

## 4.4 因特网的路由选择协议

### 4.4.1 有关路由选择协议的几个基本概念

---

#### 1. 理想的路由算法

- 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。



# 关于“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。



# 从路由算法的自适应性考虑

- **静态**路由选择策略——即非自适应路由选择，其特点是简单和开销较小，但不能及时适应网络状态的变化。
- **动态**路由选择策略——即自适应路由选择，其特点是能较好地适应网络状态的变化，但实现起来较为复杂，开销也比较大。



## 2. 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。



# 自治系统 AS (Autonomous System)

- 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。



# 因特网有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。

# 自治系统和 内部网关协议、外部网关协议



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



## 这里要指出两点

- 因特网的早期 RFC 文档中未使用“路由器”而是使用“网关”这一名词。但是在新的 RFC 文档中又使用了“路由器”这一名词。应当把这两个属于当作同义词。
- IGP 和 EGP 是协议类别的名称。但 RFC 在使用 EGP 这个名词时出现了一点混乱，因为最早的一个外部网关协议的协议名字正好也是 EGP。因此在遇到名词 EGP 时，应弄清它是指旧的协议 EGP 还是指外部网关协议 EGP 这个类别。





# 因特网的路由选择协议

---

- 内部网关协议 IGP：具体的协议有多种，如 RIP 和 OSPF 等。
- 外部网关协议 EGP：目前使用的协议就是 BGP。

## 4.4.2 内部网关协议 RIP (Routing Information Protocol)

### 1. 工作原理

- 路由信息协议 RIP 是内部网关协议 IGP 中最先得到广泛使用的协议。
- RIP 是一种分布式的基于**距离向量**的路由选择协议。
- 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。



# “距离”的定义

---

- 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。
- 这里的“距离”实际上指的是“**最短距离**”，



# “距离”的定义

---

- 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。可见 RIP 只适用于小型互联网。
- 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。



# RIP 协议的三个要点

- 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。
- 为加快协议的收敛速度，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（即**触发更新**）。



# 路由表的建立

- 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。
- 一般情况下RIP 协议的**收敛**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。

## 2. 距离向量算法

收到相邻路由器（其地址为 X）的一个路由更新报文：

(1) 先修改此报文中的所有项目：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。每一个项目都有三个关键数据，即：到目的网络N，距离是d，下一跳路由器是X。

(2) 若原路由表中没有目的网络N，则把该项目添加到路由表中。否则，查看路由表中目的网络为N的表项，若其下一跳是X，则把收到的项目替换原项目。

否则，若收到的项目中的距离d小于路由表中的距离，则进行更新，否则什么也不做。

(3) 若180秒（默认）没有收到某条路由项目的更新报文，则把该路由项目记为无效，即把距离置为16（距离为16表示不可达），若再过一段时间，如120秒，还没有收到该路由项目的更新报文，则将该路由项目从路由表中删除。

(4) 若路由表发生变化，向所有相邻路由器发送路由更新报文。

(5) 返回。



# 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。



正常情况



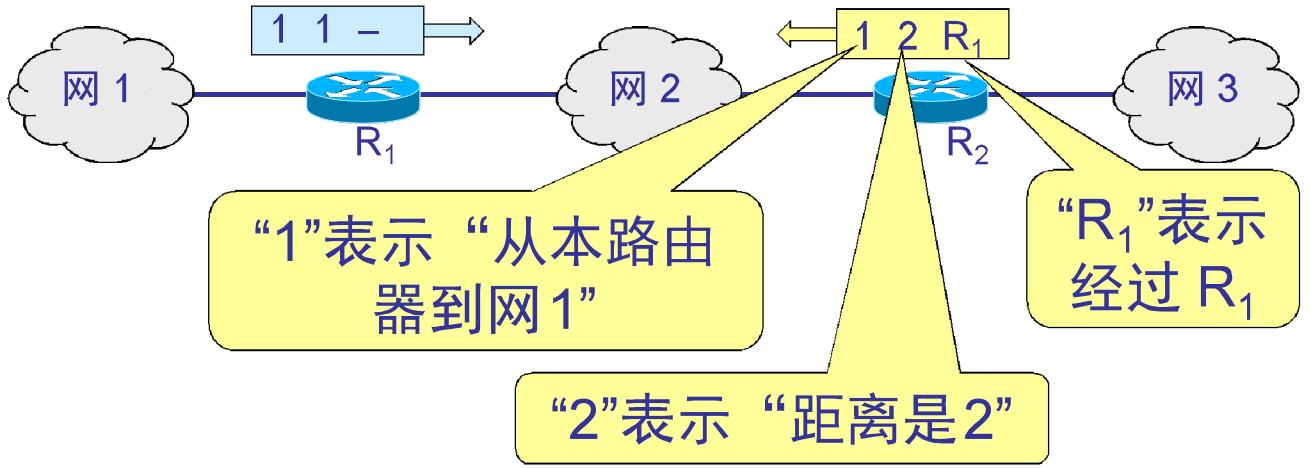
“1”表示“从本路由器到网1”

“-”表示“直接交付”

“1”表示“距离是1”

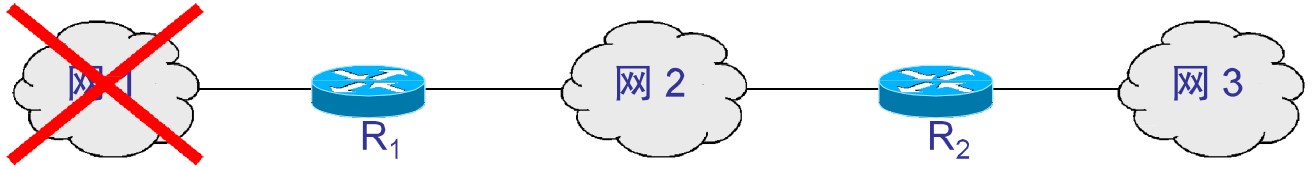
R<sub>1</sub> 说：“我到网1 的距离是1，是直接交付。”

正常情况



R<sub>2</sub> 说：“我到网 1 的距离是 2，是经过 R<sub>1</sub>。”

