

第4章 网络层：数据平面

目的:

- ❑ 了解网络层服务背后的原则，专注于数据平面：
 - 网络层服务模型
 - 转发 vs 路由
 - 路由器的工作原理
 - 网际协议
- ❑ 因特网的网络层的实现

第4章 主要内容

□ 4.1 网络层概述

- 转发和路由选择：数据平面和控制平面
- 网络服务模型

□ 4.2 路由器工作原理

- 输入端口和基于目的地的转发
- 交换
- 输出端口处理
- 何处出现排队
- 分组调度

□ 4.3 网际协议：IPv4、寻址、IPv6

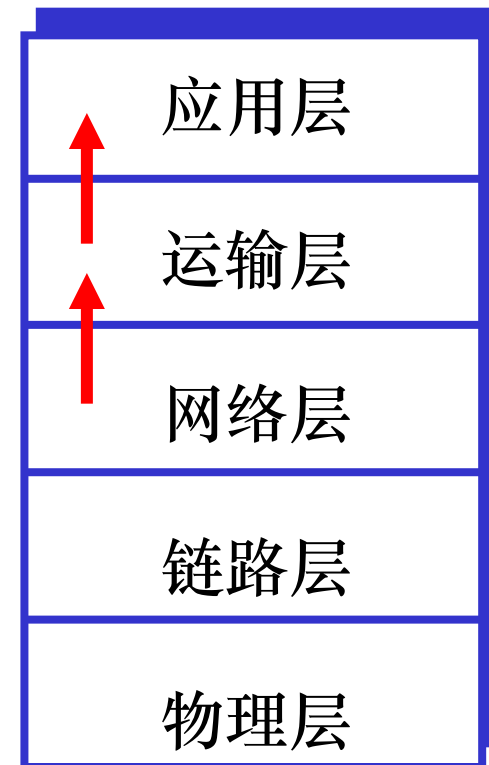
- IPv4数据报格式
- IPv4数据报分片
- IPv4编址
- 网络地址转换(NAT)
- IPv6

4.1 网络层概述

□ 运输层和网络层提供服务的区别

- 运输层为主机上的**两个进程**之间提供通信服务；
- 运输层只在网络的**主机**中出现；
- 网络层提供的**主机到主机**的通信服务；
- 网络层在网络的**主机和路由器**中出现。

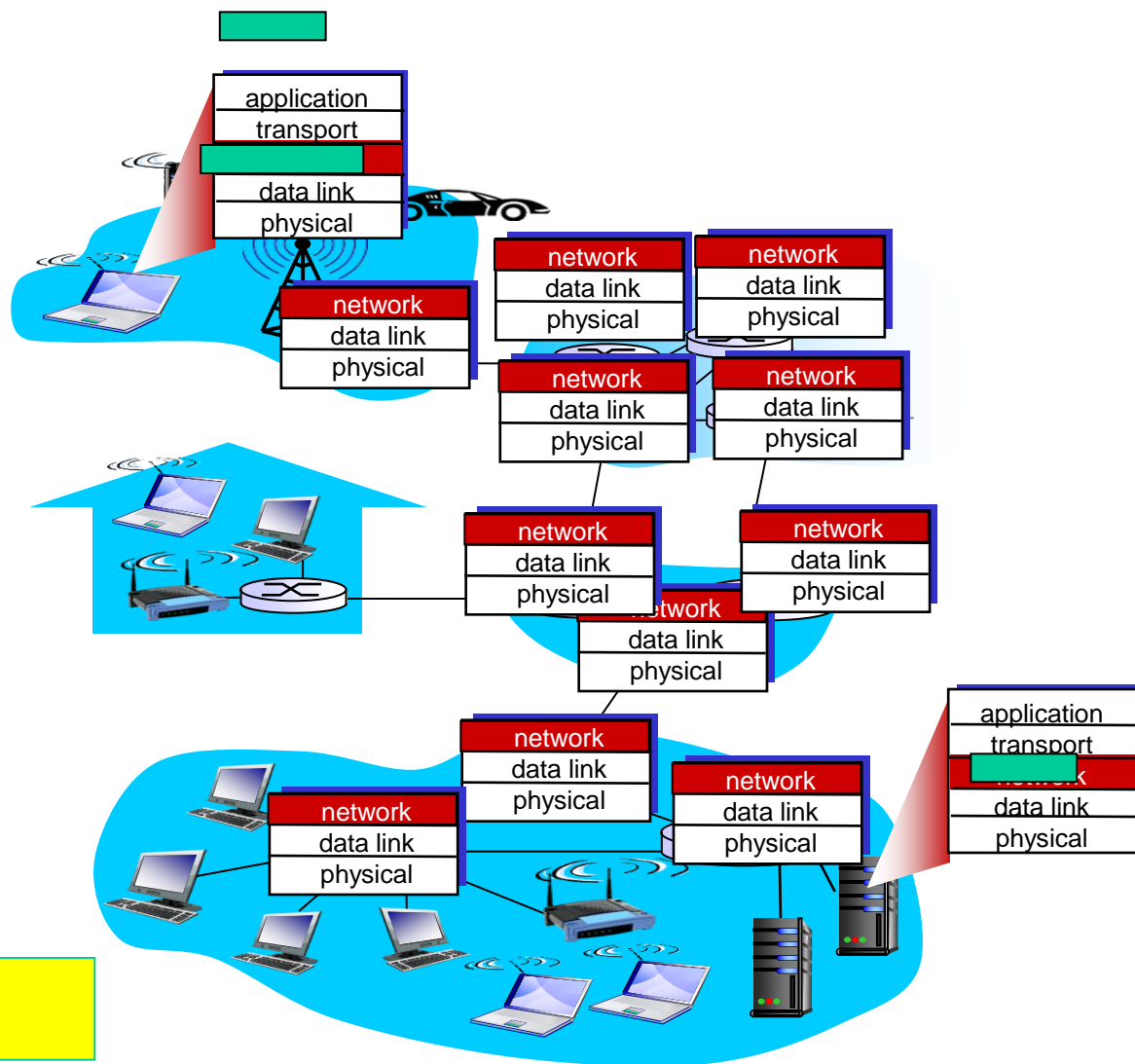
本章主要讨论网络层如何实现主机到主机的通信服务



4.1.1 转发和路由选择

- ❑ 传输报文段(segment)从发送主机到接收主机
- ❑ 发送方封装报文段(segments)为数据报(datagrams)
- ❑ 接收方递交数据报给运输层
- ❑ 网络层协议在每个主机(host)和路由器(router)上实现
- ❑ 路由器检查所有IP数据报的头部字段, 根据目的IP地址转发数据报

通常路由器只实现下三层部分



路由器的两个核心功能

网络层功能:

- 转发(forwarding): 将分组(packet)从路由器的输入端口转移到输出端口
- 路由(routing): 确定分组从发送方到接收方的路由(route)
 - 路由算法

类比: *taking a trip*

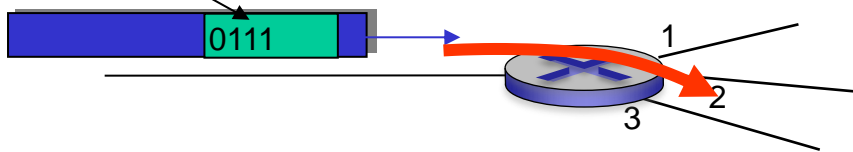
- 路由器的转发过程**类比于**通过立交桥；单个路由器
- 路由过程**类比于**从电子科大沙河校区开车到清水河校区的旅程，需要通过若干道路和立交桥（确定怎么走的线路）；网络中的所有路由器，集体经选路协议交互，决定分组从源到目的地的路径。

数据平面和控制平面

数据平面

- 本地、每个路由器的功能
- 将数据报从路由器的输入端口转发到路由器的输出端口
- forwarding功能

values in arriving
packet header

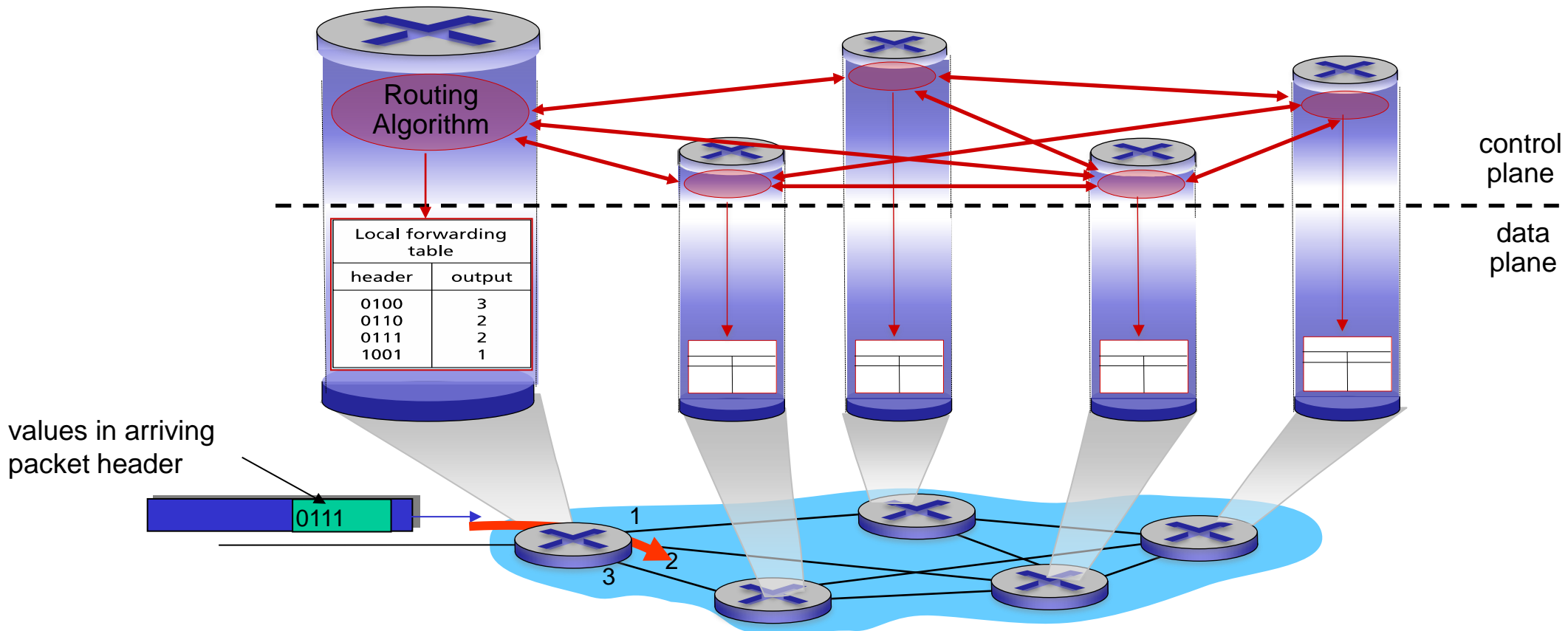


控制平面

- 网络范围的逻辑
- 确定数据报从发送主机到接收主机的端到端路径上通过的路由器是哪些
- 两个控制平面的方案：
 - 传统路由算法：在路由器中实现
 - 软件定义网络(software-defined networking, SDN)：在远程的服务器上实现

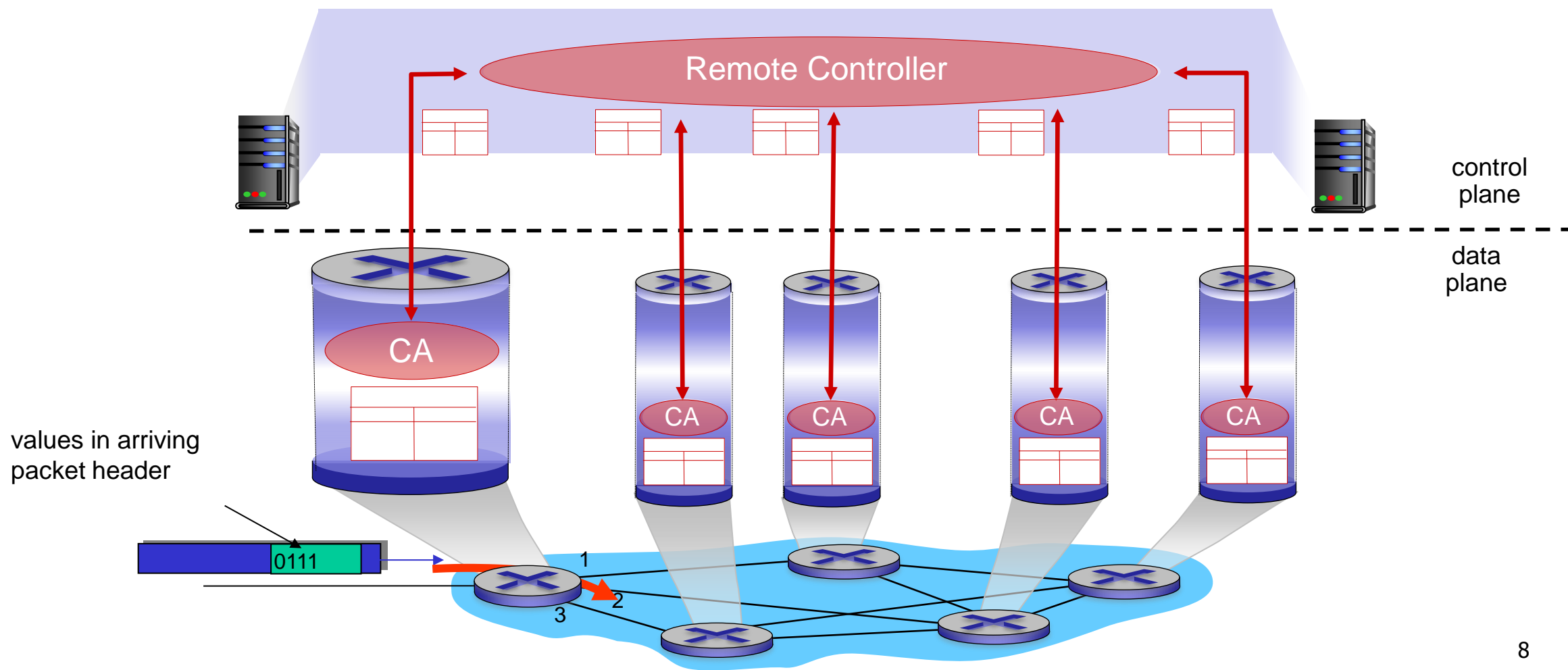
每个路由器的控制平面

每个路由器中的各个路由算法组件在控制平面中交互



逻辑上集中式的控制平面

一个特别的(通常是远程的)控制器与本地控制代理(CA)交互

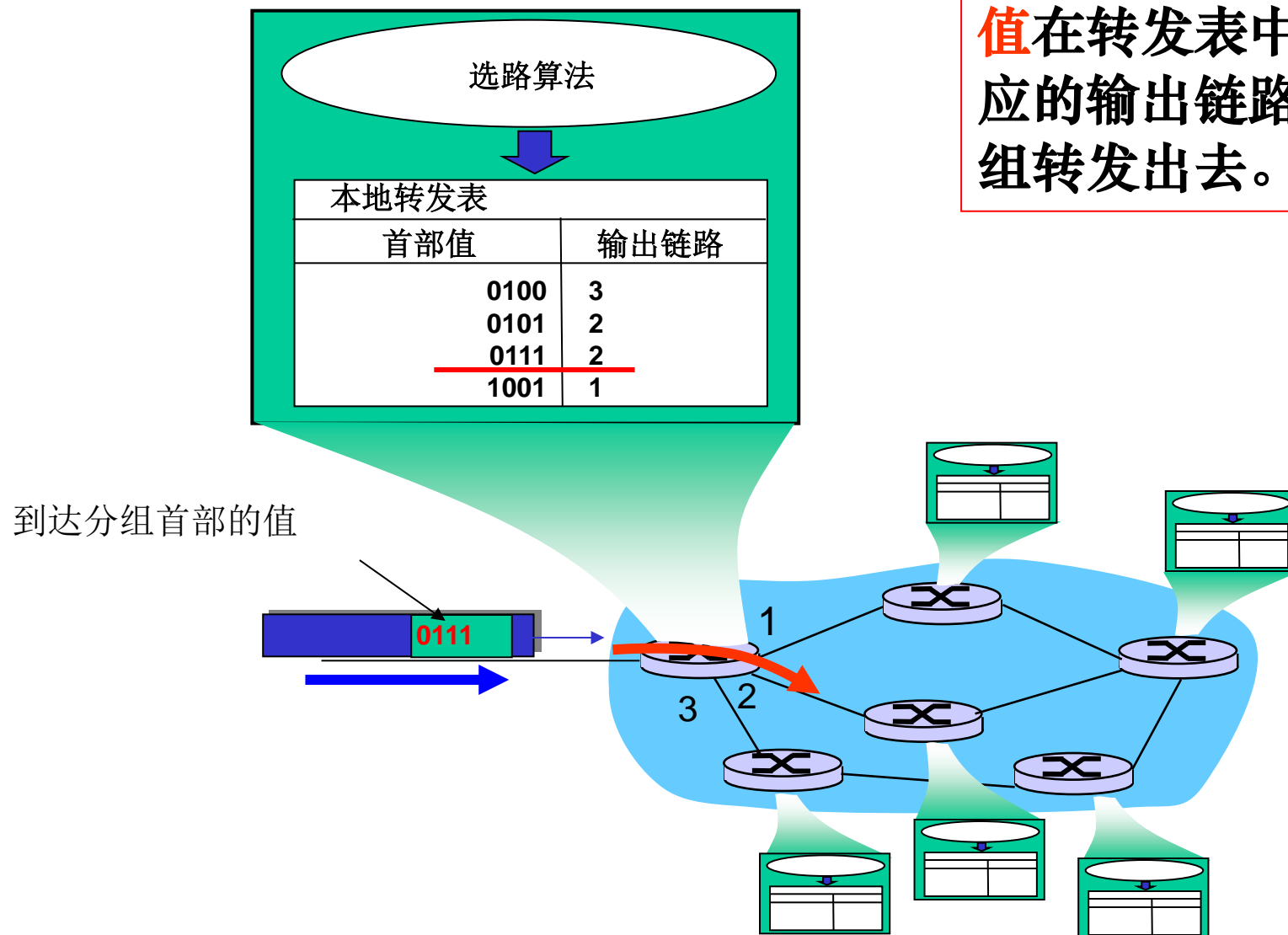


路由器如何转发分组

- **转发表**：每台路由器有一张转发表
 - **分组首部**(目的地址或某个连接标识)和**相应输出链路的对照表**;
 - 转发表的**内容由选路算法决定**(集中式或分布式)。
- **分组交换机(packet switcher)**：一台分组交换设备，根据**分组首部值**，从输入链路接口到输出链路接口传送分组。
 - **链路层交换机**：根据**链路层字段值**作转发决定的分组交换机。
 - **路由器**：根据**网络层字段值**作转发决定的分组交换机。

转发方法

路由器根据到达分组的首部**值**在转发表中查询，找到相应的输出链路接口，并将分组转发出去。



讨论： 是否在网络层建立连接

- 运输层连接：如TCP协议，在数据实际传输之前，需要发送方和接收方经过三次握手建立所需的**状态信息**。
- 网络连接：指网络层数据分组开始传输前，在所选择的**从源到目的地路径上的各路由器之间相互握手**，建立连接状态。
- 在网络层建立连接的例子：ATM、帧中继、MPLS
- **因特网的网络层不建立连接！路由器上不需要维护状态**

4.1.2 网络服务模型

问题:什么样的服务模型可以用于将数据报从发送方传输到接收方?

- 网络层可能提供的服务
 - **确保交付:** 确保分组到达目的地。
 - **具有时延上界的确保交付:** 主机到主机的时延。
 - **有序分组交付:** 按发送顺序到达。
 - **确保最小带宽:** 当发送主机以低于特定比特率的速率发送比特, 分组不会丢失, 在一定时延到达。
 - **确保最大时延抖动:** 发送方发送两个连续分组的时间间隔与接收到的间隔相同。
- 因特网的网络层提供的服务
 - 单一服务, 即尽力而为服务(best-effort service)。
 - 分组间的定时不能被保证;
 - 分组的接收顺序与发送顺序不一定相同;
 - 传送的分组不能保证最终交付, 即网络可能未向目的地交付分组。

其他网络服务模型

- ATM网络体系结构提供了确保按序时延、有界时延和确保最小带宽。
- 因特网的基本尽力而为服务模型与适当带宽供给相结合已经被证明超过“足够好”，能够满足各种网络应用。
- **ATM是基于虚电路网络**，在发送数据前需要在网络层建立连接(路径上的路由器之间)；[虚电路可参考本书第6版]
- **因特网是基于数据报网络**，在发送数据前不需要在网络层建立连接。

虚电路网络

- 数据报网络提供网络层的 **无连接** 服务
- 虚电路网络提供网络层的 **连接** 服务
- 任何网络中的网络层只提供两种服务之一，不会同时提供
 - **虚电路网络**：提供连接服务。
 - **数据报网络**：提供无连接服务。
- **传输层**的面向连接服务是在网络边缘的**端系统**中实现
- 网络层的面向连接服务是在**端系统及路由器**中实现。

虚电路(Virtual Circuits)

“源主机-目的主机路径的行为类似于电话网络的行为”

- 性能上类似
- 沿着源-目的路径的网络行为类似

- 在数据传输之前，需要为每个呼叫建立连接
- 每个分组携带**VC标识符**(不是目的主机地址)
- 位于“源-目的路径”上的**每个路由器会维护经过它的每条连接的“状态”**
- 链路和路由器的资源(带宽、缓存)可以被分配给VC(专用资源)

虚电路的实现

- 一条虚电路(VC)包括:
 - 一条从发送方到接收方的**路径**;
 - VC号, 沿路径的每条链路都有一个**VC号**;
 - 沿路径的每个路由器上都有**转发表条目**;
- 属于同一个VC的分组携带相同的VC号(而不是目的地址)
- 在每条链路上VC号可以被修改
 - 新的VC号来自于转发表

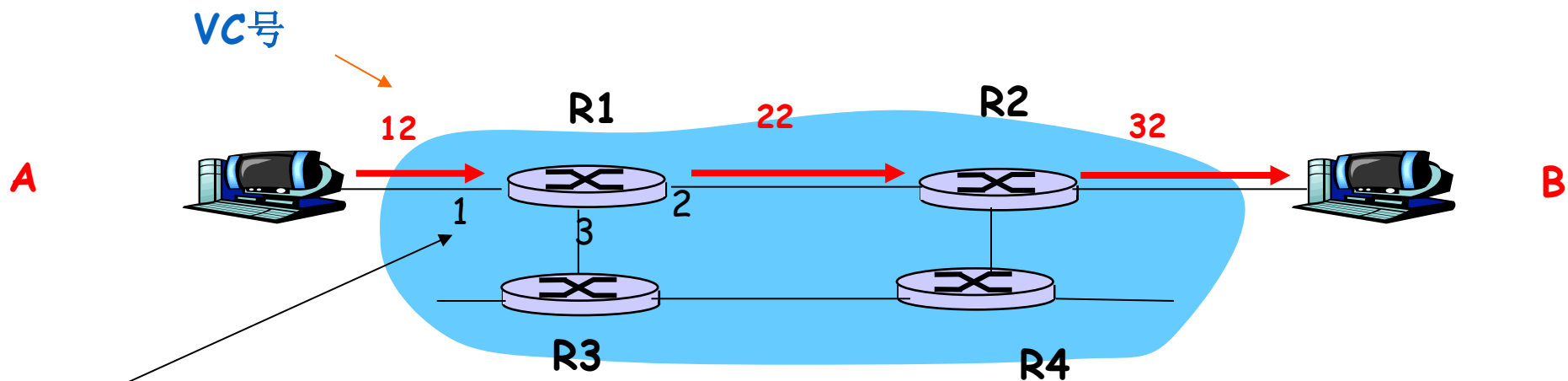
主机A与主机B通信过程

□ 主机A与主机B之间创建一条VC:

