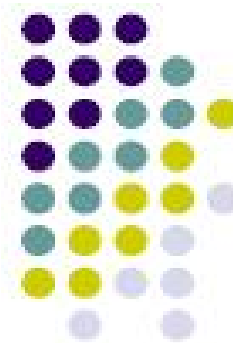


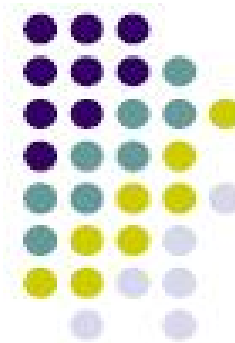
第八章 通信网络与服务

- TCP/IP 体系结构
- 网际协议
- 传输层协议
- 互联网路由协议
- DHCP与NAT
- IPv6



第八章 通信网络与服务

TCP/IP 体系结构



网络化的必要性



- 建立网络间联系

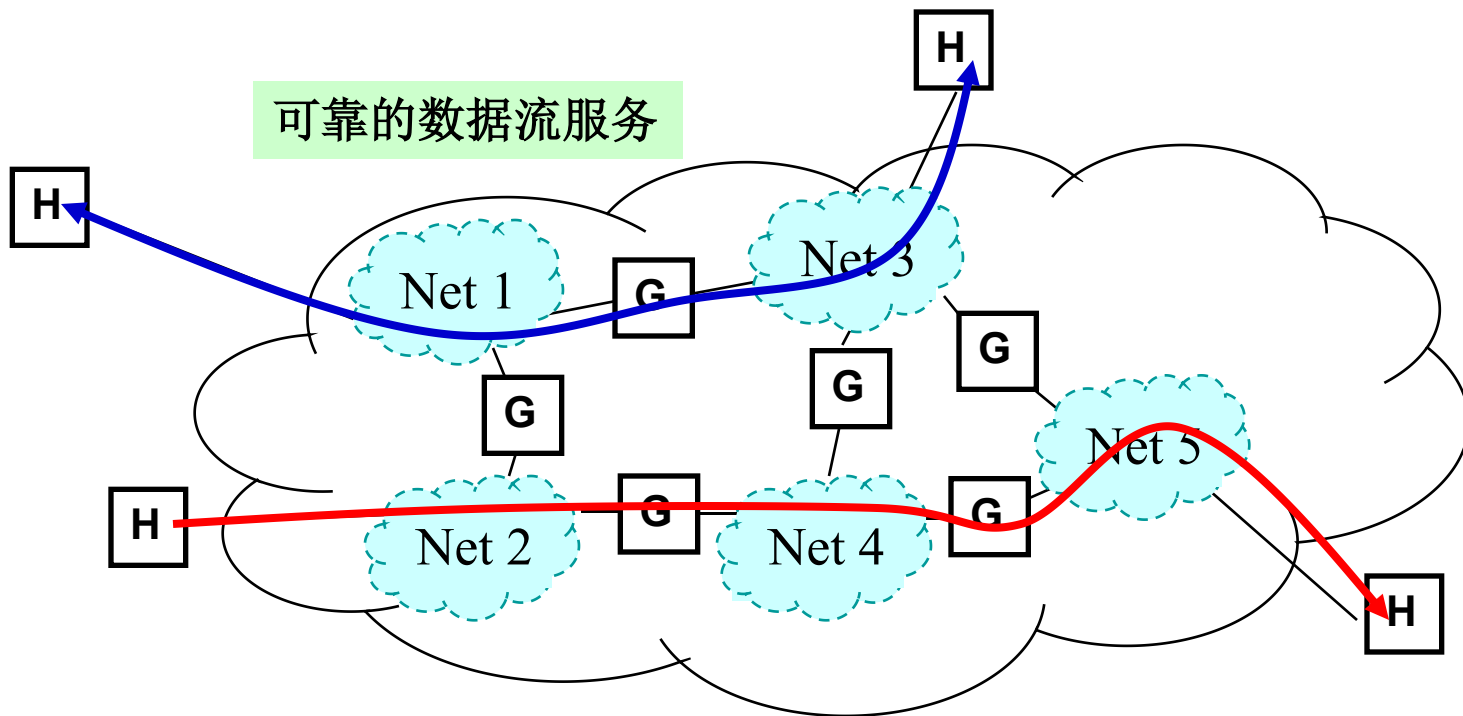
- 不同网络技术并行
- 通过IP数据包通信

- 提供通信服务

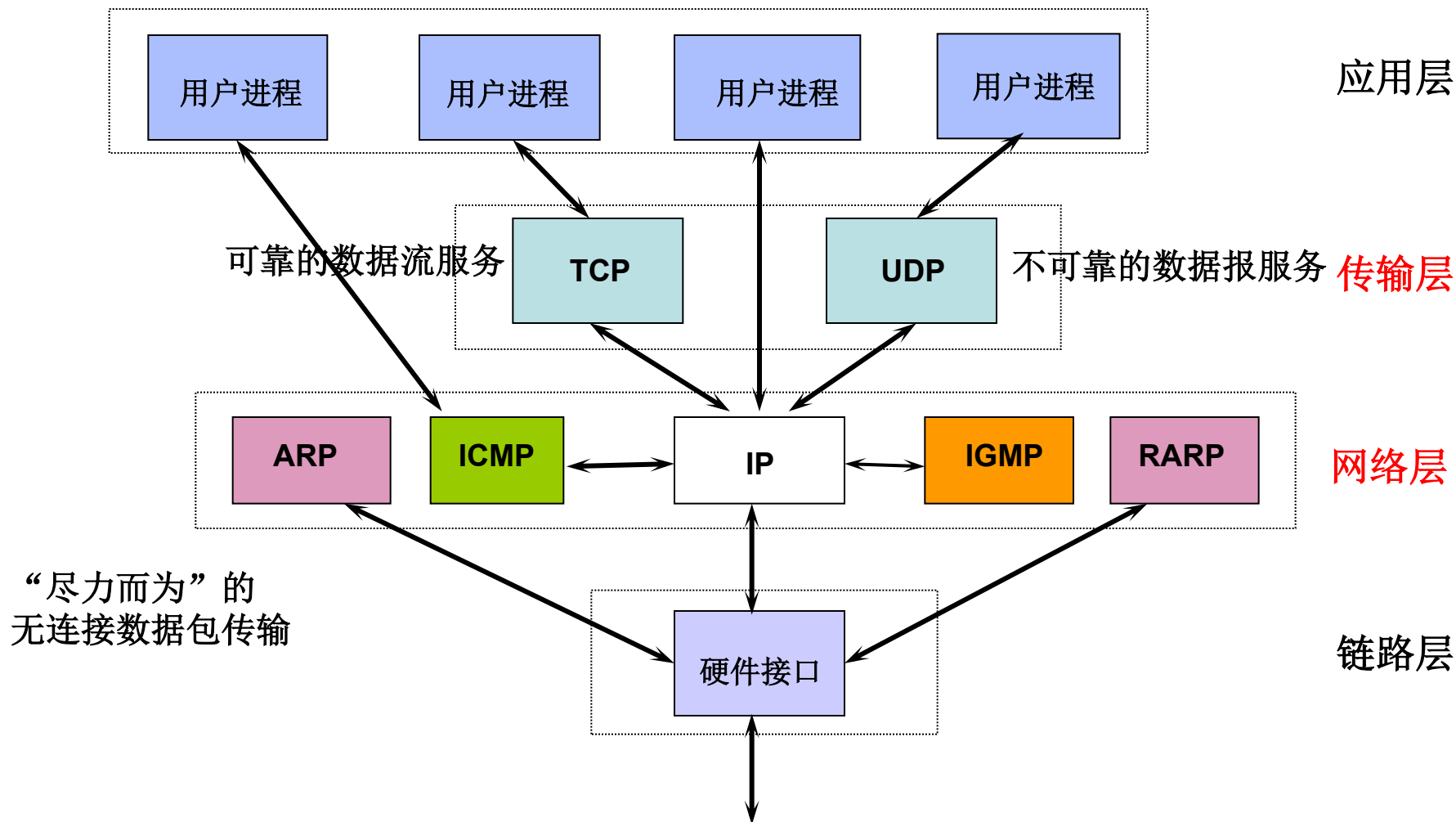
- 为用户提供通用接口

- 支持多样化应用服务

- 支持运行任何基于互联网通信的应用
- 便于新应用的网络部署



网际协议结构



封装技术



TCP报头:

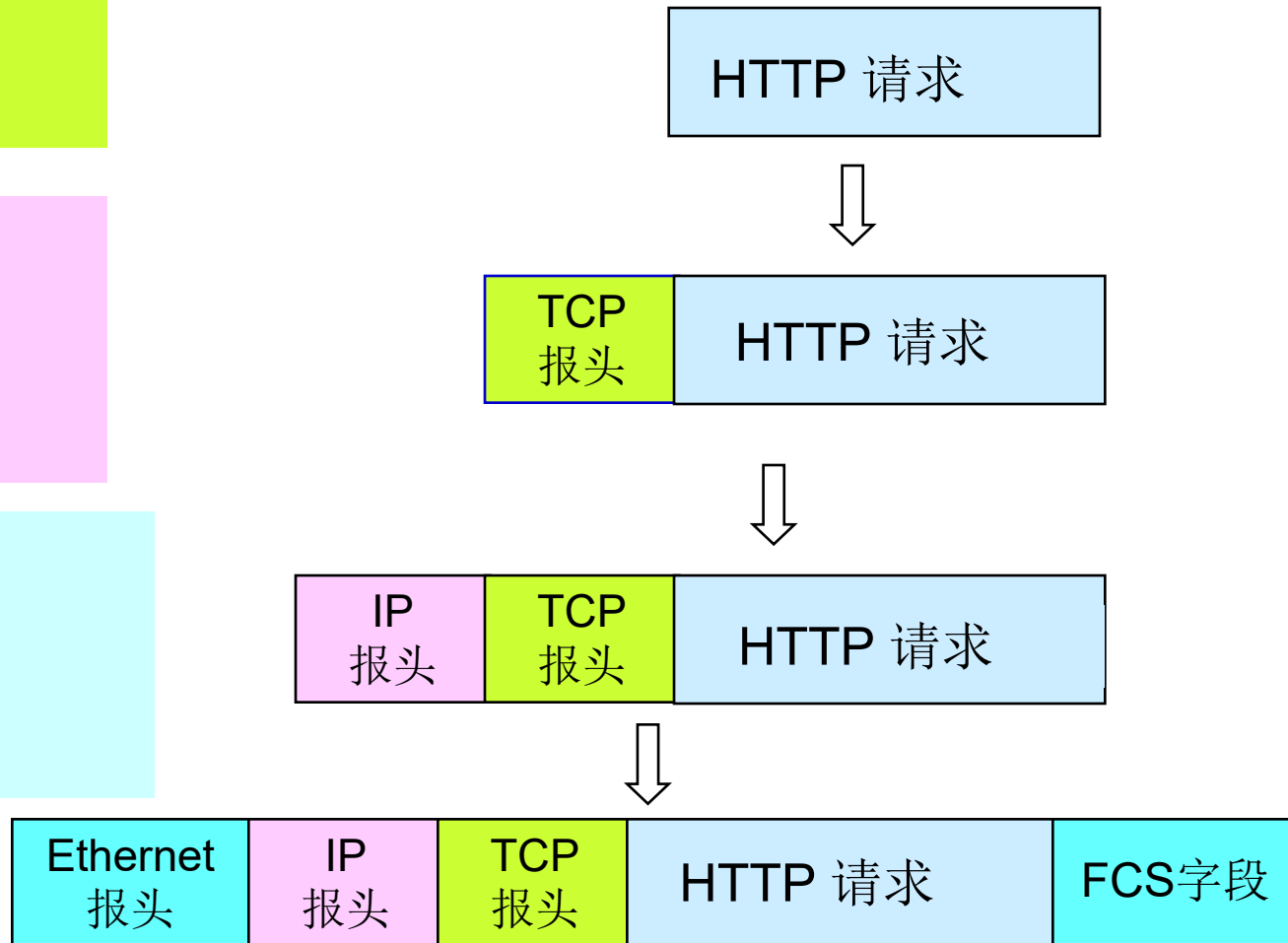
- 源、目的端口号

IP报头:

- 源、目的IP地址
- 传输协议类型

Ethernet报头:

- 源、目的MAC地址
- 网络协议类型



网络层



监督主机到主机的数据包传输（主机可能被几个物理网络分开）

- 数据链路层提供节点到节点的传输
- 传输层提供进程对进程的传输

网络层的主要任务：

- **寻址：** 唯一地识别每个设备，以允许全球通信
- **路由：** 确定从一个主机向另一个主机发送数据包的最佳路线
- **分组：** 封装从上层协议收到的数据包
- **分片：** 与每个局域网的最大传输单元（MTU）相匹配

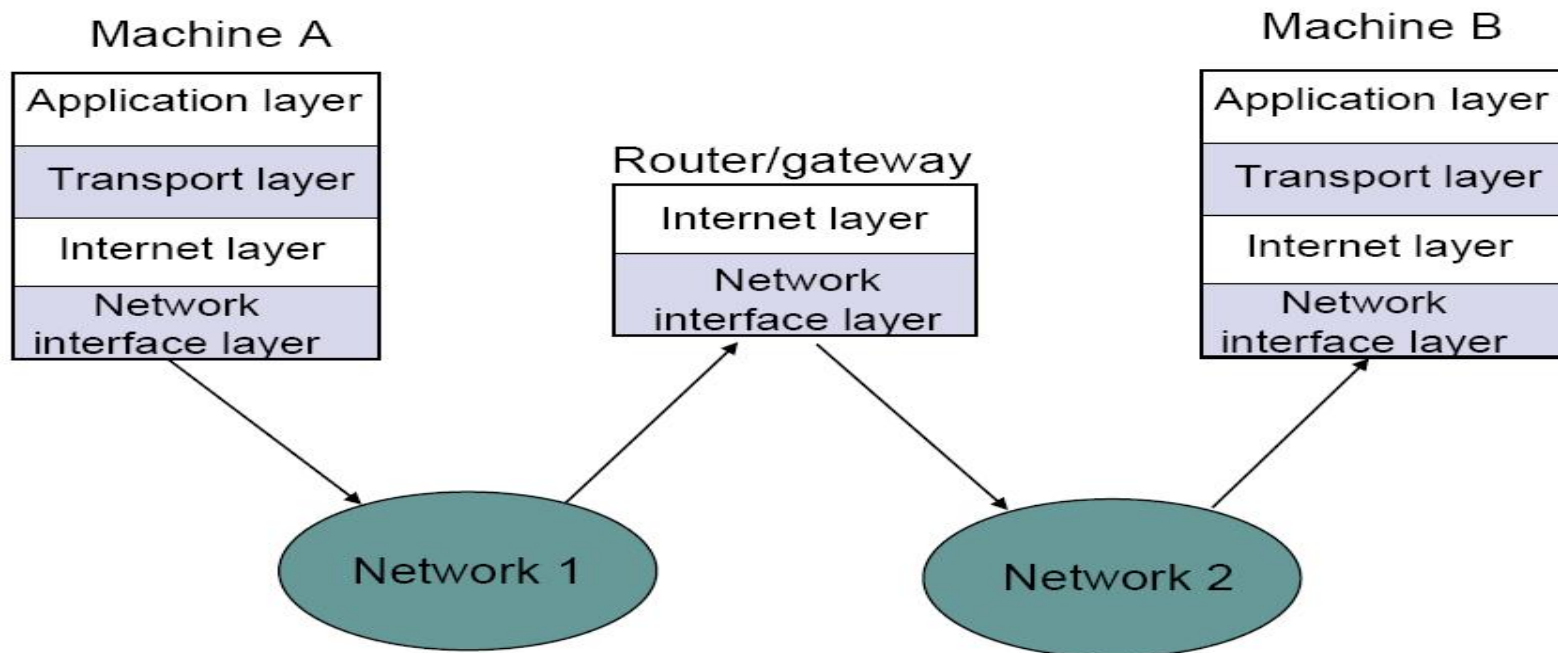
网络层



- IP数据包在互联网上传输信息

主机A的IP（源）→路由器→路由器...→路由器→主机B的IP（目的）

- 每个路由器的IP层决定下一跳（路由器）
- 网络接口在网络上传输IP数据包



网络层

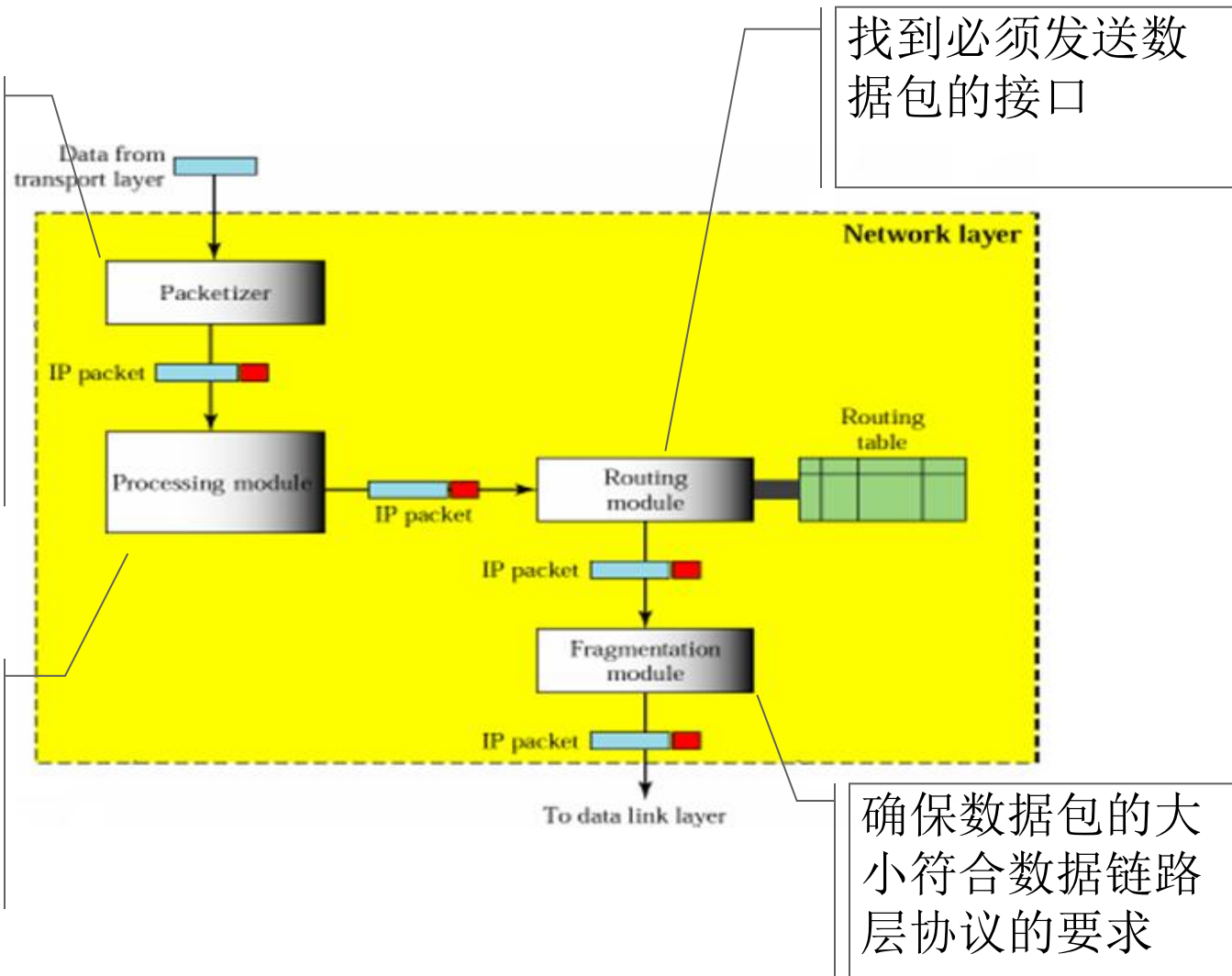


例子：源节点处互联网中的网络层职责

对来自上层的数据包进行封装，即添加报头。

- 1) 添加通用的源地址和目的地址。
- 2) 添加错误控制的字段，等等。

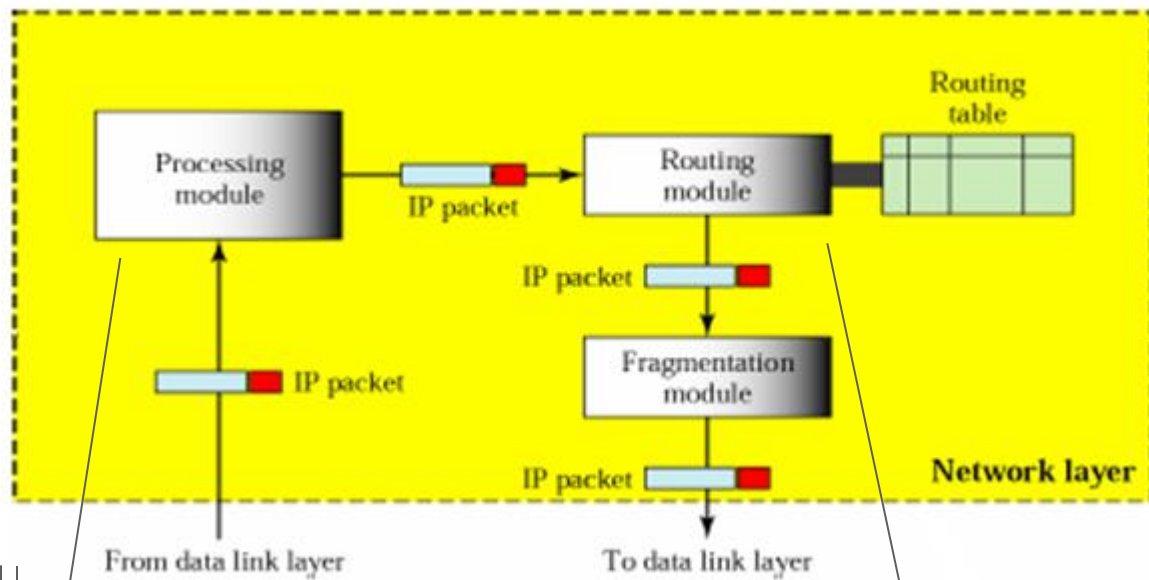
验证是否目的地地址是主机地址！
如果是，则不需要路由！



网络层



例子：中间路由节点处互联网中的网络层职责



检查数据包是否
已到达其最终目
的地或需要转发

TTL!+ 头部错误
检查!!!

