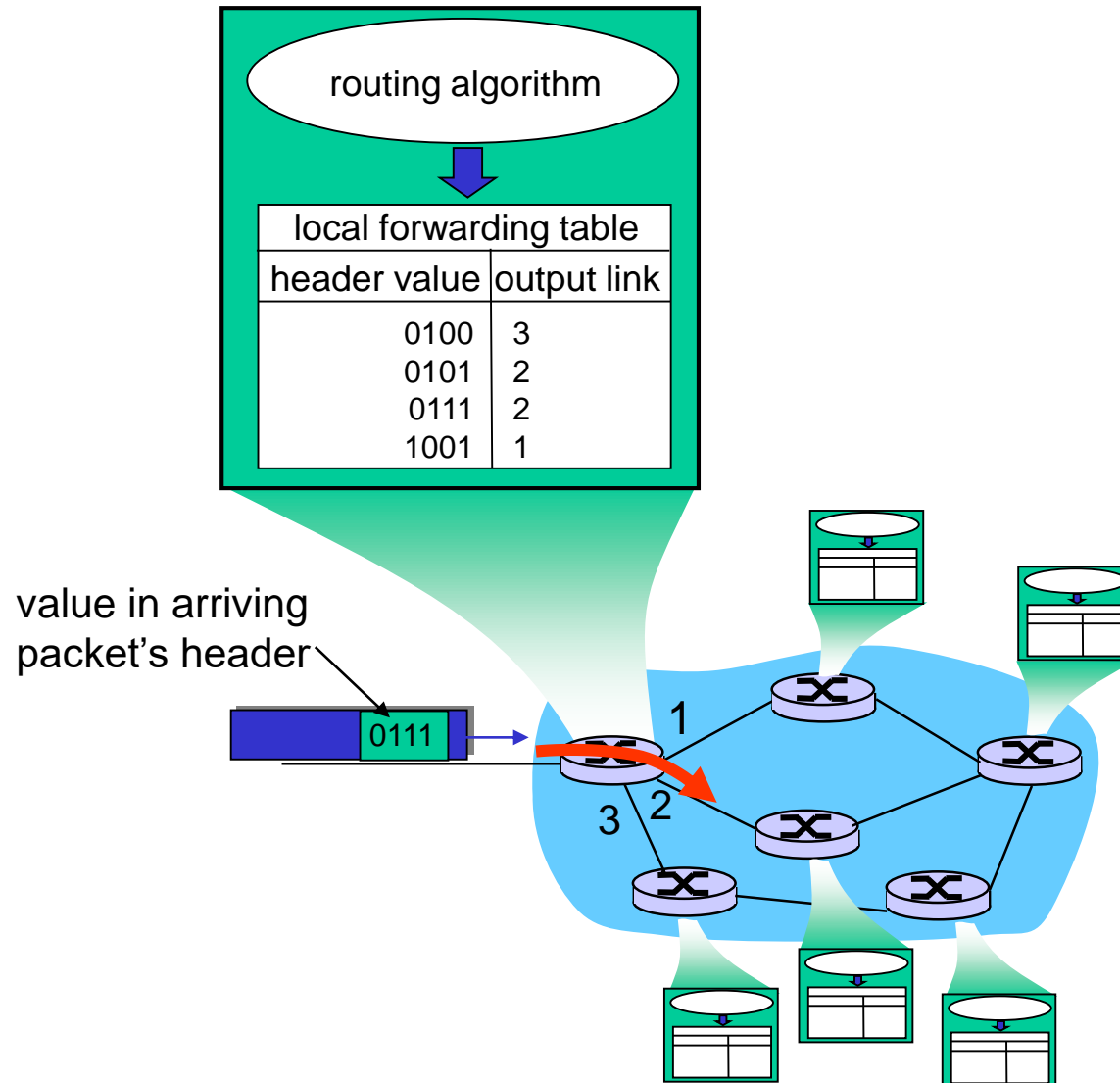


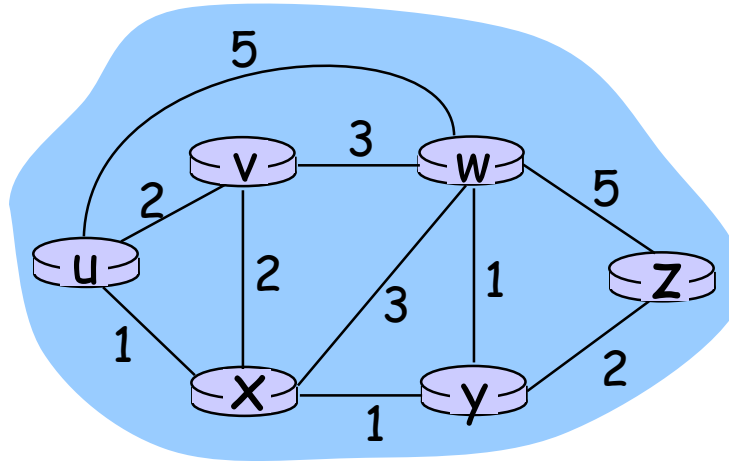
# Chapter 4: Network Layer

- ❑ 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

# Interplay between routing, forwarding



# Graph abstraction



Graph:  $G = (N, E)$

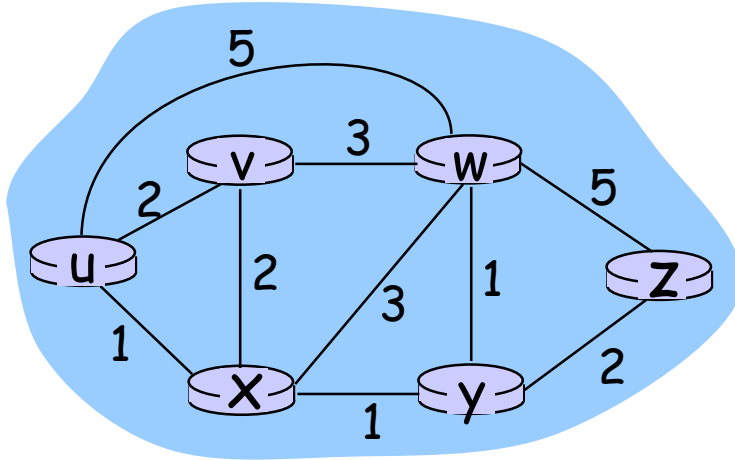
$N$  = set of routers =  $\{ u, v, w, x, y, z \}$

$E$  = set of links =  $\{ (u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) \}$

Remark: Graph abstraction is useful in other network contexts

Example: P2P, where  $N$  is set of peers and  $E$  is set of TCP connections

# Graph abstraction: costs



- $c(x,x') = \text{cost of link } (x,x')$ 
  - e.g.,  $c(w,z) = 5$
- cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

Cost of path  $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

Question: What's the least-cost path between u and z ?

Routing algorithm: algorithm that finds least-cost path

# Routing Algorithm classification

## Global or decentralized information?

### Global:

- ❑ all routers have complete topology, link cost info
- ❑ "link state" algorithms

### Decentralized:

- ❑ router knows physically-connected neighbors, link costs to neighbors
- ❑ iterative process of computation, exchange of info with neighbors
- ❑ "distance vector" algorithms

## Static or dynamic?

### Static:

- ❑ routes change slowly over time

### Dynamic:

- ❑ routes change more quickly
  - periodic update
  - in response to link cost changes

# Chapter 4: Network Layer

- ❑ 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

# A Link-State Routing Algorithm

## Dijkstra's algorithm

- ❑ net topology, link costs known to all nodes
  - accomplished via "link state broadcast"
  - all nodes have same info
- ❑ computes least cost paths from one node ('source') to all other nodes
  - gives **forwarding table** for that node
- ❑ iterative: after  $k$  iterations, know least cost path to  $k$  dest.'s

## Notation:

- ❑  $c(x,y)$ : link cost from node  $x$  to  $y$ ;  $= \infty$  if not direct neighbors
- ❑  $D(v)$ : current value of cost of path from source to dest.  $v$
- ❑  $p(v)$ : predecessor node along path from source to  $v$
- ❑  $N'$ : set of nodes whose least cost path definitively known

# Dijkstra's Algorithm

1 **Initialization:**

2  $N' = \{u\}$

3 for all nodes  $v$

4 if  $v$  adjacent to  $u$

5 then  $D(v) = c(u,v)$

6 else  $D(v) = \infty$

7

8 **Loop**

9 find  $w$  not in  $N'$  such that  $D(w)$  is a minimum

10 add  $w$  to  $N'$

11 update  $D(v)$  for all  $v$  adjacent to  $w$  and not in  $N'$  :

12  $D(v) = \min( D(v), D(w) + c(w,v) )$

13 /\* new cost to  $v$  is either old cost to  $v$  or known

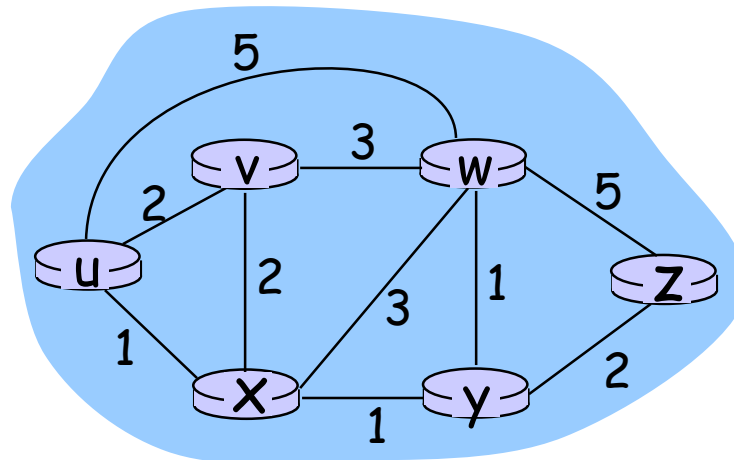
14 shortest path cost to  $w$  plus cost from  $w$  to  $v$  \*/

15 **until all nodes in  $N'$**



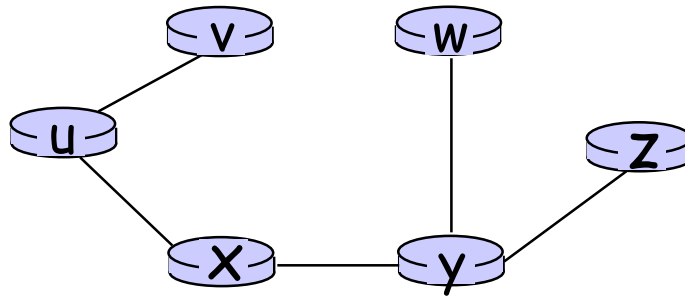
# Dijkstra's algorithm: example

Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	$\infty$	$\infty$
1	ux	2,u	4,x		2,x	$\infty$
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					



# Dijkstra's algorithm: example (2)

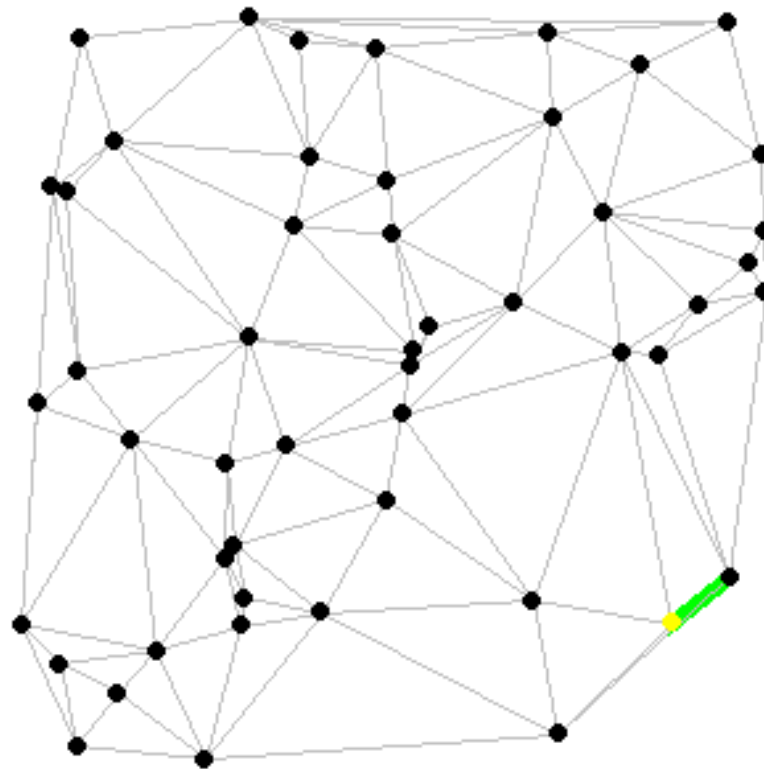
Resulting shortest-path tree from u:



Resulting forwarding table in u:

destination	link
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

Dijkstra's algorithm



[www.combinatorica.com](http://www.combinatorica.com)

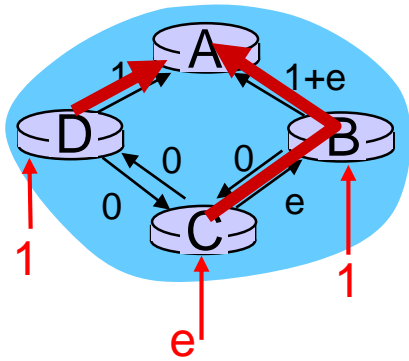
# Dijkstra's algorithm, discussion

**Algorithm complexity:**  $n$  nodes

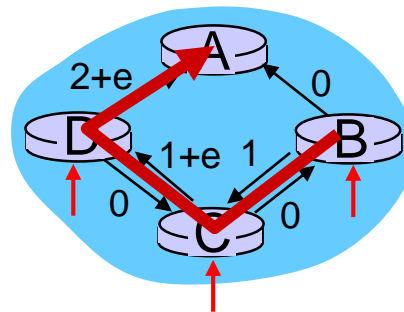
- each iteration: need to check all nodes,  $w$ , not in  $N$
- $n(n+1)/2$  comparisons:  $O(n^2)$
- more efficient implementations possible:  $O(n \log n)$

**Oscillations possible:**

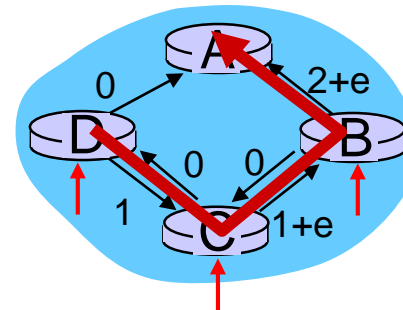
- e.g., link cost = amount of carried traffic



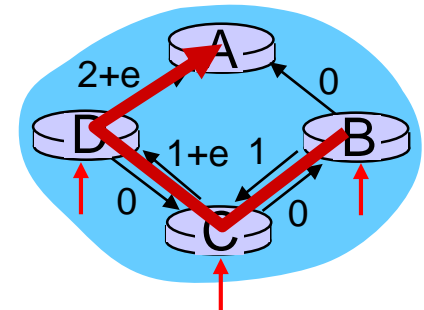
initially



given these costs,  
find new routing....  
resulting in new costs



given these costs,  
find new routing....  
resulting in new costs



given these costs,  
find new routing....  
resulting in new costs

# Chapter 4: Network Layer

- ❑ 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

# Distance Vector Algorithm

## Bellman-Ford Equation (dynamic programming)

Define

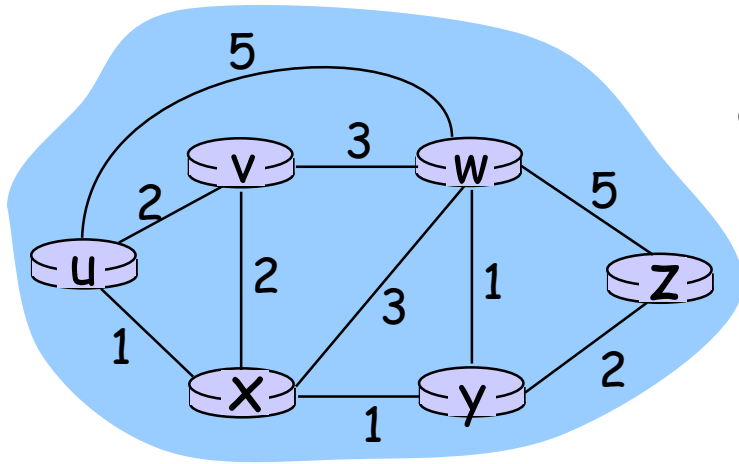
$d_x(y) :=$  cost of least-cost path from  $x$  to  $y$