路径为A → R1 → R2 → B，3条链路分配VC号12、22和32

□ 传输时，**分组VC号变化过程**:

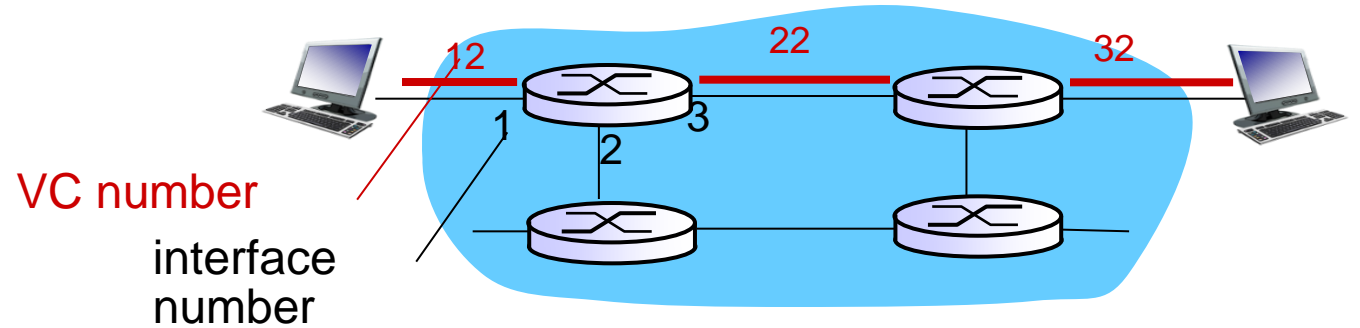
12 22 32
A → R1 → R2 → B



接口号: 路由器连接链路的接口号码

VC转发表

*forwarding table in
northwest router:*

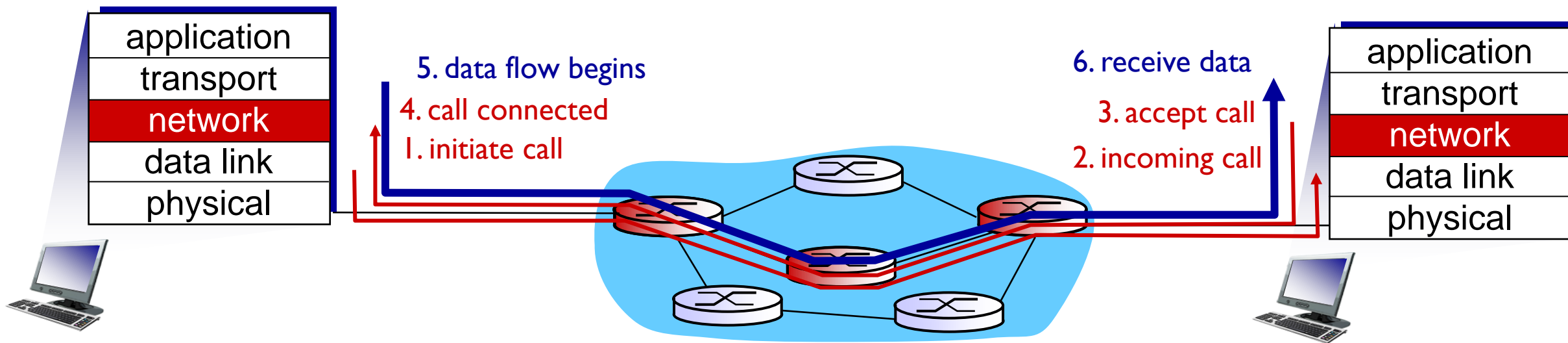


Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	3	22
2	63	1	18
3	7	2	17
1	97	3	87
...

VC routers maintain connection state information!

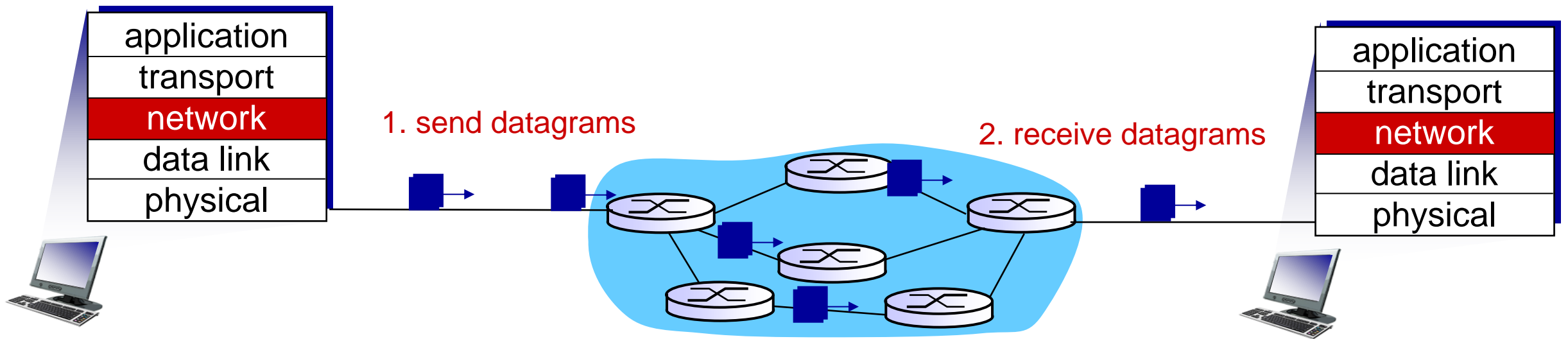
虚电路：信令协议

- 用于建立，维护，拆除虚电路
- 用于ATM, frame-relay, X.25网络
- 不能用于目前的因特网



数据报网络

- ❑ 在网络层无呼叫的过程
- ❑ 路由器： 不需要维护端到端连接的状态
- ❑ 没有网络等级的“连接”的概念
- ❑ 使用目的主机的地址进行分组转发

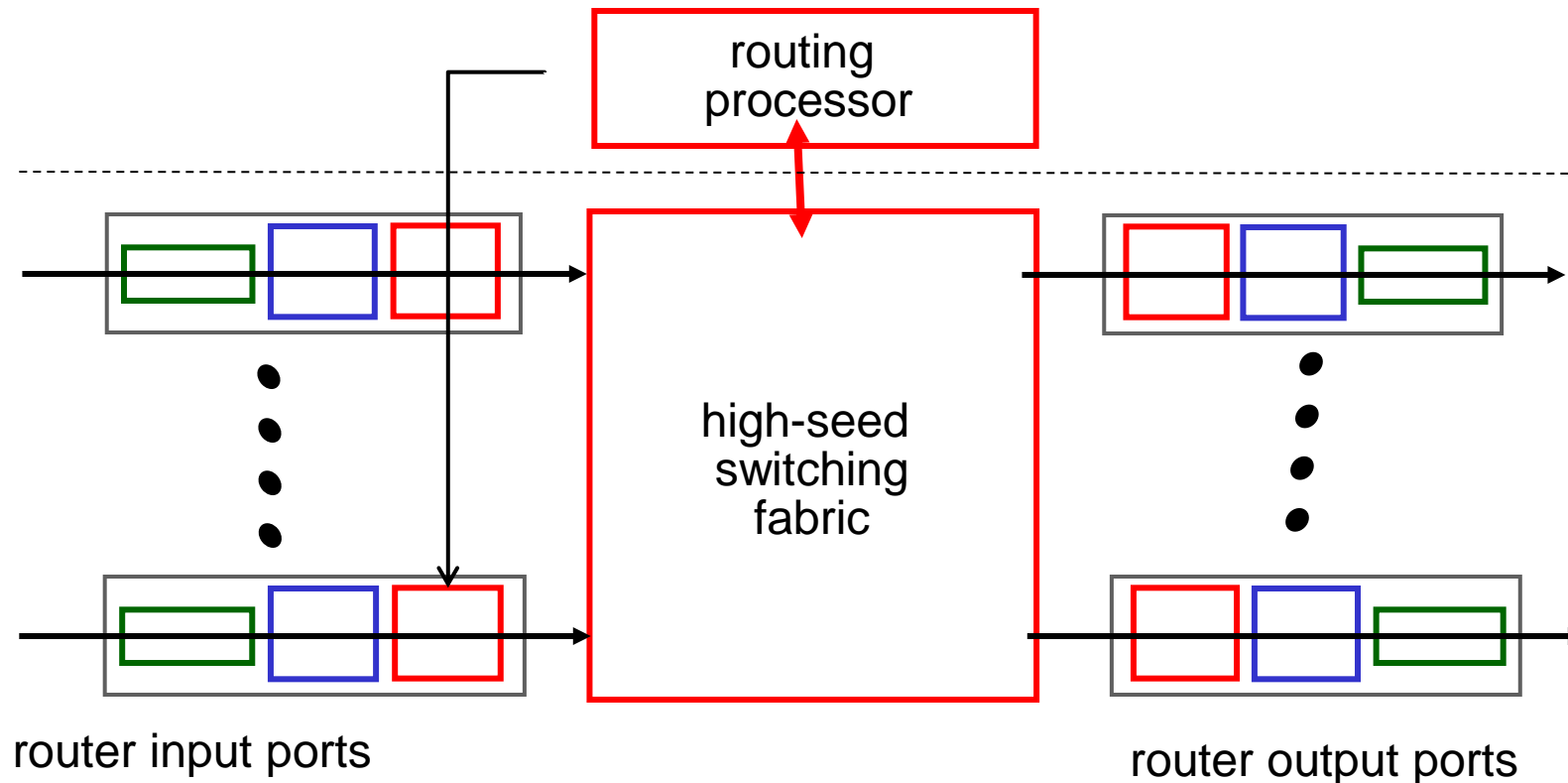


4.2 路由器工作原理

□ 路由器的两个核心功能：

- 运行路由算法/协议(RIP, OSPF, BGP)
- 将分组从路由器的输入链路传送到正确的输出链路。

□ 路由器的体系结构：

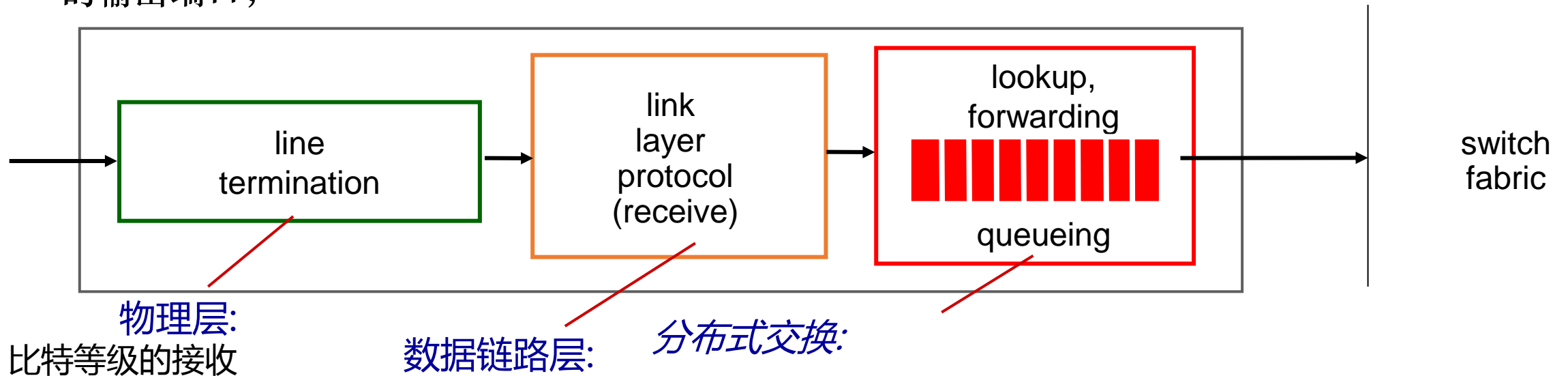


*routing, management
control plane* (software)
operates in millisecond
time frame

forwarding data plane
(hardware) operates in
nanosecond timeframe

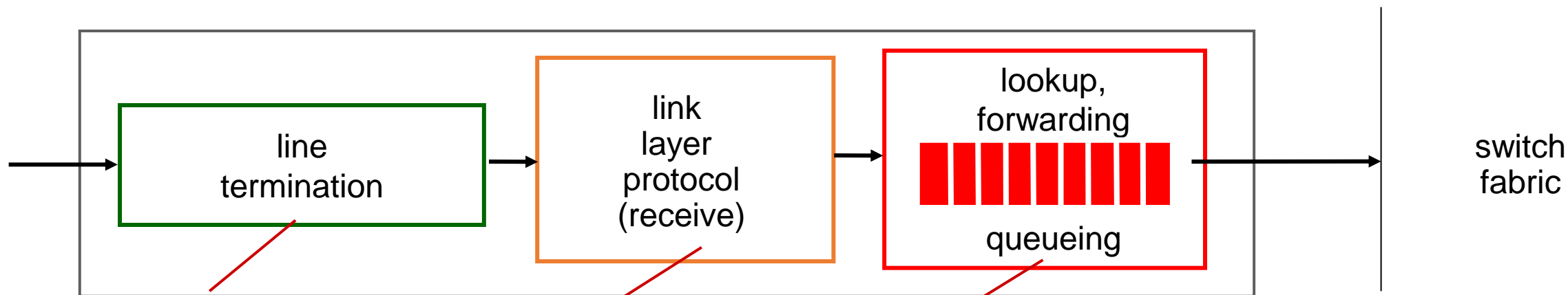
4.2.1 输入端口处理和基于目的地的转发

- ❑ 第一个**线路端接**模块：将一条**物理链路端接到路由器的物理层**；
- ❑ 第二个**数据链路处理**模块：实现路由器的**数据链路层功能**；
- ❑ 第三个**查找与转发**模块：实现**查找与转发功能**，以便分组通过路由器交换结构转发到适当的输出端口；



- 根据数据报的头部字段的值，使用输入端口内存中缓存的转发表查找输出端口
- 目标：以线速完成输入端口处理
- 排队：当数据报到达的速率快于转发到交换结构的速率

4.2.1 输入端口处理和基于目的地转发



物理层:
比特等级的接收

数据链路层:
解封装进行分析
e.g., Ethernet

分布式交换:

- 根据数据报的头部字段的值，使用输入端口内存中缓存的转发表查找输出端口
- 基于目的IP地址的转发（传统）
- 通用转发：基于任意设置的头部字段的值的组合

基于目的IP地址的转发

- 根据分组的目的IP地址在转发表中查询，找到相应的输出链路接口，并将分组转发出去
- 转发表：每台路由器维护一张路由表，目的地址与链路接口的映射表。
 - 转发表中的表项数与地址位数有关，如果每个可能的IP地址对应一项，转发表的条目可以达到40亿条(目的地址32位)

转发表

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

问题: 如果地址范围分配没有很合理, 会发生什么问题?

最长前缀匹配规则

longest prefix matching

对于给定的目的地址，使用**最长**地址前缀匹配来完成输出端口的查找

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

例子:

DA: 11001000 00010111 00010**110** 10100001

which interface?

DA: 11001000 00010111 00011000 **10101010**

which interface?

路由器查表方法

用目的地址前缀与转发表的前缀匹配：

❑ 存在匹配：向对应链路转发。

如，目的地址

11001000 00010111 00010110 10100001 →

0#

❑ 不存在匹配：选择“其他”项对应的链路转发。

❑ 存在多个匹配：使用最长前缀匹配规则，即向与最长前缀匹配的链路接口转发分组。

如，目的地址

11001000 00010111 00011000 10101010 →

1#

前缀匹配

链路接口

11001000 00010111 00010

0

11001000 00010111 00011000

1

11001000 00010111 00011

2

其他

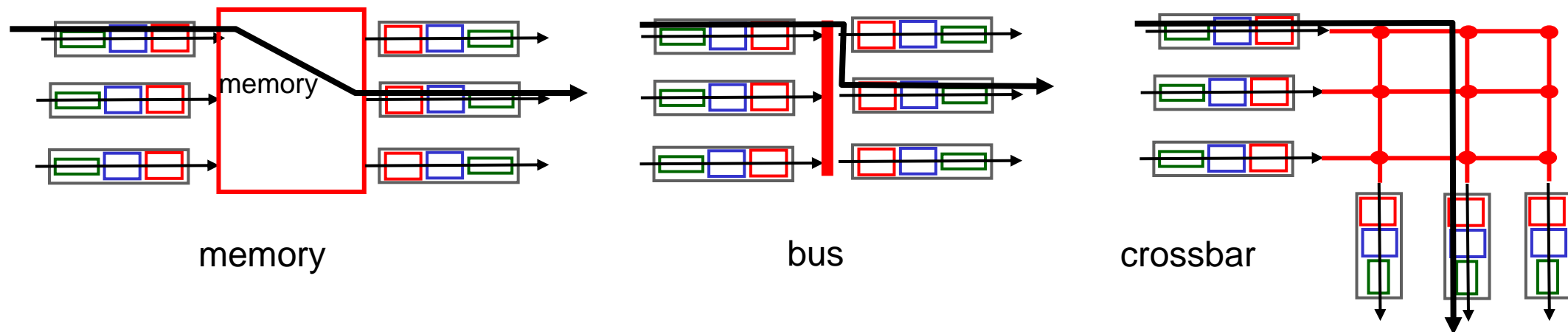
3

4.2.1 输入端口处理和基于目的地转发

- 确定将一个到达的分组通过交换结构转发给哪个输出端口。通过查找转发表实现，这里的转发表是存储在输入端口的内存中。
- 分布式交换：
 - 选路处理器计算转发表，给每个输入端口存放一份转发表拷贝。
 - 在每个输入端口本地做出交换决策，无须激活中央选路处理器。
 - 可避免在路由器中某个单点产生转发处理瓶颈。
 - 目的：以线速完成输入端口的处理
 - 排队：如果数据报到达输入端口的速度快于输入端口将数据报转发到交换结构的速度，就会发生排队

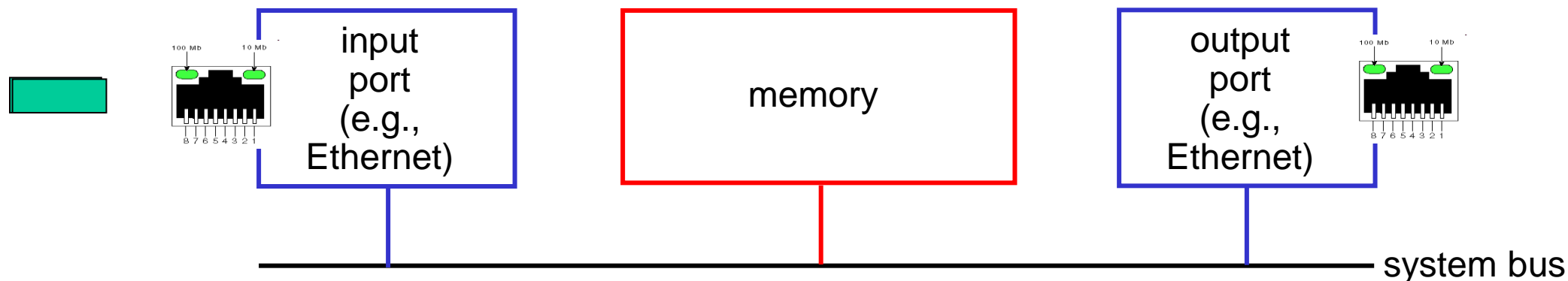
4.2.2 交换结构

- ❑ 将分组从输入端口缓存交换（转发）到恰当的输出口缓存中
- ❑ 交换速率：分组从输入端口到输出端口的速率
 - 通常用输入/输出线速的倍数来表示
 - N inputs: switching rate N times line rate desirable
- ❑ 三种类型的交换结构



经内存交换

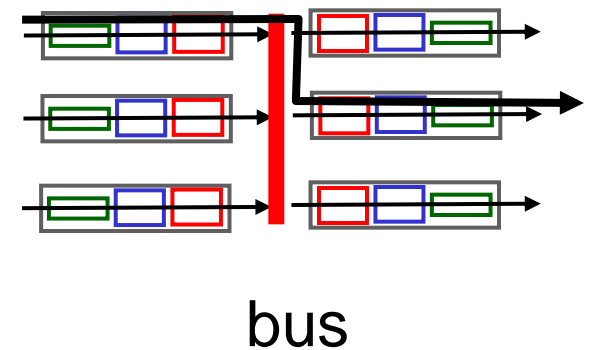
- 早期用计算机作为路由器时采用的结构(第一代)
 - 输入端口与输出端口之间的交换由CPU(选路处理器)控制完成;
 - 输入端口与输出端口类似I/O设备:
 - 当分组到达输入端口时,通过中断向选路处理器发出信号,将分组拷贝到处理器内存中;
 - 选路处理器根据分组中的目的地址查表找出适当的输出端口,将该分组拷贝到输出端口的缓存中。



交换速度受总线带宽的速度限制 (每个分组穿过两次总线)

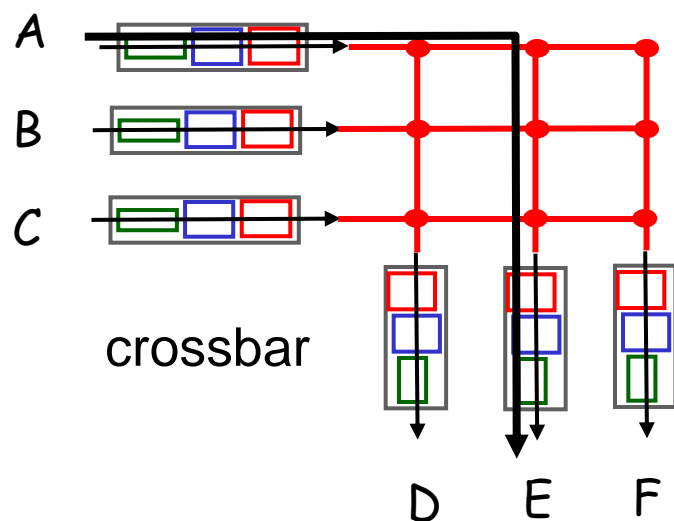
经总线交换

- 输入端口通过一条共享总线将分组直接传送到输出端口，不需要选路处理器的干预
- 每次只能有一个分组通过总线传送，分组可以被总线上的每个端口接收，只有标签正确的端口将其存储；
- 分组到达一个输入端口时，若总线正忙，会被暂时阻塞，在输入端口排队
- 路由器交换带宽受总线速率限制。