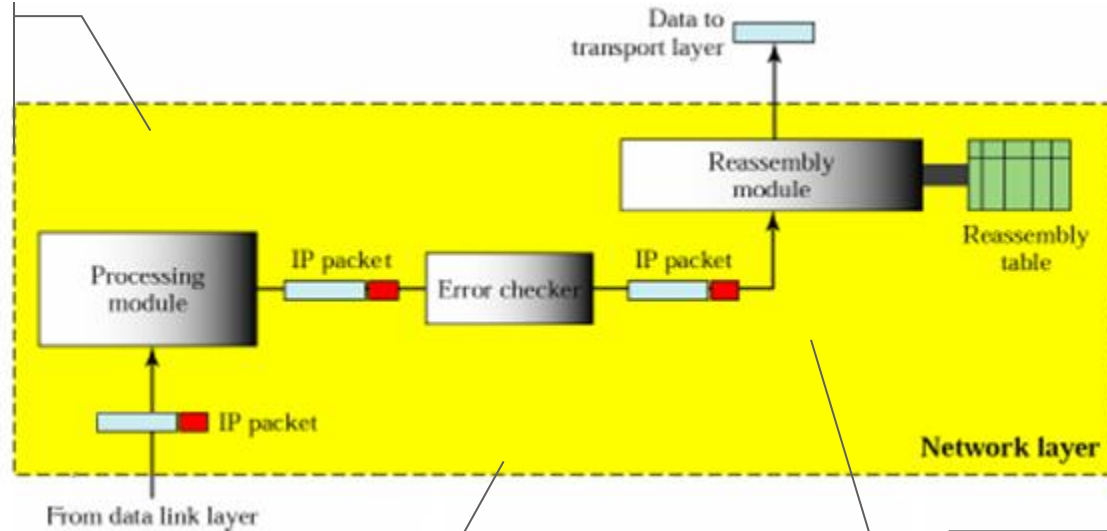
查找接口
数据包从正确
接口发送

网络层



例子：目的节点处互联网中的网络层职责

验证目的地址是否是主机地址



检查数据包在传输过程中是否有被破坏

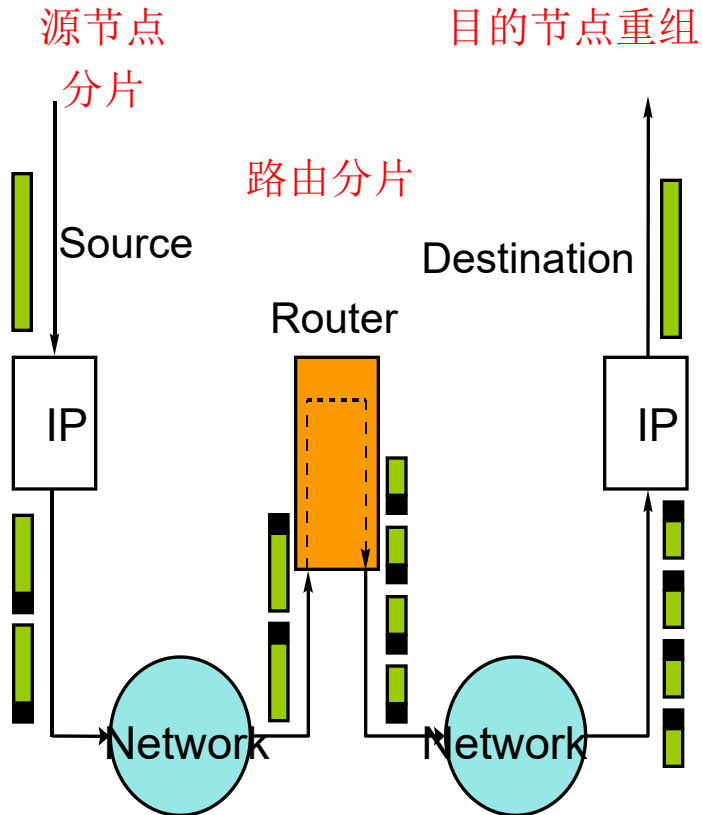
如果数据包已被分片，等待直到所有片段到达，重新组合然后发送数据包到传输层

网络层



分片和重组

—将数据报分割成更小的片段，以满足底层数据链路层协议的**MTU**要求



- **分片**：源主机或路径中的任何其他路由器
- **重组**：只由目的主机进行，为什么？
- 分片后的数据报可能会被进一步分割，为什么？
- 分片数据报的主机或路由器必须改变三个字段的值：
 - 标记
 - 分片偏移量
 - 总长度

网络层



标识字段： 16位字段 - 唯一识别来自源主机的数据报

- 为了保证唯一性，IP使用计数器来标记每个数据报。
- 当数据报被分割时，识别字段被复制到所有片段中。
- 识别号码有助于目的地重新组装数据报。

标记字段： 3位字段

- 第1位是保留位
- 第2位被称为 "不分片 "位

如果值为1，则机器不得将数据报分片

如果片段不能通过物理网络，路由器会丢弃数据包，并向源主机发送**ICMP**错误信息

- 第3位被称为 "更多片段 "位

如果值是1，数据报不是最后一个片段

如果值为0，这是最后一个或唯一的片段

D: Do not fragment

M: More fragments





片偏移字段：13位字段

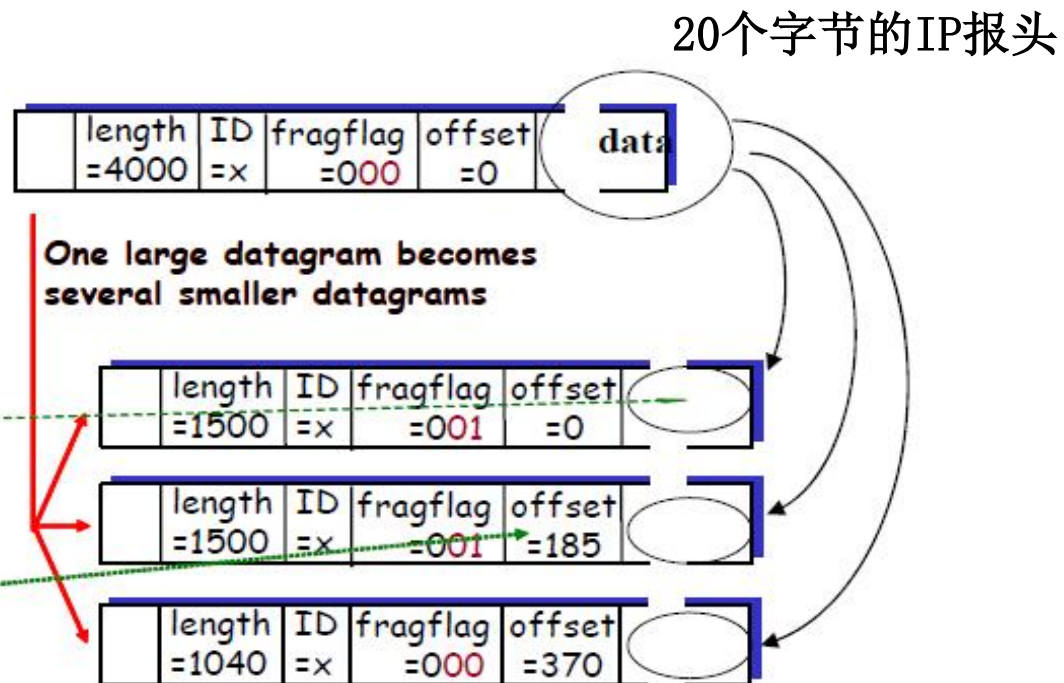
- 识别数据包中的片段位置
- 以8字节为单位测量：只有13位长，否则不能代表大于8191的序列片段
→ 大小应能被8整除

Example

- 4000 byte datagram
- MTU = 1500 bytes

1480 bytes in
data field

offset =
 $1480/8$





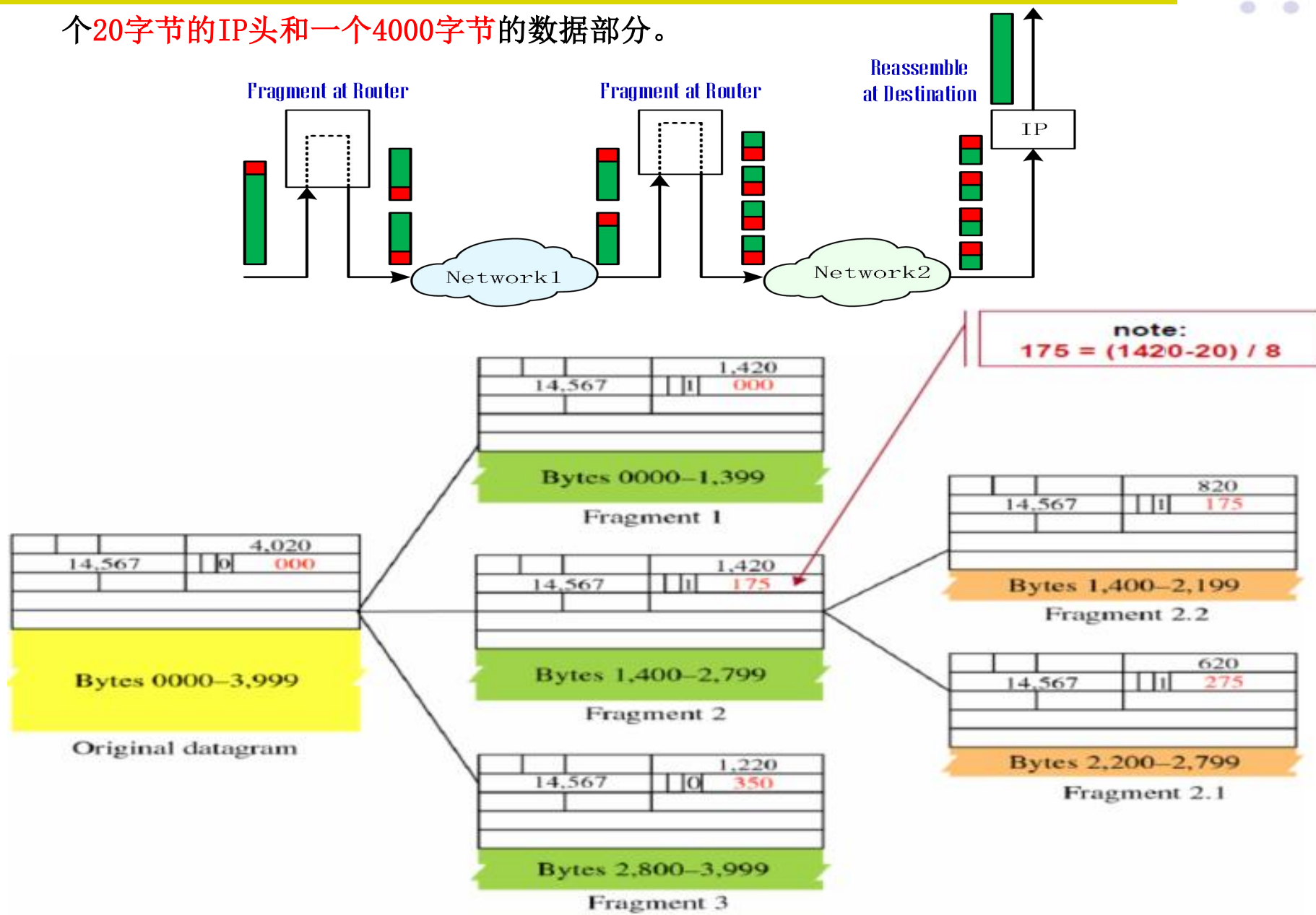
分片实例：

- 一个数据包将被转发到一个MTU为576字节的网络。该数据包有一个20字节的IP头和一个1484字节的数据部分。
- 每个片段的最大数据长度=576-20=556字节。
- 设置最大数据长度为552字节，以获得8的倍数。

	Total Length	Id	MF	Fragment Offset
Original packet	1504	x	0	0
Fragment 1	572	x	1	0
Fragment 2	572	x	1	69 = 552/8
Fragment 3	400	x	0	138=552*2/8

二次分片实例：

- 一个数据包将被转发到MTU为1420字节的网络1和MTU为820字节的网络2。该数据包有一个20字节的IP头和一个4000字节的数据部分。





二次分片实例:

网络1中的分片, MTU为1420字节

	Total Length	Id	MF	Fragment Offset
Orig. packet	4020	x	0	0
Frag. 1	1420	x	1	0
Frag. 2	1420	x	1	175 = 1400/8
Frag. 3	1220	x	0	350=1400*2/8

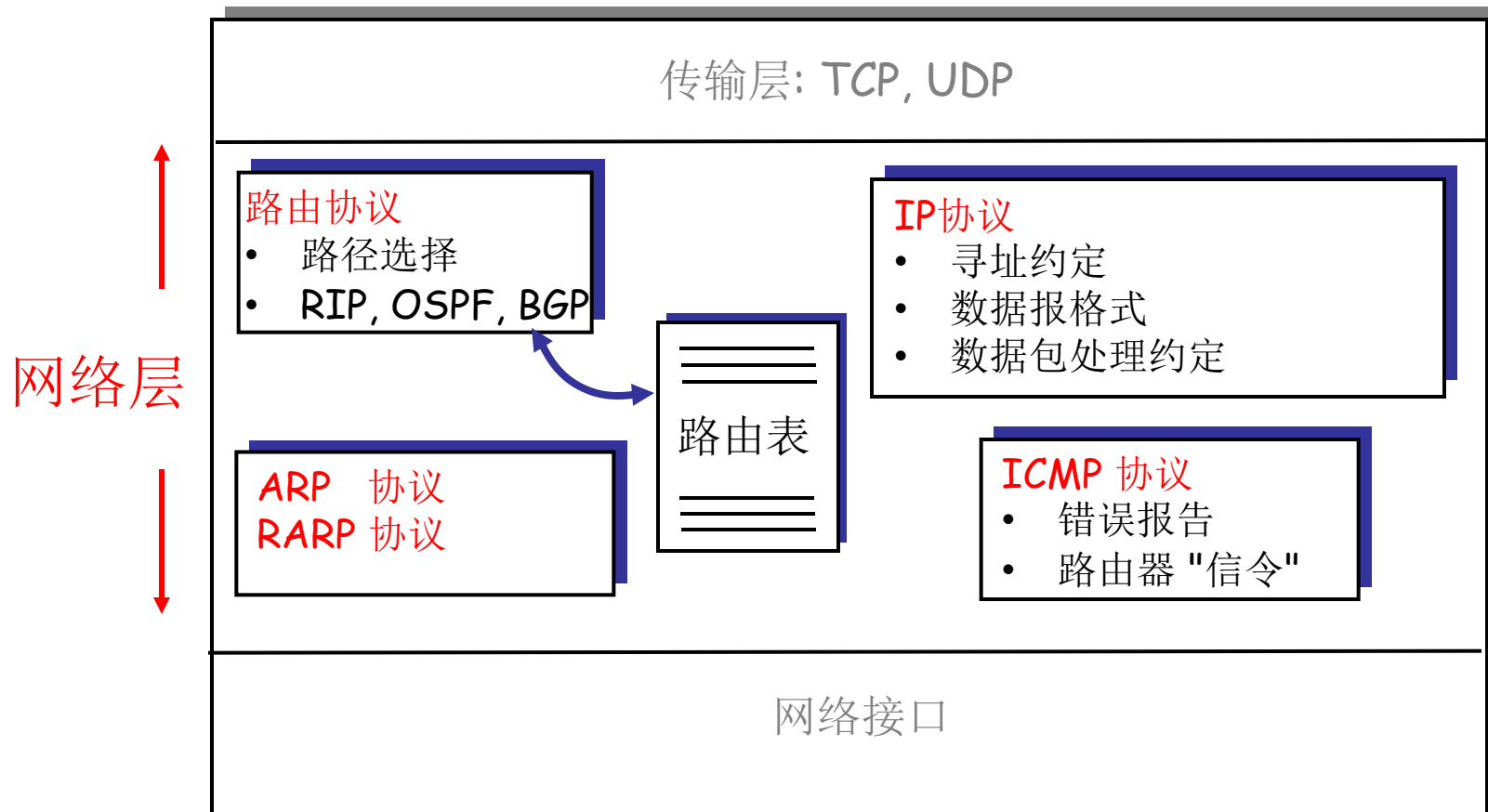
网络2中的分片, MTU为820字节

	Total Length	Id	MF	Fragment Offset
Orig. packet	4020	x	0	0
Frag. 1-1	820	x	1	0
Frag. 1-2	620	x	1	100=800/8
Frag. 2-1	820	x	1	175 = 1400/8
Frag. 2-2	620	x	1	275=2200/8
Frag. 3-1	820	x	1	350=1400*2/8
Frag. 3-2	420	x	0	450=(1400*2+800)/8



网络层

主机、路由器网络层功能。



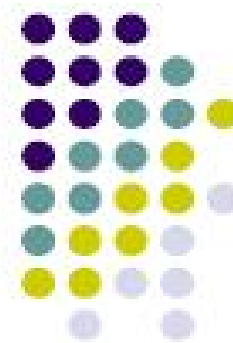
网络层



• IP	<ul style="list-style-type: none">• 主要协议• 负责 "尽最大努力 "的主机到主机的传输
ARP	<ul style="list-style-type: none">• 将下一跳的IP地址映射为其MAC/物理地址• 在将数据包传递到较低的数据链路层时使用
RARP	<ul style="list-style-type: none">• 将MAC/物理地址映射到IP地址• 在无盘机器上使用，用于恢复IP地址
ICMP	<ul style="list-style-type: none">• 被主机和路由器用来处理异常情况，如IP包头错误、无法到达的主机和网络等。
IGMP	<ul style="list-style-type: none">• 主机和路由器用于实现高效的网络层组播
Routing Protocols	<ul style="list-style-type: none">• 负责路由表的维护

第八章 通信网络与服务

网际协议
- IP 包结构



网际协议(IP)



- 提供尽力、无连接的数据包传输
 - 保持路由器的简单性，并适应网络元素的故障
 - 数据包可能丢失、失序，甚至重复。
 - 必要时，高层协议必须处理这些问题
- RFCs 791, 950, 919, 922, and 2474.
- IP是互联网STD编号5的一部分，其中还包括：
 - 互联网控制消息协议（ICMP），RFC 792
 - 互联网组管理协议（IGMP），RFC 1112

注：STD:5- 定义了互联网协议版本4。[RFC 791] 互联网协议。更新的方式。

[RFC 1349] 互联网协议套件中的服务类型

[RFC 795] 服务映射

[RFC 796] 地址映射

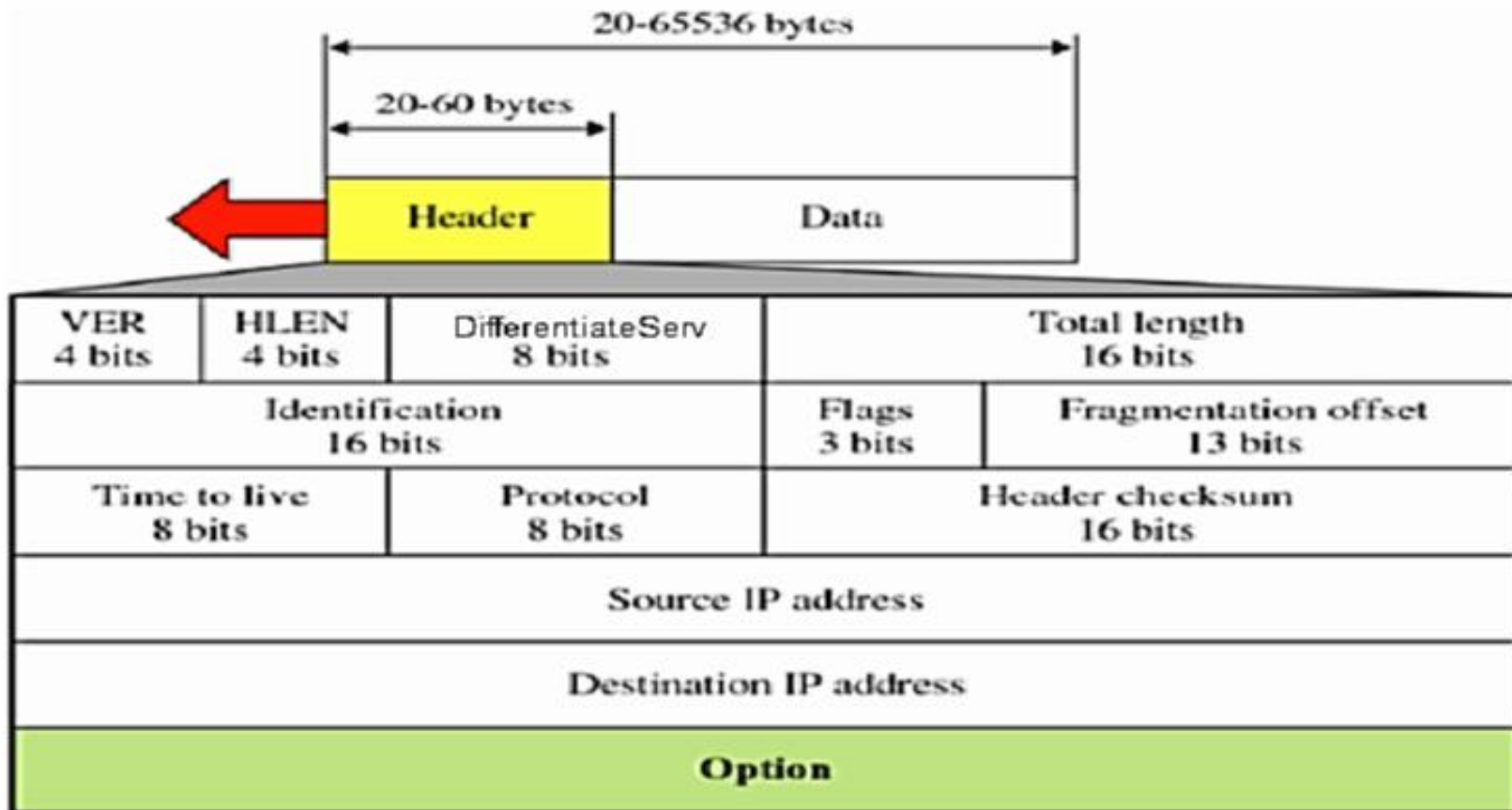
[RFC 932] [RFC 936] [RFC 940] [RFC 950] ~。

IP 报头



IP数据包=由头和数据组成的可变长度数据包

- 头部 - 长度为20至60字节，包含对路由和传输至关重要的信息
- 数据 - 长度由链路层协议的最大传输单元（MTU）决定（理论上在20至65536字节之间）



版本 - 4位字段 - 指定数据报的IP协议版本（IPv4或IPv6）。

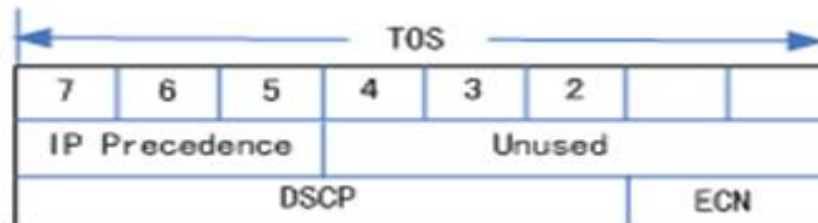
互联网报头长度（IHL）--4位字段--报头的长度（单位4字节）。

IP 报头



Differentiated Services (8bit)

替换头字段 DS字段，取代现有IPv4 TOS 和 IPv6 Traffic Class



Differentiated Services Code Point (DSCP) – 6 bit field – 区分服务代码点，有**数字**和**关键字**两种形式。

- 数字：DSCP使用6比特，十进制区间是0-63，可以定义64个优先级；
- 关键字：逐跳行为（PHB），分别是尽力服务（BE或DSCP 0）、确保转发（AF_{xy}）和加速转发（EF）。

Explicit Congestion Notification (ECN) – 2 bits -显式拥塞通知

- 00表示该数据包不使用ECN。
- 01表示该数据包是一个具有ECN功能的传输流的一部分。
- 10表示该数据包是一个实验性的ECN-capable传输流的一部分。
- 11表示该数据包经历了拥堵。

IP 报头



TOS, Type of Service - 8 bits

- 指定要求的服务类型的参数。
- 延迟/吞吐量/可靠性/货币成本

00	01	02	03	04	05	06	07
Precedence			D	T	R	M	0

优先级6和7一般保留给网络控制数据使用，如路由。

优先级5推荐给语音数据使用。

优先级4由视频会议和视频流使用。

优先级3给语音控制数据使用。

优先级1和2给数据业务使用。

优先级0为默认标记值。

优先级Precedence. 3 bits.

