

计算机网络

第四章：网络层

阮晓龙

13938213680 / rxl@hactcm.edu.cn
<http://ke.51xueweb.cn>

河南中医学院管理信息工程学科
河南中医学院网络信息中心

2014.3

本章教学计划

- 网络层提供的两种服务
- 网际协议（IP）
- 划分子网和构建超网（IP地址管理）
- 网际控制报文协议（ICMP）
- 路由选择协议（RIP、OSPF、BGP）

基础内容

- IP多播（IP多播和IGMP）
- 虚拟专用网（VPN）
- 网络地址转换（NAT）

扩展部分

本章教学计划

- 讨论多个网络通过路由器互连成为一个互联网的问题。
- 重点内容：
 - 虚拟互连网络
 - IP地址与物理地址的关系
 - IP地址分类与管理的方法
 - 路由选择协议
 - VPN的基本概念
 - NAT的基本原理

1.网络层提供的两种服务

- 在计算机网络领域，网络层应该向运输层提供怎样的服务（“面向连接”还是“无连接”）曾引起了长期的争论。
- 争论焦点的实质就是：在计算机通信中，可靠交付应当由谁来负责？是网络还是端系统？

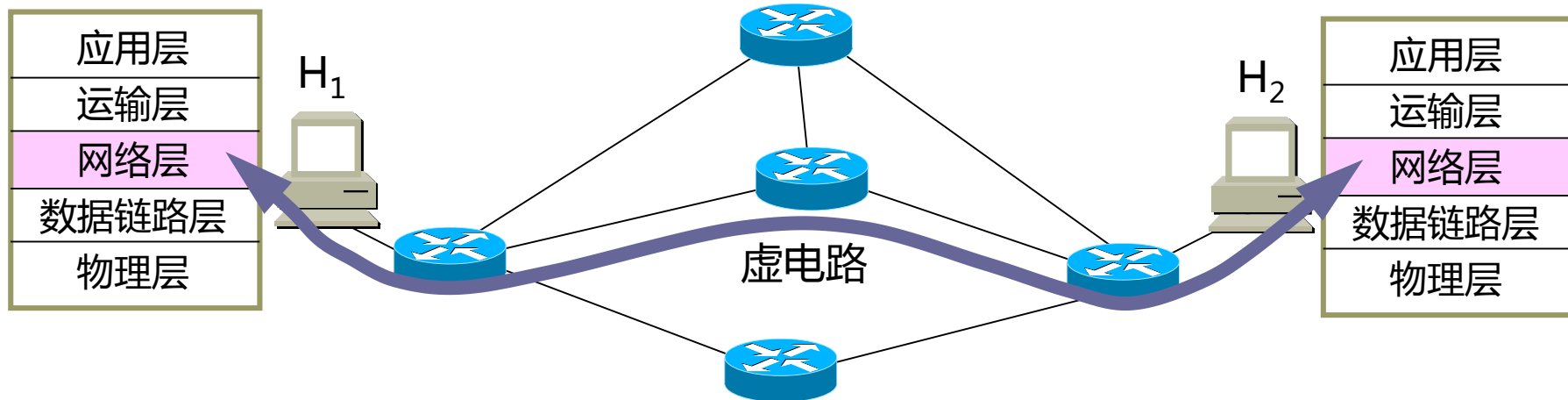
1.网络层提供的两种服务

1.1虚电路

- 电信网的成功经验：让网络负责可靠交付。
- 电信网使用昂贵的程控交换机，用**面向连接**的通信方式，使电信网络能够向用户提供可靠传输的服务。
- 虚电路VC就是当两个计算机进行通信时，应当先建立连接，以**保证双方通信所需的一切网络资源**，双方就沿着已建立的虚电路发送分组。
- 如果使用可靠传输的网络协议，就可使所发送的分组**无差错**按序到达终点。

1. 网络层提供的两种服务

1.1 虚电路



H_1 发送给 H_2 的所有分组都沿着同一条虚电路传送

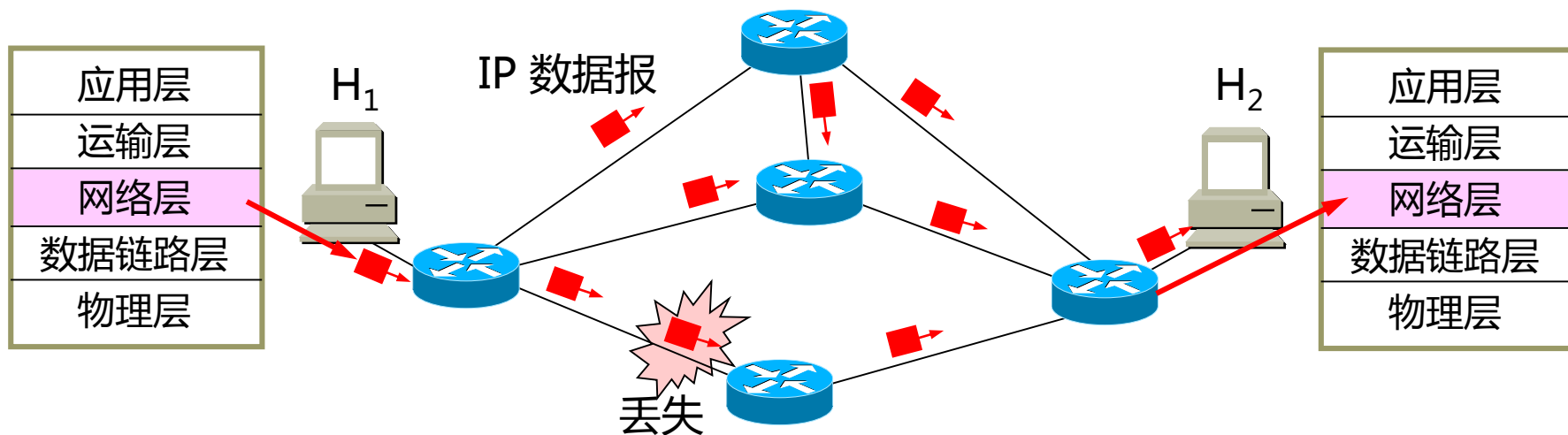
1.网络层提供的两种服务

1.2数据报

- 因特网不提供端到端的可靠服务的优势：
 - 网络中的路由器可以做得比较简单，而且价格低廉。
 - 如果主机（即端系统）中的进程之间的通信需要是可靠的，那么就由运输层负责（包括差错处理、流量控制等）。
 - 网络的造价大大降低，运行方式灵活，能够适应多种应用。
- 因特网能够发展到今日的规模，充分证明了当初采用这种设计思路的正确性。

1. 网络层提供的两种服务

1.2 数据报



H₁发送给H₂的分组可能沿着不同路径传送

1. 网络层提供的两种服务

1.3 虚电路服务与数据报服务对比

对比的方面	虚电路服务	数据报服务
设计思路	可靠通信应当由网络来保证	可靠通信应当由用户主机来保证
连接的建立	必须有	不需要
终点地址	仅在连接建立阶段使用，每个分组使用短的虚电路号	每个分组都有终点的完整地址
分组的转发	属于同一条虚电路的分组均按照同一路由进行转发	每个分组独立选择路由进行转发
当结点出故障时	所有通过出故障的结点的虚电路均不能工作	出故障的结点可能会丢失分组，一些路由可能会发生变化
分组的顺序	总是按发送顺序到达终点	到达终点时不一定按发送顺序
端到端的差错处理和流量控制	可以由网络负责，也可以由用户主机负责	由用户主机负责

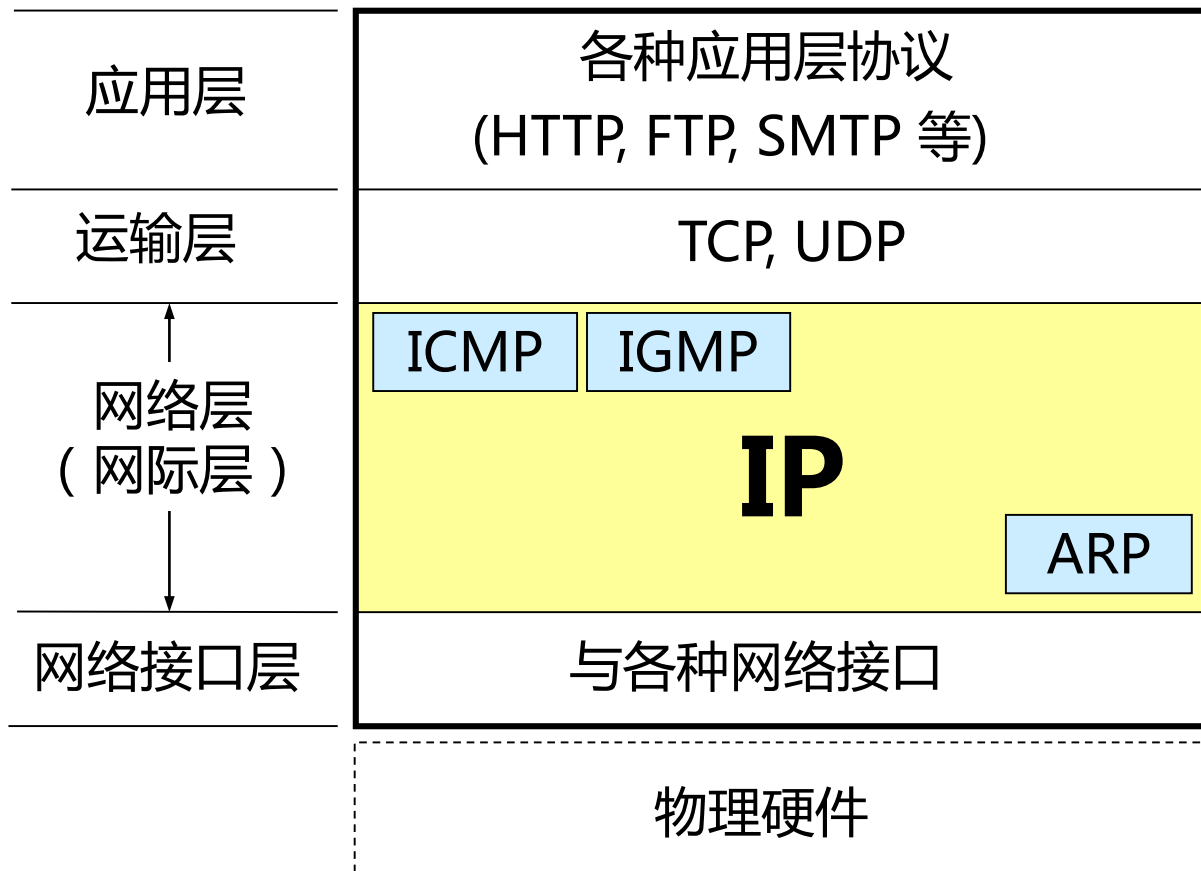
2.网际协议（IP）

- 网际协议（Internet Protocol，IP），或称互联网协议，是用于报文交换网络的一种面向数据的协议。
- IP是在TCP/IP协议中网络层的主要协议，任务是仅仅根据源主机和目的主机的地址传送数据。为此目的，IP定义了寻址方法和数据报的封装结构。
- IP的第一个架构的主要版本，现在称为IPv4，仍然是最主要的互联网协议。尽管世界各地正在积极部署IPv6，但IPv4仍然是最重要的网际协议。

2.网际协议 (IP)

- 网际协议 IP 是 TCP/IP 体系中两个最主要的协议之一。与 IP 协议配套使用的还有三个协议：
 - 地址解析协议 **ARP**
(Address Resolution Protocol)
 - 网际控制报文协议 **ICMP**
(Internet Control Message Protocol)
 - 网际组管理协议 **IGMP**
(Internet Group Management Protocol)

2.网际协议 (IP)



2.网际协议（IP）

2.1虚拟互连网络

- 互连在一起的网络要进行通信，会遇到许多问题需要解决，如：
 - 不同的寻址方案
 - 不同的最大分组长度
 - 不同的网络接入机制
 - 不同的超时控制
 - 不同的差错恢复方法
 - 不同的状态报告方法
 - 不同的路由选择技术
 - 不同的用户接入控制
 - 不同的服务（面向连接服务和无连接服务）
 - 不同的管理与控制方式

2.网际协议（IP）

2.1虚拟互连网络

- 因为用户的需求是多种多样的，所以没有一种单一的网络能够适应所有用户的需求。
- 网络技术是不断发展的，网络的制造厂家也要不停止的推出新产品，以获得更加的市场份额和持续利润。
- 市场上有不同性能、不同网络协议的网络，分布在不同的位置，由不同的组织和人员来管理。

2.网际协议（IP）

2.1虚拟互连网络

- 从一般概念上来讲，将网络互相连接起来要使用一些中间设备。根据中间设备所在的层次，可以有以下四种不同的中间设备：
 - 转发器（repeater）：物理层中继系统。
 - 网桥（桥接器，bridge）：数据链路层中继系统。
 - 路由器（router）：网络层中继系统。
 - 桥路由器（brouter）：数据链路层和网络层混合中继系统。
 - 网关（gateway）：网络层以上使用的中继系统。

2.网际协议（IP）

2.1虚拟互连网络

- 在市场上，主要的网络中间设备都有比较典型的产品名称。

中间设备	工作层次	主要产品
转发器（repeater）	物理层	集线器
网桥（桥接器，bridge）	数据链路层	交换机，二层交换机
路由器（router）	网络层	路由器
桥路器（brouter）	数据链路层和网络层	路由交换机，三层交换机
网关（gateway）	网络层以上	网关、七层交换机

2.网际协议（IP）

2.1虚拟互连网络

- 当中继系统是转发器或网桥时，一般并不称之为网络互连，因为这仅仅是把一个网络扩大了，而这仍然是一个网络。
- 网关由于比较复杂，目前使用得较少。
- 我们讨论网络互联都是指用路由器进行网络互联和路由选择。
- 路由器就是一台专用计算机，用来在互联网中进行路由选择。
 - 由于历史的原因，许多有关 TCP/IP 的文献将网络层使用的路由器称为网关。

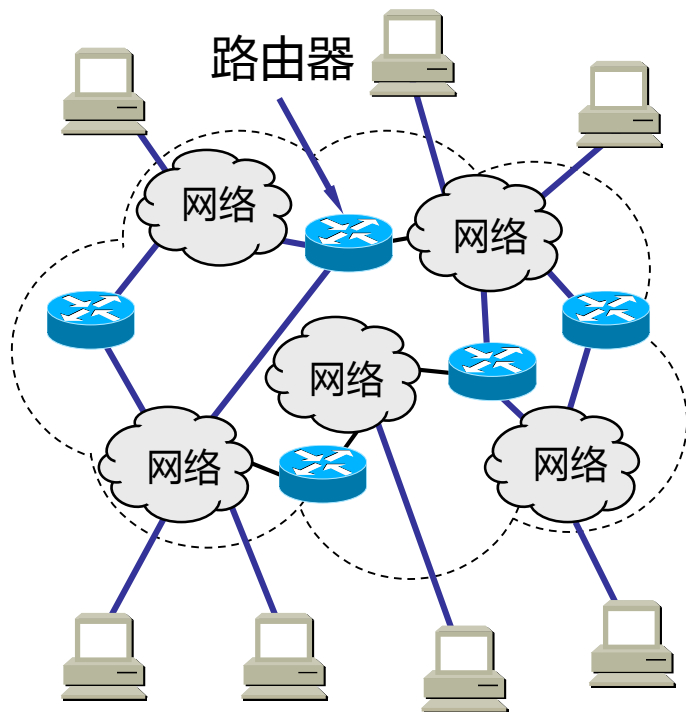
2.网际协议（IP）

2.1虚拟互连网络

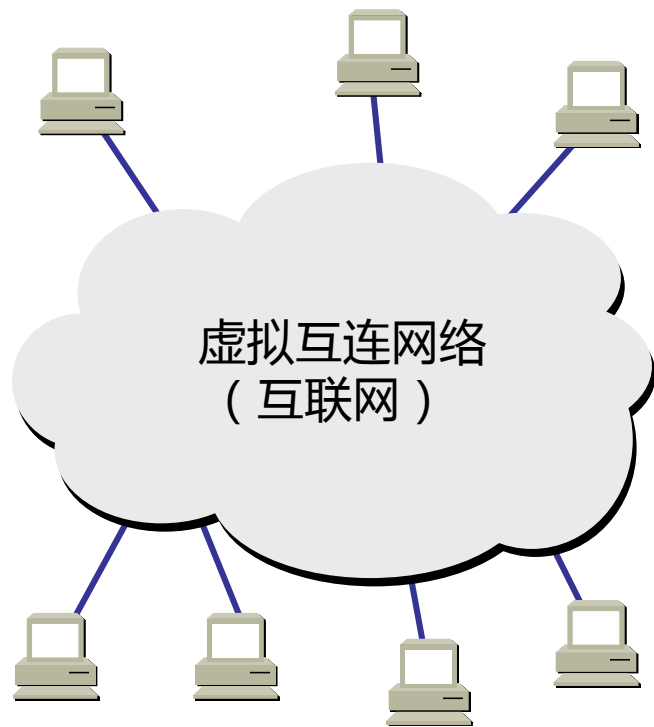
- TCP/IP体系在网络互连上采用的做法是在网络层使用标准化协议，但相互连接的的网络则可以是异构的。
- 参加互连的计算机网络都采用相同的网际协议（IP），因此可以把互连以后的计算机网络看成为一个虚拟互连网络（internet）。

2.网际协议 (IP)

2.1虚拟互连网络



(a) 互连网络



(b) 虚拟互连网络

2.网际协议（IP）

2.1虚拟互连网络

- 所谓**虚拟互连网络**也就是**逻辑互连网络**，它的意思就是互连起来的各种物理网络的异构性本来是客观存在的，但是我们利用IP协议就可以使这些性能各异的网络从用户看起来好像是一个统一的网络。
- 使用IP协议的虚拟互连网络可简称为IP网。
- 使用虚拟互连网络的好处是：当互联网上的主机进行通信时，就好像在一个网络上通信一样，而看不见互连的各具体的网络异构细节。

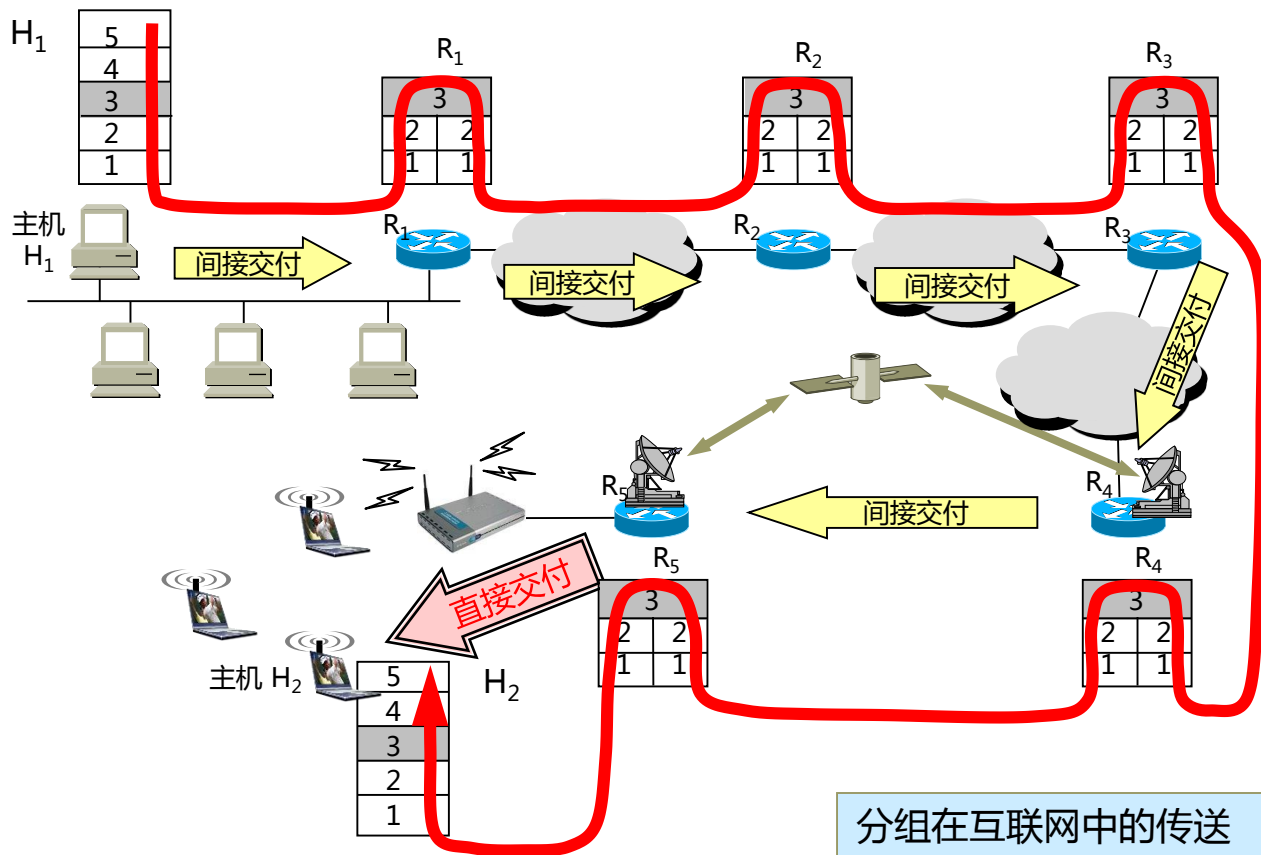
2.网际协议（IP）

2.1虚拟互连网络

- 互联网可以由多种异构网络互连组成。
- 在网络上，两台主机通信分类两种形式：
 - **直接交付**：在一个物理网络上，数据报被源主机直接传送到目标主机上。
 - **间接交付**：当源主机和目标主机分别处于不同的物理网络上时，数据报由源主机通过网络上的路由器间接的传送到目标主机上。

2.网际协议 (IP)

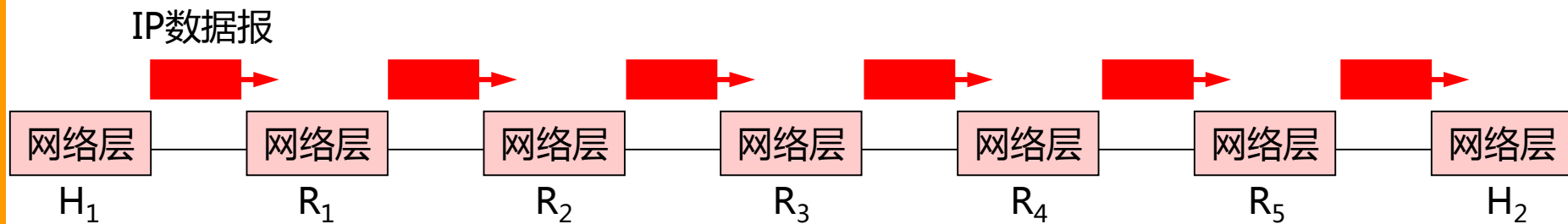
2.1虚拟互连网络

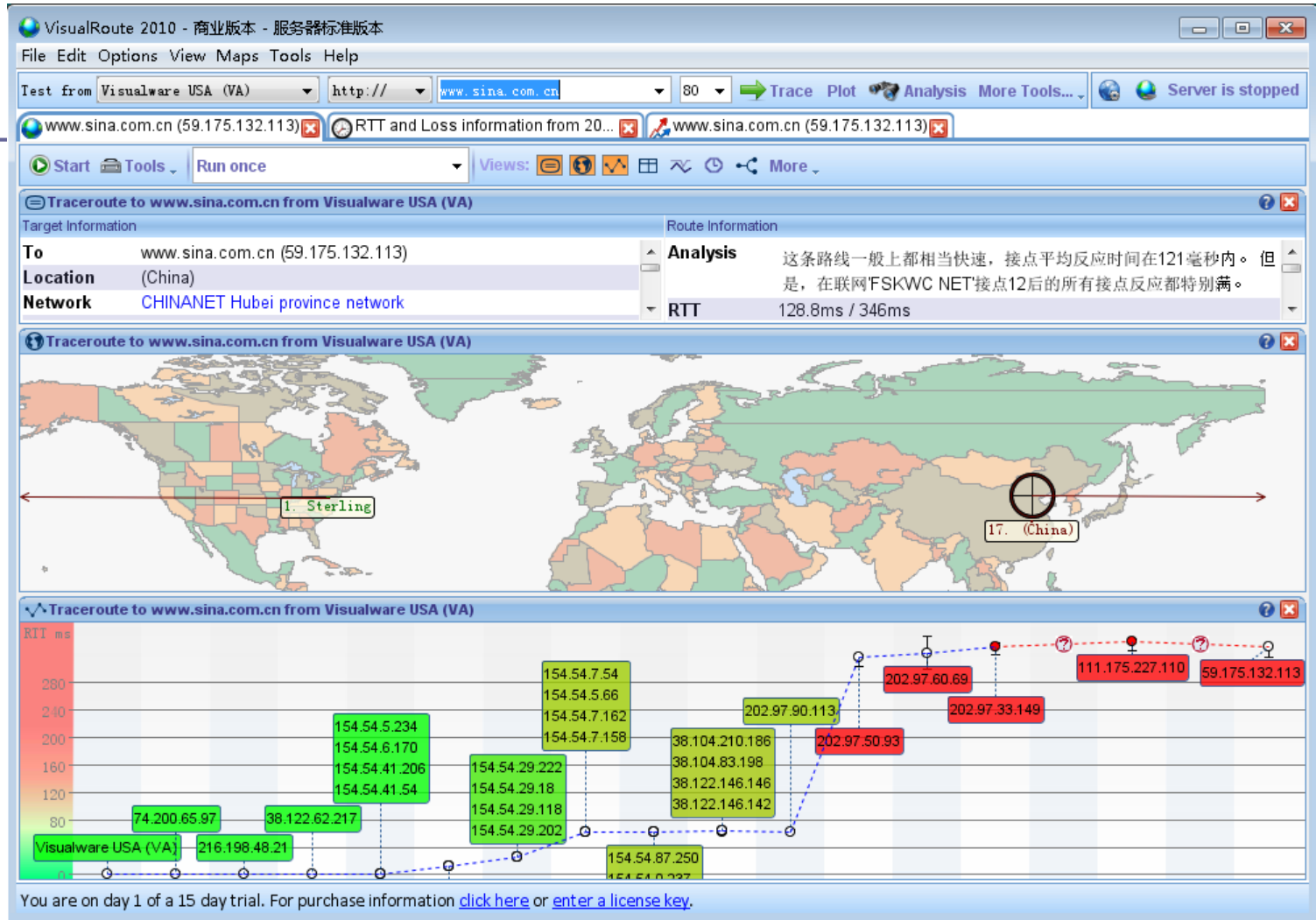


分组在互联网中的传送

2.网际协议 (IP)

2.1虚拟互连网络





2.网际协议 (IP)

2.2分类的IP地址

□ IP地址及其表示方法

- 把整个因特网看成为一个单一的、抽象的网络。
- IP 地址就是给每个连接在因特网上的主机 (或路由器) 分配一个在全世界范围是唯一的 32 位的标识符。
- IP 地址现在由因特网名字与号码指派公司ICANN (Internet Corporation for Assigned Names and Numbers)进行分配。

2.网际协议（IP）

2.2分类的IP地址

□ IP地址及其表示方法

- IP地址的编址方法共经过了三个历史阶段：

分类的IP地址：这是最基本的编址方法，在1981年就通过了相应的标准协议。

子网的划分：这是对最基本的编址方法的改进，其标准[RFC 950]在1985年通过。

构成超网：这是比较新的无分类编址方法。1993年提出后很快就得到推广应用。

2.网际协议（IP）

2.2分类的IP地址

□ 分类的IP地址

- 分类的IP地址就是将IP地址划分为若干个固定类，每一类地址都有两个字段组成：网络号、主机号。
- **网络号（net-id）**：标志主机（或路由器）所连接到的网络。一个网络号在整个因特网范围内必须是唯一的。
- **主机号（host-id）**：标志该主机（或路由器）。一个主机号在它前面的网络号所指明的网络范围内必须是唯一的。
- 一个IP地址在整个因特网范围内必须是唯一的。

2.网际协议 (IP)

2.2分类的IP地址

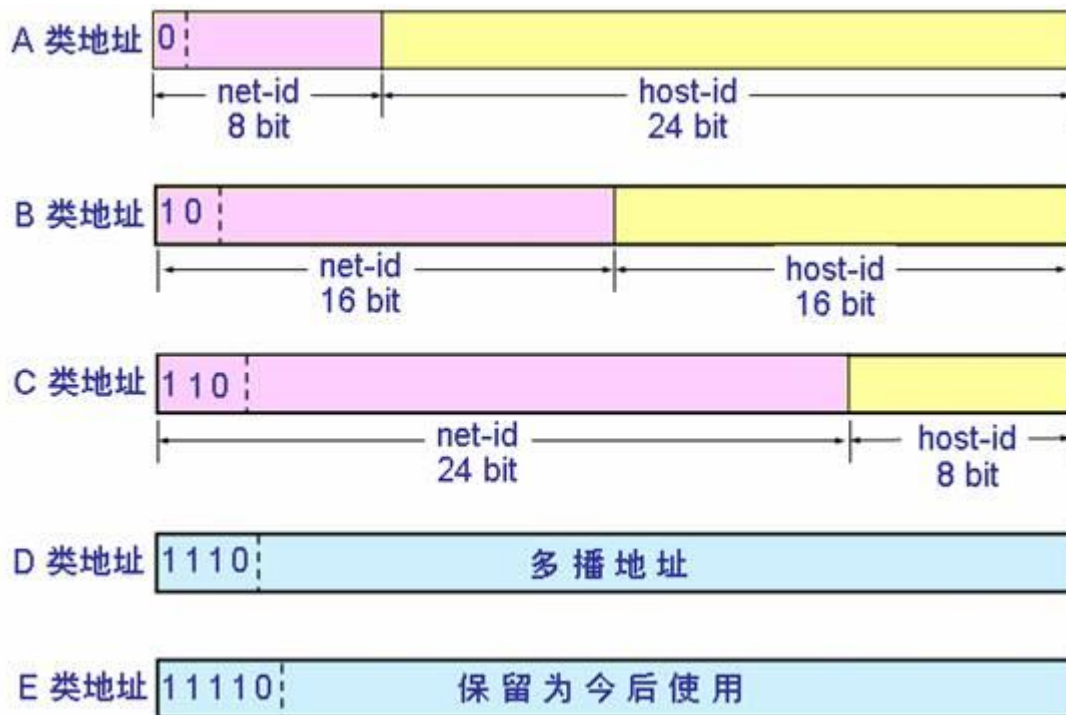
- 分类的IP地址
 - 两级的IP地址可以“定义为”：

IP地址 ::= { <网络号>, <主机号> }

2.网际协议 (IP)

