

# 第六章 传输层

6.1 传输层服务

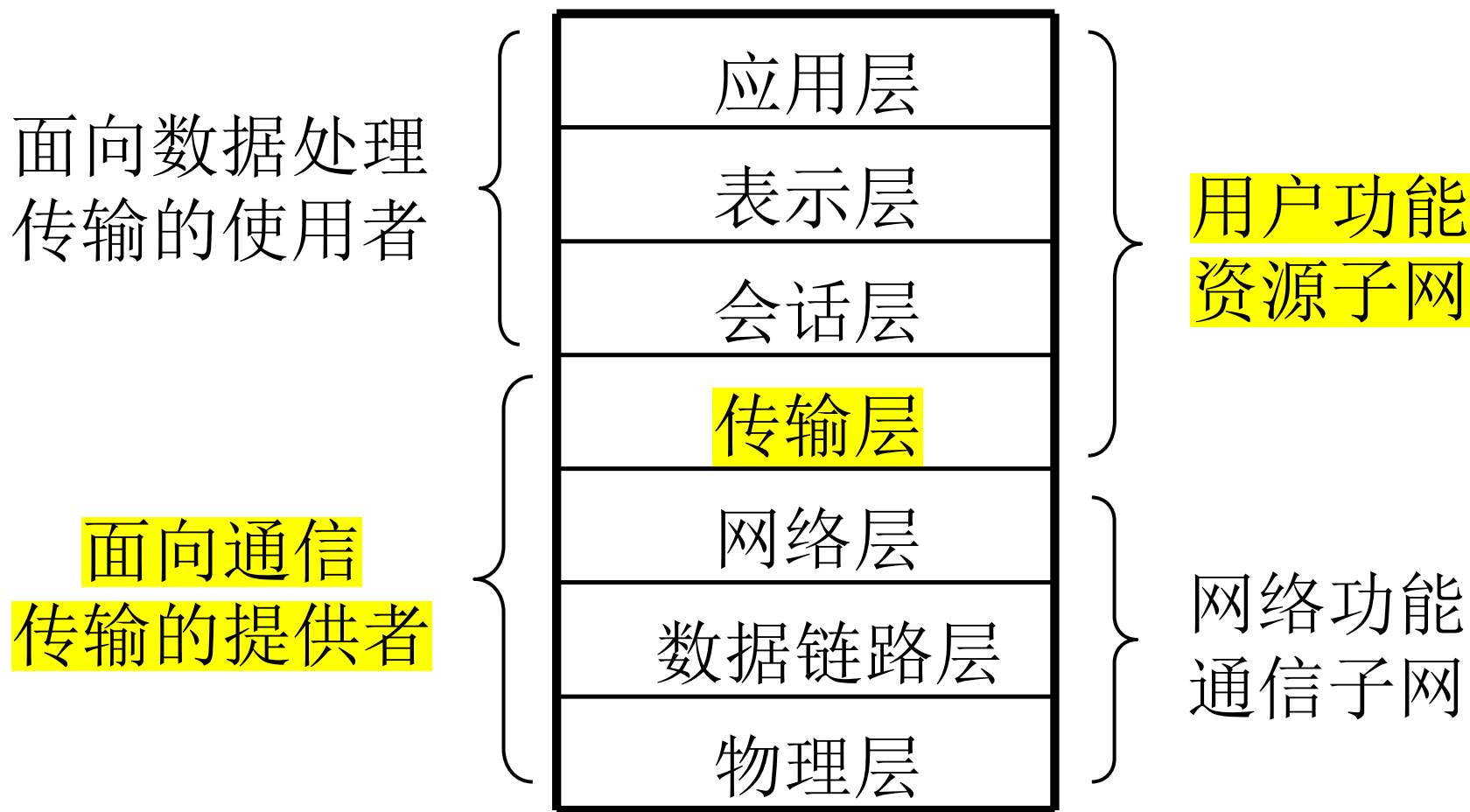
6.2 用户数据报协议

6.3 传输控制协议

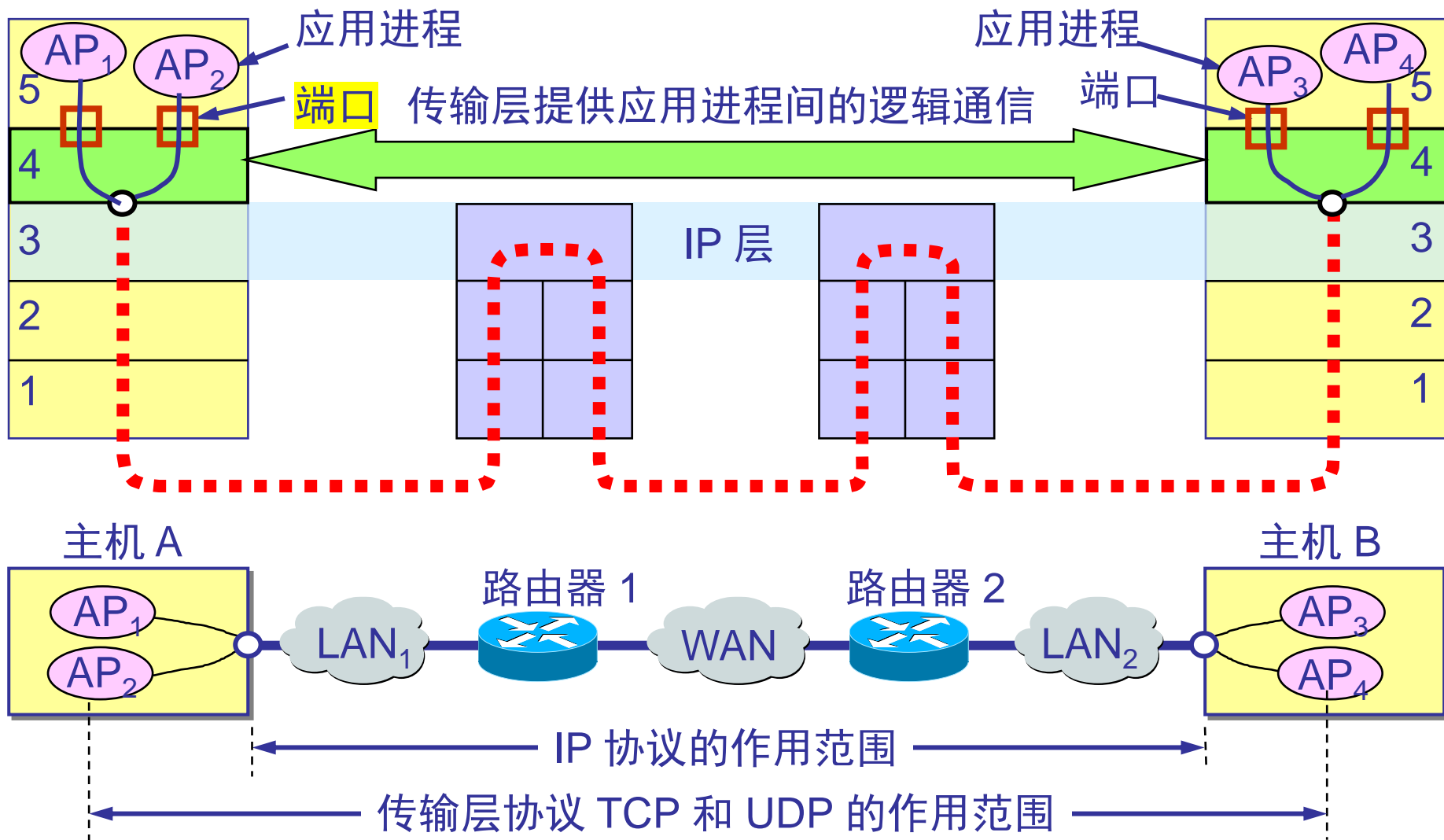
## 6.1 传输层服务

- 在OSI 参考模型中，传输层位于通信子网和资源子网之间，是整个协议层次中最核心的一层。
- 传输层为源主机上的进程和目的主机上的进程之间提供可靠的透明数据传送，使高层用户在相互通信时不必关心通信子网实现的细节。