Then

$$d_x(y) = \min_v \{c(x,v) + d_v(y)\}$$

where min is taken over all neighbors  $v$  of  $x$

# Bellman-Ford example



Clearly,  $d_v(z) = 5$ ,  $d_x(z) = 3$ ,  $d_w(z) = 3$

B-F equation says:

$$\begin{aligned} d_u(z) &= \min \{ c(u,v) + d_v(z), \\ &\quad c(u,x) + d_x(z), \\ &\quad c(u,w) + d_w(z) \} \\ &= \min \{ 2 + 5, \\ &\quad 1 + 3, \\ &\quad 5 + 3 \} = 4 \end{aligned}$$

Node that achieves minimum is next  
hop in shortest path → forwarding table

# Distance Vector Algorithm

- $D_x(y)$  = estimate of least cost from  $x$  to  $y$
- Node  $x$  knows cost to each neighbor  $v$ :  
 $c(x,v)$
- Node  $x$  maintains distance vector  $D_x = [D_x(y): y \in N]$
- Node  $x$  also maintains its neighbors' distance vectors
  - For each neighbor  $v$ ,  $x$  maintains  $D_v = [D_v(y): y \in N]$



# Distance vector algorithm (4)

## Basic idea:

- ❑ From time-to-time, each node sends its own distance vector estimate to neighbors
- ❑ Asynchronous
- ❑ When a node  $x$  receives new DV estimate from neighbor, it updates its own DV using B-F equation:

$$D_x(y) \leftarrow \min_v \{c(x,v) + D_v(y)\} \quad \text{for each node } y \in N$$

- ❑ Under minor, natural conditions, the estimate  $D_x(y)$  converge to the actual least cost  $d_x(y)$

# Distance Vector Algorithm (5)

## Iterative, asynchronous:

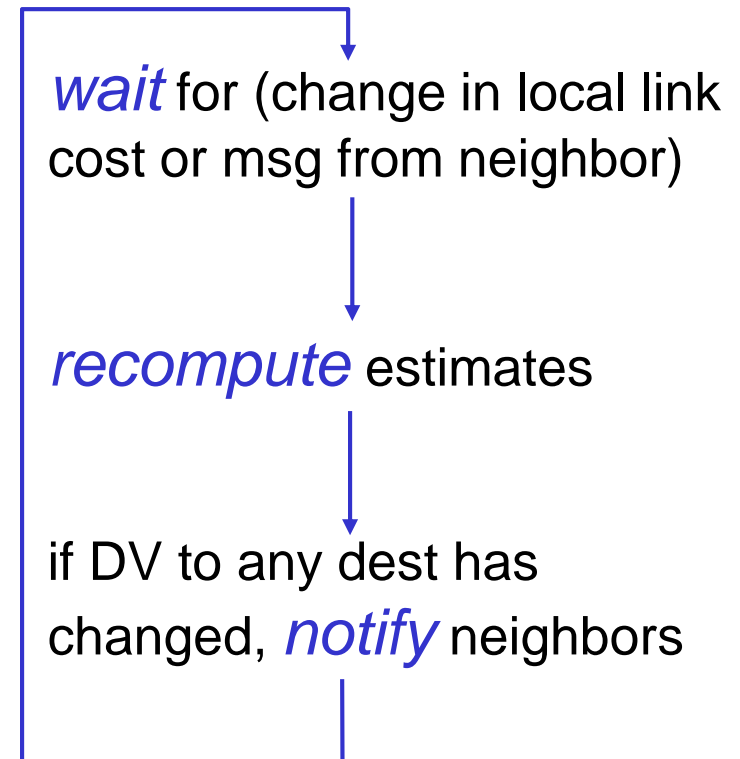
each local iteration caused by:

- ❑ local link cost change
- ❑ DV update message from neighbor

## Distributed:

- ❑ each node notifies neighbors *only* when its DV changes
  - neighbors then notify their neighbors if necessary

## Each node:



$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\} \\ = \min\{2+0, 7+1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\} \\ = \min\{2+1, 7+0\} = 3$$

**node x table**

		cost to		
from		x	y	z
	x	0	2	7
	y	∞	∞	∞
	z	∞	∞	∞

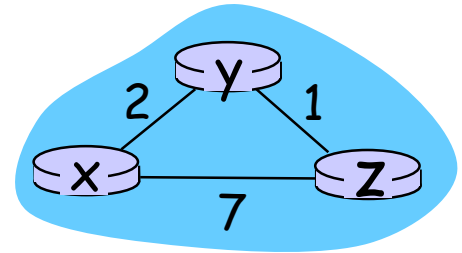
**node y table**

		cost to		
from		x	y	z
	x	∞	∞	∞
	y	2	0	1
	z	∞	∞	∞

**node z table**

		cost to		
from		x	y	z
	x	∞	∞	∞
	y	∞	∞	∞
	z	7	1	0

		cost to		
from		x	y	z
	x	0	2	3
	y	2	0	1
	z	7	1	0



time

$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\} \\ = \min\{2+0, 7+1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\} \\ = \min\{2+1, 7+0\} = 3$$

### node x table

		cost to		
		x	y	z
from	x	0	2	7
	y	$\infty$	$\infty$	$\infty$
	z	$\infty$	$\infty$	$\infty$

### node y table

		cost to		
		x	y	z
from	x	$\infty$	$\infty$	$\infty$
	y	2	0	1
	z	$\infty$	$\infty$	$\infty$

### node z table

		cost to		
		x	y	z
from	x	∞	∞	∞
	y	∞	∞	∞
	z	7	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	7	1	0

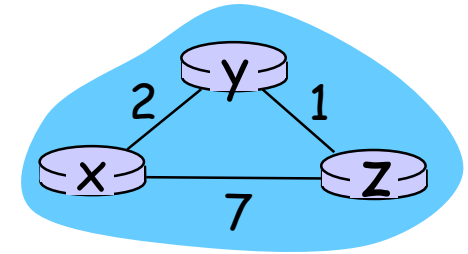
		cost to		
from		x	y	z
	x	0	2	7
	y	2	0	1
	z	7	1	0

		cost to		
		x	y	z
from	x	0	2	7
	y	2	0	1
	z	3	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	3	1	0

		cost to		
from		x	y	z
	x	0	2	3
	y	2	0	1
	z	3	1	0

		cost to		
		x	y	z
from	x	0	2	3
	y	2	0	1
	z	3	1	0

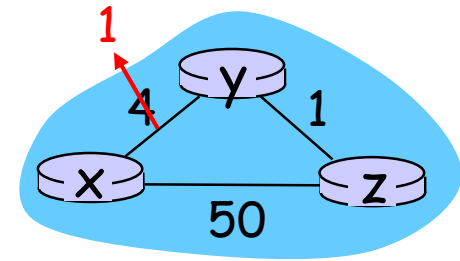


time

# Distance Vector: link cost changes

## Link cost changes:

- ❑ node detects local link cost change
- ❑ updates routing info, recalculates distance vector
- ❑ if DV changes, notify neighbors



“good  
news  
travels  
fast”

At time  $t_0$ , y detects the link-cost change, updates its DV, and informs its neighbors.

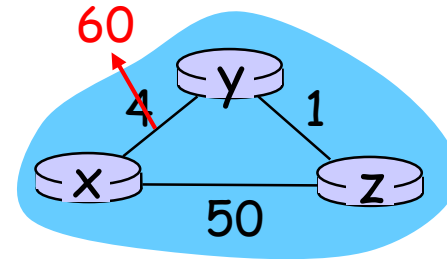
At time  $t_1$ , z receives the update from y and updates its table. It computes a new least cost to x and sends its neighbors its DV.

At time  $t_2$ , y receives z's update and updates its distance table. y's least costs do not change and hence y does not send any message to z.

# Distance Vector: link cost changes

## Link cost changes:

- ❑ good news travels fast
- ❑ bad news travels slow - "count to infinity" problem!
- ❑ 44 iterations before algorithm stabilizes: see text



## Poisoned reverse:

- ❑ If Z routes through Y to get to X :
  - Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- ❑ will this completely solve count to infinity problem?

# Comparison of LS and DV algorithms

## Message complexity

- ❑ LS: with  $n$  nodes,  $E$  links,  $O(nE)$  msgs sent
- ❑ DV: exchange between neighbors only
  - convergence time varies

## Speed of Convergence

- ❑ LS:  $O(n^2)$  algorithm requires  $O(nE)$  msgs
  - may have oscillations
- ❑ DV: convergence time varies
  - may be routing loops
  - count-to-infinity problem

**Robustness:** what happens if router malfunctions?

## LS:

- node can advertise incorrect *link* cost
- each node computes only its own table

## DV:

- DV node can advertise incorrect *path* cost
- each node's table used by others
  - error propagate thru network

# DV与LS的对比

## □ 三句话总结

- DV以最短路径算法为核心，动态更新路由信息，更新信息来自局部（邻居之间），存在“计数到无穷”问题
- LS以最短路径算法为核心，动态更新路由信息，更新信息来自全局（整个网络），存在路由震荡问题
- 蓝色是共性，红色是特性



# 思考题



- 如何解决路由震荡的问题？
- 毒性逆转能否解决“计数到无穷”的问题？
  - 可以提供几个简单的例子来说明
  - 现实协议是如何解决的？

# Chapter 4: Network Layer

- ❑ 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

# Hierarchical Routing

Our routing study thus far - idealization

- ❑ all routers identical
- ❑ network “flat”

... *not* true in practice

**scale:** with 200 million destinations:

- ❑ can't store all dest's in routing tables!
- ❑ routing table exchange would swamp links!

**administrative autonomy**

- ❑ internet = network of networks
- ❑ each network admin may want to control routing in its own network

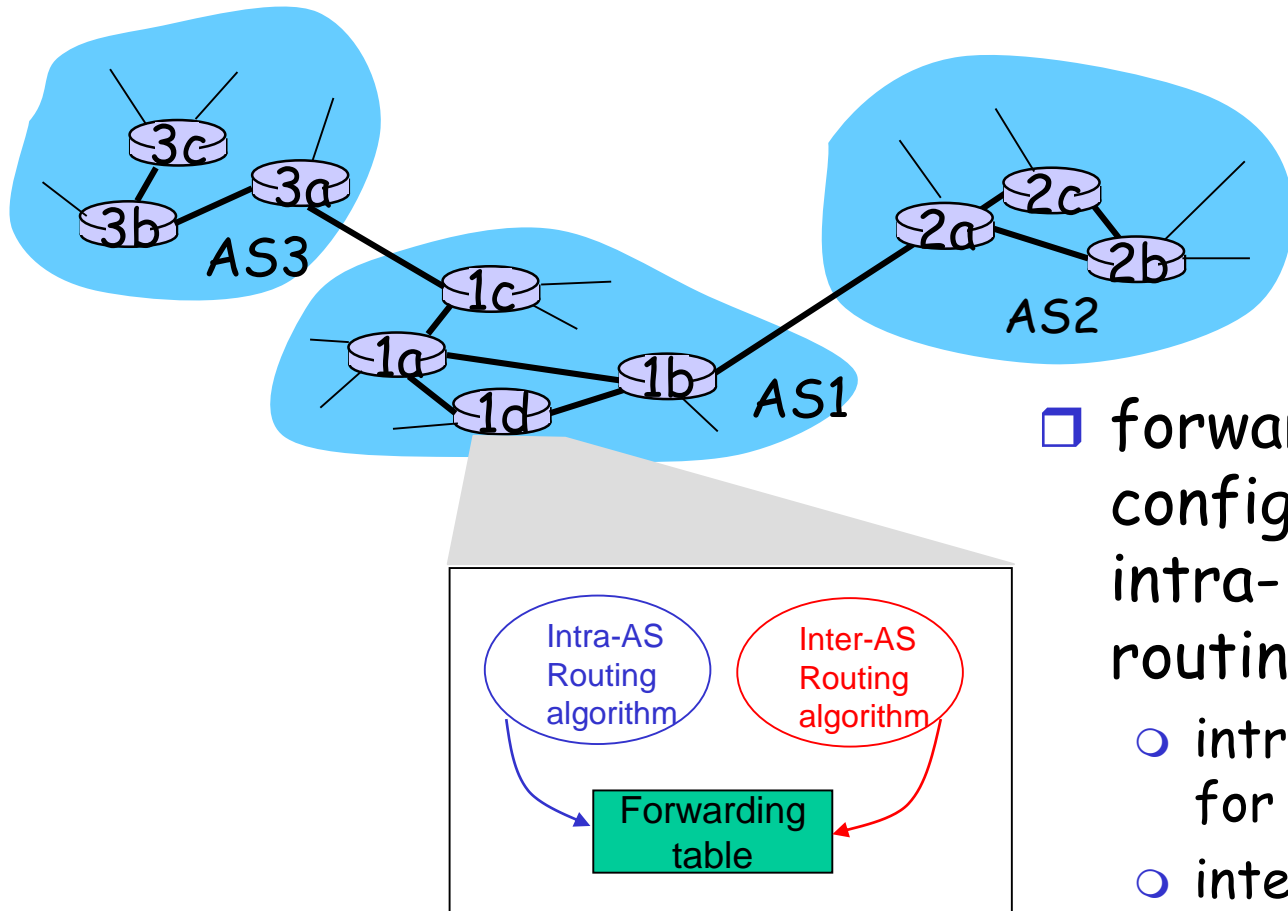
# Hierarchical Routing

- ❑ aggregate routers into regions, "autonomous systems" (AS)
- ❑ routers in same AS run same routing protocol
  - "intra-AS" routing protocol
  - routers in different AS can run different intra-AS routing protocol

## Gateway router

- ❑ Direct link to router in another AS

# Interconnected ASes



- ❑ forwarding table configured by both intra- and inter-AS routing algorithm
  - intra-AS sets entries for internal dests
  - inter-AS & intra-AS sets entries for external dests

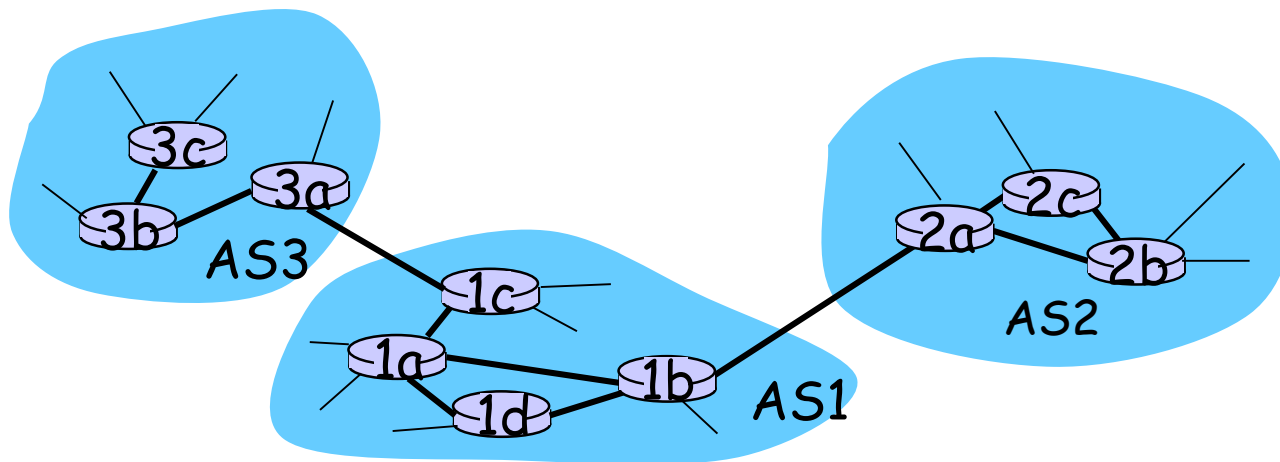
# Inter-AS tasks

- suppose router in AS1 receives datagram destined outside of AS1:
  - router should forward packet to gateway router, but which one?