2.2分类的IP地址

□ 分类的IP地址



2.网际协议（IP）

2.2分类的IP地址

- 从IP地址的结构来看，IP地址并不仅仅指明一个主机，而是还指明了主机所连接到的网络。
- 当初，把IP地址划分为A、B、C三个类别，是因为在现实中，有的网络拥有很多主机，有的网络上的主机很少。把IP地址划分为A、B、C三个类别，能够更好的满足用户的不同需求。
- 当某单位申请到一个IP地址时，实际上获得了具有同样网络号的一组地址。其中主机号是自行分配的，只要无重复即可。

2.网际协议 (IP)

2.2分类的IP地址

□ 点分十进制记法：提高IP地址的可读性

计算机存放的IP地址是连续的二进制代码

1101001101000101001000000010010

每隔8位进行分割,以便于阅读

11010011 01000101 00100000 00010010

将八位二进制数转为十进制数

211

69

32

18

点分十进制记录IP地址

211.69.32.18

2.网际协议 (IP)

2.2分类的IP地址

□ 常用的三种类别的IP地址

IP 地址的指派范围

网络类别	最大可指派的网络数	第一个可指派的网络号	最后一个可指派的网络号	每个网络中的最大主机数
A	$126 (2^7 - 2)$	1	126	16777214
B	$16383 (2^{14} - 1)$	128.1	191.255	65534
C	$2097151 (2^{21} - 1)$	192.0.1	223.255.255	254

2.网际协议 (IP)

2.2分类的IP地址

□ 常用的三种类别的IP地址

一般不使用的特殊 IP 地址

网络号	主机号	源地址使用	目的地址使用	代表的意义
0	0	可以	不可	在本网络上的本主机 (见 6.6 节 DHCP 协议)
0	host-id	可以	不可	在本网络上的某个主机 host-id
全 1	全 1	不可	可以	只在本网络上进行广播 (各路由器均不转发)
net-id	全 1	不可	可以	对 net-id 上的所有主机进行广播
127	非全 0 或全 1 的任何数	可以	可以	用作本地软件环回测试之用

2.网际协议（IP）

2.2分类的IP地址

□ IP地址具有的重要特点：

- **每一个IP地址都由网络号和主机号两部分组成。**IP地址是一种分等级的地址结构。分等级的两个好处是：第一，IP地址管理机构在分配IP地址时只分配网络号，而剩下的主机号则由得到该网络号的单位自行分配。这样就方便了IP地址的管理。第二，路由器仅根据目的主机所连接的网络号来转发分组（而不考虑目的主机号），这样就可以使路由表中的项目数大幅度减少，从而减小了路由表所占的存储空间。

2.网际协议（IP）

2.2分类的IP地址

□ IP地址具有的重要特点：

- 实际上IP地址是标志一个主机（或路由器）和一条链路的接口。当一个主机同时连接到两个网络上时，该主机就必须同时具有两个相应的IP地址，其网络号net-id必须是不同的。这种主机称为多归属主机(multihomed host)。由于一个路由器至少应当连接到两个网络（这样它才能将IP数据报从一个网络转发到另一个网络），因此一个路由器至少应当有两个不同的IP地址。

2.网际协议（IP）

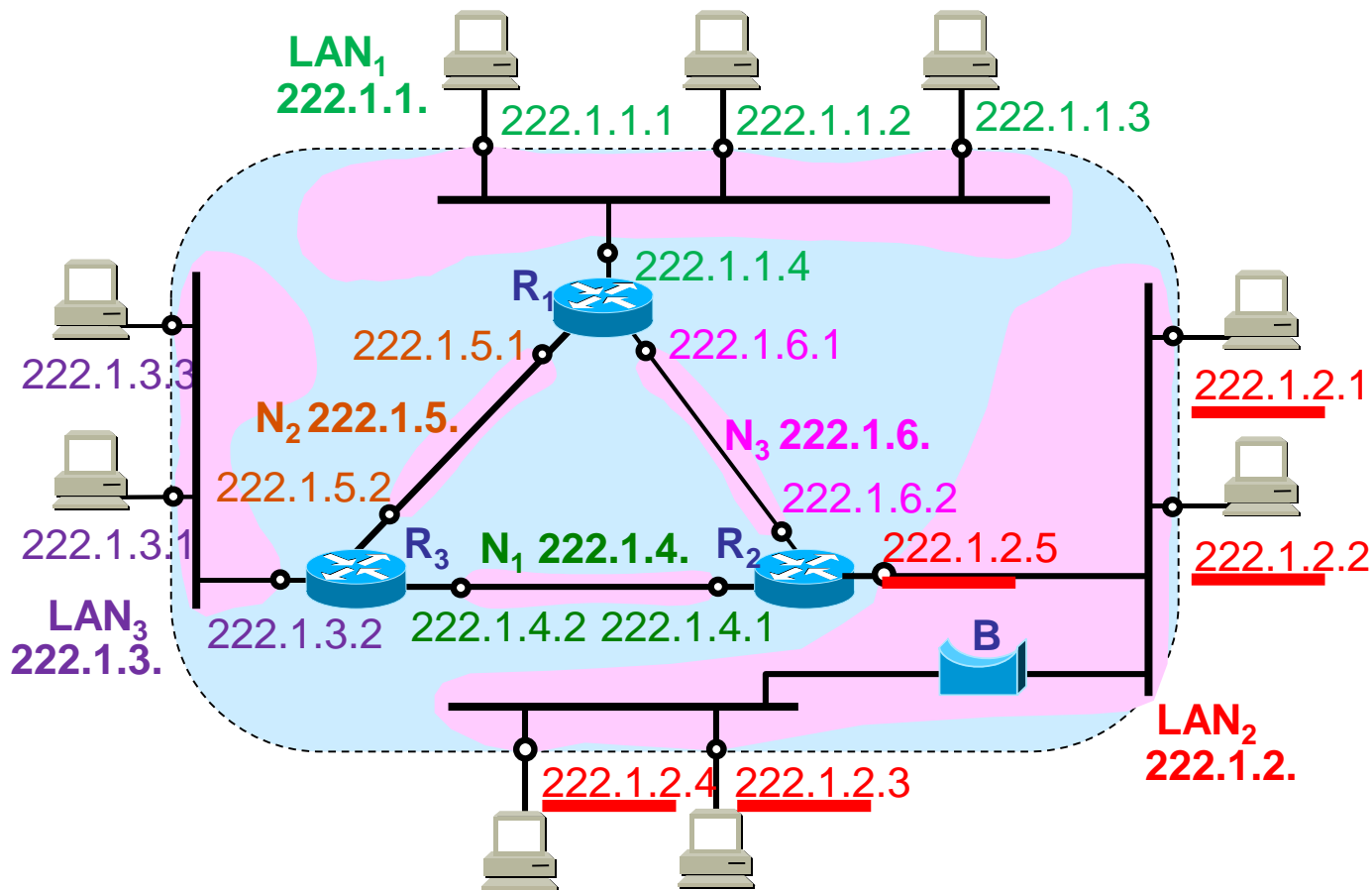
2.2分类的IP地址

□ IP地址具有的重要特点：

- 一个网络是指具有相同网络号net-id的主机的集合。用转发器或网桥连接起来的若干个局域网仍为一个网络，因此这些局域网都具有同样的网络号net-id。
- 在IP地址中，所有分配到网络号的网络都是平等的。所有分配到网络号net-id的网络，范围很小的局域网，还是可能覆盖很大地理范围的广域网，都是平等的。

2.网际协议 (IP)

2.2分类的IP地址



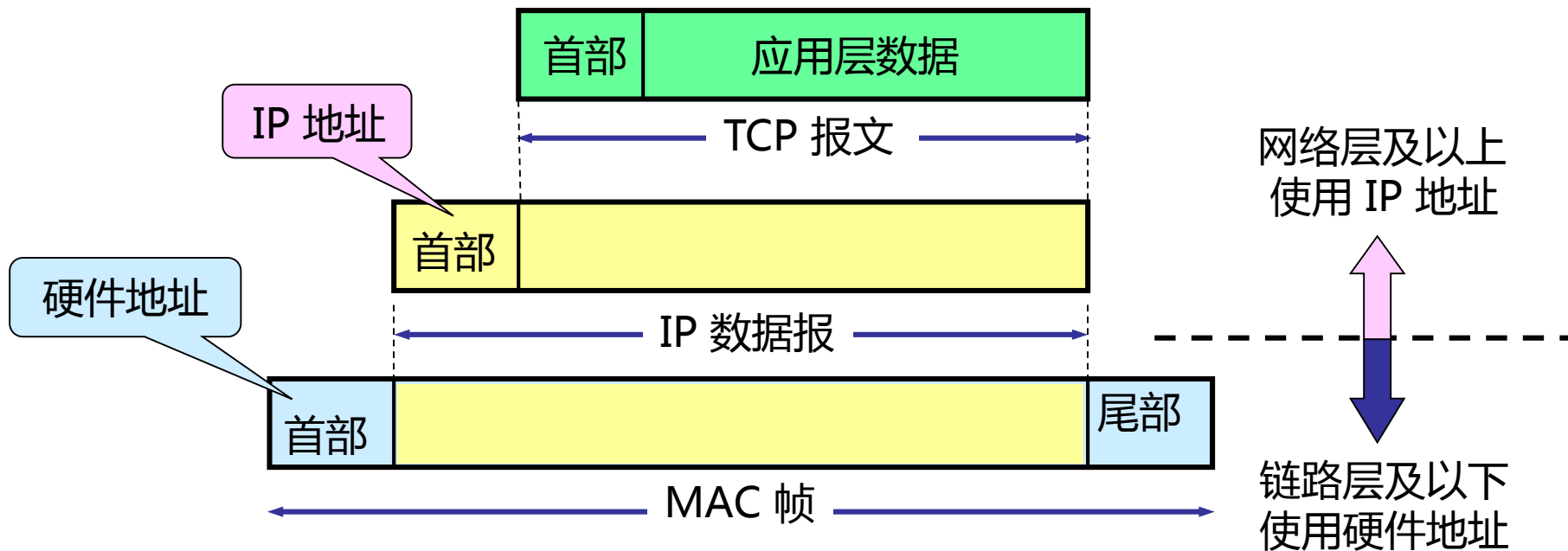
2.网际协议（IP）

2.3 IP地址与硬件地址

- IP地址和硬件地址的区别：
 - 物理地址是数据链路层和物理层使用的地址。
 - IP地址是网络层和以上各层使用的地址，IP地址是一种逻辑地址。
- 在发送数据时，数据从高层下到低层，然后才到通信链路上传输。使用IP地址的IP数据报一旦交给了数据链路层，就被封装成MAC帧。MAC帧在传送时使用的源地址和目的地址都是硬件地址。
- 连接在通信链路上的设备在接受MAC帧时，看不到IP地址。只有把数据帧提交给网络层，才能够识别到IP地址。

2. 网际协议 (IP)

2.3 IP地址与硬件地址



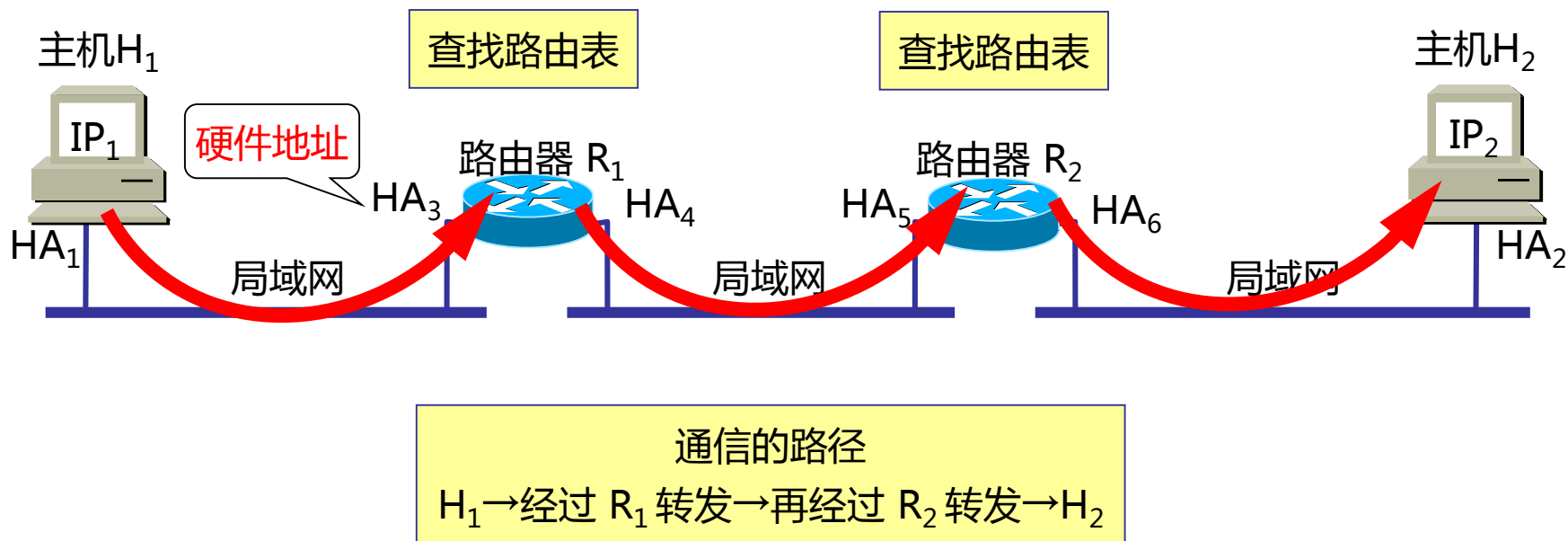
2.网际协议（IP）

2.3 IP地址与硬件地址

- IP地址放在IP数据报的首部，硬件地址则放在MAC帧的首部。
- 在网络层和网络层以上使用IP地址，在数据链路层使用硬件地址。
- 当IP数据报放入数据链路层的MAC帧中以后，整个的IP数据报就成为MAC帧的数据，因而在数据链路层看不到数据报的IP地址信息。

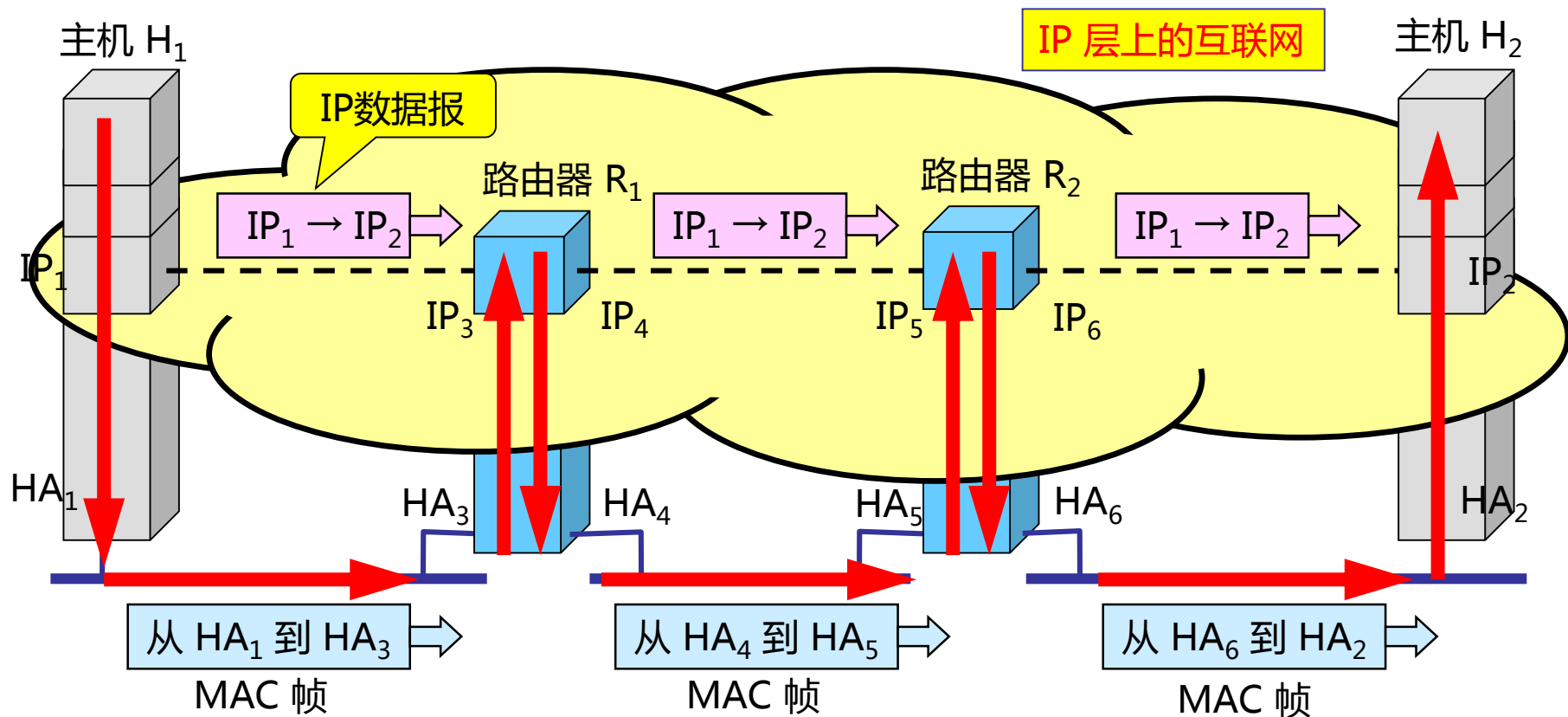
2.网际协议（IP）

2.3 IP地址与硬件地址



2. 网际协议 (IP)

2.3 IP地址与硬件地址



2.网际协议（IP）

2.3 IP地址与硬件地址

不同层次、不同区间的源地址和目的地址

网络层次中 地址信息 不同的 通信阶段	网络层 写入 IP 数据报首部的地址		数据链路层 写入 MAC 帧首部的地址	
	源地址	目的地址	源地址	目的地址
从 H1 到 R1	IP1	IP2	HA1	HA3
从 R1 到 R2	IP1	IP2	HA4	HA5
从 R2 到 H2	IP1	IP2	HA6	HA2

2.网际协议 (IP)

2.3 IP地址与硬件地址

□ 几个重点和总结：

- 在IP层抽象的互联网上只能看到IP数据报。
- 虽然在IP数据报中有源站IP地址，但是路由器只根据目的站的IP地址的**网络号**进行路由选择。
- 在局域网的链路层，只能看见MAC帧。IP数据报被封装到MAC帧中作为数据部分。
- 尽管互连在一起的网络的硬件地址体系各不相同，但IP层抽象的互联网却屏蔽了下层很复杂的细节。只要在网络层上讨论问题，就能够使用统一的、抽象的IP地址研究主机和主机或路由器之间的通信。

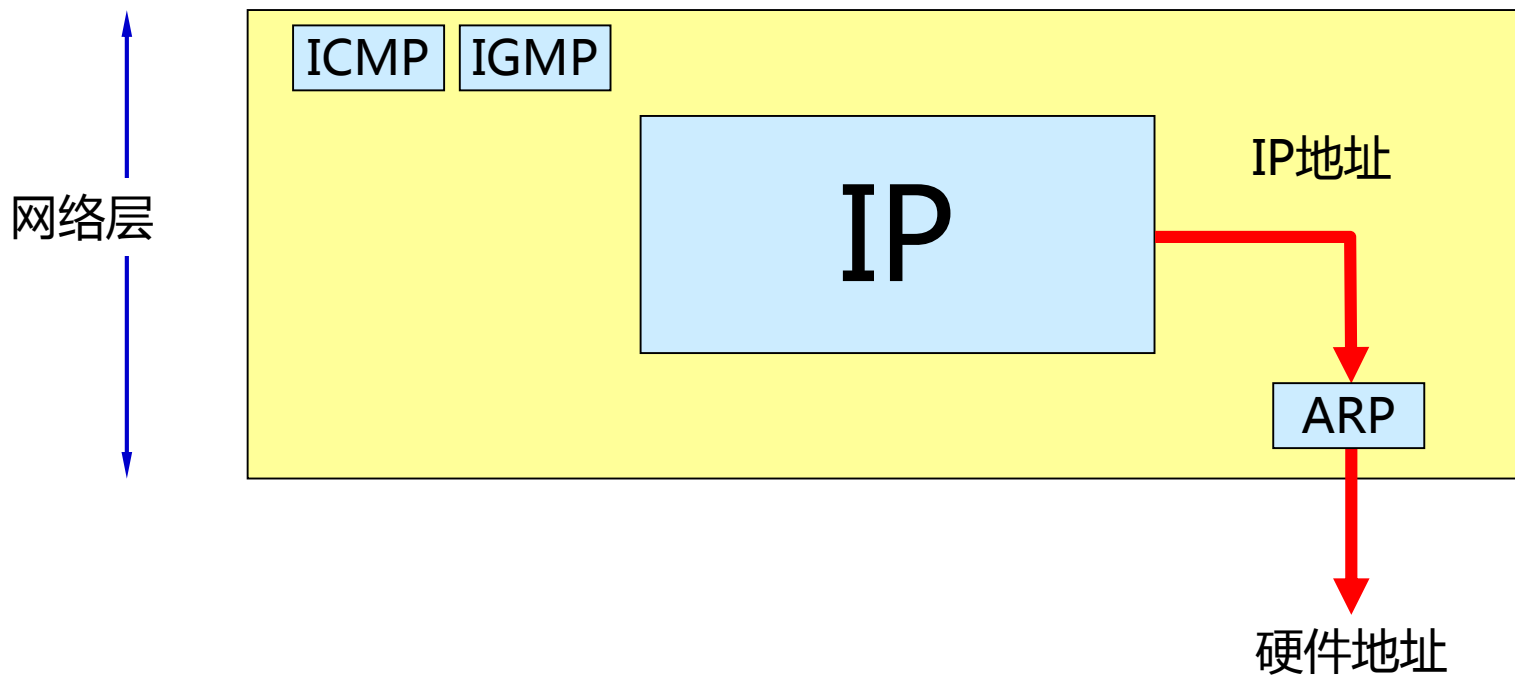
2.网际协议 (IP)

2.4地址解析协议 (ARP)

- ❑ 地址解析协议 (ARP) 的作用是：在知道一个IP地址时，查找到该IP地址对应的硬件地址。
- ❑ 由于IP协议用到了ARP，因此把ARP放到网络层来介绍。有些技术文档和书籍考虑到ARP最终解析的是硬件地址，把ARP放到数据链路层。
- ❑ 不管网络层使用的是什么协议，在实际网络的链路上传送数据帧时，最终还是必须使用硬件地址。

2.网际协议 (IP)

2.4地址解析协议 (ARP)



2.网际协议 (IP)

2.4地址解析协议 (ARP)

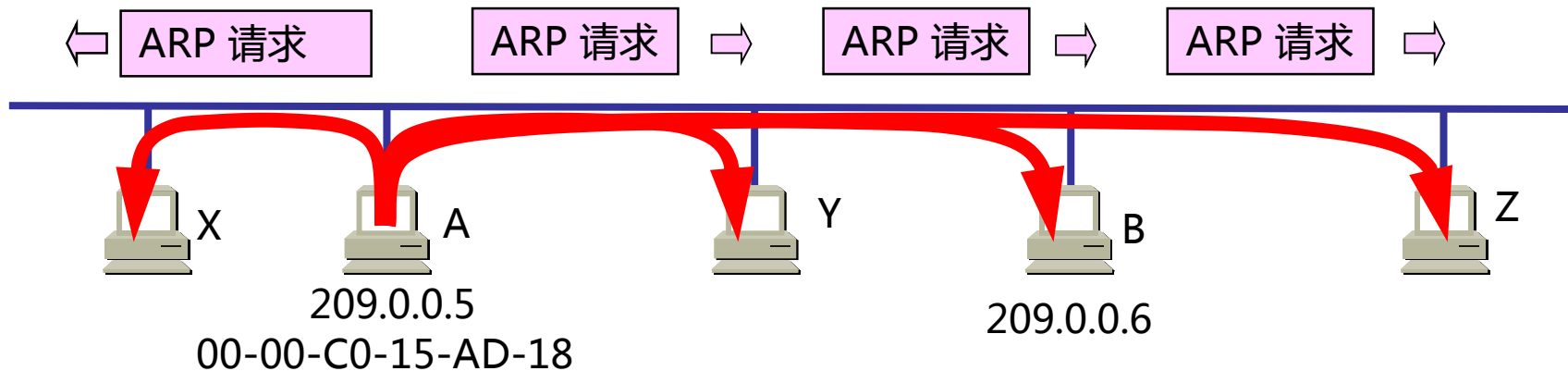
- 每一个主机都设有一个ARP高速缓存(ARP cache), 里面有所在的局域网上的各主机和路由器的IP地址到硬件地址的映射表。
- 当主机A欲向本局域网上的某个主机B发送IP数据报时, 就先在其ARP高速缓存中查看有无主机B的IP地址。如有, 就可查出其对应的硬件地址, 再将此硬件地址写入MAC帧, 然后通过局域网将该MAC帧发往此硬件地址。

2.网际协议 (IP)

2.4地址解析协议 (ARP)

主机A广播发送ARP请求分组

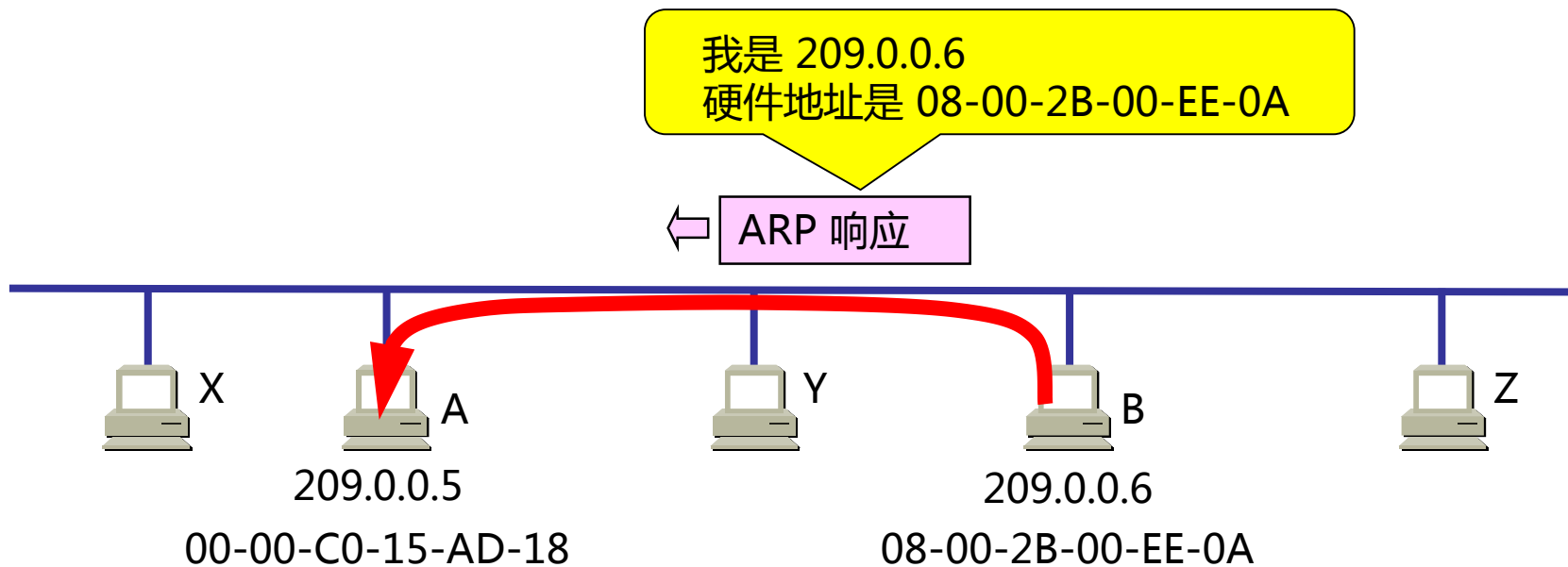
我是 209.0.0.5，硬件地址是 00-00-C0-15-AD-18
我想知道主机 209.0.0.6 的硬件地址



2.网际协议 (IP)

2.4地址解析协议 (ARP)

主机B向A发送ARP响应分组



2.网际协议 (IP)

2.4地址解析协议 (ARP)

- ❑ 为了减少网络上的通信量，主机A在发送其ARP请求分组时，就将自己的IP地址到硬件地址的映射写入ARP请求分组。
- ❑ 当主机B收到A的ARP请求分组时，就将主机A的这一地址映射写入主机B自己的ARP高速缓存中。这对主机B以后向A发送数据报时就更方便了。

2.网际协议（IP）

2.4地址解析协议（ARP）

- ARP是解决同一个局域网上的主机或路由器的IP地址和硬件地址的映射问题。
- 如果所要找的主机和源主机不在同一个局域网，那么就要通过ARP找到一个位于本局域网上的某个路由器的硬件地址，然后把分组发送给这个路由器，让这个路由器把分组转发给下一个网络。剩下的工作就由下一个网络来做。

2.网际协议（IP）

2.4地址解析协议（ARP）

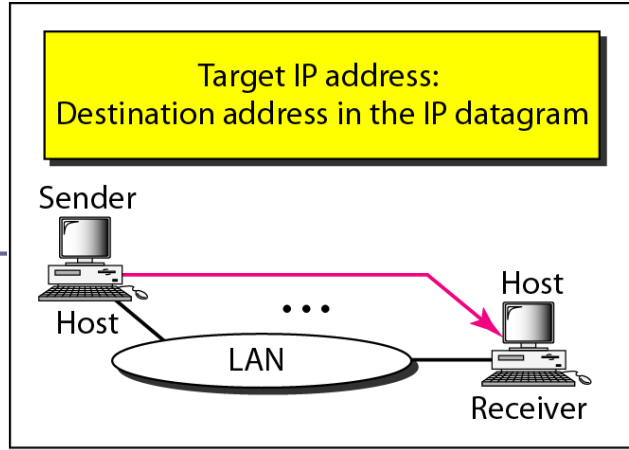
- 从IP地址到硬件地址的解析是自动进行的，主机的用户对这种地址解析过程是不知道的。
- 只要主机或路由器要和本网络上的另一个已知IP地址的主机或路由器进行通信，ARP协议就会自动地将该IP地址解析为链路层所需要的硬件地址。

2.网际协议（IP）

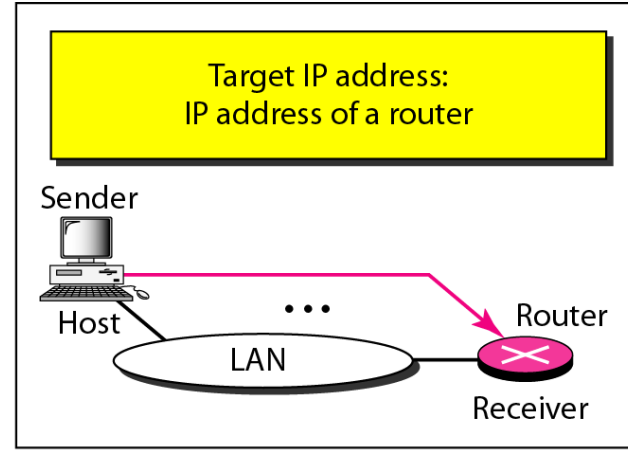
2.4地址解析协议（ARP）

□ 使用ARP的四种典型情况：

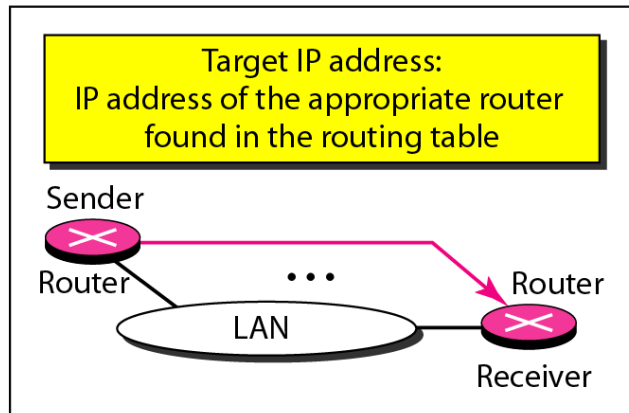
- 发送方是主机，要把IP数据报发送到本网络上的另一个主机。这时用ARP找到目的主机的硬件地址。
- 发送方是主机，要把IP数据报发送到另一个网络上的一个主机。这时用ARP找到本网络上的一个路由器的硬件地址。剩下的工作由这个路由器来完成。
- 发送方是路由器，要把IP数据报转发到本网络上的一个主机。这时用ARP找到目的主机的硬件地址。
- 发送方是路由器，要把IP数据报转发到另一个网络上的一个主机。这时用ARP找到本网络上另一个路由器的硬件地址。剩下的工作由这个路由器来完成。



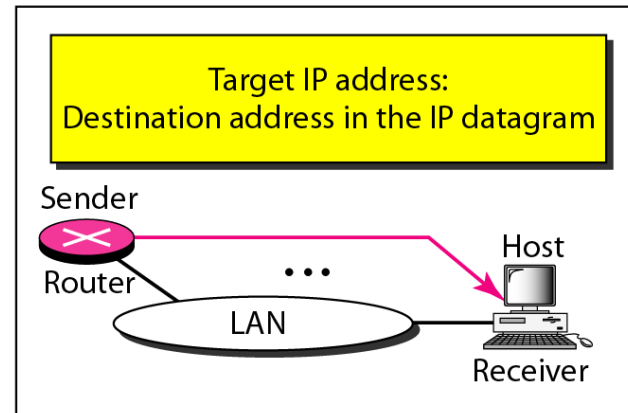
Case 1. A host has a packet to send to another host on the same network.