正常情况



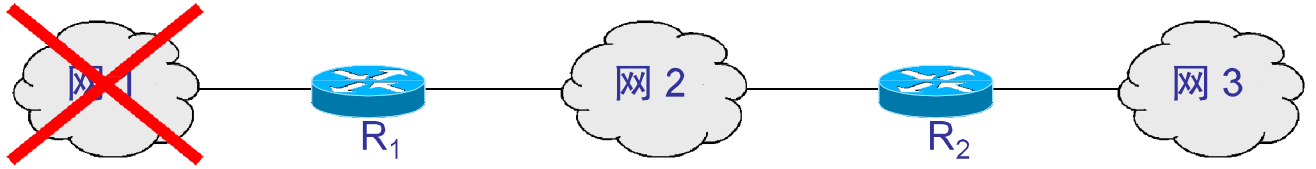
网1出了故障



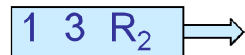
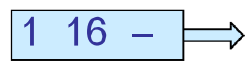
R<sub>1</sub> 说：“我到网1的距离是16（表示无法到达），是直接交付。”

但 R<sub>2</sub> 在收到 R<sub>1</sub> 的更新报文之前，还发送原来的报文，因为这时 R<sub>2</sub> 并不知道 R<sub>1</sub> 出了故障。

正常情况

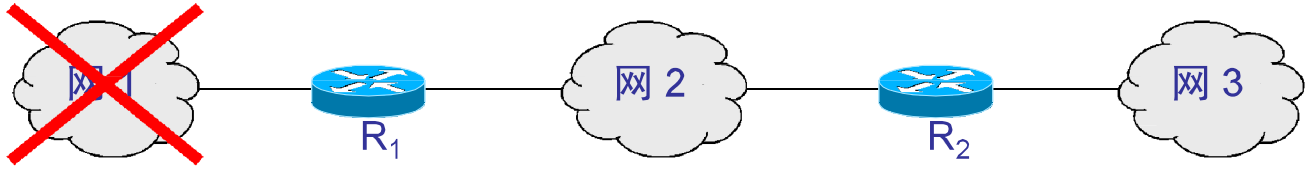


网 1 出了故障



R<sub>1</sub> 收到 R<sub>2</sub> 的更新报文后，误认为可经过 R<sub>2</sub> 到达网 1，于是更新自己的路由表，说：“我到网 1 的距离是 3，下一跳经过 R<sub>2</sub>”。然后将此更新信息发送给 R<sub>2</sub>。

正常情况



网 1 出了故障

1 16 -

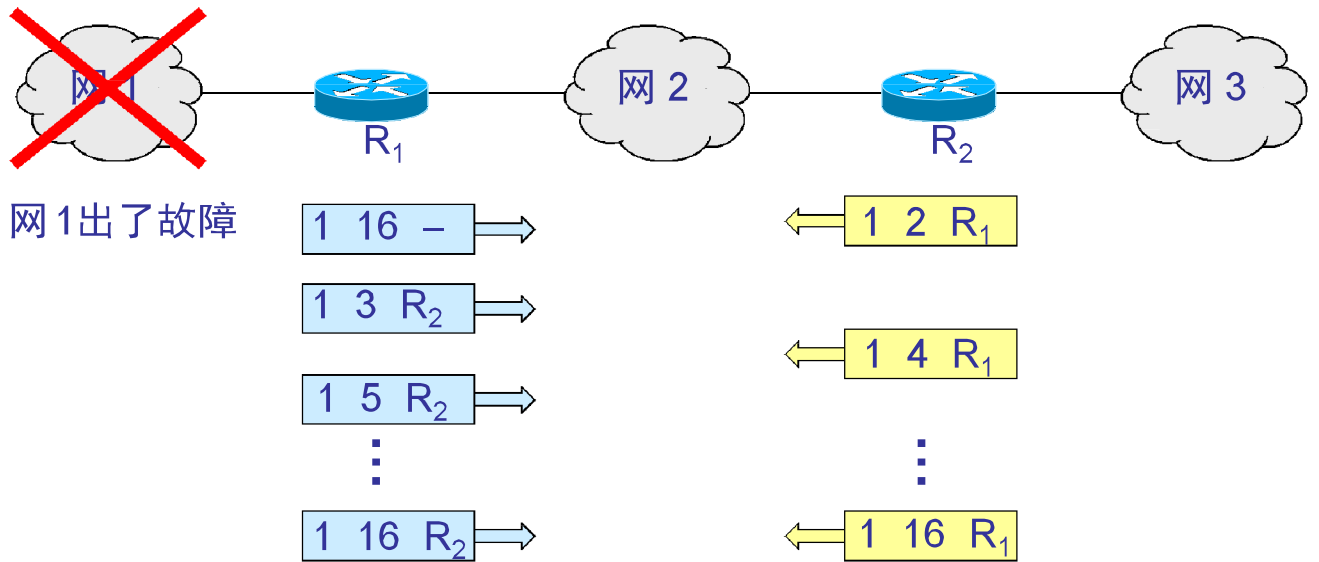
1 2 R<sub>1</sub>

1 3 R<sub>2</sub>

1 4 R<sub>1</sub>

R<sub>2</sub> 以后又更新自己的路由表为 “1, 4, R<sub>1</sub>”, 表明 “我到网 1 距离是 4, 下一跳经过 R<sub>1</sub>”。

