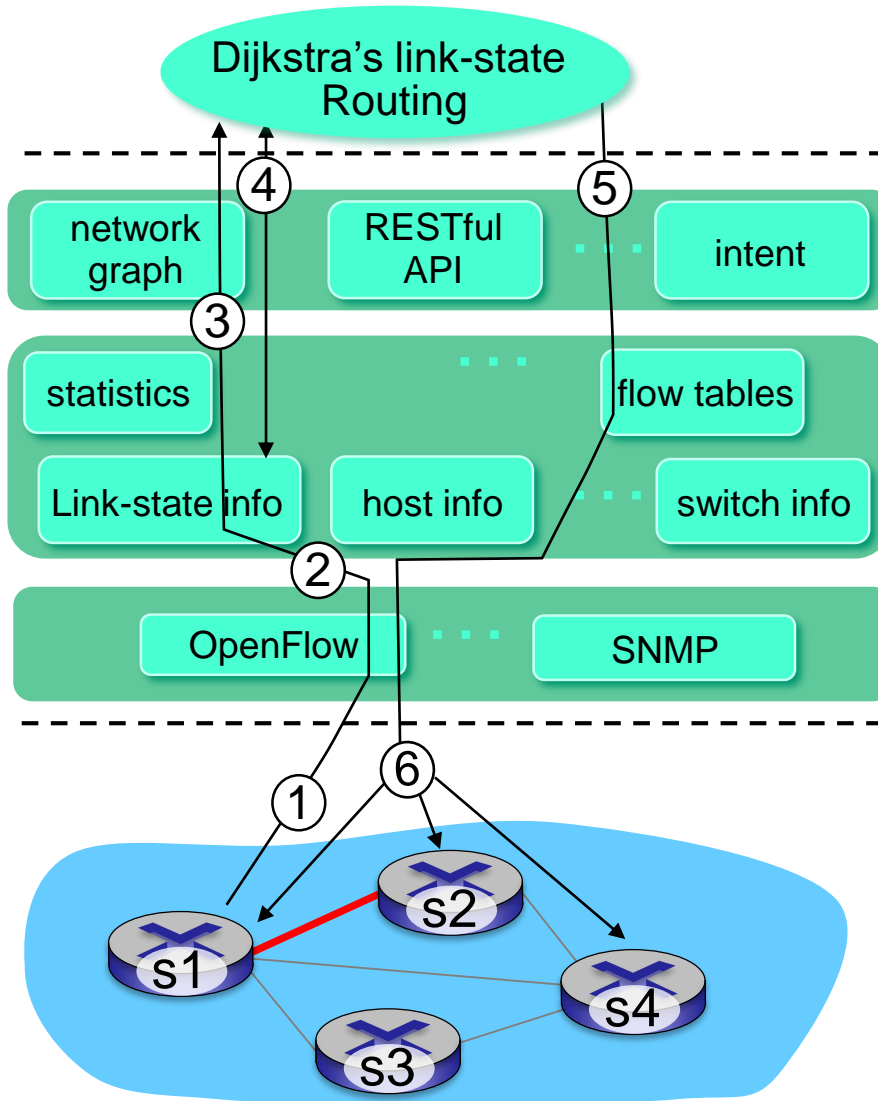- *packet-out:* controller can send this packet out of specific switch port

OpenFlow Controller

# OpenFlow: switch-to-controller messages

*Key switch-to-controller messages*

- *packet-in:* transfer packet (and its control) to controller. See packet-out message from controller

- *flow-removed:* flow table entry deleted at switch

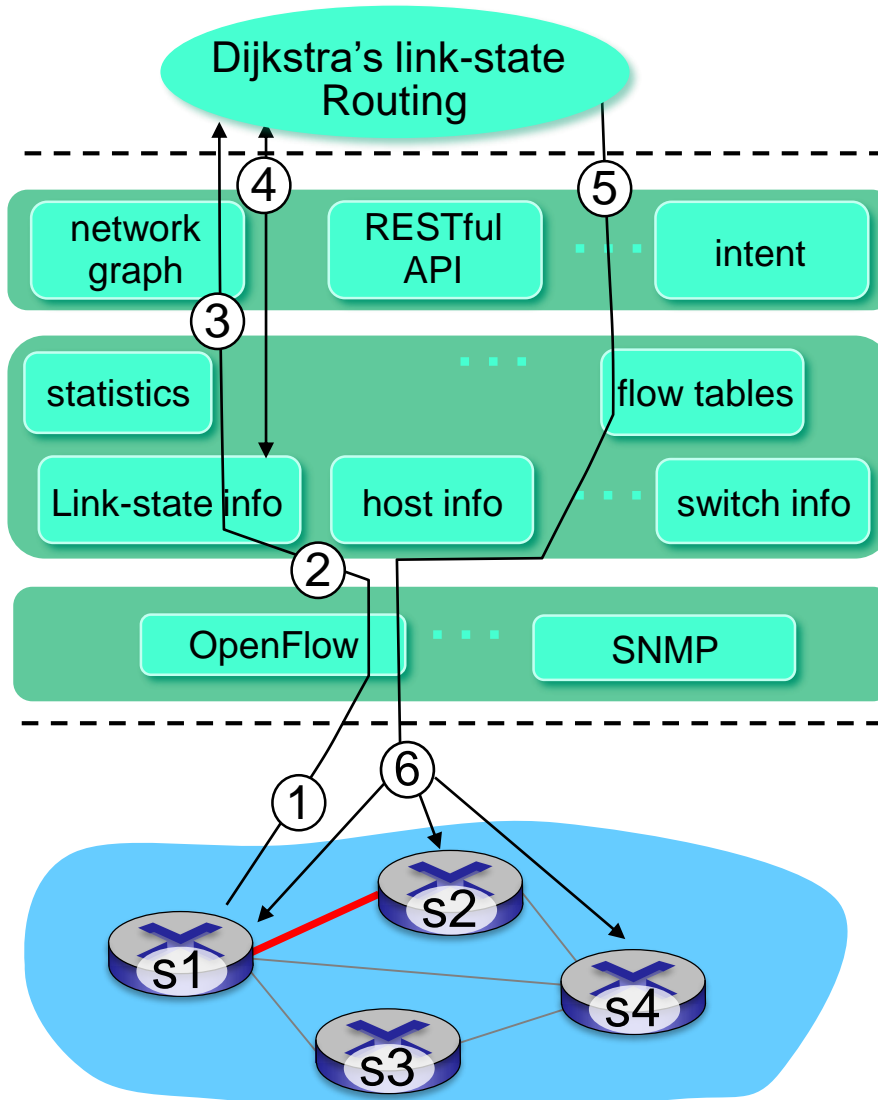- *port status:* inform controller of a change on a port.

OpenFlow Controller

# SDN: control/data plane interaction example



① S1, experiencing link failure using OpenFlow *port-status* message to notify controller

② SDN controller receives OpenFlow message, updates link status info

③ Dijkstra's routing algorithm application has previously registered to be called when ever link status changes. It is called.