## 6.1.1 传输层功能与服务概述



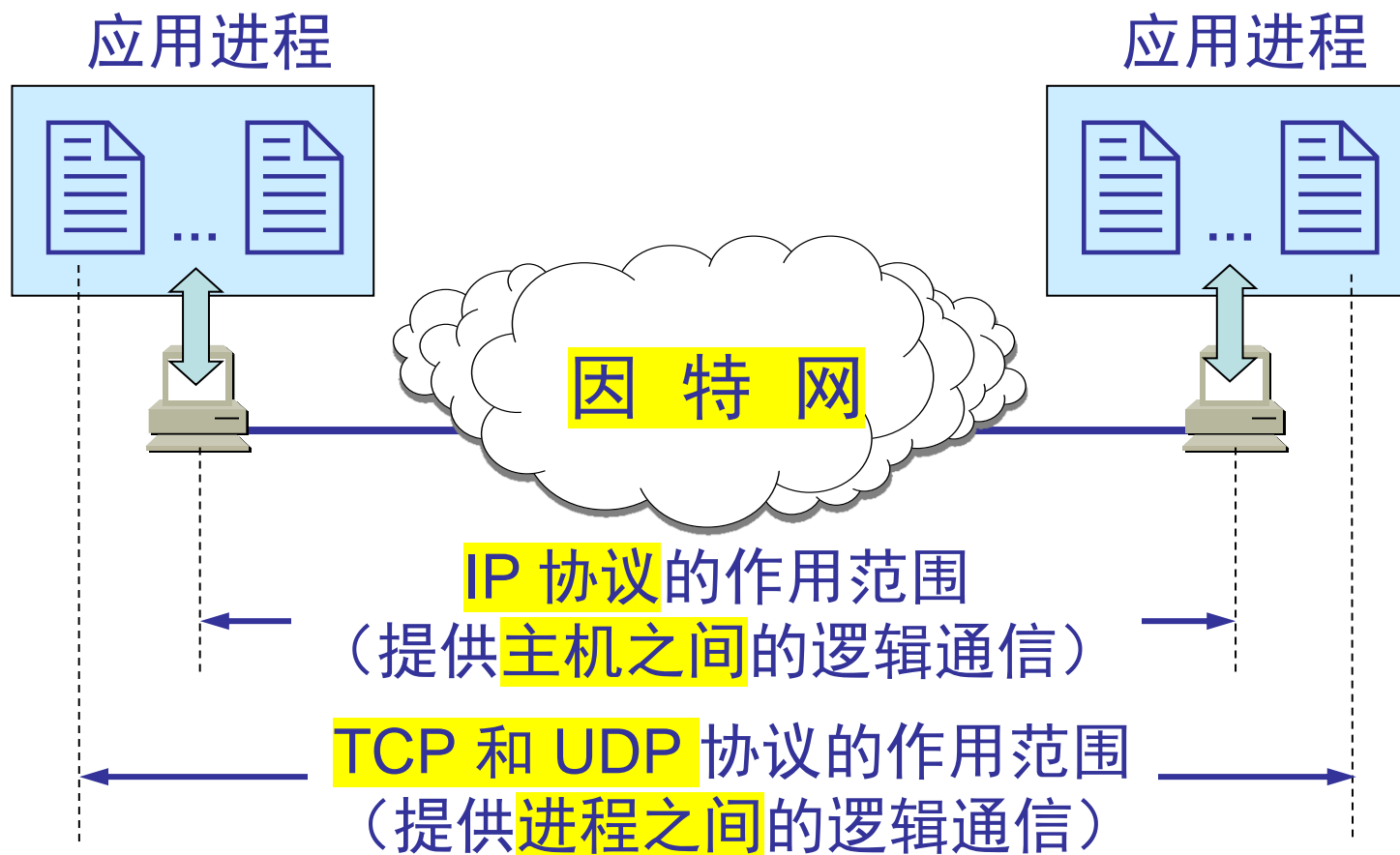
# 传输层提供逻辑通信



# 传输协议

- 传输服务是通过传输层实体间使用传输协议来实现的。
- 通信不仅仅是发生在从源计算机到目的计算机，而且是从端应用程序到端应用程序。

# 传输层和网络层协议的区别



# 传输层功能

## ■ 具体功能包括：

- 端到端的报文传递(把数据可靠地从一方用户进程送到另一方用户进程，传输层通信与通信子网工作方式无关)
- 服务点的寻址
- 拆分和组装
- 连接控制

## ■ 传输层为上层提供两种类型服务：

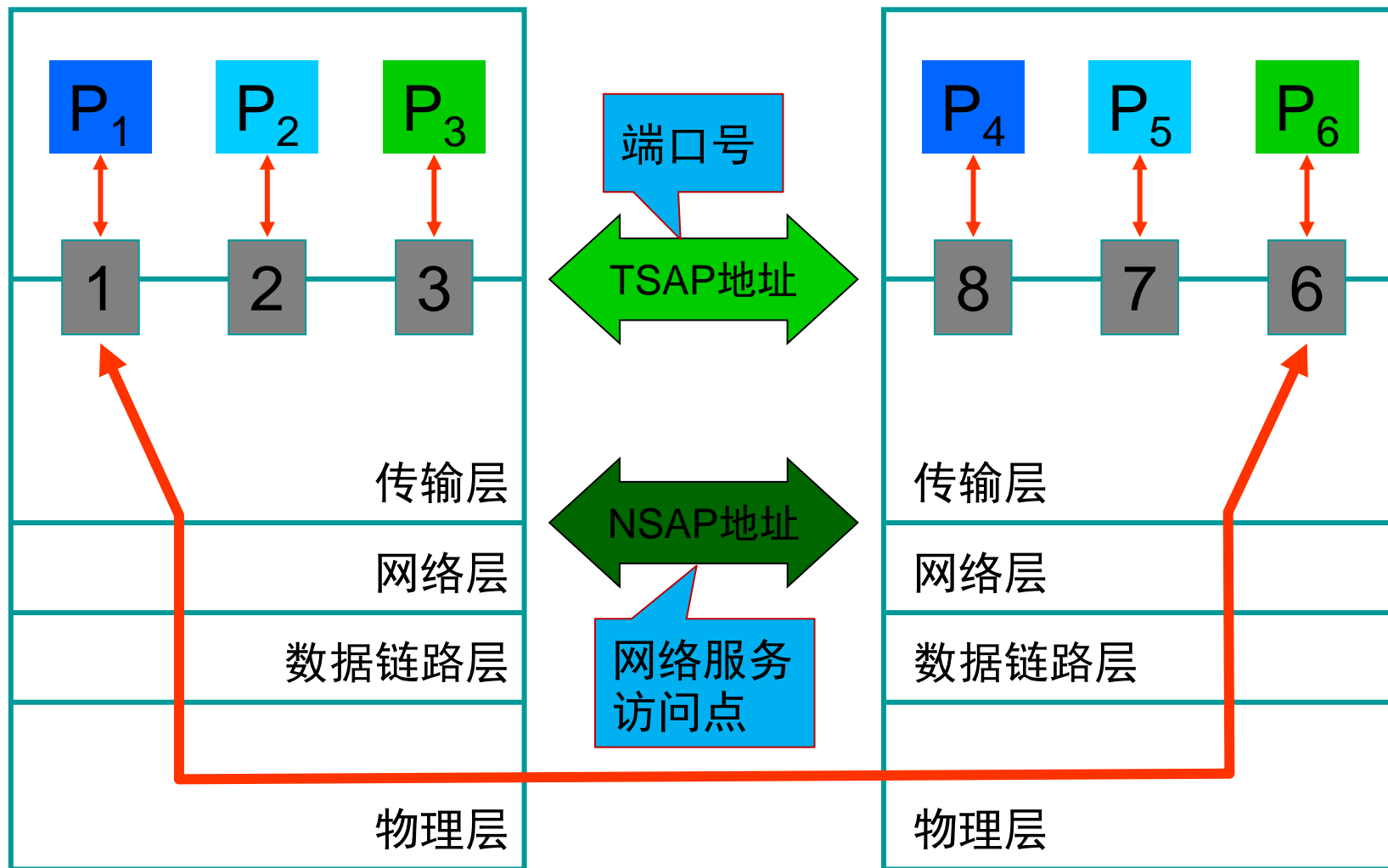
- 面向连接的传输服务
- 面向无连接的传输服务

## 6.1.2 传输层地址

- 由一台计算机上的应用程序所产生的数据不仅必须被另外一台计算机所接收，而且必须被这台计算机上正确的应用程序所接收。
- 一个应用程序同一个远程应用程序通信时，它必须知道两个地址：
  - TSAP地址(传输访问服务点，称为端口号)
  - NSAP地址(网络服务访问点)