Value	Description
0	Routine. (普通)
1	Priority. (优先)
2	Immediate. (快速)
3	Flash. (闪速)
4	Flash override. (疾速)
5	CRITIC/ECP. (关键)
6	Internetwork control. (网间控制)
7	Network control. (网络控制)

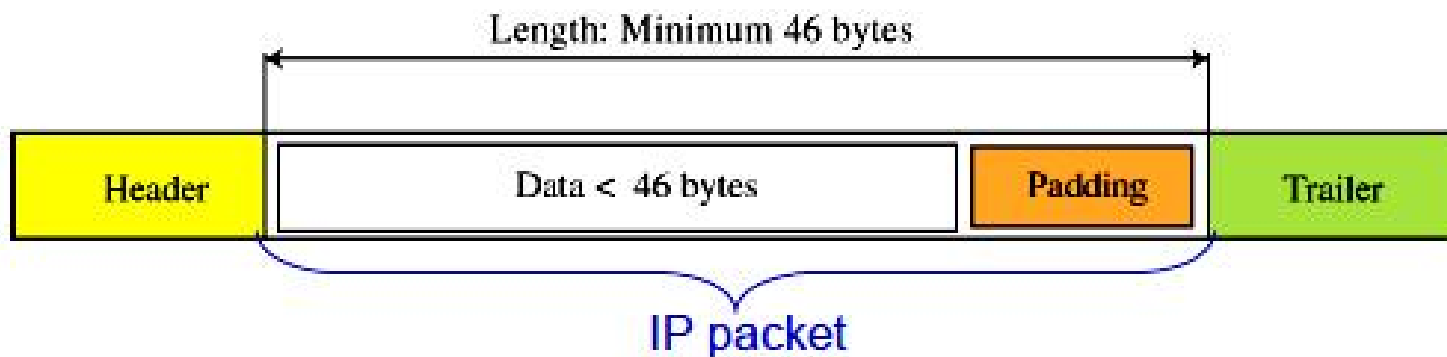
IP 报头



总长(16-bit)

--以字节为单位的数据报总长度，包括报头

- 16位 \Rightarrow 最大尺寸 = 65,535 字节
- 大数据报必须被分割成碎片，以便能够通过有MTU限制的网络
- 小数据报在以太网协议中必须被填充，以便能够通过对数据最小尺寸有限制的网络





标识、标记和片偏移字段

分片中使用的3个字段（16+3+13位）

- IPv6不允许在路由器上进行分片，因为操作很耗时
→ 如果一个IPv6数据包太大，它就会被简单地丢弃，然后一个ICMP消息被发回给源节点

Time-to-live (TTL) (8-bit)

控制数据报访问的最大跳数和/或在网络中停留的时间

- 每次路由器处理数据报时，该字段都会减去一个，当TTL达到0时，数据报被丢弃
- 目的
 - 1) 数据报不会永远循环传输
 - 2) 限制其行程（例如，仅局域网：TTL=1）

IP 报头



协议 (8-bit)

- 表示具体的传输层协议，该IP数据报的数据部分应被传递给该协议
 - 只在最终目的地使用，以便于解复用过程
 - 值：1 - ICMP, 2 - IGMP, 6 - TCP, 17 - UDP, 89 - OSPF

源和目的IP地址 - (例如，IPv4 32位)

- 在IP数据报到达最终目的地之前保持不变

填充

- 用来使报头成为32位的倍数

IP 报头



头部校验和 (16-bit)

- 仅检测头部中的错误!
- 必须重新计算并存储在每个路由器上, 因为某些字段 (例如 TTL) 可能会改变
- 路由器丢弃检测到错误的数据报
- 校验和计算:
 - 1) 将报头分成 16 位 (2 字节) 部分——校验和字段本身设置为 0
 - 2) 使用 1s 补码算法对所有部分求和

每个中间路由器必须:

- 1) 验证/重新计算每个传入数据包的校验和
- 2) 计算每个传出数据包的校验和

4	5	0	28	
1		0	0	
4	17	0		
10.12.14.5				
12.6.7.9				
4, 5, and 0	→ 0100010100000000			
28	→ 00000000000011100			
1	→ 00000000000000001			
0 and 0	→ 00000000000000000			
4 and 17	→ 0000010000010001			
0	→ 00000000000000000			
10.12	→ 0000101000001100			
14.5	→ 0000111000000101			
12.6	→ 0000110000000110			
7.9	→ 0000011100001001			
<hr/>				
Sum	→ 0111010001001110			
Checksum	→ 1000101110110001			

错误检测/纠正不是网络层的责任。
那为什么要对 IP 报头执行错误检测? !

IP 报头

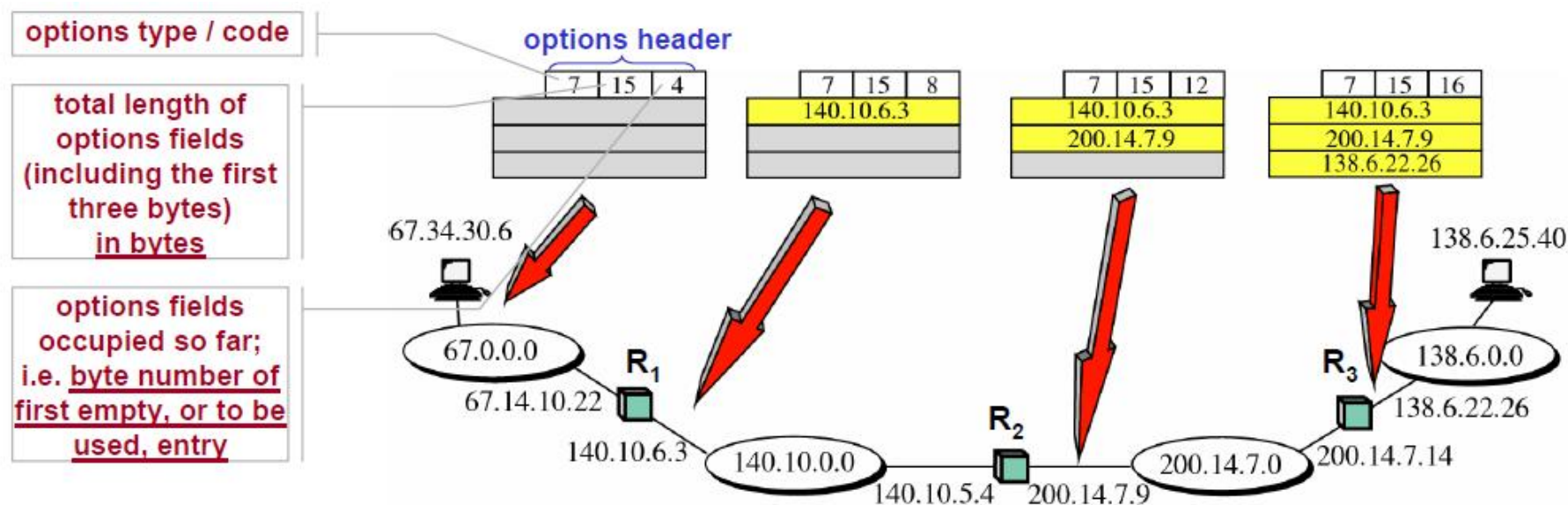


可选项 (32 位的整数倍)

- 不是每个数据报都需要！
- 允许扩展 IP 报头用于特殊目的

(a) 记录路由选项 —— 用于跟踪数据报采用的路由

- 源为 IP 地址创建空字段 - 最多 9 个 4 字节 IP 地址 (40 字节选项 - 4 字节选项标头)
- 每个处理数据报的路由器插入其 **传出** IP 地址



IP 报头



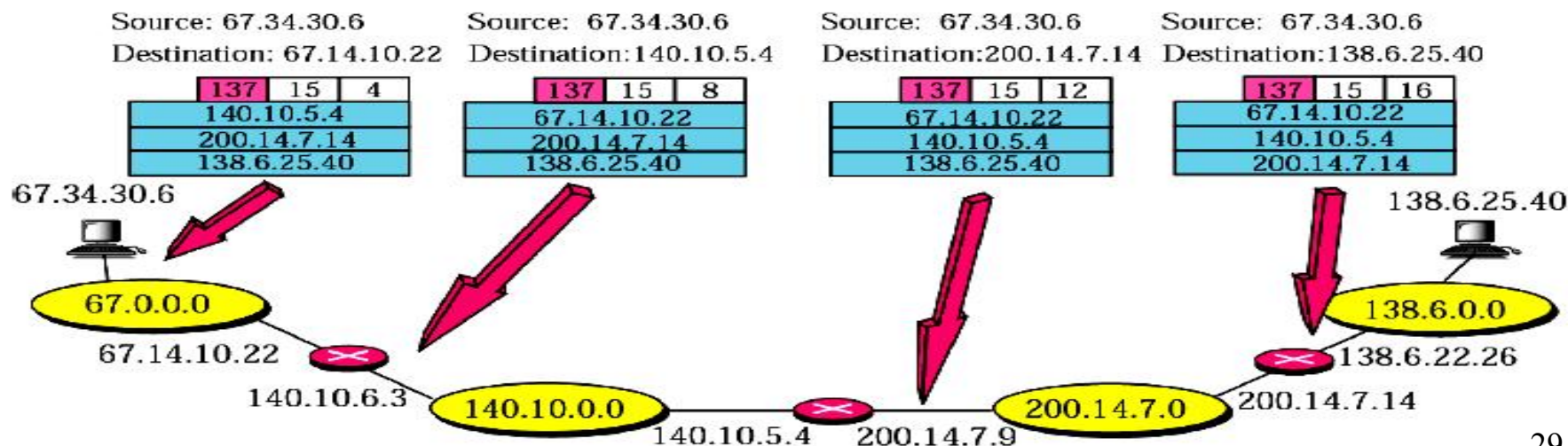
可选项

(b) 时间戳选项—类似于 (a)，加上记录每个路由器的数据报结束处理时间，以毫秒为单位

(c) 严格源路由选项—源节点用于预先确定数据报的路由

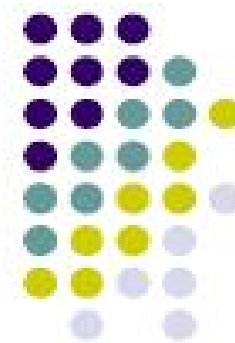
源节点提供一个IP地址列表（路由器序列），数据报在前往目的地的途中必须访问该地址。

(d) 松散源路由选项—必须访问列表中的每个路由器，尽管数据报也可以访问其他路由器



第八章 通信网络与服务

网际协议 - IPV4 地址



互联网名称和地址



互联网名称

- 每个主机都有一个独特的名字
 - 与物理位置无关
 - 便于人类记忆
 - 域名
 - 在单一行政单位下的组织
- 主机名称
 - 给予主机的名称
- 用户名称
 - 分配给用户的名称

互联网地址

- 每个主机都有全球唯一的逻辑32位IP地址
- 每个与网络的物理连接都有单独的地址
- IP地址有两部分。 *netid*、 *hostid*
 - *netid*是唯一的
 - *netid*有利于路由选择
- 点状十进制记法。
 - int1.int2.int3.int4*
 - 141.223.1.2

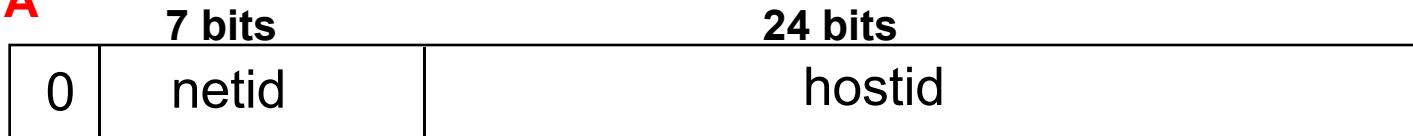
DNS将域名解析为IP地址

IP 地址



- RFC 1166
- 互联网上的每个主机都有唯一的32位IP地址
- 每个地址有两个部分：netid和hostid
- netid是唯一的，由以下机构管理
 - American Registry for Internet Numbers (ARIN)
 - Reseaux IP Europeens (RIPE)
 - Asia Pacific Network Information Centre (APNIC)
- IP地址分配给网络接口而不是主机
 - ⇒主机与网络的每个物理连接都需要单独的地址；例如，“多宿主”主机
- 点分十进制符号
 - IP 地址 10000000 10000111 01000100 00000101
 - 点分十进制 128.135.68.5

Class A



- 126个网络——多达1600万主机 **1.0.0.0 to 127.255.255.255**

Class B



- 16,382个网络——多达64,000个主机 **128.0.0.0 to 191.255.255.255**

Class C



- 200万个网络——多达254个主机 **192.0.0.0 to 223.255.255.255**

IP地址分类



Class D

28 bits

1	1	1	0	multicast address
---	---	---	---	-------------------

**224.0.0.0 to
239.255.255.255**

- 组播地址
- 多达2.5亿个组播群同时进行

Class E

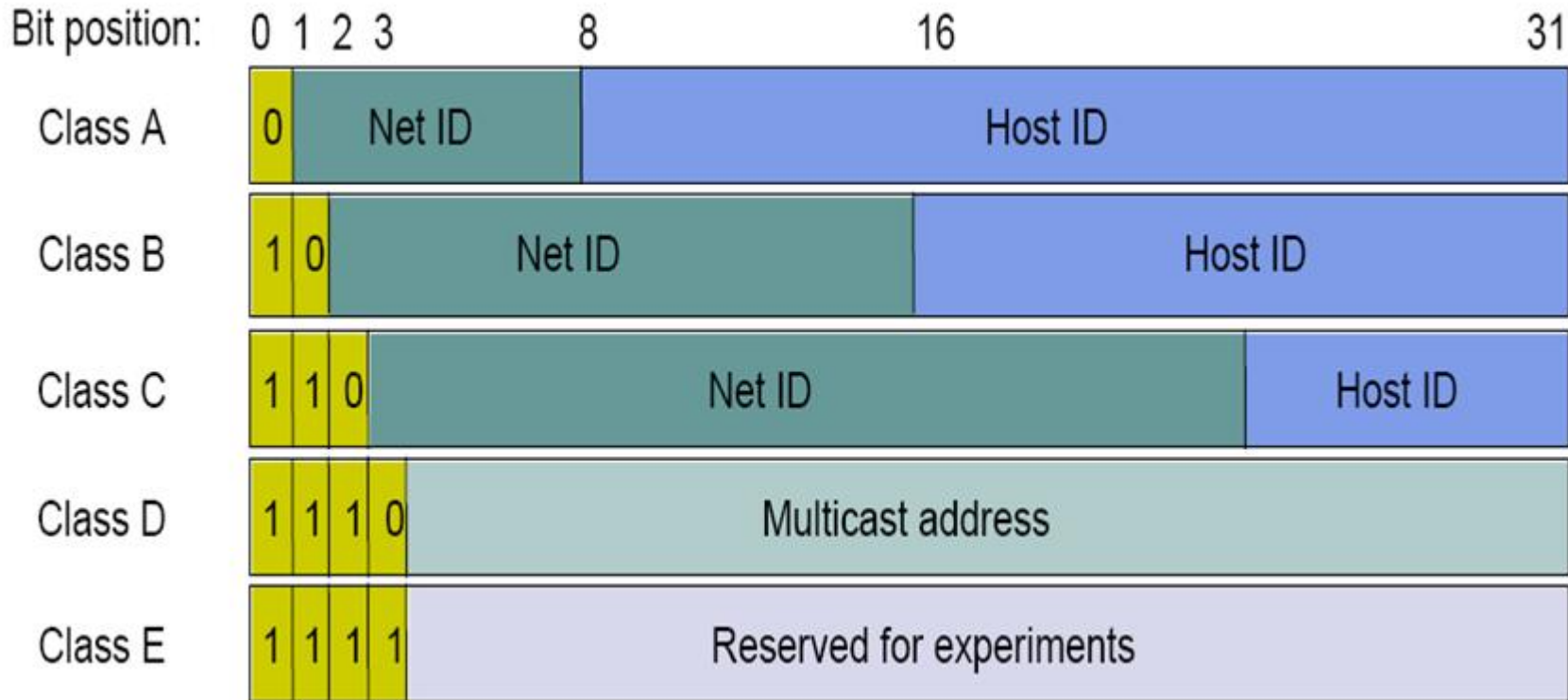
28 bits

1	1	1	1	multicast address
---	---	---	---	-------------------

240.0.0.0~~

- 留作实验用

IP地址分类



A/B/C IP 地址范围



类别	最大网络数量	第一个可用的网络ID	最后可用的网络ID	每个网络中的最大主机数
A	126 ($2^7 - 2$)	1	126	16,777,214
B	16,384 (2^{14})	128.0	191.255	65,534
C	2,097,152 (2^{21})	192.0.0	223.255.255	254

0.0.0.0是保留的，意味着 "未知地址"

特殊IP地址



- 主机ID “全0”是保留的，指网络编号
166.111.0.0 (掩码255.255.0.0)
202.119.22.0 (掩码255.255.255.0)
- 主机ID “全1”被保留用来向特定网络上的所有主机广播
166.111.255.255 (掩码255.255.0.0)
202.119.22.255 (掩码255.255.255.0)



特殊IP地址

— 保留地址

- **0.0.0.0**是保留的，意味着 "未知地址"。通常用于启动无盘工作站
- **255.255.255.255**是保留地址，用于向本地网络中的每台主机广播。
- **127.x.x.x** 表示 "这个节点"（本地环回）。
- 发送到这个地址的信息将永远不会离开本地主机。它的目的是测试网络软件
- **24.x.x.x** 有线电视网络配置
- **169.254.X.X** 链接本地。

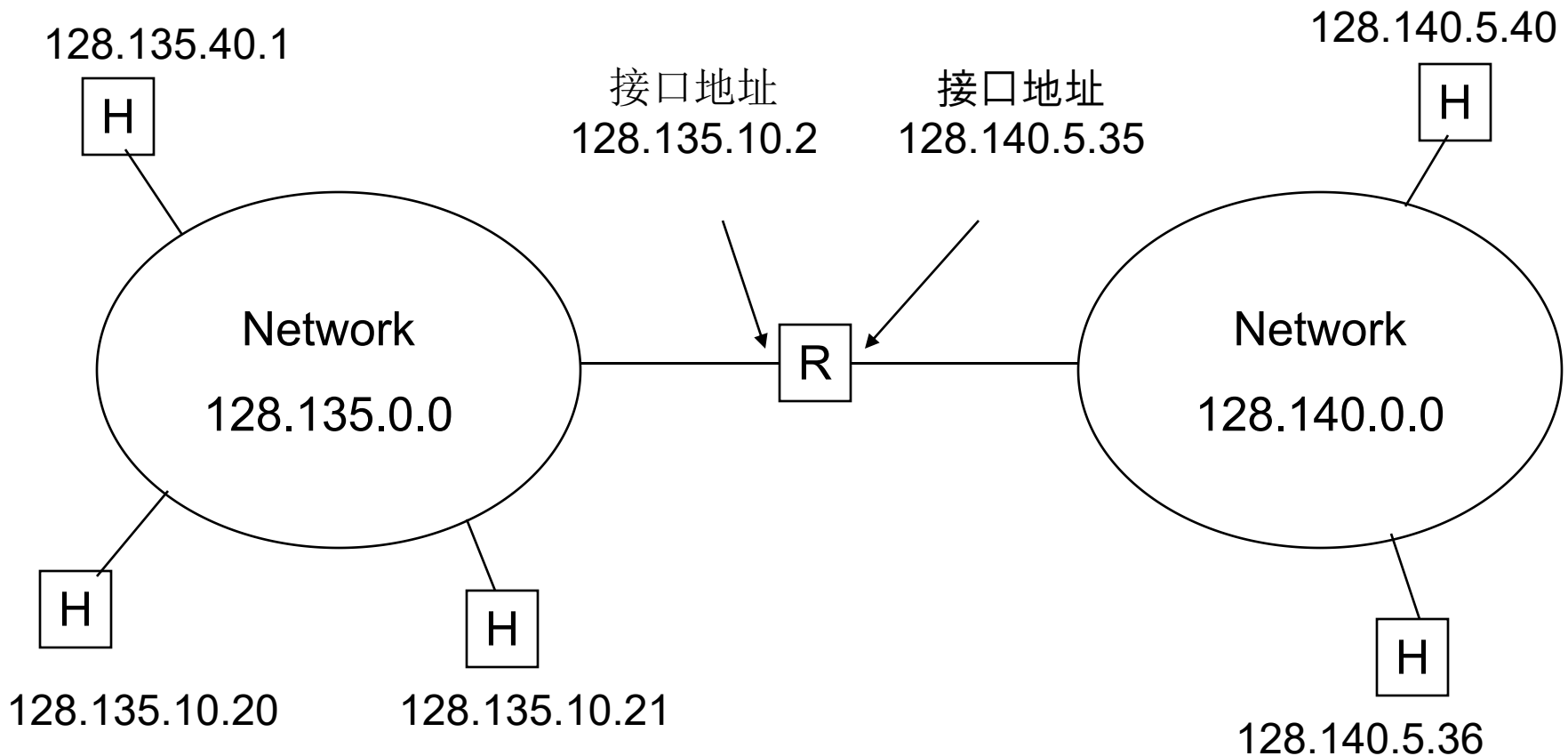
它被分配用于单一链路上的主机之间的通信。主机通过自动配置获得这些地址，例如在可能找不到DHCP服务器的时候
- **191.255.x.x**仍由IANA保留。



特殊IP地址：私人IP地址

- 私人网络使用的特定范围的IP地址（RFC 1918）
- 只限于在私人网络中使用；公共互联网的路由器会丢弃带有这些地址的数据包
 - 范围1：10.0.0.0至10.255.255.255
 - 范围2：172.16.0.0至172.31.255.255
 - 范围3：192.168.0.0至192.168.255.255
- 私有地址可以通过一个公共IP地址的代理访问互联网
- 网络地址转换（NAT）用于转换私人 and 全球IP地址

IP地址实例



主机ID=所有0的地址指的是网络。

主机ID=所有1的地址指的是一个广播包

R = 路由

H = 主机

IP地址：如何获得？



主机部分

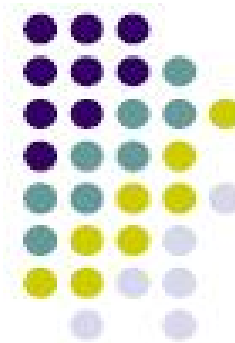
用户由互联网服务提供商（ISP）分配IP地址，通过

- 烧进系统(如引导程序)
- DHCP: Dynamic Host Configuration Protocol, 动态主机配置协议：动态地获得地址，“即插即用”