这就是好消息传播得快，而坏消息传播得慢。网络出故障的传播时间往往需要较长的时间(例如数分钟)。这是 RIP 的一个主要缺点。



这样不断更新下去，直到  $R_1$  和  $R_2$  到网 1 的距离都增大到 16 时， $R_1$  和  $R_2$  才知道网 1 是不可达的。



# 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。

## 4.4.3 内部网关协议 OSPF (Open Shortest Path First)

1. OSPF 协议的基本特点
  - “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”。
  - 是分布式的**链路状态协议**。





# 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。

# 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。



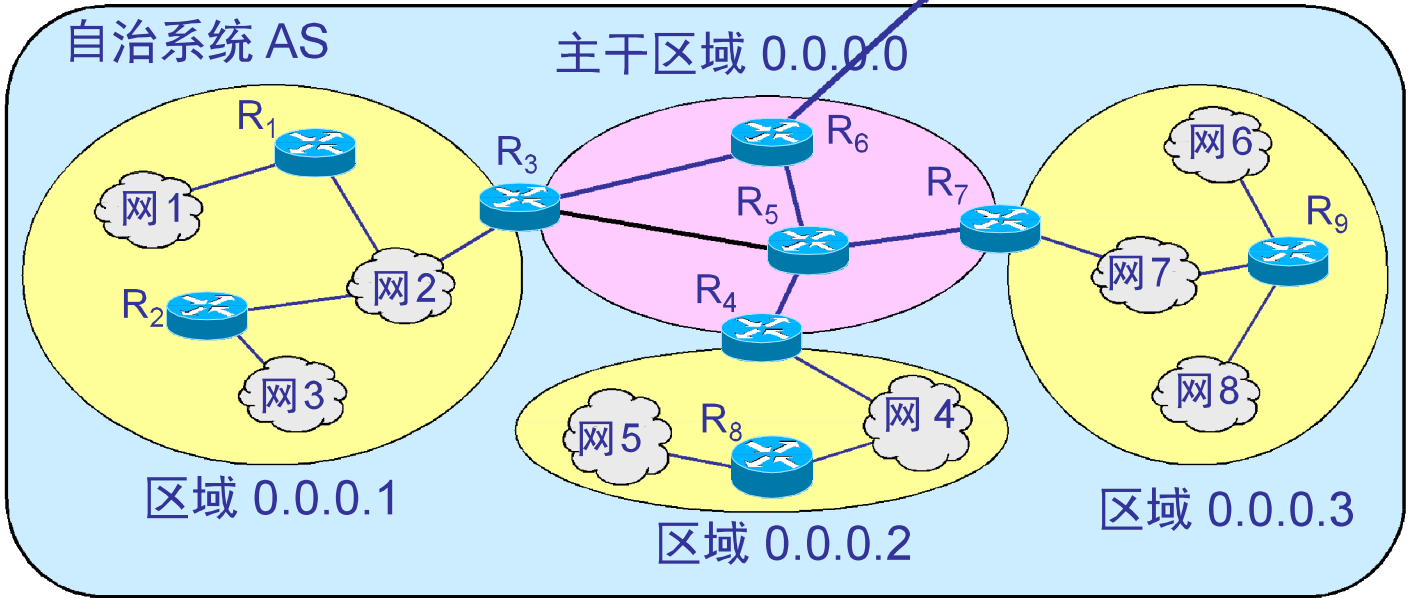
# OSPF 的区域(area)

---

- 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。

# OSPF 划分为两种不同的区域

至其他自治系统



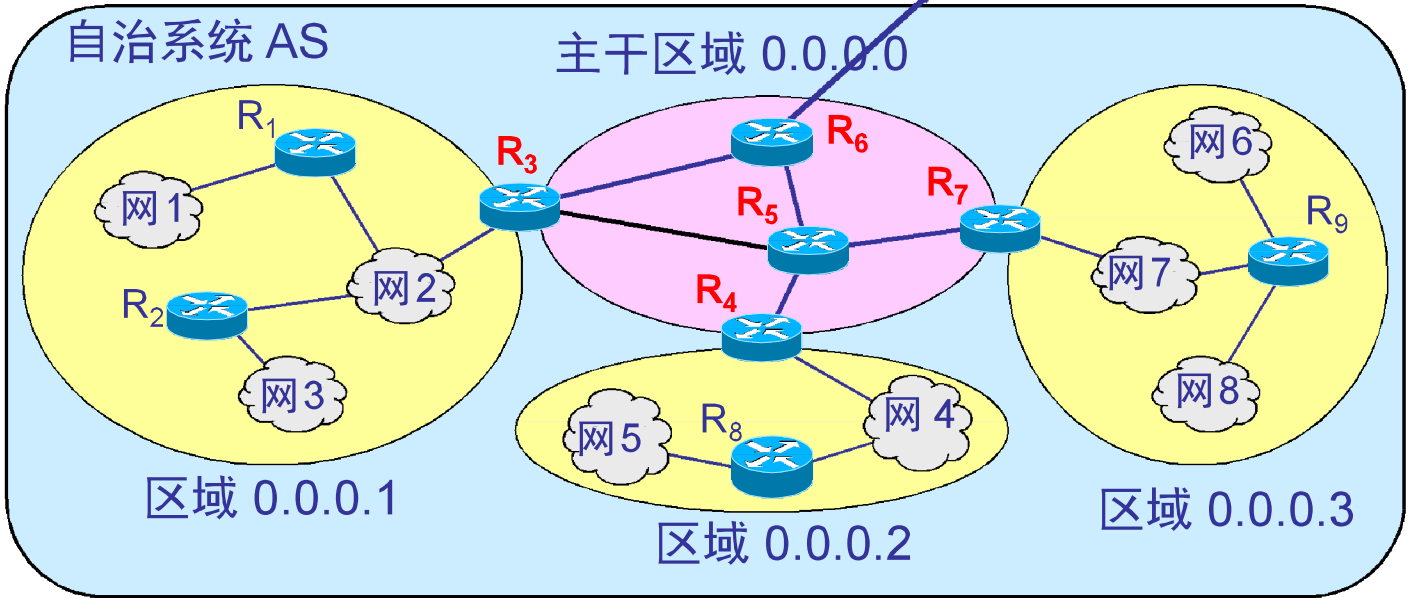


# 划分区域

- 划分区域的好处就是将利用洪泛法交换链路状态信息的范围局限于每一个区域而不是整个的自治系统，这就减少了整个网络上的通信量。
- 在一个区域内部的路由器只知道本区域的完整网络拓扑，而不知道其他区域的网络拓扑的情况。
- OSPF 使用层次结构的区域划分。在上层的区域叫作**主干区域**(backbone area)。主干区域的标识符规定为0.0.0.0。主干区域的作用是用来连通其他在下层的区域。