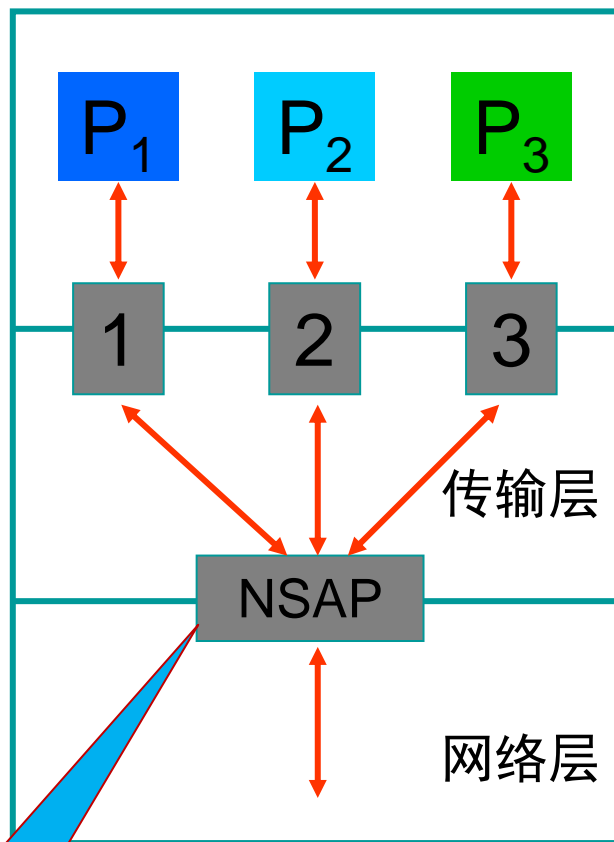
# 访问点通信



应用程序使用服务访问点通信

## 6.1.3 传输层复用

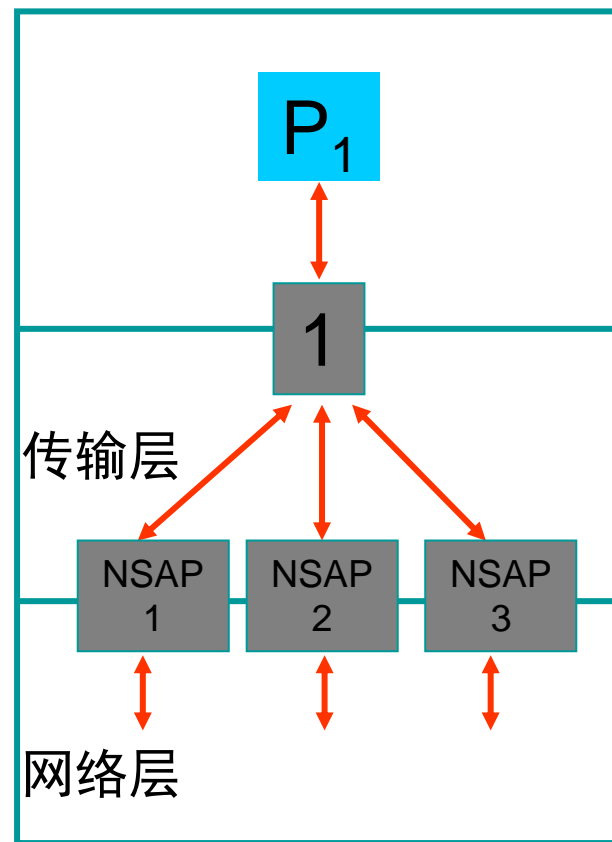
- 提高效率：向上复用、向下复用



向上复用

层网络结构是虚电路

网络服务  
访问点