Case 2. A host wants to send a packet to another host on another network. It must first be delivered to a router.



Case 3. A router receives a packet to be sent to a host on another network. It must first be delivered to the appropriate router.



Case 4. A router receives a packet to be sent to a host on the same network.

2.网际协议（IP）

2.4地址解析协议（ARP）

□ 为什么不直接使用硬件地址进行通信？

- 由于全世界存在着各式各样的网络，它们使用不同的硬件地址。要使这些异构网络能够互相通信就必须进行非常复杂的硬件地址转换工作，因此几乎是不可能的事。
- 连接到因特网的主机都拥有统一的IP地址，它们之间的通信就像连接在同一个网络上那样简单方便，因为调用ARP来寻找某个路由器或主机的硬件地址都是由计算机软件自动进行的，对用户来说是看不见这种调用过程的。

2.网际协议 (IP)

2.4地址解析协议 (ARP)

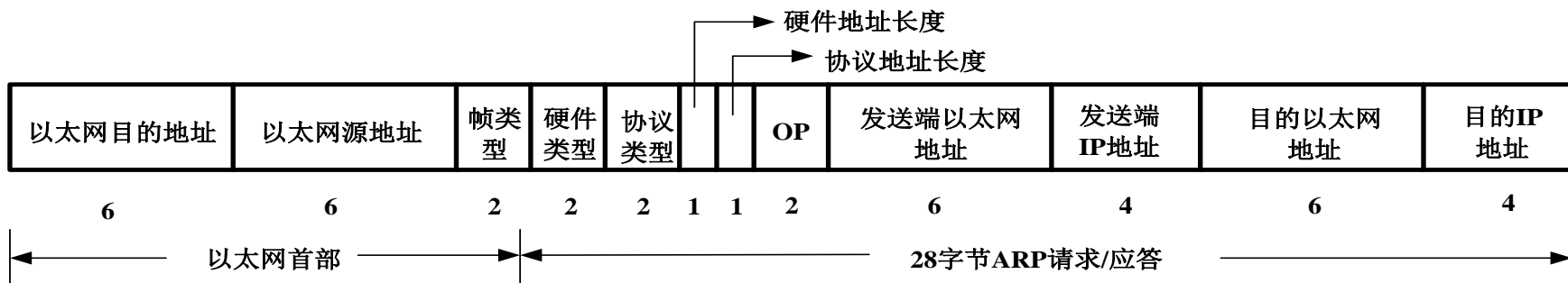
- 如何查看本地主机的ARP高速缓存？
 - 在Linux操作系统中，通过在shell环境下输入“arp”查看，也可以通过arp命令进行更多操作。
 - 在Windows操作系统中，通过【运行】【cmd】，在命令窗体中输入“arp -a”查看，也可以通过arp命令进行更多操作。
- 现场演示arp命令的使用，并介绍arp高速缓存。

```
ruanxiaolong@Teach-Ubuntu-Server-VMs:~$ arp
Address          HWtype  HWaddress      Flags Mask    Iface
211.69.32.1     ether   e4:68:a3:a3:fa:7d  C             eth0
211.69.32.122   ether   00:0c:29:16:25:97  C             eth0
HACTCM-DNS-2    ether   bc:ae:c5:07:7a:82  C             eth0
```


2.网际协议 (IP)

2.4地址解析协议 (ARP)

□ ARP的数据帧格式：



□ 通过Wireshark进行ARP数据报的分析。

2.1

(ARP)

```

3144 39.966249000    HuaweiTe_a3:fa:7d Broadcast ARP 60 who has 211.69.32.179? Tell 211.69.32.1
3146 39.966249000    HuaweiTe_a3:fa:7d Broadcast ARP 60 who has 211.69.32.179? Tell 211.69.32.1
3150 40.123107000    HuaweiTe_a3:fa:7d Broadcast ARP 60 who has 211.69.32.128? Tell 211.69.32.1
-----
Frame 3146: 60 bytes on wire (480 bits), 60 bytes captured (480 bits) on interface 0
Interface id: 0
Encapsulation type: Ethernet (1)
Arrival Time: Apr 3, 2014 15:28:56.006516000 [#####]
[Time shift for this packet: 0.000000000 seconds]
Epoch Time: 1396510136.006516000 seconds
[Time delta from previous captured frame: 0.015001000 seconds]
[Time delta from previous displayed frame: 0.041086000 seconds]
[Time since reference or first frame: 39.966249000 seconds]
Frame Number: 3146
Frame Length: 60 bytes (480 bits)
Capture Length: 60 bytes (480 bits)
[Frame is marked: False]
[Frame is ignored: False]
[Protocols in frame: eth:arp]
[Coloring Rule Name: ARP]
[Coloring Rule String: arp]
Ethernet II, Src: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d), Dst: Broadcast (ff:ff:ff:ff:ff:ff)
  Destination: Broadcast (ff:ff:ff:ff:ff:ff)
    Address: Broadcast (ff:ff:ff:ff:ff:ff)
      .... ..1. .... .. = LG bit: Locally administered address (this is NOT the factory default)
      .... ..1 .... .. = IG bit: Group address (multicast/broadcast)
  Source: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d)
    Address: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d)
      .... ..0. .... .. = LG bit: Globally unique address (factory default)
      .... ..0 .... .. = IG bit: Individual address (unicast)
  Type: ARP (0x0806)
  Padding: 00000000000000000000000000000000
  Address Resolution Protocol (request)
    Hardware type: Ethernet (1)
    Protocol type: IP (0x0800)
    Hardware size: 6
    Protocol size: 4
    Opcode: request (1)
    Sender MAC address: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d)
    Sender IP address: 211.69.32.1 (211.69.32.1)
    Target MAC address: 00:00:00_00:00:00 (00:00:00:00:00:00)
    Target IP address: 211.69.32.179 (211.69.32.179)

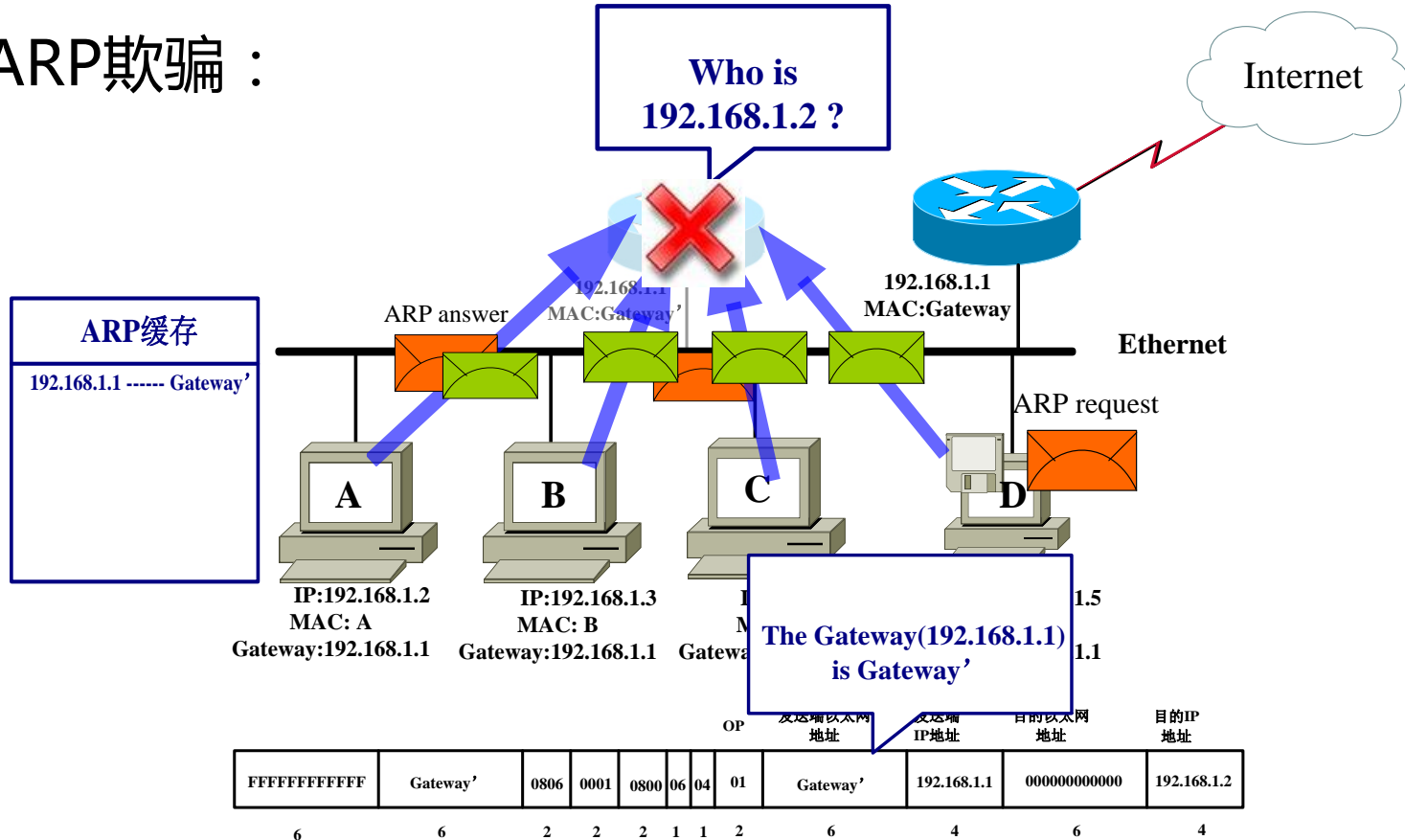
0000 ff ff ff ff ff ff e4 68 a3 a3 fa 7d 08 06 00 01 .....h ...}....
0010 08 00 06 04 00 01 e4 68 a3 a3 fa 7d d3 45 20 01 .....h ...}.E .
0020 00 00 00 00 00 00 d3 45 20 b3 00 00 00 00 00 00 .....E .....
0030 00 00 00 00 00 00 00 00 00 00 00 00 .....

```

2.网际协议 (IP)

2.4地址解析协议 (ARP)

ARP欺骗 :



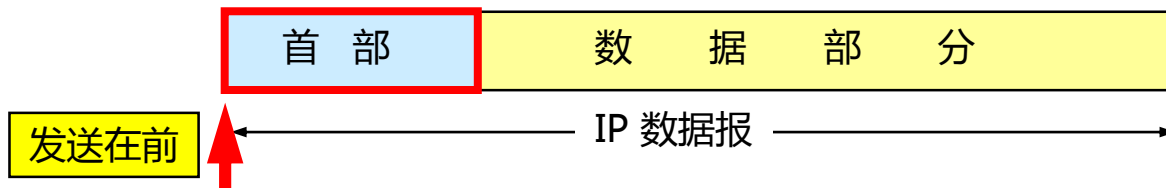
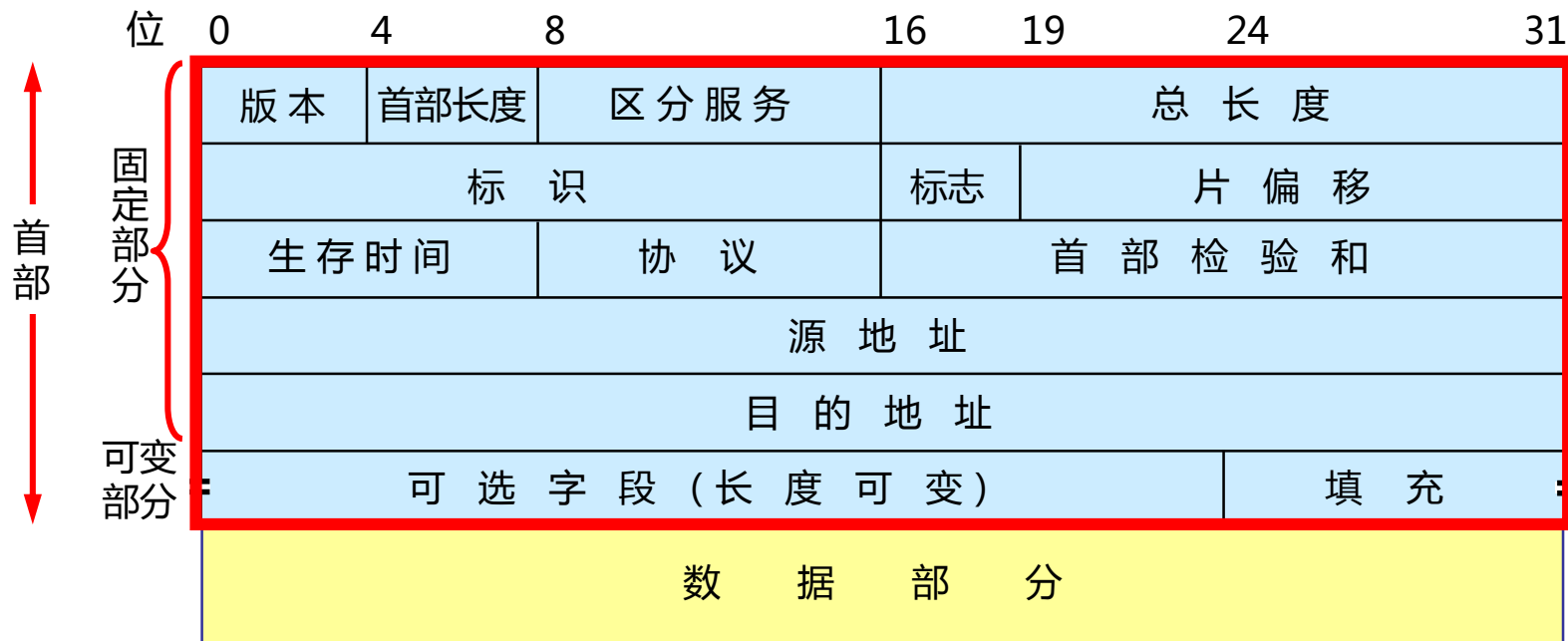
2.网际协议（IP）

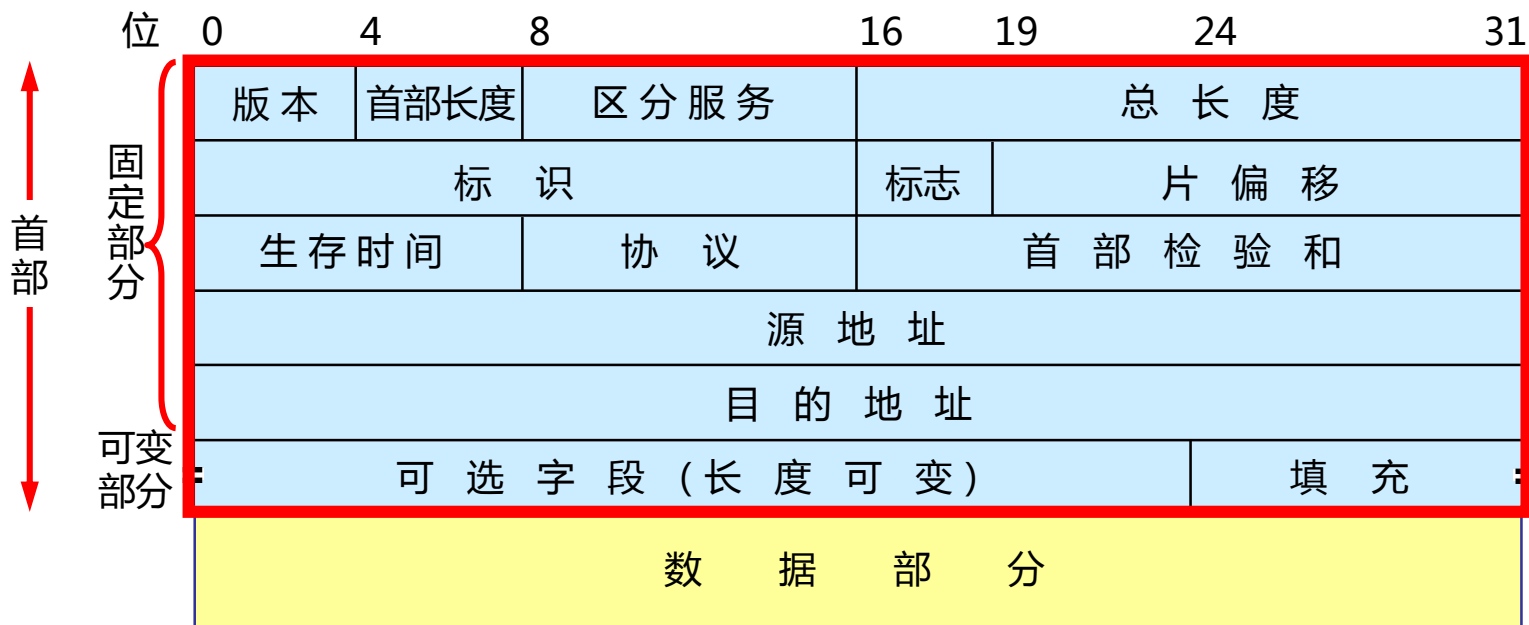
2.5 IP数据报的格式

- 一个IP数据报由首部和数据两部分组成。
- 首部的前一部分是固定长度，共20字节，是所有IP数据报必须具有的。
- 在首部固定部分的后面是一些可选字段，其长度是可变的。
- IP数据报的格式也能够说明了IP协议的功能。

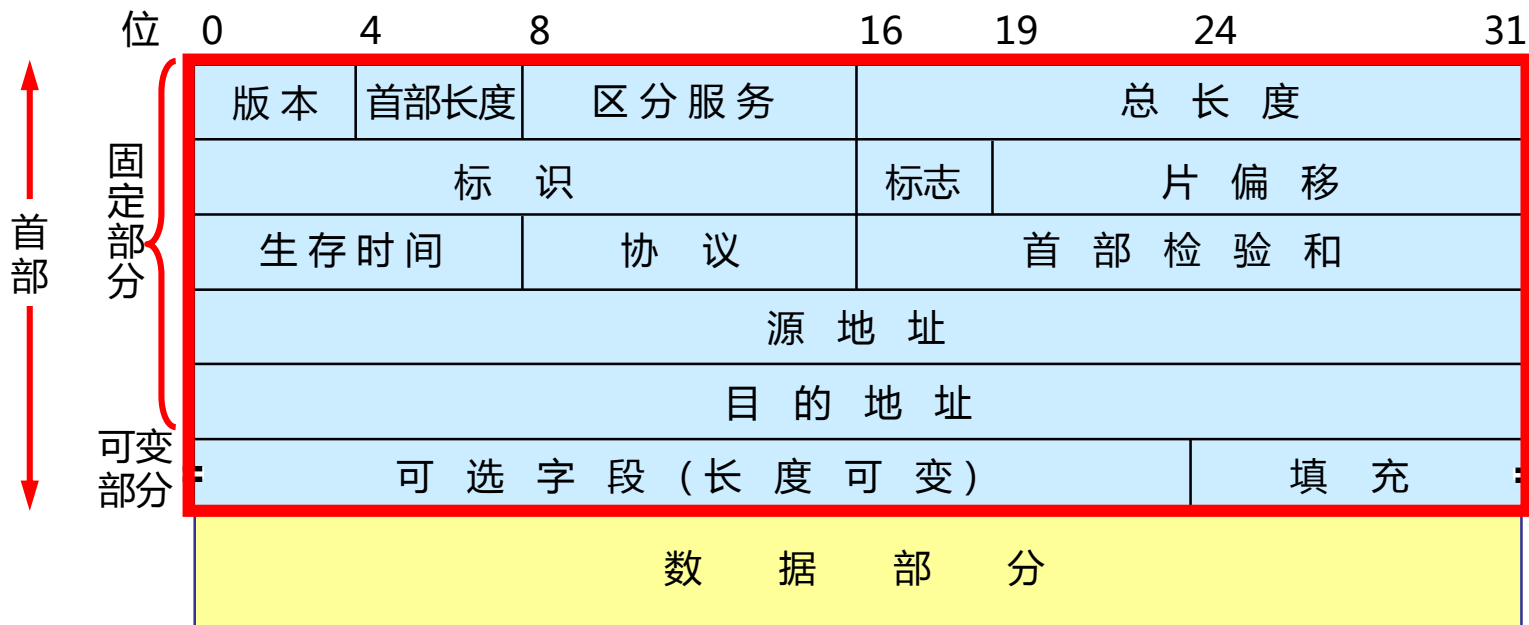
2. 网际协议 (IP)

2.5 IP数据报的格式

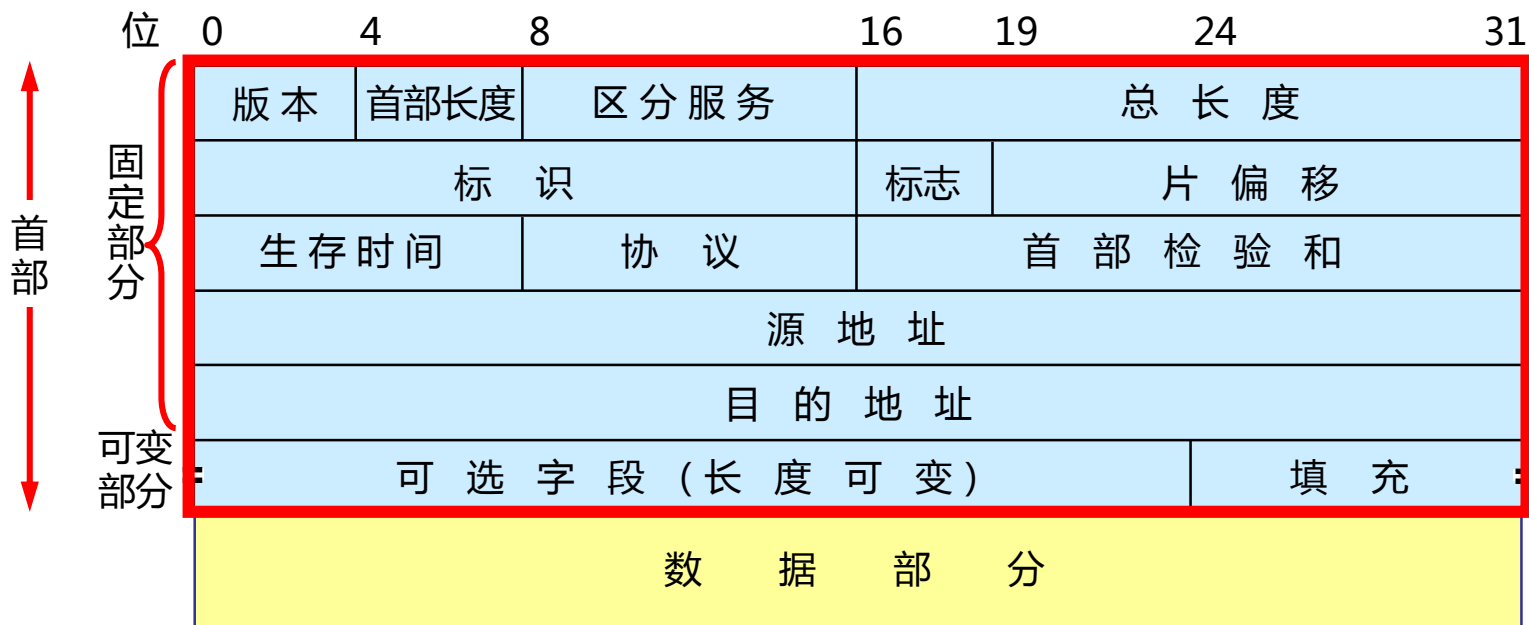




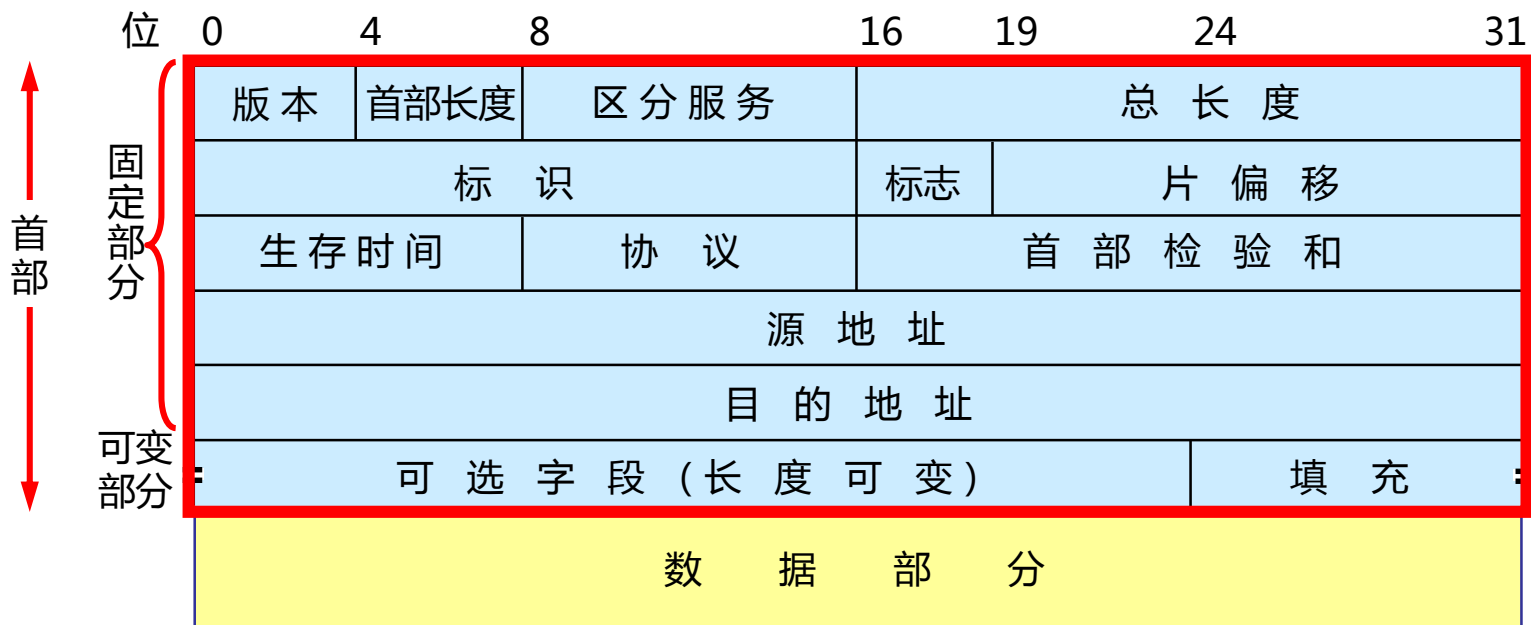
- **版本**：4位，指IP协议的版本，目前的IP协议版本号为4(即 IPv4)。
- **首部长度**：4位，可表示的最大数值是15个单位(一个单位为4字节)，因此IP的首部长度的最大值是60字节。
- **区分服务**：8位，用来获得更好的服务。在旧标准中叫做服务类型，但实际上一直未被使用过。1998年这个字段改名为区分服务。只有在使用区分服务 (DiffServ) 时，这个字段才起作用。在一般的情况下都不使用这个字段