## AS1 must:

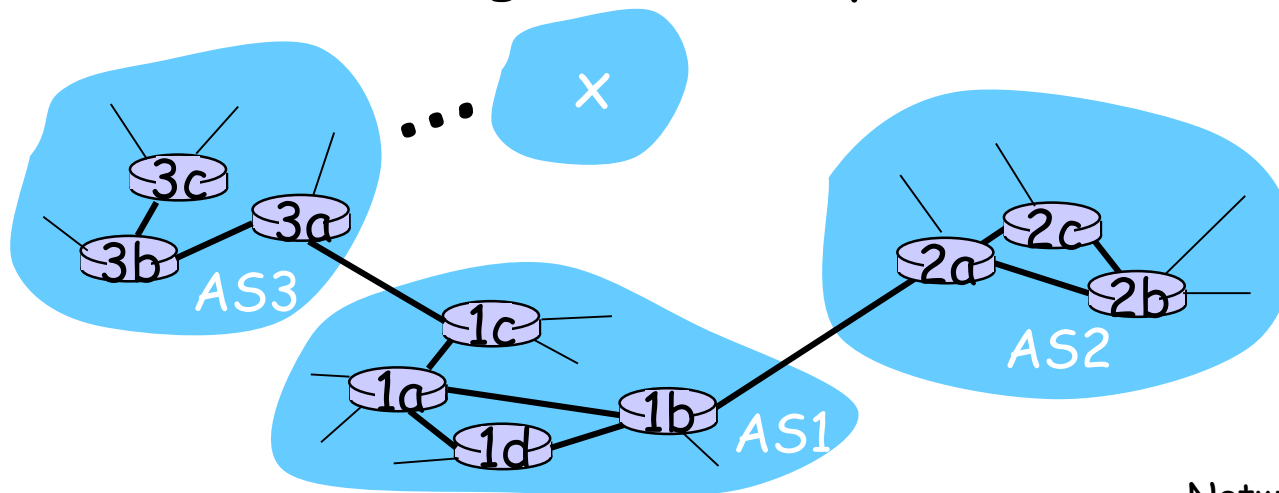
1. learn which destds are reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

**Job of inter-AS routing!**



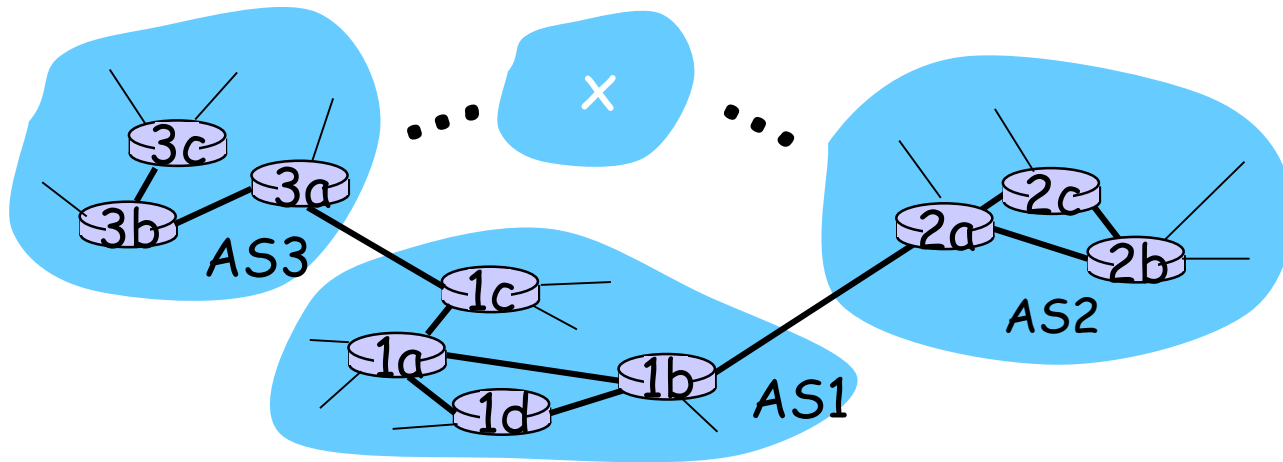
## Example: Setting forwarding table in router 1d

- suppose AS1 learns (via inter-AS protocol) that subnet **x** reachable via AS3 (gateway 1c) but not via AS2.
- inter-AS protocol propagates reachability info to all internal routers.
- router 1d determines from intra-AS routing info that its interface **I** is on the least cost path to 1c.
  - installs forwarding table entry **(x,I)**



# Example: Choosing among multiple ASes

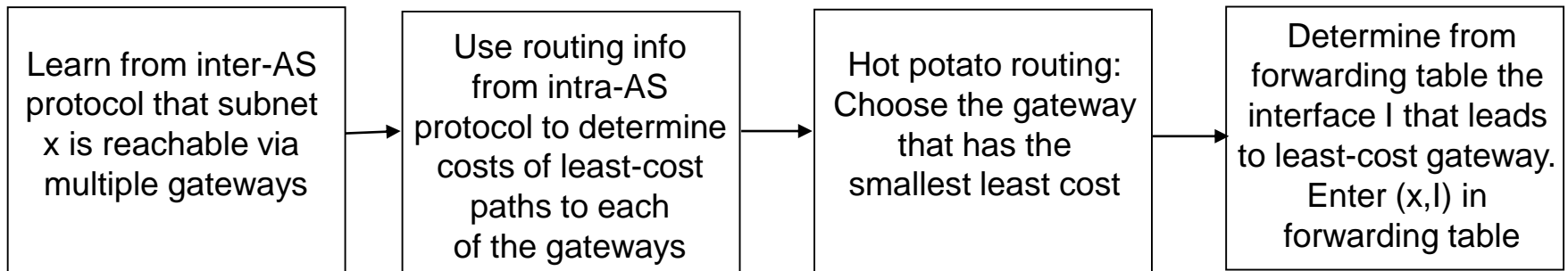
- now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 and from AS2.
- to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**.
  - this is also job of inter-AS routing protocol!





# Example: Choosing among multiple ASes

- ❑ now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 and from AS2.
- ❑ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**.
  - this is also job of inter-AS routing protocol!
- ❑ **hot potato routing**: send packet towards closest of two routers.



# Chapter 4: Network Layer

- ❑ 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

# Intra-AS Routing

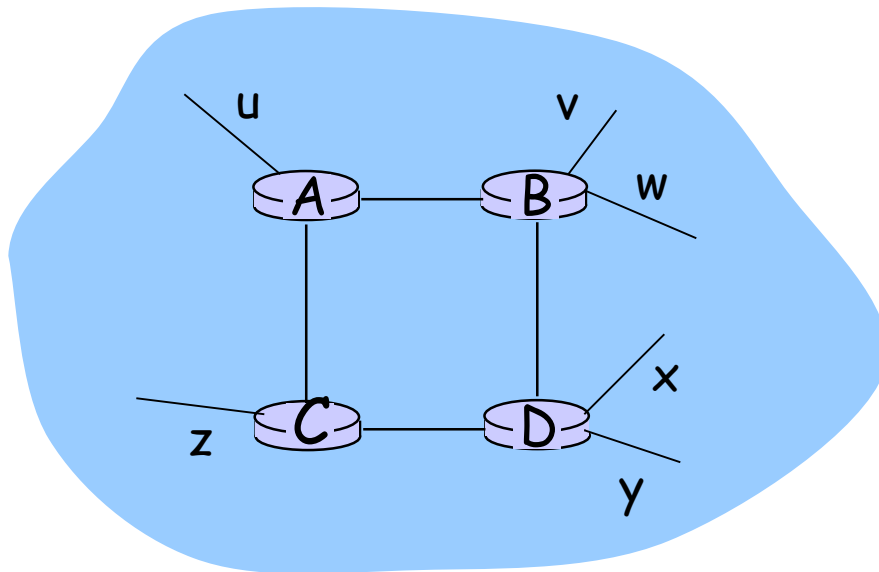
- ❑ also known as **Interior Gateway Protocols (IGP)**
- ❑ most common Intra-AS routing protocols:
  - RIP: Routing Information Protocol
  - OSPF: Open Shortest Path First
  - IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

# Chapter 4: Network Layer

- ❑ 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

# RIP (Routing Information Protocol)

- ❑ distance vector algorithm
- ❑ included in BSD-UNIX Distribution in 1982
- ❑ distance metric: # of hops (max = 15 hops)



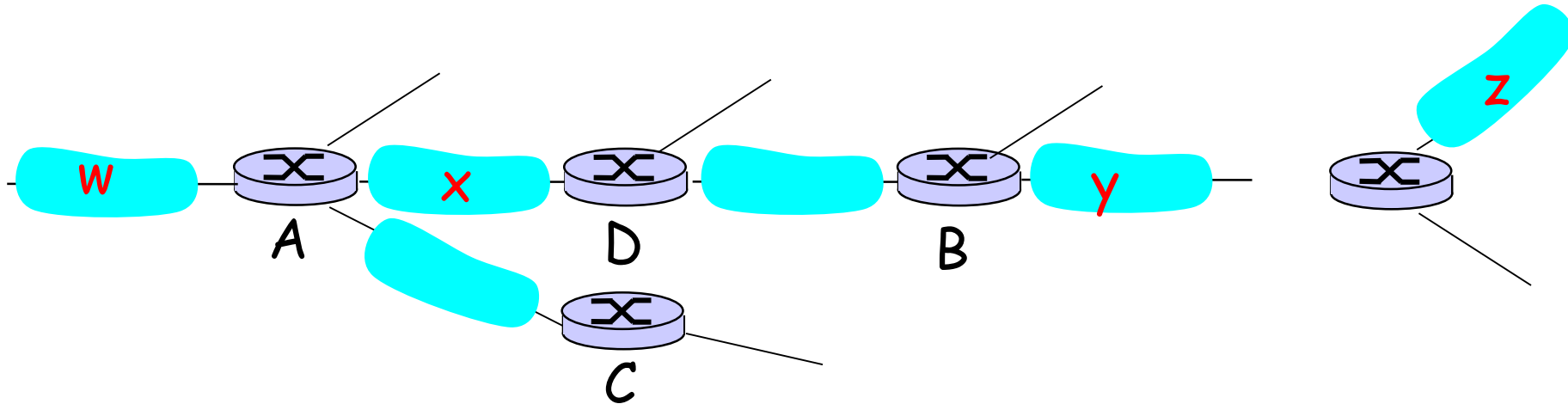
From router A to subnets:

<u>destination</u>	<u>hops</u>
u	1
v	2
w	2
x	3
y	3
z	2

# RIP advertisements

- ❑ distance vectors: exchanged among neighbors every 30 sec via Response Message (also called **advertisement**)
- ❑ each advertisement: list of up to 25 destination subnets within AS

# RIP: Example



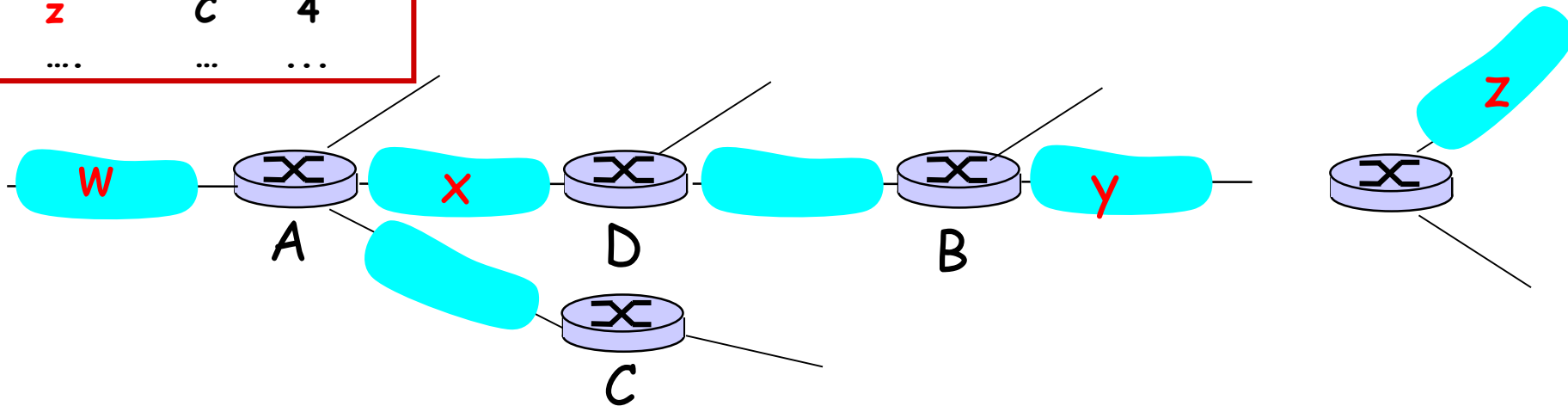
Destination Network	Next Router	Num. of hops to dest.
W	A	2
Y	B	2
Z	B	7
X	--	1
....	....	....

Routing/Forwarding table in D

# RIP: Example

Dest	Next	hops
w	-	1
x	-	1
z	C	4
...	...	...

Advertisement  
from A to D



Destination Network	Next Router	Num. of hops to dest.
w	A	2
y	B	2
z	<del>B</del> A	<del>7</del> 5
x	--	1
...	...	...

Routing/Forwarding table in D

Network Layer 4-40



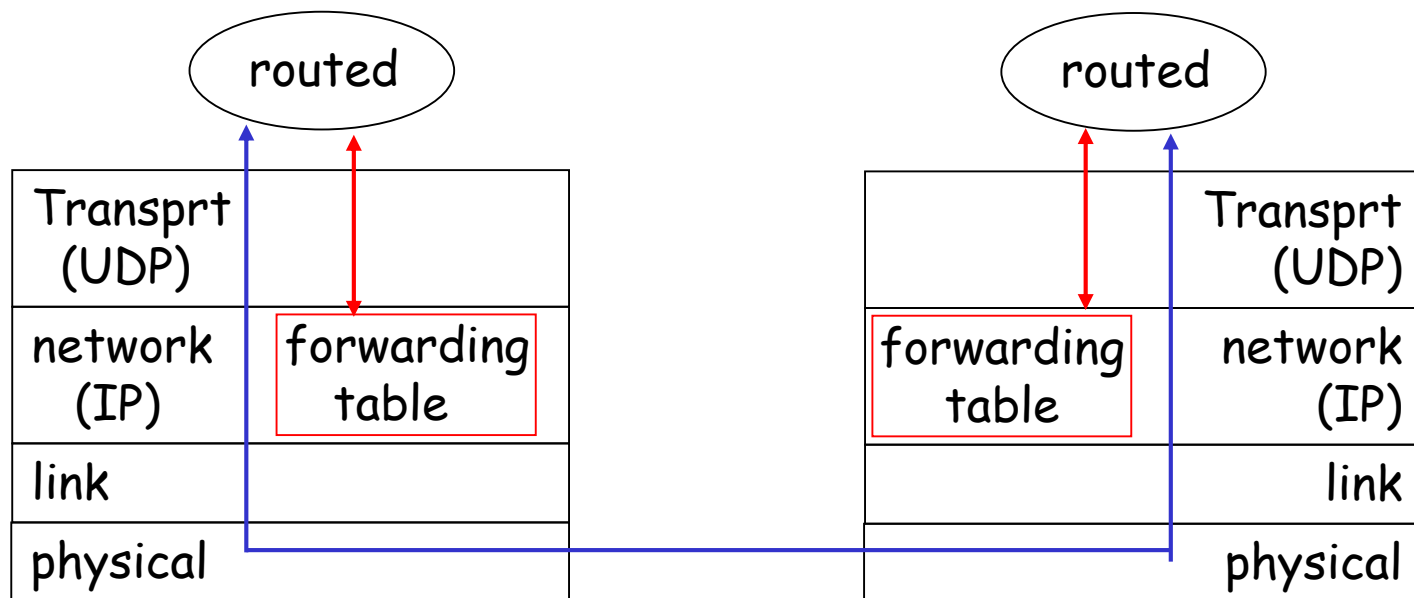
# RIP: Link Failure and Recovery

If no advertisement heard after 180 sec -->  
neighbor/link declared dead

- routes via neighbor invalidated
- new advertisements sent to neighbors
- neighbors in turn send out new advertisements (if tables changed)
- link failure info quickly (?) propagates to entire net
- *poison reverse* used to prevent ping-pong loops (infinite distance = 16 hops)

# RIP Table processing

- ❑ RIP routing tables managed by **application-level** process called route-d (daemon)
- ❑ advertisements sent in UDP packets, periodically repeated