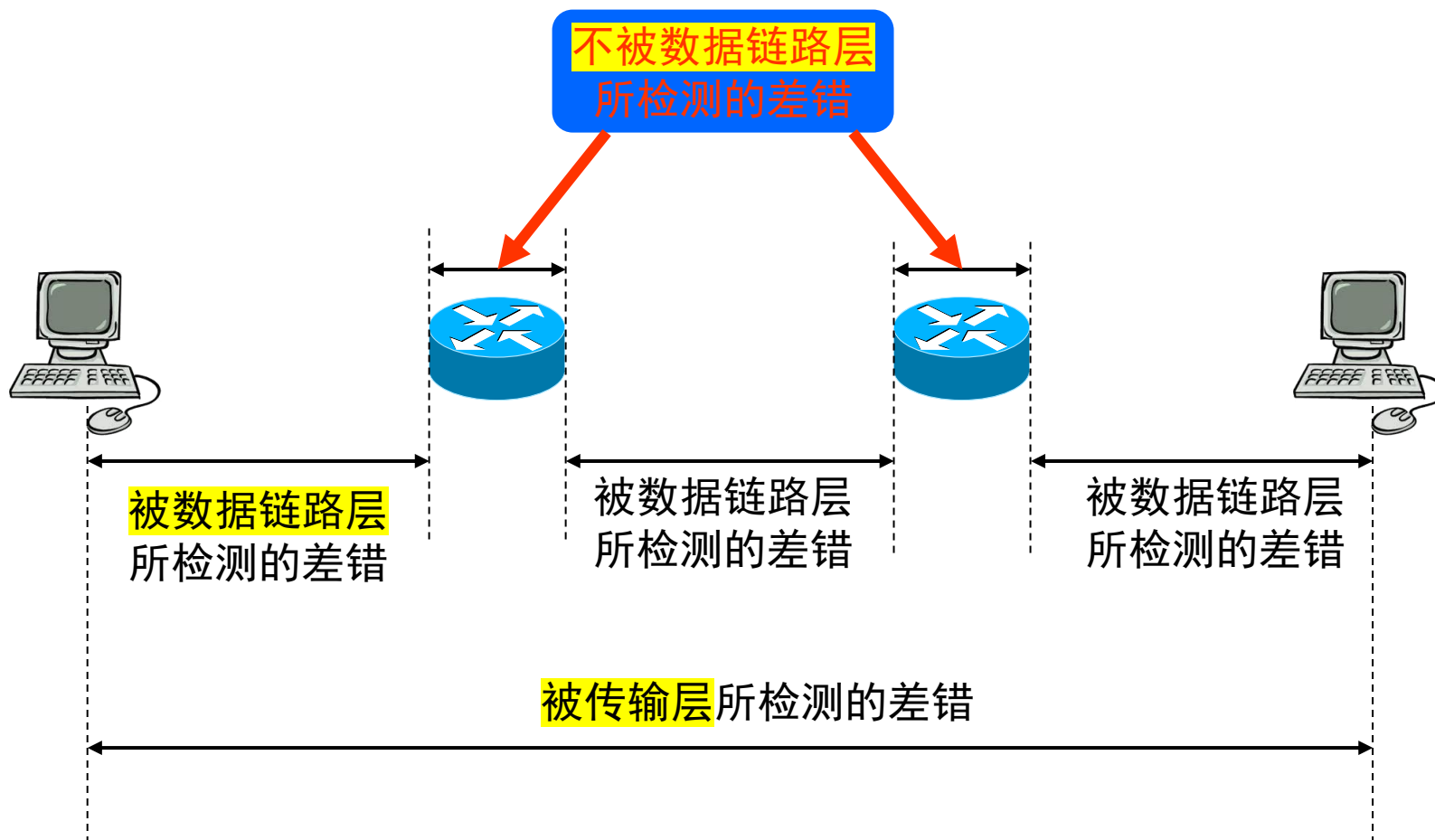


向下复用

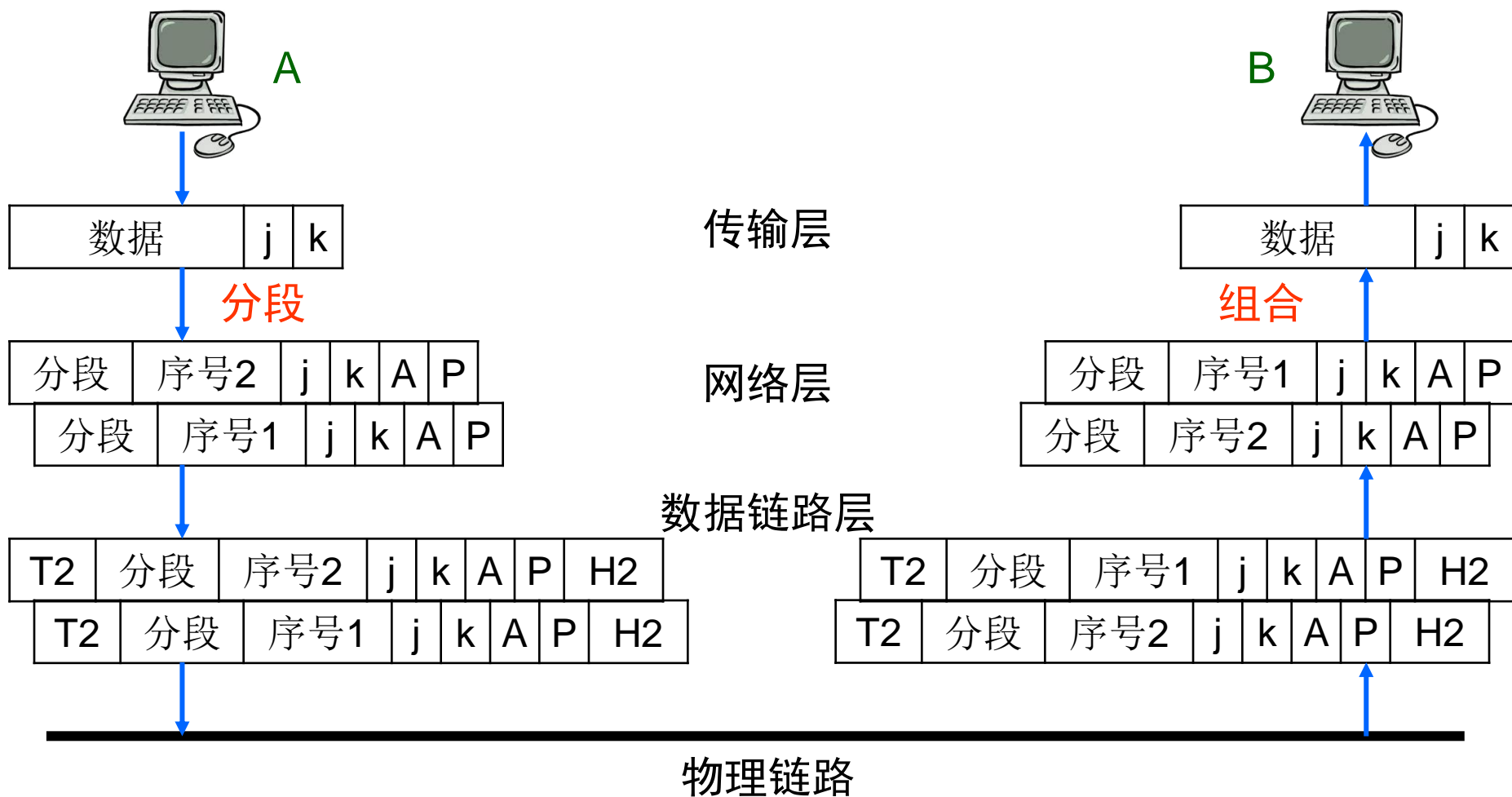
## 6.1.4 可靠传输

- 可靠传输包括以下4个方面：
  - 差错控制：保证可靠性
  - 次序控制：分段和连接、序列编号
  - 丢失控制：丢失重传
  - 重复控制：通过序列编号完成