经互联网络交换

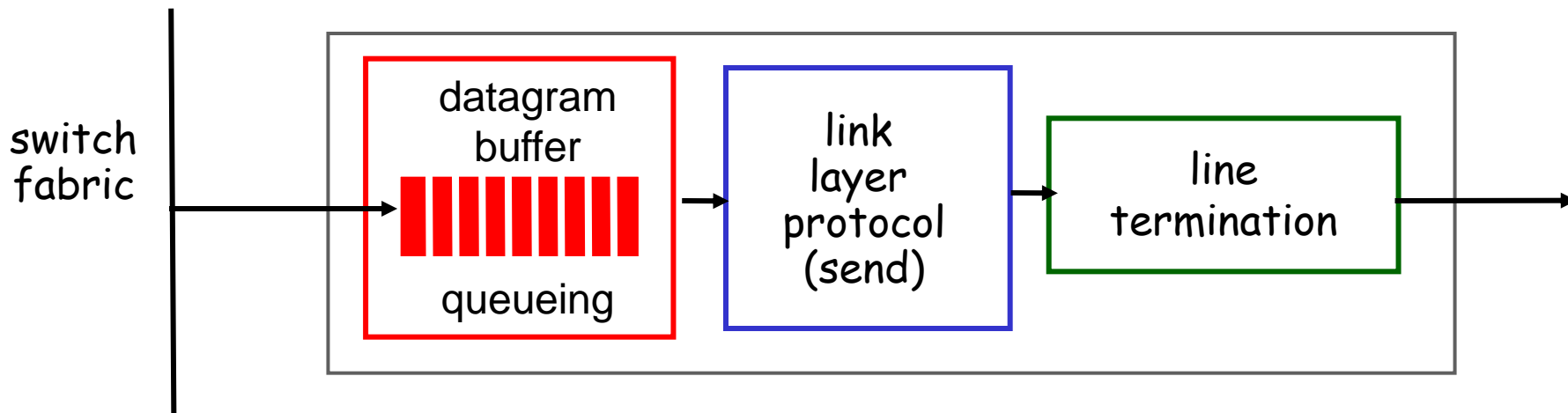
- 纵横式(crossbar)交换机：由 $2n$ 条总线组成， n 个输入端口与 n 个输出端口连接。
- 到达输入端口的分组沿水平总线穿行，直至与所希望的输出端口的垂直总线交叉点：
 - 若该条垂直总线空闲，则分组被传送到输出端口；
 - 否则，该到达的分组被阻塞，必须在输入端口排队。



从A到E端口转发分组，从B到D端口转发分组可以同时进行，因此速率高于经总线的交换结构

4.2.3 输出端口处理

- ❑ 取出存放在输出端口内存中的分组，并将其传输到输出链路上。
- ❑ 当交换结构将分组交付给输出端口的速率超过输出链路速率，就需要排队与缓存管理功能。当输出端口的缓冲区溢出时，就会出现延时和丢包。

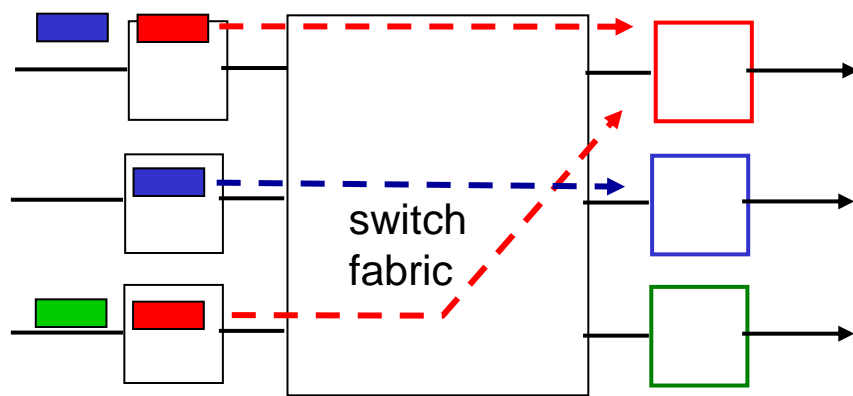


4.2.4 何时出现排队

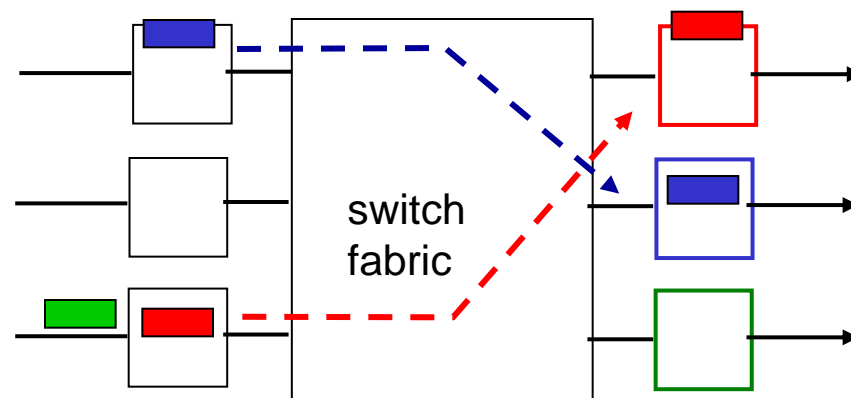
□ 输入端口排队

- 当交换结构的速度慢于输入端口的速度，就会在输入端口的缓冲区发生排队
- 会导致排队延时和由于输入缓冲区溢出导致的丢包！

□ 线头阻塞（Head-of-the-Line (HOL) blocking）：在队列前面的被阻塞的数据报会阻止队列中的其他数据报被转发。



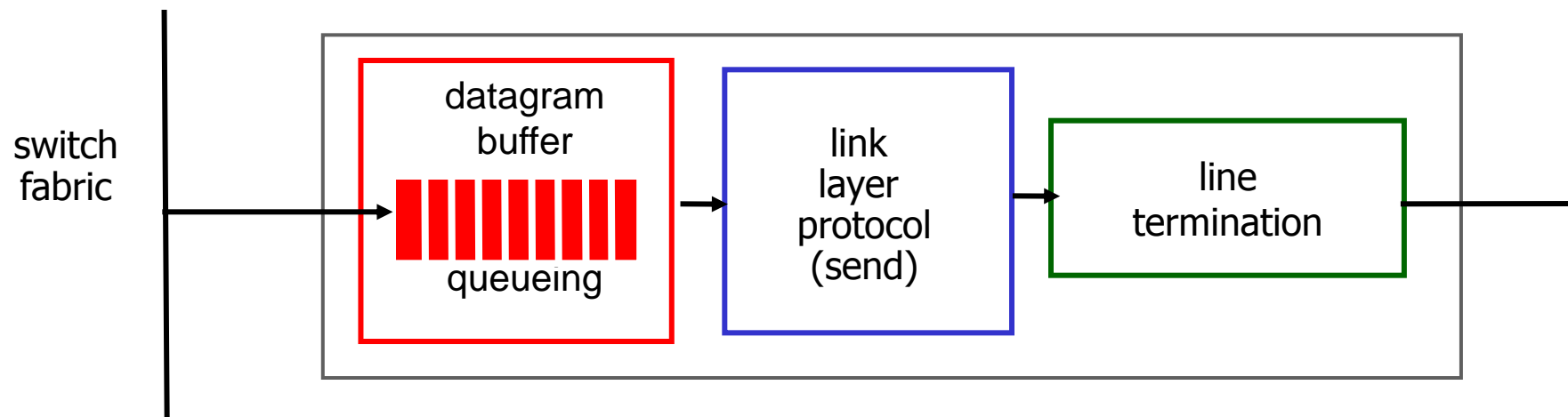
输出端口竞争：
只有一个红色的数据报可以被转移
左下角的红色分组被阻塞



一个分组时间之后：绿色分组经历了HOL阻塞

输出端口排队

This slide is HUGEY important!



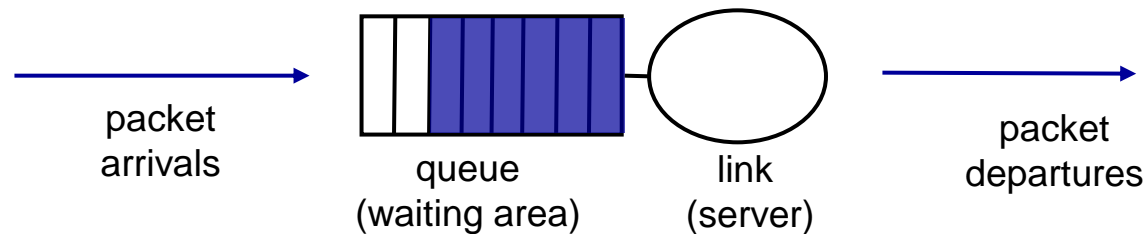
Datagram (packets) can be lost due to congestion,
lack of buffers

- ❑ 当数据报从交换结构到达输出端口的速率超过输出端口的传输速率，则需要缓存数据报
- ❑ 调度机制用来在队列中选择传输的数据报

Priority scheduling – who gets best performance,
network neutrality

4.2.5 分组调度

- 调度：选择该链路上下一个发送的分组
- FIFO (first in first out) 调度：根据到达队列的顺序发送
 - real-world example?
 - *discard policy*: if packet arrives to full queue: who to discard?
 - *tail drop*: drop arriving packet
 - *priority*: drop/remove on priority basis
 - *random*: drop/remove randomly



三种调度策略的实例——FIFO

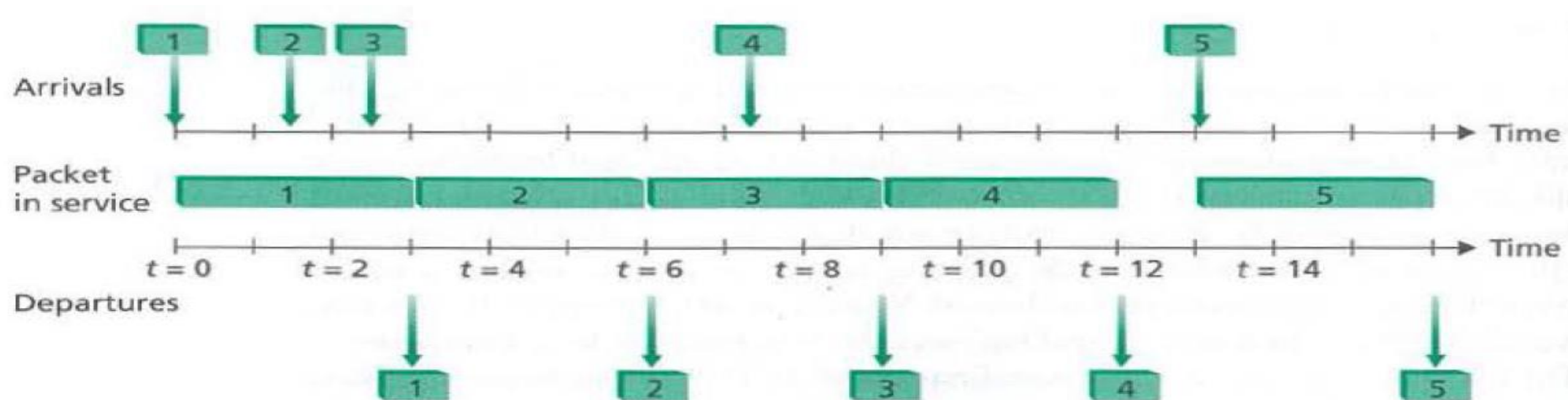


Figure 4.11 ♦ The FIFO queue in operation

三种调度策略的实例——优先级

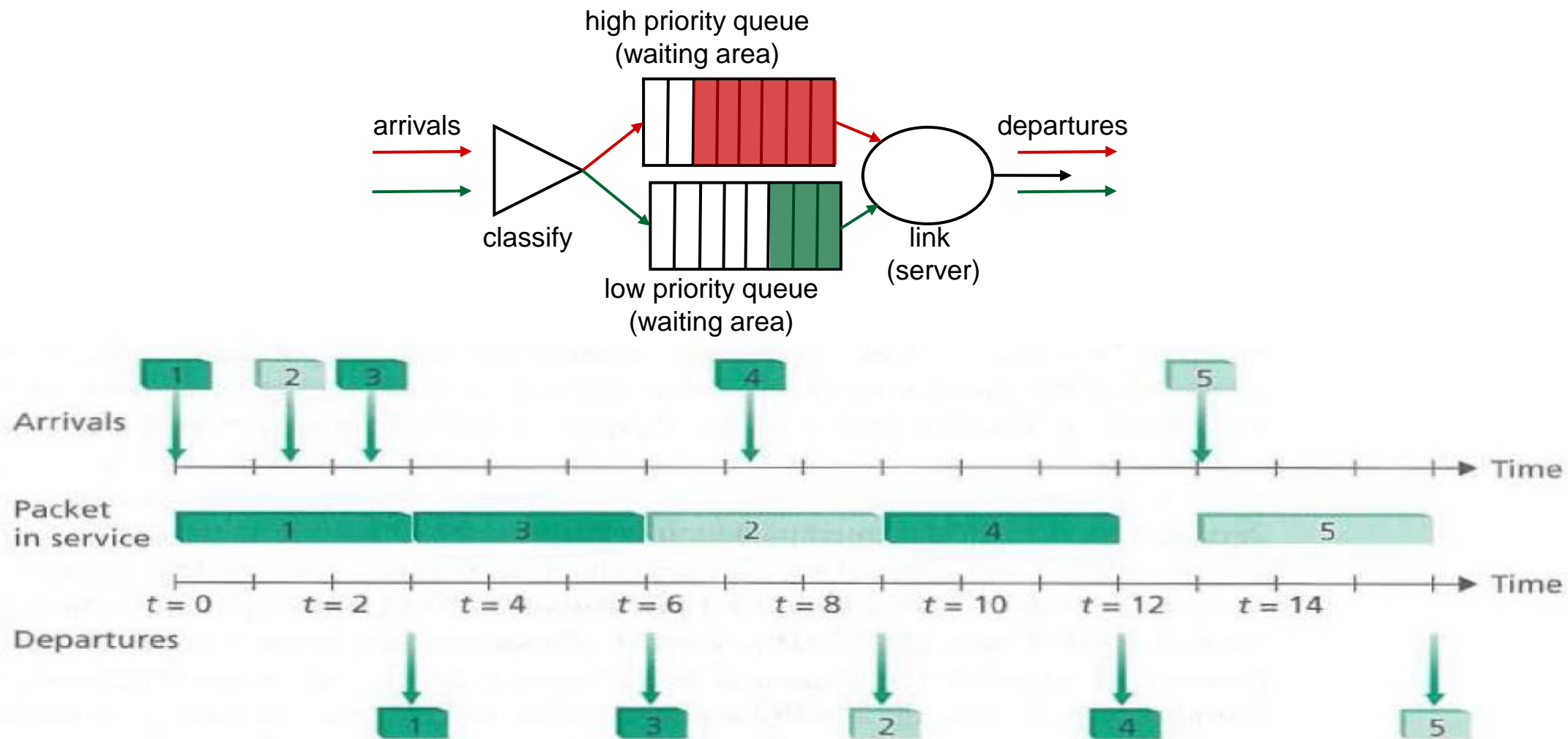


Figure 4.13 ♦ The priority queue in operation

三种调度策略的实例—Round Robin

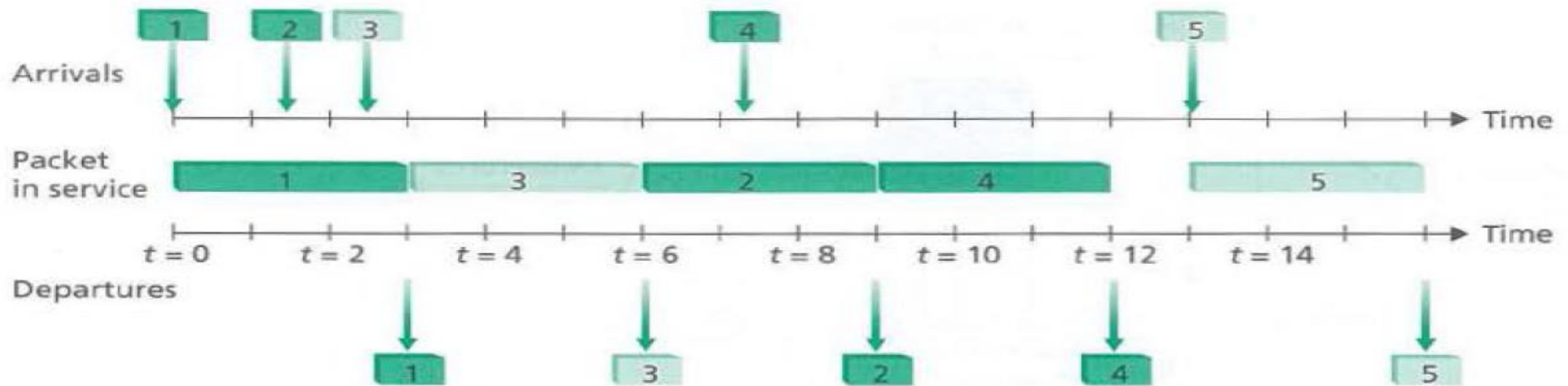
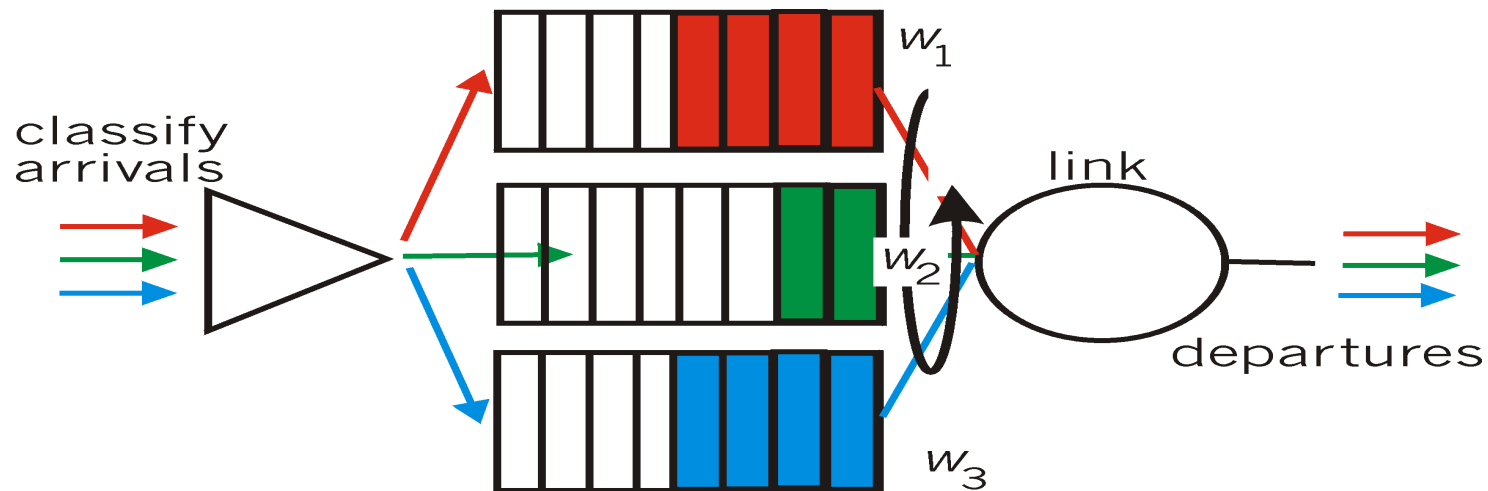


Figure 4.14 ♦ The two-class robin queue in operation

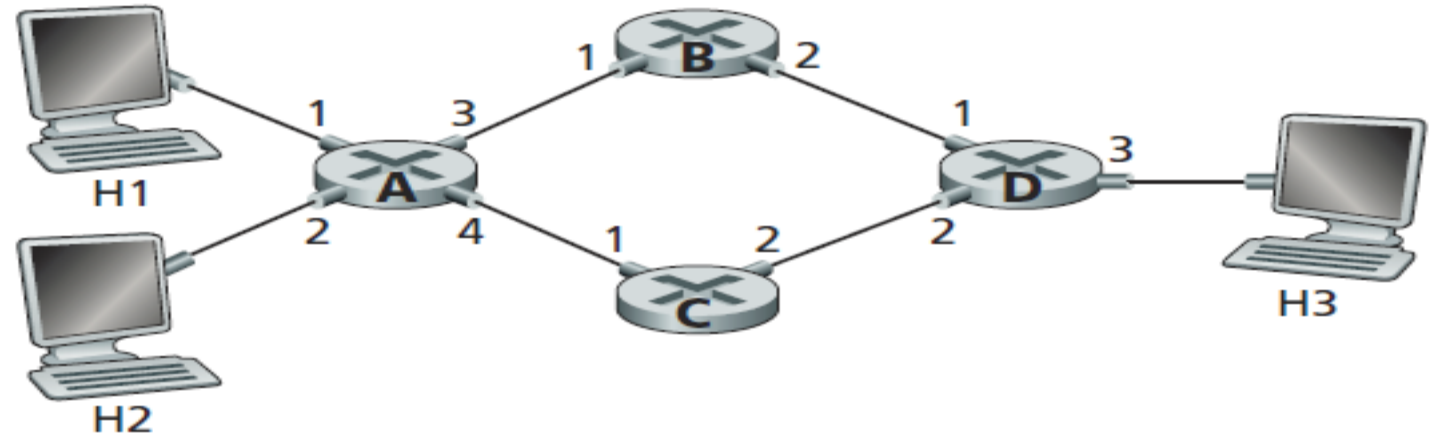
三种调度策略的实例——加权公平队列

Weighted Fair Queuing (WFQ):

- generalized Round Robin
- 在循环中每个队列有自己的权重
- real-world example?



课堂练习



P4. Consider the network below.

- Suppose that this network is a datagram network. Show the forwarding table in router A, such that all traffic destined to host H3 is forwarded through interface 3.
- Suppose that this network is a datagram network. Can you write down a forwarding table in router A, such that all traffic from H1 destined to host H3 is forwarded through interface 3, while all traffic from H2 destined to host H3 is forwarded through interface 4? (Hint: this is a trick question.)
- Now suppose that this network is a virtual circuit network and that there is one ongoing call between H1 and H3, and another ongoing call between H2 and H3. Write down a forwarding table in router A, such that all traffic from H1 destined to host H3 is forwarded through interface 3, while all traffic from H2 destined to host H3 is forwarded through interface 4.
- Assuming the same scenario as (c), write down the forwarding tables in nodes B, C, and D.