④ Dijkstra's routing algorithm access network graph info, link state info in controller, computes new routes
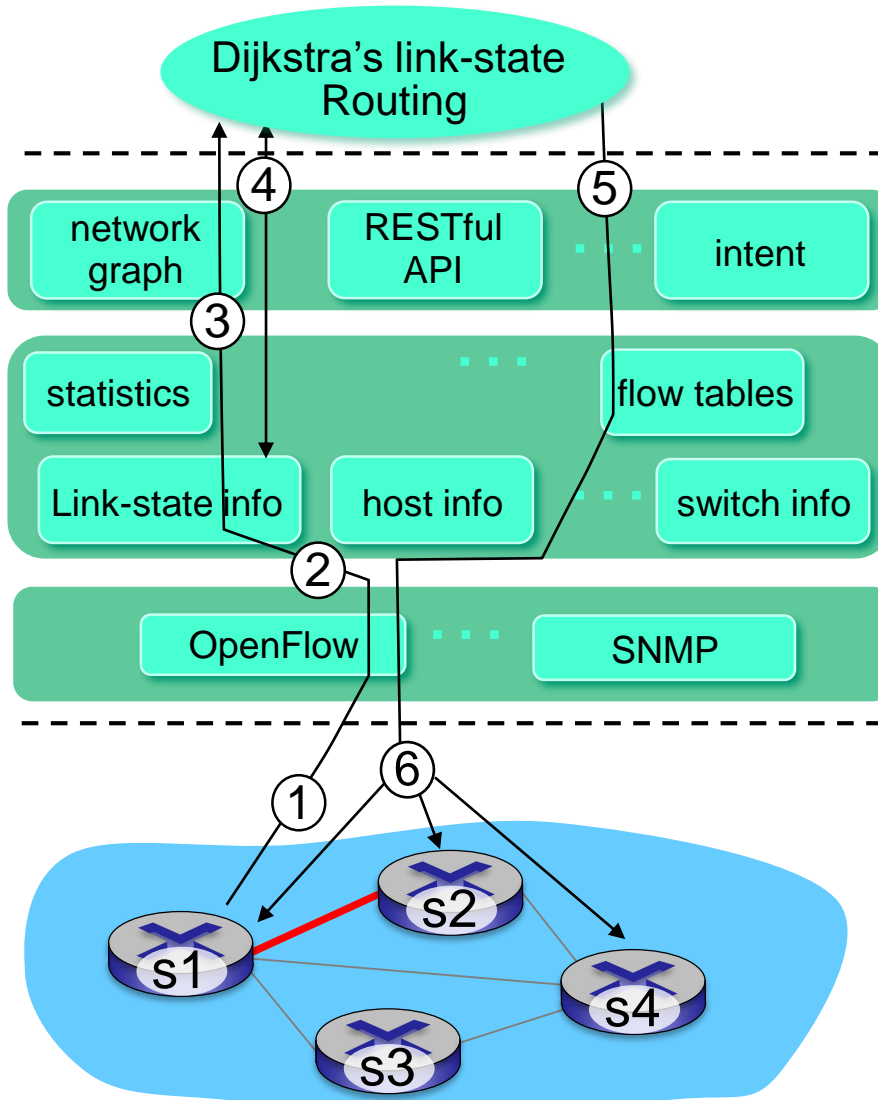
# SDN: control/data plane interaction example



Two important differences from the earlier per-router-control scenario:

- Dijkstra's algorithm is executed as a separate application, outside of the packet switches.

- Packet switches send link updates to the SDN controller and not to each other.

# SDN: control/data plane interaction example



⑤ link state routing app interacts with flow-table-computation component in SDN controller, which computes new flow tables needed

⑥ Controller uses OpenFlow to install new tables in switches that need updating

# Chapter 5: outline

# ICMP: internet control message protocol

- **used by hosts & routers to communicate network-level information**
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- **network-layer "above" IP:**
  - ICMP msgs carried in IP datagrams
- **ICMP message:**
  - Type + code + the header and the first 8 bytes of IP datagram causing error

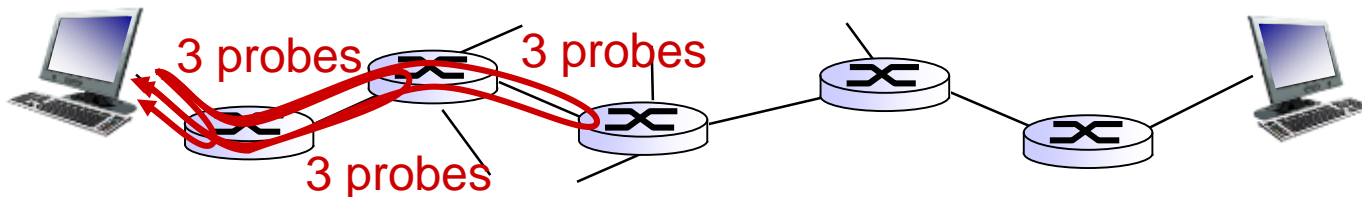| Type | Code | description |
|------|------|-------------|
| 0 | 0 | echo reply (ping) |
| 3 | 0 | dest. network unreachable |
| 3 | 1 | dest host unreachable |
| 3 | 2 | dest protocol unreachable |
| 3 | 3 | dest port unreachable |
| 3 | 6 | dest network unknown |
| 3 | 7 | dest host unknown |
| 4 | 0 | source quench (congestion control - not used) |
| 8 | 0 | echo request (ping) |
| 9 | 0 | route advertisement |
| 10 | 0 | router discovery |
| 11 | 0 | TTL expired |
| 12 | 0 | bad IP header |

# Traceroute and ICMP

- source sends series of UDP segments to destination
  - first set has TTL =1
  - second set has TTL=2, etc.
  - unlikely port number
- when datagram in $n$th set arrives to nth router:
  - router discards datagram and sends source ICMP message (type 11, code 0)
  - ICMP message include name of router & IP address

- when ICMP message arrives, source records RTTs

*stopping criteria:*

- UDP segment eventually arrives at destination host
- destination returns ICMP "port unreachable" message (type 3, code 3)
- source stops