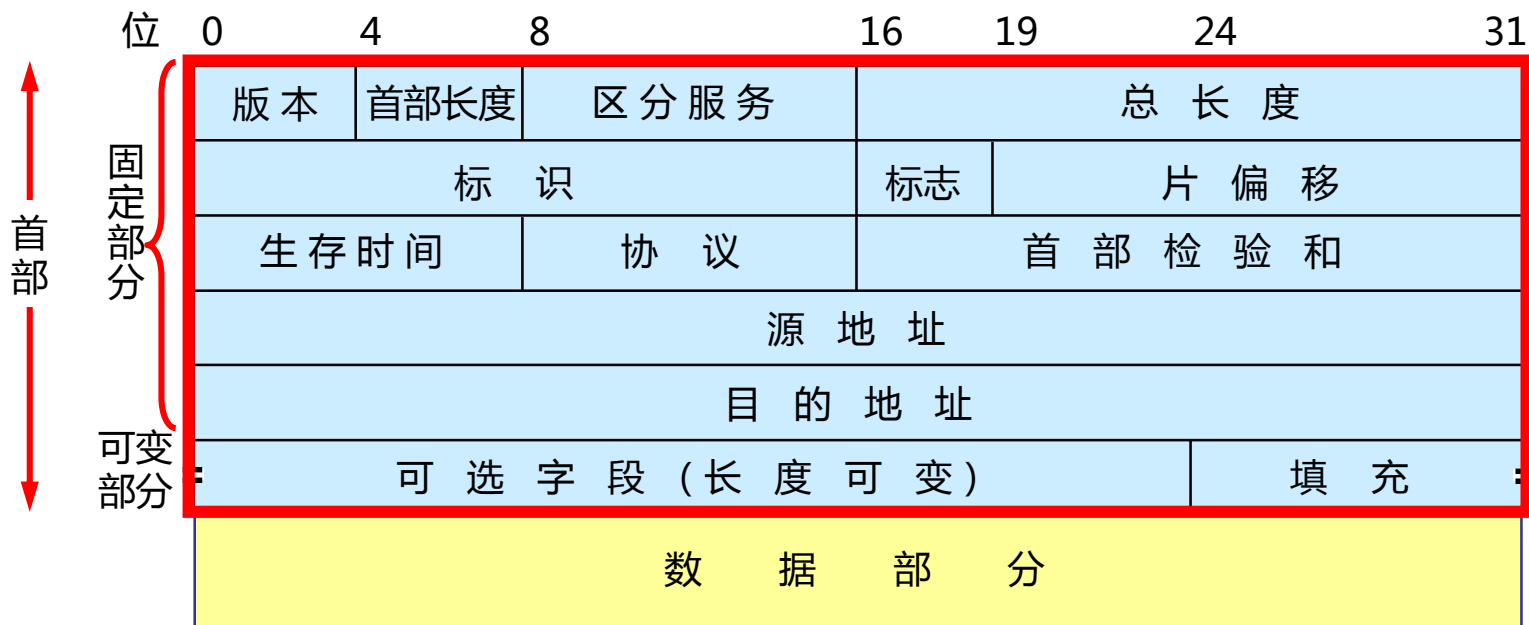
- **总长度**：16位，指首部和数据之和的长度，单位为字节，因此数据报的最大长度为65535字节。总长度必须不超过最大传送单元MTU，如果超过了，就需要把过长的数据报进行分片处理。
- **标识(identification)**：16位，它是一个计数器，用来产生数据报的标识。
- **标志(flag)**：3位，目前只有前两位有意义。
 - 标志字段的最低位是MF(More Fragment)，MF=1表示后面“还有分片”，MF=0表示最后一个分片。
 - 标志字段中间的一位是DF(Don't Fragment)，只有当DF=0时才允许分片。



- **片偏移**：13位。片偏移指出：较长的分组在分片后，某片在原分组中的相对位置。也就是说，相对于用户数据字段的起点，该片从何处开始。片偏移以8个字节为偏移单位，每个分片的长度一定是8个字节的整数倍。
- **生存时间**：8位，记为TTL(Time To Live)，数据报在网络中可通过的路由器数的最大值，TTL限制的为“跳数限制”，路由器转发数据报之前把TTL减1，如果TTL为0，路由器就丢弃这个数据报。因此TTL的单位是“跳数”。
- **协议**：8位，协议字段指出此数据报携带的数据使用何种协议，以方便目的主机的IP层将数据部分上交给哪个处理程序进行处理。



- **首部校验和**：16位。这个校验只检测数据报的首部，不对数据报的数据部分进行校验。首部校验和没有使用CRC这样复杂的计算，而是采用了更加简单的方法。（算法参考教材的内容，并进行介绍）
- **源地址**：32位，发送数据报的IP地址。
- **目的地址**：32位，接收数据报的IP地址。

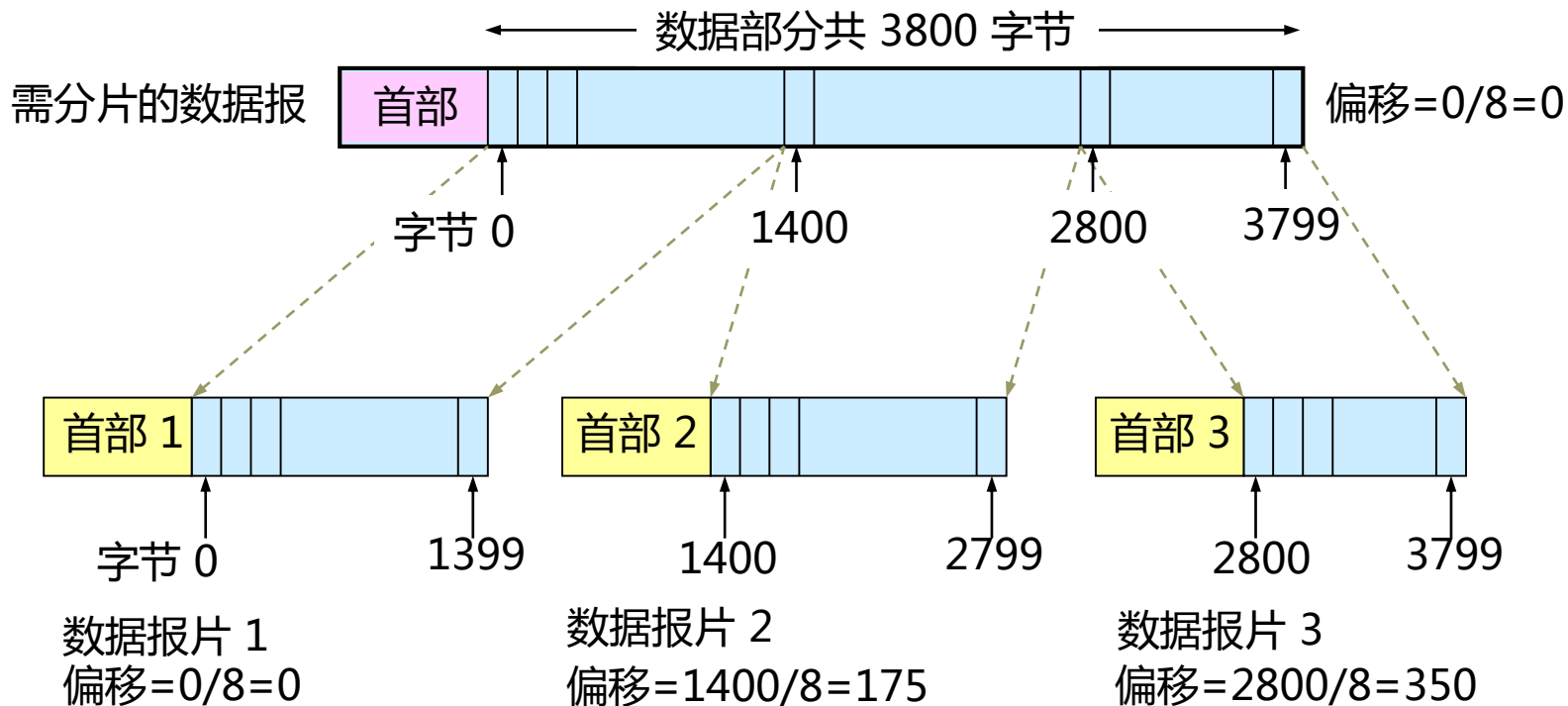


- IP首部的可变部分就是一个选项字段。
- 选项字段用来支持排错、测量以及安全等措施，内容很丰富。
- 此字段的长度为1-40个字节不等，取决于所选择的项目内容。
- 增加首部的可变部分是为了提高IP数据报的功能，但同时也使得IP数据报的首部变为可变的，增加了路由器的开销。
- IPv6将IP数据报的首部长度定为为固定的。

2.网际协议 (IP)

2.5 IP数据报的格式

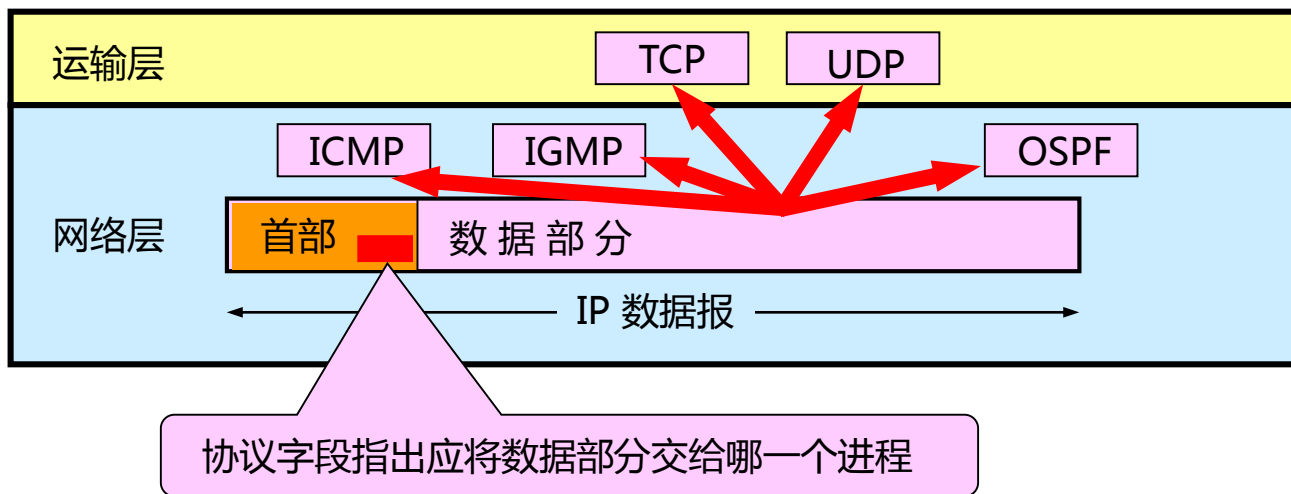
□ IP数据报分片的分析：



2.网际协议（IP）

2.5 IP数据报的格式

□ 协议字段的作用：



2.网际协议（IP）

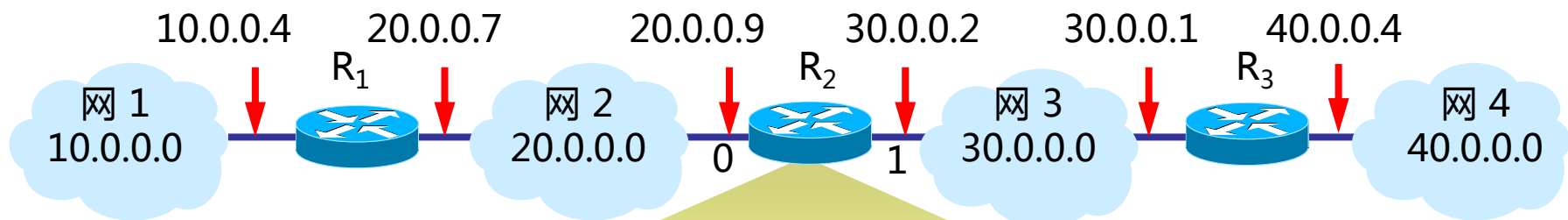
2.6 IP层转发分组的流程

□ 举个极端的例子：

- 有四个A类网络通过三个路由器连接在一起。每一个网络上都可能有成千上万个主机。
- 可以想像，若按目的主机号来制作路由表，则所得出的路由表就会过于庞大。
- 但若按主机所在的网络地址来制作路由表，那么每一个路由器中的路由表就只包含4个项目。这样就可使路由表大大简化。

2.网际协议 (IP)

2.6 IP层转发分组的流程

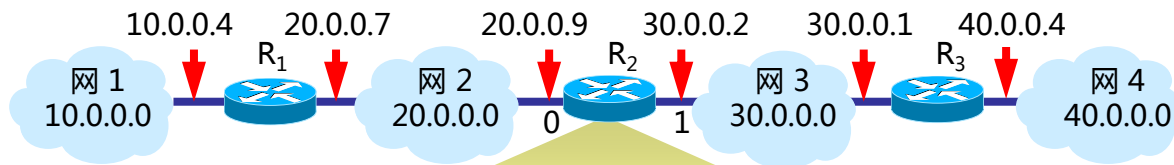


路由器 R₂ 的路由表

目的主机所在的网络	下一跳地址
20.0.0.0	直接交付, 接口 0
30.0.0.0	直接交付, 接口 1
10.0.0.0	20.0.0.7
40.0.0.0	30.0.0.1

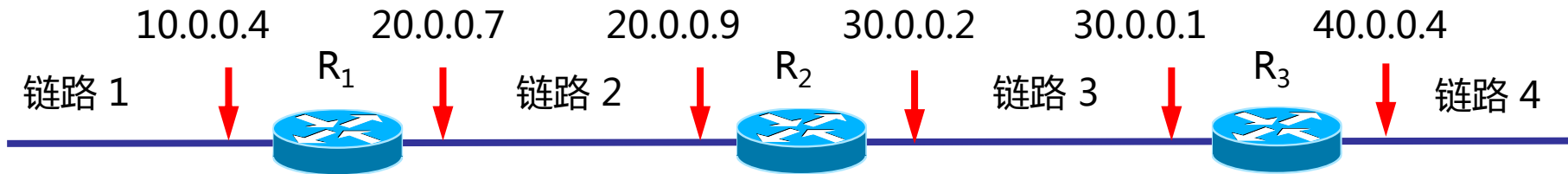
2.网际协议 (IP)

2.6 IP层转发分组的流程



路由器 R₂ 的路由表

目的主机所在的网络	下一跳地址
20.0.0.0	直接交付, 接口 0
30.0.0.0	直接交付, 接口 1
10.0.0.0	20.0.0.7
40.0.0.0	30.0.0.1



2.网际协议（IP）

2.6 IP层转发分组的流程

- 在互联网上转发分组时，是从一个路由器转发到下一个路由器。
- 在路由表中，对每一条路由最主要的是两个信息：
（目的网络地址，下一跳地址）
- 根据目的网络地址就能确定下一跳路由器，最终结果是：
 - IP数据报最终一定可以找到目的主机所在目的网络上的路由器（可能要通过多次的间接交付）。
 - 只有到达最后一个路由器时，才试图向目的主机进行直接交付。

2.网际协议（IP）

2.6 IP层转发分组的流程

□ 特定主机路由：

- 这种路由是为特定的目的主机指明一个路由。
- 采用特定主机路由可使网络管理人员能更方便地控制网络和测试网络，同时也可在需要考虑某种安全问题时采用这种特定主机路由。

□ 默认路由：

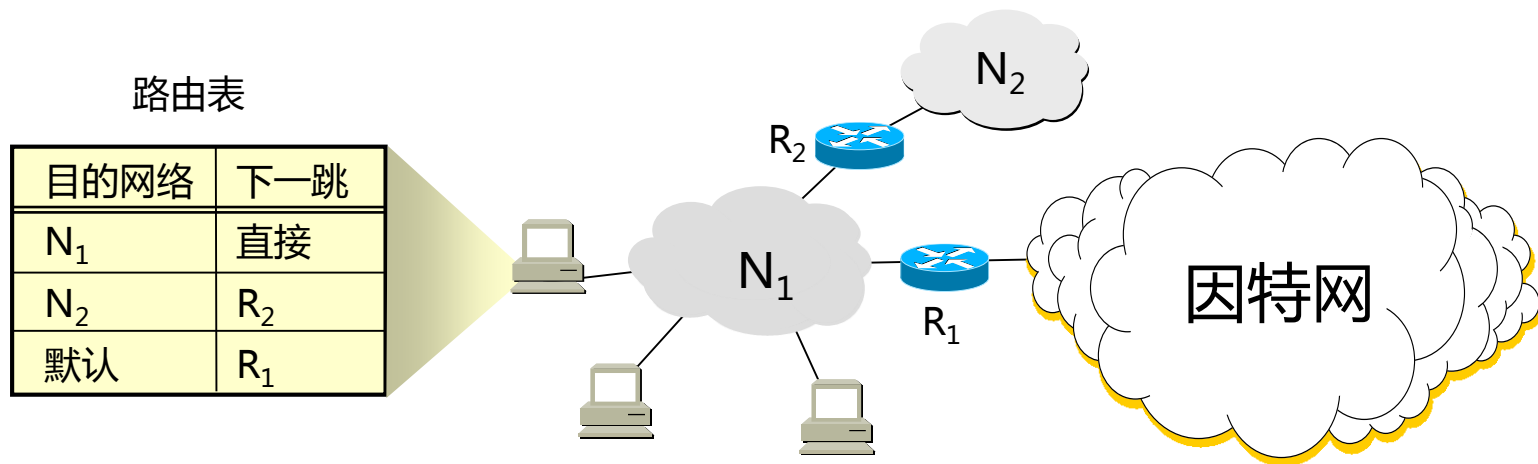
- 路由器还采用默认路由以减少路由表所占用的空间和搜索路由表所用的时间。这种转发方式在一个网络只有很少的对外连接时是很有用的。
- 默认路由在主机发送IP数据报时往往更能显示出它的好处。
- 如果一个主机连接在一个小网络上，而这个网络只用一个路由器和因特网连接，那么在这种情况下使用默认路由是非常合适的。

2.网际协议（IP）

2.6 IP层转发分组的流程

□ 默认路由：

- 只要目的网络不是N1和N2，就一律选择默认路由，把数据报先间接交付路由器R1，让R1再转发给下一个路由器。



2.网际协议（IP）

2.6 IP层转发分组的流程

□ 需要注意的是：

- IP数据报的首部中没有地方可以用来指明“下一跳路由器的IP地址”。
- 当路由器收到待转发的数据报，不是将下一跳路由器的IP地址填入IP数据报，而是送交下层的网络接口软件。
- 网络接口软件使用ARP负责将下一跳路由器的IP地址转换成硬件地址，并将此硬件地址放在链路层的MAC帧的首部，然后根据这个硬件地址找到下一跳路由器。

2.网际协议（IP）

2.6 IP层转发分组的流程

□ 分组转发算法

- ① 从数据报的首部提取目的主机的IP地址D, 得出目的网络地址为N。
- ② 若网络N与此路由器直接相连, 则把数据报直接交付目的主机D; 否则是间接交付, 执行③。
- ③ 若路由表中有目的地址为D的特定主机路由, 则把数据报传送给路由表中所指明的下一跳路由器; 否则, 执行④。
- ④ 若路由表中有到达网络N的路由, 则把数据报传送给路由表指明的下一跳路由器; 否则, 执行⑤。
- ⑤ 若路由表中有一个默认路由, 则把数据报传送给路由表中所指明的默认路由器; 否则, 执行⑥。
- ⑥ 报告转发分组出错。

2. 网

组的流程

The screenshot displays the OSPF configuration for a Cisco router (10.0.2.253) in the Observium interface. The main table shows the router's status and OSPF areas. Below each area, a sub-table lists the configured ports, their types, and states.

Router Id	Status	ABR	ASBR	Areas	Ports	Neighbours
10.0.2.253	enabled	yes	yes	3	14(14)	13

Area Id	Status	Ports
0.0.0.0	enabled	5(5)

Port	Status	Port Type	Port State
enabled	enabled	broadcast	designatedRouter
enabled	enabled	broadcast	designatedRouter
enabled	enabled	broadcast	designatedRouter
enabled	enabled	pointToPoint	loopback
enabled	enabled	broadcast	designatedRouter

Area Id	Status	Ports
0.0.0.1	enabled	1(1)

Port	Status	Port Type	Port State
enabled	enabled	broadcast	designatedRouter

Area Id	Status	Ports
0.0.0.50	enabled	8(8)

Port	Status	Port Type	Port State
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint
enabled	enabled	pointToPoint	pointToPoint

Router Id	Device	IP Address	Status
10.0.2.252	unknown	10.0.1.1	full
10.0.2.250	unknown	10.0.1.17	full

3.划分子网与构建超网

3.1划分子网：三级IP地址

□ 举例讨论：

- 某单位获得A类地址，却只有300台主机。
- 路由器里面的路由表应该有多少条记录？
- 某单位紧急扩展网络后，如何让新扩展网络接入互联网？

3.划分子网与构建超网

3.1划分子网：三级IP地址

- 两级IP地址的不足：
 - IP 地址空间的利用率有时很低。
 - 给每一个物理网络分配一个网络号会使路由表变得太大因而使网络性能变坏。
 - 两级的 IP 地址不够灵活。
- 解决思路：
 - 申请到地址后，可以再进行二次划分，把一个A类或B类地址划分为多个网络。
 - 让地址段更小、更灵活。

3.划分子网与构建超网

3.1划分子网：三级IP地址

- 从1985年起，在IP地址中又增加了一个新的字段：“子网号字段”，两级的IP地址变成为三级的IP地址。
- 从两级的IP地址变为三级的IP地址后，解决了分类IP地址管理的不足，让IP管理和使用变得更加灵活。
- 将两级的IP地址变为三级的IP地址的做法，叫作划分子网(subnetting)，或叫子网划分、子网路由选择。
- 目前，划分子网已成为因特网的正式标准协议。

3.划分子网与构建超网

3.1划分子网：三级IP地址

□ 划分子网的思路：

- 划分子网纯属一个单位内部的事情。单位对外仍然表现为没有划分子网的网络。
- 从主机号借用若干个位作为子网号 subnet-id，而主机号 host-id 也就相应减少了若干个位。

IP地址 ::= {<网络号>, <子网号>, <主机号>}

3.划分子网与构建超网

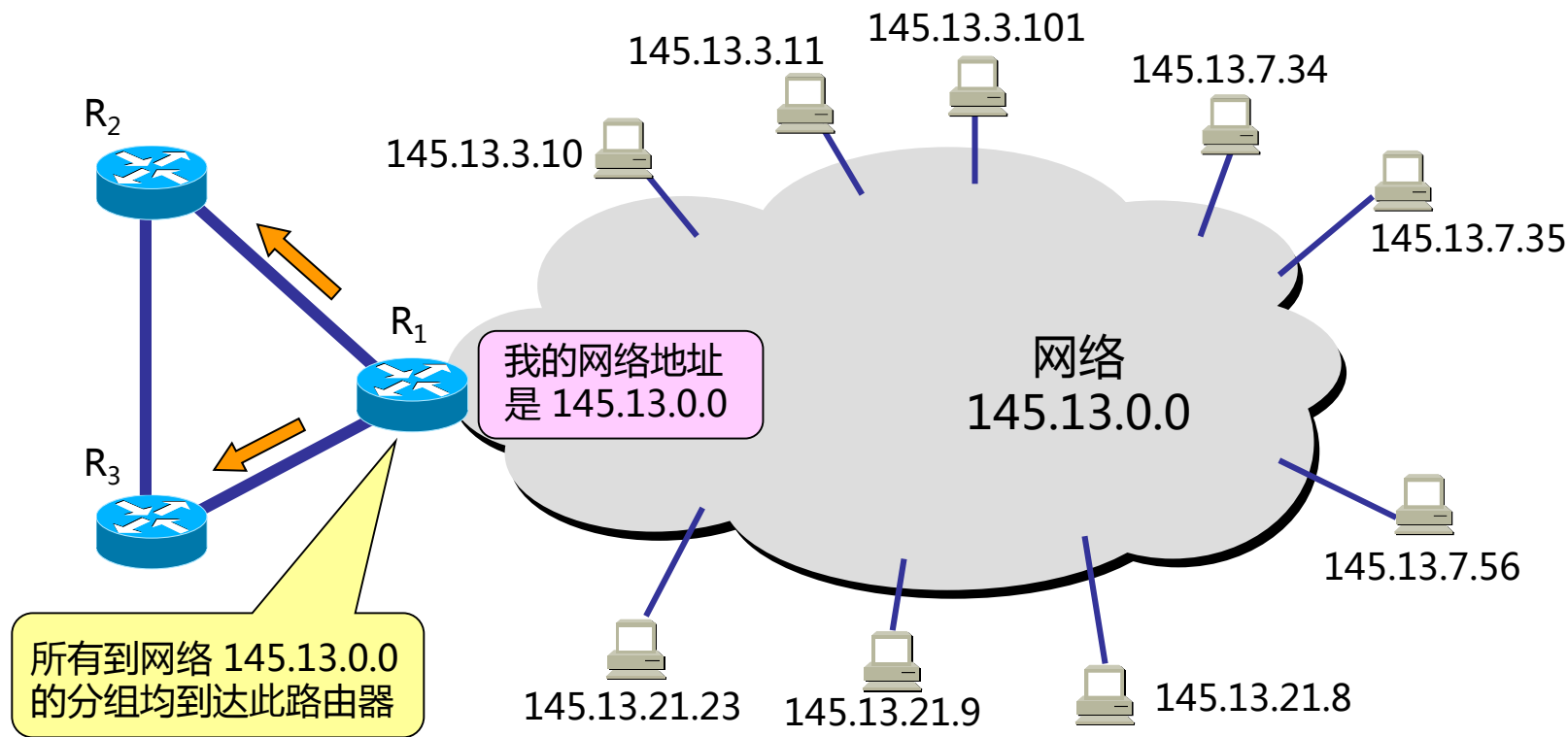
3.1划分子网：三级IP地址

□ 划分子网的思路：

- 凡是从其他网络发送给本单位某个主机的IP数据报，仍然是根据IP数据报的目的网络号net-id，先找到连接在本单位网络上的路由器。
- 然后此路由器在收到IP数据报后，再按目的网络号net-id和子网号subnet-id找到目的子网。
- 最后就将IP数据报直接交付目的主机。

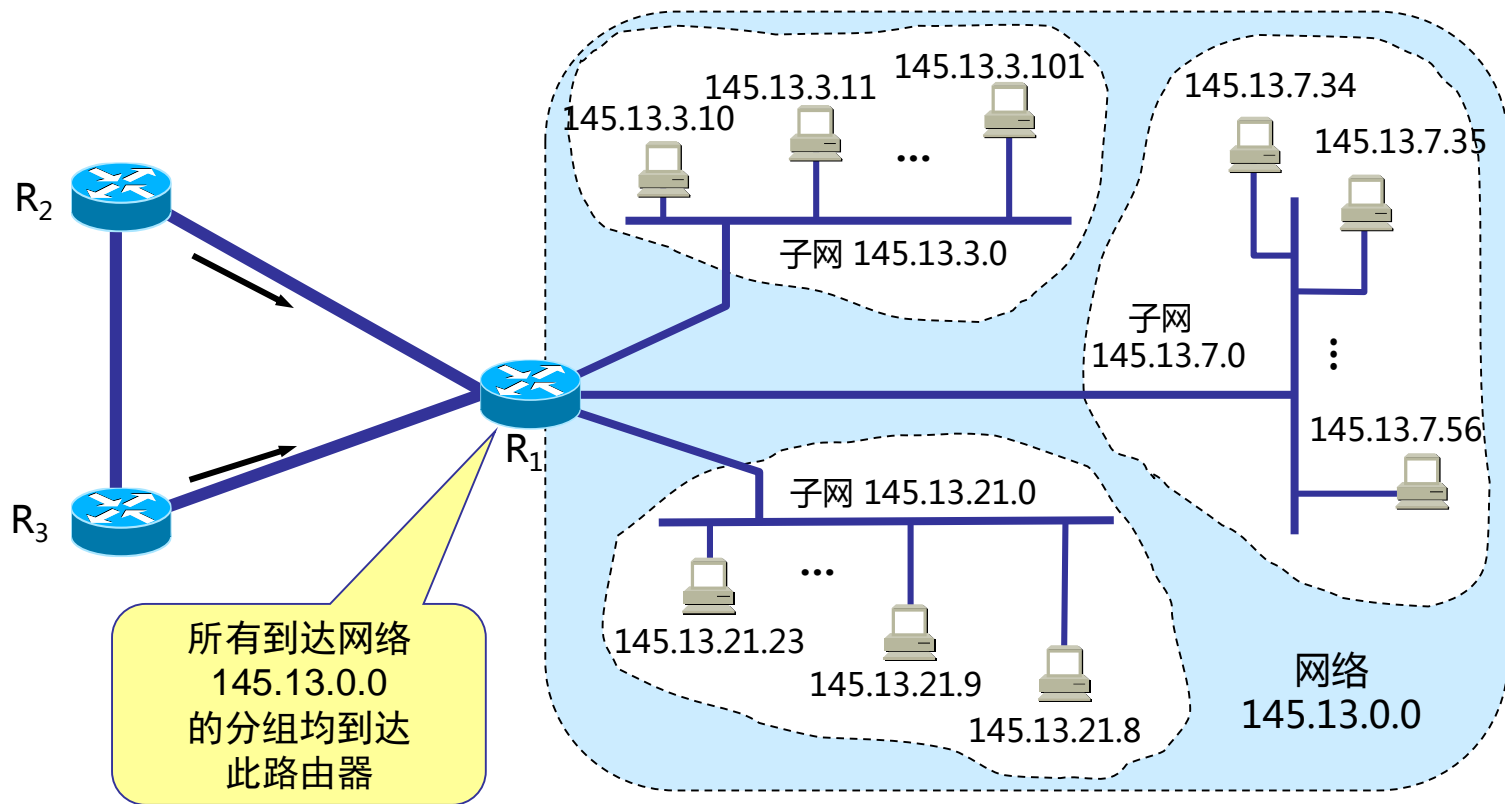
3.划分子网与构建超网

3.1划分子网：三级IP地址



3.划分子网与构建超网

3.1划分子网：三级IP地址



3.划分子网与构建超网

3.1划分子网：三级IP地址

- 当没有划分子网时，IP地址是两级结构。
- 划分子网后IP地址就变成了三级结构。
- 划分子网只是把IP地址的主机号host-id这部分进行再划分，而不改变IP地址原来的网络号 net-id。
- 划分子网是单位内部为了管理而自行开展的工作，对于上层接入网络而讲没有任何影响。

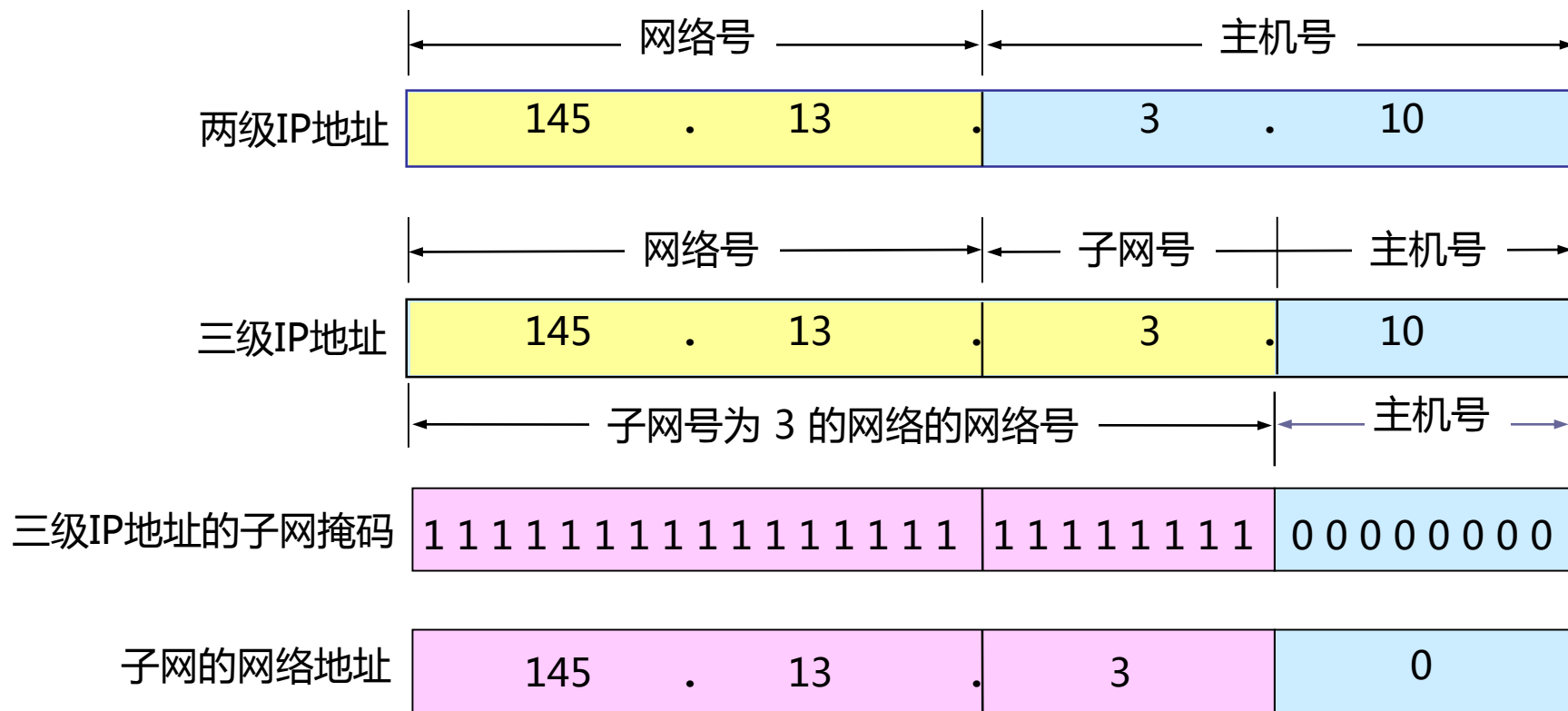
3.划分子网与构建超网

3.1划分子网：子网掩码

- 从IP数据报的首部无法看出源主机或目的主机所连接的网络是否进行了子网划分。
- 为了看到一个IP地址的网络号和主机号，于是使用了子网掩码(subnet mask)来作为辅助手段，以计算某一个IP地址的网络号。