课堂练习

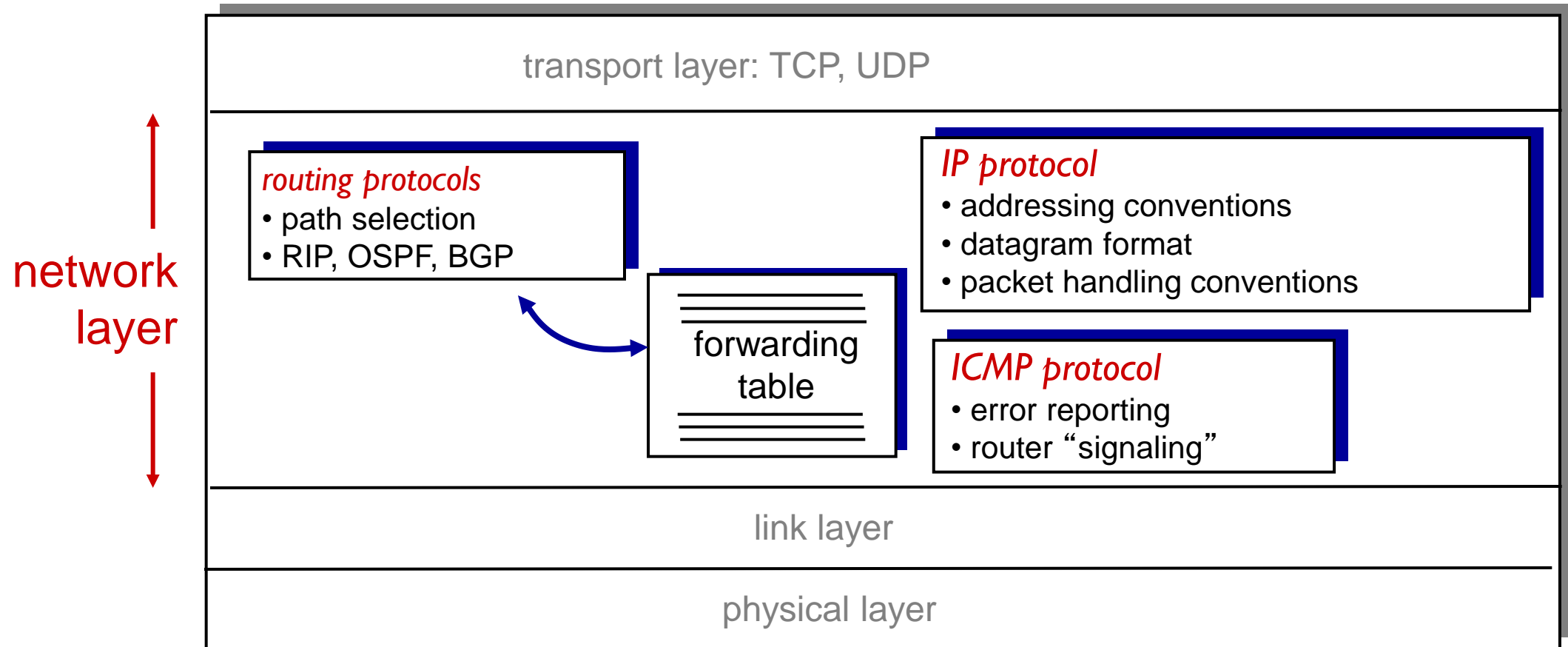
Destination Address Range	Link Interface
11100000 00000000 00000000 00000000 through 11100000 00111111 11111111 11111111	0
11100000 01000000 00000000 00000000 through 11100000 01000000 11111111 11111111	1
11100000 01000001 00000000 00000000 through 11100001 01111111 11111111 11111111	2
otherwise	3

- Provide a forwarding table that has five entries, uses longest prefix matching, and forwards packets to the correct link interfaces.
- Describe how your forwarding table determines the appropriate link interface for datagrams with destination addresses:

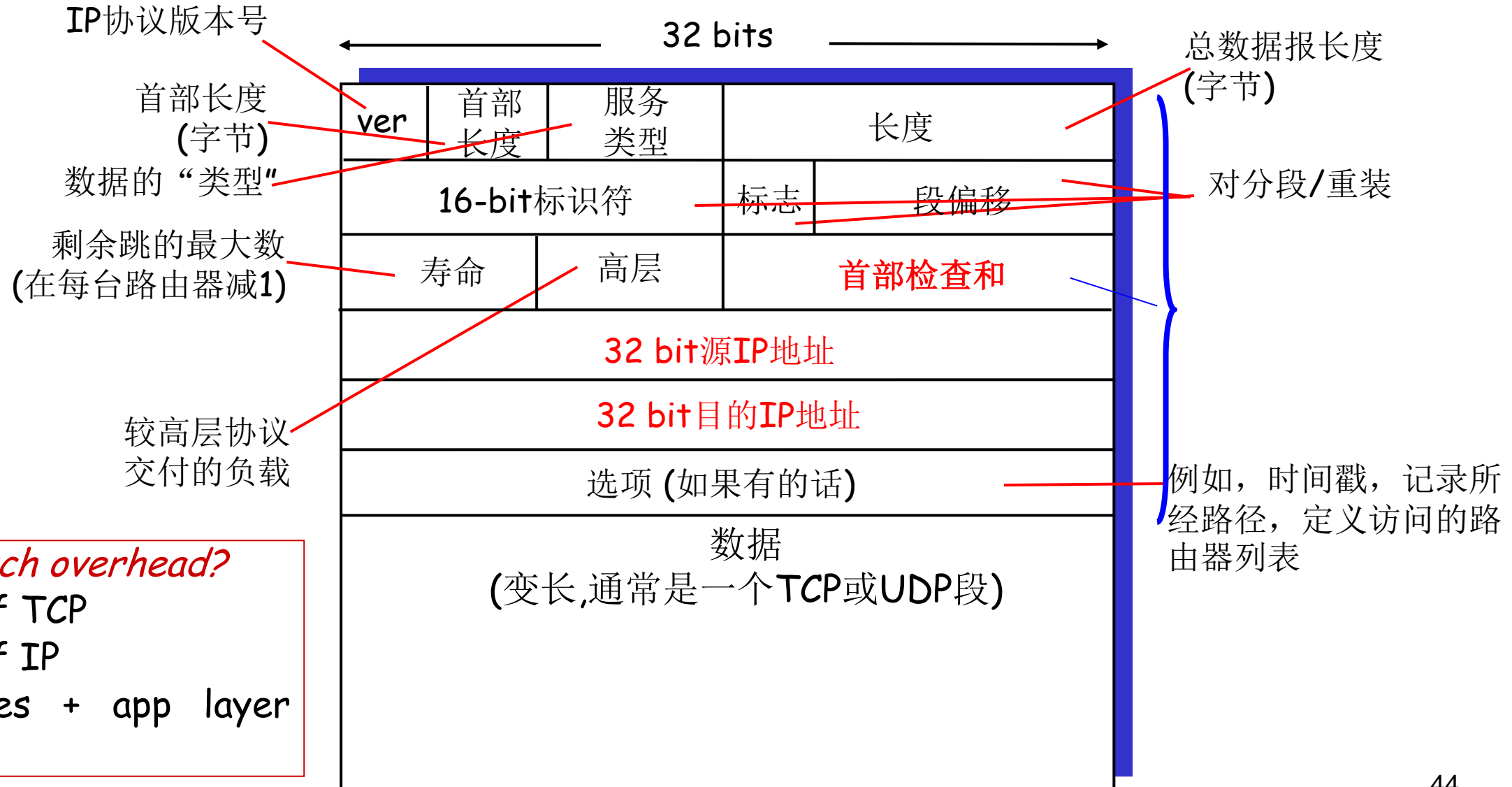
11001000 10010001 01010001 01010101
11100001 01000000 11000011 00111100
11100001 10000000 00010001 01110111

4.3: 网际协议: IPv4、寻址、IPv6

host, router network layer functions:

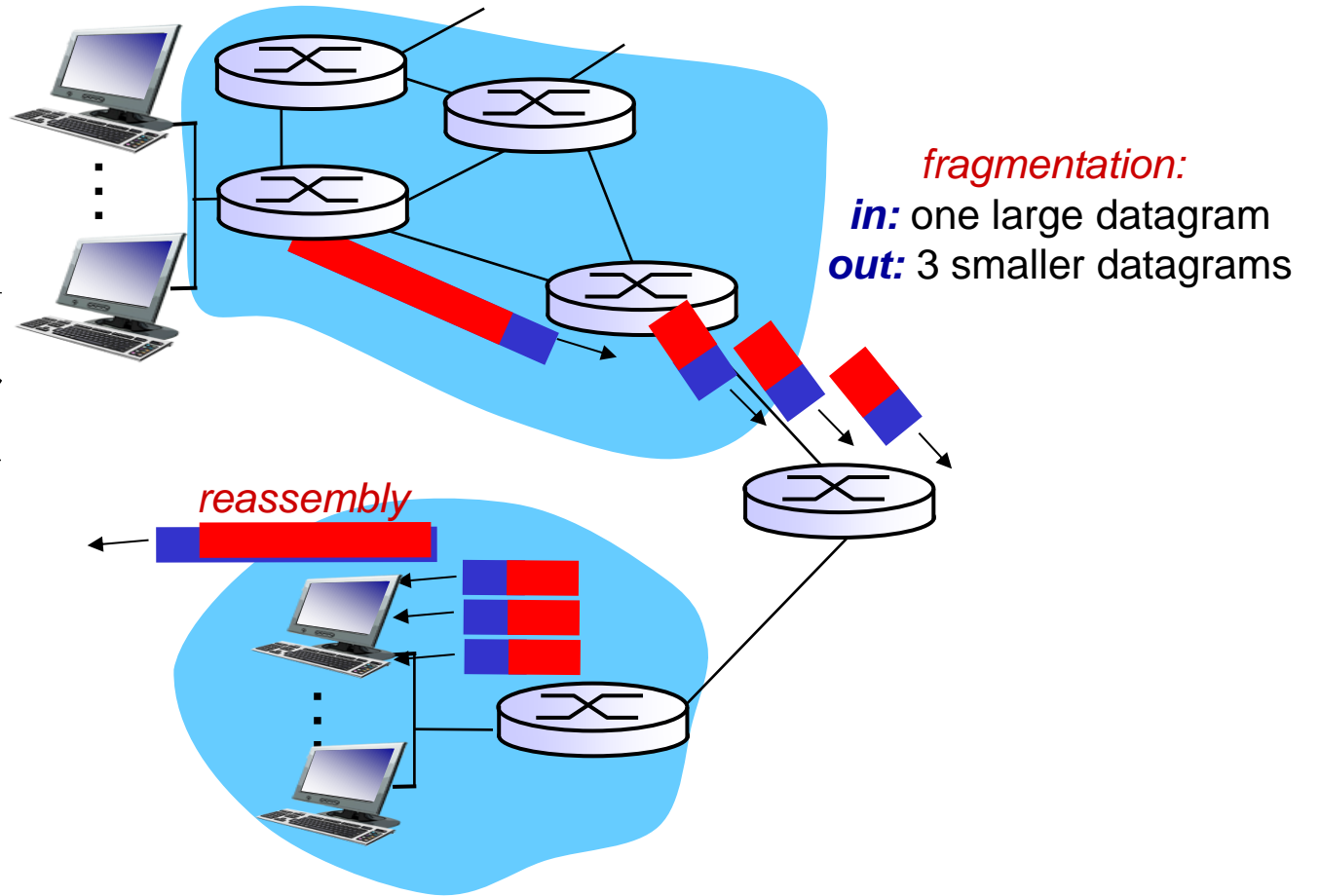


4.3.1 IP 数据报格式 (IPv4)



4.3.2 IPv4数据报分片和重组

- ❑ 每个数据链路有自己的MTU，链路类型不同，MTU的值也不同，这里MTU指的是数据链路帧的数据区的最大字节数
- ❑ 在因特网中，一个大的分组可能在路由器中被分割为几个分片，在最终的目的主机上，将这些分片重新组装成一个大的分组
- ❑ 为了进一步识别出这些分组，需要对分片进行标识



分片的例子

example:

- ❖ 4000 byte datagram
- ❖ MTU = 1500 bytes



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

一个大的数据报分割为若干小的数据报

1480 bytes in
data field

$$\text{offset} = 1480/8 = 185$$

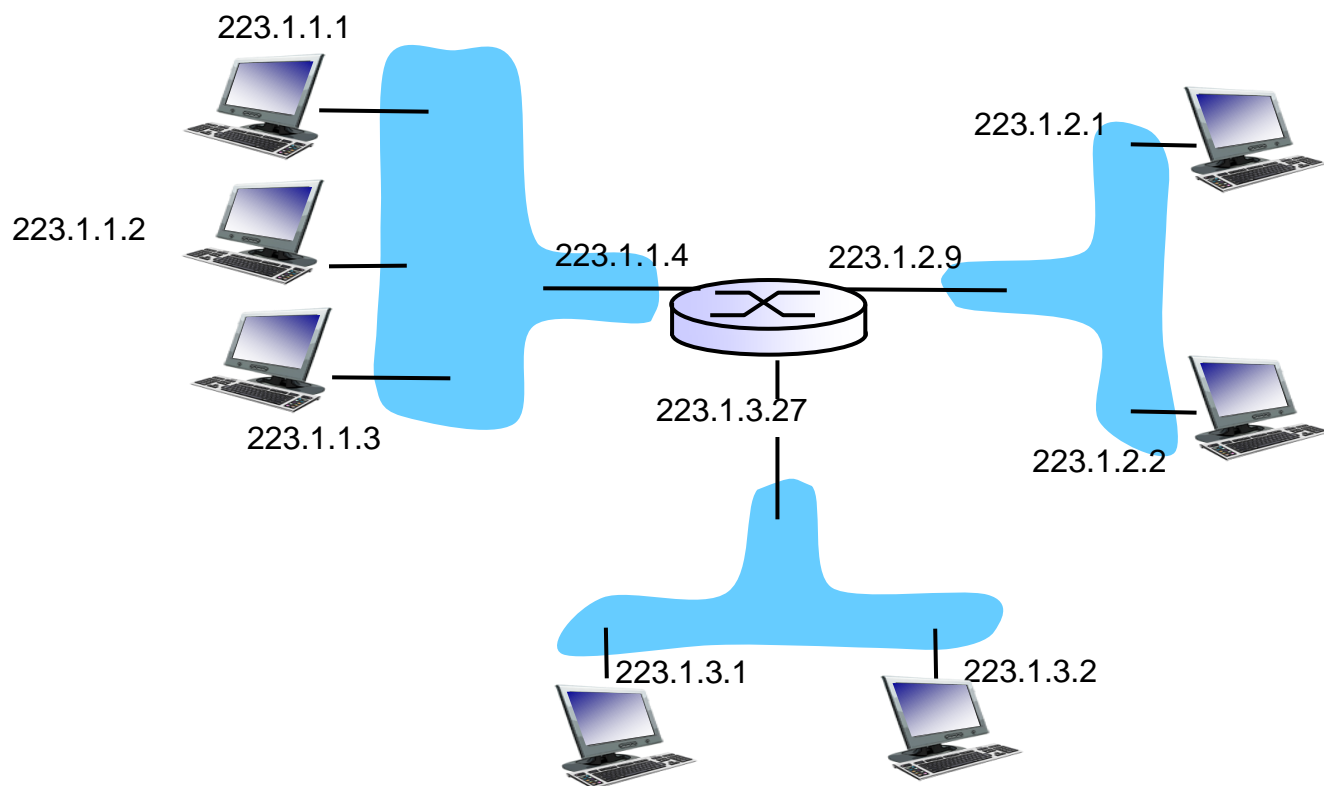
length =1500	ID =x	fragflag =1	offset =0
length =1500	ID =x	fragflag =1	offset =185
length =1040	ID =x	fragflag =0	offset =370

原始的数据4000字节的数据报，包含20字节的IP头部，实际数据载荷为3980字节，现在分为3个报文段：
3980=1480+1480+1020

$$\text{offset} = 1480*2/8 = 370$$

4.3.3 IPv4编址

- IPv4地址: 分配给主机或路由器接口(interface)的32-bit 标识符
- 接口: 主机/路由器与物理链路之间的边界
 - 路由器有多个接口
 - 主机可以有多个接口
 - 每个接口有一个IP地址



$$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_1 \underbrace{00000001}_1 \underbrace{00000001}_1$$

IPv4地址(补充)

- 早期的分类编址：根据不同的取值范围，早期将IP地址分为五类
 - IP地址中前5位用于标识IP地址的类别
 - A类地址的第一位为“0”，
 - B类地址的前两位为“10”，
 - C类地址的前三位为“110”，
 - D类地址的前四位为“1110”，（组播）
 - E类地址的前五位为“11110”。（实验）
 - A类、B类与C类地址为基本的IP地址。

IPv4地址(补充)

- IP地址由两部分组成：
 - 网络号：指明主机所在物理网络的编号。
 - 主机号：主机在物理网络中的编号。

	1	8	16	24	32	主机地址范围
A类地址	0	网络地址 (7位)	主机地址 (24位)			1. 0. 0. 0到 127. 255. 255. 255
B类地址	10	网络地址 (14位)		主机地址 (16位)		128. 0. 0. 0到 191. 255. 255. 255
C类地址	110	网络地址 (21位)			主机地址 (8位)	192. 0. 0. 0到 223. 255. 255. 255
D类地址	1110	多目的广播地址 (28位)				224. 0. 0. 0到 239. 255. 255. 255
E类地址	11110	保留用于实验和将来使用				240. 0. 0. 0到 247. 255. 255. 255

早期分类地址导致的问题

- 一个A类的IP地址，可以有24bit用于分配主机地址，因此可以支持(2的24次方)个主机，但是一个家庭或者组织往往不需要这么多的地址空间，造成浪费。
- 一个C类的IP地址，只有8bit用于分配主机地址，因此只能支持256个主机，又不太够用。
- 因此，逐渐按类别进行地址分配被CIDR技术取代

202.38.75.11 (C类, 24比特网络号, 8比特主机号)

113.12.23.06 (A类, 8比特网络号, 16比特主机号)

常用的特殊IP地址

- ❑ 127.0.0.1/8-127.255.255.254/8 (127.0.0.0和127.255.255.255除外)**环回地址**
 - 这是预留的一组IP地址，主要是用来识别主机本身的地址。也叫做“localhost”一般用来测试使用的。做开发的人比较熟悉。
- ❑ 本地链路地址 (169.254.0.0 / 16 或 169.254.0.1 到 169.254.255.254)
 - 通常被称为自动专用 IP 地址 (APIPA) 地址。Windows DHCP 客户端在无可用 DHCP 服务器时，使用它们进行自我配置。在对等连接中非常有用。
- ❑ 10.x.x.x, 172.16.x.x-172.31.x.x, 192.168.x.x.
 - 这三个地址段主要是我们私有的内网地址。也就是我们平时企业或者家里局域网所使用的地址段，我们比较熟悉的应该就是192.168.x.x 这个地址段了

常用的特殊IP地址

□ 0.0.0.0

这个地址严格上来说都不是真正意义上的IP地址。主要是用来标识不清楚的网络和主机的。系统遇到无法识别的网络或主机的时候会统一的归纳到这个地址

□ 255.255.255.255

这个地址是受限的广播地址。主要指一个网段内的所有主机

常用的特殊IP地址

□ 0.0.0.0

这个地址严格上来说都不是真正意义上的IP地址。主要是用来标识不清楚的网络和主机的。系统遇到无法识别的网络或主机的时候会统一的归纳到这个地址

□ 255.255.255.255

这个地址是受限的广播地址。主要指一个网段内的所有主机

子网的概念

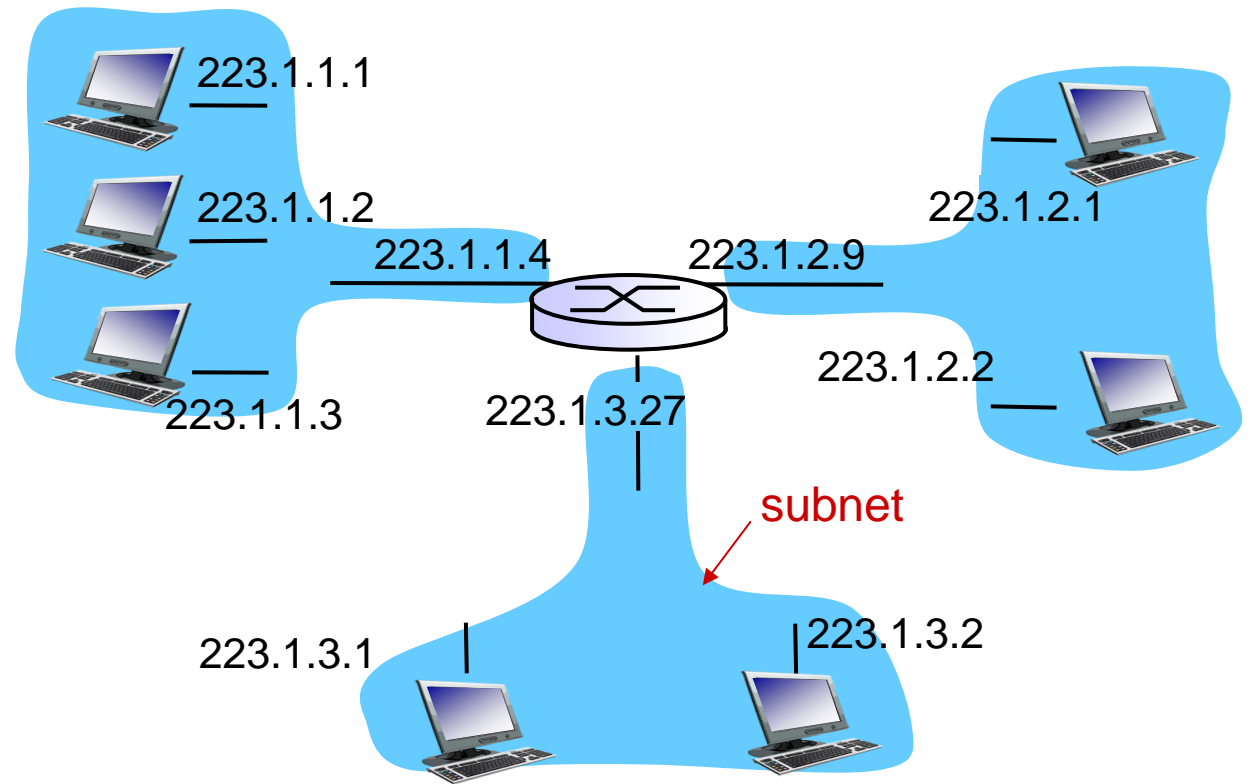
□ IP 地址:

- 子网部分 (高位 bits)
- 主机部分 (低位 bits)

□ 什么是一个子网？

(从IP地址的观点来看)