3 probes   3 probes

3 probes

# Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state

- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol
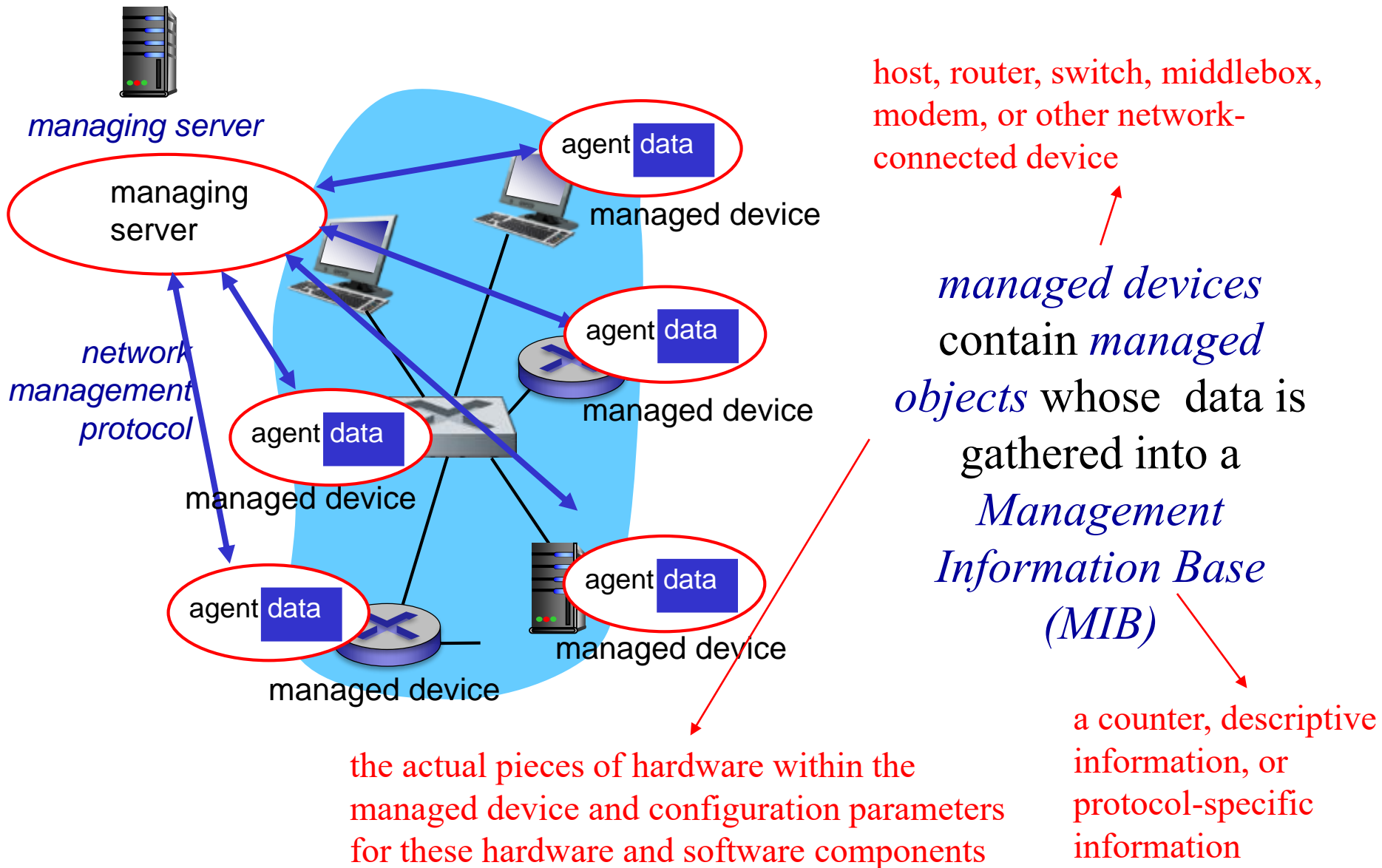
5.7 Network management and SNMP

# What is network management?

- **autonomous systems (aka "network"):** 1000s of interacting hardware/software components
- other complex systems requiring monitoring, control:
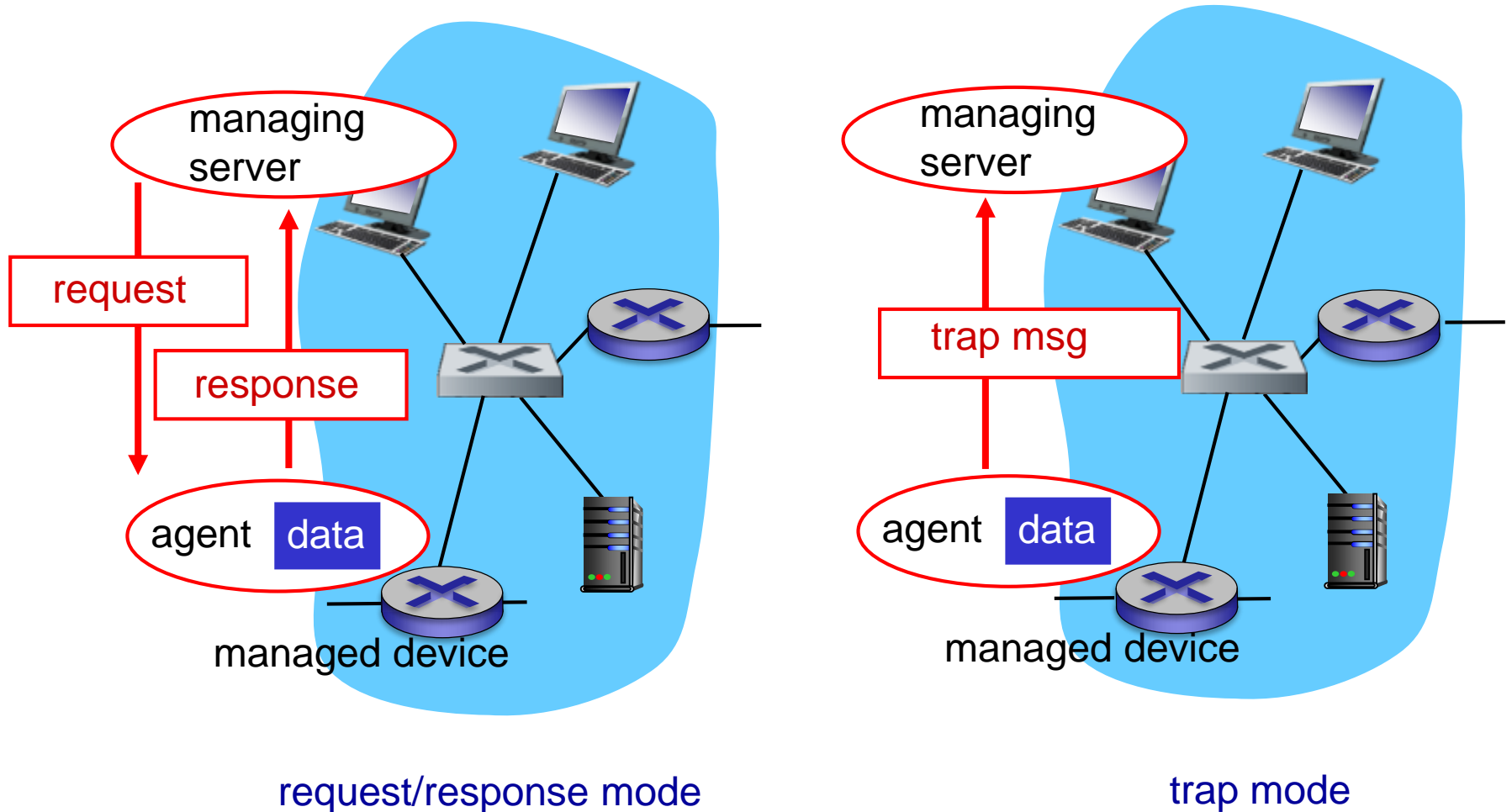  - jet airplane
  - nuclear power plant
  - others?

"Network management includes the deployment, integration and coordination of the hardware, software, and human elements to monitor, test, poll, configure, analyze, evaluate, and control the network and element resources to meet the real-time, operational performance, and Quality of Service requirements at a reasonable cost."