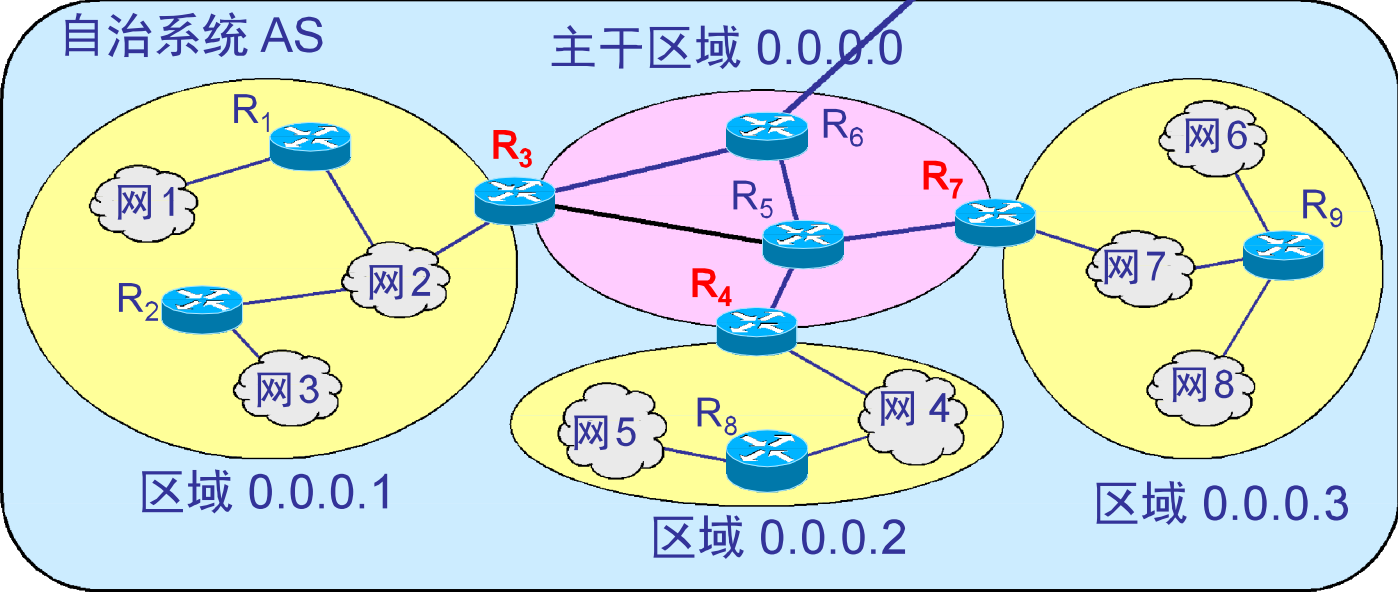
# 主干路由器

至其他自治系统



# 区域边界路由器

至其他自治系统





# OSPF 直接用 IP 数据报传送

- 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。





# OSPF 的其他特点

- 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。

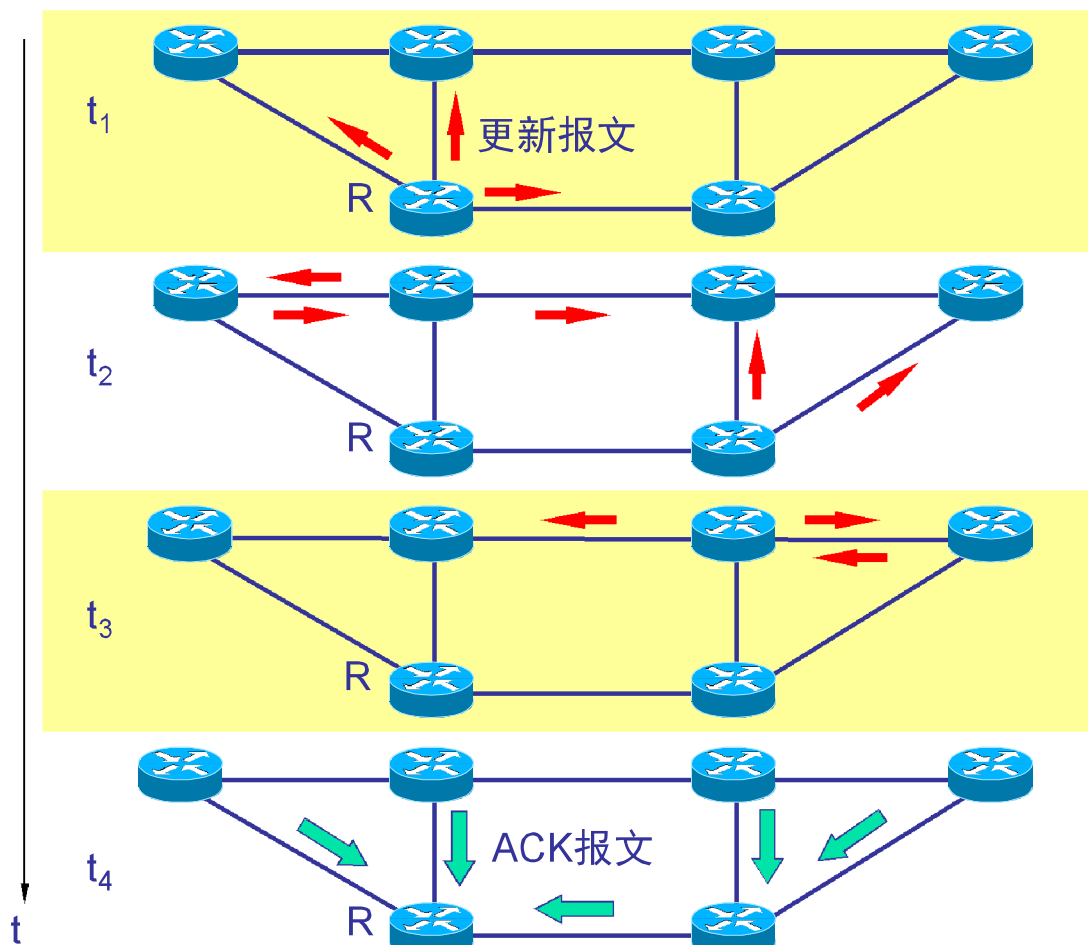


## 2. OSPF 的五种分组类型

---

- 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。

# OSPF 使用的是可靠的洪泛法





## OSPF 的其他特点

---

- OSPF 还规定每隔一段时间，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 100 ms。



# 指定的路由器 (designated router)

- 多点接入的局域网采用了指定的路由器的方法，使广播的信息量大大减少。
- 指定的路由器代表该局域网上所有的链路向连接到该网络上的各路由器发送状态信息。



## 4.4.4 外部网关协议 BGP

---

- 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。
- 可以将 BGP-4 简写为 BGP。



# BGP 使用的环境却不同

- 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。
- 因此，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**。