# Chapter 4: Network Layer

- ❑ 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

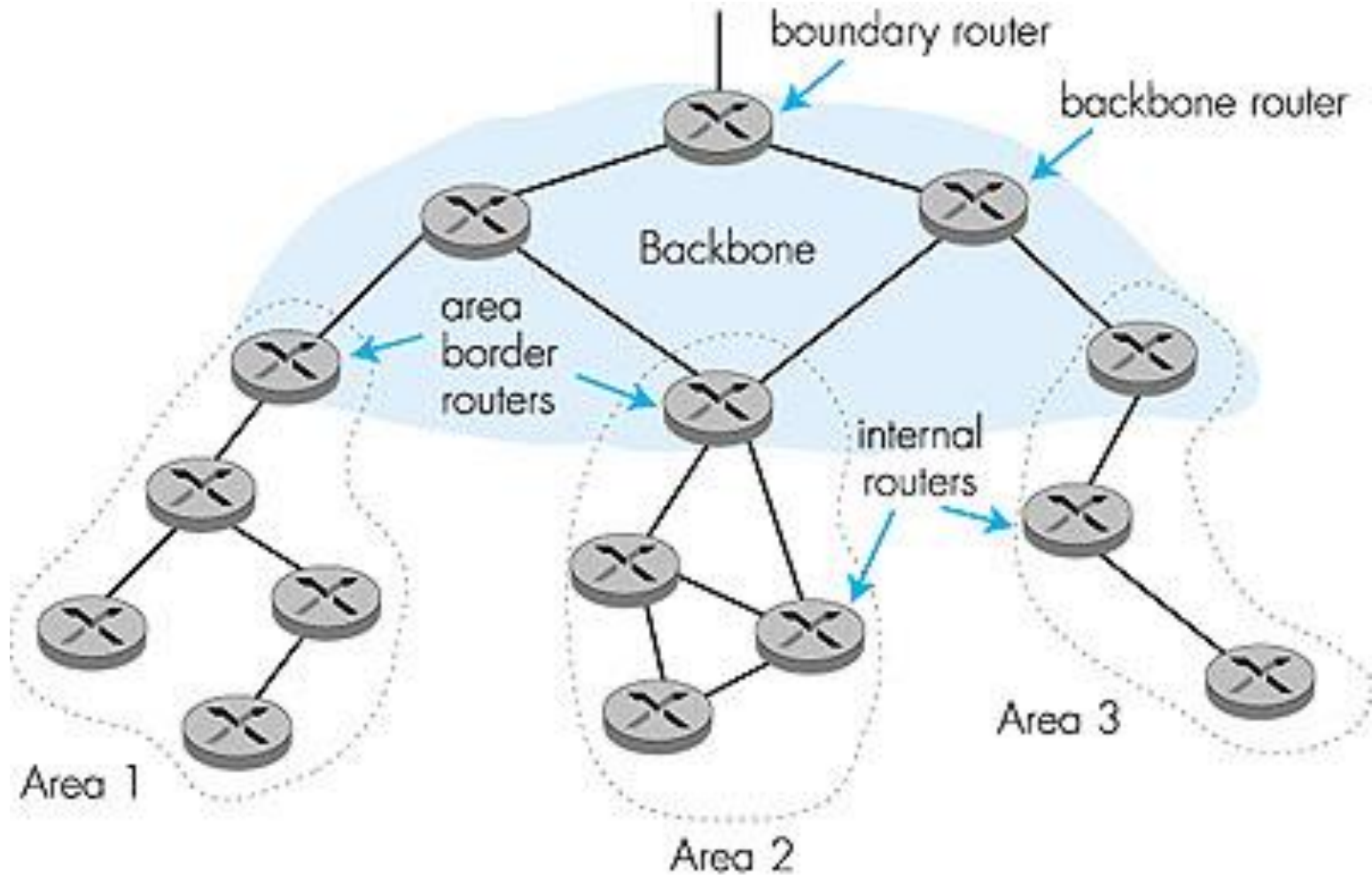
# OSPF (Open Shortest Path First)

- ❑ “open”: publicly available
- ❑ uses Link State algorithm
  - LS packet dissemination
  - topology map at each node
  - route computation using Dijkstra's algorithm
- ❑ OSPF advertisement carries one entry per neighbor router
- ❑ advertisements disseminated to **entire** AS (via flooding)
  - carried in OSPF messages directly over IP (rather than TCP or UDP)

# OSPF "advanced" features (not in RIP)

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

# Hierarchical OSPF



# Hierarchical OSPF

- ❑ **two-level hierarchy:** local area, backbone.
  - Link-state advertisements only in area
  - each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.
- ❑ **area border routers:** "summarize" distances to nets in own area, advertise to other Area Border routers.
- ❑ **backbone routers:** run OSPF routing limited to backbone.
- ❑ **boundary routers:** connect to other AS's.

# Chapter 4: Network Layer

- ❑ 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



# Internet inter-AS routing: BGP

- ❑ **BGP (Border Gateway Protocol):** *the de facto standard*
- ❑ 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.
- ❑ allows subnet to advertise its existence to rest of Internet: *"I am here"*

# Internet inter-AS routing: BGP

## □ BGP (Border Gateway Protocol):

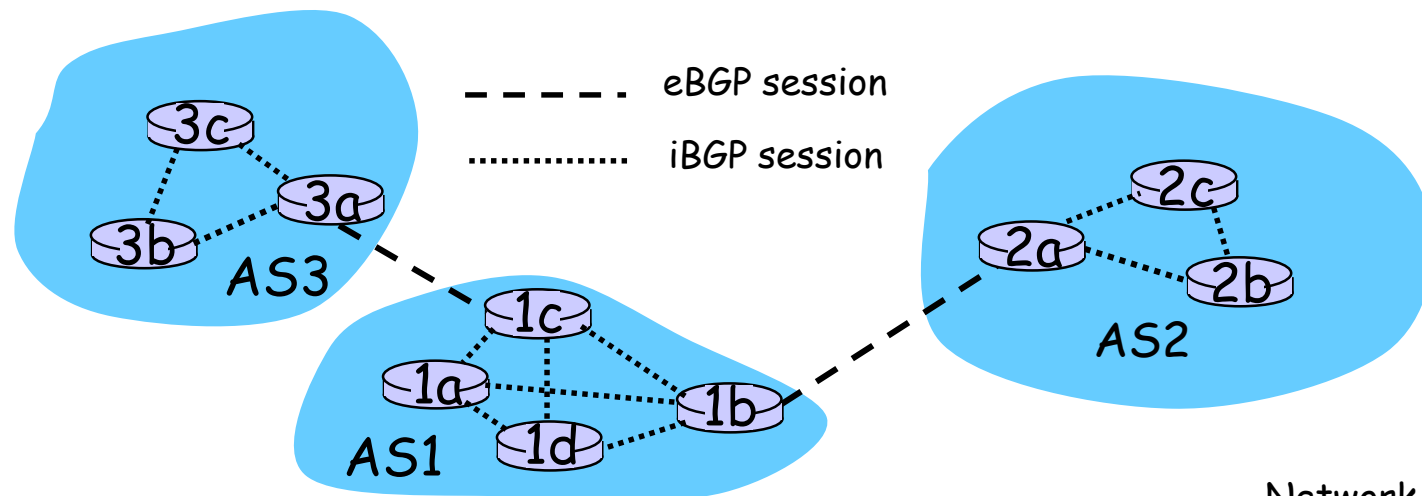
- BGP 是**唯一**一个用来处理像因特网大小的网络协议
- 也是**唯一**能够妥善处理好不相关路由域间的多路连接协议。

## □ 在BGP网络中，可以将一个网络分成多个自治系统。

- 自治系统间使用**eBGP**广播路由
- 自治系统内使用**iBGP**在自己的网络内广播路由。

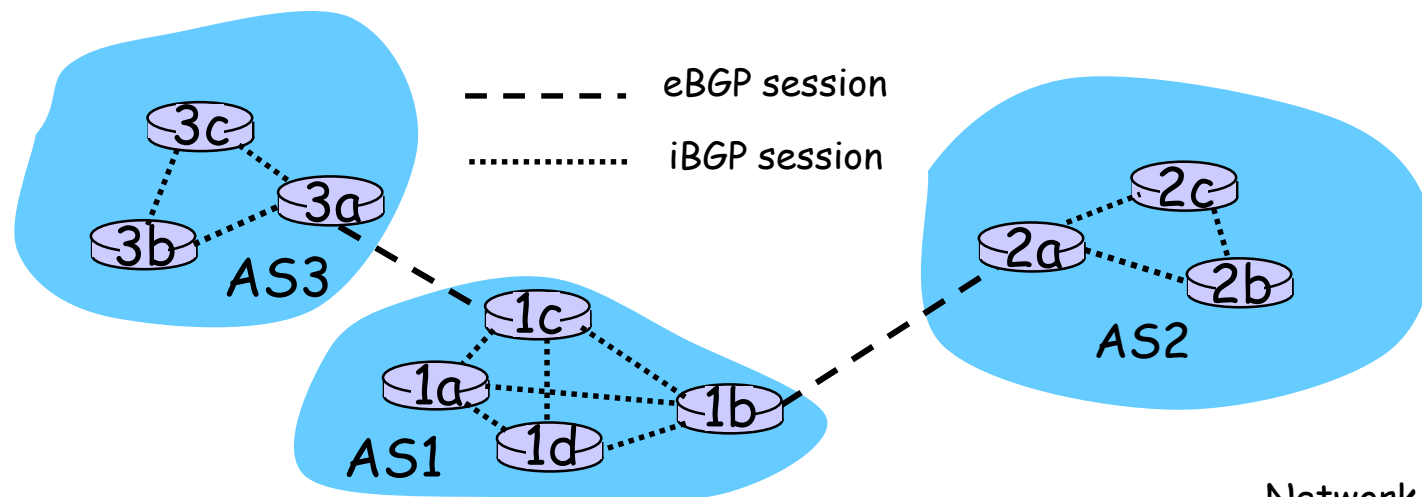
# BGP basics

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



# Distributing reachability info

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

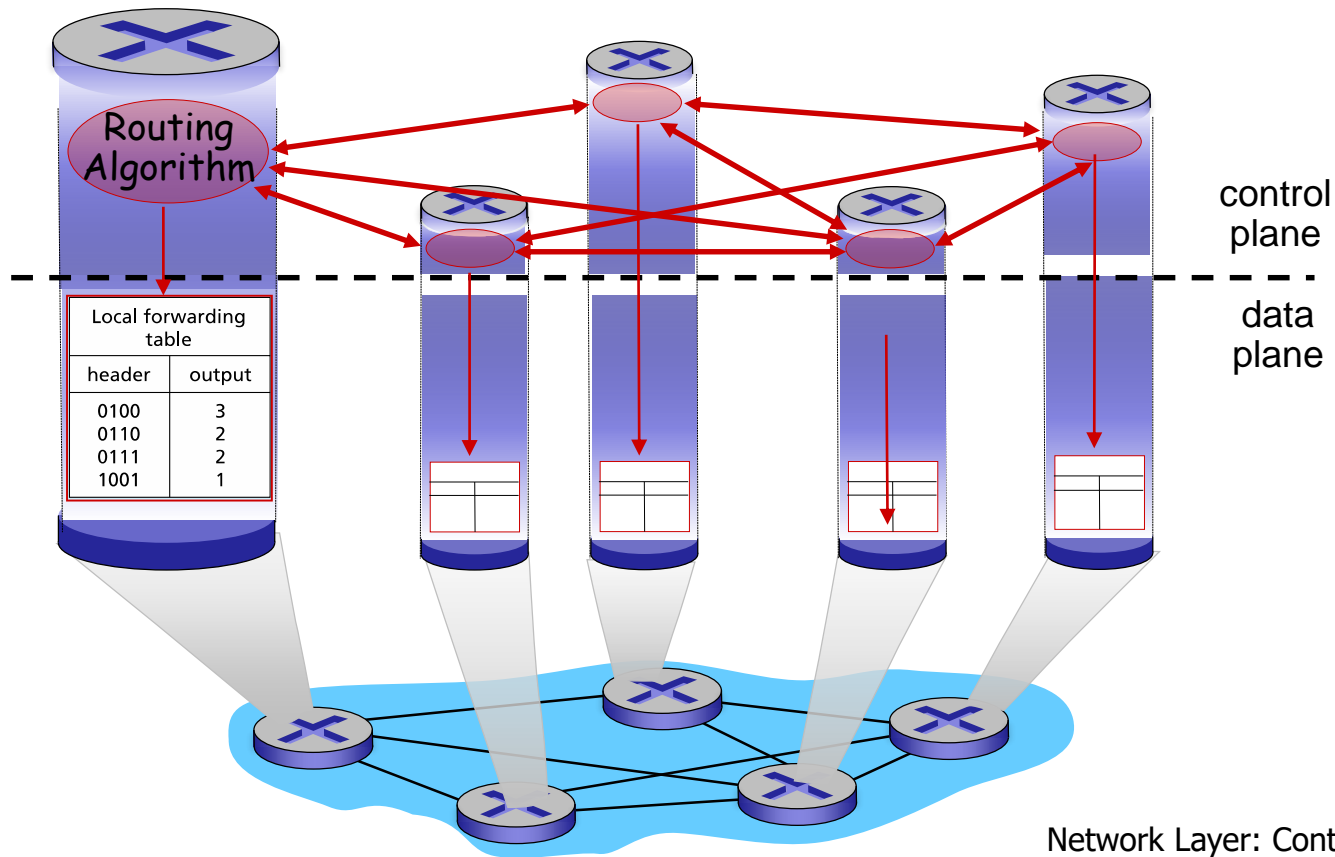


# Software defined networking (SDN)

- ❑ Internet network layer: historically has been implemented via distributed, per-router approach
  - *monolithic* router contains switching hardware, runs proprietary implementation of Internet standard protocols (IP, RIP, IS-IS, OSPF, BGP) in proprietary router OS (e.g., Cisco IOS)
  - different "middleboxes" for different network layer functions: firewalls, load balancers, NAT boxes, ..
- ❑ ~2005: renewed interest in rethinking network control plane

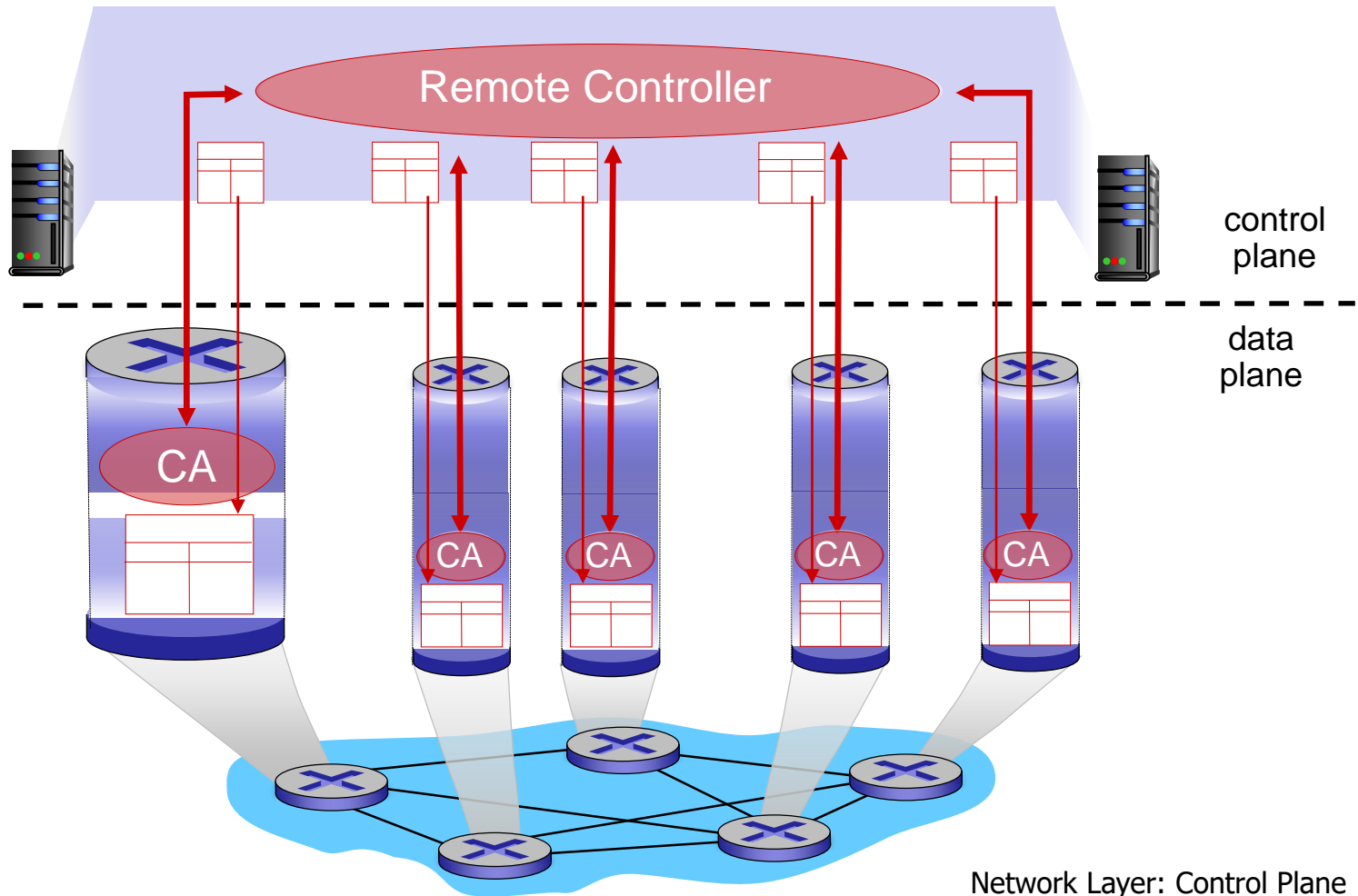
# Recall: per-router control plane

Individual routing algorithm components *in each and every router* interact with each other in control plane to compute forwarding tables