# Infrastructure for network management



*managing server*

managing server

*network management protocol*

agent data

managed device

agent data

managed device

agent data

managed device

agent data

managed device

agent data

managed device

host, router, switch, middlebox, modem, or other network-connected device

*managed devices* contain *managed objects* whose data is gathered into a *Management Information Base (MIB)*

the actual pieces of hardware within the managed device and configuration parameters for these hardware and software components

a counter, descriptive information, or protocol-specific information

# Simple Network Management Protocol (SNMP)

Two usages of SNMP



request/response mode

trap mode

# SNMP protocol: message types

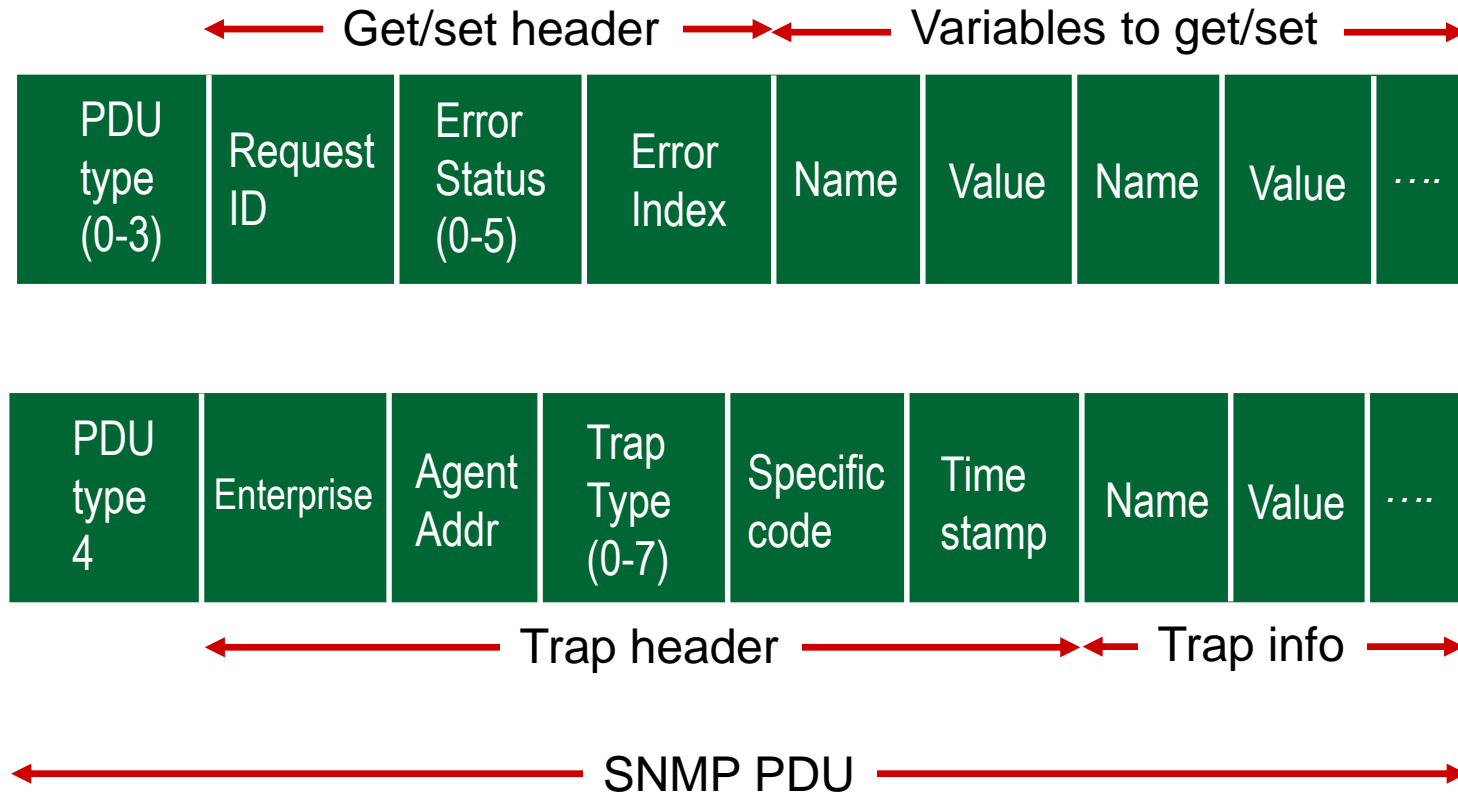| Message type | Function |
|---|---|
| GetRequest<br>GetNextRequest<br>GetBulkRequest | manager-to-agent: "get me data"<br>(data instance, next data in list, block of data) |
| InformRequest | manager-to-manager: here's MIB value |
| SetRequest | manager-to-agent: set MIB value |
| Response | Agent-to-manager: value, response to Request |
| Trap | Agent-to-manager: inform manager of exceptional event |

# SNMP protocol: message formats

Get/set header ⟷ Variables to get/set

| PDU type (0-3) | Request ID | Error Status (0-5) | Error Index | Name | Value | Name | Value | .... |
|---|---|---|---|---|---|---|---|---|

| PDU type 4 | Enterprise | Agent Addr | Trap Type (0-7) | Specific code | Time stamp | Name | Value | .... |
|---|---|---|---|---|---|---|---|---|

Trap header ⟷ Trap info

SNMP PDU

*More on network management:* see earlier editions of text!

# Chapter 5: summary

*we've learned a lot!*

- approaches to network control plane
  - per-router control (traditional)
  - logically centralized control (software defined networking)
- traditional routing algorithms
  - implementation in Internet: OSPF, BGP
- SDN controllers
  - implementation in practice: ODL, ONOS
- Internet Control Message Protocol
- network management

*next stop:  link layer!*

# CS 305: Computer Networks
## Fall 2022

**Link Layer**

**Ming Tang**

Department of Computer Science and Engineering
Southern University of Science and Technology (SUSTech)

# Chapter 6: Link layer and LANs

*our goals:*

- understand principles behind link layer services:
  - error detection, correction
  - sharing a broadcast channel: multiple access
  - link layer addressing
  - local area networks: Ethernet, VLANs
- instantiation, implementation of various link layer technologies

# Link layer, LANs: outline

# Link layer: introduction

*terminology:*

- hosts and routers: nodes
- communication channels that connect adjacent nodes along communication path: links
  - wired links
  - wireless links
- layer-2 packet: frame, encapsulates datagram

*link layer* has responsibility of transferring datagram from one node to *physically adjacent* node over a link

# Link layer: introduction

# Link layer: context

- datagram transferred by different link protocols over different links:
  - e.g., Ethernet on first link, PPP on intermediate links, 802.11 on last link
- each link protocol provides different services
  - e.g., may or may not provide rdt over link

*transportation analogy:*

- trip from SUSTech to Tsinghua
  - metro: SUSTech to SZ North
  - High speed train: SZ North to Beijing West
  - taxi: Beijing West to Tsinghua
- tourist = datagram
- transport segment = communication link
- transportation mode = link layer protocol
- travel agent = routing algorithm