# 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。



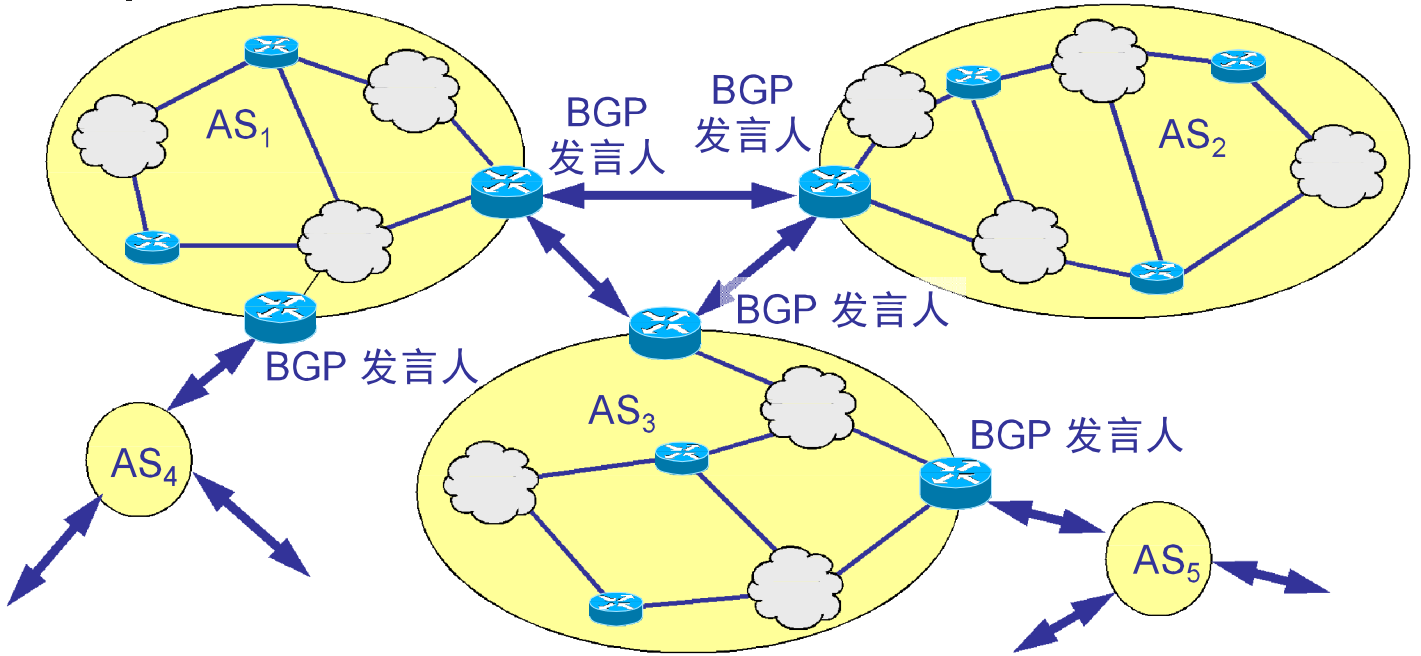


# BGP 交换路由信息

---

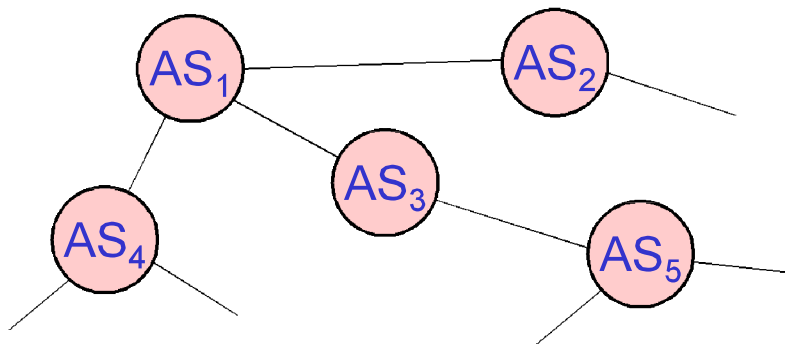
- 一个 BGP 发言人与其他自治系统中的 BGP 发言人要交换路由信息，就要先建立 TCP 连接，然后在此连接上交换 BGP 报文以建立 BGP 会话(session)，利用 BGP 会话交换路由信息。
- 使用 TCP 连接能提供可靠的服务，也简化了路由选择协议。
- 使用 TCP 连接交换路由信息的两个 BGP 发言人，彼此成为对方的邻站或对等站。

# BGP 发言人和 自治系统 AS 的关系



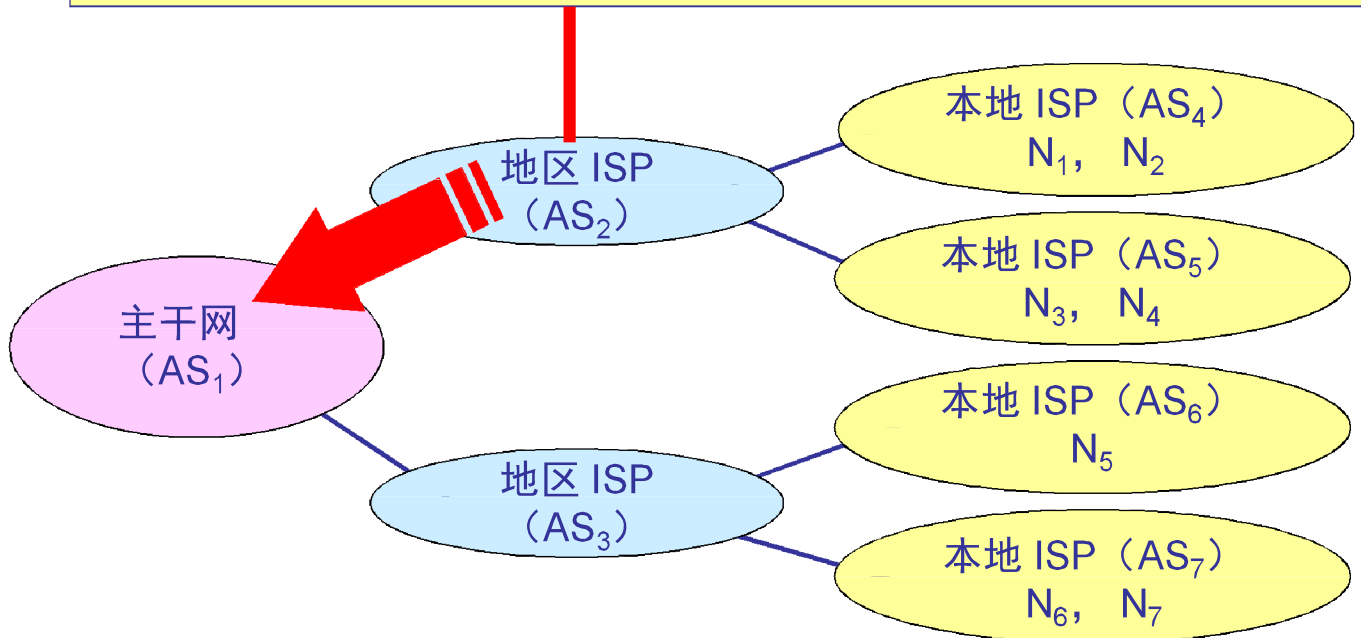
# AS 的连通图举例

- BGP 所交换的网络可达性的信息就是要到达某个网络所要经过的一系列 AS。
- 当 BGP 发言人互相交换了网络可达性的信息后，各 BGP 发言人就根据所采用的策略从收到的路由信息中找出到达各 AS 的较好路由。