# Recall: 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?

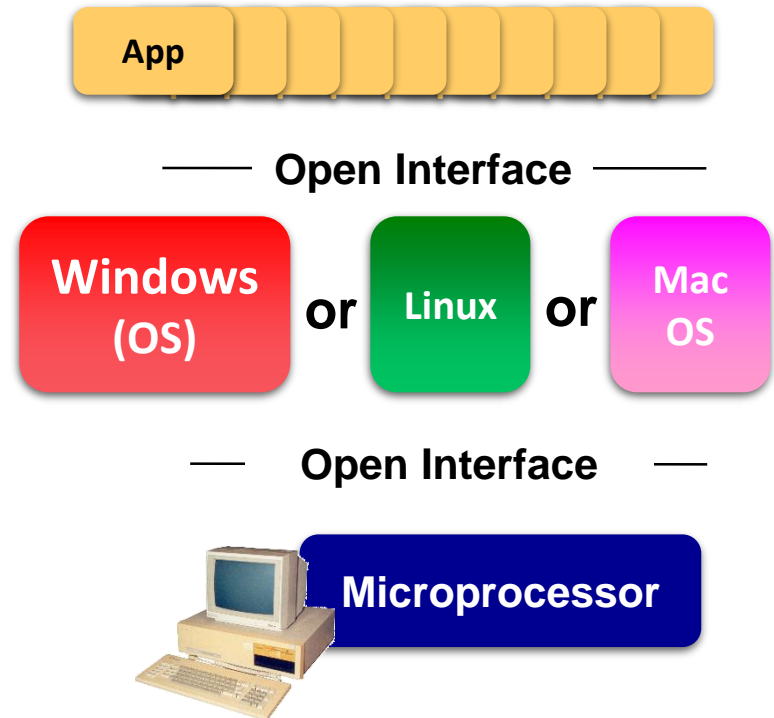
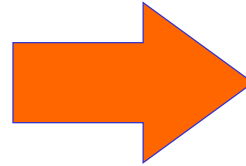
- easier network management: avoid router misconfigurations, greater flexibility of traffic flows
- table-based forwarding (recall OpenFlow API) allows “programming” routers
  - centralized “programming” easier: compute tables centrally and distribute
  - distributed “programming”: more difficult: compute 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\*

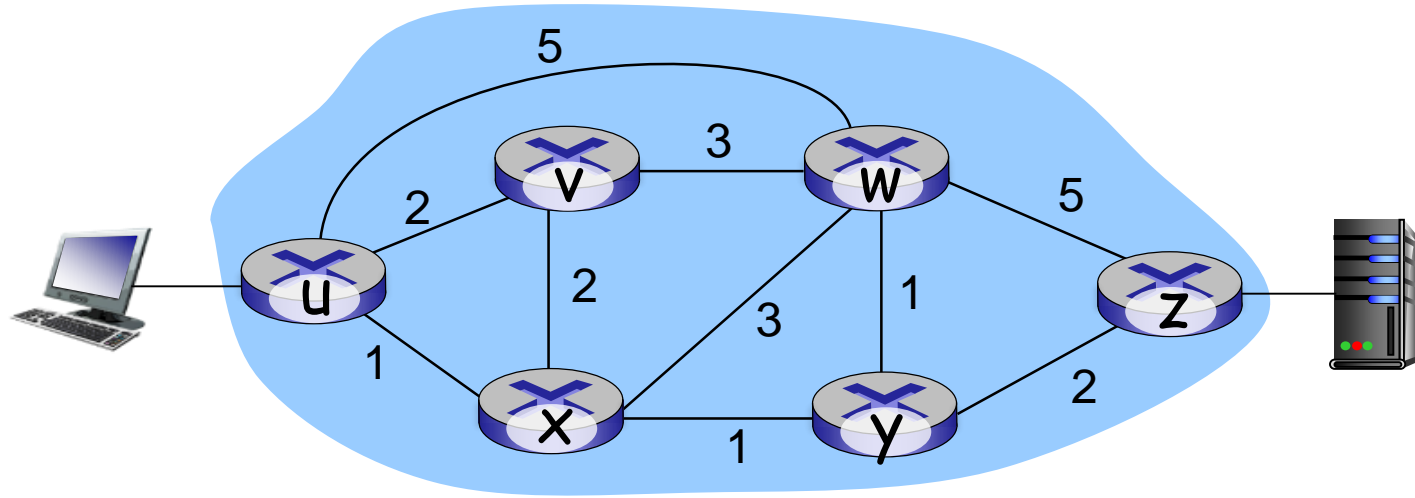


Vertically integrated  
Closed, proprietary  
Slow innovation  
Small industry



Horizontal  
Open interfaces  
Rapid innovation  
Huge industry

# Traffic engineering: difficult traditional routing

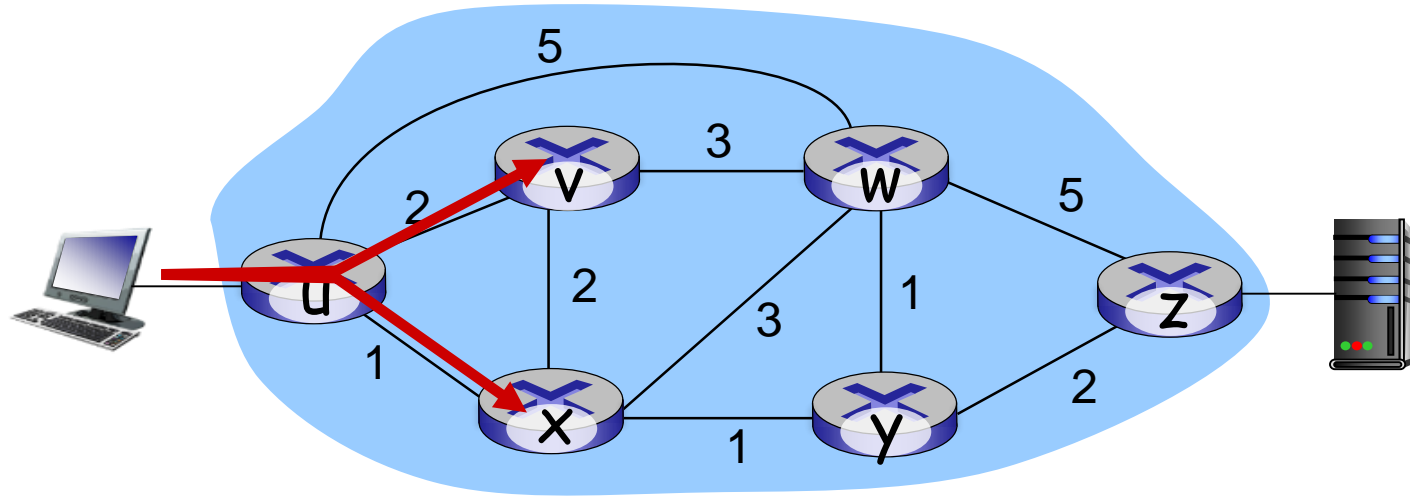


Q: what if network operator wants u-to-z traffic to flow along uvwz, x-to-z traffic to flow xwyz?

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

*Link weights are only control "knobs": wrong!*

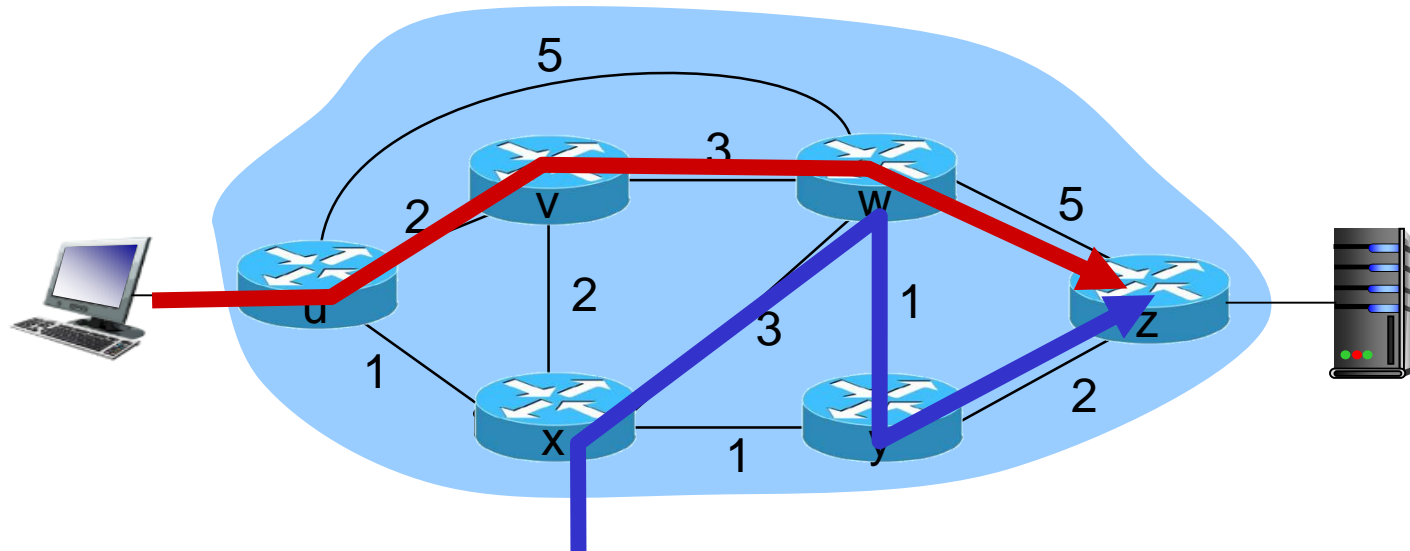
# Traffic engineering: difficult



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

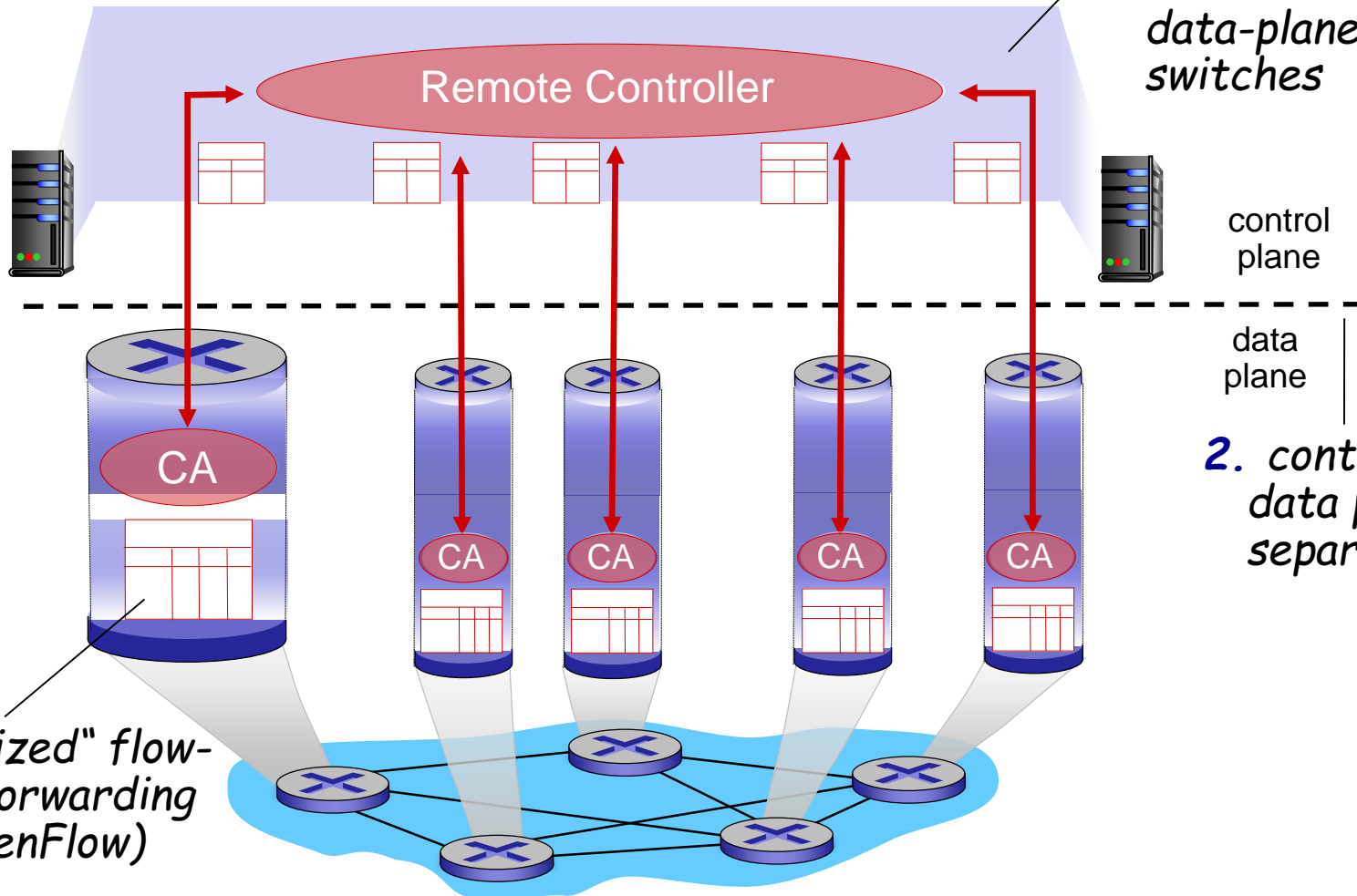
routing

access control

...

load balance

3. control plane functions external to data-plane switches



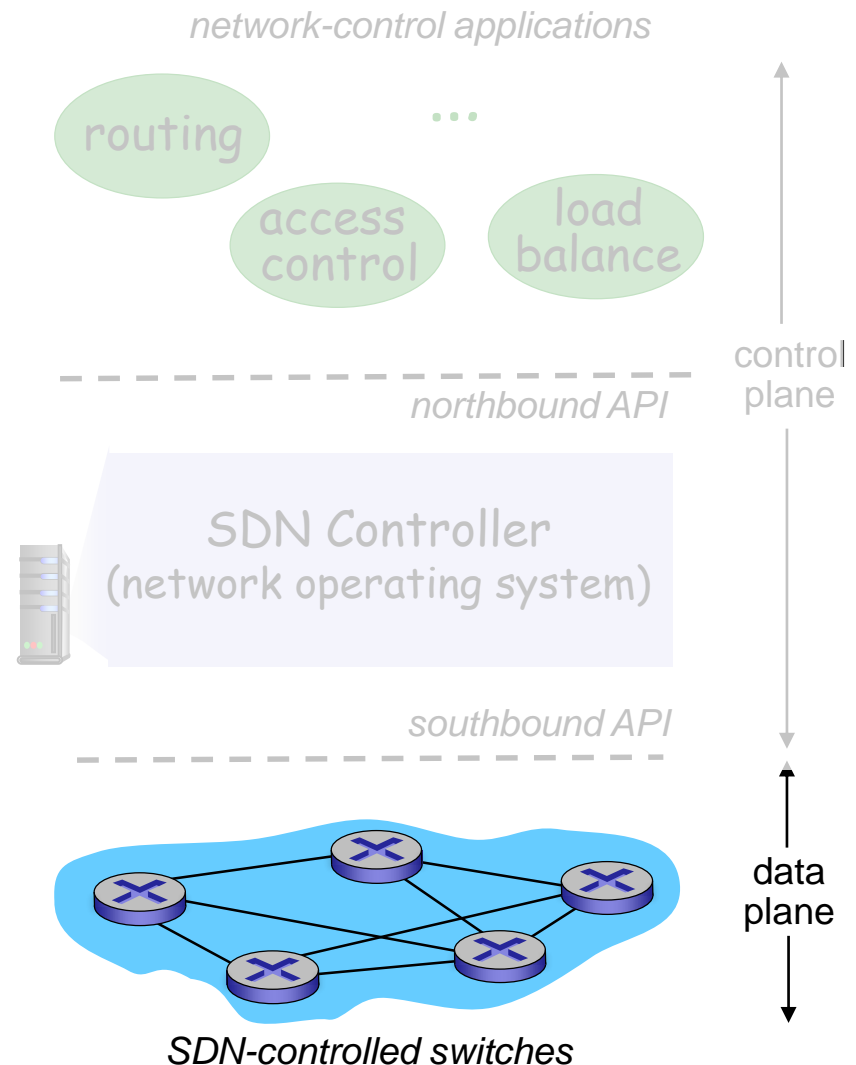
1: generalized flow-based forwarding (e.g., OpenFlow)