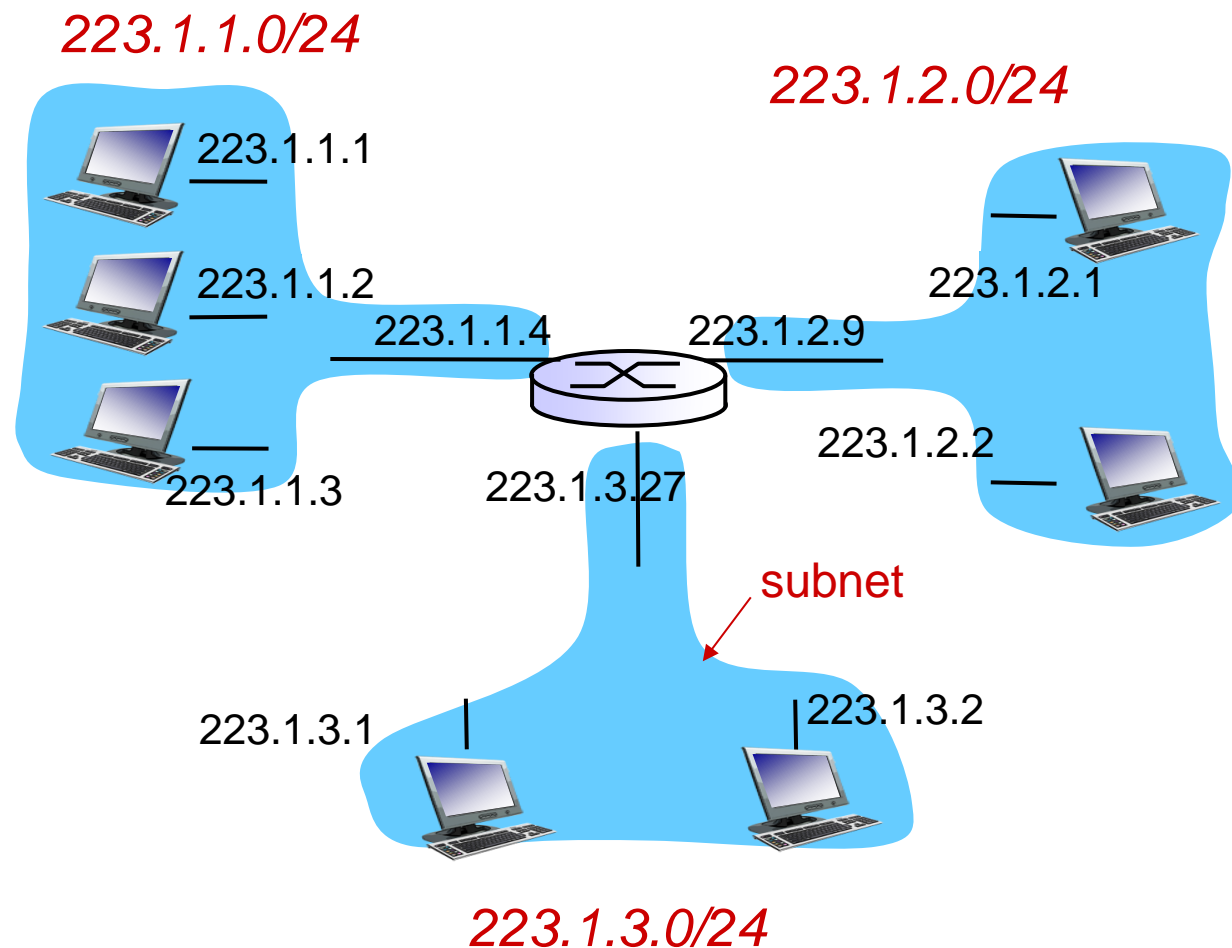
- 设备接口的IP地址具有同样的网络部分
- 没有路由器的介入，物理上能够相互到达



network consisting of 3 subnets

子网的概念

- 定义：为了确定子网，分开主机和路由器的每个接口，从而产生了几个分离的网络岛，这些独立网络中的每一个叫做一个子网(subnet).



subnet mask: /24

图中有多少个子网？

6个子网，请同学们分别列出这6个子网

223.1.1.0/24

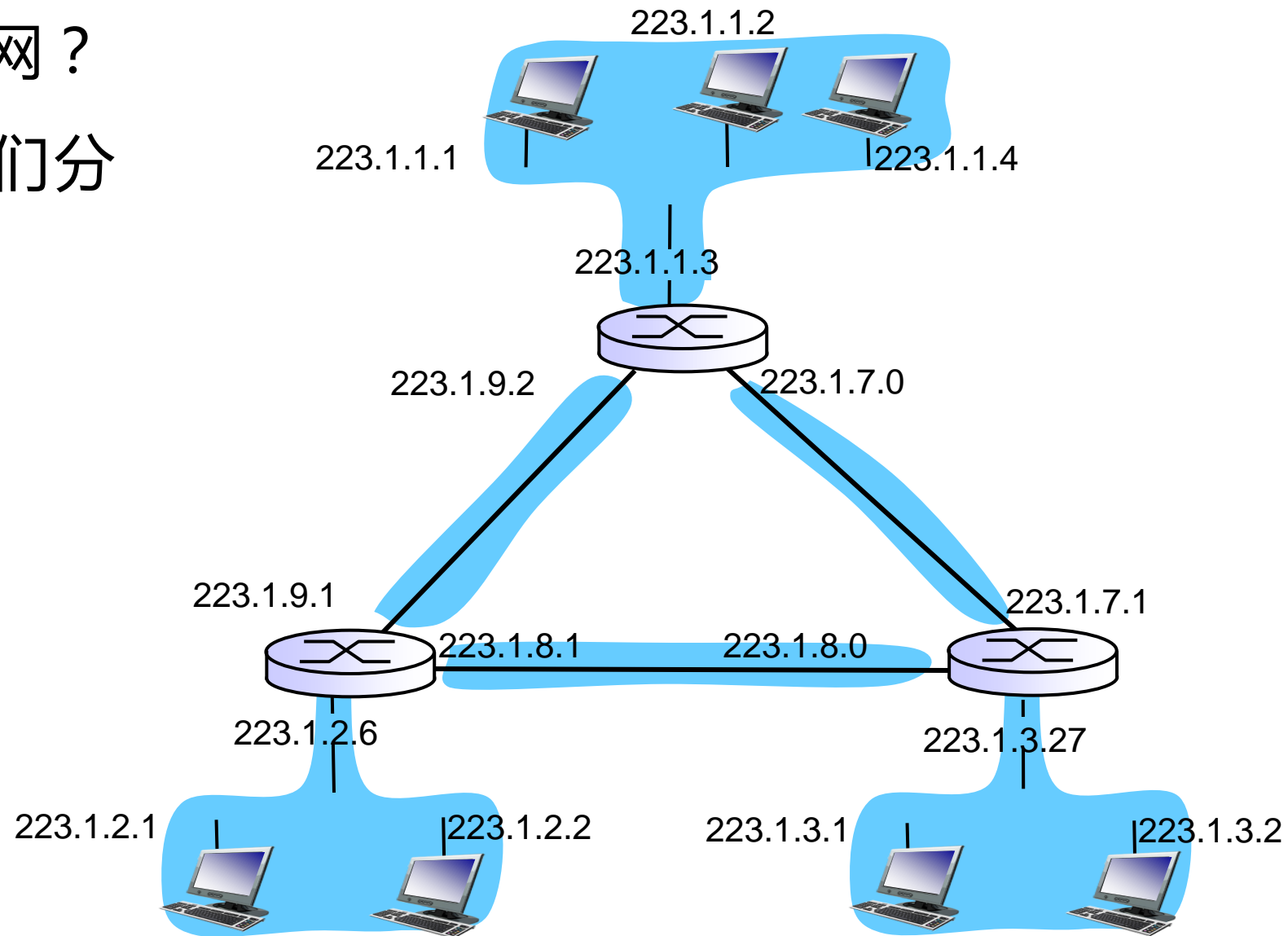
223.1.9.0/24

223.1.7.0/24

223.1.8.0/24

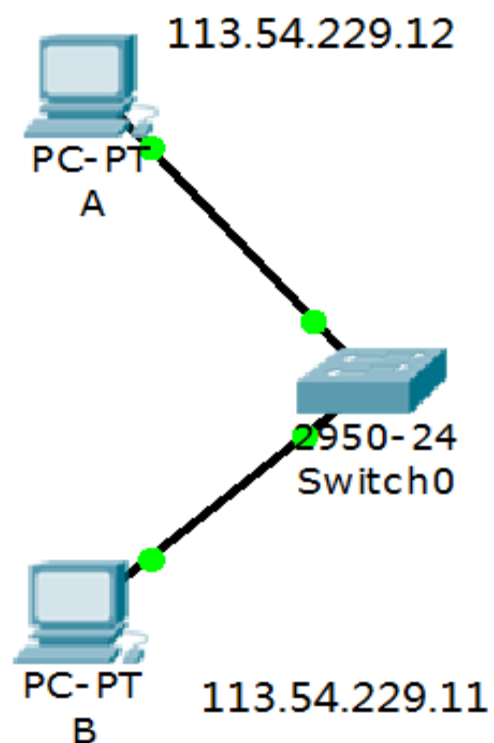
223.1.2.0/24

223.1.3.0/24



例子：

主机A的网络号：113.54.224.0/19



子网的例子，主机A和主机B通过二层交换机相连，可以相互ping通



网络前缀为/19
这个说法等价于子网掩码：
255.255.224.0 (11100000)

主机A能够ping通主机B

如果把B主机的IP地址改为
113.54.224.11/19

主机A仍然能够ping通主机B
请说明理由。

位于同一个子网的两台主机能够
相互通信

CIDR的实例

□ 请同学们写出这个网络接口的相关信息

- IPv4地址
- 子网掩码
- IPv4的默认网关
- IPv4地址对应的网络号

113.54.224.0/19

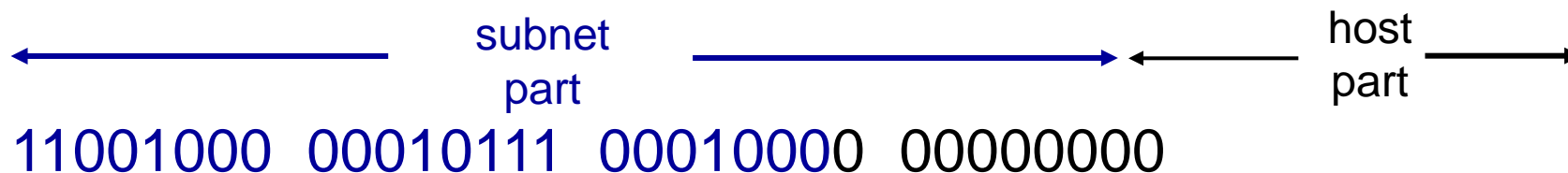
课堂练习



无类别域间路由选择CIDR

□ CIDR: Classless InterDomain Routing

- 地址中的网络部分可以任意长
- 地址格式: a.b.c.d/x, 这里 x 是地址网络部分的bit数
- 子网掩码: 用全1表示网络位, 全0表示主机位



子网掩码

200.23.16.0/23

11111111 11111111 11111110 00000000

255

255

254

0

子网掩码和网络前缀

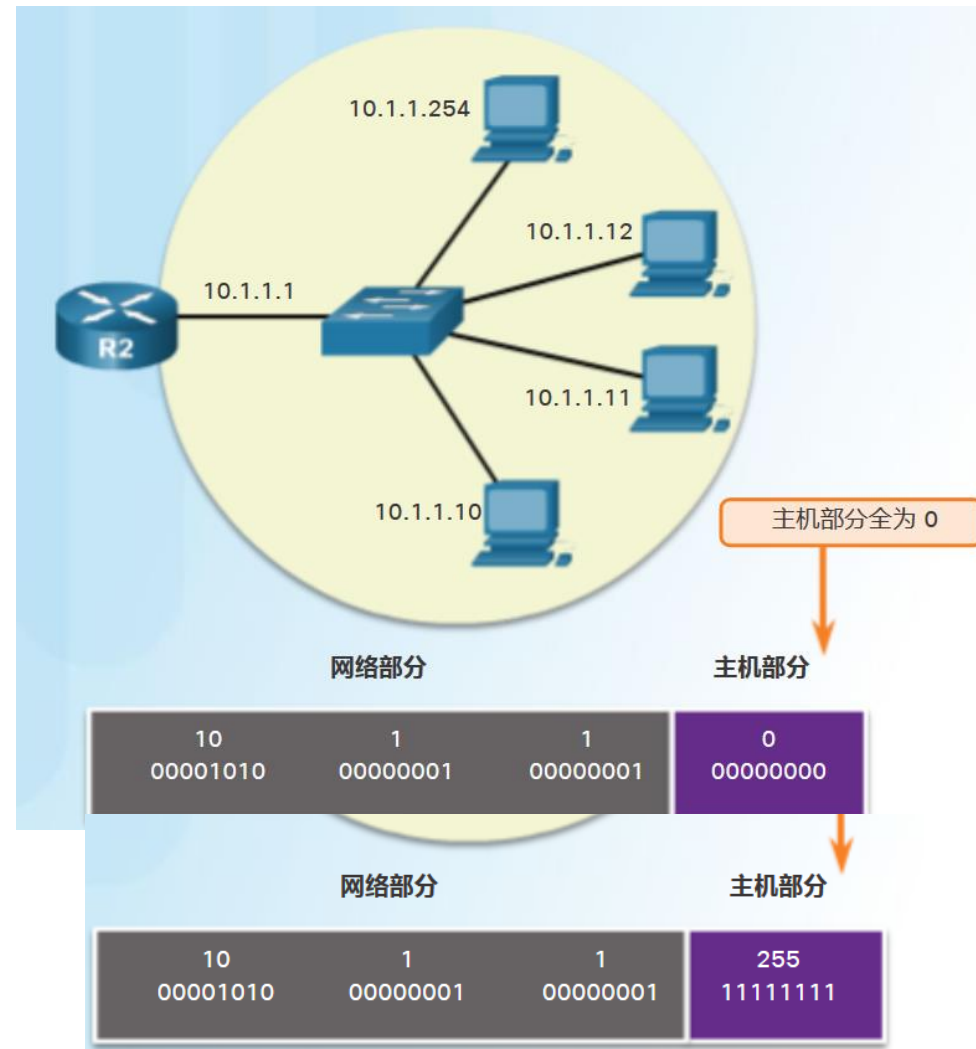
比较子网掩码和前缀长度

子网掩码	32 位地址	前缀长度
255.0.0.0	11111111.00000000.00000000.00000000	/8
255.255.0.0	11111111.11111111.00000000.00000000	/16
255.255.255.0	11111111.11111111.11111111.00000000	/24
255.255.255.128	11111111.11111111.11111111.10000000	/25
255.255.255.192	11111111.11111111.11111111.11000000	/26
255.255.255.224	11111111.11111111.11111111.11100000	/27
255.255.255.240	11111111.11111111.11111111.11110000	/28
255.255.255.248	11111111.11111111.11111111.11111000	/29
255.255.255.252	11111111.11111111.11111111.11111100	/30

网络地址、主机地址、广播地址

- ❑ 网络号：主机位全0
 - 同一个网络中的所有主机都具有相同的网络地址
 - 例子：网络地址10.1.1.0/24
- ❑ 网络广播地址：主机位全1
 - 例子：广播地址10.1.1.255/24
- ❑ 网络地址：主机位除了全0和全1都可以分配给一个网络接口

课堂练习：202.38.75.11/16，对应的网络地址和网络的广播地址

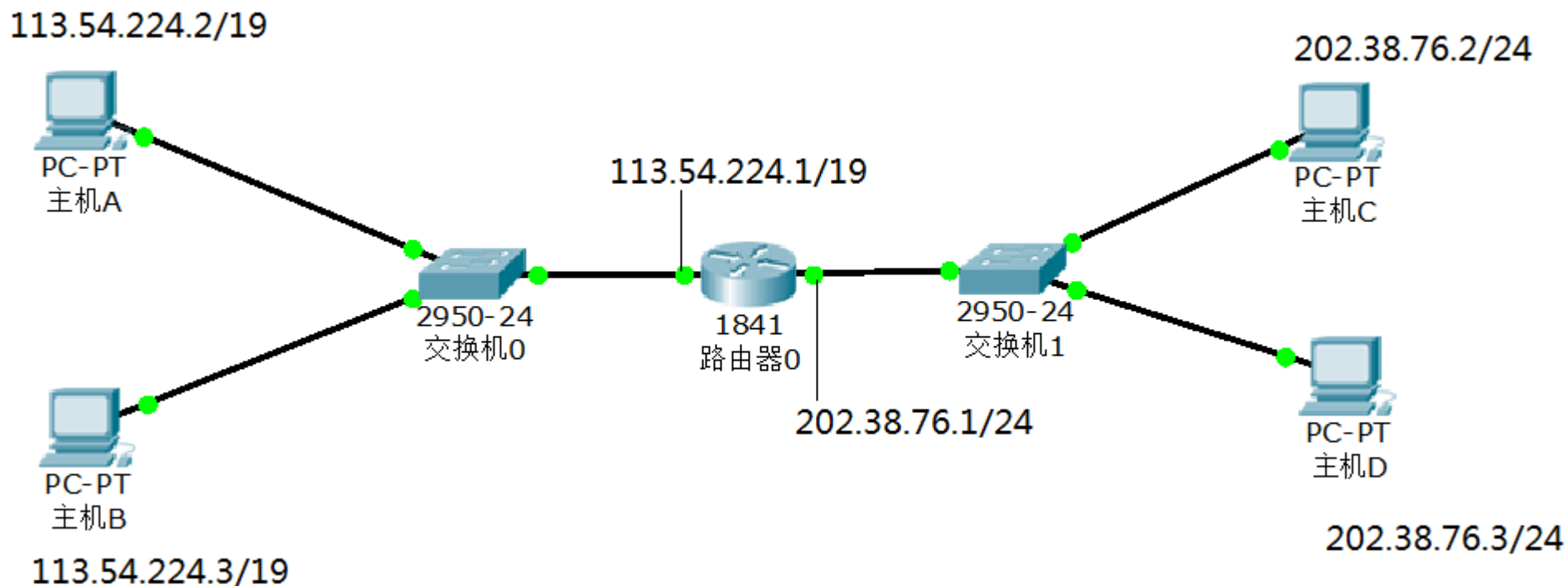


IPv4寻址

- 首先，我们需要理解CIDR技术，理解网络前缀的含义，能够利用子网掩码计算给定IP地址的网络号；很多同学喜欢根据IP地址直接判断，这是不对的，应该根据掩码操作后再比较；
- 理解因特网中IP寻址的概念；
 - 属于相同子网的主机通信，可以直接发送分组，不需要路由器参与；
 - 属于不同子网的主机通信，发送主机首先将分组发送到它的网关路由器，由路由器帮助它转发分组到接收主机；
 - 路由器将根据分组的目的地IP地址所属的子网地址，来查找其对应的输出口，将分组输出到接收主机或下一跳路由器，直到分组被接收主机接收。

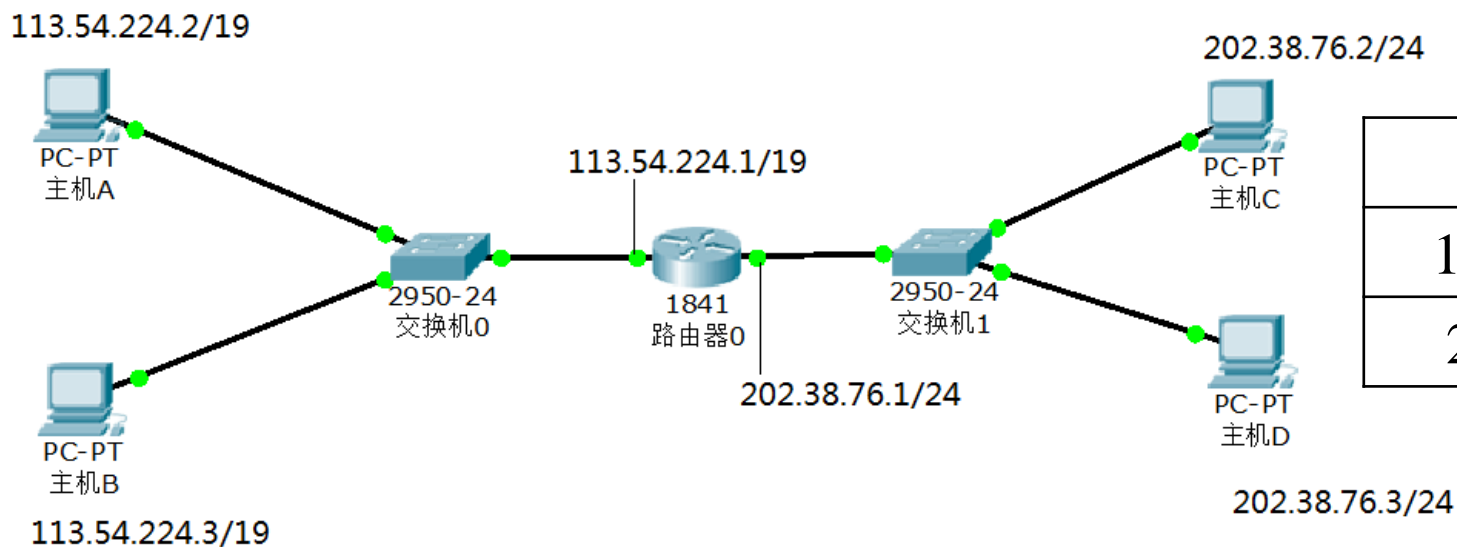
IPv4寻址

- 属于相同子网的两台主机通信，发送主机可以直接将分组发送给接收主机，不需要路由器参与；
 - 例如主机A(113.54.224.2/19)发送分组给主机B(113.54.224.3/19)；



IPv4寻址

- 属于不同子网的两台主机通信，发送主机首先将分组发送给它的网关路由器，由网关路由器帮助转发分组到接收主机；
 - 例如主机A(113.54.224.2/19)发送分组给出C(202.38.76.2/24)；
- 路由器是根据分组携带的目的IP地址的**网络号(子网地址)**，在转发表中查找其对应的输出端口，然后将分组转发到该网络号对应的输出端口输出。