第八章 通信网络与服务

网际协议
-子网

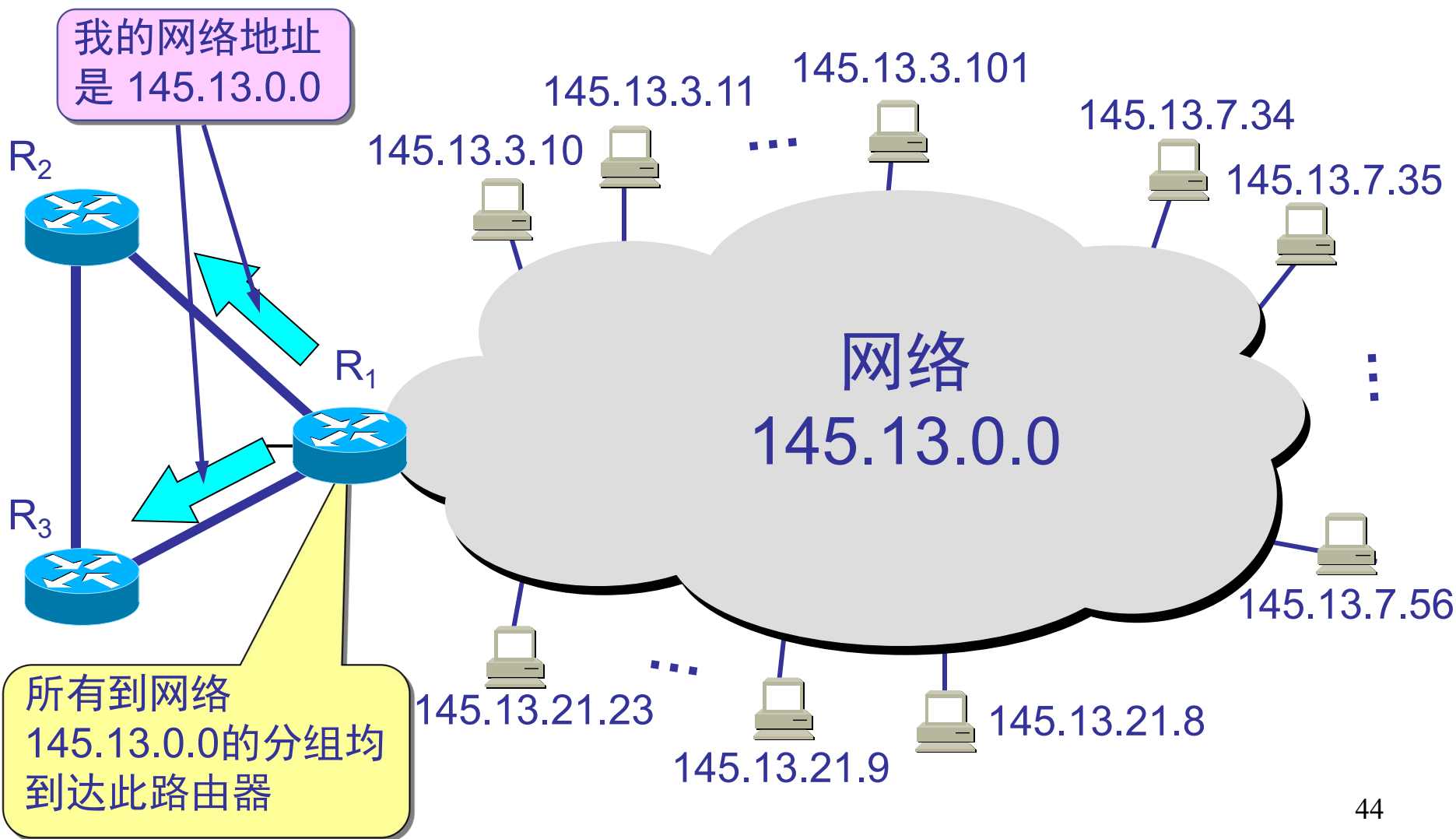




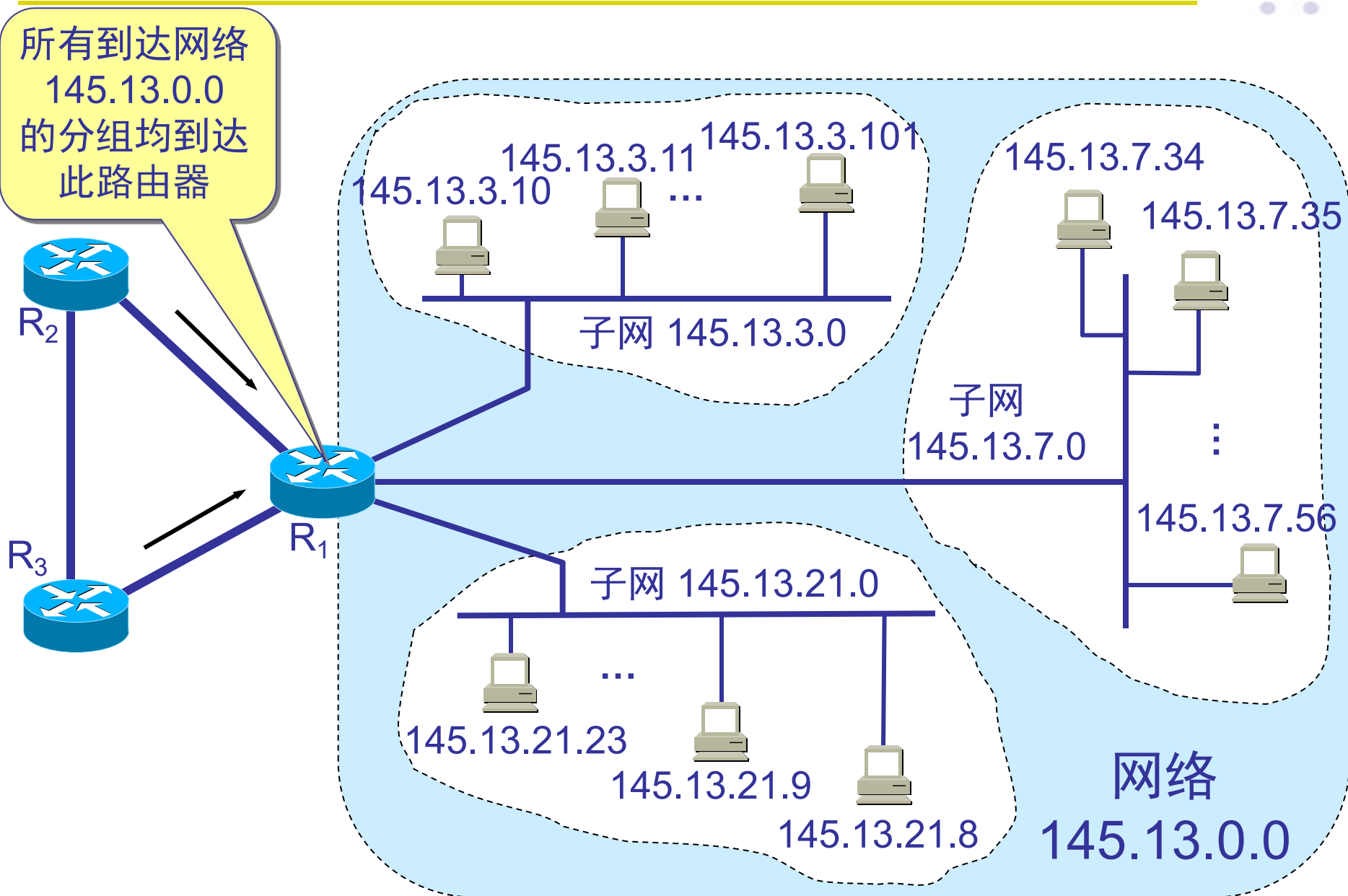
- 为什么需要子网？
 - A/B/C中最大主机数量彼此不同
 - A/B/C类不灵活
 - 简化了对局域网的管理
- 子网寻址引入了另一个层次（如141. 223. 1. 2）
- 子网形式
 - 仅在域内
 - 对远程网络是透明的
 - 将部分主机位改为子网位。
 - 引入网络掩码以了解内部如何进行子网。
 - 为了规范使用，将传统的A/B/C地址视为特殊的子网网络（只有一个子网）



一个未划分子网的 B 类网络 145.13.0.0



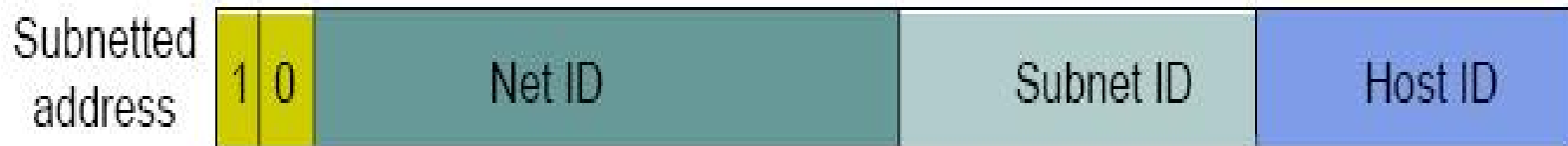
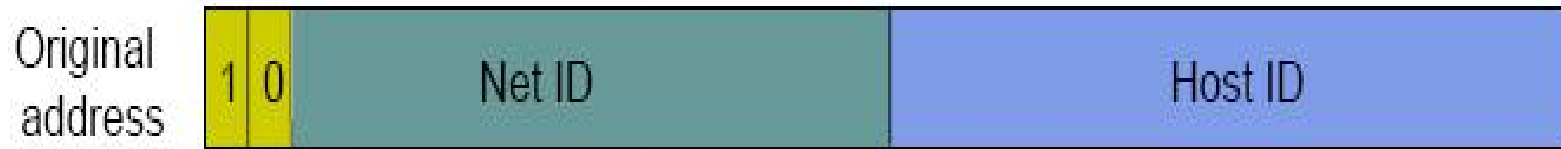
划分为三个子网后对外仍是一个网络



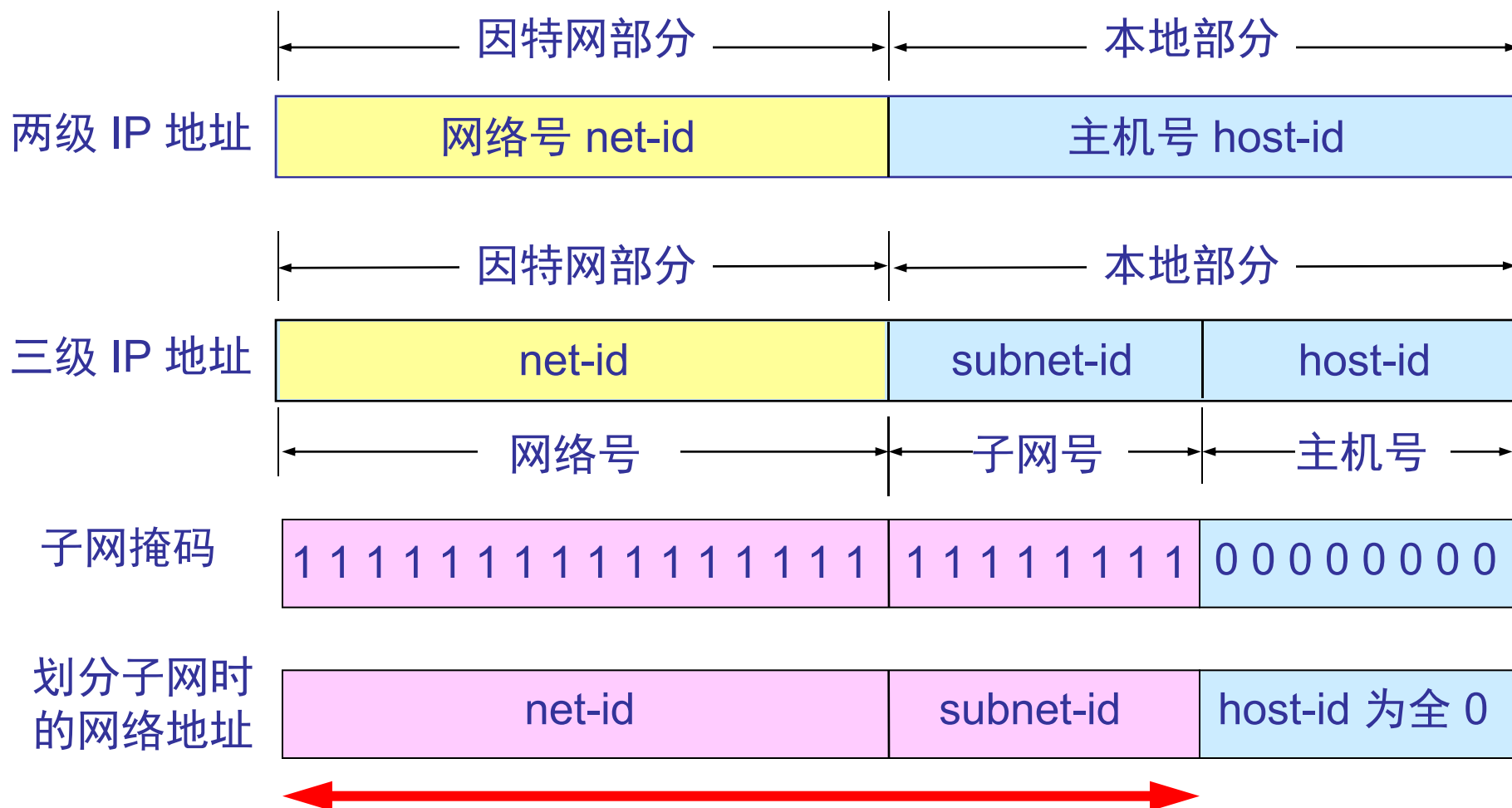
子网地址



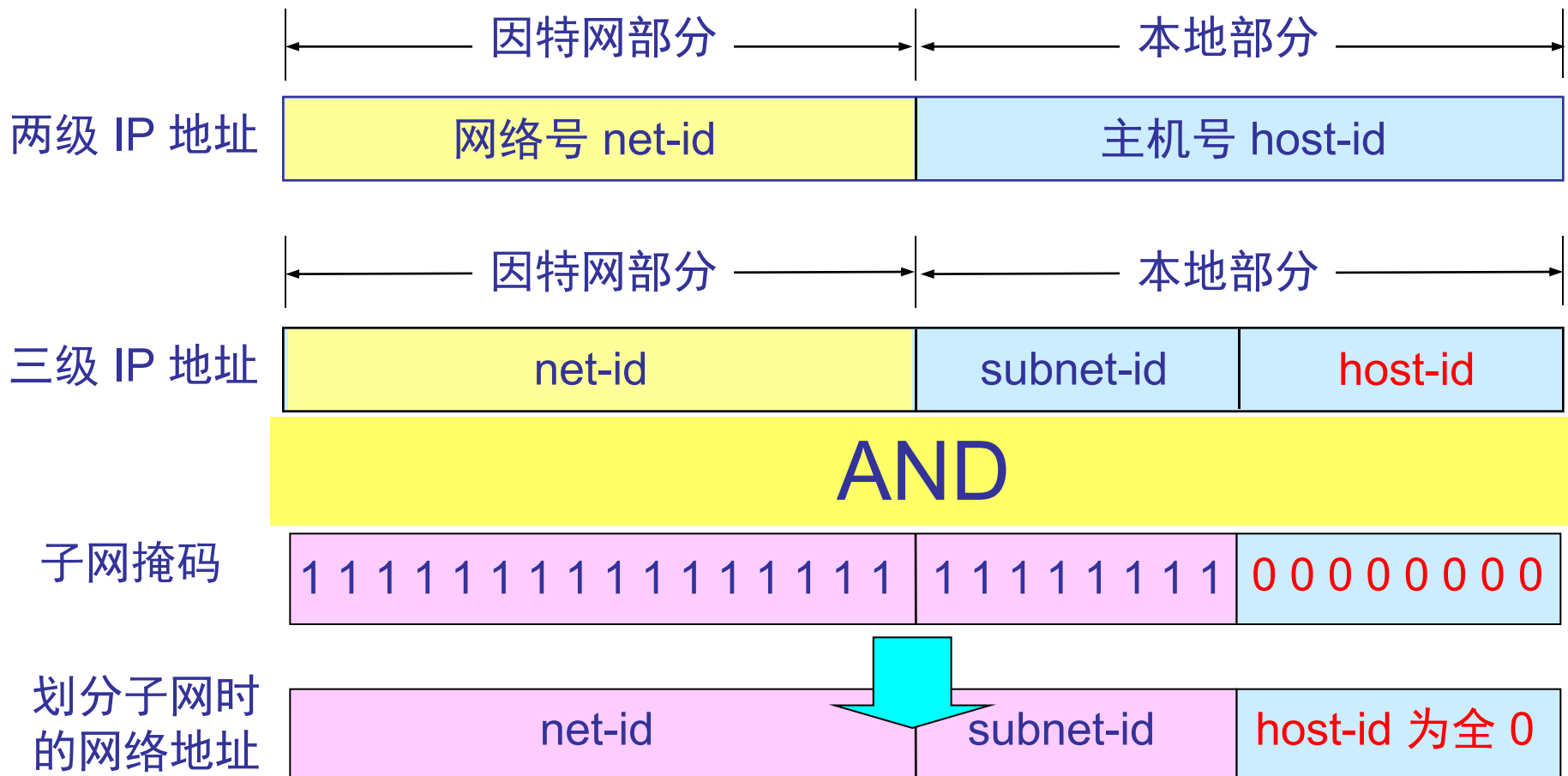
- 用来查找子网号码的掩码



IP地址和子网掩码



IP地址和子网掩码



网络掩码



掩码中的所有位等于1表示IP地址的对应位是网络部分，0表示主机部分。

- IP 地址 202.119.22.5
掩码 255.255.255.0

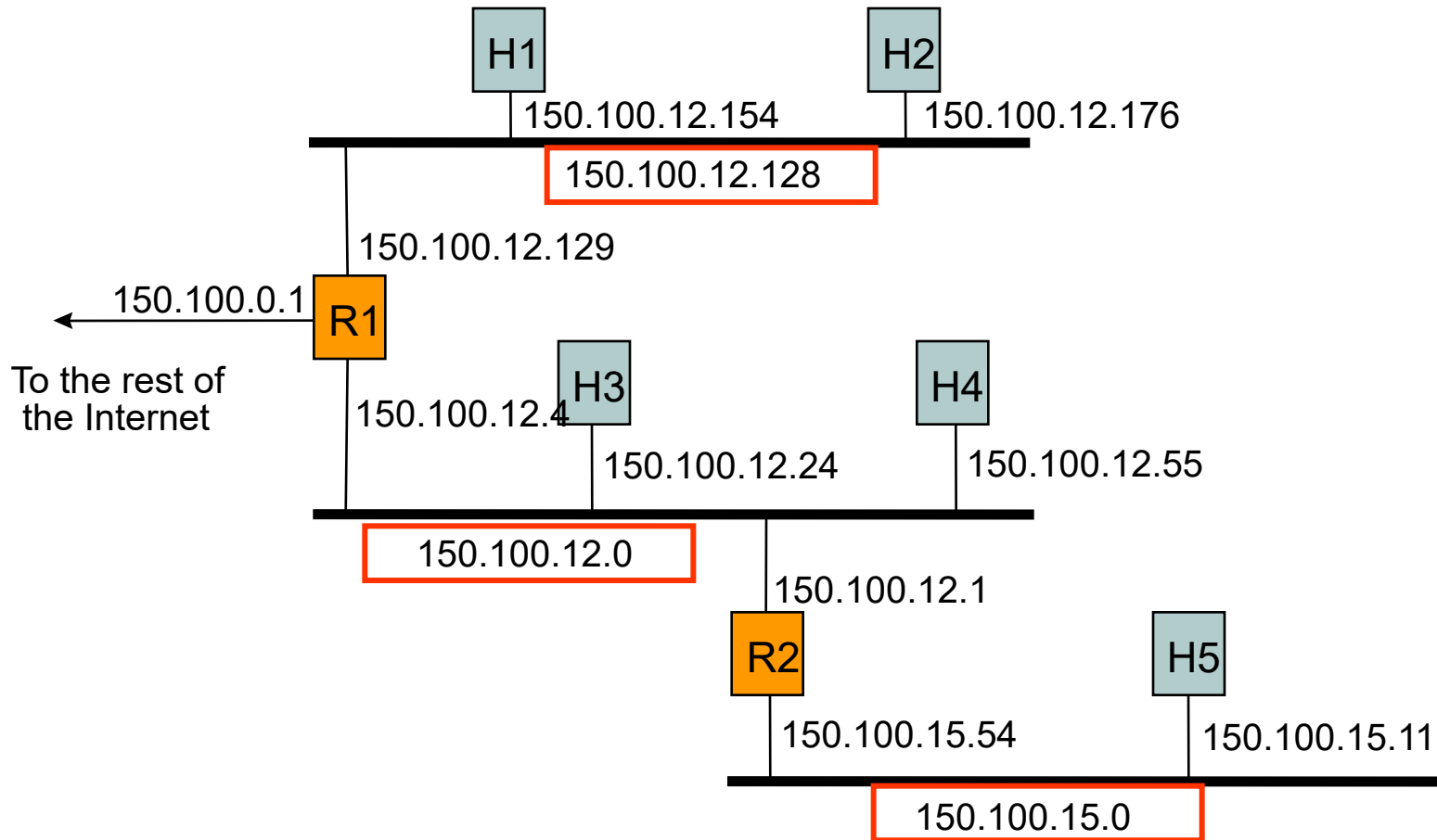
1100 1010	0111 0111	0001 0110	00000101
& 1111 1111	1111 1111	1111 1111	00000000
1100 1010	0111 0111	0001 0110	00000000
(202)	(119)	(22)	(0)

- 地址掩码相同，处于同一网络

A/B/C类网络默认掩码

A => 255.0.0.0 B=> 255.255.0.0 C=> 255.255.255.0

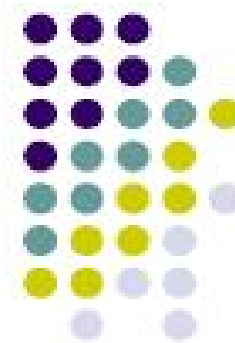
子网例子



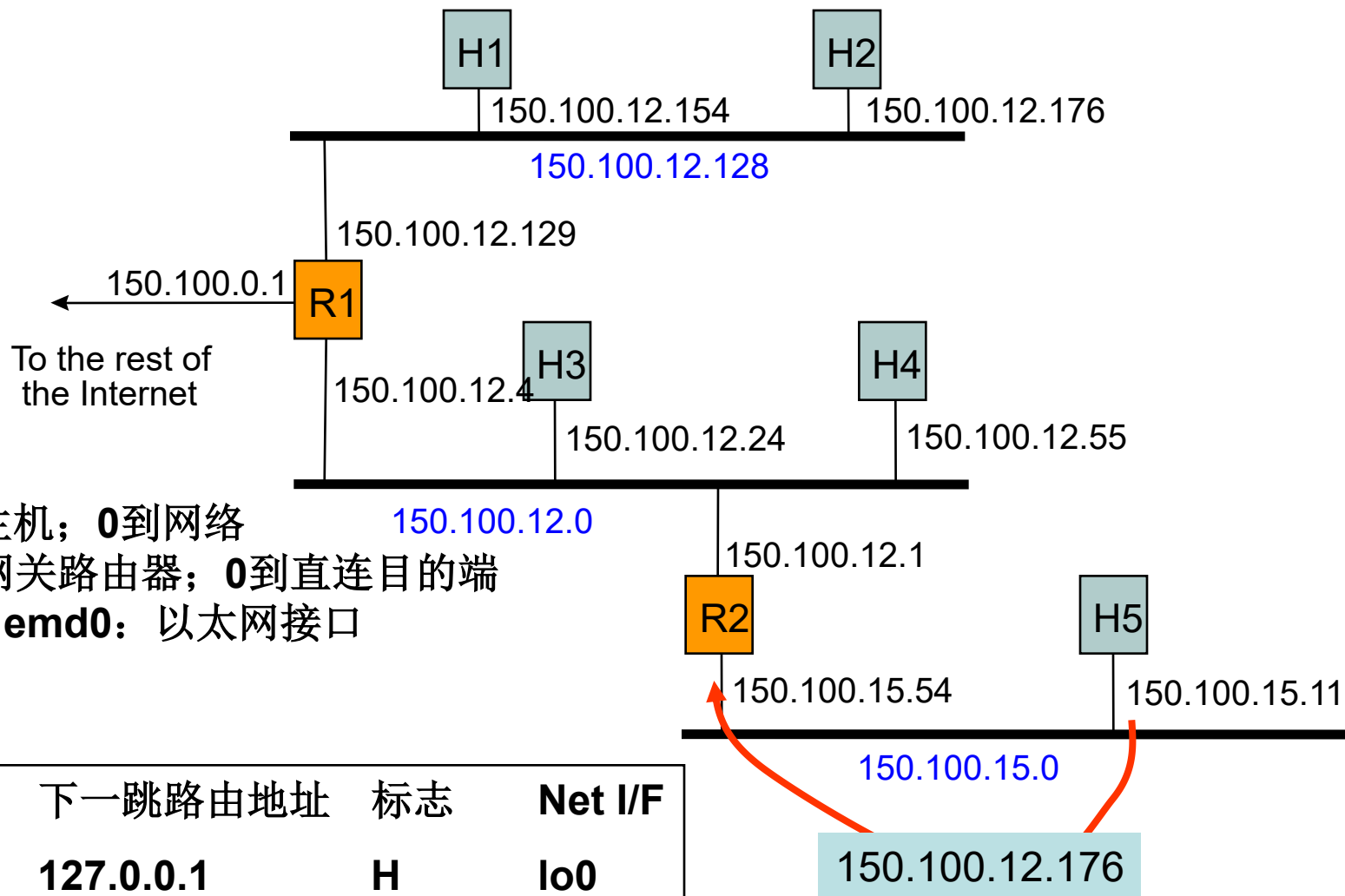
第八章 通信网络与服务



网际协议
-IP 路由



例子：主机H5向主机H2发送数据包



H: 1, 路由到主机; **0**到网络

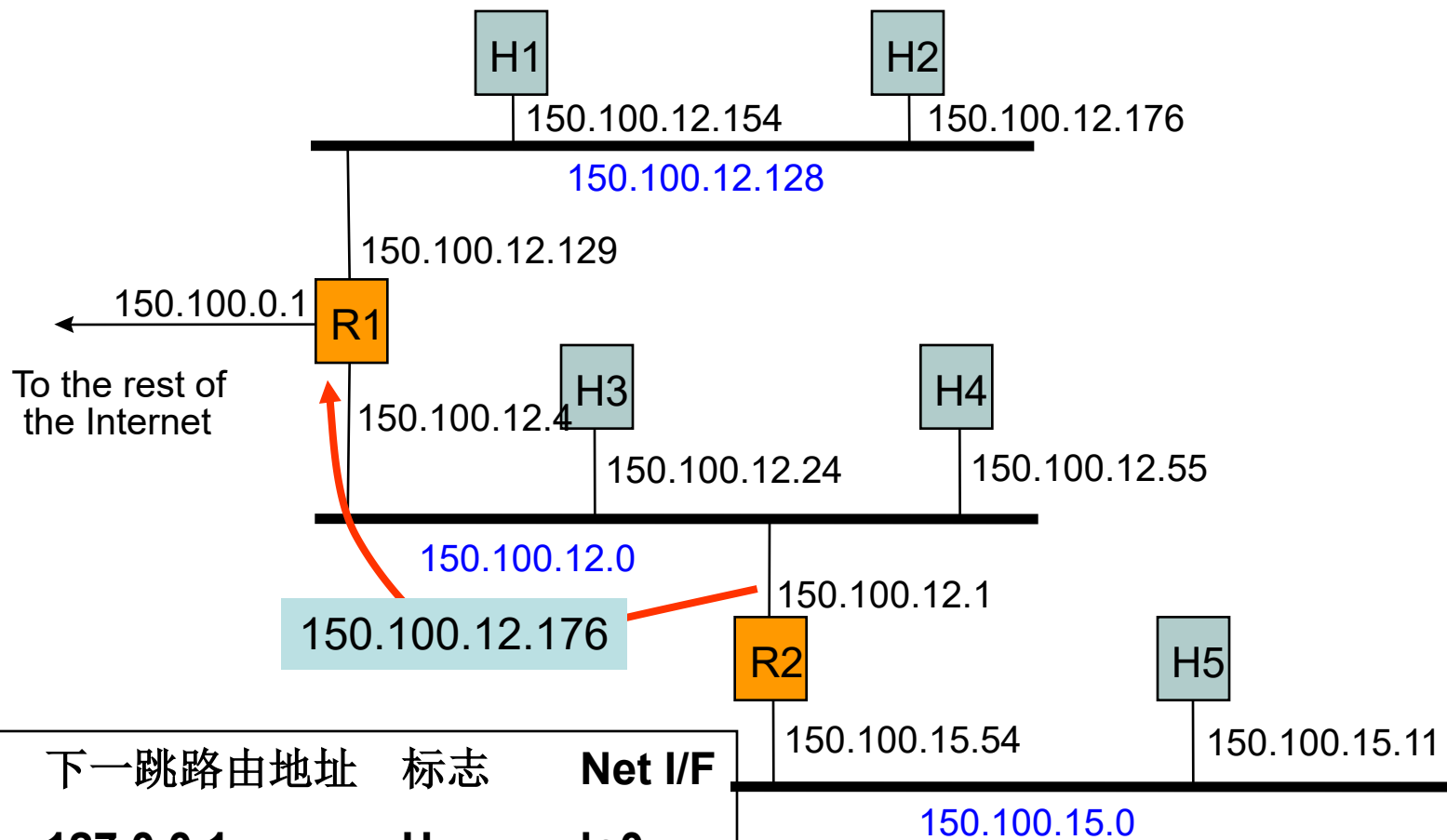
G: 1, 路由到网关路由器; **0**到直连目的端

lo0: 环回接口; **emd0**: 以太网接口

H5路由表

目的地址	下一跳路由地址	标志	Net I/F
127.0.0.1	127.0.0.1	H	lo0
default	150.100.15.54	G	emd0
150.100.15.0	150.100.15.11		emd0