# 差错控制

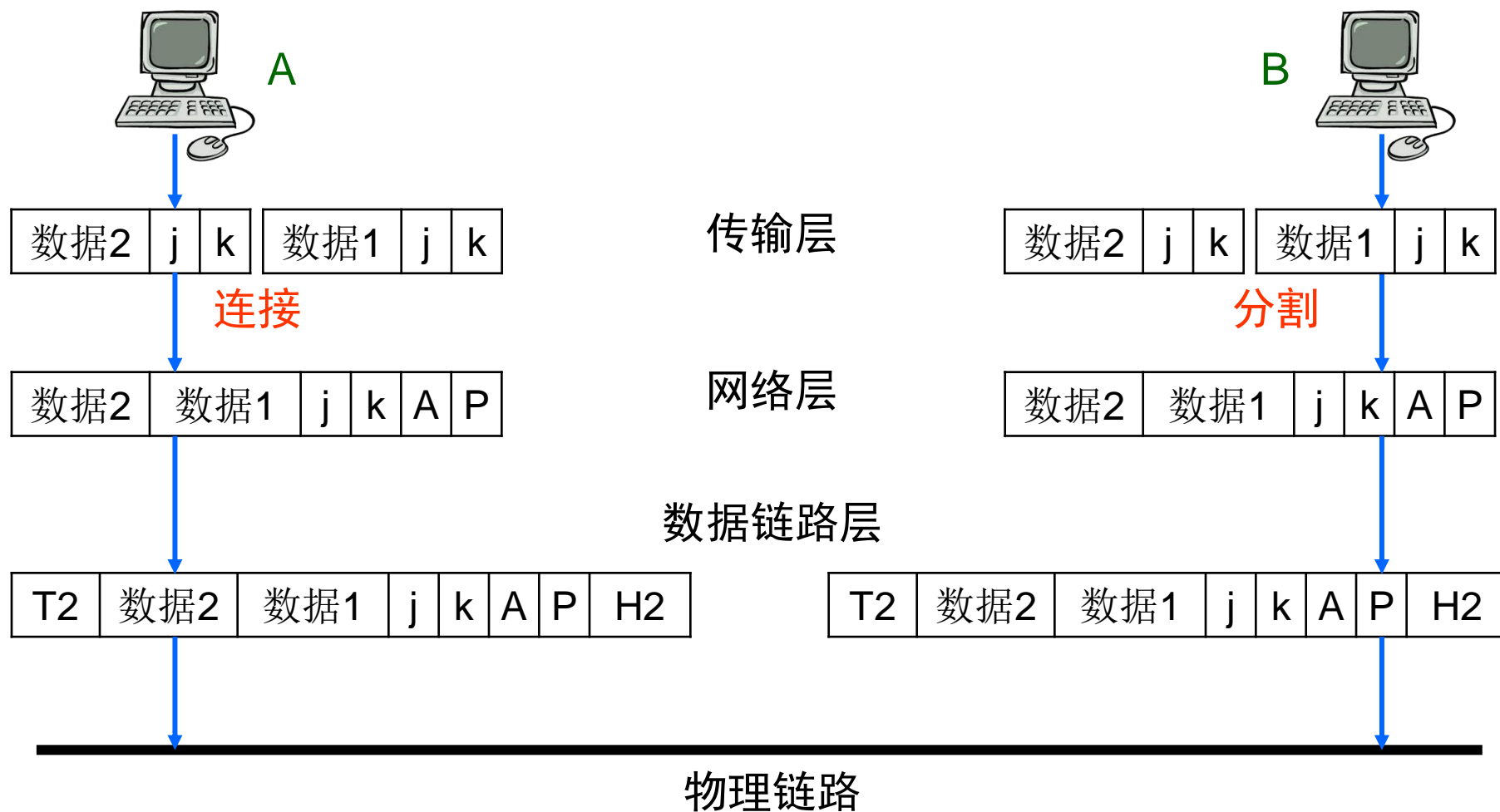


# 次序控制——分段和组合



H2、T2: 数据链路层帧首尾, A、P: 网络地址, j、k: 端口号

# 次序控制——连接和分割



H2、T2: 数据链路层帧首尾, A、P: 网络地址, j、k: 端口号

# 丢失控制

- 传输层要确保一次传输的所有片段都到达目的地
- 如果发生了丢失要重传

# 重复控制

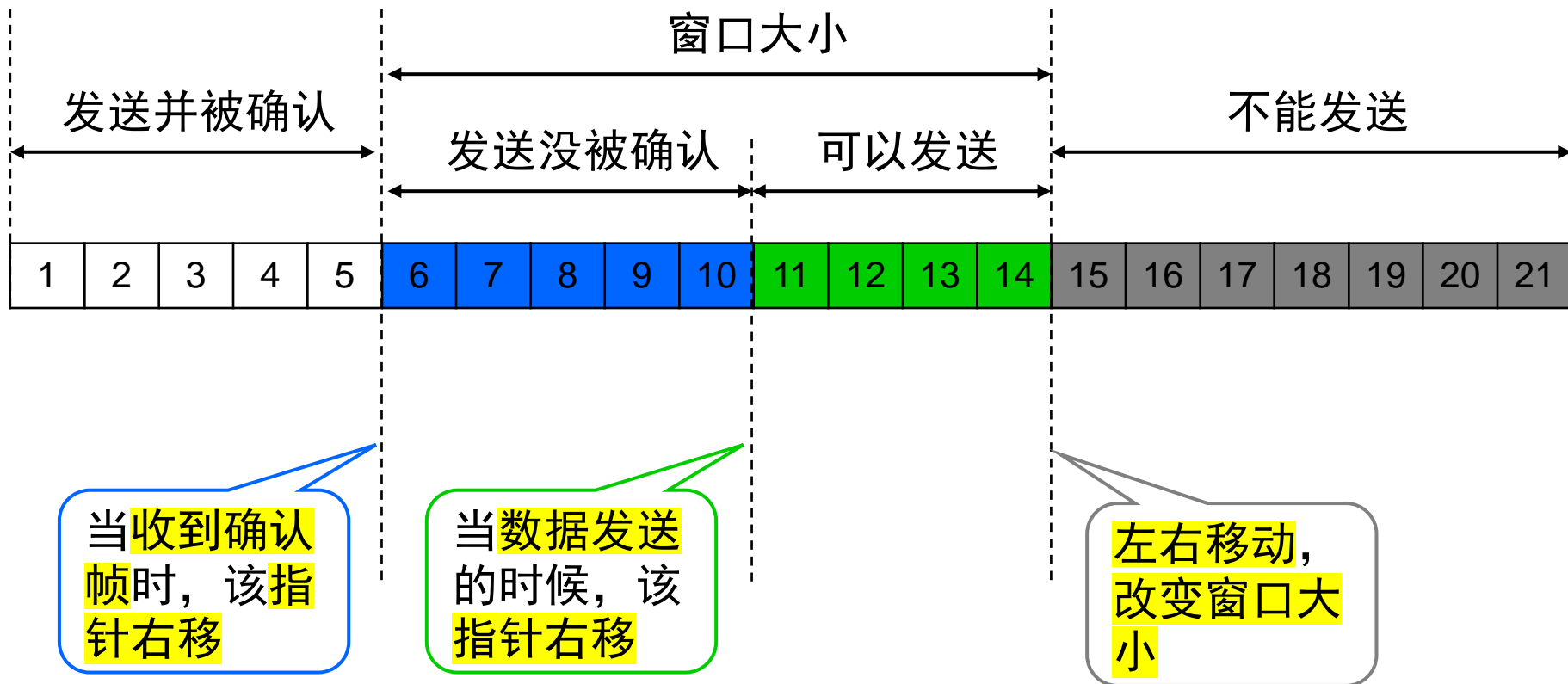
- 使用序列编号保证接收方丢弃重复数据



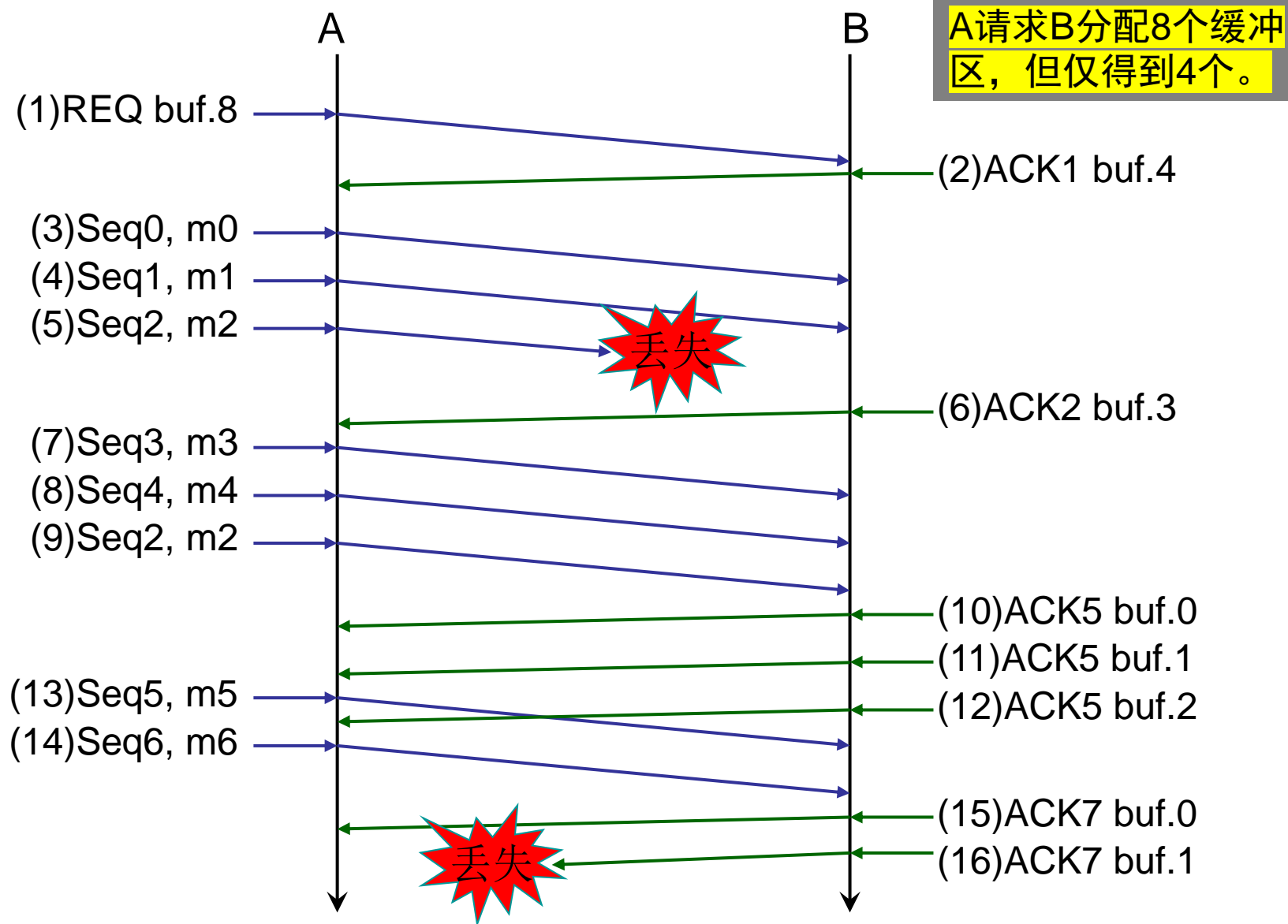