2. control, data plane separation

# SDN perspective: data plane switches

## *Data plane switches*

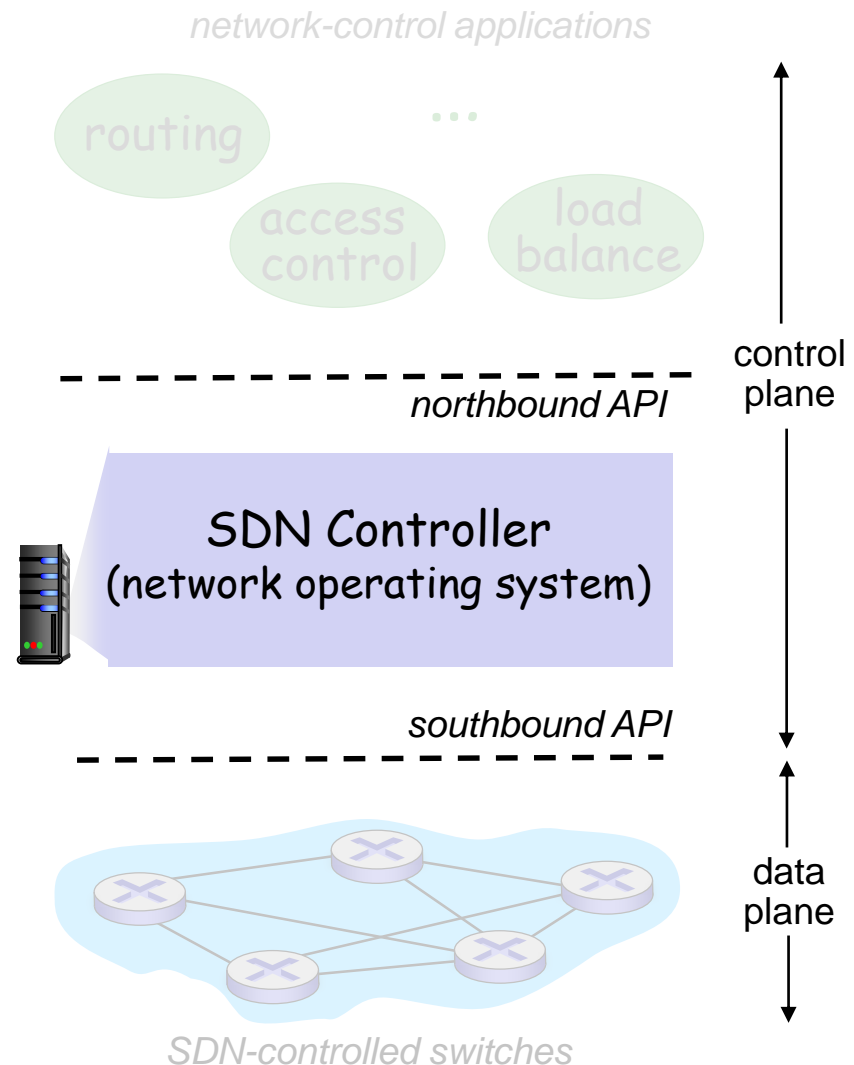
- ❑ fast, simple, commodity switches implementing generalized data-plane forwarding (Section 4.4) in hardware
- ❑ switch flow table computed, installed by controller
- ❑ API for table-based switch control (e.g., OpenFlow)
- ❑ protocol for communicating with controller (e.g., OpenFlow)



# SDN perspective: SDN controller

## *SDN controller (network OS):*

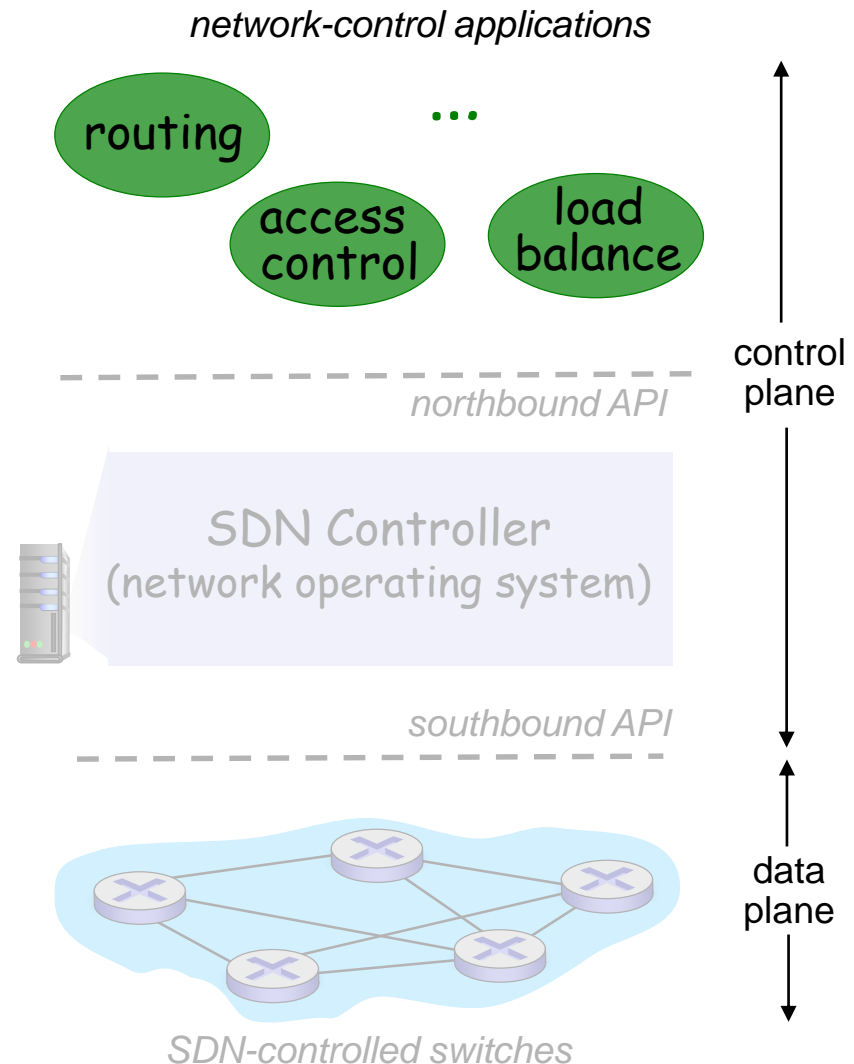
- maintain network state information
- interacts with network control applications "above" via northbound API
- interacts with network switches "below" via southbound API
- implemented as distributed system for performance, scalability, fault-tolerance, robustness



# SDN perspective: control applications

## *network-control apps:*

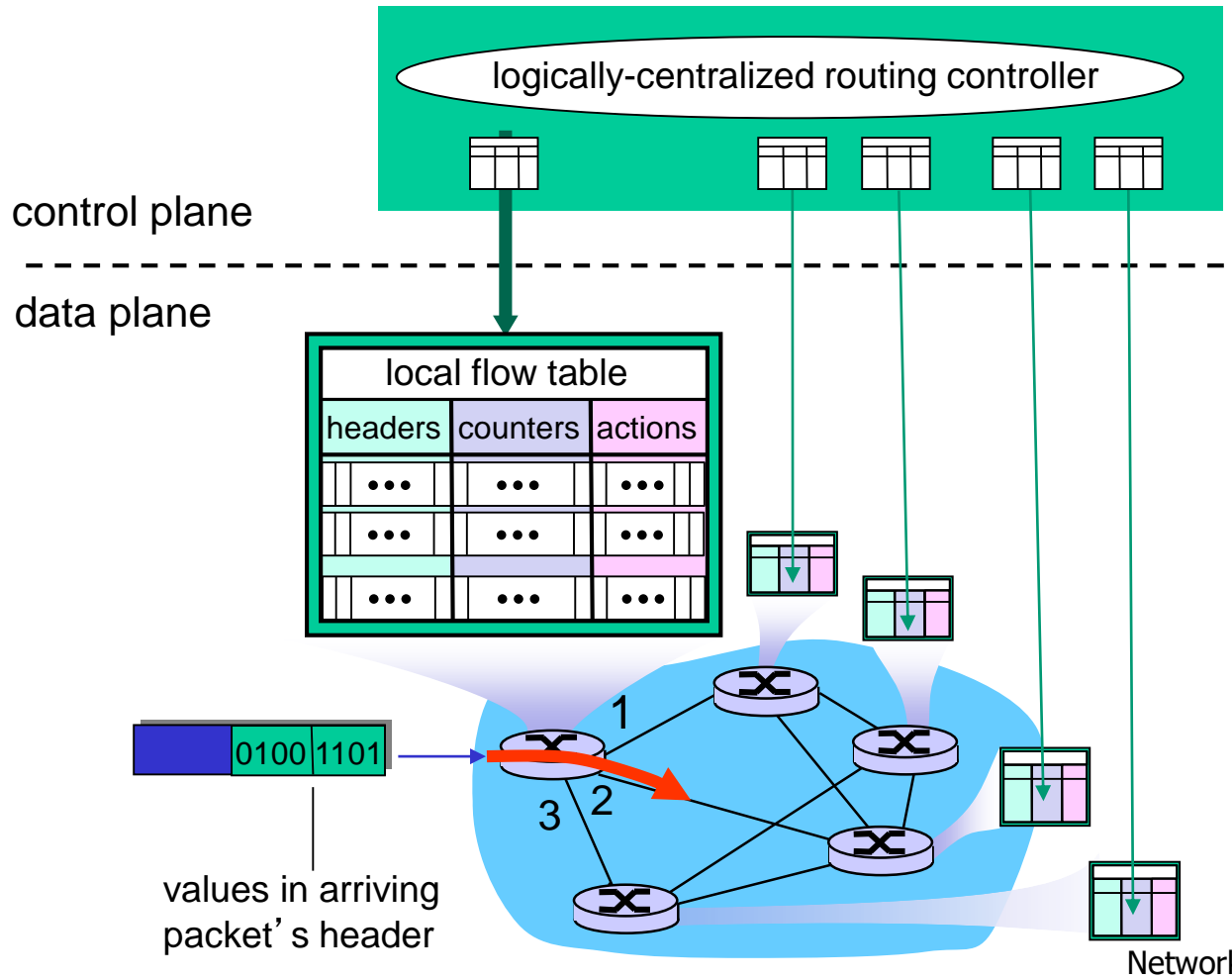
- “brains” of control: implement control functions using lower-level services, API provided by SDN controller
- *unbundled*: can be provided by 3<sup>rd</sup> party: distinct from routing vendor, or SDN controller





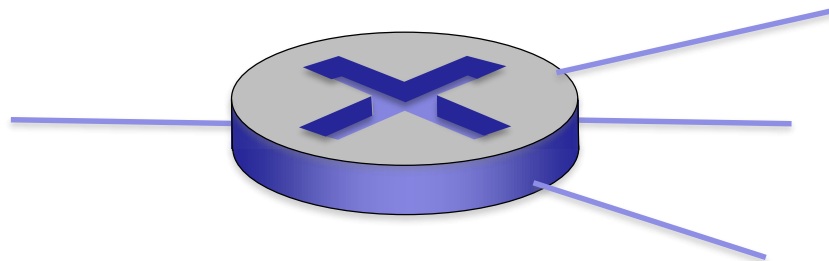
# Generalized Forwarding and SDN

Each router contains a *flow table* that is computed and distributed by a *logically centralized routing controller*



# OpenFlow data plane abstraction

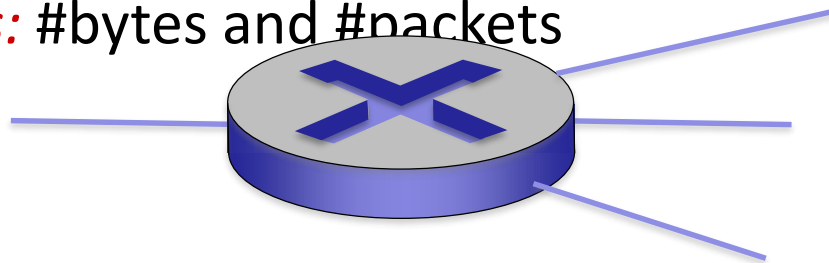
- ❑ *flow*: defined by header fields
- ❑ generalized forwarding: simple packet-handling rules
  - *Pattern*: match values in packet header fields
  - *Actions: for matched packet*: drop, forward, modify, matched packet or send matched packet to controller
  - *Priority*: disambiguate overlapping patterns
  - *Counters*: #bytes and #packets



*Flow table in a router (computed and distributed by controller) define router's match+action rules*

# OpenFlow data plane abstraction

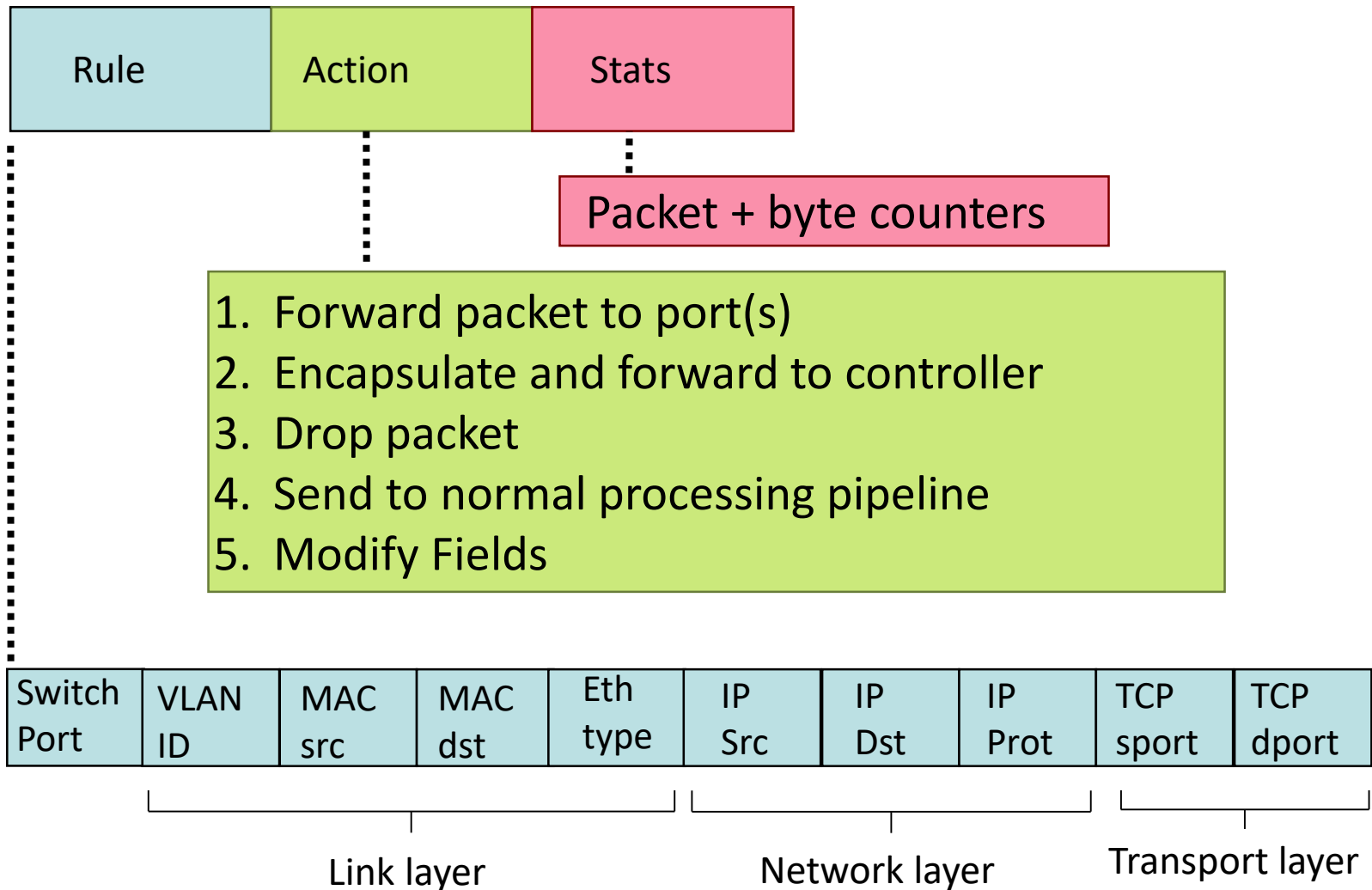
- ❑ *flow*: defined by header fields
- ❑ generalized forwarding: simple packet-handling rules
  - *Pattern*: match values in packet header fields
  - *Actions: for matched packet*: drop, forward, modify, matched packet or send matched packet to controller
  - *Priority*: disambiguate overlapping patterns
  - *Counters*: #bytes and #packets



\* : wildcard

1. src=1.2.\*.\*, dest=3.4.5.\* → drop
2. src = \*.\*.\*.\*, dest=3.4.\*.\* → forward(2)
3. src=10.1.2.3, dest=\*.\*.\*.\* → send to controller