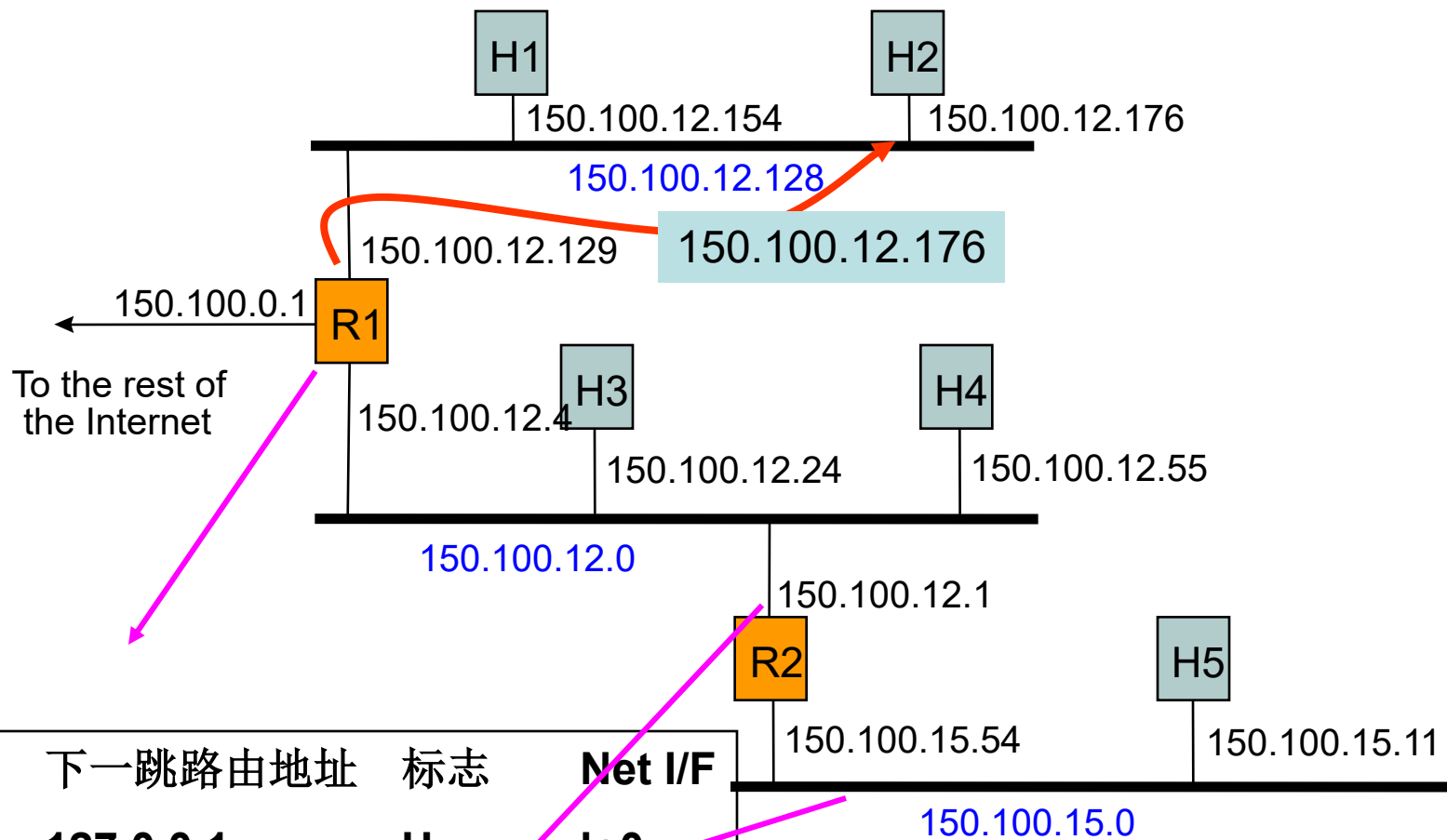
例子： 主机H5向主机H2发送数据包



R2路由表

目的地址	下一跳路由地址	标志	Net I/F
127.0.0.1	127.0.0.1	H	lo0
default	150.100.12.4	G	em0
150.100.15.0	150.100.15.54		em1
150.100.12.0	150.100.12.1		em0

例子：主机H5向主机H2发送数据包



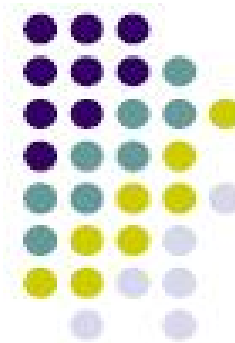
R1路由表

目的地址	下一跳路由地址	标志	Net I/F
127.0.0.1	127.0.0.1	H	lo0
150.100.12.176	150.100.12.176		emd0
150.100.12.0	150.100.12.4		emd1
150.100.15.0	150.100.12.1	G	emd1

第八章 通信网络与服务



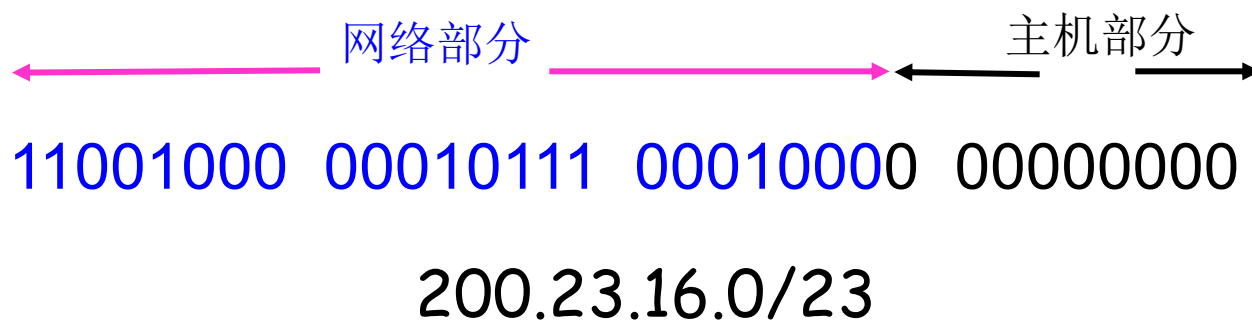
网际协议
- CIDR



CIDR: Classless Inter Domain Routing



- 分级寻址:
 - 地址空间使用效率低下，地址空间耗尽
 - 例如，B类网络为65K个主机分配了足够的地址，即使该网络中只有2K个主机。
- CIDR:无分类域间路由
 - 地址的网络部分具有任意的长度
 - 地址格式：a.b.c.d/x，其中x是地址的网络部分的#位数



CIDR 例子



128.14.32.0/20 表示的地址 (2^{12} 个地址)

最小地址 →

所有地址的 20 bit 前缀都是一样的

最大地址 →

10000000	00001110	00100000	00000000
10000000	00001110	00100000	00000001
10000000	00001110	00100000	00000010
10000000	00001110	00100000	00000011
10000000	00001110	00100000	00000100
10000000	00001110	00100000	00000101
...			
10000000	00001110	00101111	11111011
10000000	00001110	00101111	11111100
10000000	00001110	00101111	11111101
10000000	00001110	00101111	11111110
10000000	00001110	00101111	11111111



一组IP地址的分配

194.0.0.0	–	195.255.255.255	Europe
198.0.0.0	–	199.255.255.255	North America
200.0.0.0	–	201.255.255.255	Central/South America
202.0.0.0	–	203.255.255.255	Asia and the Pacific

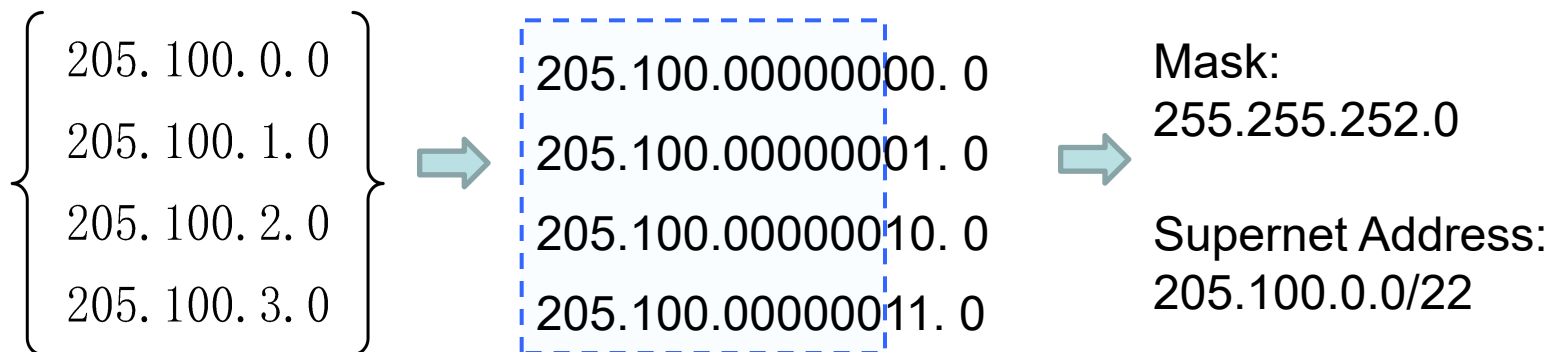
University	First address	Last address	How many	Written as
Cambridge	194.24.0.0	194.24.7.255	2048	194.24.0.0/21
Edinburgh	194.24.8.0	194.24.11.255	1024	194.24.8.0/22
(Available)	194.24.12.0	194.24.15.255	1024	194.24.12/22
Oxford	194.24.16.0	194.24.31.255	4096	194.24.16.0/20

CIDR



- CIDR不是一个独立的路由算法。它必须与其他算法一起工作。
- 标准A/B/C类IP地址可以描述为/8 /16 /24。
- CIDR和网络掩码可以进行转换。CIDR /n表示在网络掩码中，有多少个高位是网络位。
 - /8<=>掩码255.0.0.0
 - /12<=>掩码255.240.0.0
 - 等等

[CIDR]例子



超网



- 使用可变长度的掩码汇总一组连续的C类地址
- 例如：205.100.0.0/22

IP地址（205.100.0.0）和掩码长度（22）。

IP地址=11001101 01100100 00000000 00000000

掩码 = 11111111 11111111 11111100 00000000

- 包含4个C类区块。

从10101101 01100100 00000000 00000000

即：205.100.0.0

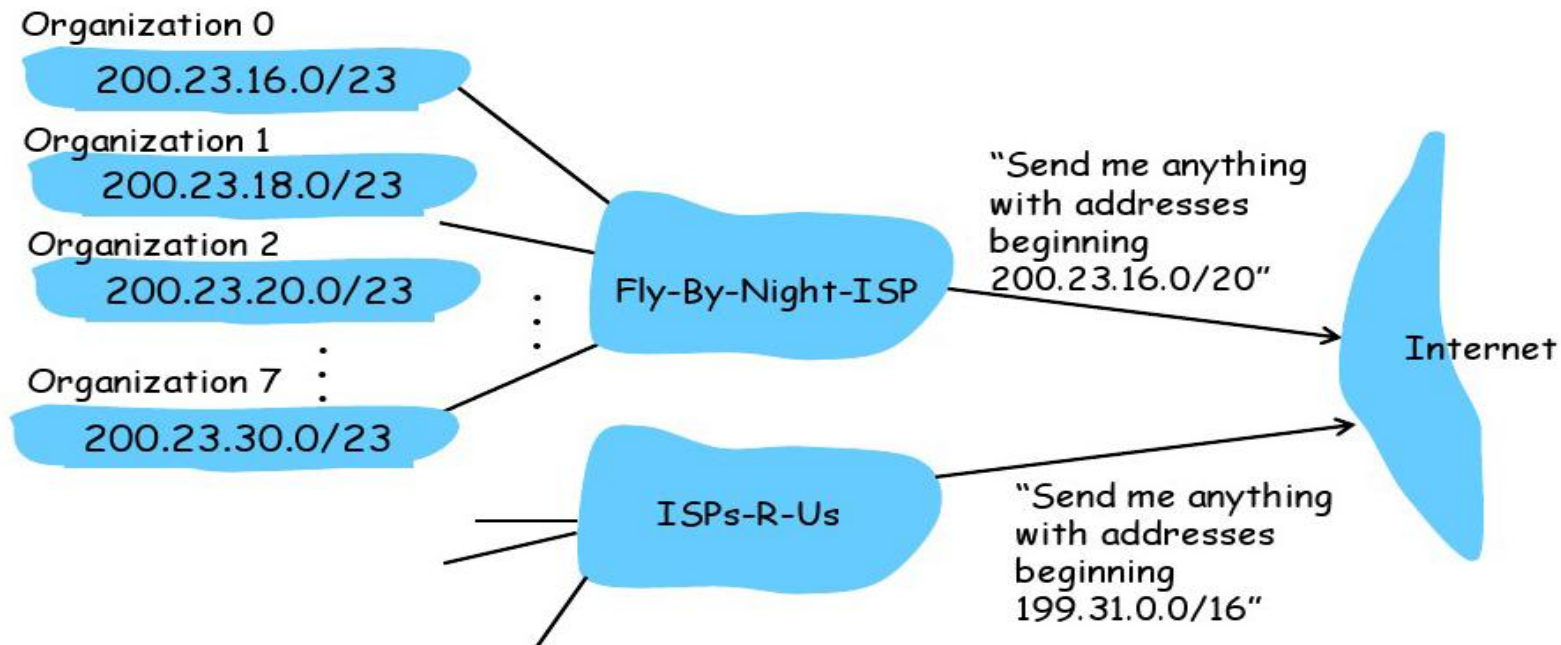
至10101101 01100100 00000011 00000000

即205.100.3.0

例子 [CIDR 路由聚合]



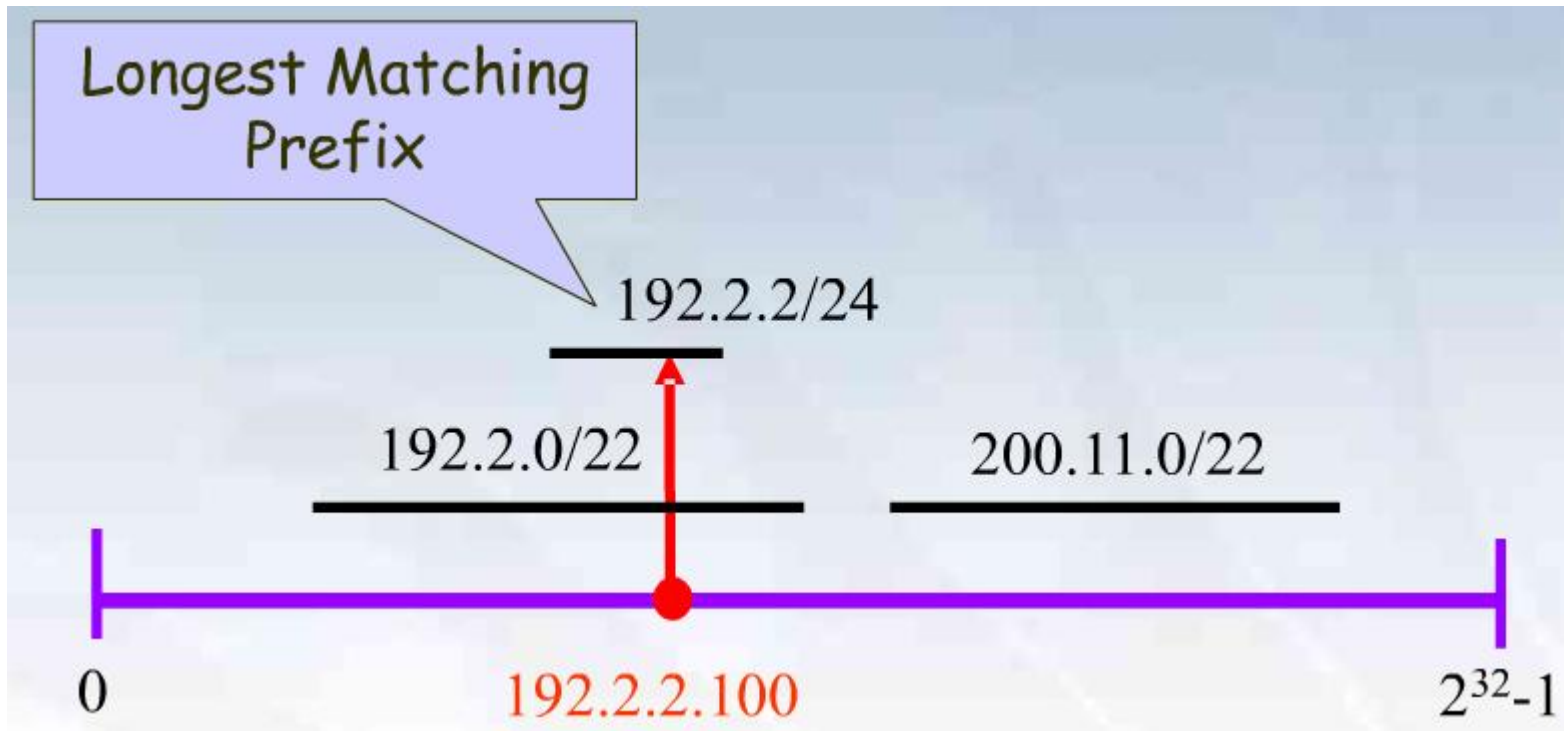
ISP's block	<u>11001000</u>	<u>00010111</u>	<u>00010000</u>	00000000	200.23.16.0/20
Organization 0	<u>11001000</u>	<u>00010111</u>	<u>00010000</u>	00000000	200.23.16.0/23
Organization 1	<u>11001000</u>	<u>00010111</u>	<u>00010010</u>	00000000	200.23.18.0/23
Organization 2	<u>11001000</u>	<u>00010111</u>	<u>00010100</u>	00000000	200.23.20.0/23
...	
Organization 7	<u>11001000</u>	<u>00010111</u>	<u>00011110</u>	00000000	200.23.30.0/23



路由表查询



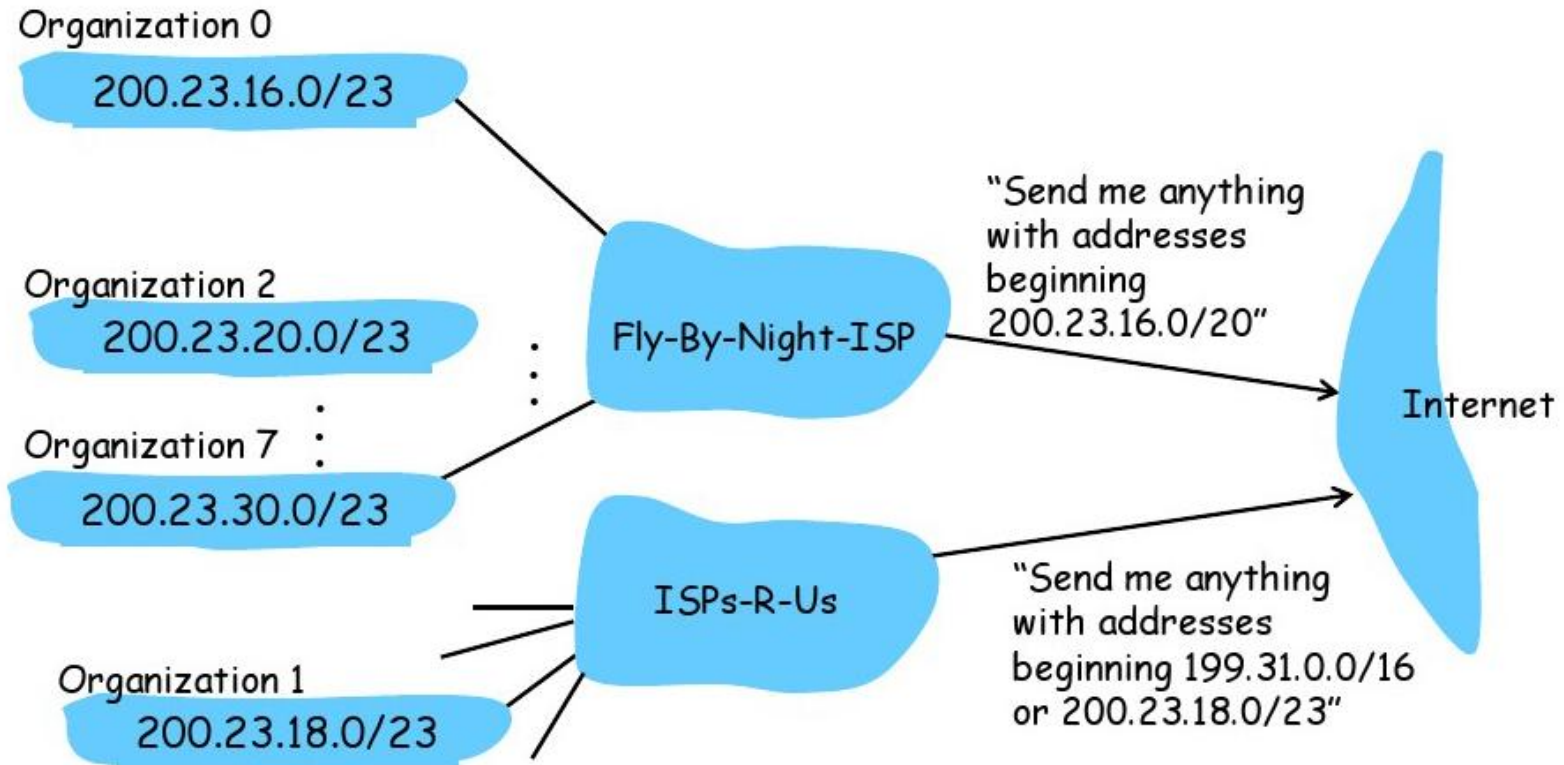
- 最长前缀匹配



最长前缀匹配



- CIDR对路由和转发的影响
- 路由表和路由协议必须携带
IP地址和掩码
- 多个条目可以匹配一个给定的IP目标地址
- 例如：路由表可能包含
 - 205.100.0.0/22，对应于一个给定的超级网。
 - 205.100.0.0/20，这是将更多的目的地汇总到一个超级网中的结果。
 - 数据包必须使用更具体的路由，即最长前缀匹配



ISPs-R-Us has a more specific route to Organization 1
200.23.18.0

路由表



- 手动设置机器为每个前缀发送数据包的位置
- 路线打印（在windows平台）
- 例子

Active Routes:

Network	Destination	Netmask	Gateway	Interface	Metric
	0.0.0.0	0.0.0.0	202.119.27.1	202.119.27.171	1
	127.0.0.0	255.0.0.0	127.0.0.1	127.0.0.1	1
	192.168.17.0	255.255.255.0	192.168.17.1	192.168.17.1	20
	192.168.17.1	255.255.255.255	127.0.0.1	127.0.0.1	20
	192.168.17.255	255.255.255.255	192.168.17.1	192.168.17.1	20
	192.168.247.0	255.255.255.0	192.168.247.1	192.168.247.1	20

...