## 6.1.5 传输层流量控制

- 和数据链路层一样，传输层也负责流量控制。但是，传输层中的流量控制是作用在端到端上的，而不是作用在单条链路上的。
- 传输层流量控制也使用滑动窗口协议，但是传输层中的窗口在大小上是可以变化的，以适应可使用的缓冲区的变化情况。
- 使用3个指针识别缓冲区

# 滑动窗口



# 传输层的流量控制

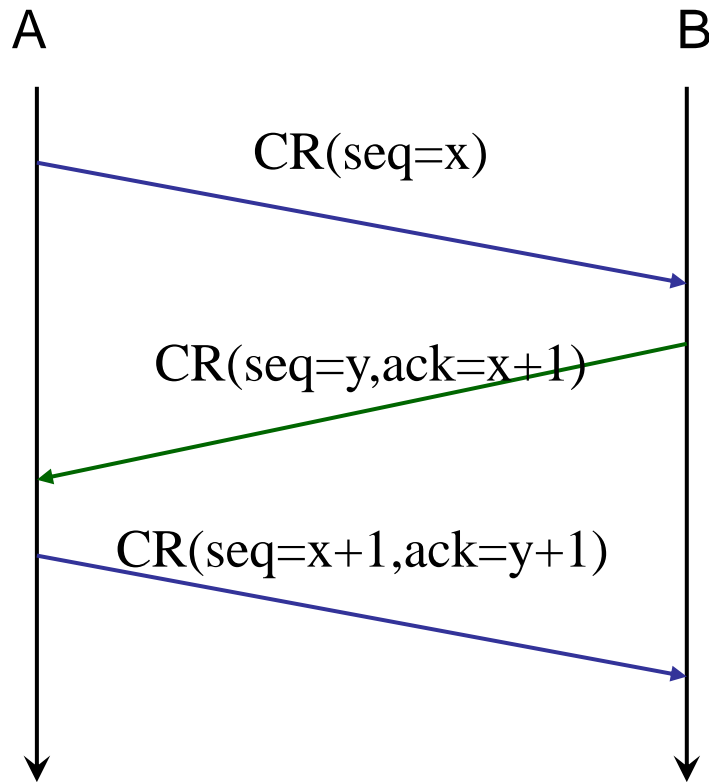


## 6.1.6 传输连接

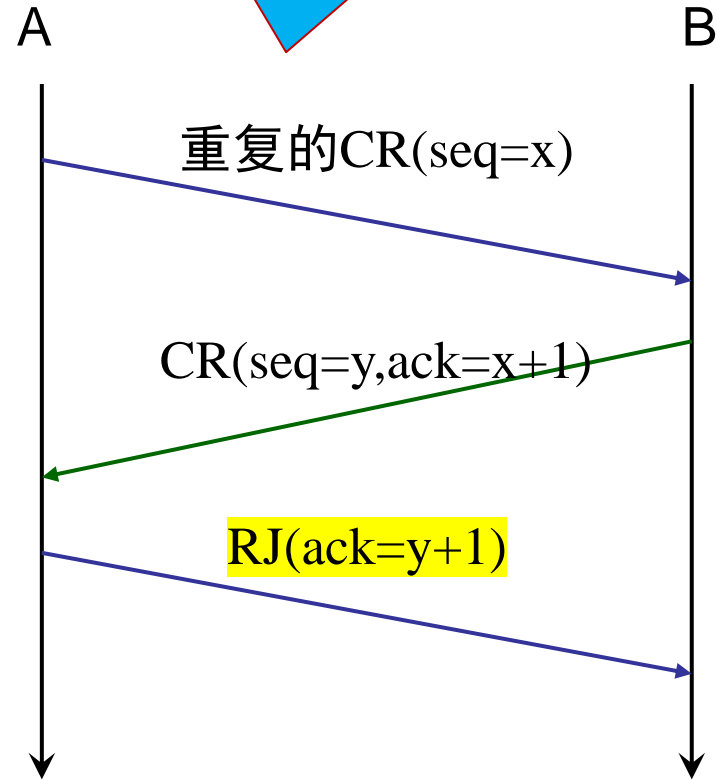
- 传输层端到端的传送可以采用两种模式来完成：
  - 面向连接
  - 面向无连接
- 面向连接传输由三个步骤：连接建立，数据传输和连接终止。