# OpenFlow: Flow Table Entries



# Examples

## Destination-based forwarding:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
-------------	---------	---------	----------	---------	--------	--------	---------	-----------	-----------	--------

\* \* \* \* \* 51.6.0.8 \* \* \* port6

*IP datagrams destined to IP address  
51.6.0.8 should be forwarded to router output  
port 6*

## Firewall:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Forward
-------------	---------	---------	----------	---------	--------	--------	---------	-----------	-----------	---------

\* \* \* \* \* 22 drop

*do not forward (block) all datagrams destined to TCP  
port 22*

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Forward
-------------	---------	---------	----------	---------	--------	--------	---------	-----------	-----------	---------

\* \* \* \* \* 128.119.1.1 \* \* \* drop

*do not forward (block) all datagrams sent by host  
128.119.1.1*

# Examples

## Destination-based layer 2 (switch) forwarding:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	22:A7:23: 11:E1:02	*	*	*	*	*	*	*	*	port3

*layer 2 frames from MAC address  
22:A7:23:11:E1:02 should be forwarded to  
output port 3*

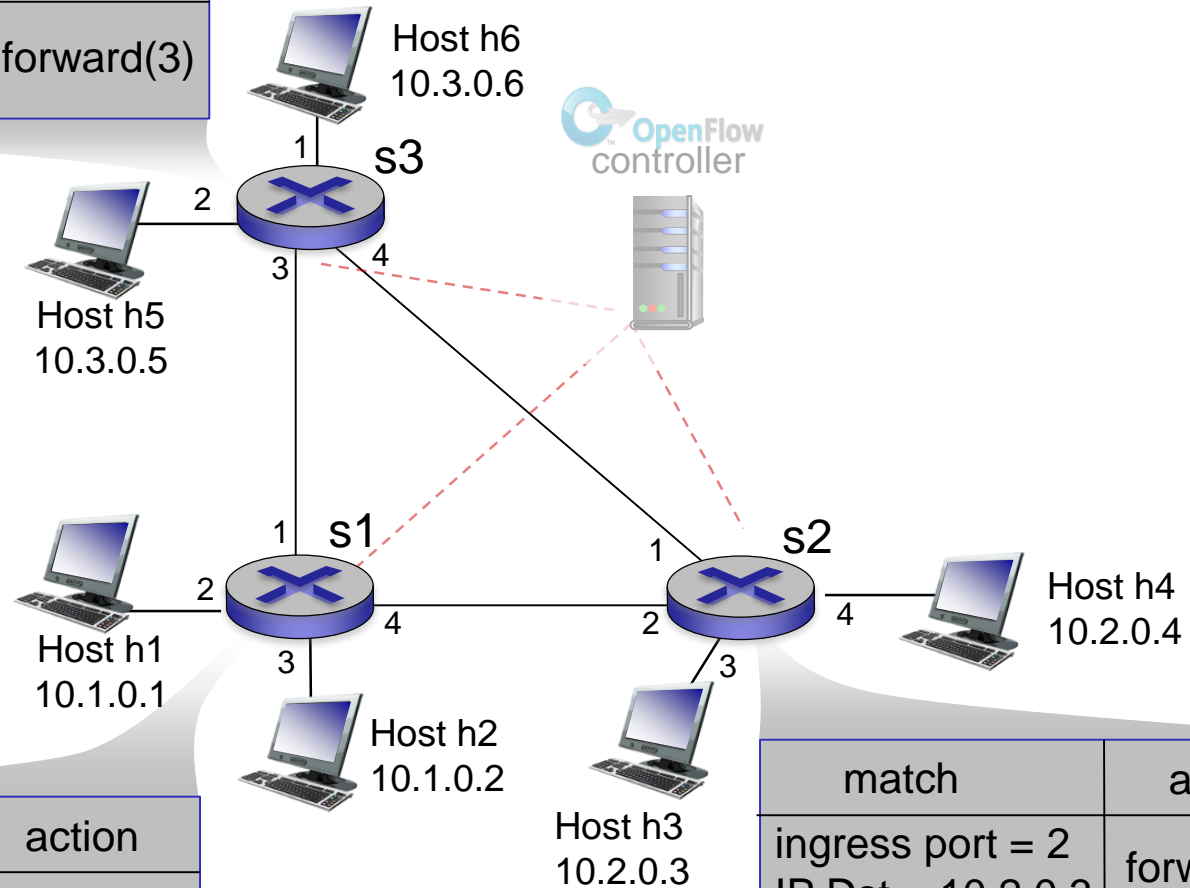
# OpenFlow abstraction

- *match+action*: unifies different kinds of devices
- Router
  - *match*: longest destination IP prefix
  - *action*: forward out a link
- Switch
  - *match*: destination MAC address
  - *action*: forward or flood
- Firewall
  - *match*: IP addresses and TCP/UDP port numbers
  - *action*: permit or deny
- NAT
  - *match*: IP address and port
  - *action*: rewrite address and port

# OpenFlow example

*Example:* datagrams from hosts h5 and h6 should be sent to h3 or h4, via s1 and from there to s2

match	action
IP Src = 10.3.*.* IP Dst = 10.2.*.*	forward(3)



match	action
ingress port = 1 IP Src = 10.3.*.* IP Dst = 10.2.*.*	forward(4)

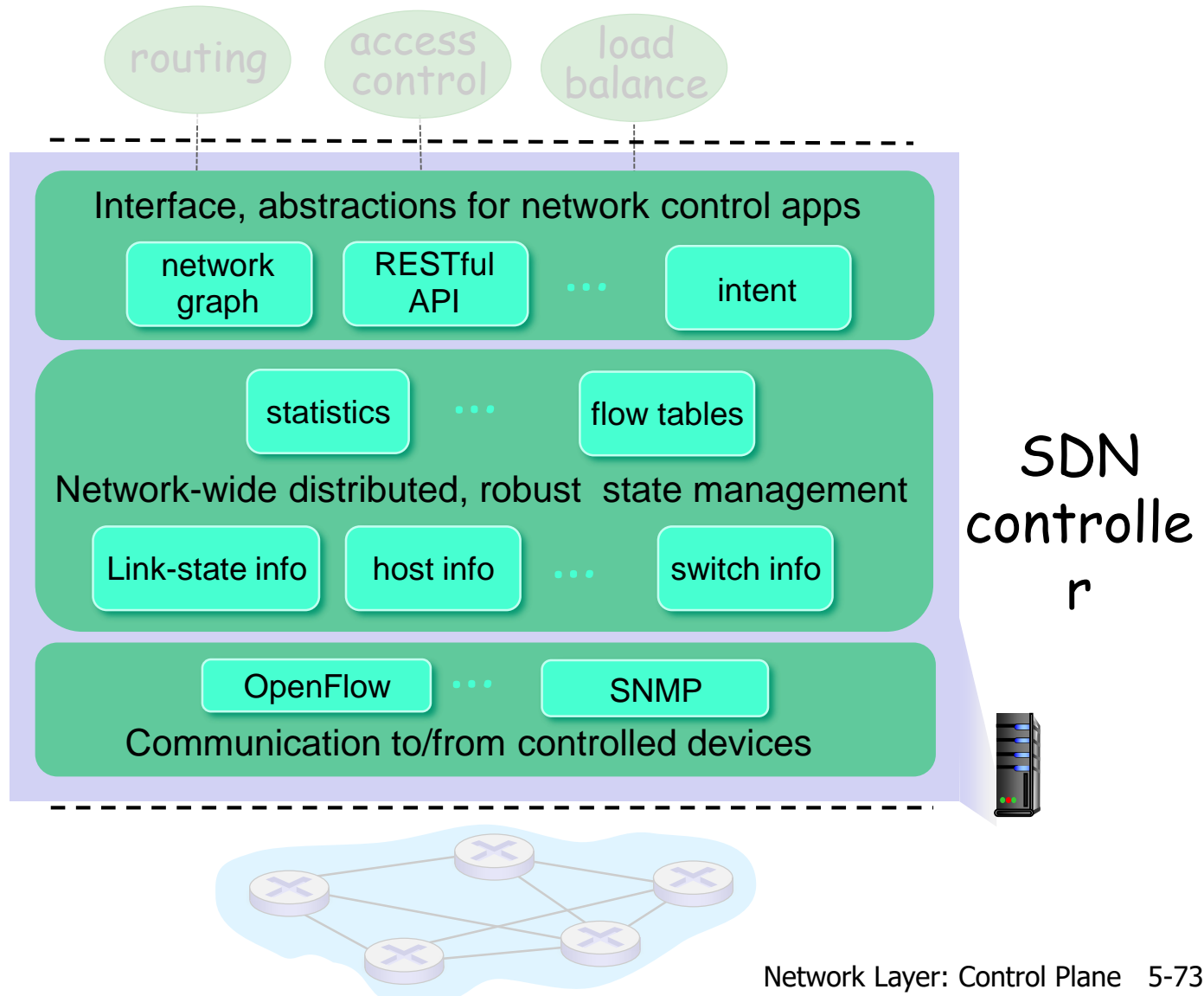
match	action
ingress port = 2 IP Dst = 10.2.0.3	forward(3)
ingress port = 2 IP Dst = 10.2.0.4	forward(4)



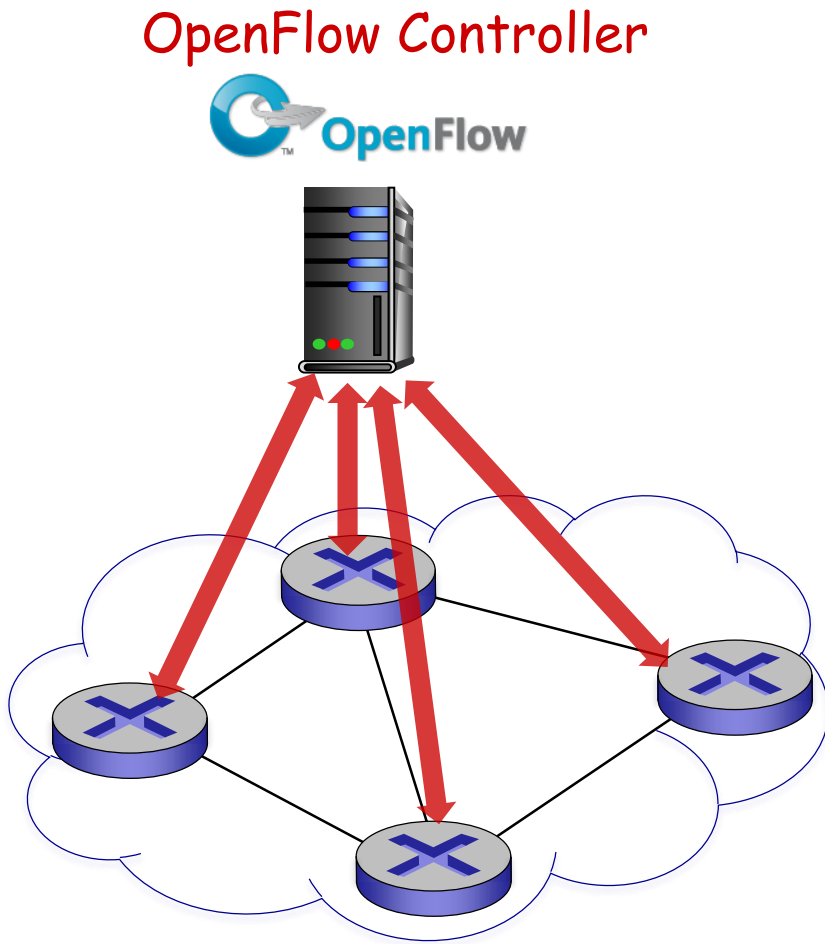
# Components of SDN controller

Interface layer  
to network  
control apps:  
abstractions API

Network-wide  
state management  
layer: state of  
networks links,  
switches, services:  
a *distributed*  
*database*  
*communication*  
layer:  
communicate  
between SDN  
controller and  
controlled  
switches



# OpenFlow protocol

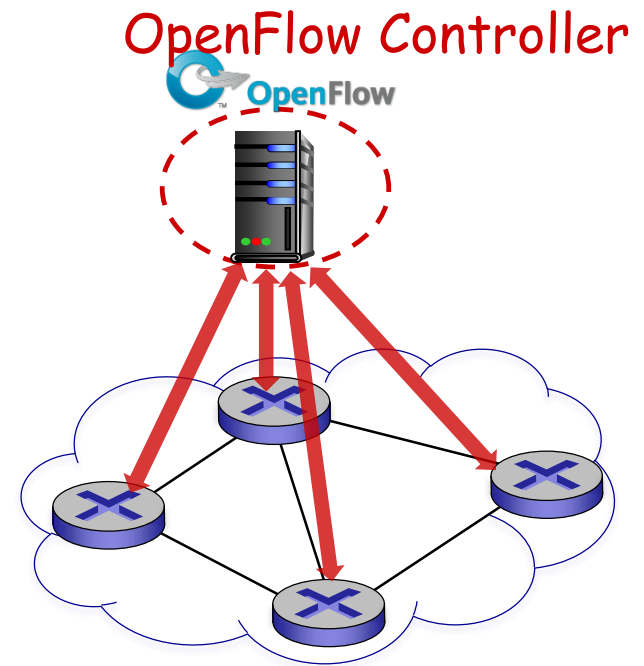


- ❑ operates between controller, switch
- ❑ TCP used to exchange messages
  - optional encryption
- ❑ three classes of OpenFlow messages:
  - controller-to-switch
  - asynchronous (switch to controller)
  - symmetric (misc)

# OpenFlow: controller-to-switch messages

## Key controller-to-switch messages

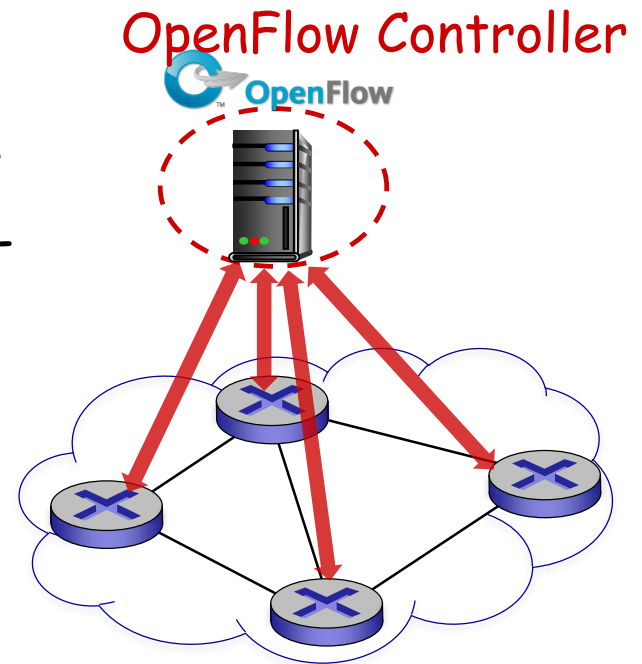
- ❑ *features*: controller queries switch features, switch replies
- ❑ *configure*: controller queries/sets switch configuration parameters
- ❑ *modify-state*: add, delete, modify flow entries in the OpenFlow tables
- ❑ *packet-out*: controller can send this packet out of specific switch port



# OpenFlow: switch-to-controller messages

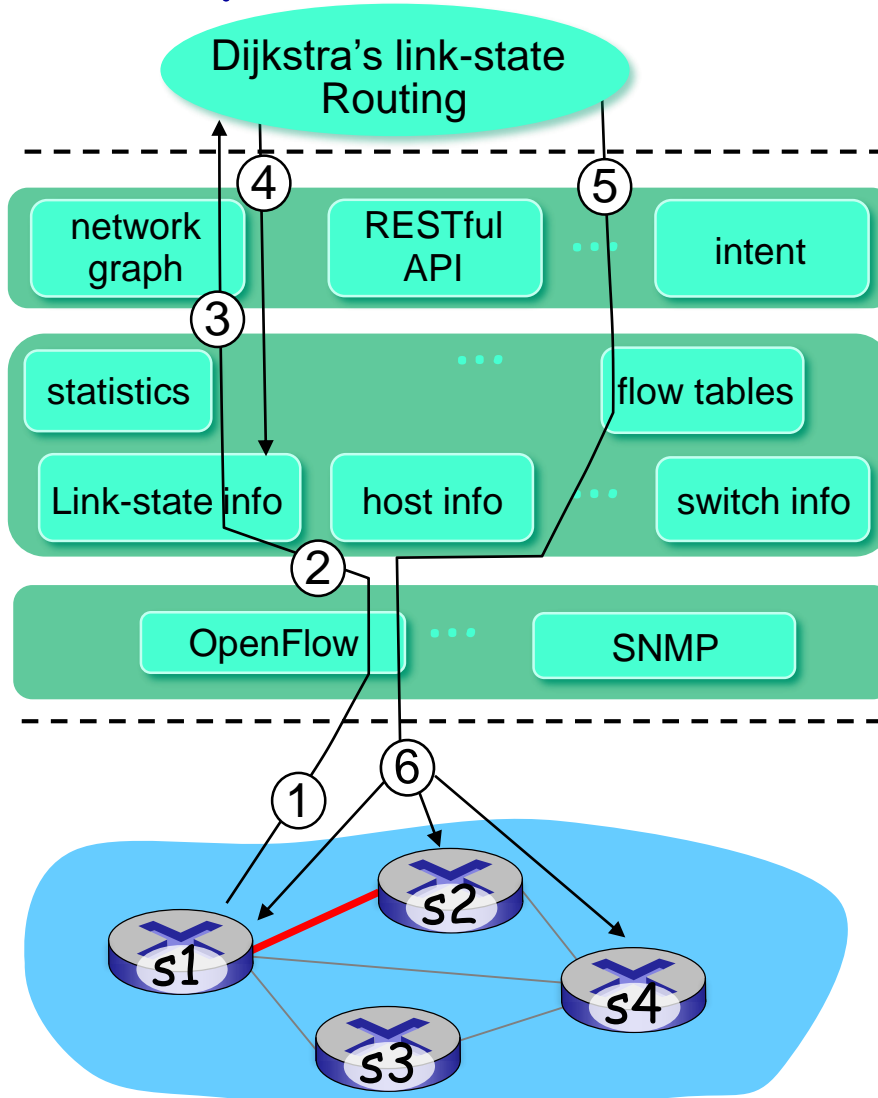
## Key switch-to-controller messages

- ❑ **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.