Default Gateway: 202.119.27.1

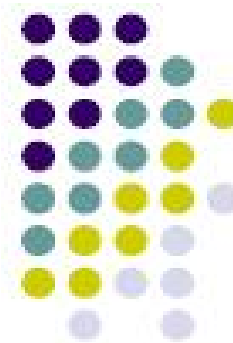
第八章 通信网络与服务



网际协议

- ARP

- ICMP



物理地址



- 局域网（和其他网络）将物理地址分配给网络的物理附件
- 网络使用自己的地址将数据包或帧传送到适当的目的地
- IP地址需要在每个IP网络接口被解析为物理地址
- 例如。 以太网使用48位地址
 - 每个以太网网络接口卡（NIC）都有全球唯一的介质访问控制（MAC）或物理地址
 - 前24位识别网卡制造商；后24位是序列号
 - 00: 90: 27: 96: 68: 07 12个十六进制数字

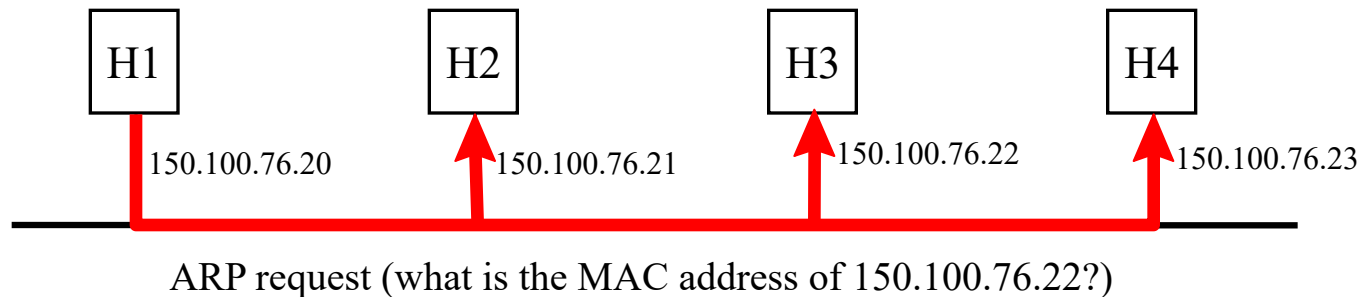
1、Address Resolution Protocol -- ARP



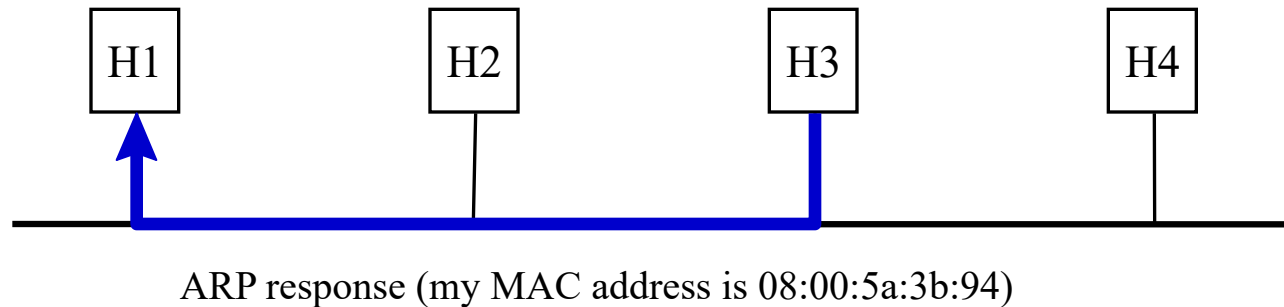
虽然IP地址标识了一个主机，但数据包是由底层网络（如以太网）物理传递的，它使用自己的物理地址（以太网的MAC地址）。如何将一个IP地址映射到一个物理地址？

H1、H2、H3、H4都在同一个网络中（如何判断？）

H1想知道H3的物理地址 -> 广播一个ARP请求



每台主机都收到请求，但只有H3用它的物理地址回复。

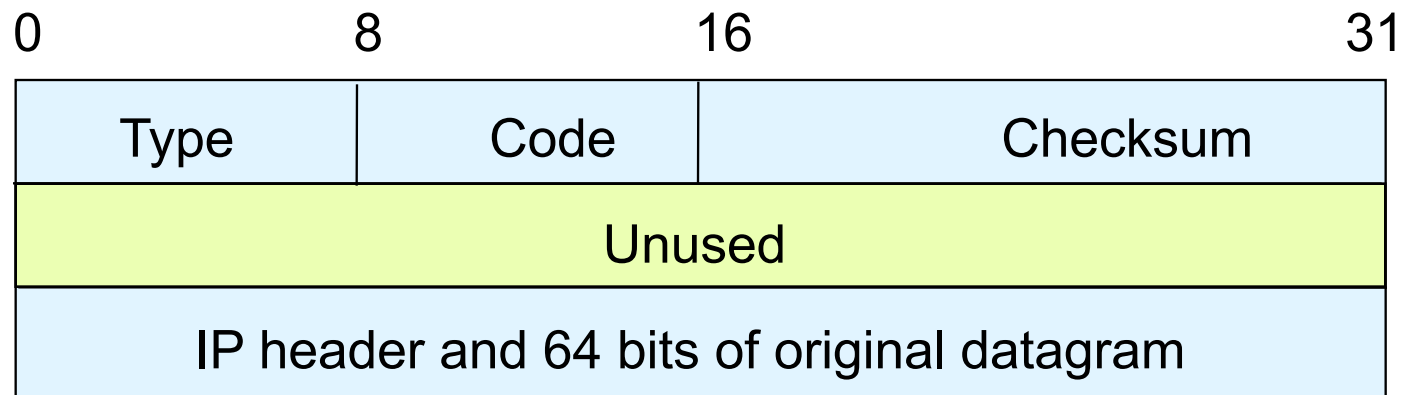


2、Internet Control Message Protocol (ICMP)



- RFC792; 封装在IP包中（协议类型=1）。
- 处理错误和控制信息
 1. 如果路由器不能传递或转发一个数据包，它会向源头发送一个ICMP “主机不可达”消息。
 2. 如果路由器收到本应发送到另一个路由器的数据包，它会向发件人发送一个ICMP “重定向”消息；发件人修改其路由表
 3. ICMP “路由器发现”消息允许主机了解其网络中的路由器并初始化和更新其路由表
 4. ICMP “回波”请求和回复有助于诊断，并用于 “ping”。

ICMP基本错误信息格式



- *Type* of message: some examples

- 0 Network Unreachable;
- 3 Port Unreachable
- 1 Host Unreachable
- 4 Fragmentation needed
- 2 Protocol Unreachable
- 5 Source route failed
- 11 Time-exceeded, code=0 if TTL exceeded

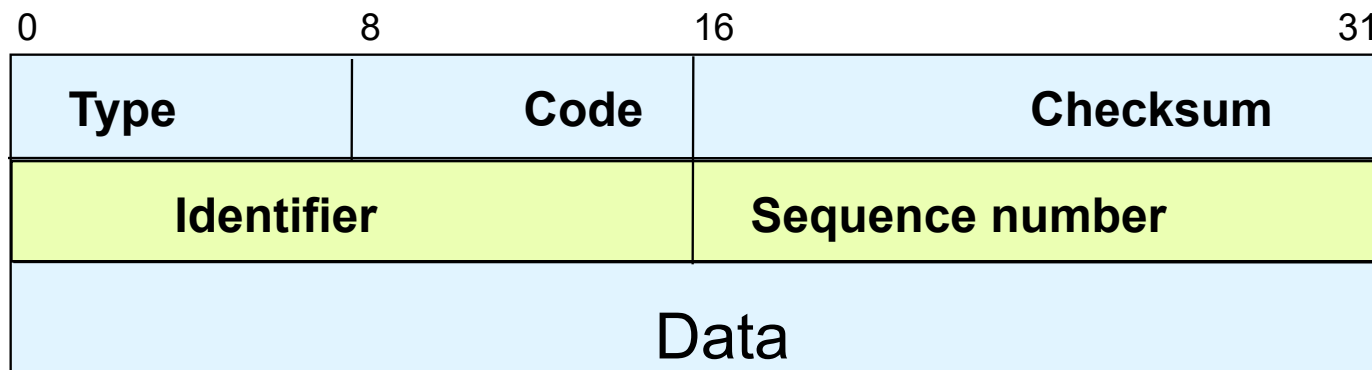
- *Code*: purpose of message

- IP header & 64 bits of original datagram

To match ICMP message with original data in IP packet

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

回声请求和回声回复消息格式



- Echo request: type=8; Echo reply: type=0
 - 目的地通过将请求中的数据复制到回复信息中，用回声回复进行回复
- 序列号用于匹配回复和请求
- 用于区分使用回声服务的不同会话的ID。
- 在PING中使用

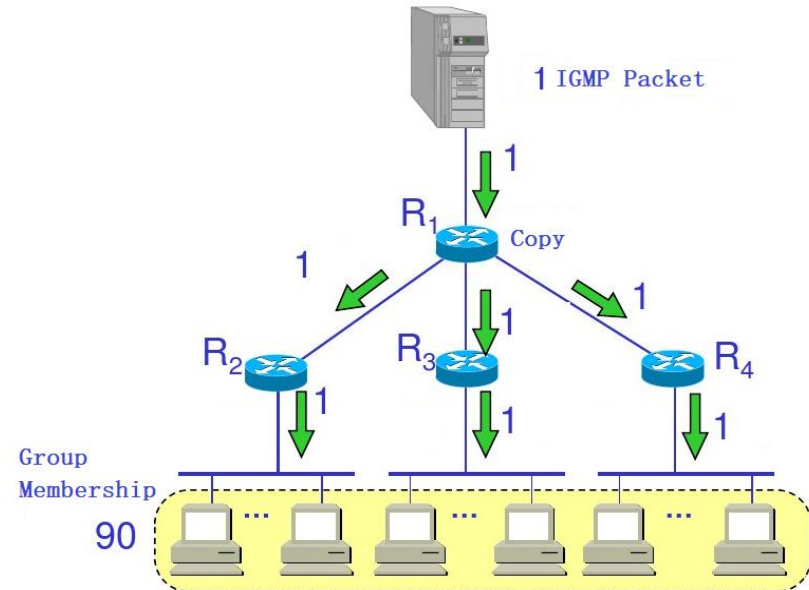
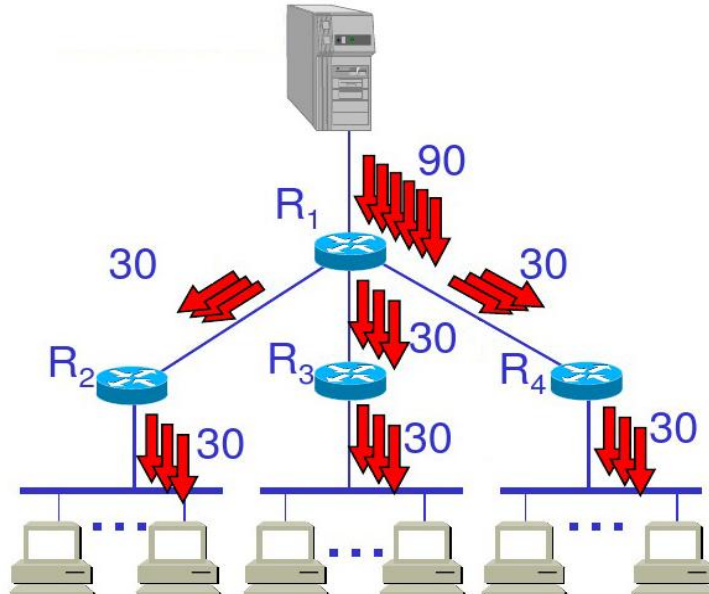
IGMP (Internet Group Management Protocol)



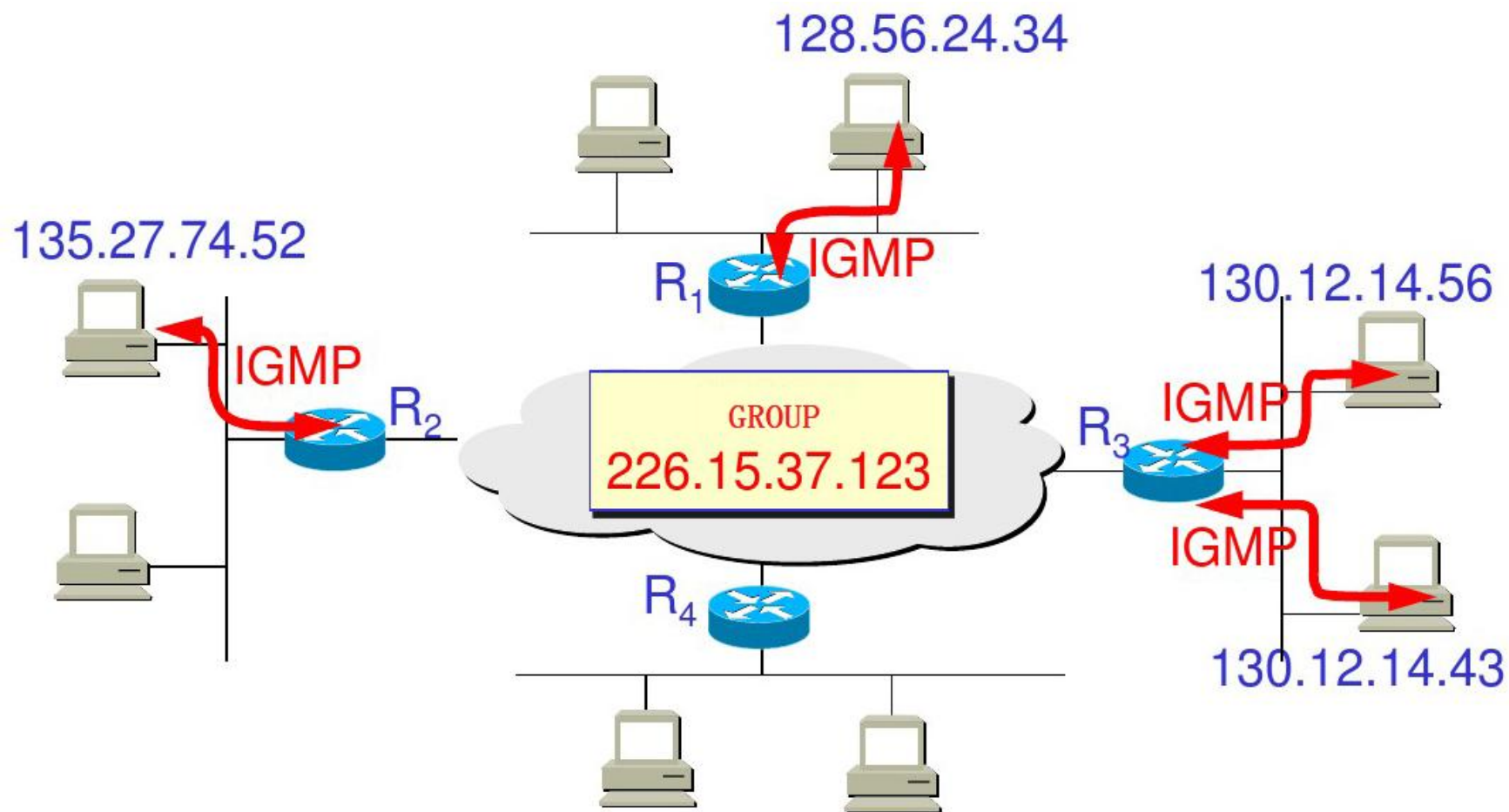
更有效地利用资源

1. IGMP_{V1} – RFC1112 in 1989
2. IGMP_{V2} – RFC2236 in 1997
3. IGMP_{V3} – RFC3376 in 2002

只能使用**D**类地址作为目的地址，而不能作为原地址！

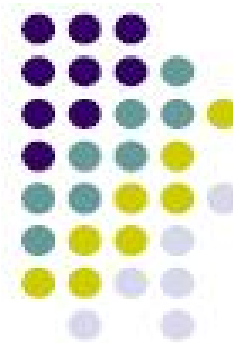


组成员



第八章 通信网络与服务

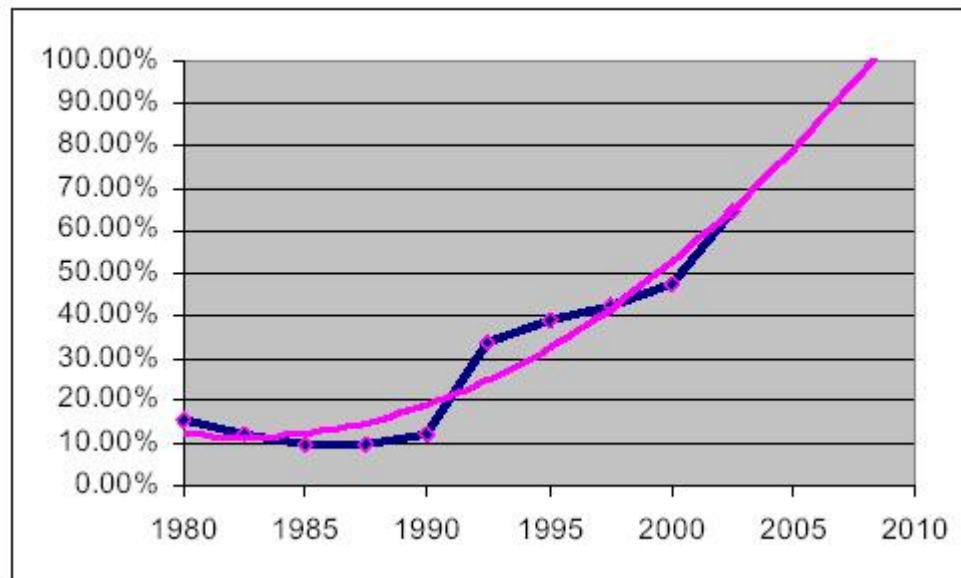
网际协议 - IPV6 地址



IP地址分配史



1981年--IPv4协议发布
1985年 ~ 总空间的1/16
1990年~占总空间的1/8
1995年~占总空间的1/3
2000年~占总空间的1/2
2002.5 ~ 占总空间的2/3
2011.2 ~ IANA将最后5个IP地址块分配给五个RIR



IANA RIR NIR ISP中的用尽情况

可能的解决方案。

- NAT网络地址转换
- DHCP地址共享/PPPoE
- CIDR（无分类域间路由）
- 加上一些地址回收
- 32位空间的理论极限：~40亿台设备
- 32位空间的实际限制：~2.5亿个设备（RFC 3194）（私有IP，广播）

更多IPv4不足



- 网络本身
 - 点对点困难
 - NAT的安全记录可疑
 - 管理困难
- 安全性
- QoS（服务质量）不足，且不是实时
- 路由表太大，处理速度慢
- 移动性
 - 但点对点的移动性是互联网的未来
- 设备自动配置困难
 - DHCP和地址所有权不能跨越组织边界工作
 - 使用外部代理进行自动配置是不可能的