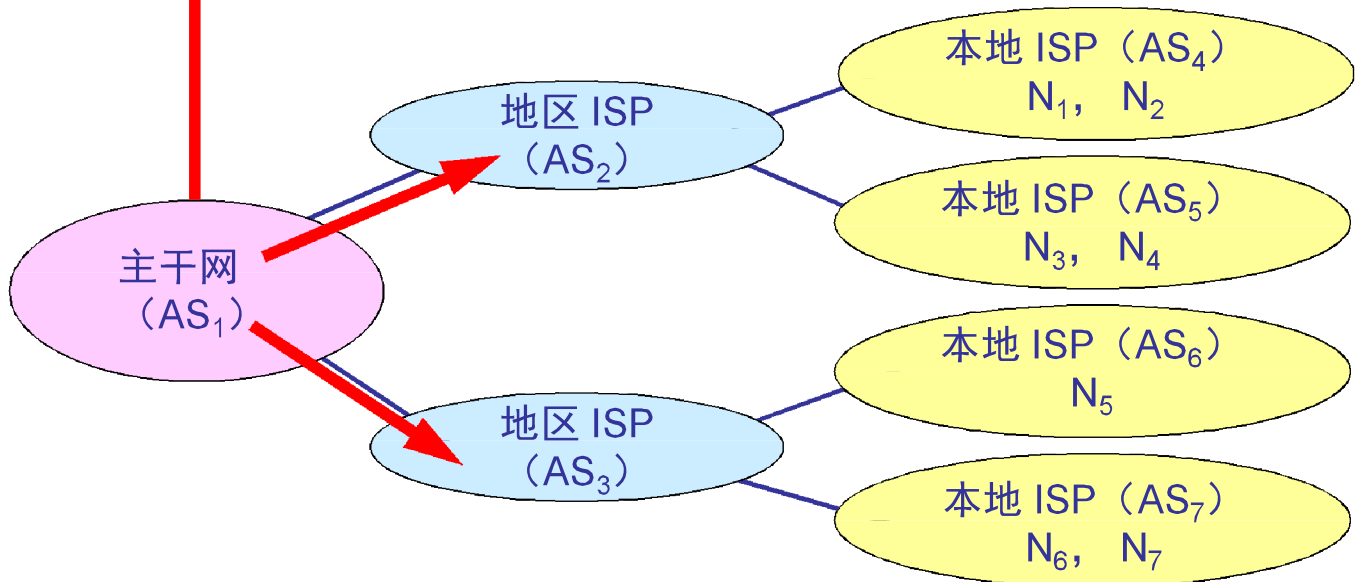
# BGP 发言人交换路径向量

自治系统  $AS_2$  的 BGP 发言人通知主干网的 BGP 发言人：“要到达网络  $N_1, N_2, N_3$  和  $N_4$  可经过  $AS_2$ 。”



# BGP 发言人交换路径向量

主干网还可发出通知：“要到达网络  $N_5$ ,  $N_6$  和  $N_7$  可沿路径  $(AS_1, AS_3)$ 。”





# 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 协议的特点

- 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。



# BGP-4 共使用四种报文

- (1) 打开(**OPEN**)报文，用来与相邻的另一个BGP发言人建立关系。
  - (2) 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。
  - (3) 保活(**KEEPALIVE**)报文，用来确认打开报文和周期性地证实邻站关系。
  - (4) 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。