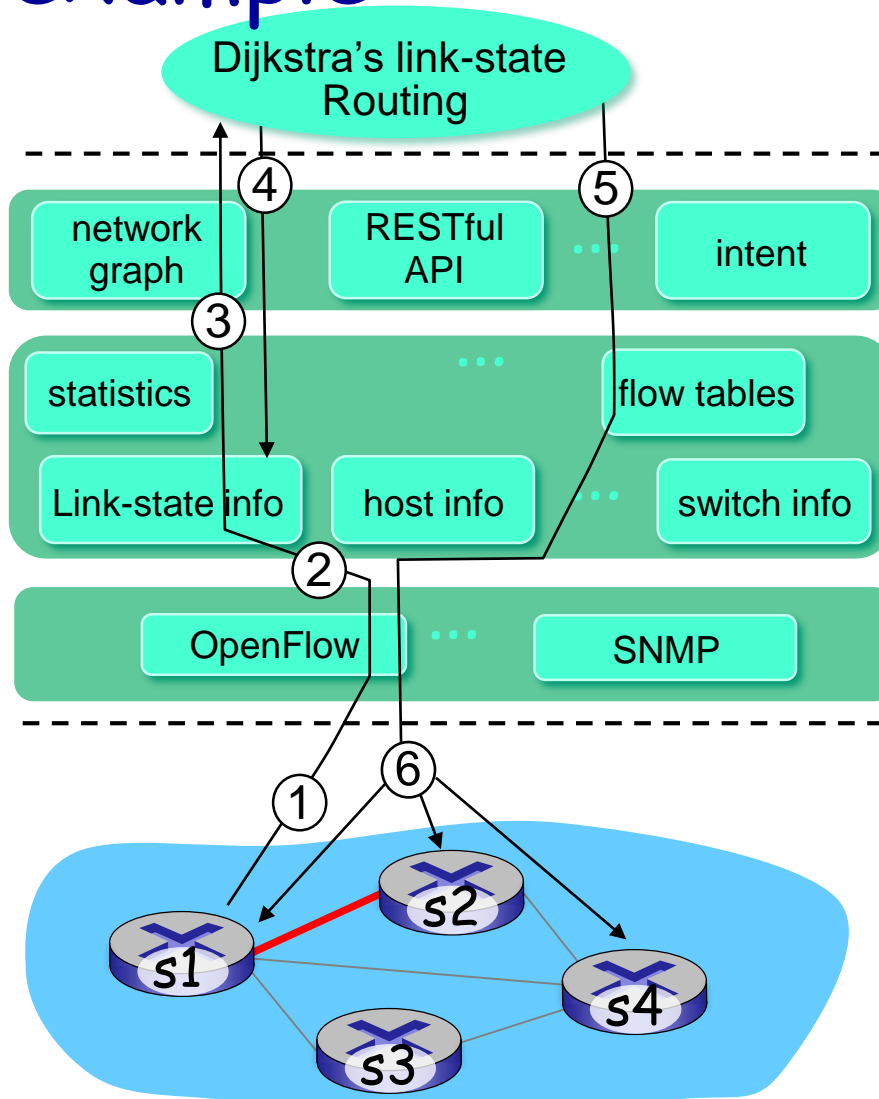
Fortunately, network operators don't "program" switches by creating/sending OpenFlow messages directly. Instead use higher-level abstraction at controller

# SDN: control/data plane interaction example



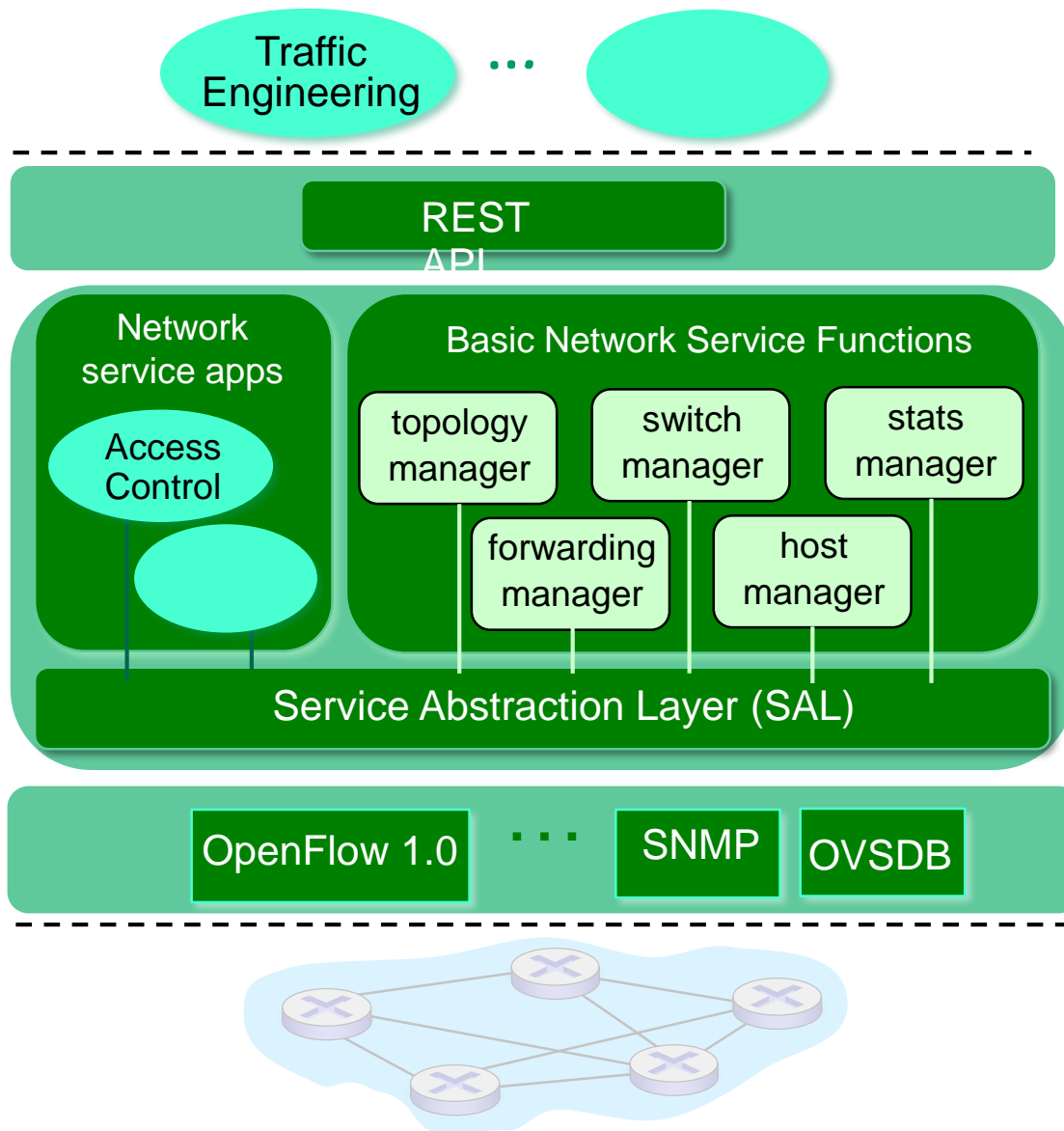
- ① S1, experiencing link failure using OpenFlow port status message to notify controller
- ② SDN controller receives OpenFlow message, updates link status info
- ③ Dijkstra's routing algorithm application has previously registered to be called when ever link status changes. It is called.
- ④ Dijkstra's routing algorithm access network graph info, link state info in controller, computes new routes

# SDN: control/data plane interaction example



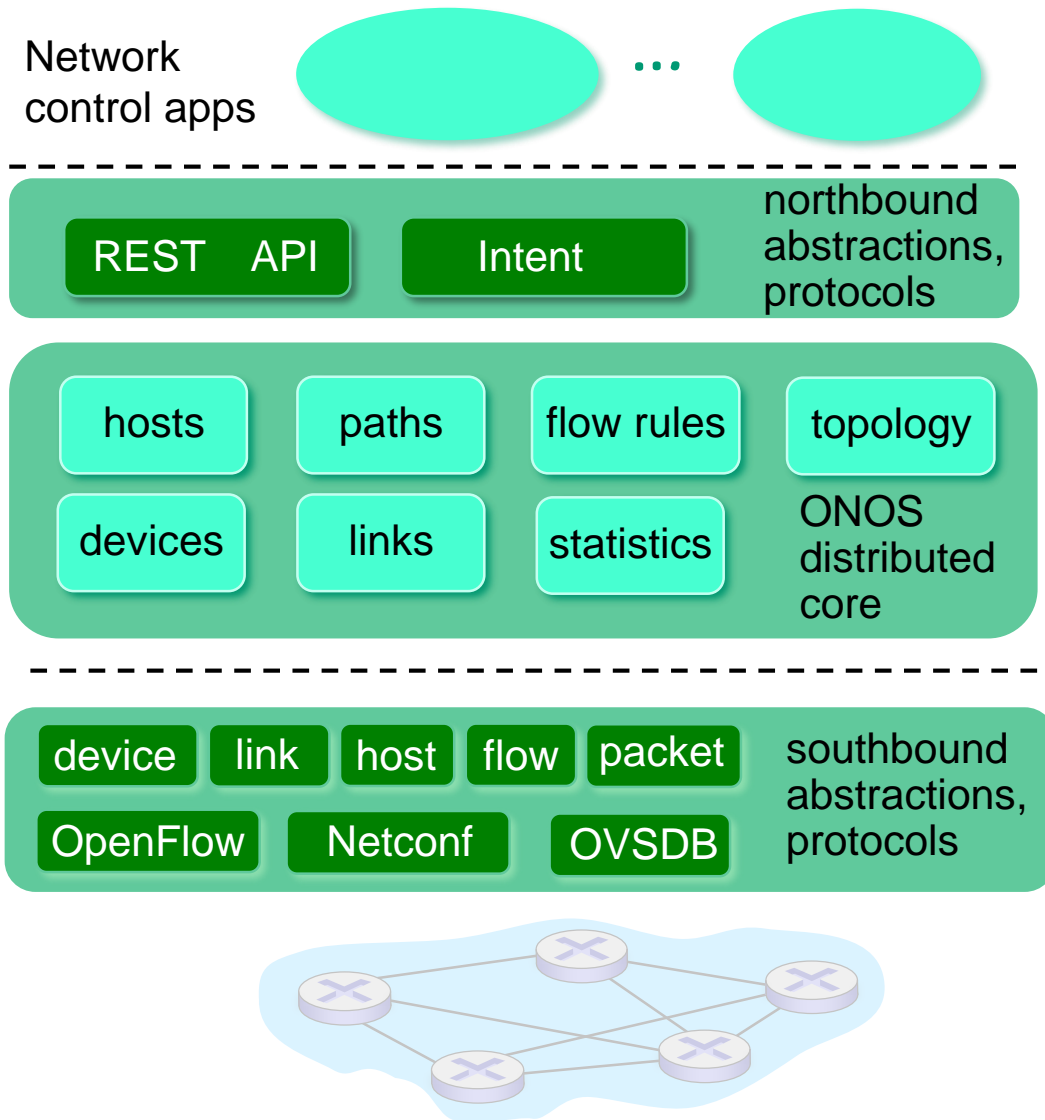
- ⑤ link state routing app interacts with flow-table-computation component in SDN controller, which computes new flow tables needed
- ⑥ Controller uses OpenFlow to install new tables in switches that need updating

# OpenDaylight (ODL) controller



- 最普及的SDN控制器
- ODL Lithium controller
- network apps may be contained within, or be external to SDN controller
- Service Abstraction Layer: interconnects internal, external applications and services

# ONOS controller



- control apps separate from controller
- intent framework: high-level specification of service: what rather than how
- considerable emphasis on distributed core: service reliability, replication performance scaling



# SDN: selected challenges

- ❑ hardening the control plane: dependable, reliable, performance-scalable, secure distributed system
  - robustness to failures: leverage strong theory of reliable distributed system for control plane
  - dependability, security: “baked in” from day one?
- ❑ networks, protocols meeting mission-specific requirements
  - e.g., real-time, ultra-reliable, ultra-secure
- ❑ Internet-scaling

# 软件定义——SDN

## □ SDN: 网络操作系统

- 软件定义的真正落地，还是SDN。
- 2011年前后，OpenFlow被用于云计算平台中进行网络管理，并被广泛接受。



S12712



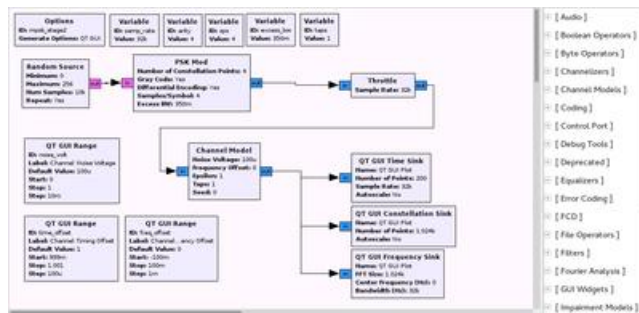
S12708

S12700 Series Agile Switches

# 软件定义——SDR

## □ SDR: 软件定义无线电

- 在通用可编程数字信号处理硬件平台上，利用软件定义来实现无线电台的各单元功能，对无线电信号进行调制或解调以及测量。
- **GNU Radio**: 开源SDR软件平台，结合通用软件无线电外设（**USRP**）开发板，进行快速开发，**WIFI**、**Bluetooth**、**RFID**、广播等



# 各种软件定义

## □ 还有很多例子

- 软件定义存储SDS
- 软件定义数据中心SDDC
- 软件定义雷达
- 软件定义汽车
- 软件定义卫星
- 。 。 。
- 软件定义制造
- 软件定义服务

# 什么是软件定义

- ❑ “软件定义”的定义
  - 用软件去定义系统的功能，用软件给硬件赋能
- ❑ 本质
  - 可编程，**API**解除了软硬件之间的耦合关系
- ❑ 特点
  - 硬件资源虚拟化
  - 系统软件平台化
  - 应用软件多样化

# 软件定义

## □ 延伸

- 解耦功能，通过**API**编程，管理下层资源，为上层提供服务。
- 这不就是操作系统做的事情嘛！
- 解耦功能，通过**ISA**编程，管理下层资源，为上层提供服务。
- 这不就是指令集做的事情嘛！
- 因特网的分层模型本身就是软件定义的好例子！

# 软件定义

## □ 延伸

- 软件定义的计算机：当前计算机（相对于大型机）
- 软件定义的手机：智能手机（相对于功能机）
- 软件定义的计算机：虚拟机
- 软件定义的计算机及环境：容器

## □ 趋势

- 将计算、存储、网络等**IT**基础资源，抽象为系统软件对虚拟资源进行管理和调用，在此基础上，用户可编写应用程序，访问资源所提供的服务，进而改变资源的行为，满足应用对资源的多样需求。

# 思考题

- 如何看待
  - 滴滴专车这样的出行服务公司
- 如果将来有
  - 软件定义的课程
  - 软件定义的公司