# Link layer services

- *framing, link access:*
  - encapsulate datagram into frame, adding header, trailer
  - channel access if shared medium
  - "MAC" addresses used in frame headers to identify source, destination
    - different from IP address!
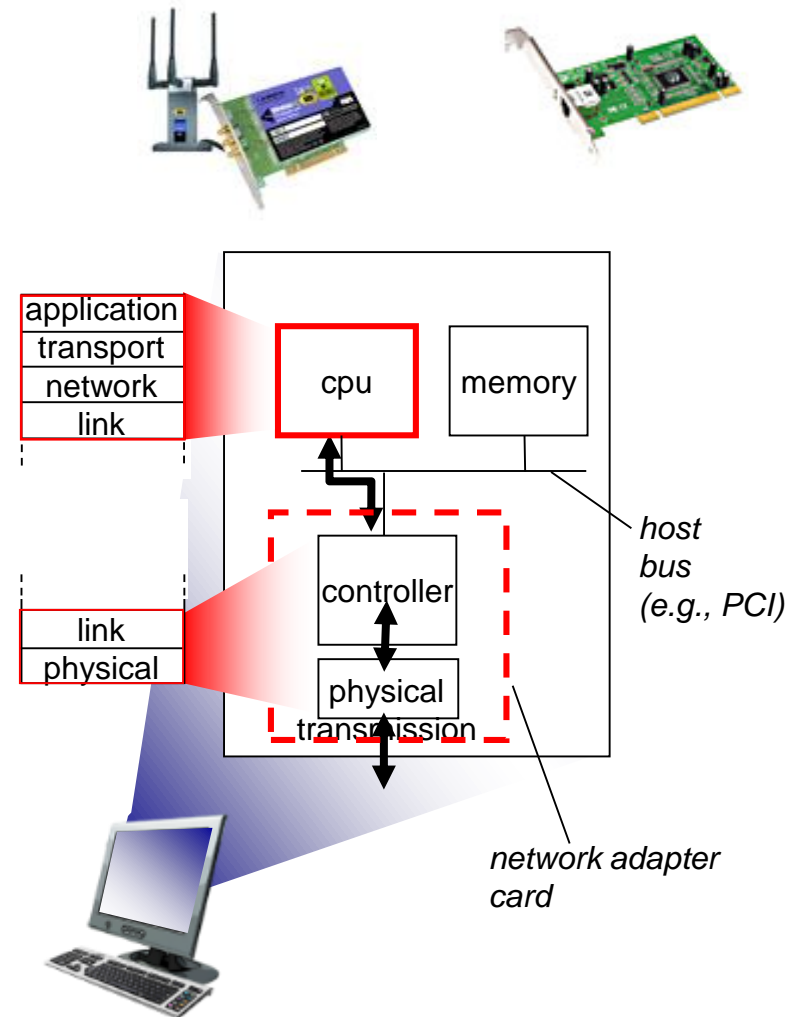- *reliable delivery between adjacent nodes*
  - we learned how to do this already (chapter 3)!
  - seldom used on low bit-error link (fiber, some twisted pair)
  - wireless links: high error rates
    - *Q:* why both link-level and end-end reliability?
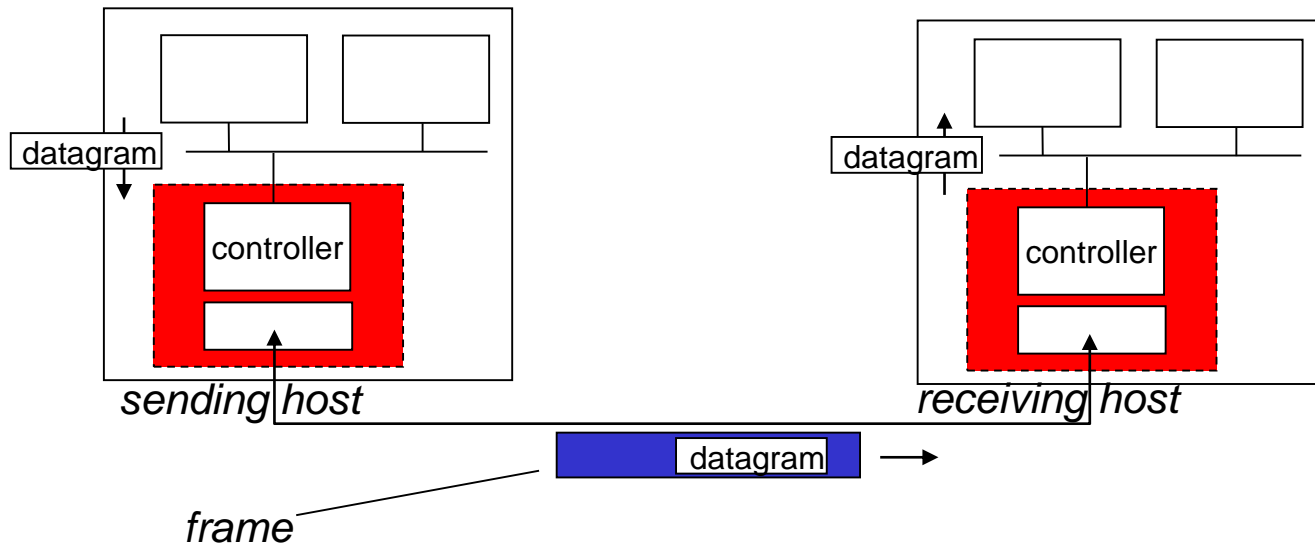
# Link layer services (more)

- *flow control:*
  - pacing between adjacent sending and receiving nodes
- *error detection*:
  - errors caused by signal attenuation, noise.
  - receiver detects presence of errors:
    - signals sender for retransmission or drops frame
- *error correction:*
  - receiver identifies *and corrects* bit error(s) without resorting to retransmission
- *half-duplex and full-duplex*
  - with half duplex, nodes at both ends of link can transmit, but not at same time

# Where is the link layer implemented?

- in each and every host
- link layer implemented in "adaptor" (aka *network interface card* NIC) or on a chip
  - Ethernet card, 802.11 card; Ethernet chipset
  - implements link, physical layer
- attaches into host's system buses
- combination of hardware, software, firmware



application
transport
network
link

cpu

memory

link
physical

controller

physical transmission

host bus (e.g., PCI)

network adapter card

# Adaptors communicating



sending host

receiving host

datagram

frame

- sending side:
  - encapsulates datagram in frame
  - adds error checking bits, rdt, flow control, etc.

- receiving side
  - looks for errors, rdt, flow control, etc.
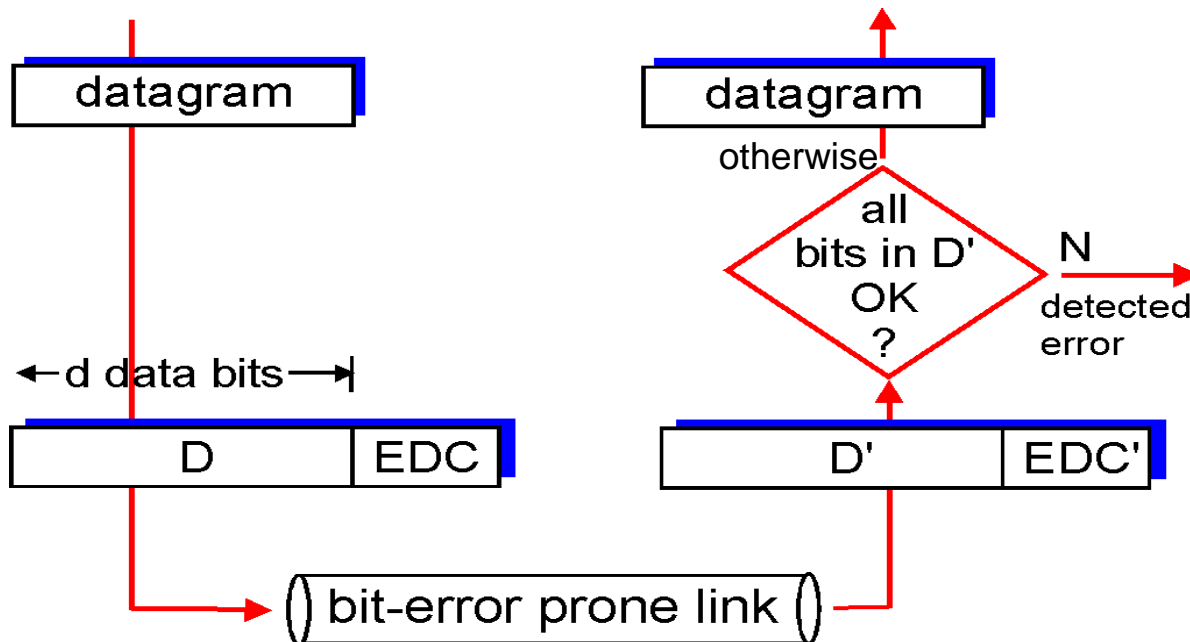  - extracts datagram, passes to upper layer at receiving side