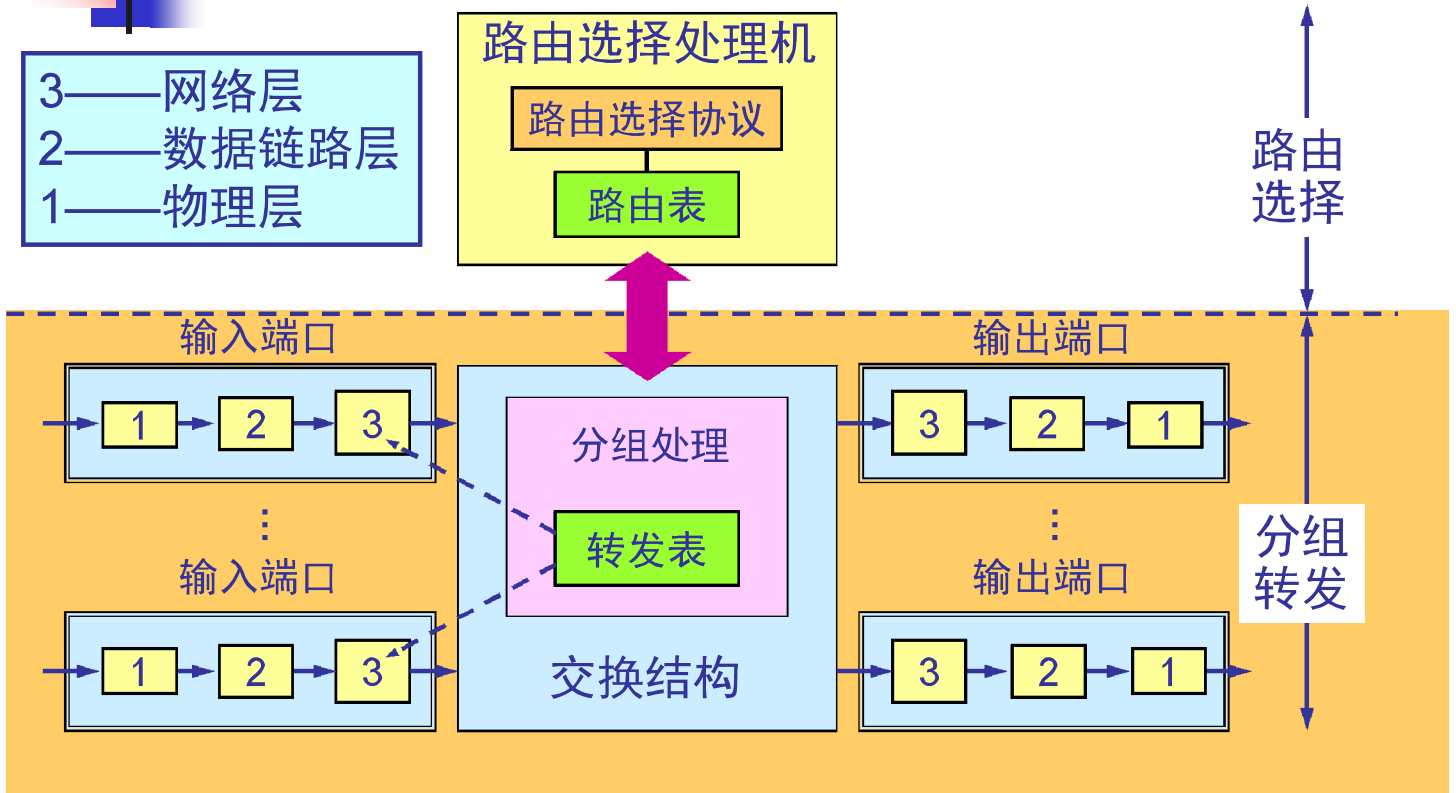


## 4.5 路由器的工作原理

### 4.5.1 路由器的构成

- 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。

# 典型的路由器的结构



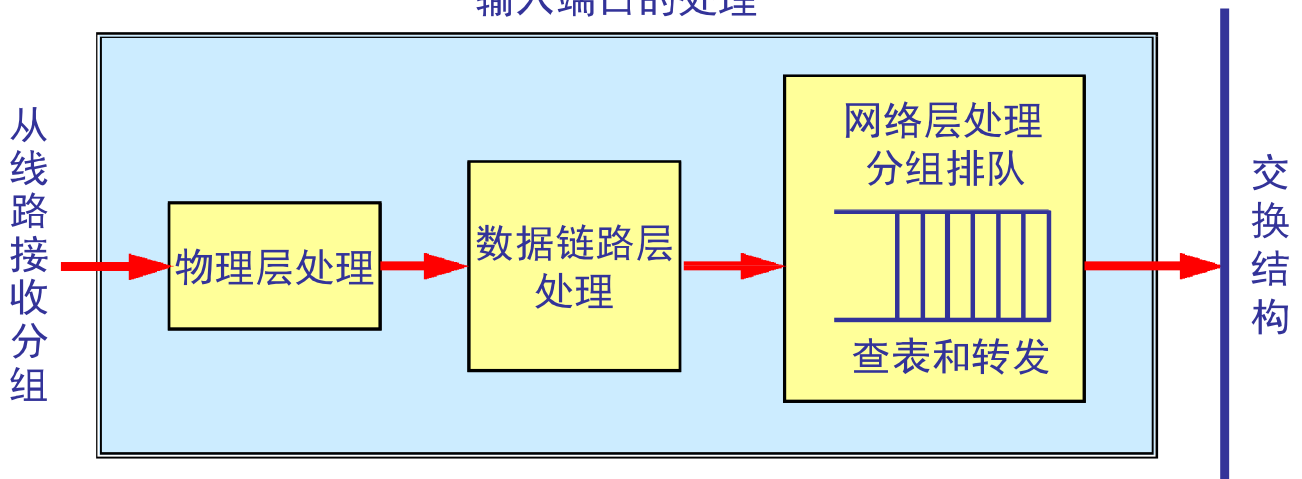
# “转发”和“路由选择”的区别

- “**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，

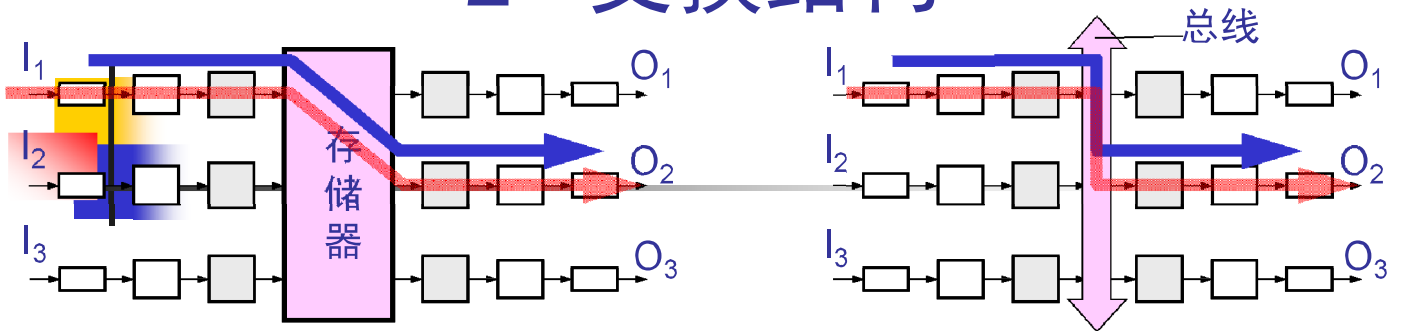
若交换结构处理分组的速率赶不上分组进入队列的速率，则会导致输入队列排队！

- 数据链路层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

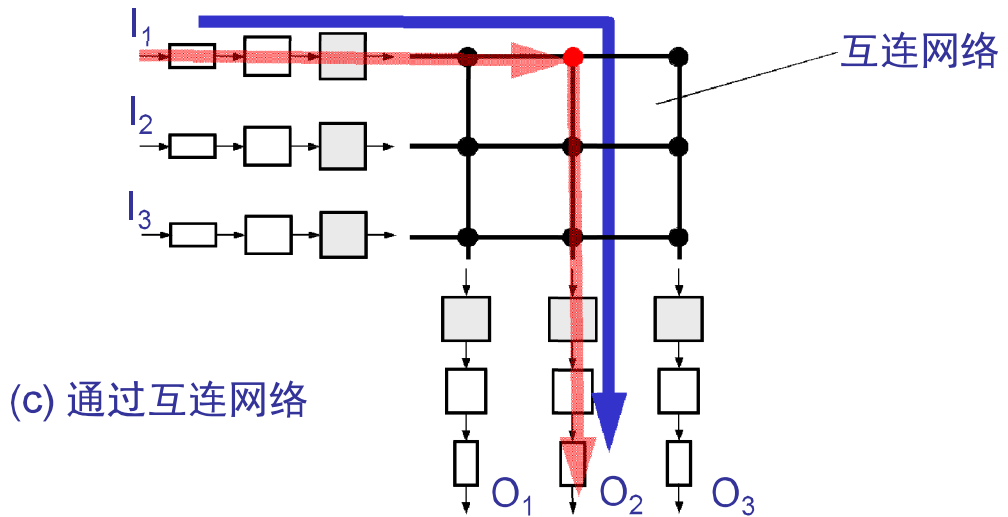


## 2 交换结构



(a) 通过存储器

(b) 通过总线



(c) 通过互连网络

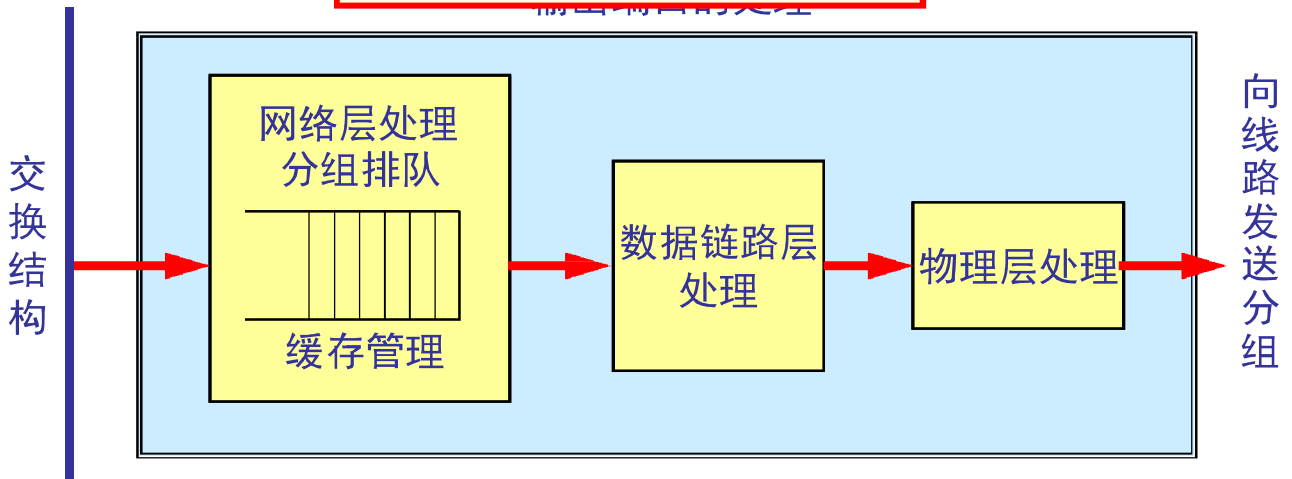