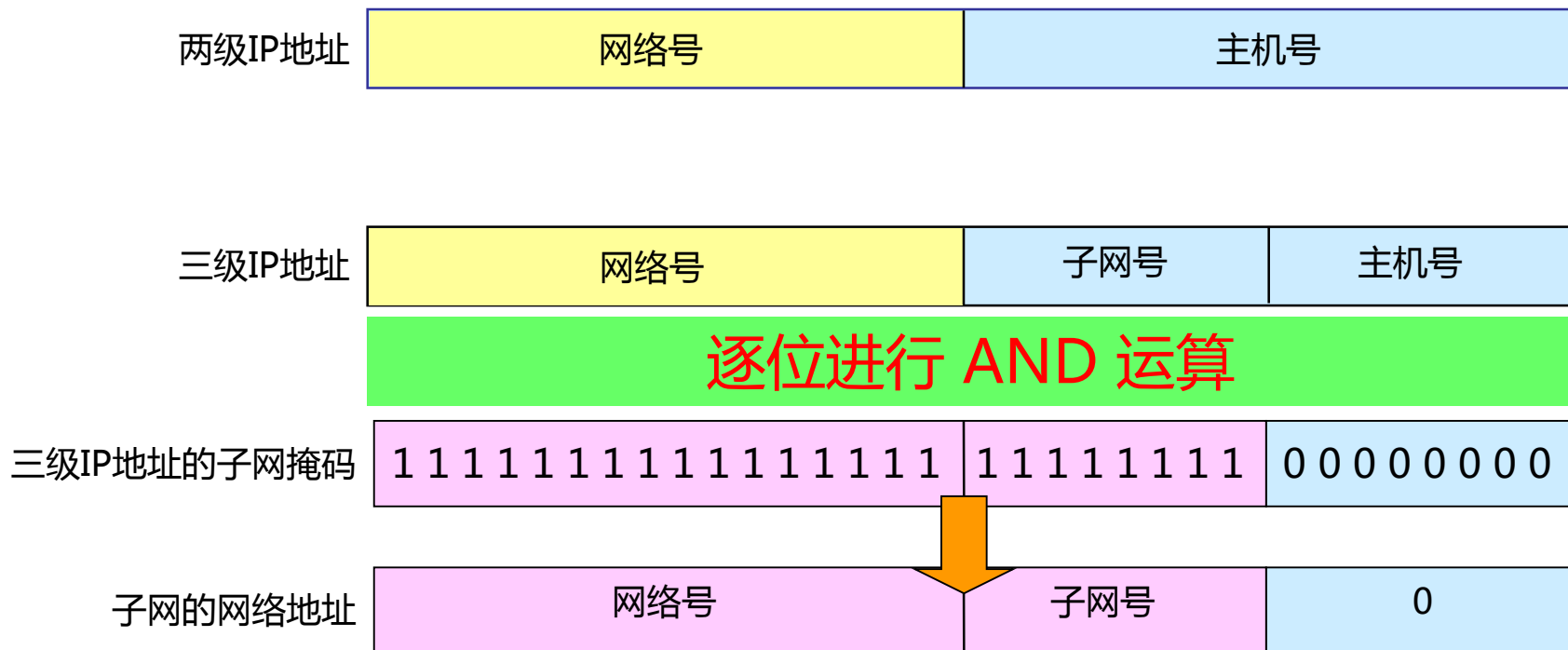
3.划分子网与构建超网

3.1划分子网：子网掩码



3.划分子网与构建超网

3.1划分子网：子网掩码



3.划分子网与构建超网

3.1划分子网：子网掩码

- 使用子网掩码之后，不管网络是否划分了子网，都可以通过子网掩码和IP地址逐位的“与”运算，方便的得出网络地址。
- 路由器在进行分组转发时，就不需要考虑是否划分子网，通过一个方法进行处理即可。

3.划分子网与构建超网

3.1划分子网：子网掩码

□ 讨论：

- IP地址：211.69.32.18，是C类地址，在具体使用时没有进行子网划分。
 - 但是在具体应用中，依然需要使用子网掩码，这是为什么？
- 不划分子网也要使用子网掩码的原因，就在于让路由器以一个方法进行分组处理，以便于以一种方法查找路由表。
- 子网掩码是一个网络或一个子网的重要属性。

3.划分子网与构建超网

3.1划分子网：子网掩码

A 类 地 址	网络地址	网络号	主机号为全 0
	默认子网掩码 255.0.0.0	1 1 1 1 1 1 1 1	0 0
B 类 地 址	网络地址	网络号	主机号为全 0
	默认子网掩码 255.255.0.0	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
C 类 地 址	网络地址	网络号	主机号为全 0
	默认子网掩码 255.255.255.0	1 1	0 0 0 0 0 0 0 0

3.划分子网与构建超网

3.1划分子网：子网掩码

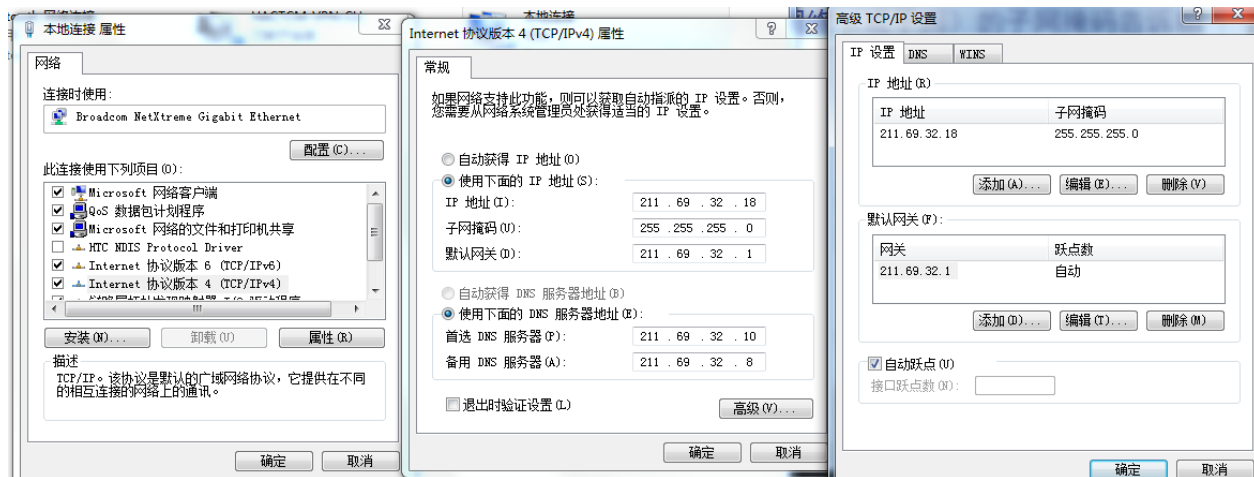
- 子网掩码是一个网络或一个子网的重要属性。
- 路由器在和相邻路由器交换路由信息时，必须把自己所在网络（或子网）的子网掩码告诉相邻路由器。
- 路由器的路由表中的每一个项目，除了要给出目的网络地址外，还必须同时给出该网络的子网掩码。
- 若一个路由器连接在两个子网上就拥有两个网络地址和两个子网掩码。

3.划分子网与构建超网

3.1划分子网：子网掩码

现场演示：

- 让一个主机同时接入到两个网络上。（为一个网络接口卡配置多个IP地址）



3.划分子网与构建超网

3.1划分子网：子网掩码

□ 计算：

- 已知IP地址211.69.32.18，子网掩码是255.255.255.0，请计算其网络地址。
- 已知IP地址211.69.32.18，子网掩码是255.255.240.0，请计算其网络地址。
- 已知IP地址211.69.32.18，子网掩码是255.255.255.252，请计算机其网络地址，并列出现该网络内的所有IP地址。

3.划分子网与构建超网

3.2使用子网掩码的分组转发过程

- 在不划分子网的两级 IP 地址下，从IP地址得出网络地址是个很简单的事。（讨论：两级IP地址下的分组转发）
- 但在划分子网的情况下，从 IP 地址却不能唯一地得出网络地址来，这是因为网络地址取决于网络所采用的子网掩码，但数据报的首部并没有提供子网掩码的信息。
- 在划分子网的情况下，分组转发的算法也必须做相应的改动。

3.划分子网与构建超网

3.2使用子网掩码的分组转发过程

- 使用子网划分后，路由表必须包含以下三项内容：**目的网络地址、子网掩码和下一跳地址。**

192.168.183.0/24	OSPF	10	82	D	10.0.1.30	Vlanif1129
192.168.254.0/24	OSPF	10	42	D	10.0.1.30	Vlanif1129
192.168.255.0/24	OSPF	10	42	D	10.0.1.30	Vlanif1129
202.4.128.0/19	Static	60	0	RD	222.21.219.73	Vlanif1222
202.38.64.0/18	Static	60	0	RD	222.21.219.73	Vlanif1222
202.38.140.0/23	Static	60	0	RD	222.21.219.73	Vlanif1222
202.38.184.0/21	Static	60	0	RD	222.21.219.73	Vlanif1222
202.38.192.0/18	Static	60	0	RD	222.21.219.73	Vlanif1222
202.112.0.0/13	Static	60	0	RD	222.21.219.73	Vlanif1222
202.120.0.0/15	Static	60	0	RD	222.21.219.73	Vlanif1222
202.127.216.0/21	Static	60	0	RD	222.21.219.73	Vlanif1222
202.127.224.0/19	Static	60	0	RD	222.21.219.73	Vlanif1222
202.179.240.0/20	Static	60	0	RD	222.21.219.73	Vlanif1222

3.划分子网与构建超网

3.2使用子网掩码的分组转发过程

□ 划分子网的情况下，路由器转发分组的算法：

①从收到的分组的首部提取目的IP地址D。

②先用各网络的子网掩码和D逐位相“与”，看是否和相应的网络地址匹配。若匹配则直接交付。否则就是间接交付，执行③。

③若路由表中有目的地址为D的特定主机路由，则将分组传送给指明的下一跳路由器；否则，执行④。

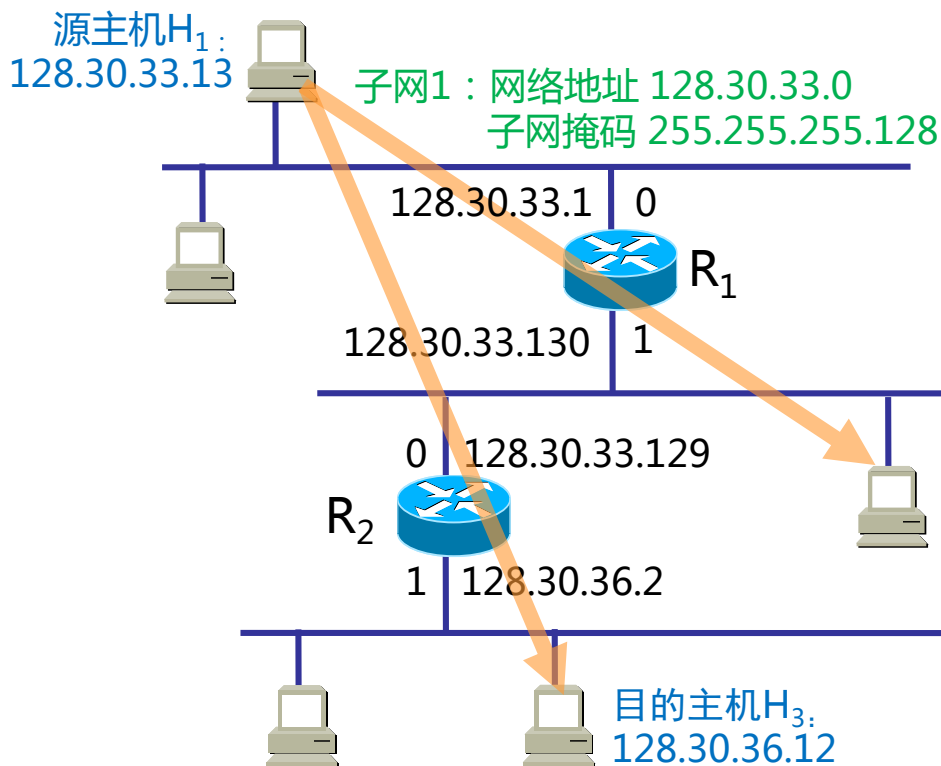
④对路由表中的每一行的子网掩码和D逐位相“与”，若其结果与该行的目的网络地址匹配，则将分组传送给该行指明的下一跳路由器；否则，执行⑤。

⑤若路由表中有一个默认路由，则将分组传送给路由表中所指明的默认路由器；否则，执行⑥。

⑥报告转发分组出错。

3.划分子网与构建超网

3.2使用子网掩码的分组转发过程



R₁的路由表（未给出默认路由器）

目的网络地址	子网掩码	下一跳
128.30.33.0	255.255.255.128	接口 0
128.30.33.128	255.255.255.128	接口 1
128.30.36.0	255.255.255.0	R ₂

3.划分子网与构建超网

3.2使用子网掩码的分组转发过程

□ 讨论：

- H1向H2发送分组数据，R1在接收到H1发送的数据报后，查找路由表的过程。
- H1向H3发送分组数据，R1、R2在接收到H1发送的数据报后，查找路由表的过程。
- R2的路由表的具体内容是什么？

3.划分子网与构建超网

3.3构建超网 (CIDR)

- 划分子网在一定程度上缓解了因特网在发展中遇到的困难。然而在 1992 年因特网仍然面临三个必须尽早解决的问题，这就是：
 - B 类地址在 1992 年已分配了近一半，眼看就要在 1994 年 3 月全部分配完毕！
 - 因特网主干网上的路由表中的项目数急剧增长（从几千个增长到几万个）。
 - 整个 IPv4 的地址空间最终将全部耗尽。

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ 网络前缀

3.划分子网与构建超网

3.3构建超网 (CIDR)

- 1987年，RFC 1009就指明了在一个划分子网的网络中可同时使用几个不同的子网掩码。
- 使用变长子网掩码**VLSM**(Variable Length Subnet Mask)可进一步提高 IP 地址资源的利用率。
- 在VLSM的基础上又进一步研究出无分类编址方法，它的正式名字是**无分类域间路由选择CIDR(Classless Inter-Domain Routing)**。

3.划分子网与构建超网

3.3构建超网 (CIDR)

- ❑ CIDR 消除了传统的A类、B类和C类地址以及划分子网的概念，因而可以更加有效地分配IPv4的地址空间。
- ❑ CIDR使用各种长度的“网络前缀” (network-prefix)来代替分类地址中的网络号和子网号。
- ❑ IP地址从三级编址（使用子网掩码）又回到了两级编址。

IP地址 ::= { <网络前缀>, <主机号> }

3.划分子网与构建超网

3.3构建超网 (CIDR)

- CIDR使用“斜线记法” (slash notation) , 它又称为CIDR记法 , 即在IP地址面加上一个斜线 “/” , 然后写上网络前缀所占的位数 (这个数值对应于三级编址中子网掩码中 1 的个数) 。
- CIDR 把网络前缀都相同的连续的IP地址组成 “CIDR 地址块” 。

IP地址 ::= { <网络前缀> , <主机号> }

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ CIDR地址表示方法讨论：

- **211.69.32.0/20**表示的地址块共有 2^{12} 个地址（因为斜线后面的20是网络前缀的位数，所以这个地址的主机号是12位）。
- 这个地址块的起始地址是 211.69.32.0。
- 在不需要指出地址块的起始地址时，也可将这样的地址块简称为“**/20 地址块**”。
- 211.69.32.0/20地址块的最小地址：211.69.32.0
- 211.69.32.0/20地址块的最大地址：211.69.47.255

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ CIDR地址表示方法讨论：

- 10.0.0.0/10表示的地址块共有 2^{22} 个地址。
- 10.0.0.0/10可简写为10/10，也就是将点分十进制中低位连续的0省略。
- 10.0.0.0/10相当于指出IP地址10.0.0.0的掩码是255.192.0.0，即
11111111 11000000 00000000 00000000
- 10.0.0.0/10可以表示为在网络前缀的后面加一个星号 * 的方法，即：00001010 00*，在星号*之前是网络前缀，而星号*表示IP地址中的主机号，可以是任意值。

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ CIDR地址表示方法讨论：

Network : 211.69.32.0
Netmask : 255.255.255.0  CIDR : 211.69.32.0/24

Network : 211.69.32.0
Netmask : 255.255.240.0  CIDR : 211.69.32.0/?

3.划分子网与构建超网

□ CIDR地址表示方法讨论：

(CIDR)

子网掩码	CIDR 名称	IP 数量	可用 IP 数量	C 类网络的数量
255.255.255.255	/32	1	1	1/256
255.255.255.254	/31	2	0	1/128
255.255.255.252	/30	4	2	1/64
255.255.255.248	/29	8	6	1/32
255.255.255.240	/28	16	14	1/16
255.255.255.224	/27	32	30	1/8
255.255.255.192	/26	64	62	1/4
255.255.255.128	/25	128	126	一半
255.255.255.0	/24	256	254	1
255.255.254.0	/23	512	510	2
255.255.252.0	/22	1024	1022	4
255.255.248.0	/21	2048	2046	8
255.255.240.0	/20	4096	4094	16
255.255.224.0	/19	8192	8190	32
255.255.192.0	/18	16,384	16,382	64
255.255.128.0	/17	32,768	32,766	128
255.255.0.0	/16	65,536	65,534	256
255.254.0.0	/15	131,072	131,070	512
255.252.0.0	/14	262,144	262,142	1024
255.248.0.0	/13	524,288	524,286	2048
255.240.0.0	/12	1,048,576	1,048,574	4096
255.224.0.0	/11	2,097,152	2,097,150	8192
255.192.0.0	/10	4,194,304	4,194,302	16,384
255.128.0.0	/9	8,388,608	8,388,606	32,768
255.0.0.0	/8	16,777,216	16,777,214	65,536
254.0.0.0	/7	33,554,432	33,554,430	131,072
252.0.0.0	/6	67,108,864	67,108,862	262,144
248.0.0.0	/5	134,217,728	134,217,726	1,048,576
240.0.0.0	/4	268,435,456	268,435,454	2,097,152
224.0.0.0	/3	536,870,912	536,870,910	4,194,304
192.0.0.0	/2	1,073,741,824	1,073,741,822	8,388,608
128.0.0.0	/1	2,147,483,648	2,147,483,646	16,777,216
0.0.0.0	/0	4,294,967,296	4,294,967,294	33,554,432

3. 划分

□ CIDR地

子网掩码	CIDR 名称	IP 数量	可用 IP 数量	C 类网络的数量
255.255.255.255	/32	1	1	1/256
255.255.255.254	/31	2	0	1/128
255.255.255.252	/30	4	2	1/64
255.255.255.248	/29	8	6	1/32
255.255.255.240	/28	16	14	1/16
255.255.255.224	/27	32	30	1/8
255.255.255.192	/26	64	62	1/4
255.255.255.128	/25	128	126	一半
255.255.255.0	/24	256	254	1
255.255.254.0	/23	512	510	2
255.255.252.0	/22	1024	1022	4
255.255.248.0	/21	2048	2046	8
255.255.240.0	/20	4096	4094	16
255.255.224.0	/19	8192	8190	32
255.255.192.0	/18	16,384	16,382	64
255.255.128.0	/17	32,768	32,766	128
255.255.0.0	/16	65,536	65,534	256
255.254.0.0	/15	131,072	131,070	512
255.252.0.0	/14	262,144	262,142	1024
255.248.0.0	/13	524,288	524,286	2048
255.240.0.0	/12	1,048,576	1,048,574	4096

超网 (CIDR)

3.划分子网与构建超网

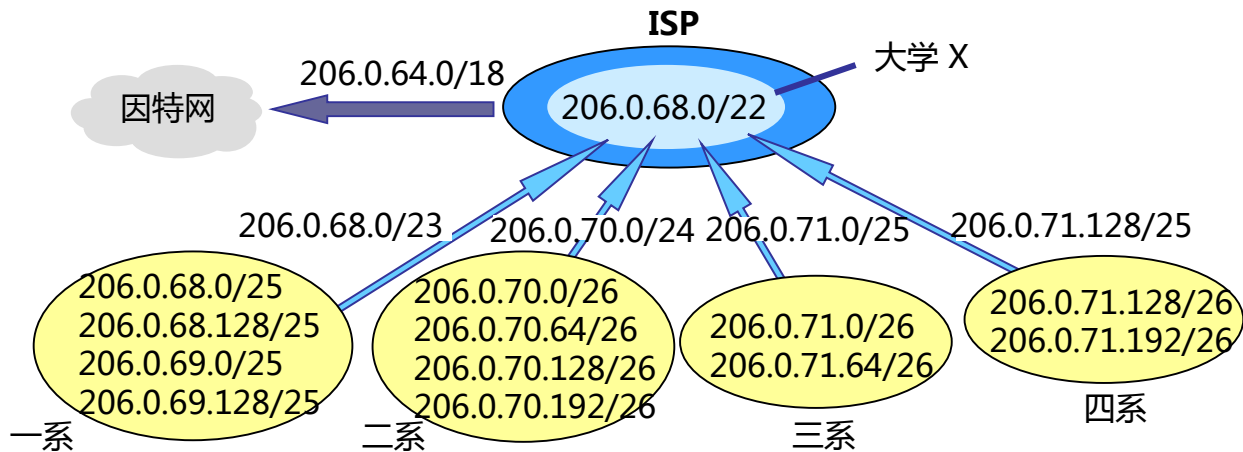
3.3构建超网 (CIDR)

- 路由聚合 (route aggregation) :
 - 一个CIDR地址块可以表示很多地址，这种地址的聚合常称为路由聚合，它使得路由表中的一个项目可以表示很多个（例如上千个）原来传统分类地址的路由。
 - 路由聚合也称为**构成超网**(supernetting)。
 - CIDR虽然不使用子网了，但仍然使用“掩码”这一名词（但不叫子网掩码）。
 - 对于**/20地址块**，它的掩码是20个连续的1。斜线记法中的数字就是掩码中1的个数。

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ 路由聚合 (route aggregation) :



单位	地址块	二进制表示	地址数
ISP	206.0.64.0/18	11001110.00000000.01*	16384
大学	206.0.68.0/22	11001110.00000000.010001*	1024
一系	206.0.68.0/23	11001110.00000000.0100010*	512
二系	206.0.70.0/24	11001110.00000000.01000110.*	256
三系	206.0.71.0/25	11001110.00000000.01000111.0*	128
四系	206.0.71.128/25	11001110.00000000.01000111.1*	128

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ 最长前缀匹配：

- 使用CIDR时，路由表中的每个项目由“网络前缀”和“下一跳地址”组成。在查找路由表时可能会得到不止一个匹配结果。
- 应当从匹配结果中选择具有最长网络前缀的路由：**最长前缀匹配 (longest-prefix matching)**。
- 网络前缀越长，其地址块就越小，因而路由就越具体 (more specific) 。
- 最长前缀匹配又称为**最长匹配**或**最佳匹配**。

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ 使用二叉线索查找路由表：

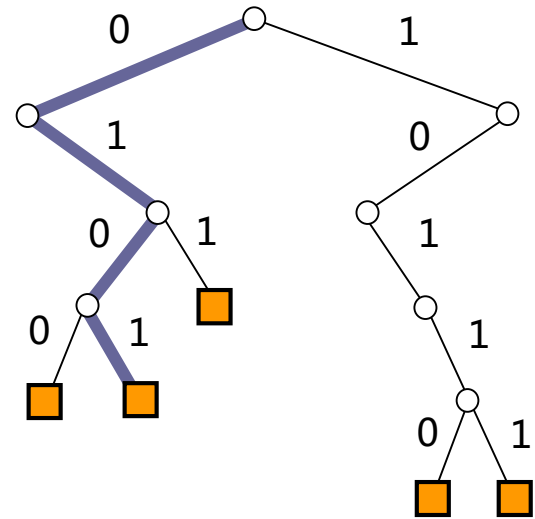
- 当路由表的项目数很大时，设法减小路由表的查找时间就成为一个非常重要的问题。为了进行更加有效的查找，通常是将无分类编址的路由表存放在一种层次的数据结构中，然后自上而下地按层次进行查找。最常用的就是**二叉线索(binary trie)**。
- IP地址中从左到右的比特值决定了从根结点逐层向下层延伸的路径，而二叉线索中的各个路径就代表路由表中存放的各个地址。
- 为了提高二叉线索的查找速度，广泛使用了各种压缩技术。

3.划分子网与构建超网

3.3构建超网 (CIDR)

□ 使用二叉线索查找路由表：

32 位的 IP 地址	唯一前缀
01000110 00000000 00000000 00000000	0100
01010110 00000000 00000000 00000000	0101
01100001 00000000 00000000 00000000	011
10110000 00000010 00000000 00000000	10110
10111011 00001010 00000000 00000000	10111



4.网际控制报文协议ICMP

4.1 ICMP

- ❑ 为了提高IP数据报交付成功的机会，在网际层使用了网际控制报文协议 ICMP (Internet Control Message Protocol)。
- ❑ ICMP是IPv4协议簇中的一个子协议，用于在IP主机、路由器之间传递控制消息。
- ❑ 控制消息是指网络通不通、主机是否可达、路由是否可用等网络本身的消息。这些控制消息虽然并不传送用户数据，但是对于用户数据的传递起着重要的作用。

4. 网际控制报文协议ICMP

4.1 ICMP

- ❑ ICMP允许主机或路由器报告差错情况和提供有关异常情况的报告。
- ❑ ICMP不是高层协议，而是IP层的协议。
- ❑ ICMP协议与ARP协议不同，ICMP依靠IP协议来完成任务，所以ICMP报文中要封装IP头部，组成IP数据报发送出去。
- ❑ ICMP一般并不用来在端系统之间传送数据，不被用户网络程序直接使用。端系统中，Ping和Traceroute等诊断网络的工具才会直接使用ICMP协议。

4.网际控制报文协议ICMP

4.1 ICMP

- ICMP报告无法传送的数据报的错误，并帮助对这些错误进行疑难解答。例如，如果IPv4不能够将数据报传送到目标主机，则路由器或目标主机上的ICMP就会向主机发送ICMP的“无法到达目标”的消息。
- 在下列情况中，通常自动发送ICMP消息：
 - IP数据报无法访问目标。
 - IP路由器（网关）无法按当前的传输速率转发数据报。
 - IP路由器将发送主机重定向为使用到达目标的更佳路由。