# Link layer, LANs: outline

# Error detection

EDC= Error Detection and Correction bits
D    = Data protected by error checking, may include header fields

- Error detection not 100% reliable!
    - protocol may miss some errors, but rarely
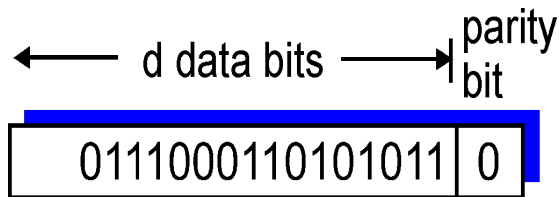    - larger EDC field yields better detection and correction, but larger overhead



- Parity checks
- Check-summing methods
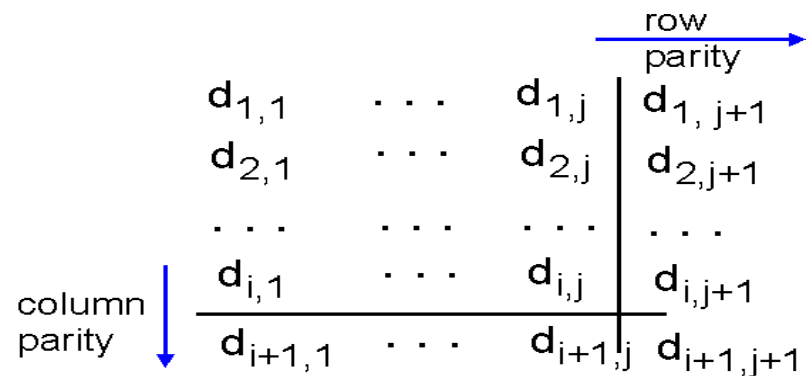- Cyclic-redundancy check

# Parity checking

## *single bit parity:*

- *d*etect single bit errors
- Even parity scheme
- Odd parity scheme



## *two-dimensional bit parity:*

- detect and correct single bit errors



no errors

parity error

correctable single bit error

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

# Parity checking

```
1 0 1 1 1 | 0
1 0 1 0 1 | 1
1 1 1 0 0 | 1
1 1 1 1 0 | 0
```

1 0 1 1 1 0 1 0 1 0 1 1 1 1 1 0 0 1 1 1 1 1 0 0

Parity Bits

**Case 1: a bit is in error.**

```
1 0 1 1 1   0
1 0 0 0 1   1       Error Detected
1 1 1 0 0   1
1 1 1 1 0   0
```

**Case 2: two bits are in error.**

Correct Bit Detect As Incorrect Bit

```
0 0 1 1 1   0
1 0 1 0 1   1       Error Detected
1 1 1 0 1   1
1 1 1 1 0   0
```

**Case 3: error not detected**

```
1 0 1 1 1   0
1 0 0 1 1   1       Not Detected so
1 1 0 1 0   1       not Corrected
1 1 1 1 0   0
```

**Many other cases …**

# Internet checksum (review)

**goal:** detect "errors" (e.g., flipped bits) in transmitted packet (note: used at transport layer only)

*sender:*
- treat segment contents as sequence of 16-bit integers
- checksum: addition (1's complement sum) of segment contents
- sender puts checksum value into UDP checksum field

*receiver:*
- compute checksum of received segment
- check if computed checksum equals checksum field value:
  - NO - error detected
  - YES - no error detected. *But maybe errors nonetheless?*

# Cyclic redundancy check

- more powerful error-detection coding
- view data bits, D, as a binary number
- choose r+1 bit pattern (generator), G
- goal: choose r CRC bits, R, such that
  - <D,R> exactly divisible by G (modulo 2)
  - receiver knows G, divides <D,R> by G. If non-zero remainder: error detected!
  - can detect all consecutive bit errors of r bits or less
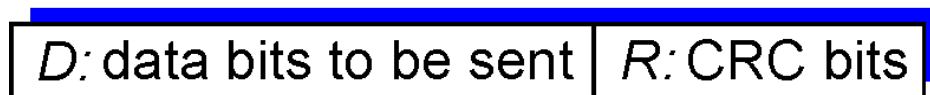- widely used in practice (Ethernet, 802.11 WiFi, ATM)

```
1011 XOR 0101 = 1110
1001 XOR 1101 = 0100

1011 - 0101 = 1110
1001 - 1101 = 0100
```



d bits ──────────→ ← r bits →

| D: data bits to be sent | R: CRC bits |

bit pattern

$$D * 2^r \quad XOR \quad R$$

mathematical formula

# Cyclic redundancy check

All CRC calculations are done in modulo-2 arithmetic without carries in addition or borrows in subtraction.

- This means that addition and subtraction are identical, and
- both are equivalent to the bitwise exclusive-or (XOR) of the operands.

```
1011 XOR 0101 = 1110          1011 - 0101 = 1110
1001 XOR 1101 = 0100          1001 - 1101 = 0100
```

Multiplication and division are the same as in base-2 arithmetic, except that any required addition or subtraction is done without carries or borrows.

```
              10001 remainder 101
      10011|100100110
            10011
            10110
            10011
              101
```
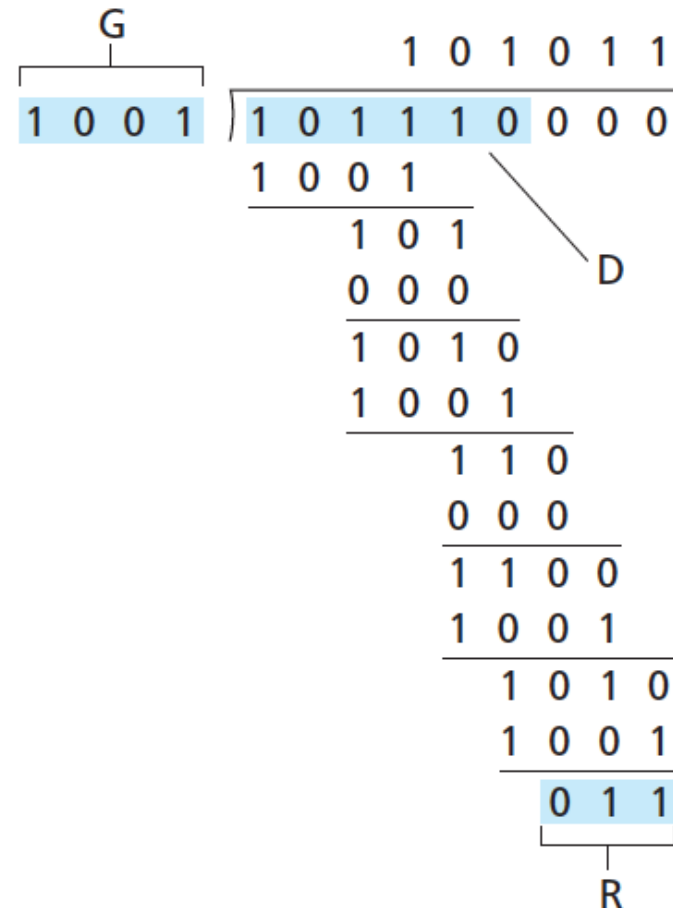
# CRC example

want:

   $D \cdot 2^r$ XOR $R = nG$

*equivalently:*

   $D \cdot 2^r = nG$ XOR $R$

*equivalently:*

if we divide $D \cdot 2^r$ by G, want remainder R to satisfy:

$$R = remainder[\frac{D \cdot 2^r}{G}]$$

G

```
                              1 0 1 0 1 1
      1 0 0 1  | 1 0 1 1 1 0 0 0 0
                 1 0 0 1
                   1 0 1
                   0 0 0
                   1 0 1 0
                   1 0 0 1
                     1 1 0
                     0 0 0
                     1 1 0 0
                     1 0 0 1
                       1 0 1 0
                       1 0 0 1
                         0 1 1
```
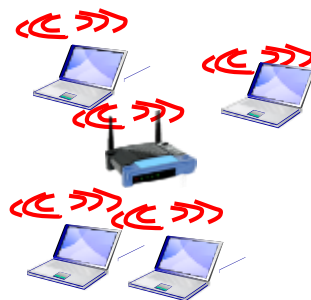
D

R

# Link layer, LANs: outline

# Multiple access links, protocols
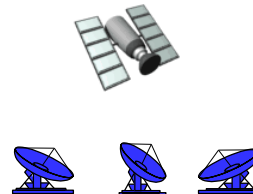
two types of "links":

- point-to-point
  - PPP for dial-up access
  - point-to-point link between Ethernet switch, host
- *broadcast (shared wire or medium)*
  - old-fashioned Ethernet
  - 802.11 wireless LAN



shared wire (e.g., cabled Ethernet)

shared RF (e.g., 802.11 WiFi)

shared RF (satellite)

humans at a cocktail party (shared air, acoustical)

# Multiple access protocols

- single shared broadcast channel
- two or more simultaneous transmissions by nodes: interference
  - *collision* if node receives two or more signals at the same time

*multiple access protocol*

- distributed algorithm that determines how nodes share channel, i.e., determine which and when node can transmit
- communication about channel sharing must use channel itself!
  - no out-of-band channel for coordination

# An ideal multiple access protocol

*given:* broadcast channel of rate R bps

*Desired properties:*

1. when one node wants to transmit, it can send at rate R.
2. when M nodes want to transmit, each can send at average rate R/M
3. fully decentralized:

   - no special node to coordinate transmissions
   - no synchronization of clocks, slots
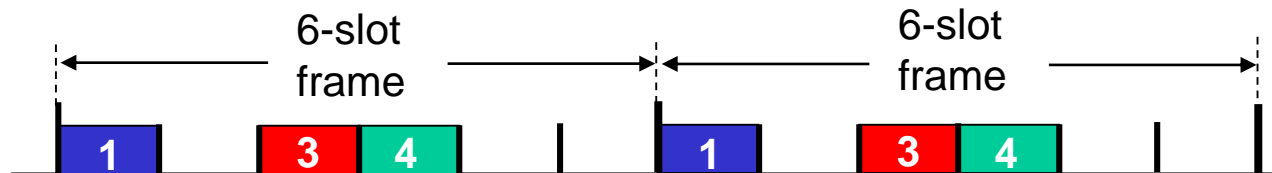4. simple

# MAC protocols: taxonomy

three broad classes:

- *channel partitioning*
  - divide channel into smaller "pieces" (time slots, frequency, code)
  - allocate piece to node for exclusive use
- *random access*
  - channel not divided, allow collisions
  - "recover" from collisions
- *"taking turns"*
  - nodes take turns, but nodes with more to send can take longer turns

# Channel partitioning MAC protocols: TDMA

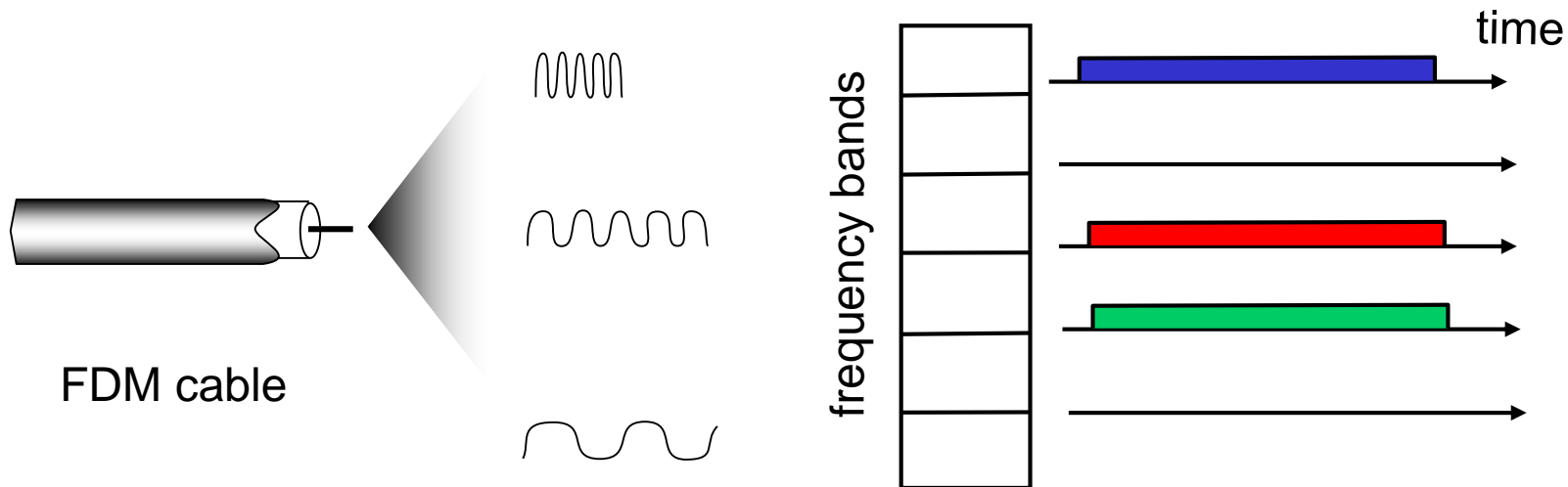## TDMA: time division multiple access

- access to channel in "rounds"
- each station gets fixed length slot (length = packet transmission time) in each round
- unused slots go idle
- example: 6-station LAN, 1,3,4 have packets to send, slots 2,5,6 idle

# Channel partitioning MAC protocols: FDMA

## FDMA: frequency division multiple access

- channel spectrum divided into frequency bands
- each station assigned fixed frequency band
- unused transmission time in frequency bands go idle
- example: 6-station LAN, 1,3,4 have packet to send, frequency bands 2,5,6 idle