若从交换结构到达队列的分组速率超过输出链路的带宽，则会

- 当交换链路层部，交

如果路由器足够快是不是就不会出现排队？

据尾

相层的目的地和链路。





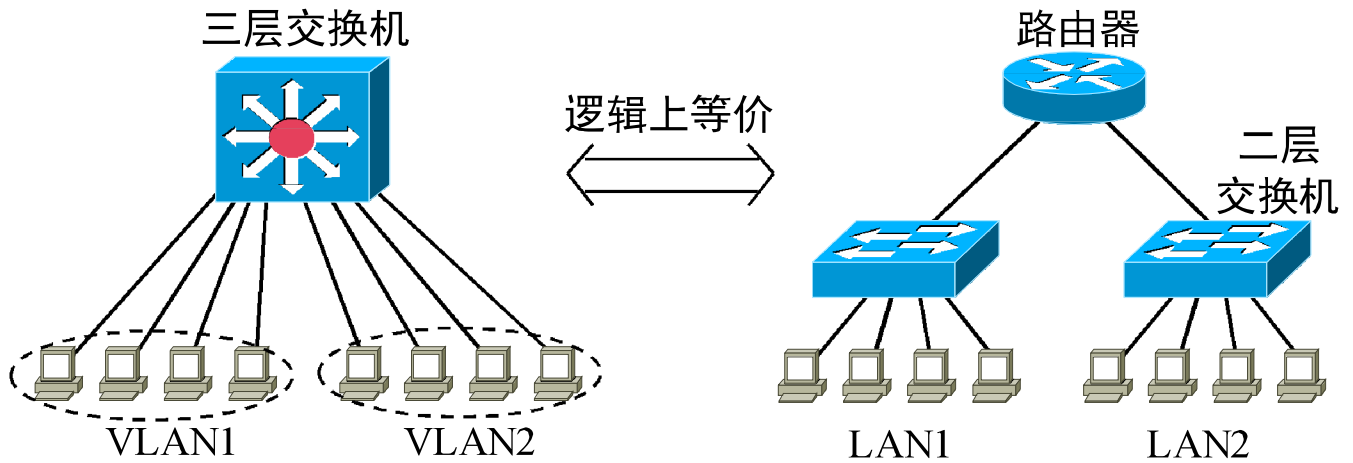
## 4.5.2 路由器与交换机的比较

- 交换机的**优点**是即插即用，并具有相对高的分组过滤和转发速度。缺点是：大型交换机网络要求交换机维护大的转发表，主机中维护大的ARP表，并可能产生广播风暴，逻辑拓扑被限制为树。
- 路由器的优点是能提供更加智能的路由选择，并能隔离广播域。缺点是：路由器不是即插即用的，对每个分组处理时间通常比交换机更长。

## 4.5.3 三层

为避免混淆，我们使用术语“**路由器**”而不使用术语“**三层交换机**”！

- “三层交换机”器和支持VLAN的二层交换机的集成体







# 三层交换机的应用

---

典型的做法是：

- 处于同一个局域网中的各个子网的互连以及局域网中VLAN间的路由，用三层交换机来代替普通路由器，实现广播域的隔离。
- 只有局域网与广域网互连，或广域网之间互连时才使用普通路由器。