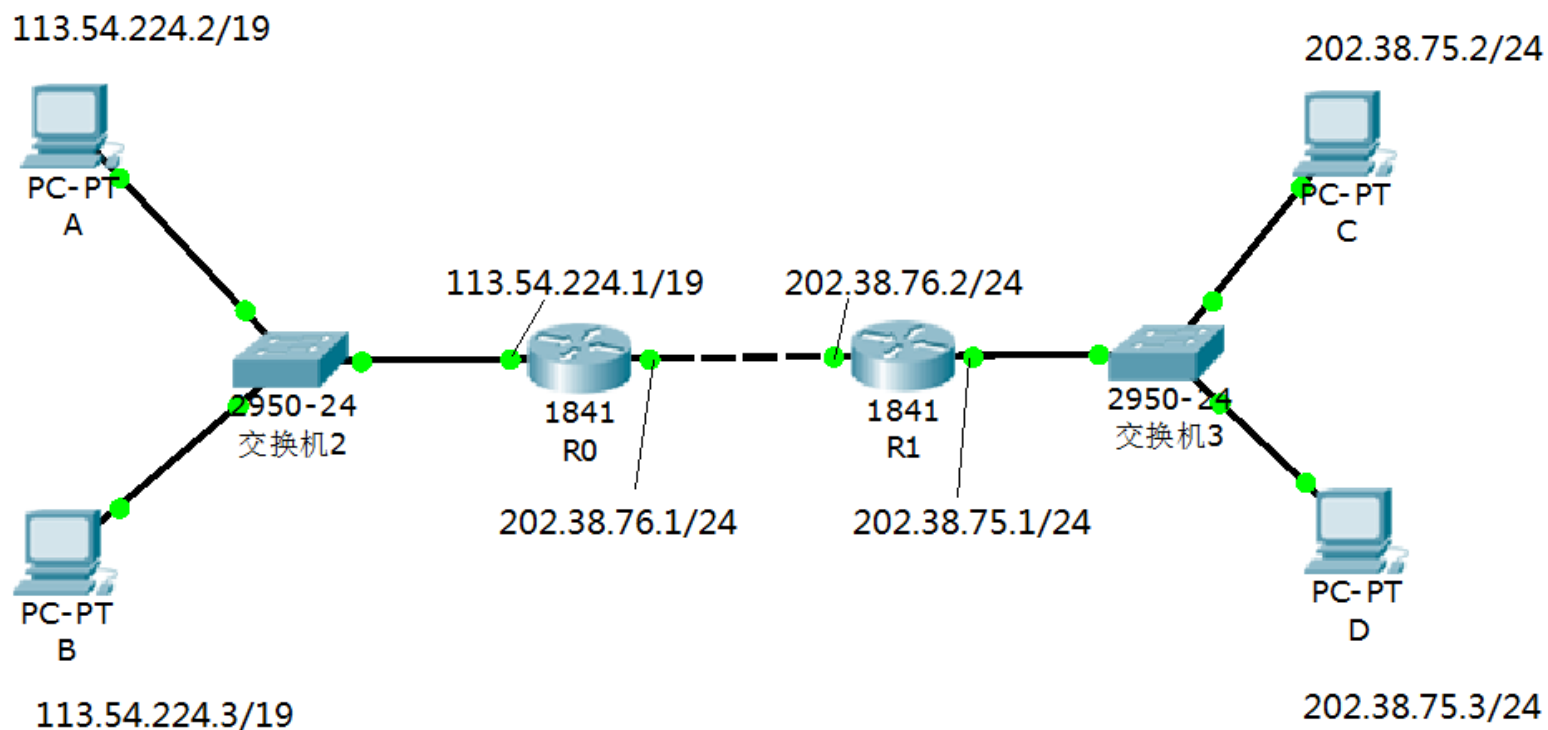


目的地址段	下一跳IP地址
113.54.224.0/19	与本地左端口直连
202.38.76.0/24	与本地右端口直连

IPv4寻址：案例分析

□ 在思科仿真软件中建立如下的拓扑结构

请同学们思考怎么配置R0和R1的转发表才能使4台主机相互通信



这里总共有3个子网：

- 子网113.54.224.0/19;
- 子网202.38.76.0/24;
- 子网202.38.75.0/24。

IPv4寻址：案例分析

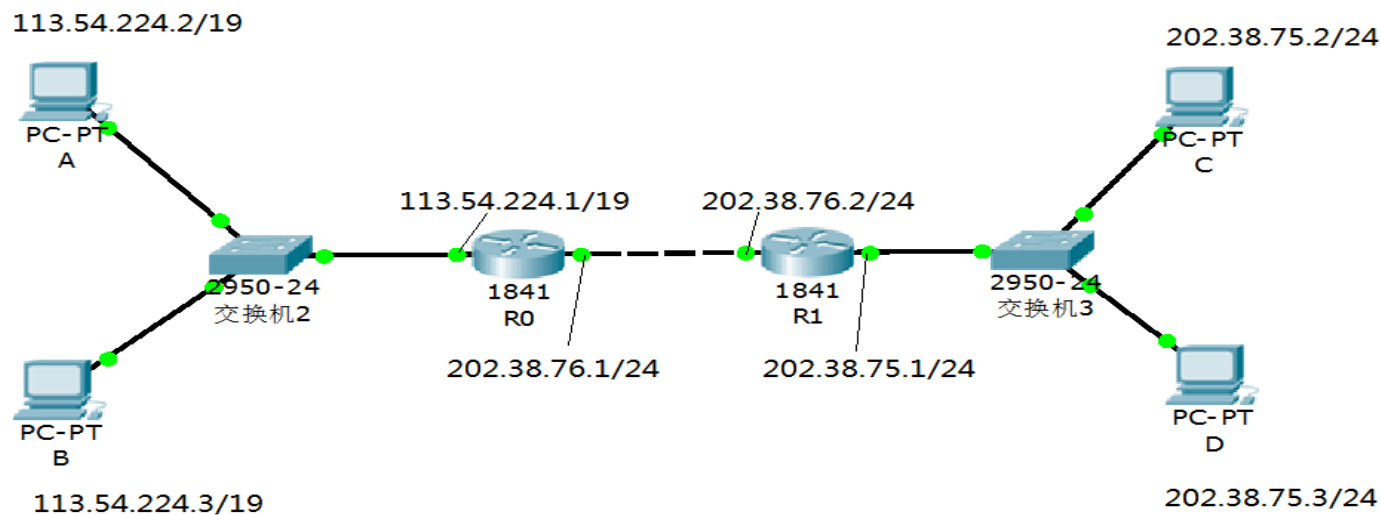
- 请同学们配置路由器R0和路由器R1的转发表，实现主机A和主机C能够相互通信
 - 假设R0和R1本身的转发表如下

R0

目的地址段	下一跳IP地址
113.54.224.0/19	与本地左端口直连
202.38.76.0/24	与本地右端口直连
202.38.75.0/24	202.38.76.2

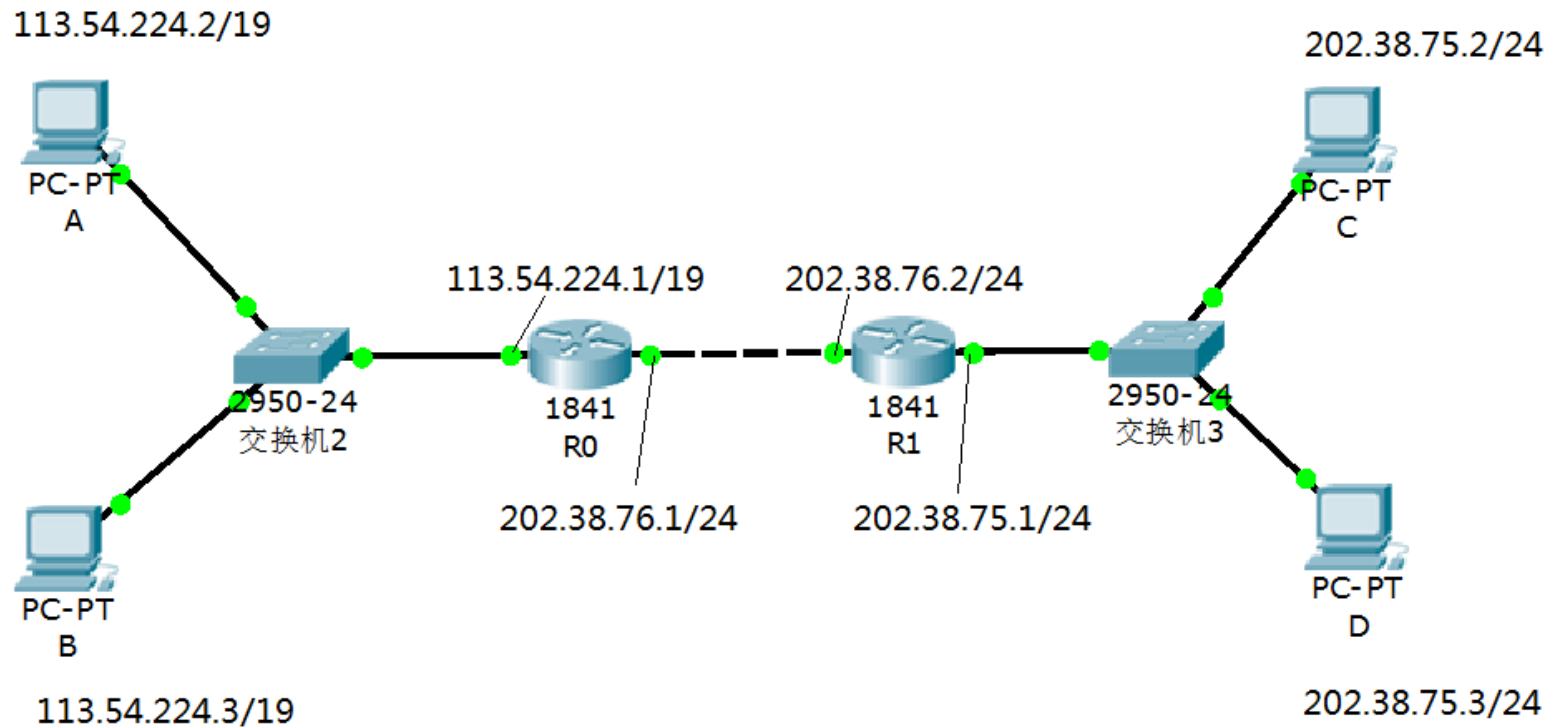
R1

目的地址段	下一跳IP地址
202.38.76.0/24	与本地左端口直连
202.38.75.0/24	与本地右端口直连
113.54.224.0/19	202.38.76.1



IPv4寻址：案例分析

- 课后思考：如果主机C的网关地址不设置，那么A主机还能和C主机通信吗？D主机还能和C主机通信吗？



提示1：A主机和C主机不属于同一子网；
提示2：C主机和D主机属于同一子网；

IPv4寻址

- ❑ IPv4寻址是根据网络号(子网部分)，首先在路由器的转发表中查找是否有该网络号对应的条目，将分组转发到该网络号对应的输出端口；
- ❑ 公司(个人)的机器如果配置了某个IP地址段，最初只有接入ISP的路由器知道这个IP地址段的信息，如何让其他路由器知道这个信息呢？（通过发送路由公告）
- ❑ 我们说路由信息的扩散是通过BGP协议(路由协议)向与该本地ISP相连的其他ISP的路由器扩散，直到所有因特网上的路由器都知道这个IP地址段的前缀信息。
- ❑ 怎么自动建立转发表条目，我们将在第5章进一步学习！

关于IP地址及域名的获得

□ ISP获得地址块的方法

ICANN: Internet Corporation for Assigned Names and Numbers

<http://www.icann.org/>

- 分配IP地址段
- 管理根DNS服务器
- 分配域名及解决纠纷

如何获得IP addresses

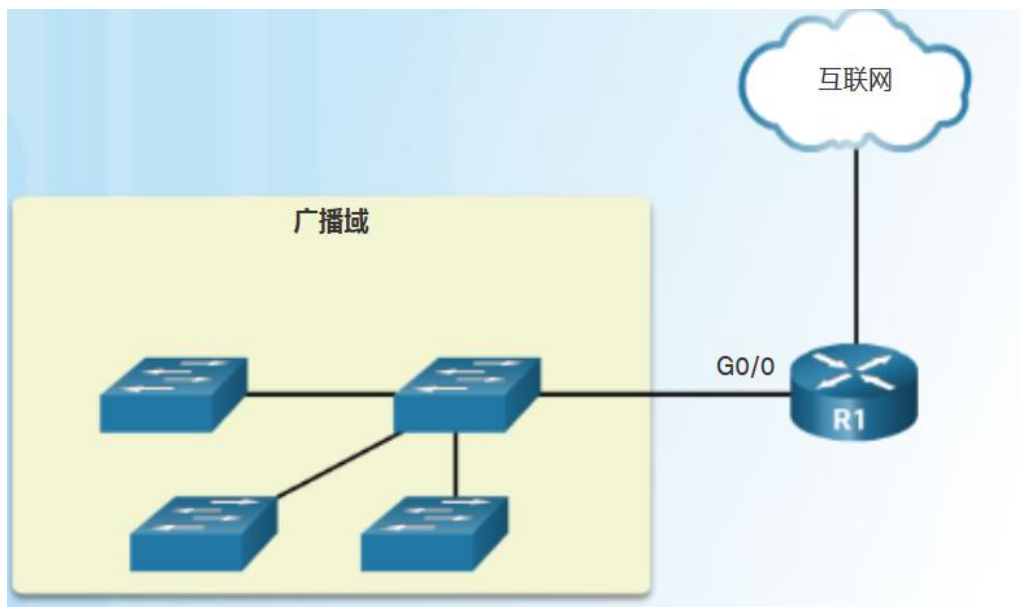
Q: 如何获得IP地址的网络部分

A: 由ISP的地址空间中分配

在分配的地址空间中划分子网

- 划分子网的原因
 - 广播域

一个企业从ISP获得地址块以后可以建立一个网络，所有主机都连入这个网络中，但是实际的操作中，企业通常将这个大的网络划分为若干小的子网，来进行组网。



左图是一个由若干交换机构成的广播域，在广播域中有许多消息(例如ARP消息、路由公告消息)是需要在所有除达到端口外的所有端口转发。因此一个大型广播域会导致：

- 网络运行缓慢，原因是它导致的流量太大
- 设备运行缓慢，原因是设备必须接受和处理每个广播数据包

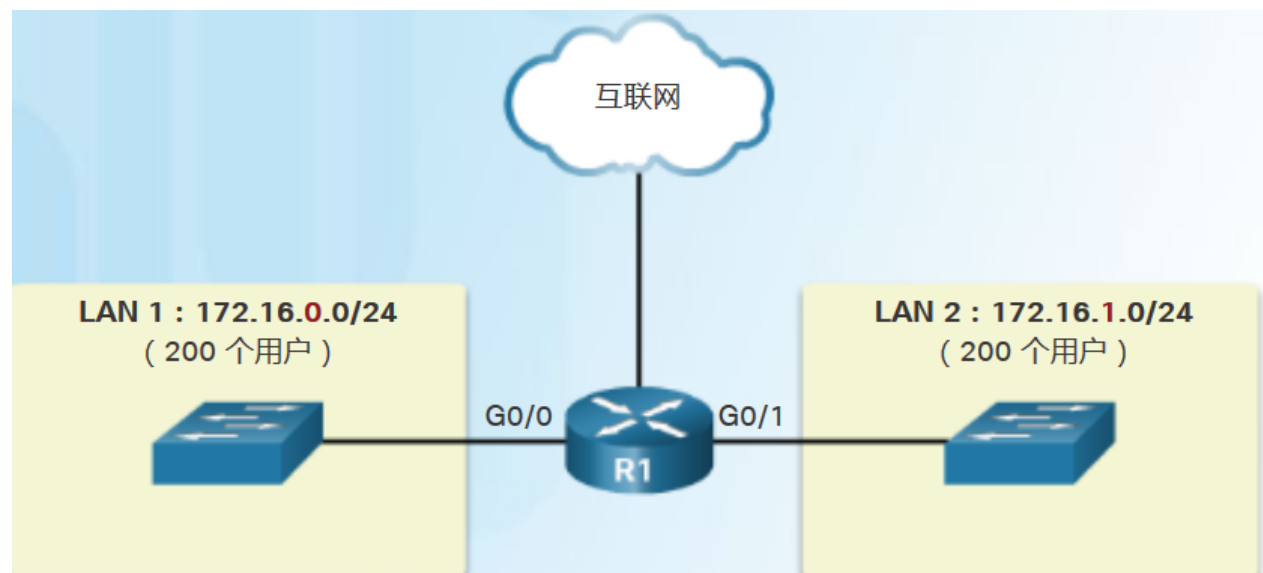
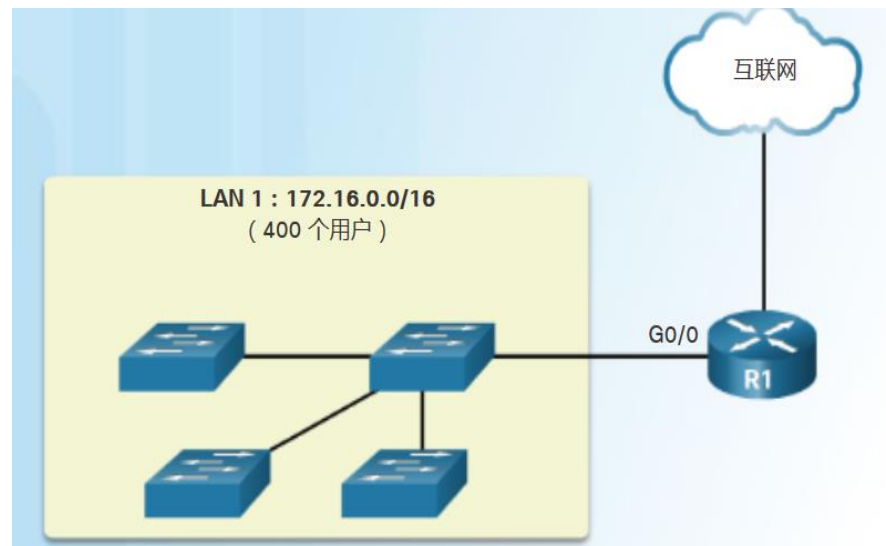
在分配的地址空间中划分子网

❑ 路由器能隔离广播域

- 路由器在收到广播消息时，它不会将其转发到其他接口。例如，当 R1 在其 Gigabit Ethernet 0/0 接口收到广播时，它不会将其转发到另一个接口。

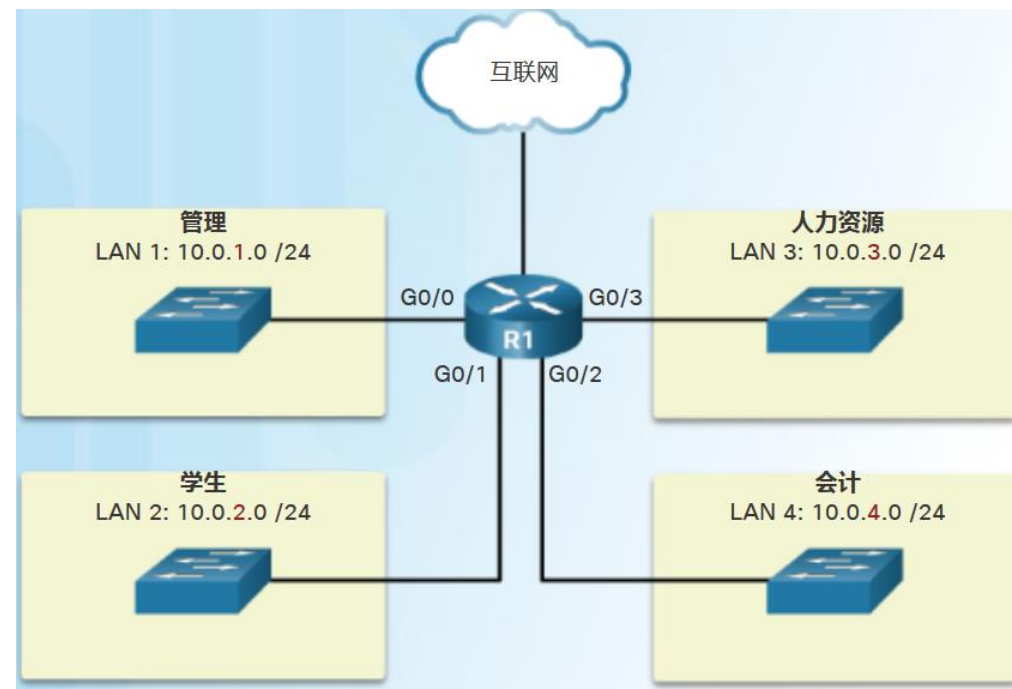
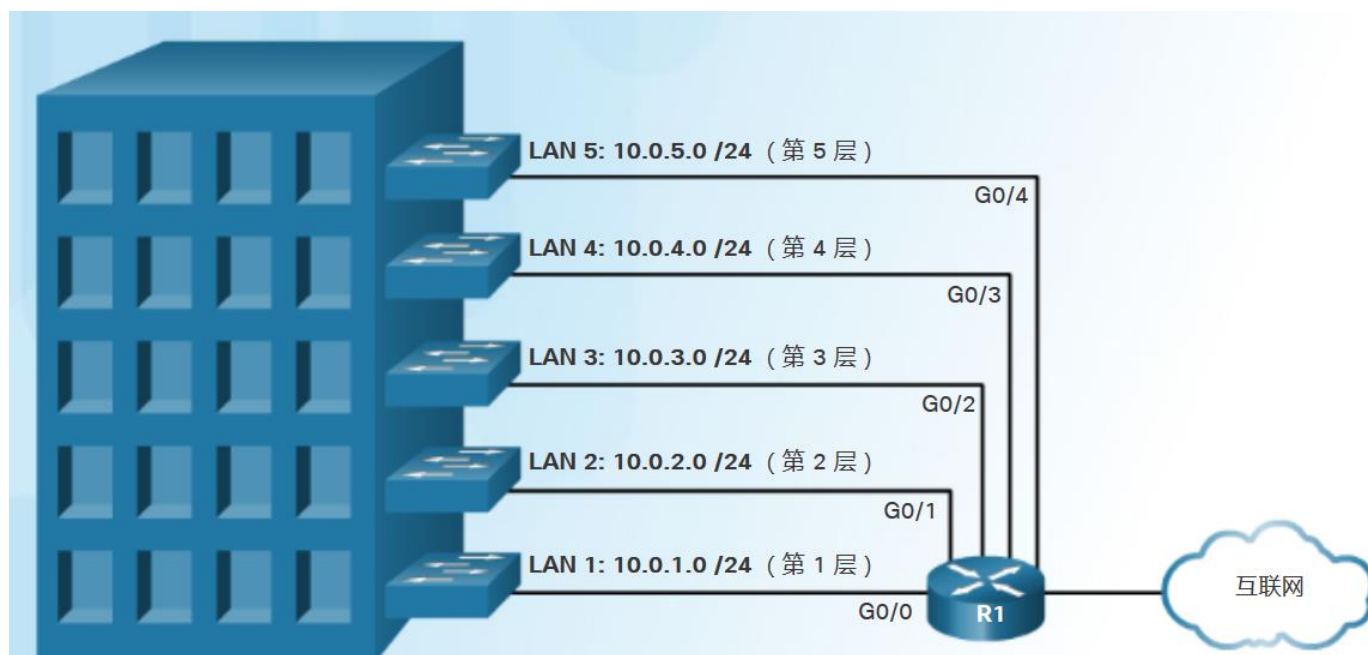
❑ 为了解决大型广播域导致的问题，可行的解决方案是，采用“**子网划分**”的方法，缩减网络规模以创建更小的广播域。

图中，网络地址为 172.16.0.0 /16 的 LAN 1 中的 400 个用户被划分到两个子网中，每个子网包含 200 个用户，网络地址分别为 172.16.0.0/24 和 172.16.1.0/24。广播消息仅在更小的网络域内传播。因此，LAN 1 中的广播不会传播到 LAN 2。



子网划分的原因

- 子网划分可以降低整体网络流量并改善网络性能。它也能让管理员实施安全策略，例如哪些子网允许或不允许进行通信。



子网划分的方法

- 创建 IPv4 子网时会将一个或多个主机位作为网络位。具体做法是延长子网掩码，从地址的主机部分借用若干位来增加网络位。借用的主机位越多，可以定义的子网也就越多。（本身网络前缀/24）

前缀长度	子网掩码	二进制子网掩码 (n = 网络, h = 主机)	子网数	主机数
/25	255.255.255.128	n n n n n n n n . n n n n n n n n . n n n n n n n n . n h h h h h h h 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 0 0 0 0 0 0 0	2	126
/26	255.255.255.192	n n n n n n n n . n n n n n n n n . n n n n n n n n . n n h h h h h h h 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 0 0 0 0 0 0	4	62
/27	255.255.255.224	n n n n n n n n . n n n n n n n n . n n n n n n n n . n n n h h h h h h 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 0 0 0 0 0	8	30
/28	255.255.255.240	n n n n n n n n . n n n n n n n n . n n n n n n n n . n n n n h h h h h 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 0 0 0 0	16	14
/29	255.255.255.248	n n n n n n n n . n n n n n n n n . n n n n n n n n . n n n n n h h h h 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 0 0 0	32	6
/30	255.255.255.252	n n n n n n n n . n n n n n n n n . n n n n n n n n . n n n n n n h h h 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 1 1 . 1 1 1 1 1 1 0 0	64	2

子网划分的方法（定长子网划分）

ISP's block	<u>11001000 00010111 0001</u> 000 0	00000000	200.23.16.0/20
Organization 0	<u>11001000 00010111 0001</u> 000 0	00000000	200.23.16.0/23
Organization 1	<u>11001000 00010111 0001</u> 001 0	00000000	200.23.18.0/23
Organization 2	<u>11001000 00010111 0001</u> 010 0	00000000	200.23.20.0/23
...
Organization 7	<u>11001000 00010111 0001</u> 111 0	00000000	200.23.30.0/23