4.网际控制报文协议ICMP

4.1 ICMP

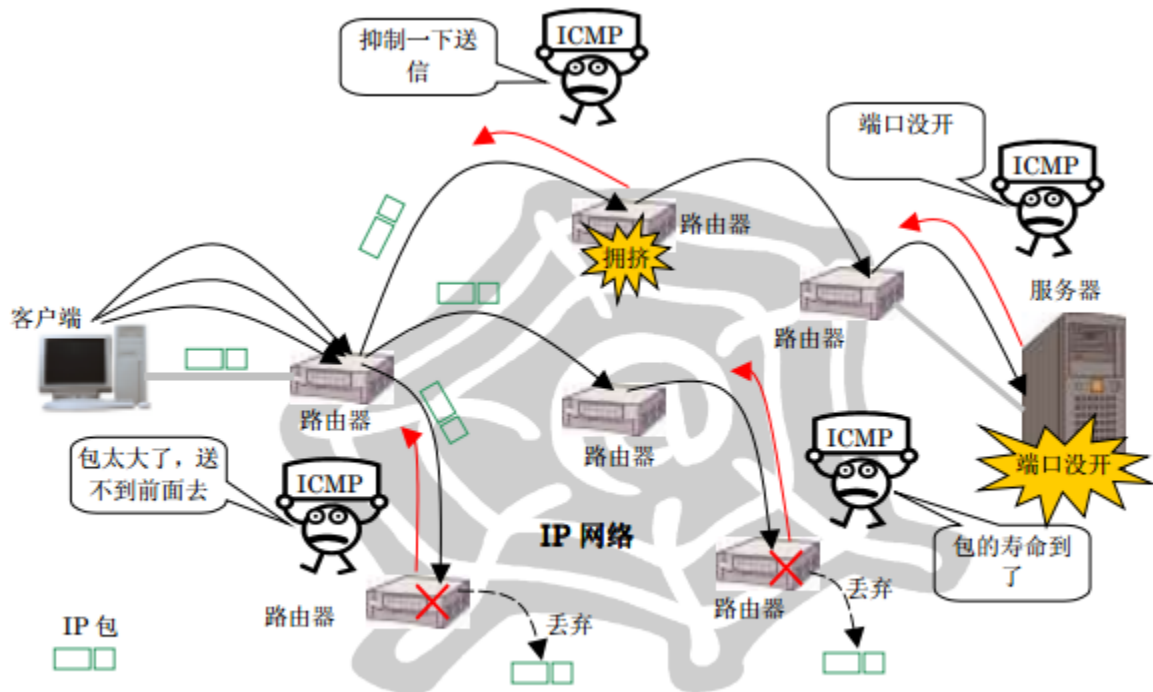
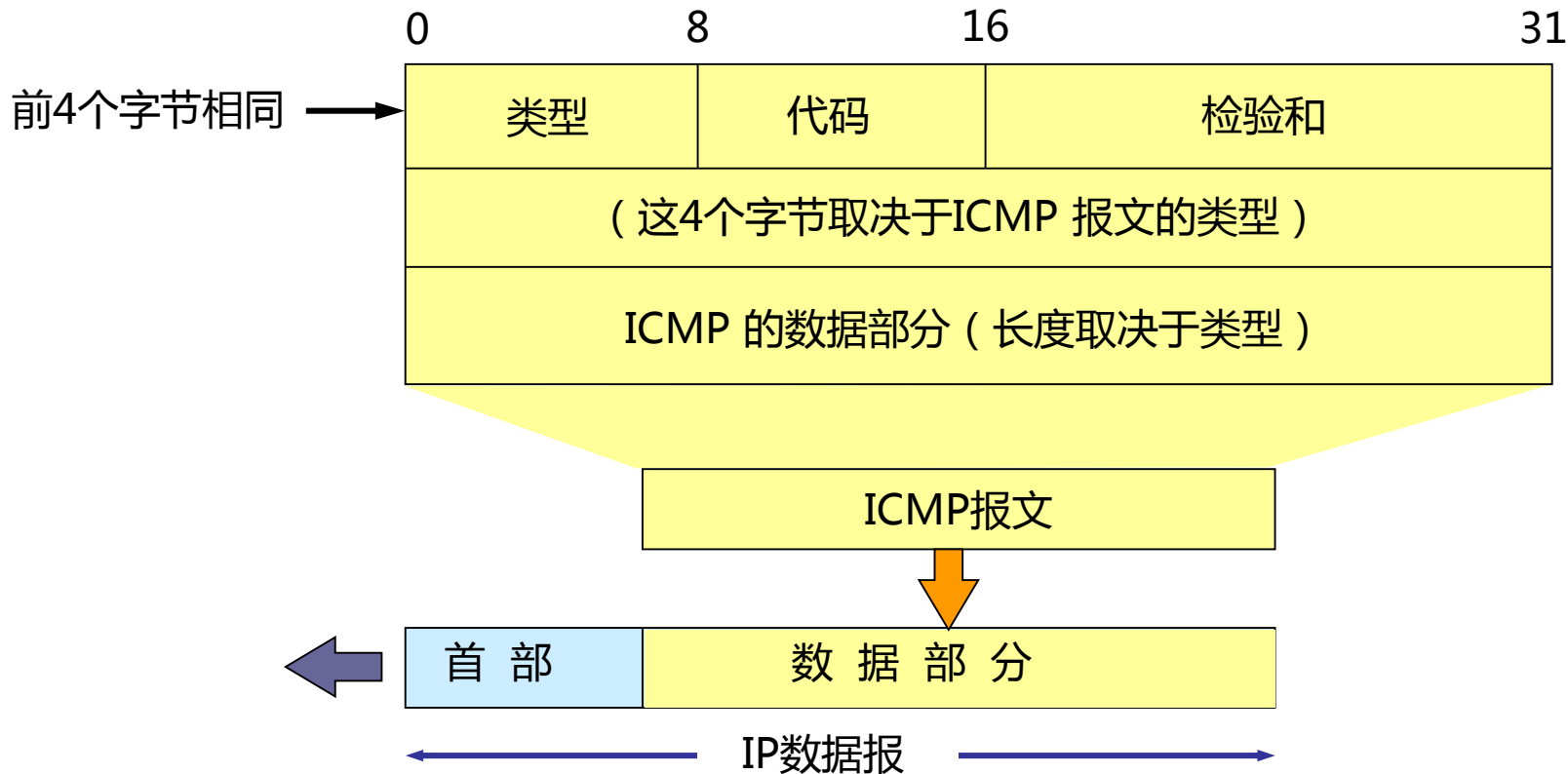


图 1 ICMP 是使 IP 通信平稳运行的辅助协议

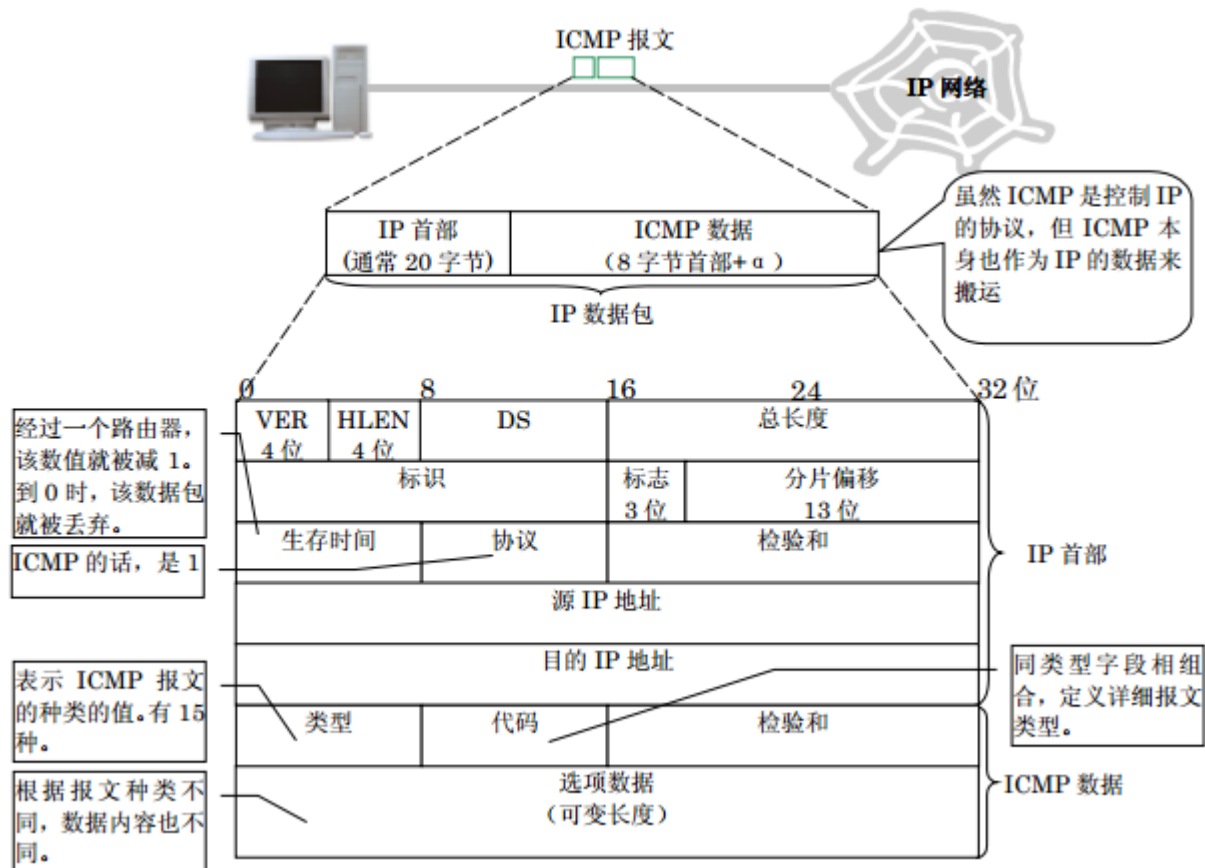
4. 网际控制报文协议ICMP

4.2 ICMP报文



4. 网际控制报文协议ICMP

4.2 ICMP报文



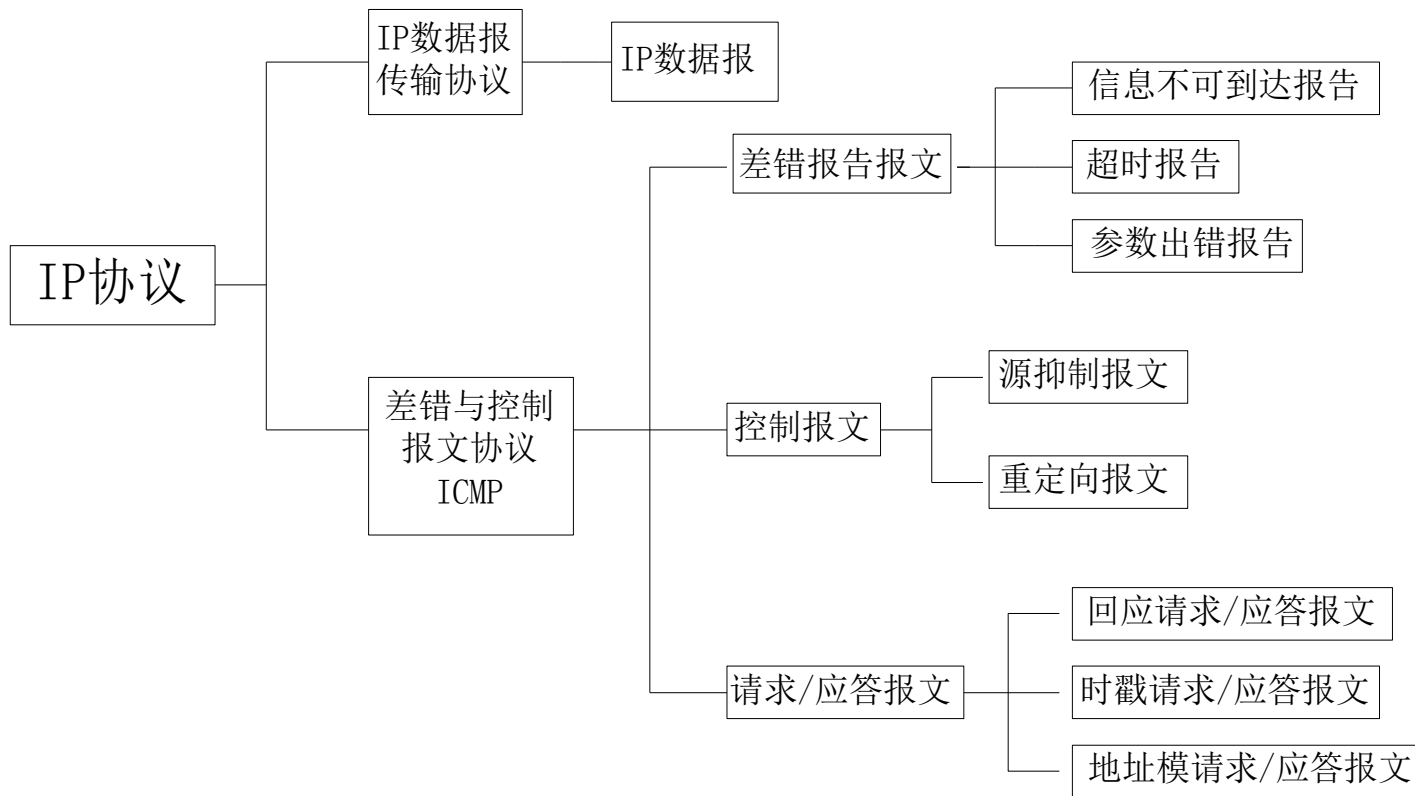
4.网际控制报文协议ICMP

4.2 ICMP报文

- ICMP报文的种类有两种，即**ICMP差错报告报文**和**ICMP询问报文**。
- ICMP报文的前4个字节是统一的格式，共有三个字段：即类型、代码和检验和。接着的4个字节的内容与ICMP的类型有关。

4. 网际控制报文协议ICMP

4.2 ICMP报文



4.网际控制报文协议ICMP

4.2 ICMP报文

- ICMP差错报告报文共有5种：
 - 终点不可达
 - 源点抑制(Source quench)
 - 时间超过
 - 参数问题
 - 改变路由（重定向）(Redirect)

4. 网际控制报文协议ICMP

4.2 ICMP报文

- 不应发送ICMP差错报告报文的几种情况：
 - 对ICMP差错报告报文不再发送ICMP差错报告报文。
 - 对第一个分片的数据报片的所有后续数据报片都不发送ICMP差错报告报文。
 - 对具有多播地址的数据报都不发送ICMP差错报告报文。
 - 对具有特殊地址（如127.0.0.0或0.0.0.0）的数据报不发送ICMP差错报告报文。

4. 网际控制报文协议ICMP

4.2 ICMP报文

- ICMP询问报文有两种：
 - 回送请求和回答报文
 - 时间戳请求和回答报文

- 停止使用的几种ICMP报有：
 - 信息请求与回答报文
 - 掩码地址请求和回答报文
 - 路由器询问和通告报文

类型	代码	名称	查询	差错
0	0	回应该答(Echo Reply)	√	
3		目的地不可达		√
	0	网路不可达		√
	1	主机不可达		√
	2	协议不可达		√
	3	端口不可达		√
	4	需要分片和不需要分片标记置位		√
	5	源路由失败		√
	6	目的网络未知		√
	7	目的主机未知		√
	8	源主机被隔离		√
	9	目的网络的通告被禁止		√
	10	目的主机的通信被禁止		√
	11	对请求的服务类型 ToS, 目的网路不可达		√
	12	对请求的服务类型 ToS, 目的主机不可达		√
	13	由于过滤,通信被强制禁止		√
	14	主机越权		√
	15	优先权中止生效		√
4	0	源抑制 (Source Quench)		√
5		重定向		√
	0	为网络 (子网) 重定向数据报		√
	1	为主机重定向数据报		√
	2	为网络和服务类型重定向数据报		√
	3	为主机和服务类型重定向数据报		√
6	0	选择主机地址		
8	0	请求回应	√	
9	0	路由器通告	√	
10	0	路由器选择请求	√	
11		超时		
	0	传输中超出 TTL=0		√
	1	分片重组 TTL=0		√
12		参数问题		
	0	指定错误的指针(坏的 IP 头部)		√
	1	缺少需要的选项		√
	2	错误长度		
13	0	时间戳请求		√
14	0	时间戳回复		√
15	0	信息请求 (已作废不用)		√
16	0	信息回复 (已作废不用)		√
17	0	地址掩码请求		√
18	0	地址掩码回复		√
30		跟踪路由		
31		数据报会话错误		
32		移动主机重定向		
33		IPv6 你在哪里		
34		IPv6 我在这里		
35		移动注册请求		
36		移动注册回复		

4.网际控制报文协议ICMP

4.3 Ping

- Ping程序是ICMP协议的最常见应用程序，由Mike Muuss 编写，可以用来测试目的主机是否可到达。
- Ping程序使用ICMP回应请求和回应应答报文来实现。

4.网际控制报文协议ICMP

4.3 Ping

- 当调用ping程序时，它发送一个包含ICMP回应请求的报文给目的地，然后等待一段很短的时间。如果没有收到应答，则重新传送请求。如果重传的请求仍没有收到应答（或收到一个ICMP目的不可达报文），ping报告该远程机器为不可达。
- 远端主机上的ICMP软件应答该回应请求报文。
- 按照协议只要收到回应请求，ICMP软件必须发送回应应答。

4. 网际控制报文协议ICMP

4.3 Ping



```
C:\Windows\system32\cmd.exe

C:\Users\RuanXiaolong>ping www.baidu.com

正在 Ping www.a.shifen.com [119.75.218.77] 具有 32 字节的数据:
来自 119.75.218.77 的回复: 字节=32 时间=12ms TTL=52
来自 119.75.218.77 的回复: 字节=32 时间=12ms TTL=52
来自 119.75.218.77 的回复: 字节=32 时间=12ms TTL=52
来自 119.75.218.77 的回复: 字节=32 时间=12ms TTL=52

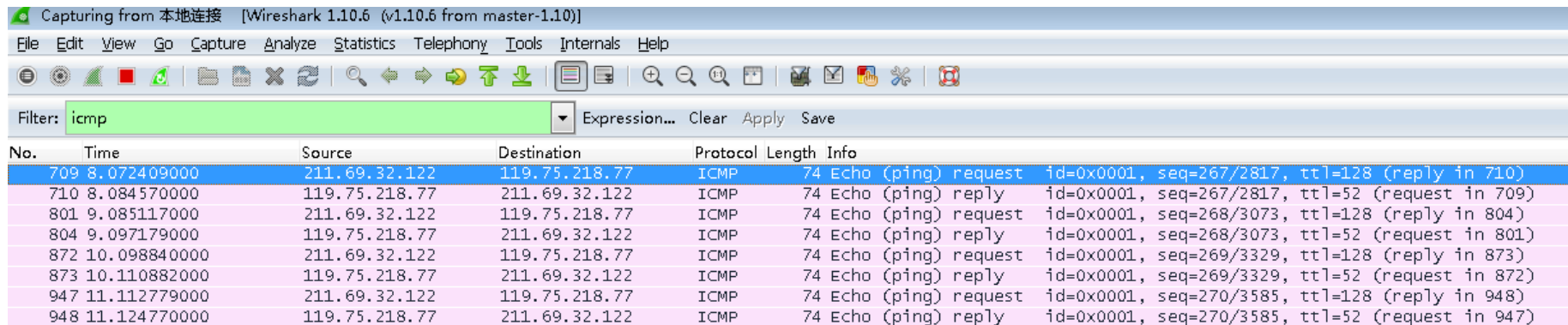
119.75.218.77 的 Ping 统计信息:
    数据包: 已发送 = 4, 已接收 = 4, 丢失 = 0 (0% 丢失),
    往返行程的估计时间(以毫秒为单位):
        最短 = 12ms, 最长 = 12ms, 平均 = 12ms

C:\Users\RuanXiaolong>
```

执行Ping操作：ping www.baidu.com

4. 网际控制报文协议ICMP

4.3 Ping



Capturing from 本地连接 [Wireshark 1.10.6 (v1.10.6 from master-1.10)]

File Edit View Go Capture Analyze Statistics Telephony Tools Internals Help

Filter: icmp

No.	Time	Source	Destination	Protocol	Length	Info
709	8.072409000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request id=0x0001, seq=267/2817, ttl=128 (reply in 710)
710	8.084570000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply id=0x0001, seq=267/2817, ttl=52 (request in 709)
801	9.085117000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request id=0x0001, seq=268/3073, ttl=128 (reply in 804)
804	9.097179000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply id=0x0001, seq=268/3073, ttl=52 (request in 801)
872	10.098840000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request id=0x0001, seq=269/3329, ttl=128 (reply in 873)
873	10.110882000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply id=0x0001, seq=269/3329, ttl=52 (request in 872)
947	11.112779000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request id=0x0001, seq=270/3585, ttl=128 (reply in 948)
948	11.124770000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply id=0x0001, seq=270/3585, ttl=52 (request in 947)

执行Ping操作：ping www.baidu.com
通过Wireshark捕获的ICMP数据报文

4.

Wireshark 1.10.6 (v1.10.6 from master-1.10)

Filter: icmp

No.	Time	Source	Destination	Protocol	Length	Info
709	8.072409000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
710	8.084570000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
801	9.085117000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
804	9.097179000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
872	10.098840000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
873	10.110882000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
947	11.112779000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
948	11.124770000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
1771	228.666000000	211.69.32.130	211.69.32.15	ICMP	132	Destination unreachable

Frame 709: 74 bytes on wire (592 bits), 74 bytes captured (592 bytes) on interface 0
 Interface id: 0
 Encapsulation type: Ethernet (1)
 Arrival Time: Apr 23, 2014 02:25:31.592503000
 [Time shift for this packet: 0.000000000 seconds]
 Epoch Time: 1398191131.592503000 seconds
 [Time delta from previous captured frame: 0.006751000 seconds]
 [Time delta from previous displayed frame: 0.000000000 seconds]
 [Time since reference or first frame: 8.072409000 seconds]
 Frame Number: 709
 Frame Length: 74 bytes (592 bits)
 Capture Length: 74 bytes (592 bits)
 [Frame is marked: False]
 [Frame is ignored: False]
 [Protocols in frame: eth:ip:icmp:data]
 [Coloring rule Name: ICMP]
 [Coloring rule string: icmp || icmpv6]

Ethernet II, Src: Vmware_16:25:97 (00:0c:29:16:25:97), Dst: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d)
 Destination: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d)
 Source: Vmware_16:25:97 (00:0c:29:16:25:97)
 Type: IP (0x0800)

Internet Protocol Version 4, Src: 211.69.32.122 (211.69.32.122), Dst: 119.75.218.77 (119.75.218.77)
 Version: 4
 Header length: 20 bytes
 Differentiated Services Field: 0x00 (DSCP 0x00: Default; ECN: 0x00: Not-ECT (Not ECN-Capable Tra
 Total Length: 60
 Identification: 0x1648 (5704)
 Flags: 0x00
 Fragment offset: 0
 Time to live: 128
 Protocol: ICMP (1)
 Header checksum: 0x0000 [validation disabled]
 Source: 211.69.32.122 (211.69.32.122)
 Destination: 119.75.218.77 (119.75.218.77)
 [Source GeoIP: Unknown]
 [Destination GeoIP: Unknown]

Internet Control Message Protocol
 Type: 8 (Echo (ping) request)
 Code: 0
 Checksum: 0x4c50 [correct]
 Identifier (BE): 1 (0x0001)
 Identifier (LE): 256 (0x0100)
 Sequence number (BE): 267 (0x010b)
 Sequence number (LE): 2817 (0x0b01)
 [Response frame: 710]

Data (32 bytes)
 Data: 6162636465666768696a6b6c6d6e6f707172737475767761...
 [Length: 32]

```

0000 e4 68 a3 a3 fa 7d 00 0c 29 16 25 97 08 00 45 00  . . . . .) . . . . .E
0010 00 3c 16 48 00 00 80 01 00 00 d3 45 20 7a 77 4b  < . H . . . . .E 2wK
0020 da 4d 08 00 4c 50 00 01 01 0b 61 62 63 64 65 66  . M . L P . . . . abcdef
0030 67 68 69 6a 6b 6c 6d 6e 6f 70 71 72 73 74 75 76  ghijklmn opqrstuv
0040 77 61 62 63 64 65 66 67 68 69                      wabcedfg hi
  
```

本地连接 <live capture in progress> File: ... Packets: 30607 · Displayed: 13 (0.0%)

Wireshark 1.10.6 (v1.10.6 from master-1.10)

Filter: icmp

No.	Time	Source	Destination	Protocol	Length	Info
709	8.072409000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
710	8.084570000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
801	9.085117000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
804	9.097179000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
872	10.098840000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
873	10.110882000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
947	11.112779000	211.69.32.122	119.75.218.77	ICMP	74	Echo (ping) request
948	11.124770000	119.75.218.77	211.69.32.122	ICMP	74	Echo (ping) reply
1771	228.666000000	211.69.32.130	211.69.32.15	ICMP	132	Destination unreachable

Frame 710: 74 bytes on wire (592 bits), 74 bytes captured (592 bytes) on interface 0
 Interface id: 0
 Encapsulation type: Ethernet (1)
 Arrival Time: Apr 23, 2014 02:25:31.604664000
 [Time shift for this packet: 0.000000000 seconds]
 Epoch Time: 1398191131.604664000 seconds
 [Time delta from previous captured frame: 0.012161000 seconds]
 [Time delta from previous displayed frame: 0.012161000 seconds]
 [Time since reference or first frame: 8.084570000 seconds]
 Frame Number: 710
 Frame Length: 74 bytes (592 bits)
 Capture Length: 74 bytes (592 bits)
 [Frame is marked: False]
 [Frame is ignored: False]
 [Protocols in frame: eth:ip:icmp:data]
 [Coloring rule Name: ICMP]
 [Coloring rule string: icmp || icmpv6]

Ethernet II, Src: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d), Dst: Vmware_16:25:97 (00:0c:29:16:25:97)
 Destination: Vmware_16:25:97 (00:0c:29:16:25:97)
 Source: HuaweiTe_a3:fa:7d (e4:68:a3:a3:fa:7d)
 Type: IP (0x0800)

Internet Protocol Version 4, Src: 119.75.218.77 (119.75.218.77), Dst: 211.69.32.122 (211.69.32.122)
 Version: 4
 Header length: 20 bytes
 Differentiated Services Field: 0x00 (DSCP 0x00: Default; ECN: 0x00: Not-ECT (Not ECN-Capable T
 Total Length: 60
 Identification: 0x1648 (5704)
 Flags: 0x00
 Fragment offset: 0
 Time to live: 52
 Protocol: ICMP (1)
 Header checksum: 0x2b21 [validation disabled]
 Source: 119.75.218.77 (119.75.218.77)
 Destination: 211.69.32.122 (211.69.32.122)
 [Source GeoIP: Unknown]
 [Destination GeoIP: Unknown]

Internet Control Message Protocol
 Type: 0 (Echo (ping) reply)
 Code: 0
 Checksum: 0x5450 [correct]
 Identifier (BE): 1 (0x0001)
 Identifier (LE): 256 (0x0100)
 Sequence number (BE): 267 (0x010b)
 Sequence number (LE): 2817 (0x0b01)
 [Request frame: 709]

Data (32 bytes)
 Data: 6162636465666768696a6b6c6d6e6f707172737475767761...
 [Length: 32]

```

0000 00 0c 29 16 25 97 e4 68 a3 a3 fa 7d 08 00 45 00  . . . . .) . . . . .E
0010 00 3c 16 48 00 00 80 01 01 0b 61 62 63 64 65 66  < . H . . . . .E +!wK.M.E
0020 70 7a 00 00 54 50 00 01 01 0b 61 62 63 64 65 66  z . TP . . . . abcdef
0030 67 68 69 6a 6b 6c 6d 6e 6f 70 71 72 73 74 75 76  ghijklmn opqrstuv
0040 77 61 62 63 64 65 66 67 68 69                      wabcedfg hi
  
```

本地连接 <live capture in progress> File: ... Packets: 32657 · Displayed: 13 (0.0%)

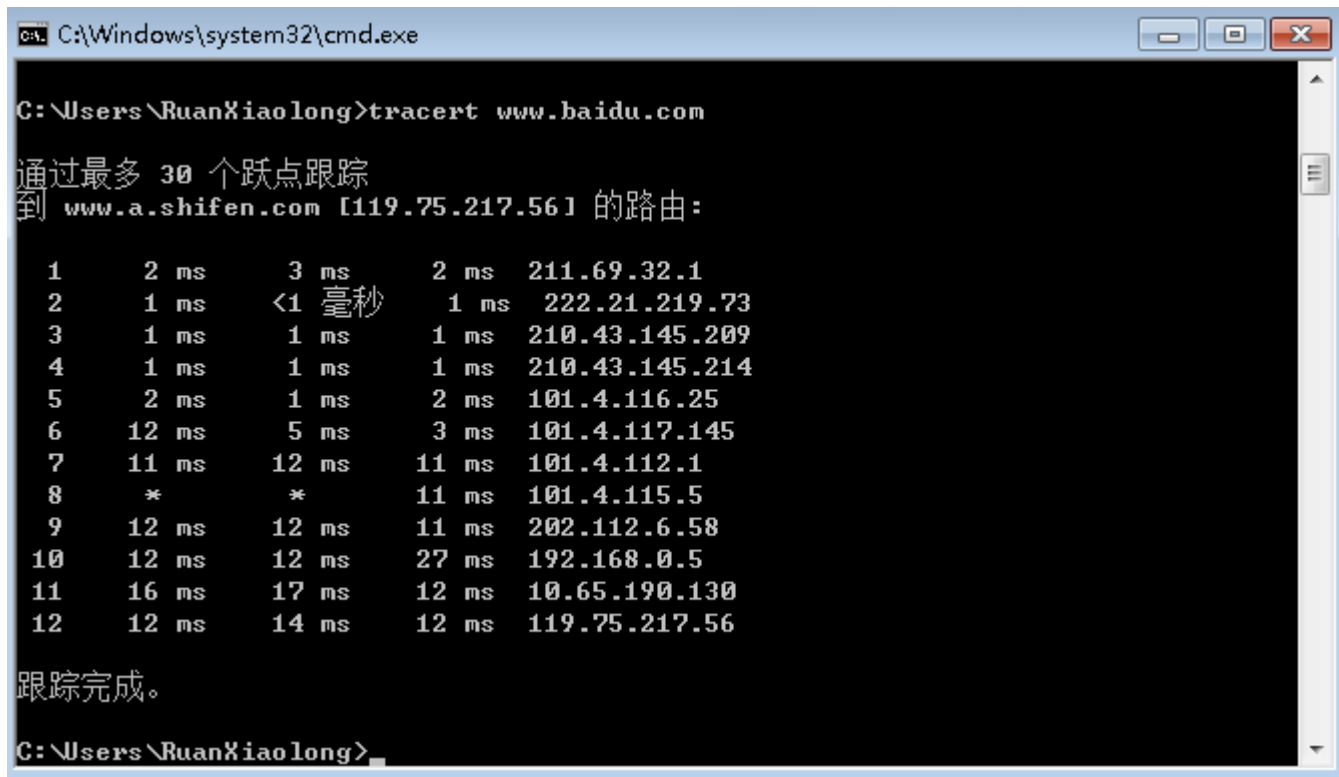
4.网际控制报文协议ICMP

4.4 tracert

- windows的tracert和linux/UNIX/router的traceroute都用于探测数据包从源到目的经过路由的IP，但两者探测的方法却有差别。
- 默认情况下，tracert是向目的地址发出ICMP请求回显数据包，而traceroute是向目的地址的某个端口（大于30000）发送UDP数据报。
- 两者用于探测的数据类型不同。但他们也有一个共同点：都是通过设置发送包的TTL的值从1开始、逐次增1的方法来探测。

4. 网际控制报文协议ICMP

4.4 tracert



```
C:\Windows\system32\cmd.exe

C:\Users\RuanXiaolong>tracert www.baidu.com

通过最多 30 个跃点跟踪
到 www.a.shifen.com [119.75.217.56] 的路由:

 1      2 ms     3 ms     2 ms    211.69.32.1
 2      1 ms     <1 毫秒  1 ms    222.21.219.73
 3      1 ms     1 ms     1 ms    210.43.145.209
 4      1 ms     1 ms     1 ms    210.43.145.214
 5      2 ms     1 ms     2 ms    101.4.116.25
 6     12 ms     5 ms     3 ms    101.4.117.145
 7     11 ms    12 ms    11 ms    101.4.112.1
 8      *      *      11 ms    101.4.115.5
 9     12 ms    12 ms    11 ms    202.112.6.58
10     12 ms    12 ms    27 ms    192.168.0.5
11     16 ms    17 ms    12 ms    10.65.190.130
12     12 ms    14 ms    12 ms    119.75.217.56

跟踪完成。

C:\Users\RuanXiaolong>
```