# 连接建立

B不知道A发的是不是一个重复的建立连接的申请，A根据B发的应答信号来判断自己发的是不是重复的建立连接的申请，如果是，则拒绝该连接请求



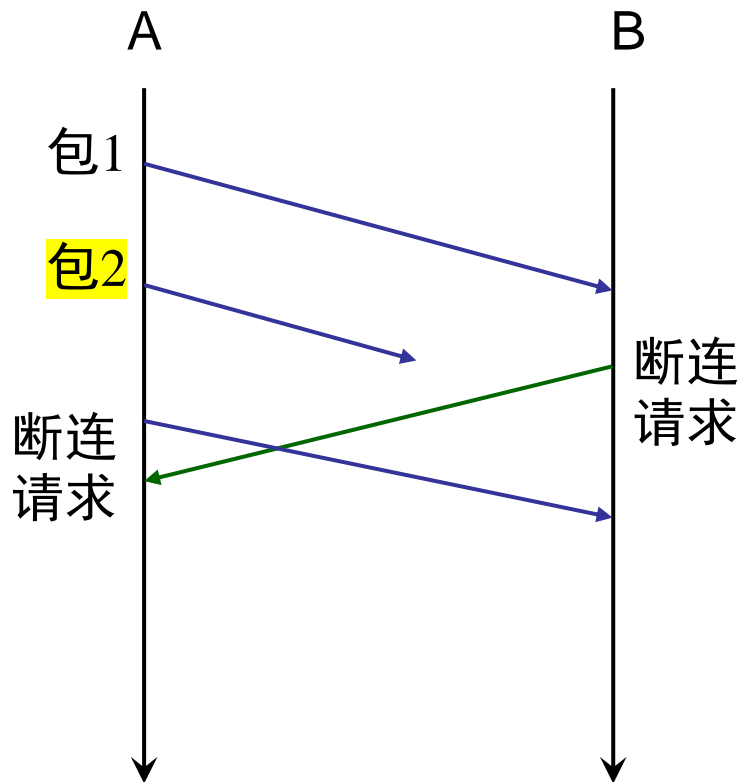
连接的建立



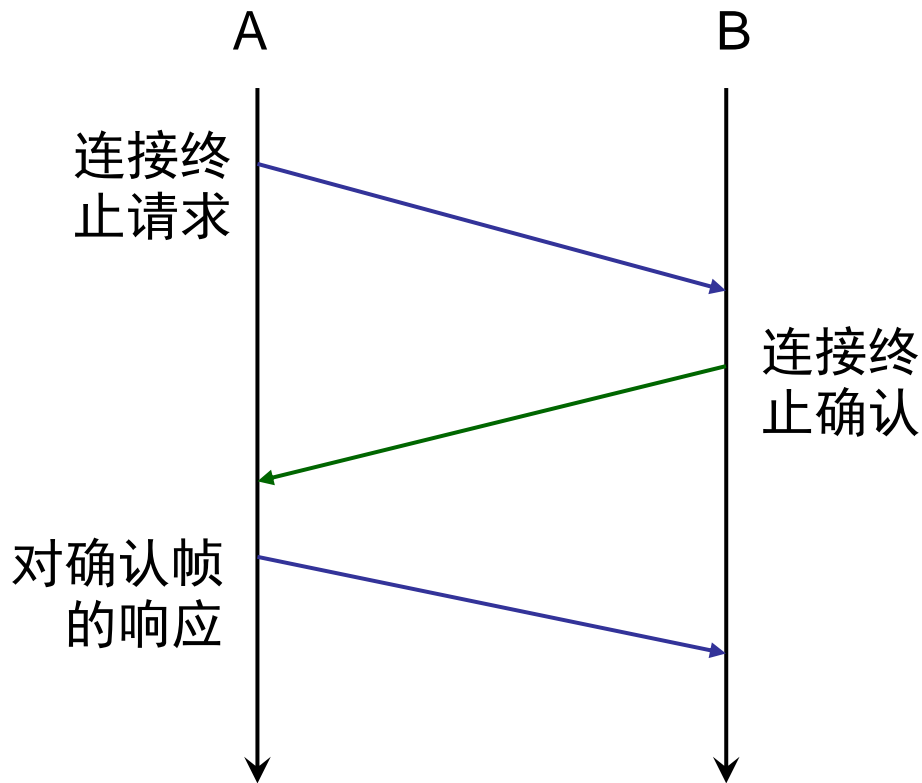
拒绝重复的连接请求

三次握手方法建立连接

# 连接终止

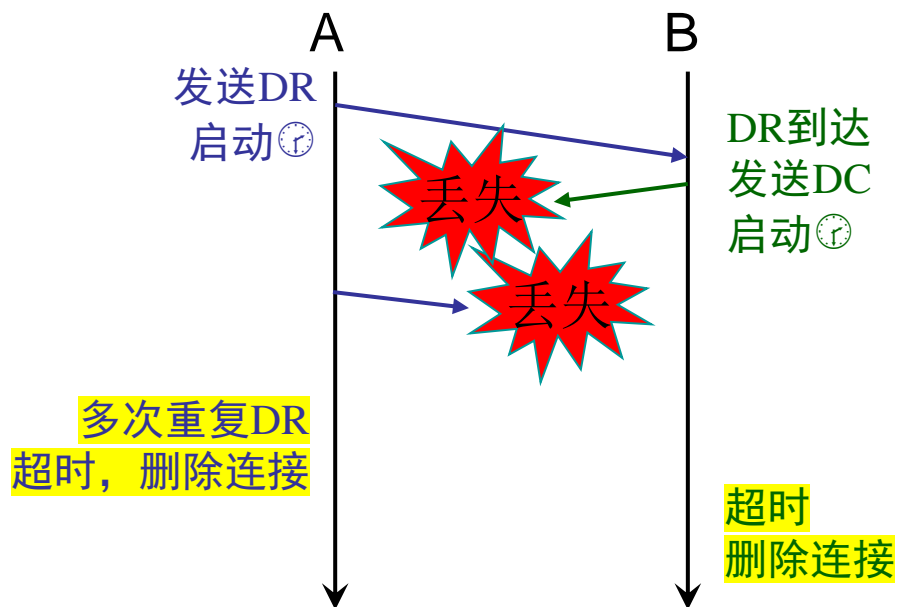
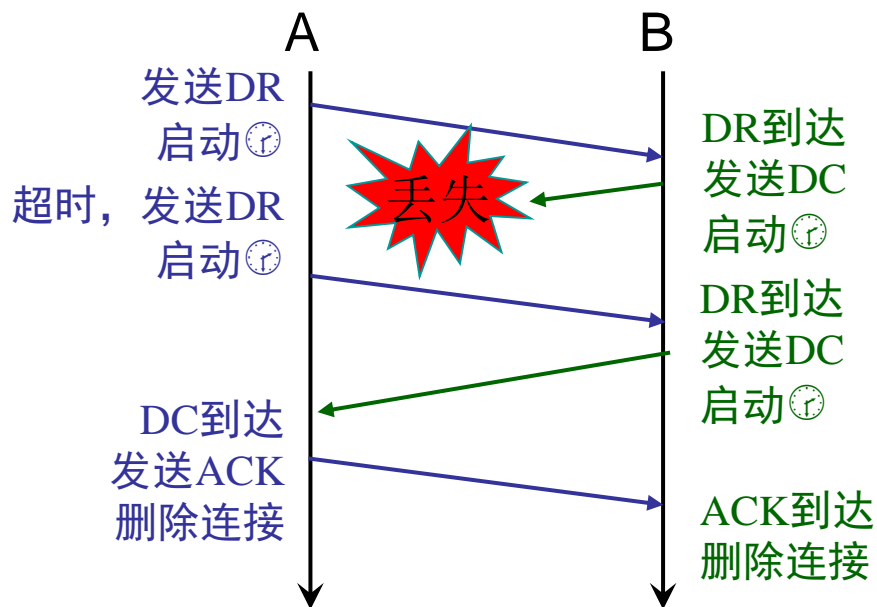
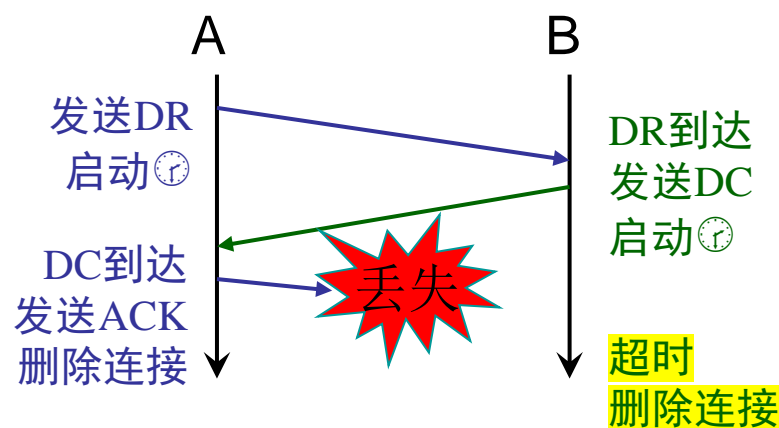
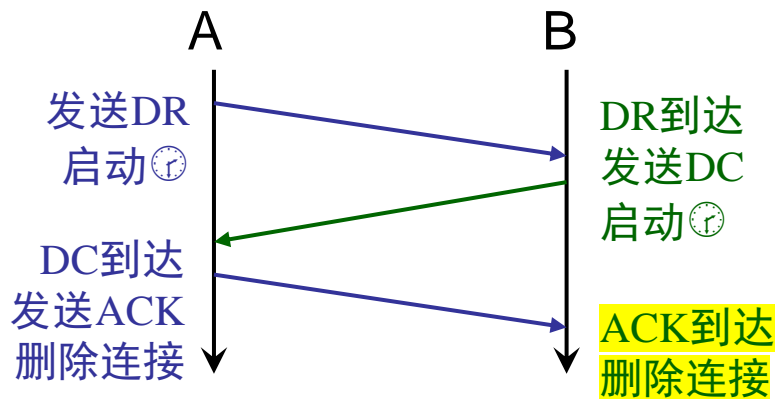