注意每个子网的网络前缀变为了**23**比特

Organization 0 11001000 00010111 0001**000**0 00000000 200.23.16.0/23

地址范围: 11001000 00010111 0001**000**0 00000001

11001000 00010111 0001**000**1 11111110

11111111.11111111.11111110.00000000

子网掩码: 255.255.254.0

如何分配子网（变长子网划分）（重要）

- 给定一个地址块，怎么进行子网划分，以及为每个子网分配地址段
- 某单位从ISP处购买了一个IP地址块202.40.70.0/24。如果该单位有3个部门，第1个部门78台计算机，第2个和第3个部门各有48台计算机，请你为3个部门划分子网，并给出各部门的网络号、子网掩码和主机可分配的IP地址范围。

子网1的最后8比特的设置：0xxxxxxx；这里后7比特表示主机号，/25；

子网2的最后8比特的设置：10xxxxxx；这里后6比特表示主机号，/26；

子网3的最后8比特的设置：11xxxxxx；这里后6比特表示主机号，/26；

如何分配子网(重要)

子网1: 202.40.70.00000001 ~ 202.40.70.01111110 (202.40.70. 1 ~ 202.40.70.126)

网络号为: 202.40.70.0 子网掩码为: 255.255.255.128

子网2: 202.40.70.10000001 ~ 202.40.70.10111110 (202.40.70. 129~202.40.70.190)

网络号为: 202.40.70.128 子网掩码为: 255.255.255.192

子网3: 202.40.70.11000001 ~ 202.40.70.11111110 (202.40.70. 193~202.40.70.254)

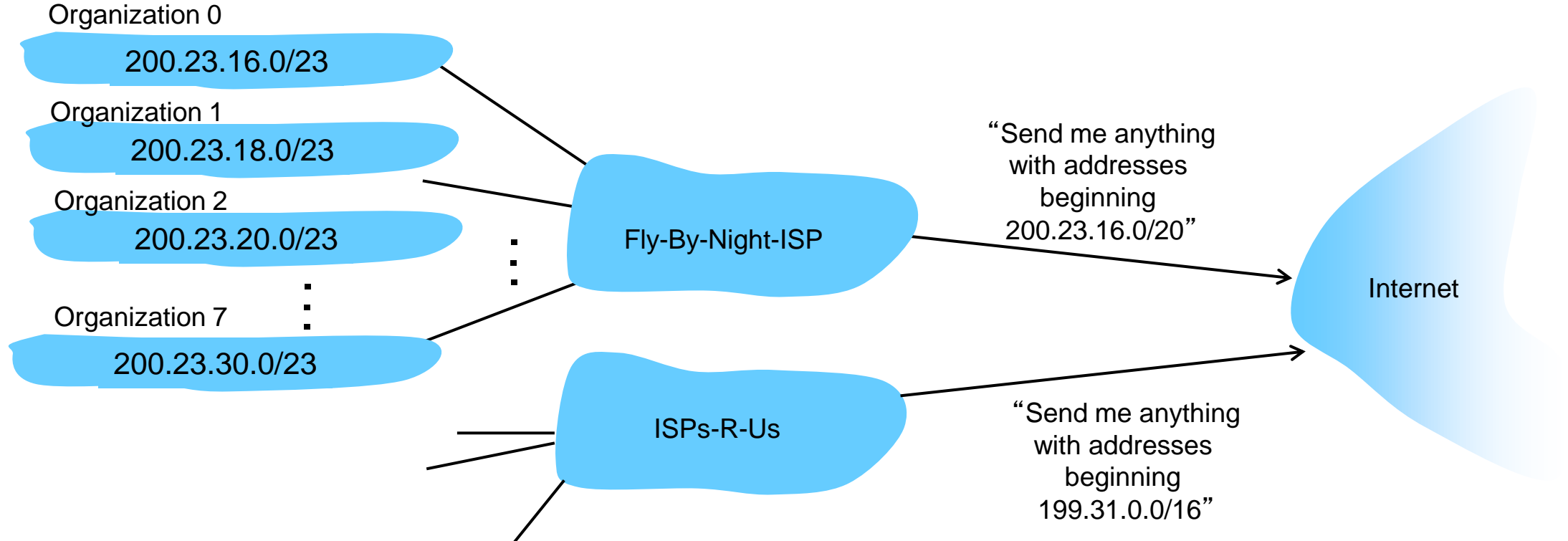
网络号为: 202.40.70.192 子网掩码为: 255.255.255.192

课堂练习

P13. Consider a router that interconnects three subnets: Subnet 1, Subnet 2, and Subnet 3. Suppose all of the interfaces in each of these three subnets are required to have the prefix 223.1.17/24. Also suppose that Subnet 1 is required to support at least 60 interfaces, Subnet 2 is to support at least 90 interfaces, and Subnet 3 is to support at least 12 interfaces. Provide three network addresses (of the form a.b.c.d/x) that satisfy these constraints.

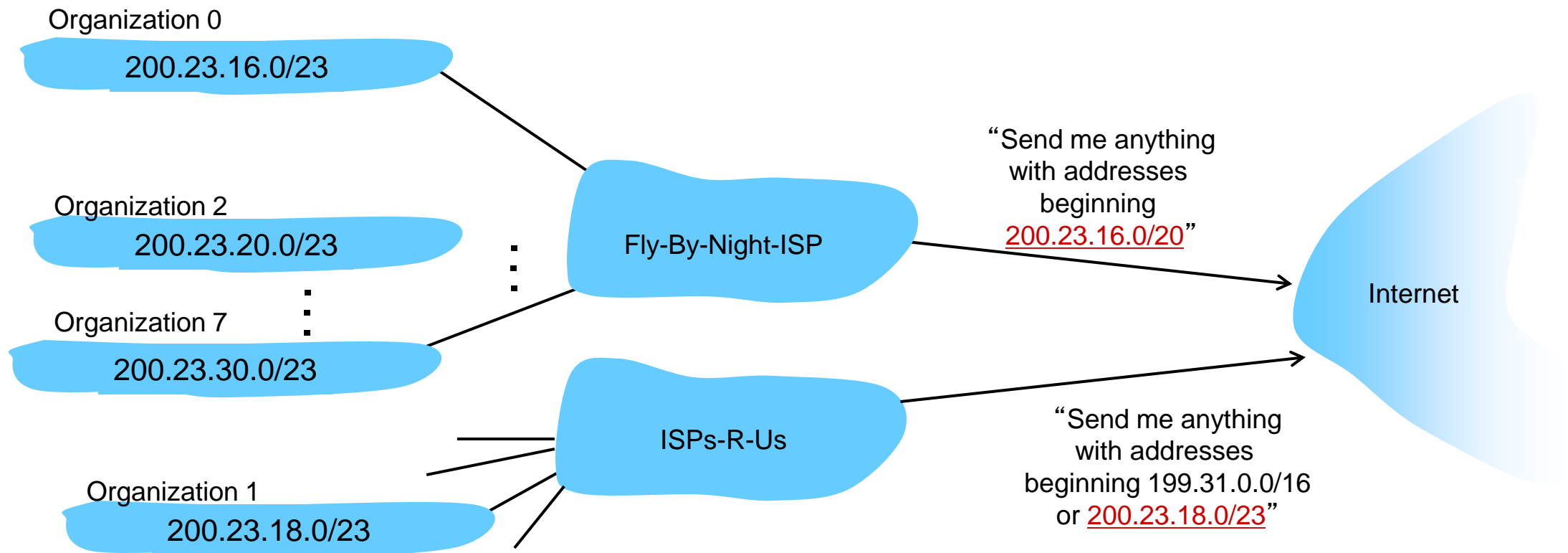
层次化地址: 路由聚合

hierarchical addressing allows efficient advertisement of routing information:



层次化地址:更详细的信息

ISPs-R-Us has a more specific route to Organization 1



配置主机IP地址

- ❑ 手工获取
- ❑ 动态主机配置协议获取（DHCP）：允许一个主机自动获取以下信息：
 - IP地址、子网掩码、默认网关、DNS地址
 - plug-and-play（即插即用）

hard-coded by system admin in a file

- Windows: control-panel->network->configuration->tcp/ip->properties
- UNIX: /etc/rc.config

DHCP: Dynamic **H**ost **C**onfiguration **P**rotocol:
dynamically get address from as server

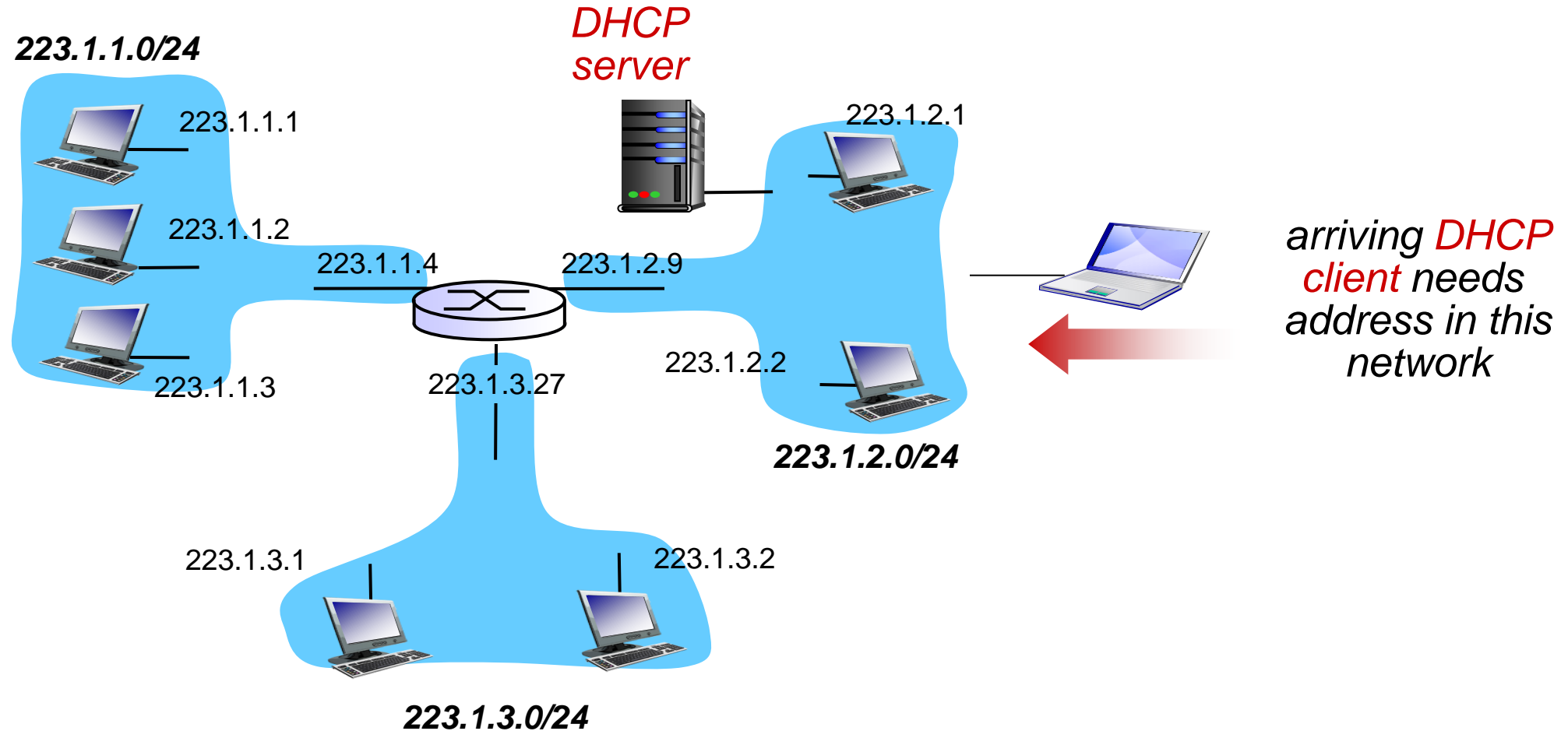
--

动态地址配置

□ DHCP概述：

- 主机广播 “DHCP discover” msg [optional]
- DHCP服务器响应 “DHCP offer” msg [optional]
- 主机请求IP地址： “DHCP request” msg
- DHCP服务器发送地址： “DHCP ack” msg

DHCP client-server scenario



DHCP客户-服务器场景

DHCP server: 223.1.2.5



DHCP discover

src : 0.0.0.0, 68
dest.: 255.255.255.255, 67
yiaddr: 0.0.0.0
transaction ID: 654

arriving
client



DHCP offer

src: 223.1.2.5, 67
dest: 255.255.255.255, 68
yiaddr: 223.1.2.4
transaction ID: 654
lifetime: 3600 secs

DHCP request

src: 0.0.0.0, 68
dest.: 255.255.255.255, 67
yiaddr: 223.1.2.4
transaction ID: 655
lifetime: 3600 secs

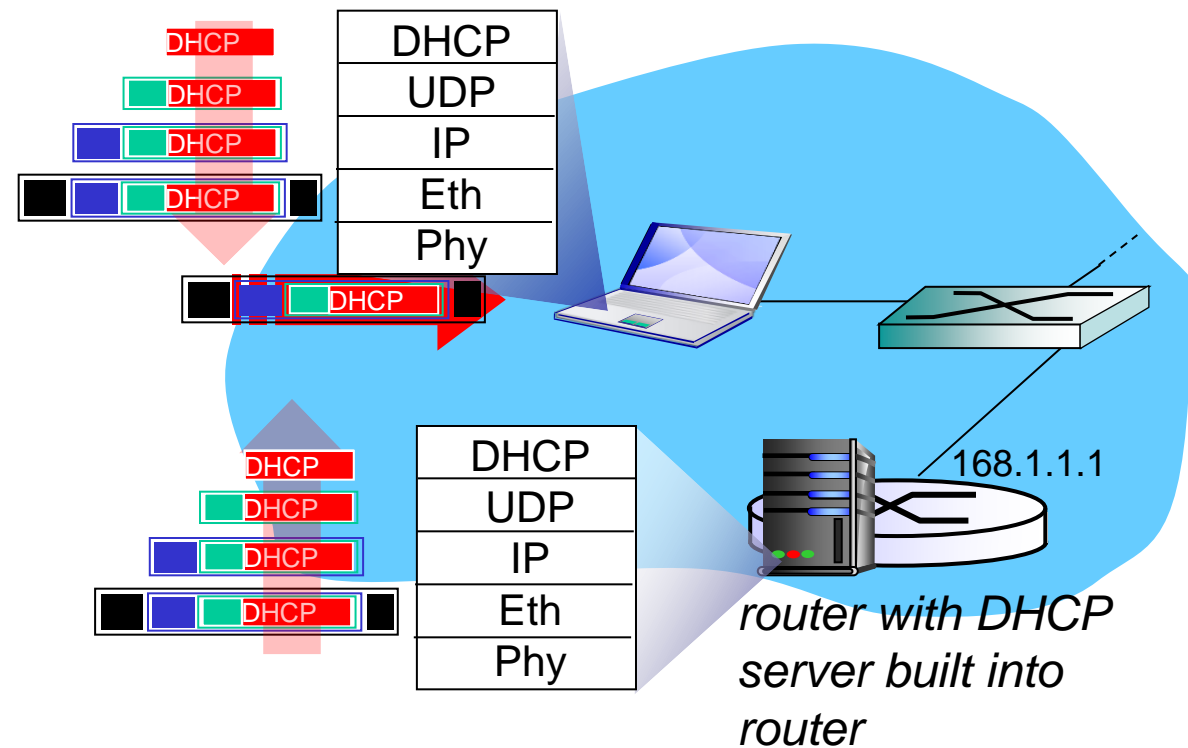
DHCP ACK

src: 223.1.2.5, 67
dest: 255.255.255.255, 68
yiaddr: 223.1.2.4
transaction ID: 655
lifetime: 3600 secs

DHCP: 不仅获得IP地址

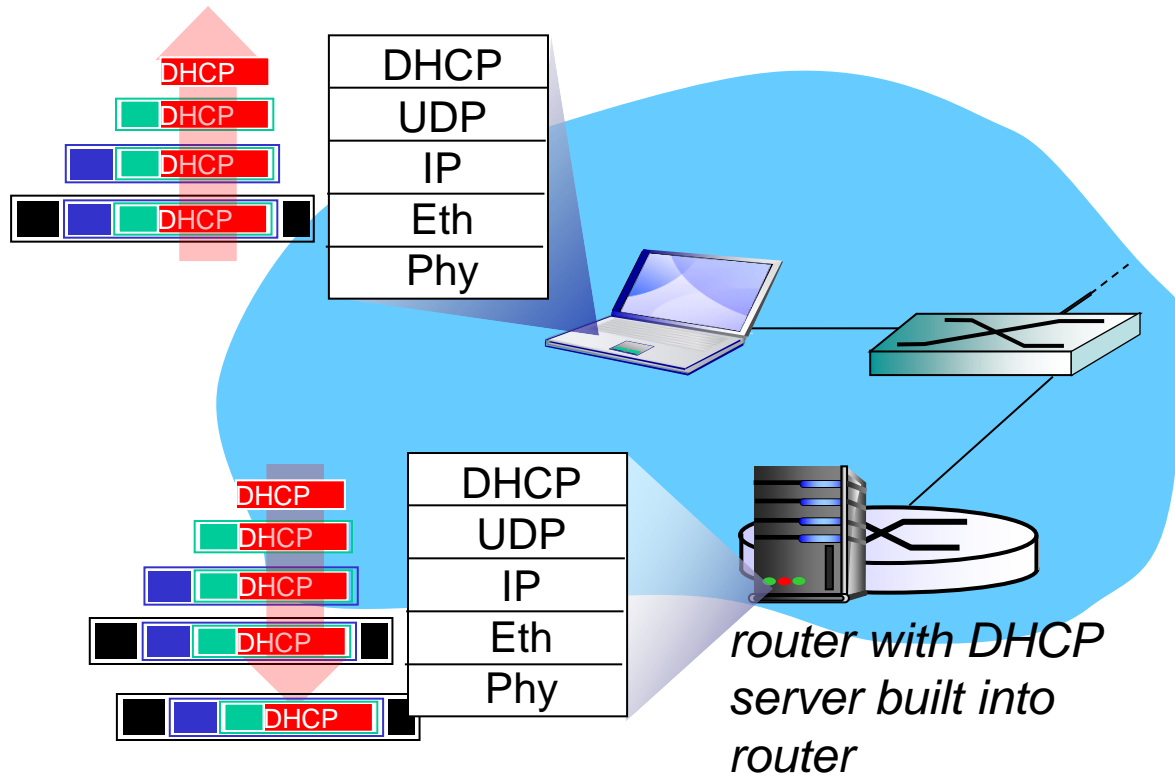
- DHCP返回的不仅仅是IP地址:
 - 客户的第一跳路由器的地址
 - DNS服务器的IP地址和域名
 - 网络掩码（用于指示网络的网络号部分和主机号部分）

DHCP: 例子



- connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use **DHCP**
- DHCP request encapsulated in **UDP**, encapsulated in **IP**, encapsulated in **802.1 Ethernet**
- Ethernet frame broadcast (dest: FFFFFFFFFFFFFFFF) on LAN, received at router running DHCP server
- Ethernet demuxed to IP demuxed, UDP demuxed to DHCP

DHCP: 例子



- DHCP server formulates DHCP ACK containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server
- encapsulation of DHCP server, frame forwarded to client, demuxing up to DHCP at client
- client now knows its IP address, name and IP address of DNS server, IP address of its first-hop router

DHCP: Wireshark output (home LAN)

Message type: **Boot Request (1)**

Hardware type: Ethernet

Hardware address length: 6

Hops: 0

Transaction ID: 0x6b3a11b7

Seconds elapsed: 0

Bootp flags: 0x0000 (Unicast)

Client IP address: 0.0.0.0 (0.0.0.0)

Your (client) IP address: 0.0.0.0 (0.0.0.0)

Next server IP address: 0.0.0.0 (0.0.0.0)

Relay agent IP address: 0.0.0.0 (0.0.0.0)

Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)

Server host name not given

Boot file name not given

Magic cookie: (OK)

Option: (t=53,l=1) **DHCP Message Type = DHCP Request**

Option: (61) Client identifier

Length: 7; Value: 010016D323688A;

Hardware type: Ethernet

Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)

Option: (t=50,l=4) Requested IP Address = 192.168.1.101

Option: (t=12,l=5) Host Name = "nomad"

Option: (55) Parameter Request List

Length: 11; Value: 010F03062C2E2F1F21F92B

1 = Subnet Mask; 15 = Domain Name

3 = Router; 6 = Domain Name Server

44 = NetBIOS over TCP/IP Name Server

.....

request

Message type: **Boot Reply (2)**

Hardware type: Ethernet

Hardware address length: 6

Hops: 0

Transaction ID: 0x6b3a11b7

Seconds elapsed: 0

Bootp flags: 0x0000 (Unicast)

Client IP address: 192.168.1.101 (192.168.1.101)

Your (client) IP address: 0.0.0.0 (0.0.0.0)

Next server IP address: 192.168.1.1 (192.168.1.1)

Relay agent IP address: 0.0.0.0 (0.0.0.0)

Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)

Server host name not given

Boot file name not given

Magic cookie: (OK)

Option: (t=53,l=1) DHCP Message Type = DHCP ACK

Option: (t=54,l=4) Server Identifier = 192.168.1.1

Option: (t=1,l=4) Subnet Mask = 255.255.255.0

Option: (t=3,l=4) Router = 192.168.1.1

Option: (6) Domain Name Server

Length: 12; Value: 445747E2445749F244574092;

IP Address: 68.87.71.226;

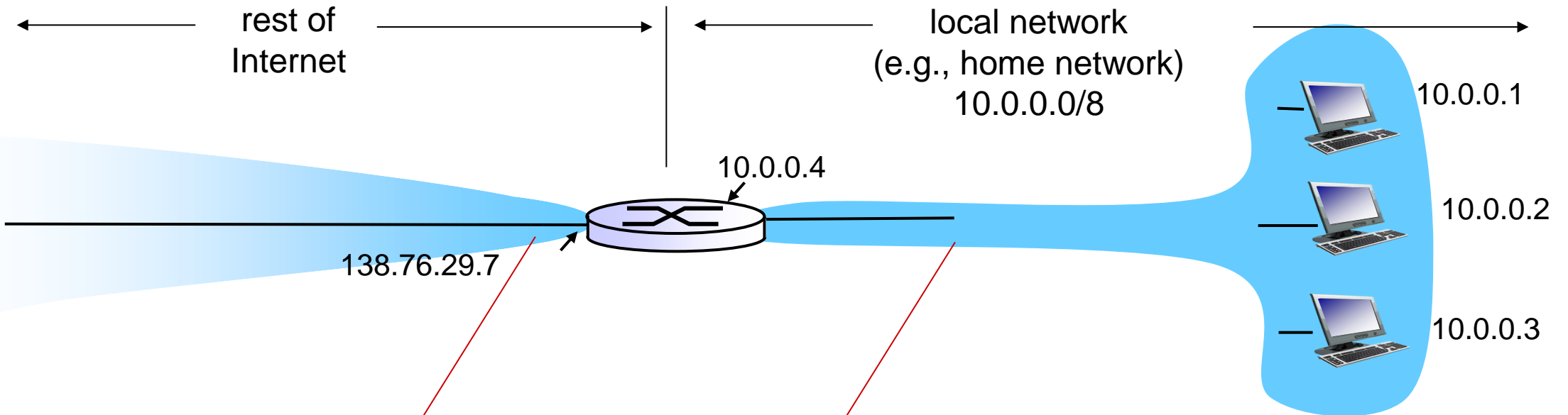
IP Address: 68.87.73.242;

IP Address: 68.87.64.146

Option: (t=15,l=20) Domain Name = "hsd1.ma.comcast.net."

reply

4.3.4 网络地址转换(NAT)



所有分组离开本地网络有同样的 单个源端 NAT IP 地址: 138.76.29.7, 不同的源端口号

源端或目的端是这个网络的分组有10.0.0/8 地址用于指示源地址或目的地址

IPv4地址枯竭问题

- ❑ IPv4地址枯竭问题(理想情况40亿个IP地址)
- ❑ 地址重用的解决方案
 - 方案1：动态分配IP地址，即不给主机配置固定IP地址，当主机不连接因特网时，IP地址可被路由器收回重新分配给其他主机使用；
 - 方案2：网络地址转换，即从32比特的IPv4地址范围中划分出若干地址段作为私有地址，配置私有IP地址的主机不能直接访问外部网络(因特网)，只能在自己的内部网络中使用，因此可以在不同的内部网络中重复使用，内网主机如果需要与外部网络通信，可以共用一个全局IP地址来解决；
 - 私有地址：10.0.0.0 ~ 10.255.255.255：172.16.0.0 ~ 172.31.255.255：192.168.0.0 ~ 192.168.255.255：

4.3.4 网络地址转换(NAT)

- 动机: 对外部网络来讲, 本地网络只用一个IP地址
 - 不需要从 ISP分配一系列地址, 只要一个IP地址用于所有设备
 - 在本地网络改变设备的IP地址不用通知外部世界
 - 可以变更 ISP 而不用改变本地网络的设备的地址
 - 本地网络内部设备不能被外部世界明确寻址,或是不可见 (增加了安全性).