No.	Time	Source	Destination	Protocol	Length	Info
776	7.100759000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=271/3841, ttl=1
777	7.103338000	211.69.32.1	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
778	7.105919000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=272/4097, ttl=1
779	7.108839000	211.69.32.1	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
780	7.110844000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=273/4353, ttl=1
781	7.113511000	211.69.32.1	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
1311	12.971001000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=274/4609, ttl=2
1312	12.971850000	222.21.219.73	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
1313	12.974177000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=275/4865, ttl=2
1314	12.975011000	222.21.219.73	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
1315	12.977696000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=276/5121, ttl=2
1316	12.978554000	222.21.219.73	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
1852	18.524059000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=277/5377, ttl=3
1853	18.524992000	210.43.145.209	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
1854	18.527417000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=278/5633, ttl=3
1855	18.528302000	210.43.145.209	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
1856	18.530291000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=279/5889, ttl=3
1857	18.531190000	210.43.145.209	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
2488	24.077668000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=280/6145, ttl=4
2489	24.079077000	210.43.145.214	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
2490	24.081326000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=281/6401, ttl=4
2491	24.082795000	210.43.145.214	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
2492	24.085184000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=282/6657, ttl=4
2493	24.086592000	210.43.145.214	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
3040	29.631898000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=283/6913, ttl=5
3041	29.633837000	101.4.116.25	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
3042	29.636127000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=284/7169, ttl=5
3043	29.637911000	101.4.116.25	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
3046	29.640450000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=285/7425, ttl=5
3047	29.642757000	101.4.116.25	211.69.32.122	ICMP	70	Time-to-live exceeded (Time to live exceeded in transit)
3657	35.184932000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=286/7681, ttl=6
3658	35.196870000	101.4.117.145	211.69.32.122	ICMP	110	Time-to-live exceeded (Time to live exceeded in transit)
3659	35.199353000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=287/7937, ttl=6
3660	35.204513000	101.4.117.145	211.69.32.122	ICMP	110	Time-to-live exceeded (Time to live exceeded in transit)
3661	35.206887000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=288/8193, ttl=6
3662	35.210351000	101.4.117.145	211.69.32.122	ICMP	110	Time-to-live exceeded (Time to live exceeded in transit)
4212	40.754192000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=289/8449, ttl=7
4213	40.765805000	101.4.112.1	211.69.32.122	ICMP	182	Time-to-live exceeded (Time to live exceeded in transit)
4214	40.768497000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=290/8705, ttl=7
4215	40.780166000	101.4.112.1	211.69.32.122	ICMP	182	Time-to-live exceeded (Time to live exceeded in transit)
4216	40.782562000	211.69.32.122	119.75.217.56	ICMP	106	Echo (ping) request id=0x0001, seq=291/8961, ttl=7

5.路由选择协议

- 网络层的主要功能是将分组从源节点路由到目的节点中，而且在大多数计算机网络中，采用的是数据报分组交换方式，数据报分组需要经过多跳（Hop）才能到达目的地。
- 路由功能是一种数据报分组交换路径选择行为，是网络层的一种基本功能。
- 路由功能和日常旅行时选择最佳路线的道理是相通的，路由选择就是综合考虑多种因素，例如线路长度、信道带宽、线路稳定性等。

5.路由选择协议

5.1路由的分类

- 路由（Routing）是把信息从源节点通过网络传送到目的节点的行为，简单的讲路由就是指网络层设备从一个接口上收到数据包，根据数据包的目的地址进行定向，并转发到另一个接口的过程。
- 路由与桥接对比的主要区别在于：桥接发生在数据链路层，连接的是同一网络或同一子网的不同网段。路由发生在网络层，连接的是不同网络或不同子网。

5.路由选择协议

5.1路由的分类

- 路由功能的实现是依靠路由器或路由交换机中的路由表进行的。
- 路由分为静态路由 (Static Routing) 和动态路由 (Dynamic Routing) 两大类。

5.路由选择协议

5.1路由的分类

- **静态路由：**
- 静态路由是手动配置的路由，特别是在小型局域网中，因为静态路由的配置和管理都比较简单。
- 静态路由的特点是：
 - 手动配置：静态路由需要管理员手动进行逐条配置。
 - 路由路径固定不变：静态路由不会随着网络的拓扑结构或链路的状态变化而变化，是静态固定的。
 - 不可通告性：静态路由信息是私有的，不会通告给其他路由器。
 - 单向性：静态路由仅为数据提供沿着下一跳的方向进行路由，不提供反向路由。

5.路由选择协议

5.1路由的分类

- **静态路由：**
- 静态路由明确的指明了到达目的网络，所以在所有同目的地址的路由中，静态路由的优先级是除“直连路由”外最高的，也就是说如果配置了到达某一网络或者某一结点的静态路由，则优先采用这条静态路由，只有当该条静态路由不可用时，才会考虑其他路由。
- 静态路由一般适合比较简单的小型网络环境，因为在这样环境中，网络管理员易于清楚的了解网络的拓扑结构，能够设置正确的路由信息。

5.路由选择协议

5.1路由的分类

- 动态路由：
- 对于较为大型的广域网来说，由于拓扑结构复杂，且网络结构可能经常变动，通常会采用更加灵活、更具自动特性的动态路由。



5.路由选择协议

5.1路由的分类

□ 动态路由的特点：

- 自动生成：路由器在启动了动态路由协议后，将会通告所直接连接的网络，则路由器间就会自动生成路由器直接连接的网络间的路由表项。
- 自动调整：当网络结构发生改变，动态路由可以随时根据网络拓扑结构的变化调整路由表项，并删除无效的路由表项。
- 自动通告：动态路由可以在相邻路由器上相互通告，以便及时反映拓扑结构的变化，生成新的动态路由表项。
- 自动生成双向路由：路由器在生成某条路由的动态路由时会自动生成回程路由表项，也就是会同时双向路由表项。

5.路由选择协议

5.1路由的分类

□ 动态路由的特点：

- 仅可生成网络间的路由表项：动态路由不能够生成到达具体节点或主机的动态路由表项。
- 不同动态路由不兼容：动态路由根据所采用的算法的不同分为不同类型，如RIP、OSPF、EIGRP、IS-IS、BGP等。不同的动态路由协议主要适用的网络环境不一样，也是不兼容的。但是支持互相重发布。

5.路由选择协议

5.2有关路由选择协议的几个基本概念

- 理想的路由算法：
- 路由选择协议的核心是路由算法，即需要何种算法来获得路由表中的各项目。
- 路由算法应具有的特点是：
 - 算法必须是正确的和完整的。
 - 算法在计算上应简单。
 - 算法应能适应通信量和网络拓扑的变化，这就是说，要有自适应性。
 - 算法应具有稳定性。
 - 算法应是公平的。
 - 算法应是最佳的

5.路由选择协议

5.2有关路由选择协议的几个基本概念

□ 最佳路由：

- 不存在一种绝对的最佳路由算法。
- 所谓“最佳”只能是相对于某一种特定要求下得出的较为合理的选择而已。
- 实际的路由选择算法，应尽可能接近于理想的算法。
- 路由选择是个非常复杂的问题：它是网络中的所有结点共同协调工作的结果。路由选择的环境往往是不断变化的，而这种变化有时无法事先知道。

5.路由选择协议

5.2有关路由选择协议的几个基本概念

□ 分层次的路由选择协议：

- 因特网采用分层次的路由选择协议。
- 因特网的规模非常大。如果让所有的路由器知道所有的网络应怎样到达，则这种路由表将非常大，处理起来也太花时间。而所有这些路由器之间交换路由信息所需的带宽就会使因特网的通信链路饱和。
- 许多单位不愿意外界了解自己单位网络的布局细节和本部门所采用的路由选择协议（这属于本部门内部的事情），但同时还希望连接到因特网上。

5.路由选择协议

5.2有关路由选择协议的几个基本概念

- 自治系统 (Autonomous System , AS) :
 - 自治系统AS的定义：在单一的技术管理下的一组路由器，而这些路由器使用一种AS内部的路由选择协议和共同的度量以确定分组在该AS内的路由，同时还使用一种AS之间的路由选择协议用以确定分组在AS之间的路由。
 - 现在对自治系统AS的定义是强调下面的事实：尽管一个AS使用了多种内部路由选择协议和度量，但重要的是一个AS对其他AS表现出的是一个单一的和一致的路由选择策略。

5.路由选择协议

5.2有关路由选择协议的几个基本概念

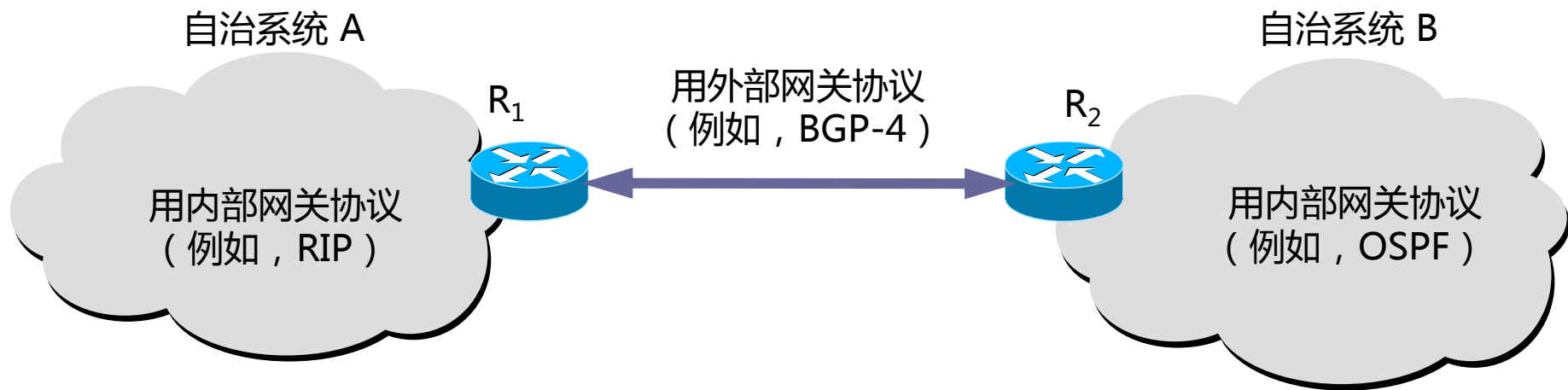
□ 因特网的两大类路由选择协议：

- 内部网关协议 IGP (Interior Gateway Protocol)：即在一个自治系统内部使用的路由选择协议。目前这类路由选择协议使用得最多，如RIP和OSPF协议。
- 外部网关协议 EGP (External Gateway Protocol)：若源站和目的站处在不同的自治系统中，当数据报传到一个自治系统的边界时，就需要使用一种协议将路由选择信息传递到另一个自治系统中。这样的协议就是外部网关协议EGP。在外部网关协议中目前使用最多的是BGP-4。

5. 路由选择协议

5.2 有关路由选择协议的几个基本概念

□ 因特网的两大类路由选择协议：



自治系统之间的路由选择叫做域间路由选择(interdomain routing)
自治系统内部的路由选择叫做域内路由选择(intradomain routing)

5.路由选择协议

5.3内部网关协议 RIP

- 路由信息协议(Routing Information Protocol,RIP)是内部网关协议IGP中最先得到广泛使用的协议。它的中文名称很少使用，叫做路由信息协议。
- RIP是一种分布式的基于距离向量的路由选择协议，是因特网的标准协议。
- RIP协议要求网络中的每一个路由器都要维护从它自己到其他每一个目的网络的距离记录。

5.路由选择协议

5.3内部网关协议 RIP

□ 距离：

- 从一路由器到直接连接的网络的距离定义为1。从一个路由器到非直接连接的网络的距离定义为所经过的路由器数加1。
- RIP协议中的“距离”也称为“跳数”(hop count)，因为每经过一个路由器，跳数就加1。
- RIP协议中的“距离”实际上指的是“最短距离”。RIP认为一个好的路由就是它通过的路由器的数目少，即“距离短”。
- RIP允许一条路径最多只能包含 15 个路由器，“距离”的最大值为16时即相当于不可达。
- RIP 不能在两个网络之间同时使用多条路由。RIP选择一个具有最少路由器的路由（即最短路由），哪怕还存在另一条高速(低时延)但路由器较多的路由。

5.路由选择协议

5.3内部网关协议 RIP

□ RIP协议的特点：

- 仅和相邻路由器交换信息。
- 交换的信息是当前本路由器所知道的全部信息，即自己的路由表。
- 按固定的时间间隔交换路由信息，例如，每隔30秒。然后路由器根据收到的路由信息更新路由表。当网络拓扑发生变化时，路由器也及时向相邻路由器通告拓扑变化后的路由信息。

5.路由选择协议

5.3内部网关协议 RIP

□ 路由表的建立过程：

- 路由器在刚刚开始工作时，只知道到直接连接的网络的距离（此距离定义为1）。
- 以后，每一个路由器也只和数目非常有限的相邻路由器交换并更新路由信息。
- 经过若干次更新后，所有的路由器最终都会知道到达本自治系统中任何一个网络的最短距离和下一跳路由器的地址。
- RIP协议的收敛(convergence)过程较快，即在自治系统中所有的结点都得到正确的路由选择信息的过程。

5.路由选择协议

5.3内部网关协议 RIP

- 距离向量算法：
- 收到相邻路由器（其地址为 X）的一个 RIP 报文：
 - 先修改此RIP报文中的所有项目：把“下一跳”字段中的地址都改为X，并把所有的“距离”字段的值加1。
 - 对修改后的RIP报文中的每一个项目，重复以下步骤：
 - 若项目中的目的网络不在路由表中，则把该项目加到路由表中。
 - 否则若下一跳字段给出的路由器地址是同样的，则把收到的项目替换原路由表中的项目。
 - 否则若收到项目中的距离小于路由表中的距离，则进行更新，
 - 否则，什么也不做。
 - 若3分钟还没有收到相邻路由器的更新路由表，则把此相邻路由器记为不可达路由器，即将距离置为16（距离为16表示不可达）。

5.路由选择协议

5.3内部网关协议 RIP

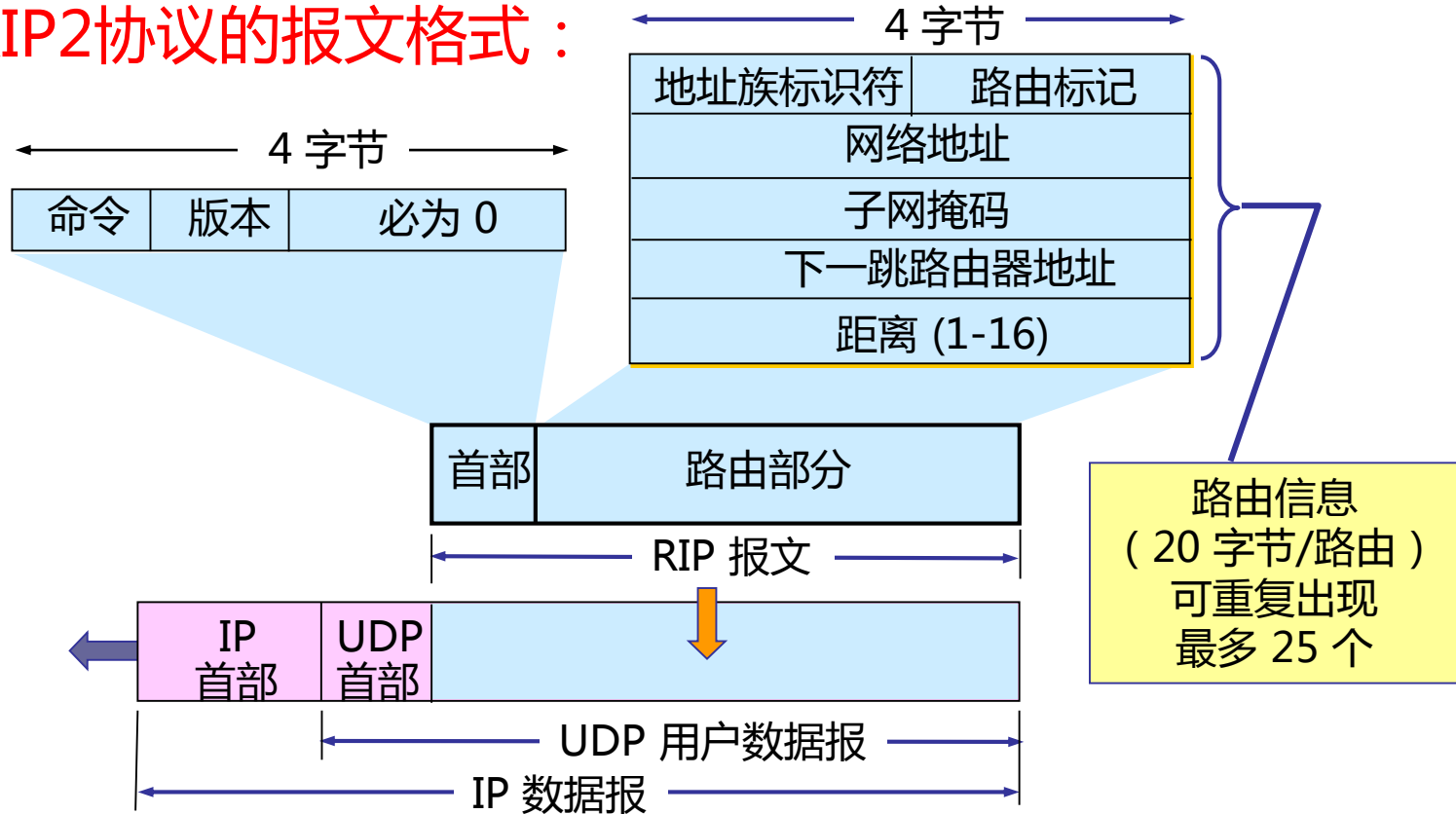
□ 路由器之间交换信息：

- RIP协议让互联网中的所有路由器都和自己的相邻路由器不断交换路由信息，并不断更新其路由表，使得从每一个路由器到每一个目的网络的路由都是最短的（即跳数最少）。
- 虽然所有的路由器最终都拥有了整个自治系统的全局路由信息，但由于每一个路由器的位置不同，它们的路由表当然也应当是不同的。

5.路由选择协议

5.3内部网关协议 RIP

□ RIP2协议的报文格式：



5.路由选择协议

5.3内部网关协议 RIP

□ RIP协议的优缺点：

- RIP存在的一个问题是当网络出现故障时，要经过比较长的时间才能将此信息传送到所有的路由器。
- RIP协议最大的优点就是实现简单，开销较小。
- RIP限制了网络的规模，它能使用的最大距离为15（16表示不可达）。
- 路由器之间交换的路由信息是路由器中的完整路由表，因而随着网络规模的扩大，开销也就增加。

5.路由选择协议

5.4内部网关协议 OSPF

- **开放最短路径优先**（Open Shortest Path First，OSPF）协议是为了克服RIP的缺点，在1989年开发出来的。
- “开放”表明OSPF协议不是受某一家厂商控制，而是公开发表的。
- “最短路径优先”是因为使用了Dijkstra提出的最短路径算法SPF。
- OSPF只是一个协议的名字，它并不表示其他的路由选择协议不是“最短路径优先”。实际上，所有的在自治系统内部使用的路由选择协议都是要寻找最短路径的。

5.路由选择协议

5.4内部网关协议 OSPF

- OSPF协议最主要的特征就是使用**分布式的链路状态协议**（Link State Protocol），而不是像RIP那样的距离向量协议。
- OSPF的三个要点：
 - 向本自治系统中所有路由器发送信息，这里使用的方法是洪泛法。
 - 发送的信息就是与本路由器相邻的所有路由器的链路状态，但这只是路由器所知道的部分信息。“链路状态”就是说明本路由器都和哪些路由器相邻，以及该链路的“度量” (metric)。
 - 只有当链路状态发生变化时，路由器才用洪泛法向所有路由器发送此信息。

5.路由选择协议

5.4内部网关协议 OSPF

- 链路状态数据库（ Link-state Database ）：
 - 由于各路由器之间频繁地交换链路状态信息，因此所有的路由器最终都能建立一个链路状态数据库。
 - 这个数据库实际上就是全网的拓扑结构图，它在全网范围内是一致的（这称为链路状态数据库的同步）。
 - OSPF的链路状态数据库能较快地进行更新，使各个路由器能及时更新其路由表。OSPF的更新过程收敛得快是其重要优点。

5.路由选择协议

5.4内部网关协议 OSPF

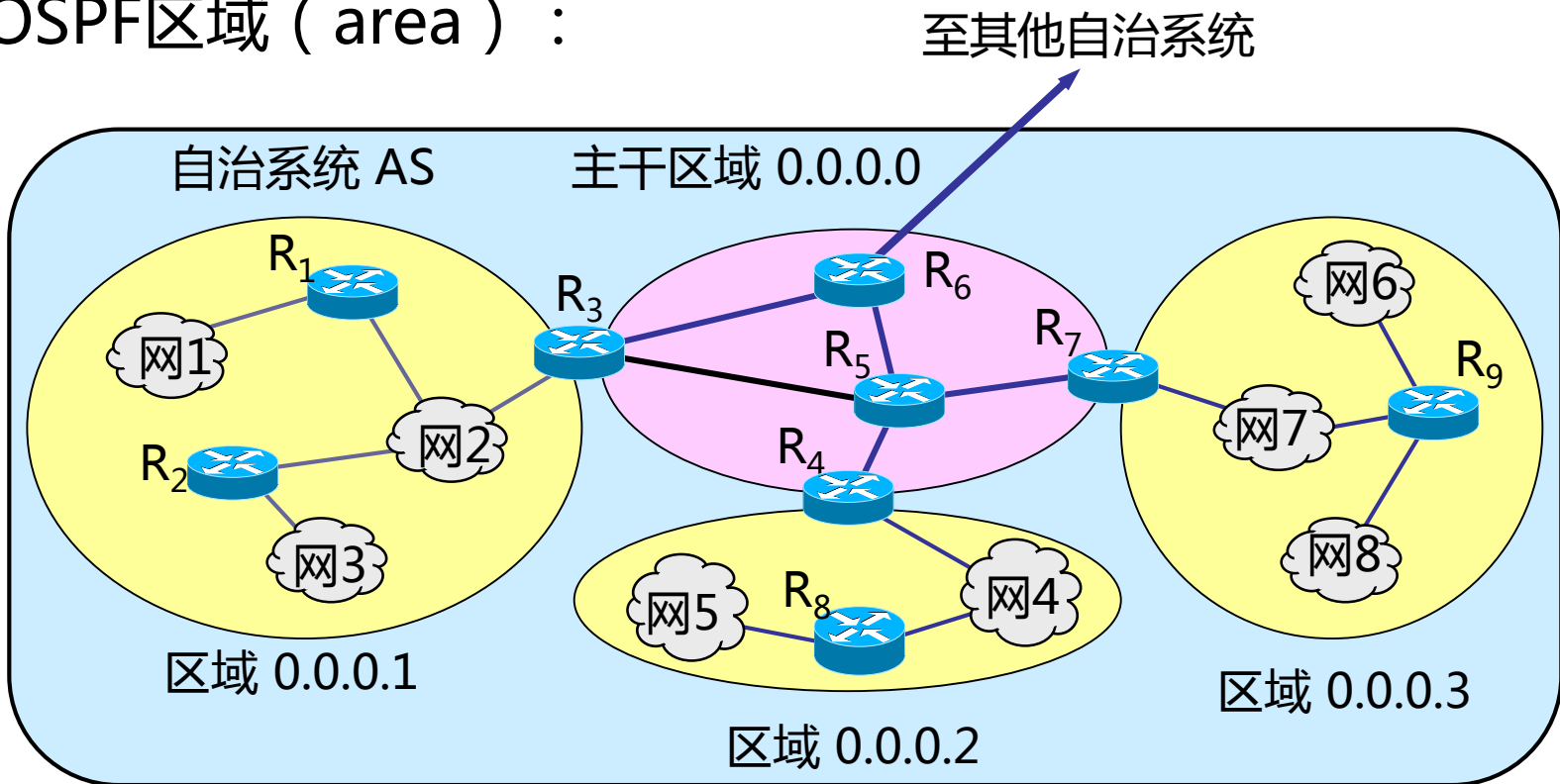
□ OSPF区域（area）：

- 为了使OSPF能够用于规模很大的网络，OSPF将一个自治系统再划分为若干个更小的范围，叫作区域。
- 每一个区域都有一个32位的区域标识符（用点分十进制表示）。
- 区域也不能太大，在一个区域内的路由器最好不超过200个。

5.路由选择协议

5.4内部网关协议 OSPF

□ OSPF区域 (area) :



5.路由选择协议

5.4内部网关协议 OSPF

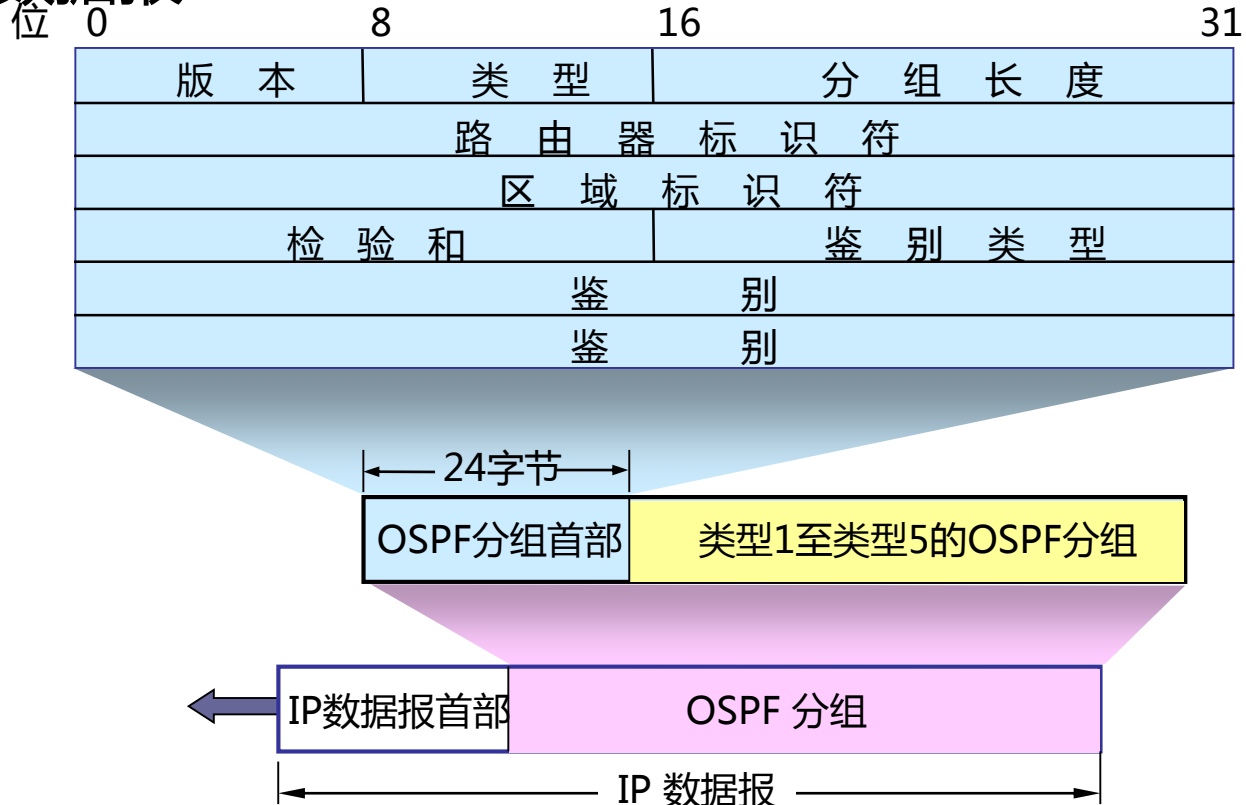
□ OSPF数据报：

- OSPF不用UDP而是直接用IP数据报传送。
- OSPF构成的数据报很短。这样做可减少路由信息的通信量。
- 数据报很短的另一好处是可以不必将长的数据报分片传送。分片传送的数据报只要丢失一个，就无法组装成原来的数据报，而整个数据报就必须重传。

5.路由选择协议

5.4内部网关协议 OSPF

□ OSPF数据报：



5.路由选择协议

5.4内部网关协议 OSPF

□ OSPF数据报：

- OSPF对不同的链路可根据IP分组的不同服务类型TOS而设置成不同的代价。因此，OSPF对于不同类型的业务可计算出不同的路由。
- 如果到同一个目的网络有多条相同代价的路径，那么可以将通信量分配给这几条路径。这叫作多路径间的负载平衡。
- 所有在OSPF路由器之间交换的分组都具有鉴别的功能。
- 支持可变长度的子网划分和无分类编址CIDR。
- 每一个链路状态都带上一个32位的序号，序号越大状态就越新。

5.路由选择协议

5.4内部网关协议 OSPF

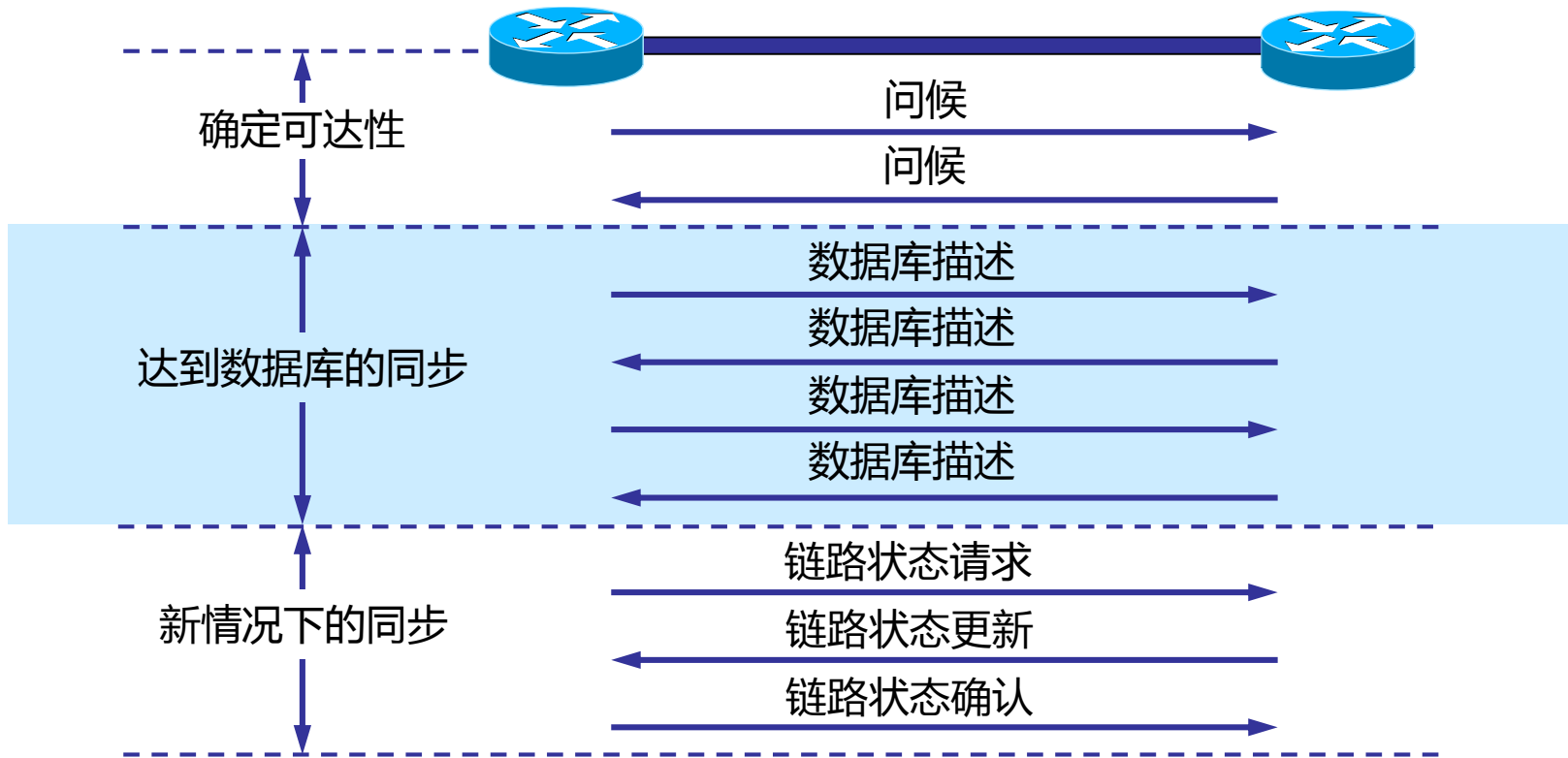
□ OSPF的五种分组类型：

- 类型1：问候(Hello)分组。
- 类型2：数据库描述(Database Description)分组。
- 类型3：链路状态请求(Link State Request)分组。
- 类型4：链路状态更新(Link State Update)分组，用洪泛法对全网更新链路状态。
- 类型5，链路状态确认(Link State Acknowledgment)分组。