IPv6相较于IPv4的优势



1. 地址耗尽问题得到解决
2. 解决了国际错配问题
3. 恢复了端到端的通信
4. 范围内的地址和地址选择
5. 更有效的转发
6. 内置安全和移动性

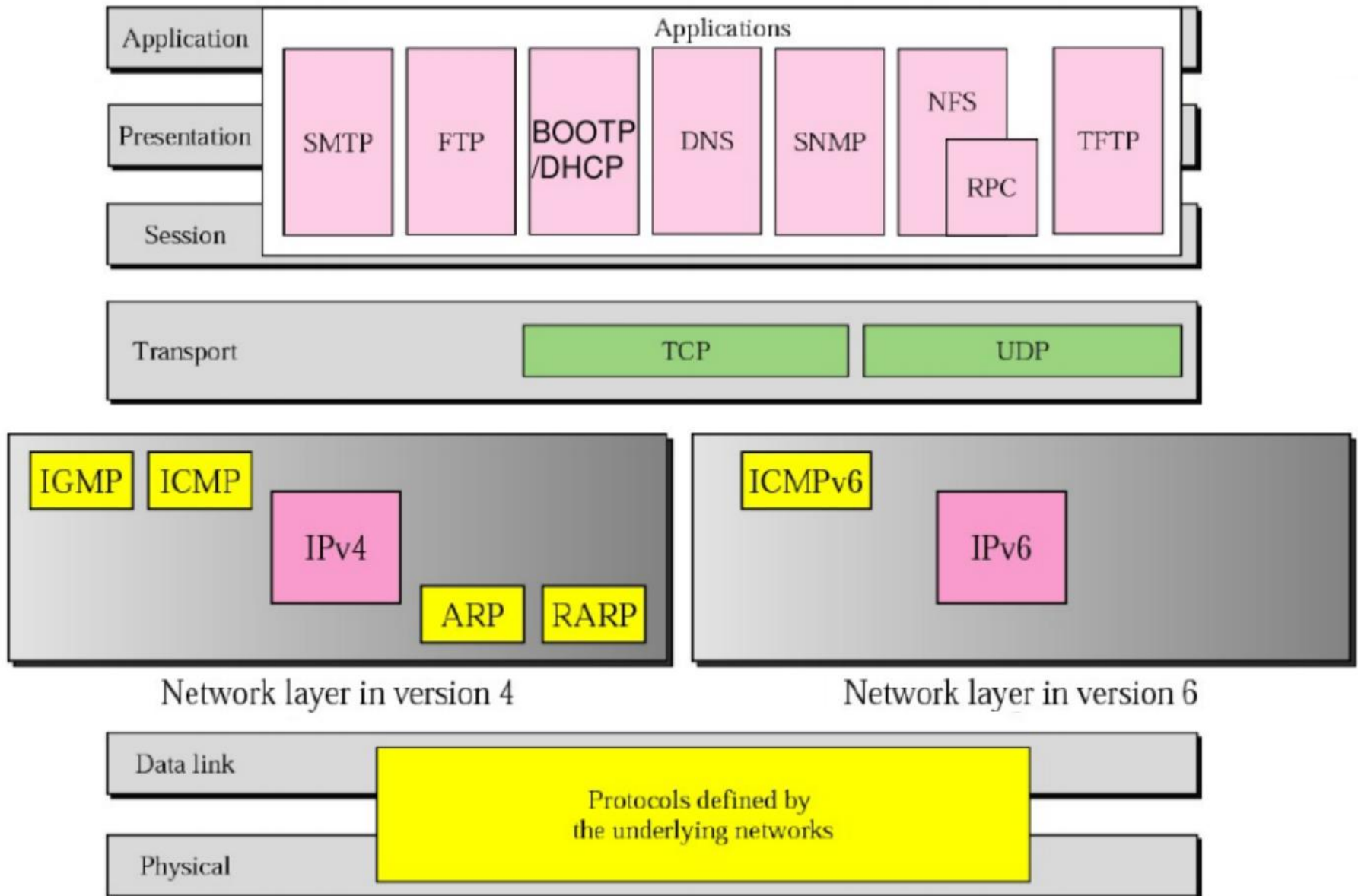
世界IPv6日 2011年6月8日

世界IPv6启动日 2012年6月6日

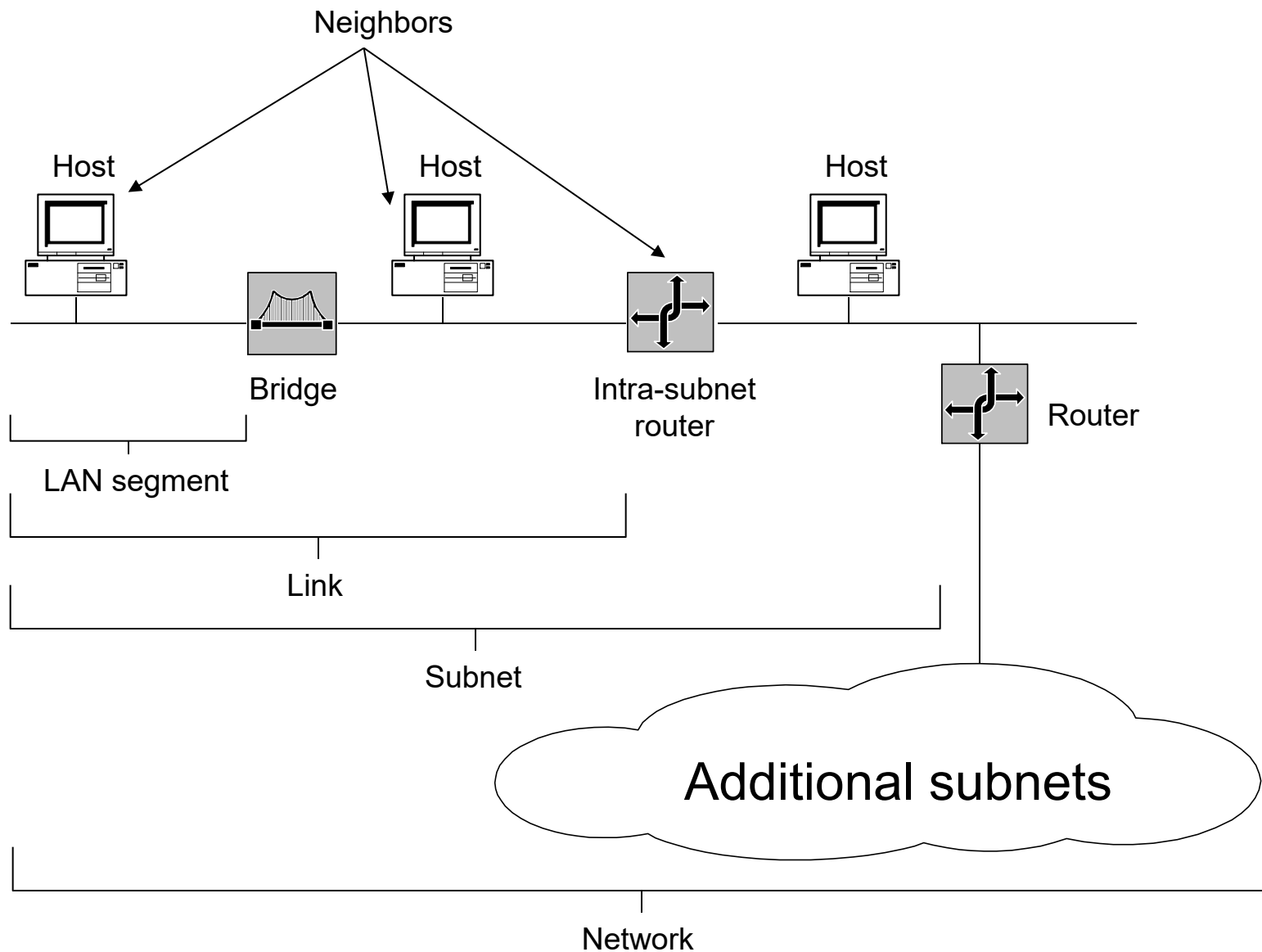
IPv6特点

1. 新的报头格式
2. 大的地址空间
3. 高效和分层的寻址和路由基础设施
4. 无状态和有状态的地址配置
5. 内置安全
6. 更好地支持QoS
7. 邻近节点互动的新协议
8. 可扩展性

TCP/IP和OSI模型



相关的基本概念



IPv4 & IPv6报头



IPv4: 20 Bytes + Options

Version (4)	Header Len. (4)	Type of Service (8)	Datagram Length (16)	
Identifier (16)			Flags (3)	Fragmentation Offset (13)
Time to Live (8)	Upper-layer Protocol (8)		Header Checksum (16)	
Source IP Address (32)				
Destination IP Address (32)				
Options (if any)				

IPv6: 40 Bytes

Version (4)	Traffic Class (8)	Flow Label (20)		
Payload Length (16)		Next Header (8)	Hop Limit (8)	
Source IP Address (128)				
Destination IP Address (128)				

IPv4 VS IPv6



特点	IPv4	IPv6
地址长度	32 bits	128 bits
支持IPSec	Optional	Required
支持QoS	Some	Better
分片	Hosts and routers	Hosts only
数据报大小	576 bytes	1280 bytes
头部校验和	Yes	No
头部选项	Yes	No
链路层地址解析	ARP (broadcast)	Multicast Neighbor Discovery Messages
组播成员	IGMP	Multicast Listener Discovery (MLD)
路由器搜索	Optional	Required
使用广播	Yes	No
配置	Manual, DHCP	Automatic, DHCP
DNS名称查询	Uses A records	Uses AAAA records
DNS反向查询	Uses IN-ADDR.ARPA	Uses IP6.INT

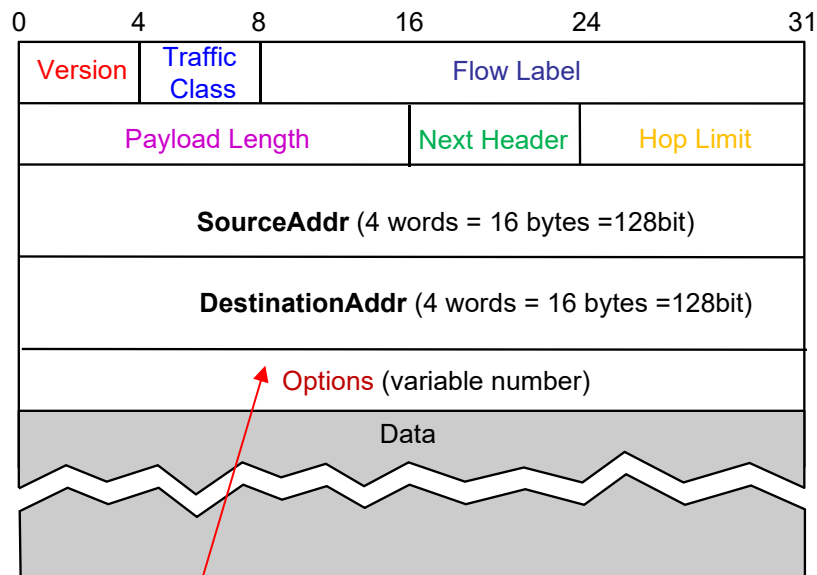
头部字段的不同



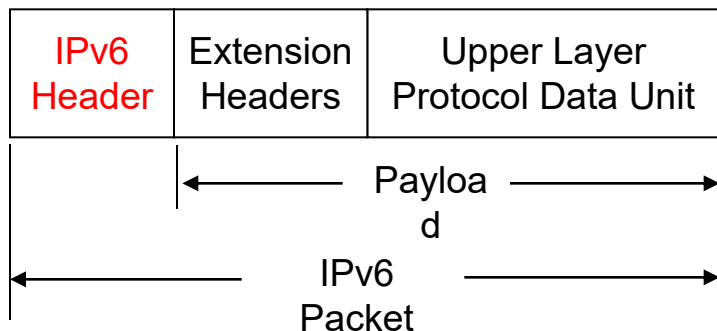
- 没有校验和
- 没有头部长度域
 - 头部没有长度变化
 - 固定格式加快了处理速度
- 没有分片和重组相关字段（也没有ID字段）
 - 过大数据包会被丢弃，并向发送者发送消息以减少数据包的尺寸
 - 主机应该做路径MTU发现
 - 主机要能够对数据包进行分割



IPv6 数据包 (PDU) 结构

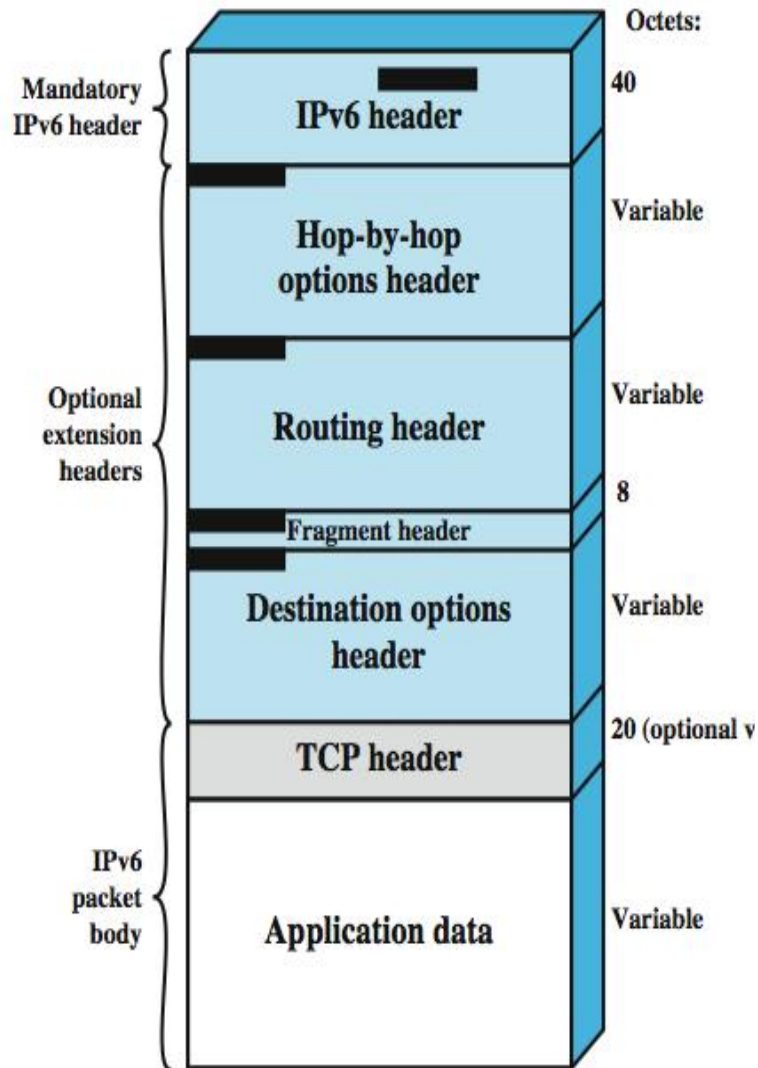


- 版本=6
- 业务类型=IPv4中的ToS
- 流标签：支持QoS和公平的带宽分配
- 有效载荷长度：不包含报头，将数据包限制在64KB以内
- 跳数限制=TTL字段
- 下一个标头：确定扩展头的类型（选项和/或上层协议）。
- 选项。“扩展头”
 - 有序的图元列表 – 6种常见类型



Header code	Header type
0	Hop-by-hop options header
43	Routing header
44	Fragment header
51	Authentication header
52	Encapsulating security payload header
60	Destination options header

下一个头字段的值



Value	Header
0	Hop-by-Hop Options Header
6	TCP
17	UDP
41	Encapsulated IPv6 Header
43	Routing Header
44	Fragment Header
50	Encapsulating Security Payload
51	Authentication Header
58	ICMPv6
59	No next header
60	Destination Options Header

Daisy-chain extension headers

分片扩展



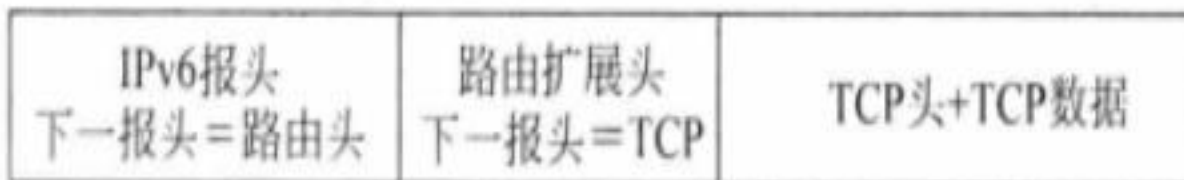
- 类似于IPv4的分片
 - 作为一个扩展头来实施
 - 放在IPv6头和数据之间（如果它是唯一使用的扩展）。
 - 13位偏移
 - 最后一个片段标记(M)
 - 比IPv4更大的片段ID字段
- 分片在终端主机上完成的

0	8	16	29	31
next header	reserved	offset	reserved	M
ID				

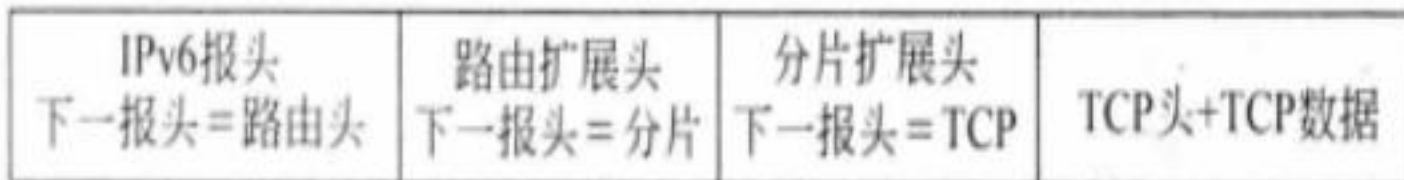
扩展头使用



(a) 0个扩展头



(b) 1个扩展头



(c) 多个扩展头

扩展头使用示例



- 128位地址空间
 - 2^{128} 可能的地址
 - 340, 282, 366, 920, 938, 463, 463, 374, 607, 431, 768, 211, 456个地址 (3.4×10^{38})
 - 每平方米的地球表面 (πD^2) 有 6.65×10^{23} 个地址
- 选择128位是为了在设计分层寻址和路由时允许多层次和灵活性
- 典型的单播IPv6地址。
64位为子网ID, 64位为接口ID

IPv6地址类型



- **单播Unicast**

单个接口的地址

一对一传送到单个接口

- **多(组)播Multicast**

一组接口的地址

一对多地传送到该地址所标识的所有接口

- **任播Anycast**

一组接口的地址

一对多地传送到此地址所标识的一组接口中距离源节点最近（根据使用的路由协议进行度量）的一个接口

- **不再有广播地址，在IPv6中的广播功能是通过组播来完成**

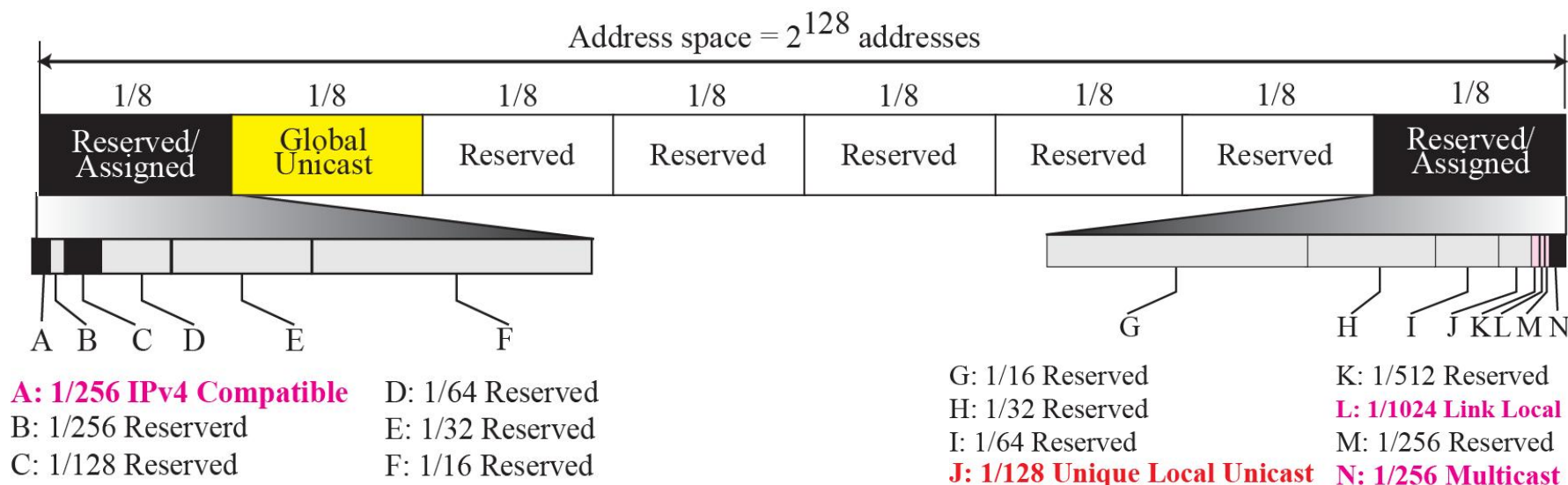




Table *Prefixes for IPv6 Addresses*

	<i>Block Prefix</i>	<i>CIDR</i>	<i>Block Assignment</i>	<i>Fraction</i>
1	0000 0000	0000::/8	Reserved (IPv4 compatible)	1/256
	0000 0001	0100::/8	Reserved	1/256
	0000 001	0200::/7	Reserved	1/128
	0000 01	0400::/6	Reserved	1/64
	0000 1	0800::/5	Reserved	1/32
	0001	1000::/4	Reserved	1/16
2	001	2000::/3	Global unicast	1/8
3	010	4000::/3	Reserved	1/8
4	011	6000::/3	Reserved	1/8
5	100	8000::/3	Reserved	1/8
6	101	A000::/3	Reserved	1/8
7	110	C000::/3	Reserved	1/8
8	1110	E000::/4	Reserved	1/16
	1111 0	F000::/5	Reserved	1/32
	1111 10	F800::/6	Reserved	1/64
	1111 110	FC00::/7	Unique local unicast	1/128
	1111 1110 0	FE00::/9	Reserved	1/512
	1111 1110 10	FE80::/10	Link local addresses	1/1024
	1111 1110 11	FEC0::/10	Reserved	1/1024
	1111 1111	FF00::/8	Multicast addresses	1/256



单播IPv6地址

- 可聚合的全局单播地址
- 链路本地地址
- 站点本地地址
- 特殊地址
- 兼容性地址
- **NSAP地址**

表 4-10 IPv6 的地址分类

地址类型	二进制前缀
未指明地址	00...0 (128 位), 可记为::/128。
环回地址	00...1 (128 位), 可记为::1/128。
多播地址	11111111 (8 位), 可记为 FF00::/8。
本地链路单播地址	1111111010 (10 位), 可记为 FE80::/10。
全球单播地址	(除上述四种外, 所有其他的二进制前缀)

结 点 地 址 (128 bit)

子网前缀 (n bit)

接口标识符 ($128 - n$) bit

全球路由选择前缀 (n bit)

子网标识符 (m bit)

接口标识符 ($128 - n - m$) bit

IPv6地址语法



- IPv6 地址二进制格式:

```
00100001110110100000000011010011000000000000000000001011110
0111011000000101010101000000000111111111111111100010100010
01110001011010
```

- 沿着16位的边界划分:

```
0010000111011010  0000000011010011  0000000000000000
0010111100111011  0000001010101010  0000000011111111
1111111000101000  1001110001011010
```

- 每个16位的数据块被转换为16进制，并以冒号为界:

21DA:00D3:0000:2F3B:02AA:00FF:FE28:9C5A

- 每个16位块内减去前导零值:

21DA:D3:0:2F3B:2AA:FF:FE28:9C5A

压缩0值



- 一些IPv6地址包含长序列的0
- 设置为0的单个连续的16位块序列可以压缩为 "::"（双冒号）
- 例子:
 - FE80:0:0:0:2AA:FF:FE9A:4CA2
FE80::2AA:FF:FE9A:4CA2
 - FF02:0:0:0:0:0:0:2
FF02::2
- 不能使用0压缩来其他16位块的一部分
FF02:30:0:0:0:0:0:5 不能压缩为 FF02:3::5

压缩0值例子



1080:0:0:0:9:802:500C:328B \Rightarrow 1080::9:802:500C:328B

FF01:0:0:0:0:0:0:101 (Multicast) \Rightarrow FF01::101

0:0:0:0:0:0:0:1 (Loopback) \Rightarrow ::1

0:0:0:0:0:0:0:0 (Unspecified) \Rightarrow ::

12AB:0000:0000:0000:0000:CD80:0000:0000 / 60 \Rightarrow 双冒号标记在一个地址中只能出现一次!
12AB::CD80:0:0 / 60 or 12AB:0:0:0:0:CD80:: / 60

例:

12AB:0:0:0:0:CD8:: / 60 (12AB:0:0:0:0:0CD8:0:0 / 60)

12AB::CD80 / 60 \Rightarrow (12AB:0:0:0:0:0:0:CD80 / 60)

12AB::CD8 / 60 (12AB:0:0:0:0:0:0:0CD8 / 60)

例子



压缩以下地址：

- a. 0000:0000:FFFF:0000:0000:0000:0000:0000
- b. 1234:2346:0000:0000:0000:0000:0000:1111
- c. 0000:0001:0000:0000:0000:0000:1200:1000
- d. 0000:0000:0000:0000:0000:FFFF:24.123.12.6

答案

- a. 0:0:FFFF::
- b. 1234:2346::1111
- c. 0:1::1200:1000
- d. ::FFFF:24.123.12.6

IPv6前缀



- 前缀是地址的一部分，其中的位有固定的值，是路由或子网的标识位
- IPv6子网或路由总是使用CIDR符号。
地址/前缀(长度)
- 例子:
 - 21DA:D3:: / 48 为一个路由
 - 21DA:D3:0:2F3B:: / 64 为一个子网
- 不再有点分十进制子网掩码!