4.3.4 网络地址转换(NAT)

□ 执行NAT，路由器必须做到：

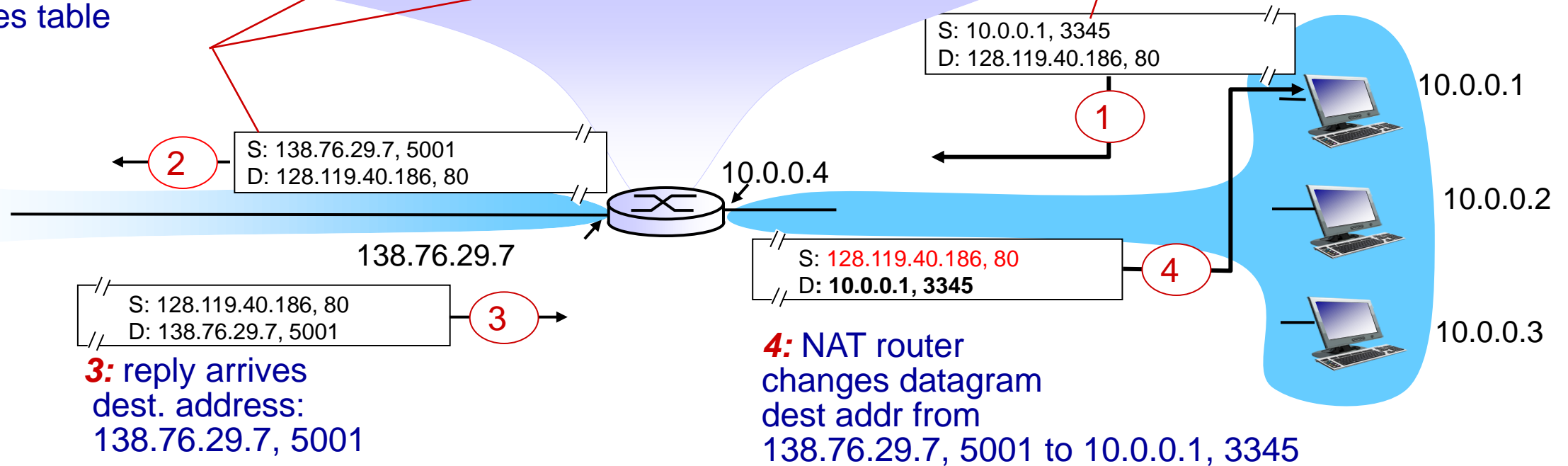
- 外出的分组：替换每个外出的分组的 (源IP 地址, 端口号) 为 (NAT IP 地址, 新端口号)，远程客户/服务器用(NAT IP地址, 新端口号)作为目的地来响应。
- (在NAT转换表中)记录每个(源IP 地址, 端口号)到 (NAT IP地址, 新端口号) 转换配对
- 进来的分组：对每个进来的分组，用保存在NAT表中的对应的(源IP 地址, 端口号) 替换分组中的目的域 (NAT IP 地址, 新端口号)

4.3.4 网络地址转换(NAT)

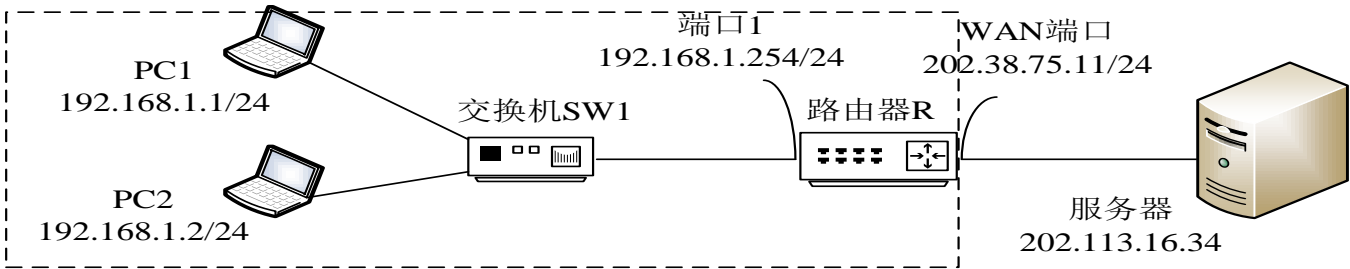
2: NAT router changes datagram source addr from 10.0.0.1, 3345 to 138.76.29.7, 5001, updates table

NAT translation table	
WAN side addr	LAN side addr
138.76.29.7, 5001	10.0.0.1, 3345
.....

1: host 10.0.0.1 sends datagram to 128.119.40.186, 80



课堂练习



假定一个家庭网络中的拓扑结构如下图所示，其内部网络为192.168.1.0/24，路由器的WAN端口的地址为202.38.75.11/24，这是一个全局IP地址，假设PC1和PC2上分别向服务器202.113.16.34/24的80端口发起2个TCP连接，则在路由器R上需要建立4个对应的NAT转换表项，假设在PC1、PC2上的端口号分别为3000~3003，NAT转换表如下所示

源IP: 202.113.16.34, 源port: 80

目的IP: 202.38.75.11, 目的port: 4000

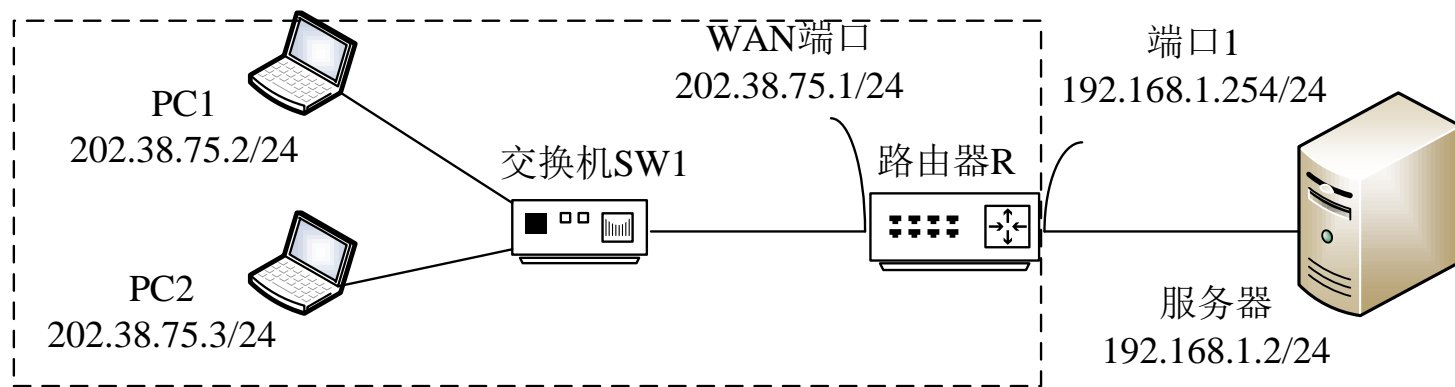
TCP连接	WAN端		LAN端	
	IP地址	端口号	IP地址	端口号
TCP1	202.38.75.11	4000	192.168.1.1	3000
TCP2	202.38.75.11	4001	192.168.1.1	3001
TCP3	202.38.75.11	4002	192.168.1.2	3002
TCP4	202.38.75.11	4003	192.168.1.2	3003

请写出服务器返回给PC1的TCP1的数据分组的源、目的地址，及源、目的端口号，写出PC1收到的数据分组的源、目的地址，及源、目的端口号。

源IP: 202.113.16.34, 源port: 80,
目的IP: 192.168.1.1, 目的port: 3000

NAT案例讨论

- 案例：如果在内网中有个服务器，这个服务器只有内部IP地址，那么需要在路由器上做怎样的配置，才能使外网的客户机访问这个服务器呢？



- 回答：只需在路由器的NAT表中手动增加一个表项：
(202.38.75.1: 80) 映射为 (192.168.1.2: 80)。外网主机访问时，只需使用http://202.38.75.1:80就可以访问。

4.3.4 网络地址转换(NAT)

□ 16-bit 端口号:

- 一个局域网地址可以同时支持60,000 个并发连接! (2的16次方)

□ NAT存在争议

- 路由器只应该处理到第三层
- 违反了端到端主张
 - 应用程序设计者在设计时不得不将NAT加以考虑
 - 如P2P应用程序
- 应使用IPv6来解决地址短缺问题
- NAT穿透: 如果客户机想要连接到NAT后面的服务器呢?

4.3.5 IPv6: 动机

- ❑ *initial motivation*: 32-bit address space soon to be completely allocated.
- ❑ additional motivation:
 - header format helps speed processing/forwarding
 - header changes to facilitate QoS

IPv6 datagram format:

- **fixed-length 40 byte header**
- no fragmentation allowed

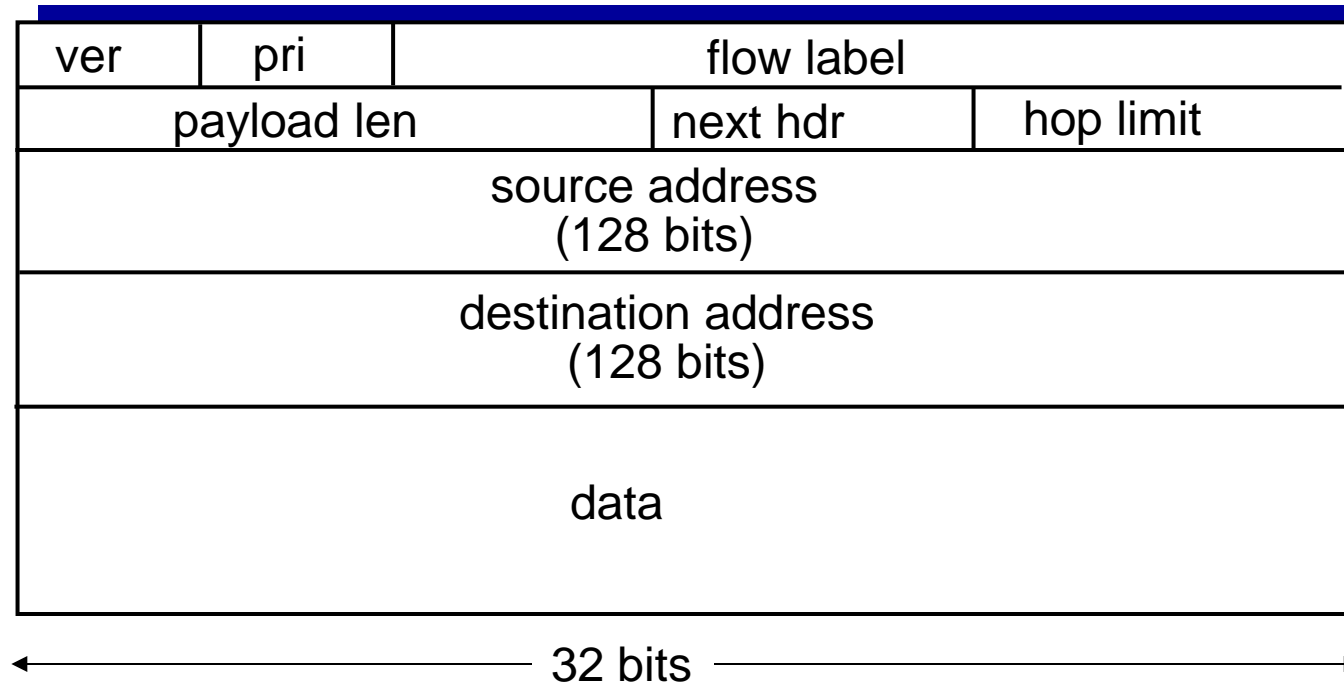
IPv6 数据报格式

priority: identify priority among datagrams in flow

flow Label: identify datagrams in same “flow.”

(concept of “flow” not well defined).

next header: identify upper layer protocol for data

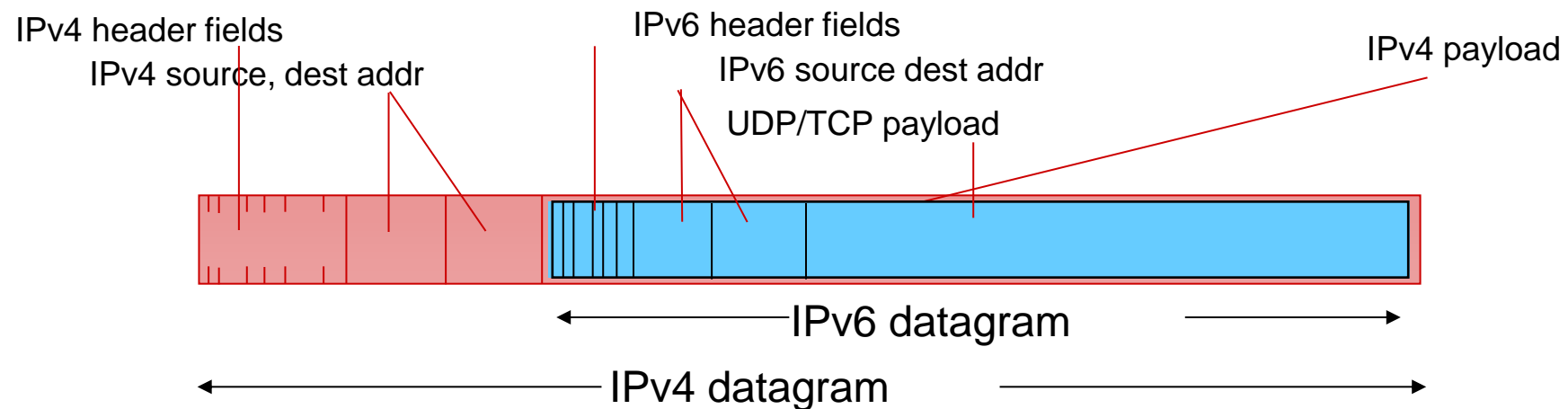


Other changes from IPv4

- ❑ *checksum*: removed entirely to reduce processing time at each hop
- ❑ *options*: allowed, but outside of header, indicated by “Next Header” field
- ❑ *ICMPv6*: new version of ICMP
 - additional message types, e.g. “Packet Too Big”
 - multicast group management functions

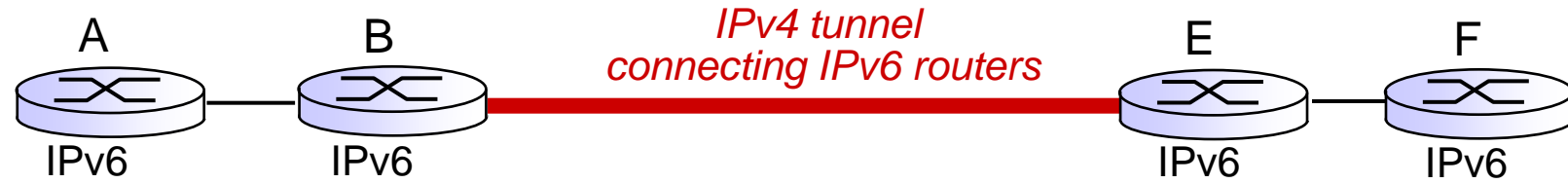
Transition from IPv4 to IPv6

- not all routers can be upgraded simultaneously
 - no “flag days”
 - how will network operate with mixed IPv4 and IPv6 routers?
- *tunneling*: IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers

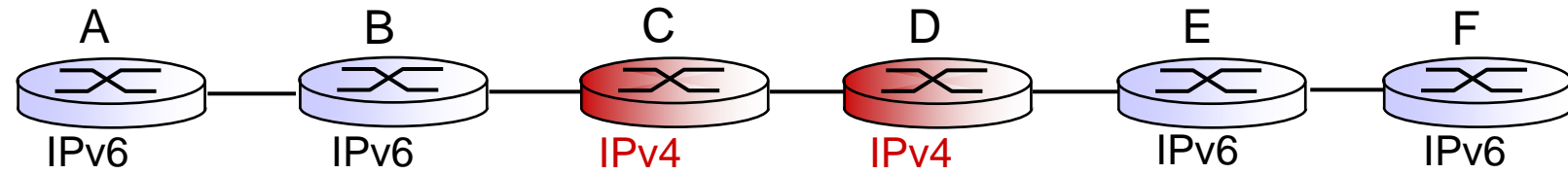


Tunneling

logical view:

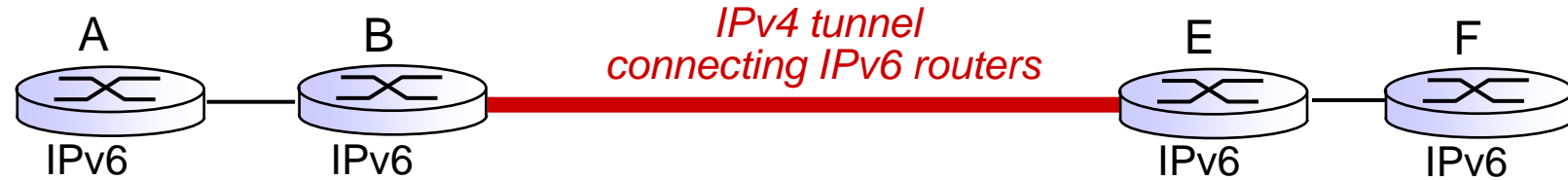


physical view:

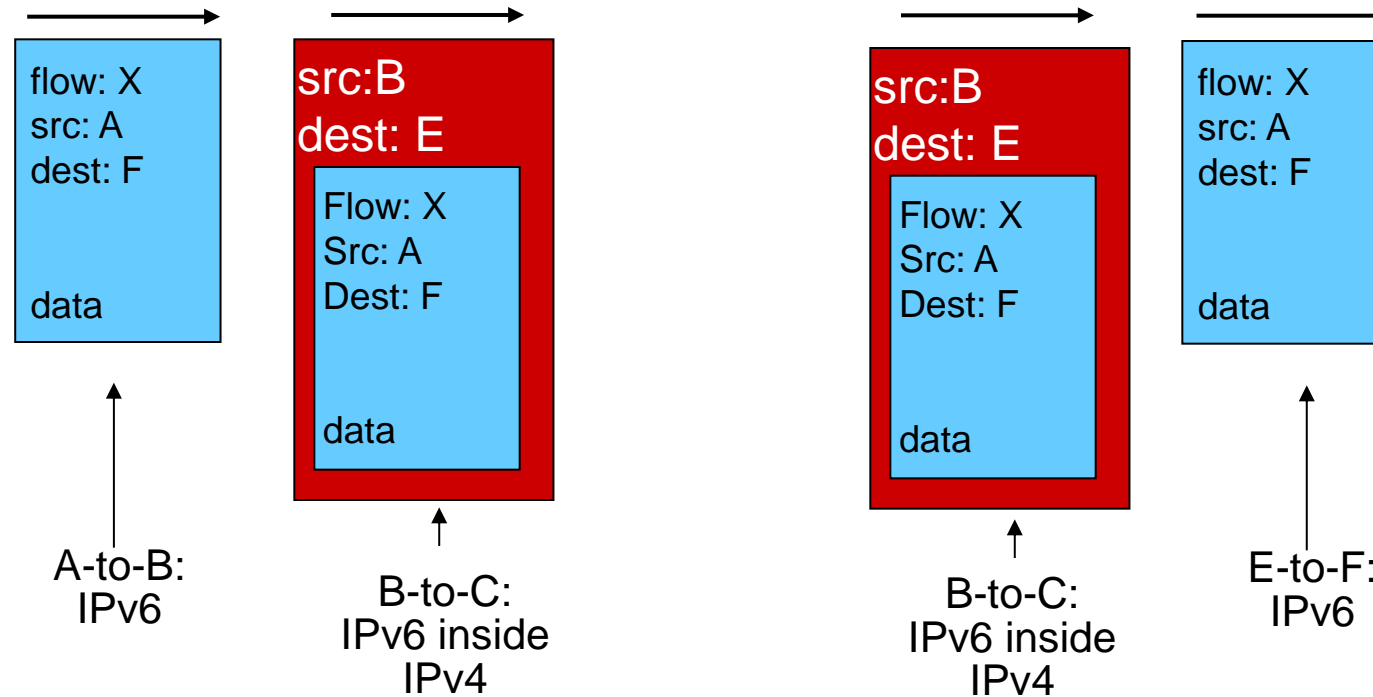
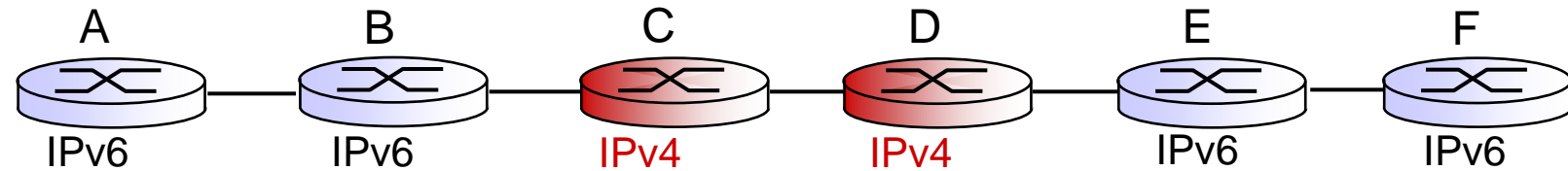


Tunneling

logical view:



physical view:



IPv6: adoption

- ❑ Google: 8% of clients access services via IPv6
- ❑ NIST: 1/3 of all US government domains are IPv6 capable
- ❑ *Long (long!) time for deployment, use*
 - 20 years and counting!
 - think of application-level changes in last 20 years: WWW, Facebook, streaming media, Skype, ...
 - *Why?*