5.路由选择协议

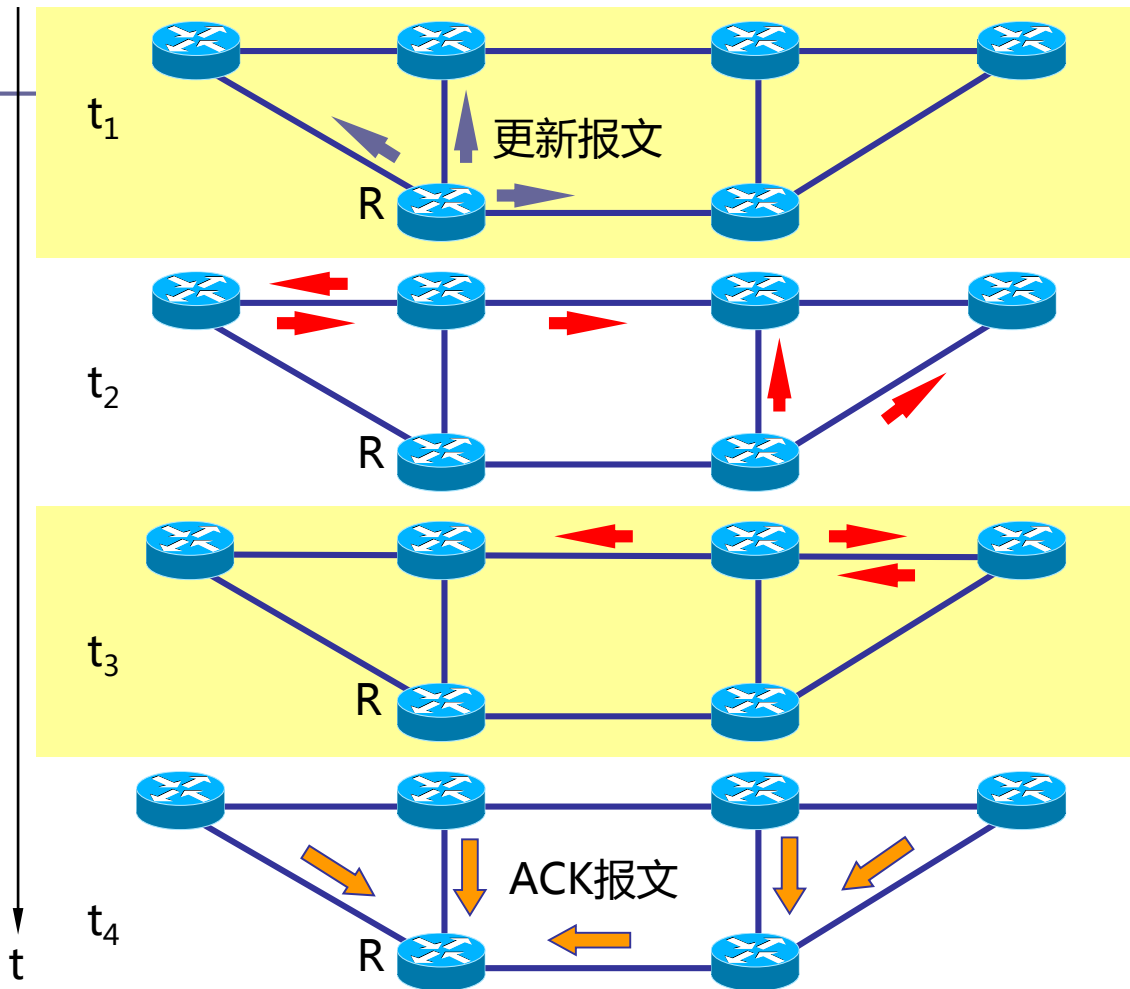
5.4内部网关协议 OSPF

□ OSPF的五种分组类型：



5. 路由选择协议

- OSPF使用的是可靠
- 的洪泛法：



5.路由选择协议

5.4内部网关协议 OSPF

□ OSPF的其他特点：

- OSPF规定每隔一段时间，如30分钟，要刷新一次数据库中的链路状态。
- 由于一个路由器的链路状态只涉及到与相邻路由器的连通状态，因而与整个互联网的规模并无直接关系。因此当互联网规模很大时，OSPF协议要比距离向量协议RIP好得多。
- OSPF没有“坏消息传播得慢”的问题，据统计，其响应网络变化的时间小于100ms。
- 多点接入的局域网采用了指定的路由器的方法，使广播的信息量大减少。

5.路由选择协议

5.4外部网关协议 BGP

- 边界网关协议BGP是不同自治系统的路由器之间交换路由信息的协议。
- BGP较新版本是2006年1月发表的BGP-4（BGP第4个版本），即RFC 4271 ~ 4278。
- 为了简单起见，将BGP-4简写为BGP。

5.路由选择协议

5.4外部网关协议 BGP

- 为什么在不同的AS之间不能够使用RIP或OSPF？
- 因特网的规模太大，使得自治系统之间路由选择非常困难。对于自治系统之间的路由选择，要寻找最佳路由是很不现实的。
 - 当一条路径通过几个不同 AS 时，要想对这样的路径计算出有意义的代价是不太可能的。
 - 比较合理的做法是在 AS 之间交换“可达性”信息。
- 自治系统之间的路由选择必须考虑有关策略。边界网关协议 BGP 只能是力求寻找一条能够到达目的网络且比较好的路由（不能兜圈子），而并非要寻找一条最佳路由。

5.路由选择协议

5.4外部网关协议 BGP

□ BGP发言人 (BGP speaker)

- 每一个自治系统的管理员要选择至少一个路由器作为该自治系统的“ BGP 发言人”。
- 一般说来，两个 BGP 发言人都是通过一个共享网络连接在一起的，而 BGP 发言人往往就是 BGP 边界路由器，但也可以不是 BGP 边界路由器。

5.路由选择协议

5.4外部网关协议 BGP

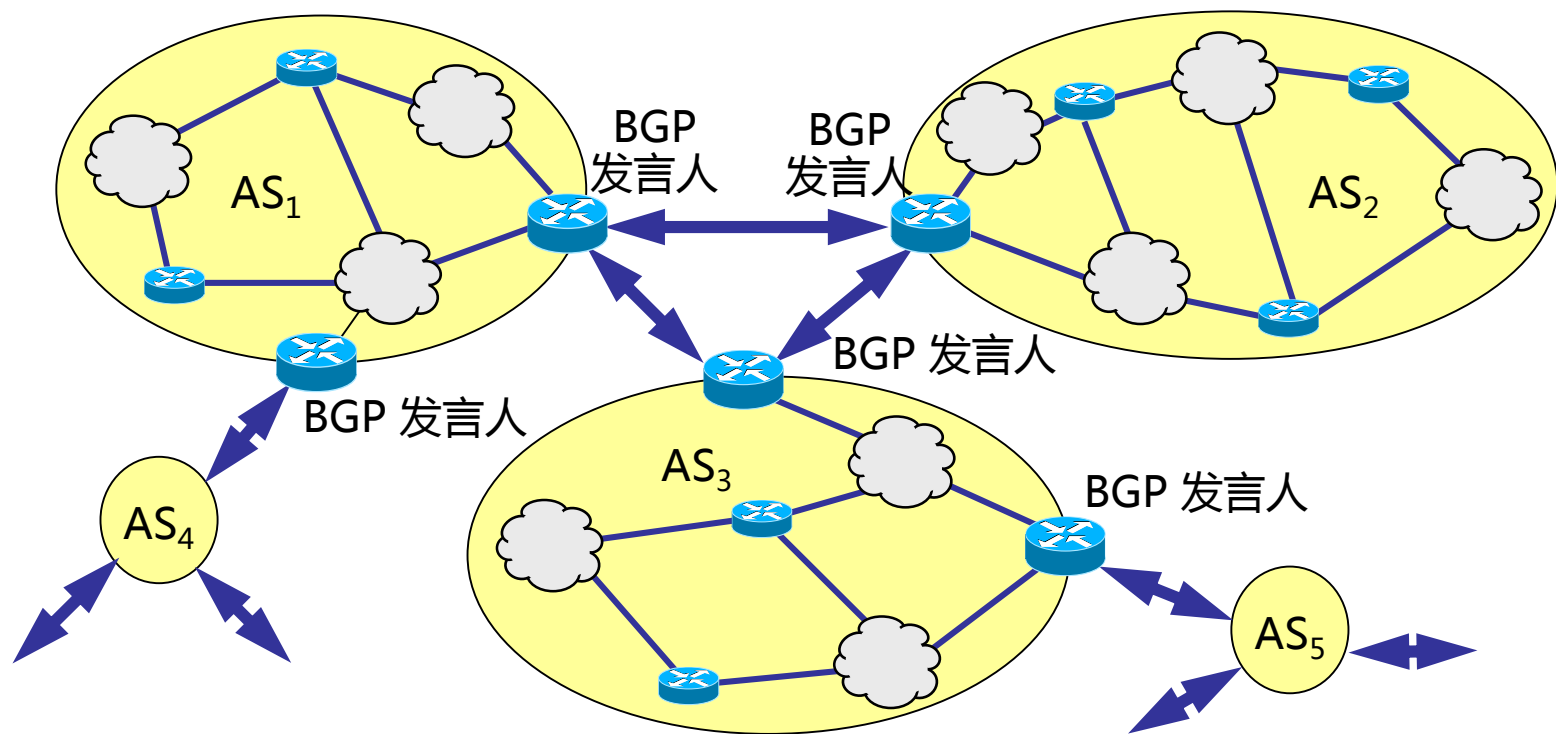
□ BGP发言人 (BGP speaker)

- 一个BGP发言人与其他自治系统中的BGP发言人要交换路由信息，就要先建立TCP连接，然后在此连接上交换BGP报文以建立BGP会话(session)，利用BGP会话交换路由信息。
- 使用TCP连接能提供可靠的服务，也简化了路由选择协议。
- 使用TCP连接交换路由信息的两个BGP发言人，彼此成为对方的邻站或对等站。

5.路由选择协议

5.4外部网关协议 BGP

□ BGP发言人自治系统AS的关系



5.路由选择协议

5.4外部网关协议 BGP

□ BGP协议的特点：

- BGP协议交换路由信息的结点数量级是自治系统数的量级，这要比这些自治系统中的网络数少很多。
- 每一个自治系统中BGP发言人（或边界路由器）的数目是很少的。这样就使得自治系统之间的路由选择不致过分复杂。
- **BGP支持CIDR**，因此BGP的路由表也就应当包括目的网络前缀、下一跳路由器，以及到达该目的网络所要经过的各个自治系统序列。
- 在BGP刚刚运行时，BGP的邻站是交换整个的BGP路由表。但以后只需要在发生变化时**更新有变化的部分**。这样做对节省网络带宽和减少路由器的处理开销方面都有好处。

5.路由选择协议

5.4外部网关协议 BGP

□ BGP-4共使用四种报文：

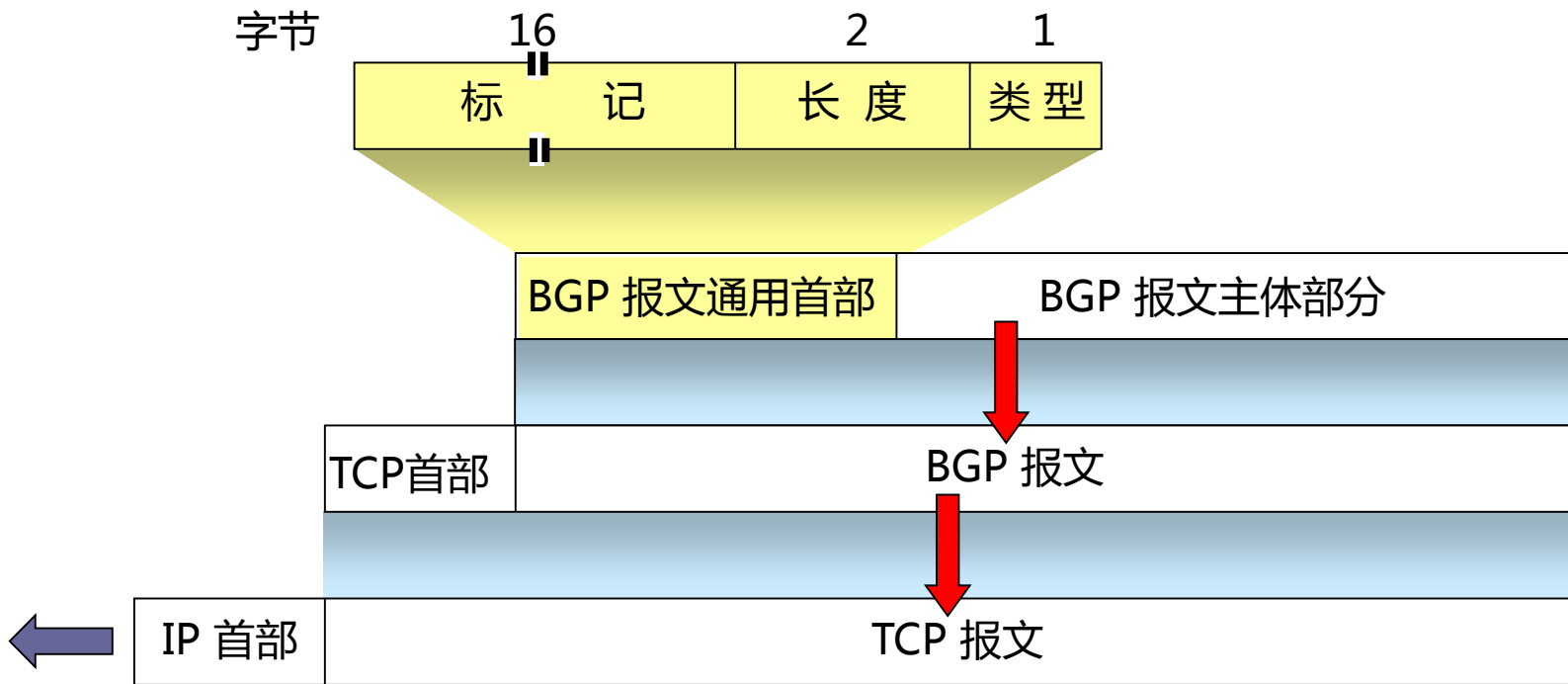
- 打开(OPEN)报文，用来与相邻的另一个BGP发言人建立关系。
- 更新(UPDATE)报文，用来发送某一路由的信息，以及列出要撤消的多条路由。
- 保活(KEEPALIVE)报文，用来确认打开报文和周期性地证实邻站关系。
- 通知(NOTIFICATION)报文，用来发送检测到的差错。

- 在RFC 2918中增加了ROUTE-REFRESH报文，用来请求对等端重新通告。

5.路由选择协议

5.4外部网关协议 BGP

□ BGP报文结构：



5.路由选择协议

5.5路由器的结构

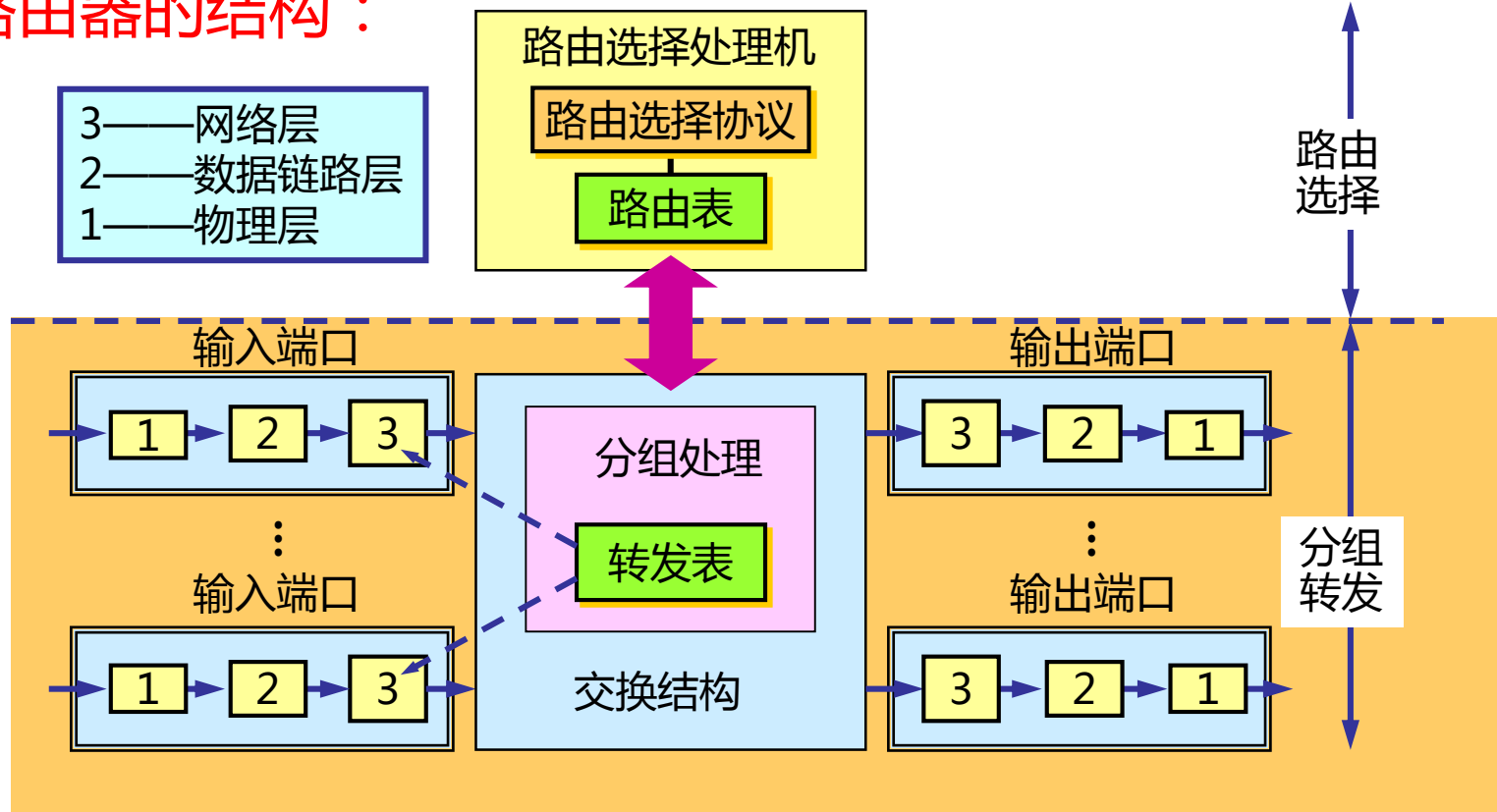
□ 路由器的结构：

- 路由器是一种具有多个输入端口和多个输出端口的专用计算机，其任务是转发分组。也就是说，将路由器某个输入端口收到的分组，按照分组要去的目的地（即目的网络），把该分组从路由器的某个合适的输出端口转发给下一跳路由器。
- 下一跳路由器也按照这种方法处理分组，直到该分组到达终点为止。

5. 路由选择协议

5.5 路由器的结构

□ 路由器的结构：



5.路由选择协议

5.5路由器的结构

- “转发” 和 “路由选择” 的区别：
 - “转发” (forwarding)就是路由器根据转发表将用户的 IP 数据报从合适的端口转发出去。
 - “路由选择” (routing)则是按照分布式算法，根据从各相邻路由器得到的关于网络拓扑的变化情况，动态地改变所选择的路由。
 - 路由表是根据路由选择算法得出的。而转发表是从路由表得出的。
 - 在讨论路由选择的原理时，往往不去区分转发表和路由表的区别。

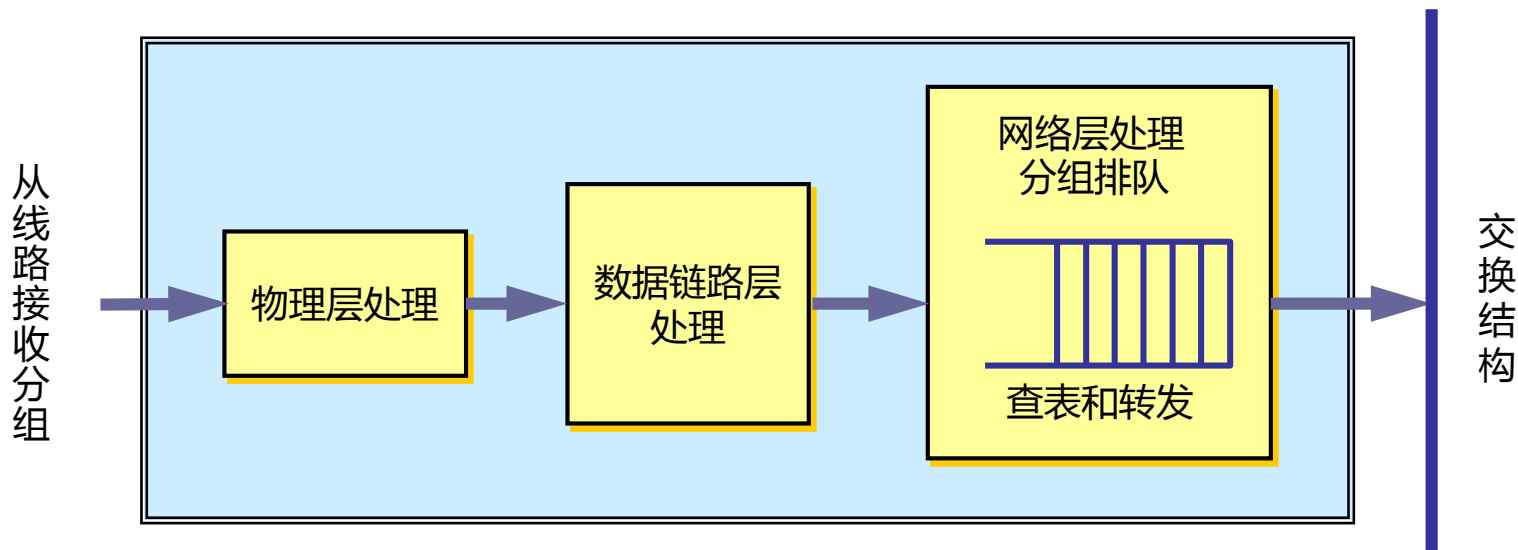
5.路由选择协议

5.5路由器的结构

□ 输入端口对线路上收到的分组的处理：

- 数据链路层剥去帧首部和尾部后，将分组送到网络层的队列中排队等待处理。这会产生一定的时延。

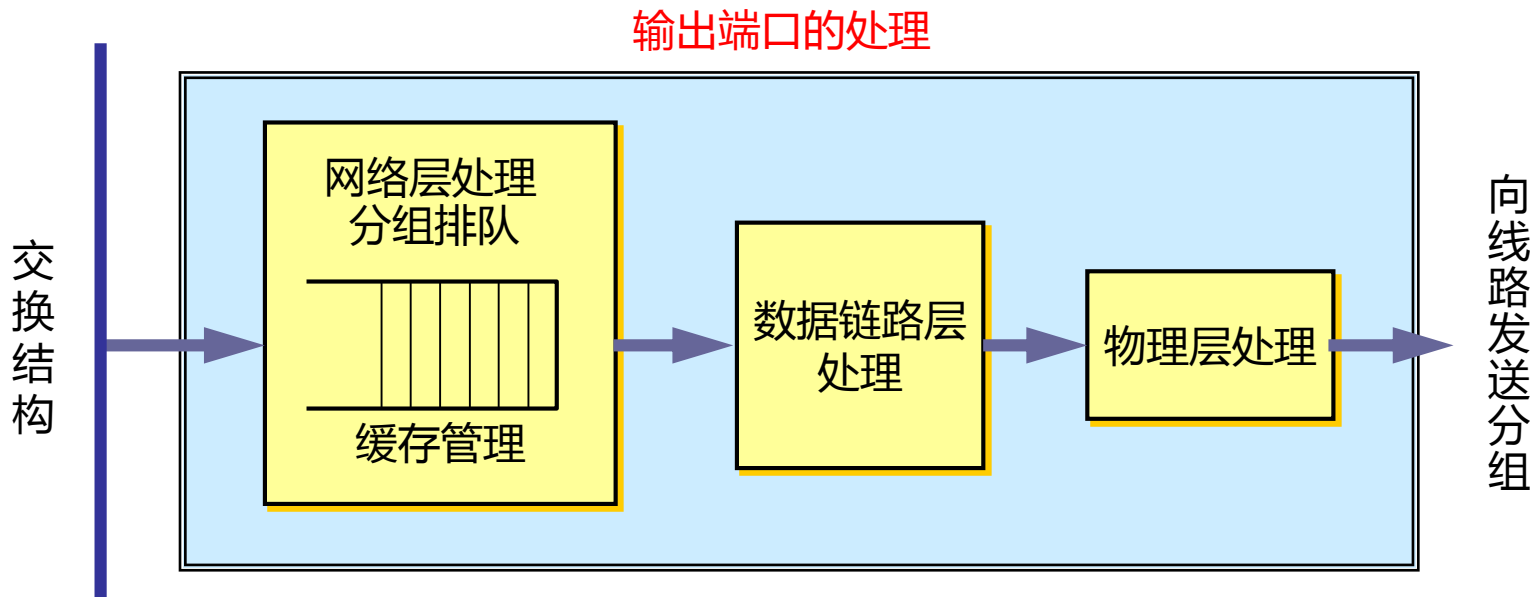
输入端口的处理



5.路由选择协议

5.5路由器的结构

- 输出端口将交换结构传送来的分组发送到线路：
 - 当交换结构传送过来的分组先进行缓存。数据链路层处理模块将分组加上链路层的首部和尾部，交给物理层后发送到外部线路。



5.路由选择协议

5.5路由器的结构

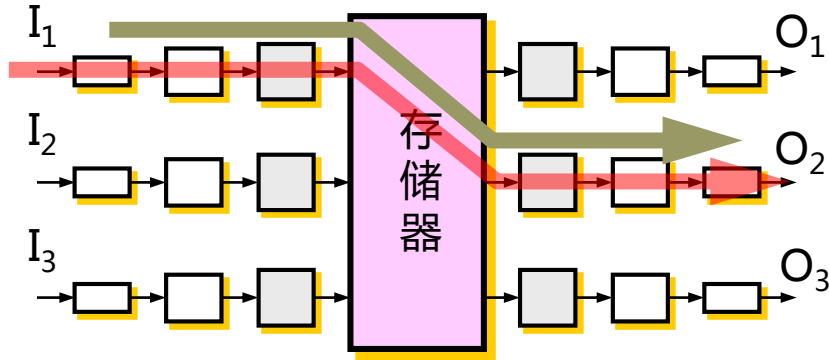
□ 分组丢弃：

- 若路由器处理分组的速率赶不上分组进入队列的速率，则队列的存储空间最终必定减少到零，这就使后面再进入队列的分组由于没有存储空间而只能被丢弃。
- 路由器中的输入或输出队列产生溢出是造成分组丢失的重要原因。

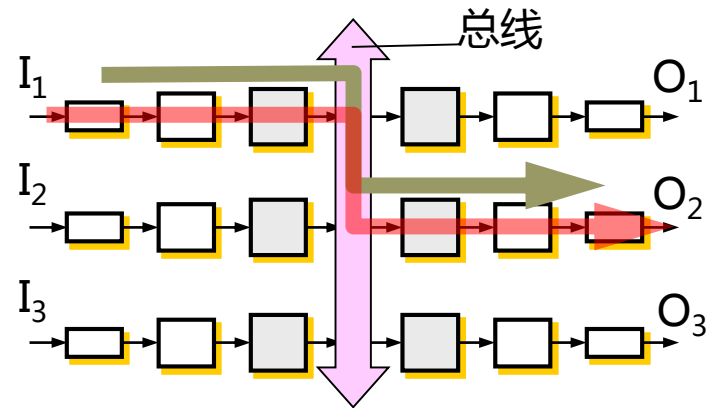
5. 路由选择协议

5.5 路由器的结构

□ 交换结构：



(a) 通过存储器

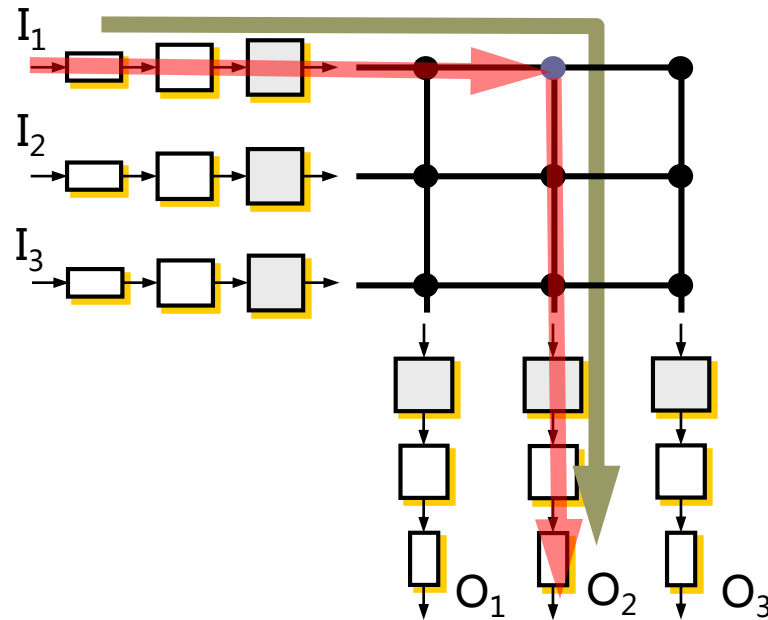


(b) 通过总线

5.路由选择协议

5.5路由器的结构

□ 交换结构：



互连网络

(c) 通过互连网络

5.路由选择协议

5.5路由器的结构



NE40E-X8



NE40E-X16



NE40E-X3

5.路由选择协议

5.5路由器的结构



AR2220



AR2240

5.路由选择协议

5.5路由器的结构



Thanks