IPv6特殊地址

- 未指明地址

- 0:0:0:0:0:0:0:0 or ::

- 环回地址

- 0:0:0:0:0:0:0:1 or ::1

- 多播地址

- FF01:0:0:0:0:0:0:101 or FF01::101

- DNS服务器通常在:

- FEC0:0:0:0:FFFF::1, FEC0:0:0:0:FFFF::2,
FEC0:0:0:0:FFFF::3

兼容地址

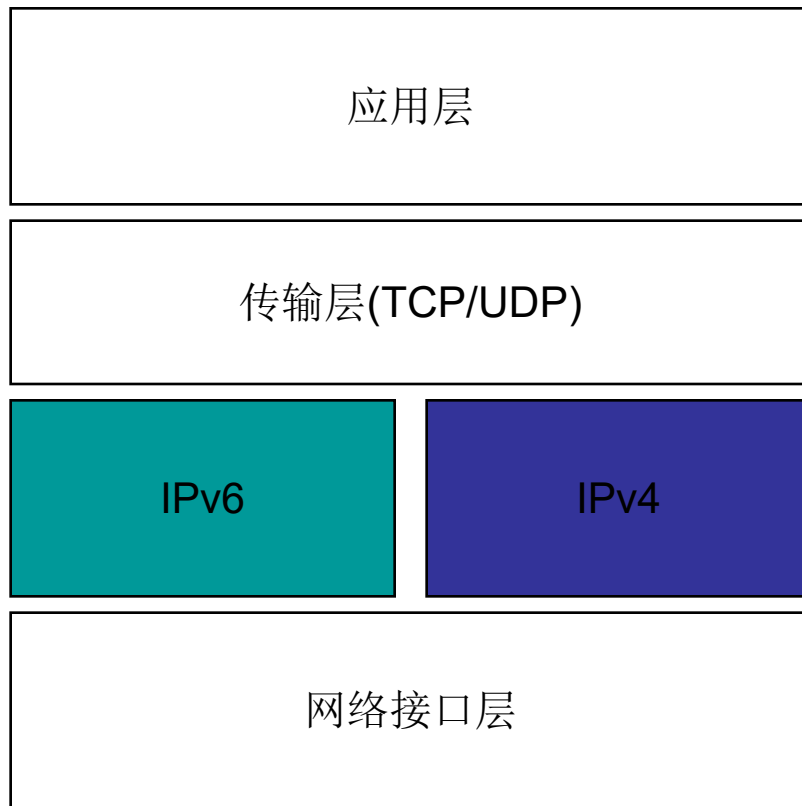


- IPv4-compatible address (兼容IPv4)
0:0:0:0:0:0:w.x.y.z or ::w.x.y.z
- IPv4-mapped address (不支持IPv6)
0:0:0:0:0:FFFF:w.x.y.z or ::FFFF:w.x.y.z
- 6over4 address
Interface ID of ::WWXX:YYZZ
- 6to4 address
Prefix of 2002:WWXX:YYZZ:: / 48
- ISATAP address
Interface ID of ::0:5EFE:w.x.y.z

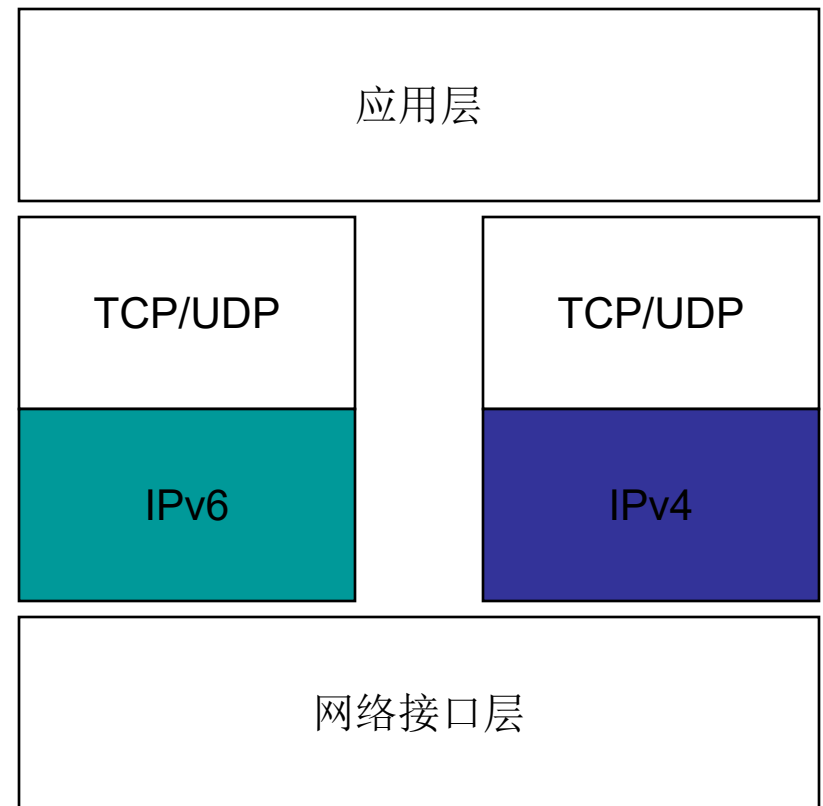
4、IPv6 层结构



双IP层架构



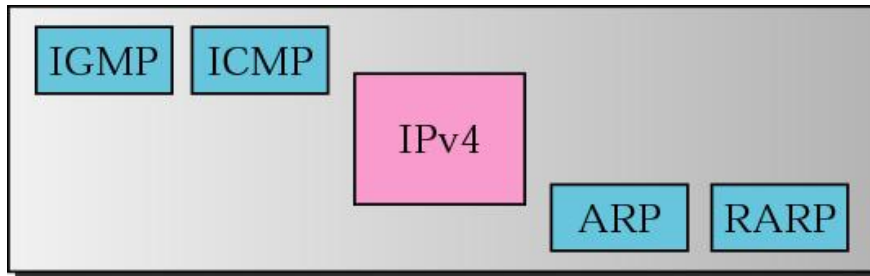
双堆栈结构



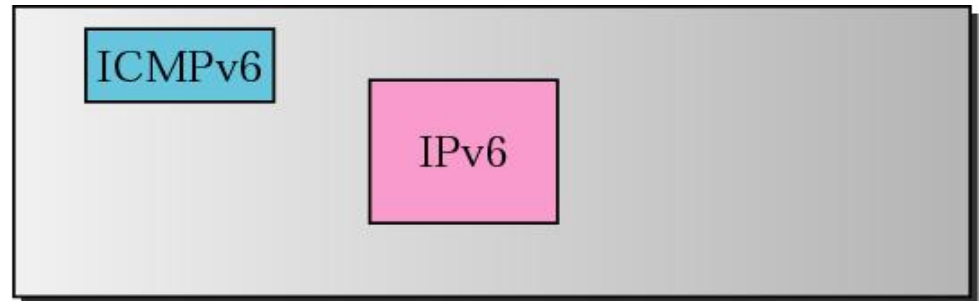
IPv6网络层



- v4和v6中的网络层比较



Network layer in version 4

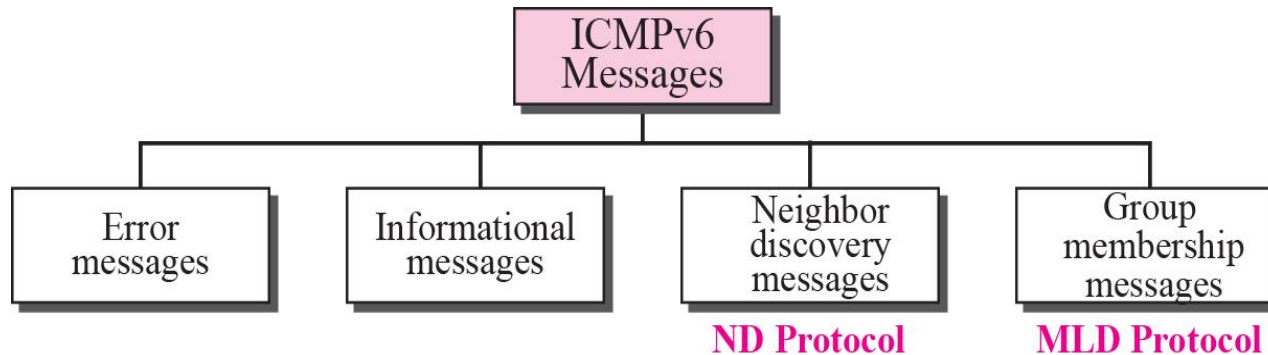


Network layer in version 6

ICMPv6



- 互联网控制消息协议第6版（ICMPv6）遵循第4版的相同策略和目的，比ICMPv4更复杂
 - 一些在IPv4的网络层中独立的协议现在是ICMPv6的一部分
 - 增加了一些新的消息，使其更加有用
- ICMPv6消息的分类法

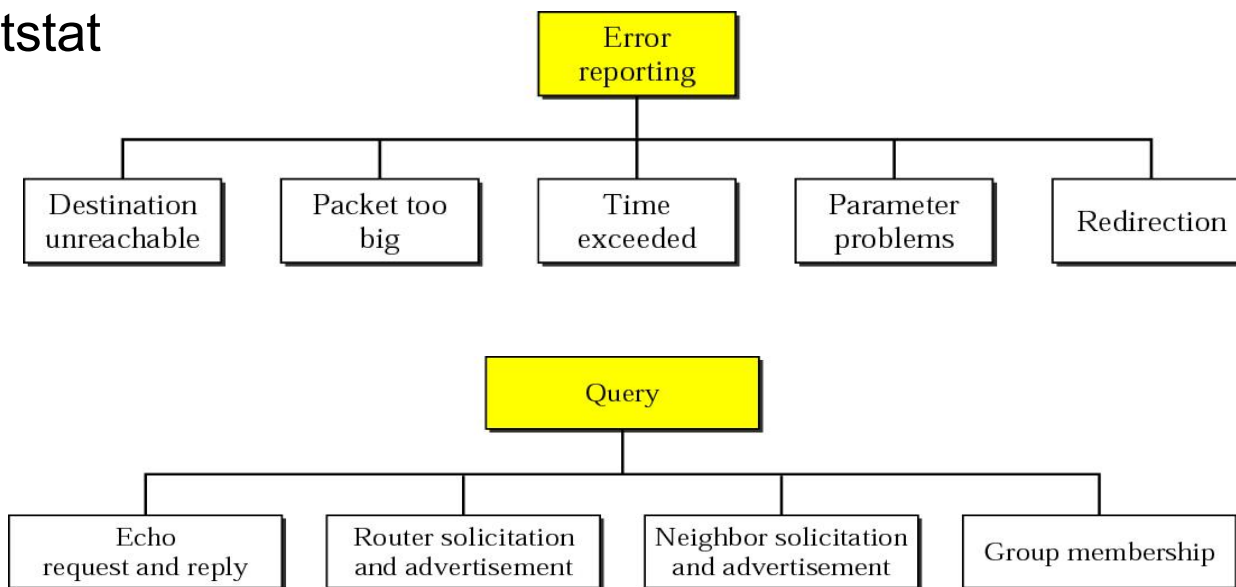


ICMPv6



- 支持IPv6的实用程序

1. 组播收听发现功能
2. 路由
3. Ping
4. 追踪器
5. Pathping
6. Netstat

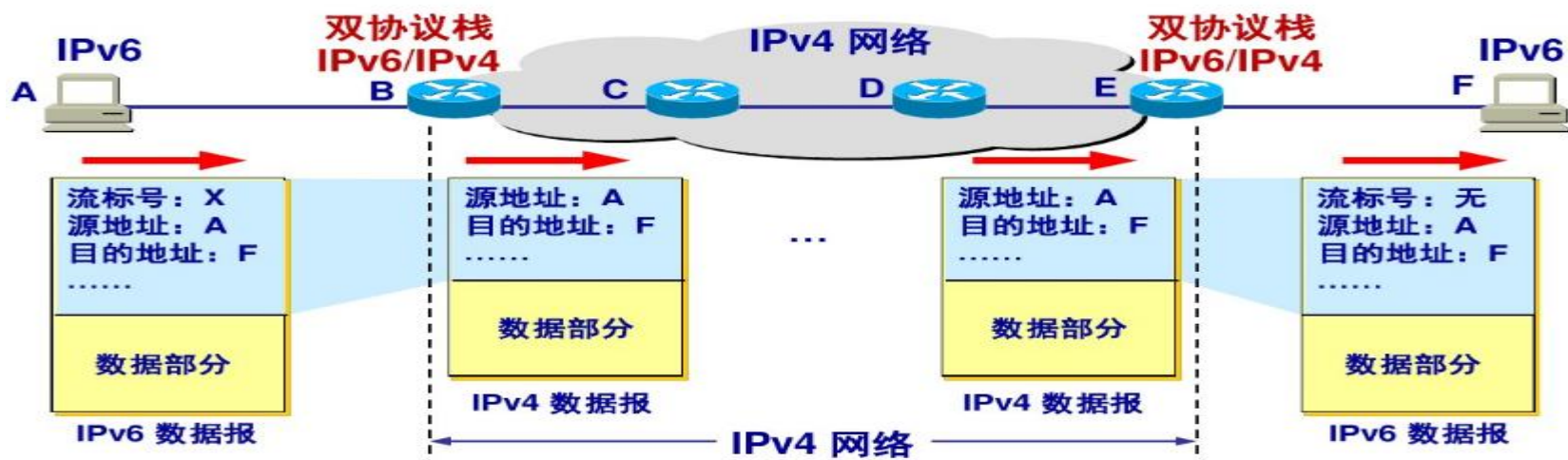


5、从V4到V6的过渡



1. 双栈操作 - v6节点在v4和v6模式下运行，并使用版本字段来决定使用哪个协议栈

- 节点可以被分配一个v4兼容的v6地址
 - 允许支持v6的主机谈论v6，即使本地路由器只谈论v4
 - 标志着对隧道的需求
 - 在32位v4地址上添加96个0（零扩展）--例如：::10.0.0.1
- 节点可以被分配一个v4映射v6的地址
 - 允许同时支持v6和v4的主机与v4主机通信
 - 在v4地址上添加2个字节的1，然后将其余部分扩展为零--例如：::ffff:10.0.0.1



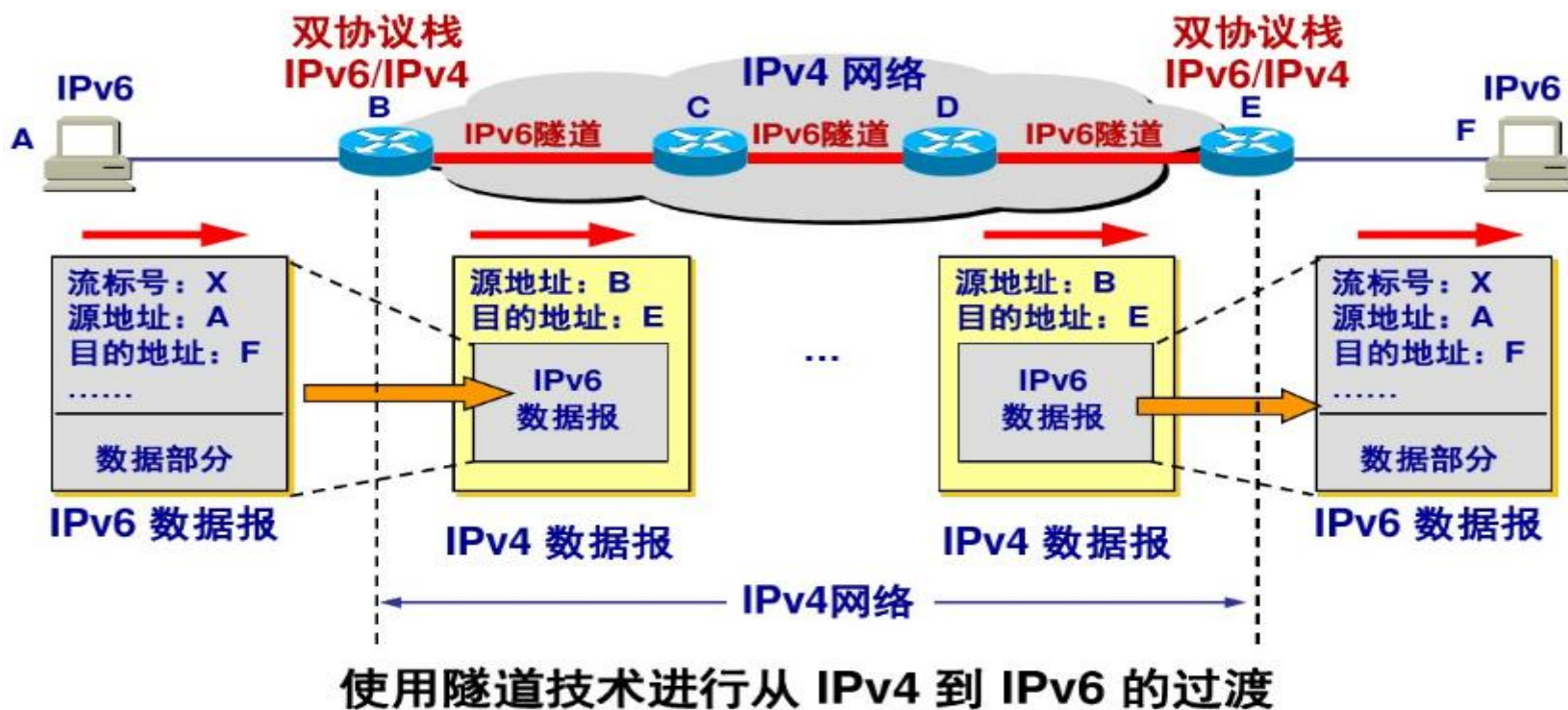
使用双协议栈进行从 IPv4 到 IPv6 的过渡

从V4到V6的过渡



2.隧道—用来处理两个v6路由器之间的v4路由器的网络

- 简单地将v6数据包和它们的所有信息封装在v4数据包中，直到你遇到下一个v6路由器





IPv6中的问题

- 地址长度：可用地址与开销的关系
- 跳数限制：是否需要65K？
- 最大数据包大小：更大的带宽要求更大的数据包
- 校验和是否必要？
- 服务器如何处理这两种类型的数据包？
- 在IP中安全是必要的吗？
如何最好地实现？
- DNS在过渡中可能非常重要--如何实现？



IPv6: 我国下一代互联网发展的机遇和挑战

- 互联网诞生于 20 世纪 60 年代，而我国于 1994 年才正式获准接入互联网，是互联网行业的后来者，在 IPv4 地址资源获取、互联网核心技术开发与标准制定方面都未能占得先机。
- 在 IPv6 新环境下，我国跟发达国家站在同一起跑线上，可以说 IPv6 为我国互联网的发展打开了一个创新空间，不仅解决了我国互联网发展中遇到的重大问题，还带来了新的机遇

逃离 IP 枯竭魔咒：



- 根据 CNNIC 2018年8月发布的第42次《中国互联网络发展状况统计报告》，截至2018年6月我国互联网用户数量已突破 8 亿，为全球最多的国家，而IPv4地址数量仅约3.39亿个。
- 2011年2月国际地址分配机构（IANA）在将其IPv4地址空间段的最后2 地址组分配出去，宣告地区性注册机构 (RIR)可用IPv4地址空间中“空闲池”的终结，亚洲地区由于分配地址少，是会最早出现地址枯竭的地区。
- 随着我国移动互联网、物联网、5G 等新技术新应用的不断推进，对IP地址的需求愈发旺盛 → IPv6的诞生就解决了这一问题！
- IPv6地址共有128位，可提供约340万亿个IP地址。因此，IPv6不仅能满足现有的地址需求，还可以为移动互联网、物联网乃至其他未知的全新业务留有地址空间。
- 我国也可积极参与IPv6 地址分配，争取更多地址资源以一劳永逸地解决地址短缺问题，这对于中国互联网经济的未来是决胜性的。

提高互联网生活与产业的安全和效率：



- 由于IPv4地址的极度缺乏，我国互联网用户目前实际使用的是运营商利用网络地址转换（NAT）生成的虚拟内部IP地址，即多个用户同时使用一个公有IP地址与外部通信。
- 这种方式会带来两大隐患：
 - 一方面，由于运营商并未为虚拟IP提供全面的端口映射，导致网络效率低下，设备互联、下载速度慢；
 - 另一方面，NAT 破坏了端对端的透明性，使个人用户无法与公有 IP 地址一对一匹配，如果发生网络安全事件将无法追溯到源地址。
- 而IPv6有海量可分配地址，所有设备都将可以有自己的单独地址，实现人一设备一IP地址的对应关系，解决网络实名制和用户身份溯源问题，实现网络精准管理。
- 此外，IPv6 的地址分配一开始就遵循聚类的原则，这使得路由器能在路由表中用一条记录表示一片子网，大大减小了路由器中路由表的长度，提高了路由器转发数据包的速度，从而提高了网络速度。

IPv6 发展仍任重道远：



- 总体上看，我国IPv6有较好的研究基础，良好的政策条件，初步形成了网络、应用和终端协同推进、齐头并进的良好发展局面，但与国际上 IPv6 发展较快的国家相比，还存在比较大差距，需要在后续工作中持续大力推进。
 - 政府已经主动发力，推进 IPv6 国家发展战略落实。我国发布的《推进互联网协议第六版（IPv6）规模部署行动计划》明确推进IPv6部署重要意义，提出部署的总体要求和主要目标，并从[互联网应用](#)、[网络和应用基础设施](#)、[网络安全](#)和[关键前沿技术](#)五方面安排实施的步骤，是我国全面推动IPv6 部署的指导文件。
 - 相关政府部门需结合行业与地方实际，制定实施方案，提出具体举措，并在落实中主动发现并解决存在的问题，成为 IPv6 部署的第一推动力。

积极开展标准制定，参与国际规则博弈：



- 我国IPv6国内外标准化已经有一定的基础，在国内标准方面，标准化工作于2001年启动，目前已发布相关邮电行业标准一百余项，涉及资源、网络、应用等各个领域；在IPv6国际标准方面，我国积极参与IETF IPv6国标准制定工作，并主导制定多项IPv6领域RFCs标准。下一步，我国应一方面根据国内外IPv6发展现状，持续跟进国内外标准制定，规范IPv6行业发展，加强国际交流，争取国际话语权；另一方面需助力标准体系落地，推动标准应用实施，使标准能够充分应用于IPv6行业，促进行业健康有序发展。

如何摆脱西方对我国互联网发展的制约



- 域名服务器是构架互联网所必需的关键基础设施，负责互联网最顶级的域名解析，被称为互联网的“中枢神经”。在IPv4体系内，美国利用先发优势主导的根服务器治理体系已延续近30年，全球13台根服务器中就有10台位于美国，其余3台分布在英国、瑞典和日本。我国互联网不得不依赖于西方国家控制的根服务器，如发现紧急情况（如2014年1月21日国内过半网站DNS瘫痪事件），缺少应变和反制手段，这对于我国互联网安全乃至国家安全都带来了巨大的隐患。IPv6部署意味着需要部署新的根服务器，中国也将有机会拥有自己的根服务器，逃离某天突然爆发的IP战困局。

如何提升我国在互联网发展中的话语权



- 现如今，标准已从传统意义上的产品互换和质量评判的依据上升为产业整体发展战略的重要组成部分，成为事关产业发展的基础性、先导性和战略性工作。我国是IETF标准化工作的后来者，IPv4 领域的标准起草绝大多数都是由西方国家主导的，我国主导的 IPv4 领域标准仅有个位数。因此在 IPv4 领域，我国不得不遵从其他国家的标准要求，发展严重受限。截至2018年8月，IETF共发布RFC标准近8500个，其中由中国主导起草虽然只有100 个左右，但几乎都是IPv6领域的标准。我国在IPv6 领域标准的制定中已经跟上了国际社会的步伐，IPv6 时代是我国互联网协议相关标准制定方面从跟跑到领跑的绝好机会。

如何激发企业积极性，实现技术安全和可控



- 互联网领域核心技术受制于人一直是制约我国互联网健康发展的瓶颈问题。随着互联网规模的不断扩大，网络和信息系統面临的问题和挑战越来越大，掌握核心技术的重要性越来越突出。应抓住全球互联网向 **IPv6** 转化这一历史机遇，实行“弯道超车”，在解决下一代互联网面临的重大技术挑战、特别是安全可信和自主可控的**IPv6** 技术体系方面，抢得国际话语权和主动权，力争使我国在基于**IPv6** 的下一代互联网研究上走在世界的前列。



中办、国办《推进 IPv6 规模部署行动计划》

2017 ~ 2018

应用系统：

- 国内用户量排名前50位的商业网站及应用，
- 省部级以上政府和中央企业外网网站系统，
- 中央和省级新闻及广播电视媒体网站系统，
- 工业互联网等新兴领域的网络与应用；
- 超大型互联网数据中心（IDC），
- 排名前5位的内容分发网络（CDN），
- 排名前10位云服务平台的50%云产品；

管道：

- 互联网骨干网、城域网和接入网，
- LTE网络及业务，
- 广电骨干网

终端：

- 新增网络设备、固定网络终端、移动终端。

2019 ~ 2020

应用系统：

- 新增网络地址不再使用私有IPv4地址，
- 国内用户量排名前100位的商业网站及应用，
- 市地级以上政府外网网站、新闻广播电视媒体网站系统；
- 大型互联网数据中心，
- 排名前10位的内容分发网络，
- 排名前10位云服务平台的全部云产品；

管道：

- 广电网络，
- 5G网络及业务，
- 国际出入口。

终端：

- 各类新增移动和固定终端，

2021 ~ 2025

- 到2025年末，我国IPv6网络规模、用户规模、流量规模位居世界第一位，
- 网络、应用、终端全面支持IPv6，
- 全面完成向下一代互联网的平滑演进升级，形成全球领先的下一代互联网技术产业体系。

到2018年末，市场驱动的良好发展环境基本形成，IPv6活跃用户数达到2亿，在互联网用户中的占比不低于20%

到2020年末，市场驱动的良好发展环境日臻完善，IPv6活跃用户数超过5亿，在互联网用户中的占比超过50%

到2025年末，我国IPv6网络规模、用户规模、流量规模位居世界第一位，形成全球领先的下一代互联网技术产业体系。