FDM cable

# Random access protocols

- when node has packet to send
  - transmit at full channel data rate R.
  - no *a priori* coordination among nodes
- two or more transmitting nodes ➔ "collision",
- random access MAC protocol specifies:
  - how to detect collisions
  - how to recover from collisions (e.g., via delayed retransmissions)
- examples of random access MAC protocols:
  - slotted ALOHA
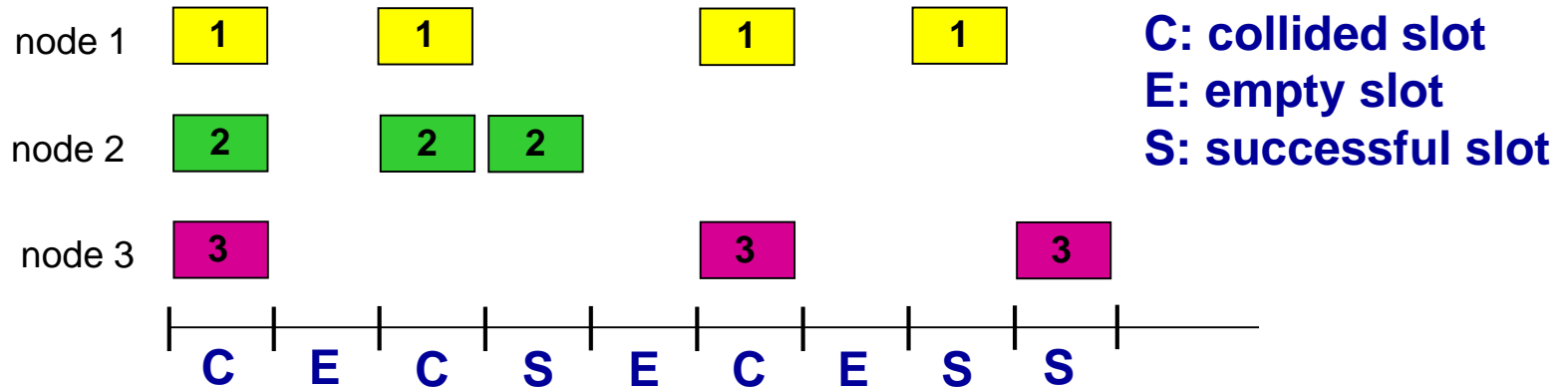  - ALOHA
  - CSMA, CSMA/CD, CSMA/CA

# Slotted ALOHA

*assumptions:*

- all frames same size
- time divided into equal size slots (time to transmit 1 frame)
- nodes start to transmit only slot beginning
- nodes are synchronized
- if 2 or more nodes transmit in slot, all nodes detect collision

*operation:*

- when node obtains fresh frame, transmits in next slot
  - *if no collision:* node can send new frame in next slot
  - *if collision:* node retransmits frame in each subsequent slot with prob. p until success

# Slotted ALOHA



node 1 | 1 | 1 | 1 | 1

node 2 | 2 | 2 | 2

node 3 | 3 | 3 | 3

**C: collided slot**
**E: empty slot**
**S: successful slot**

C  E  C  S  E  C  E  S  S

*Pros:*

- single active node can continuously transmit at full rate of channel
- highly decentralized: only slots in nodes need to be in sync
- simple

*Cons:*

- collisions, wasting slots
- idle slots
- clock synchronization

# Slotted ALOHA: efficiency

*efficiency*: long-run fraction of successful slots (many nodes, all with many frames to send)

- *suppose:* N nodes with many frames to send, each transmits in slot with probability $p$
- prob that given node has success in a slot $= p(1-p)^{N-1}$
- prob that *any* node has a success $= Np(1-p)^{N-1}$

- max efficiency: find $p*$ that maximizes $Np(1-p)^{N-1}$
- for many nodes, take limit of $Np*(1-p*)^{N-1}$ as $N$ goes to infinity, gives:
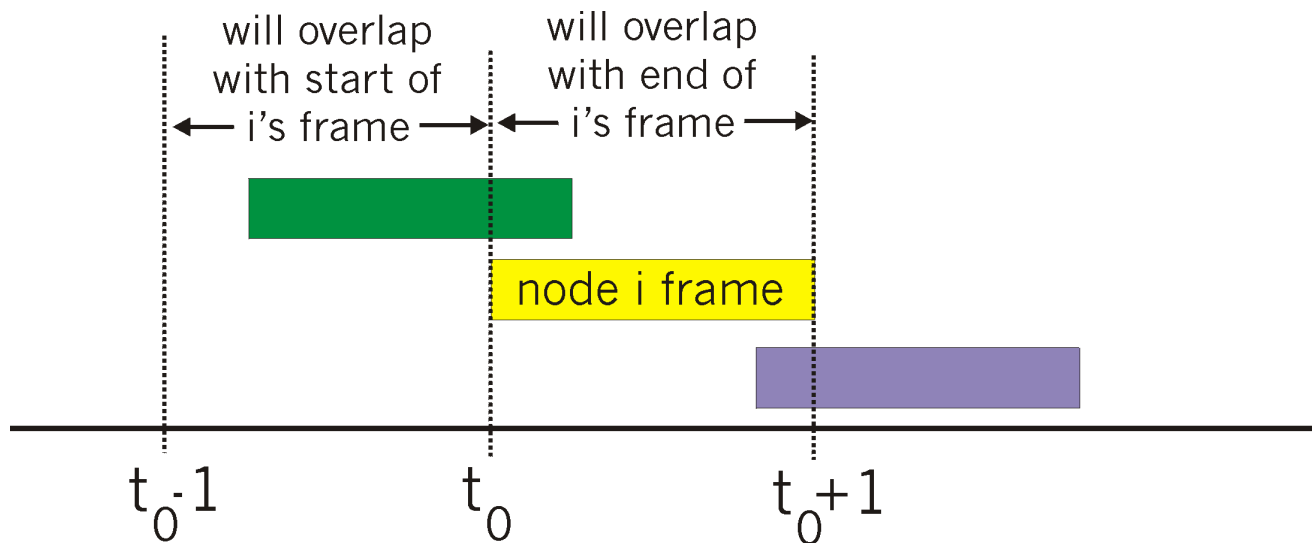
*max efficiency = 1/e = .37*

*at best:* channel used for useful transmissions 37% of time!

*!*

# Pure (unslotted) ALOHA

- unslotted Aloha: simpler, no synchronization
- when frame first arrives
  - transmit immediately
- collision probability increases:
  - frame sent at $t_0$ collides with other frames sent in $[t_0-1, t_0+1]$

will overlap with start of i's frame ← → ← will overlap with end of i's frame → 

node i frame

$t_0-1$      $t_0$      $t_0+1$

# Pure ALOHA efficiency

P(success by given node) = P(node transmits) .

$\quad$ P(no other node transmits in $[t_0-1, t_0]$ .

$\quad$ P(no other node transmits in $[t_0-1, t_0]$

$$= p \cdot (1-p)^{N-1} \cdot (1-p)^{N-1}$$

$$= p \cdot (1-p)^{2(N-1)}$$

… choosing optimum p and then letting $n$ $\longrightarrow$ $\infty$

$$= 1/(2e) = .18$$

even *worse* than slotted Aloha!