错误断连，数据丢失



正确的断连方式

# 三次握手方法的连接释放



DR: 断连请求 DC: 断连证实 ACK: 确认

## 6.2 用户数据报协议UDP

- UDP采用非连接的方式提供网络应用层的事务处理
- UDP不提供可靠性，即UDP协议不提供端到端的确认和重传功能，它不保证数据包一定能到达目的地，因此称为不可靠协议。
  - 不可靠：丢失、重复、延迟、乱序和损坏
- UDP必需在IP上运行。它的下层协议是以IP作为前提的。



# UDP特征

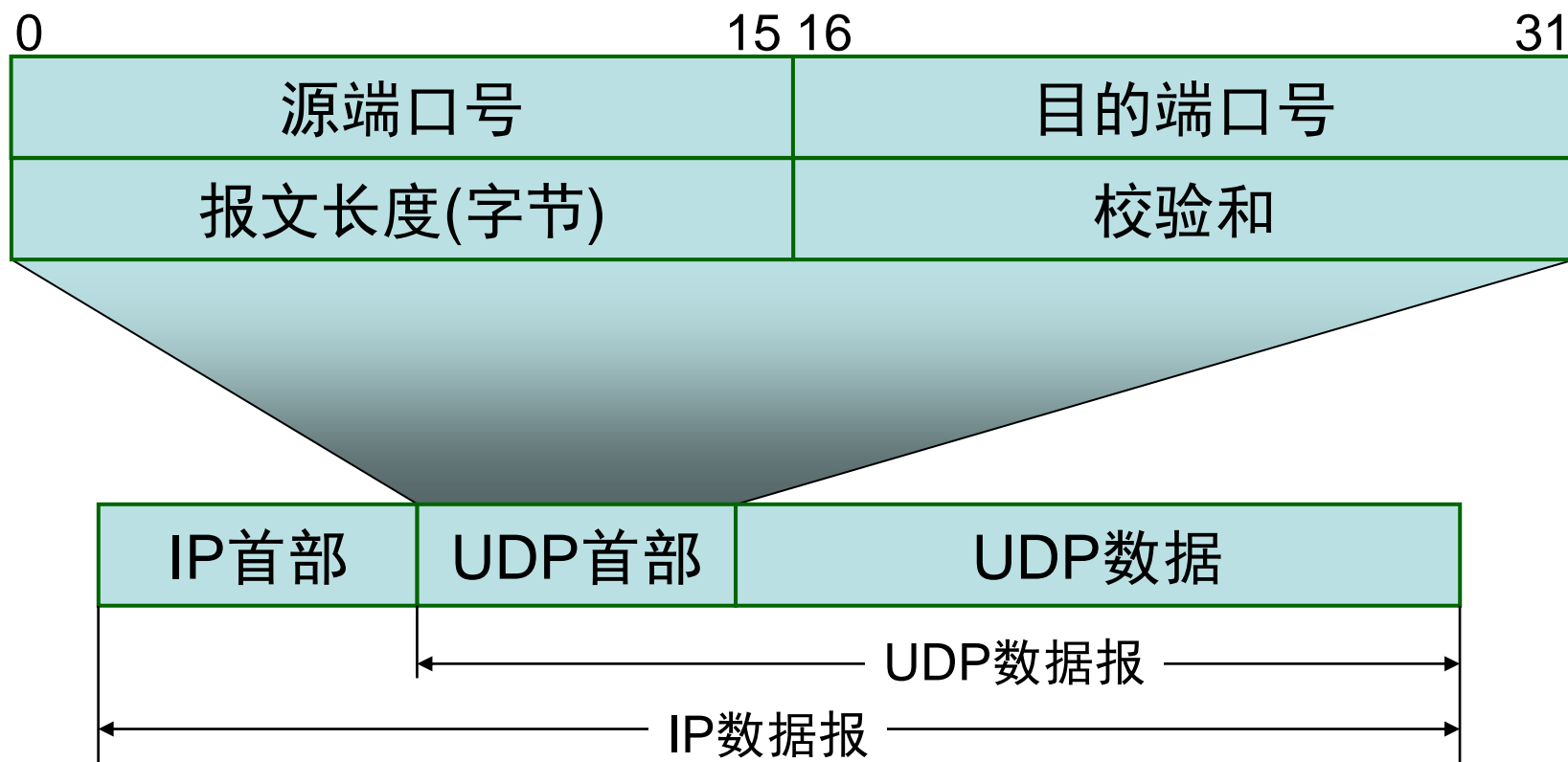
- UDP提供端到端服务，允许应用进程发送和接收单个报文，每个报文被装进单个数据报中进行传输。
- UDP可以被表征为：
  - 端到端
  - 无连接
  - 面向报文
  - 尽力而为
  - 任意交互
  - 操作系统无关

# 采用的通讯方式

- UDP允许采用4种交互通信方式
  - 一对一
  - 一对多
  - 多对一
  - 多对多
- 使用协议端口号标识端点
  - 因为UDP必须跨越异构计算机，为了避免含糊性，UDP定义一个标识符抽象集——协议端口号(protocol port number)。
  - 在每台计算机内提供端口号与操作系统所用的程序标识符之间的映射关系。

## 6.2.1 UDP报头格式

- UDP协议直接利用IP协议进行UDP数据报的传输。
- UDP数据报封装在IP数据报中



# UDP报头格式(1)

0	15	16	31
源端口号		目的端口号	
报文长度(字节)		校验和	

- **16位UDP源端口号(Source Port):** 该端口号作为接收进程返回数据时的目的端口。
  - 可选字段。若不选用, 其值为0。
- **16位UDP目的端口号(Destination Port):** 该端口号是作为接收主机内与特定应用进程相关联的地址。

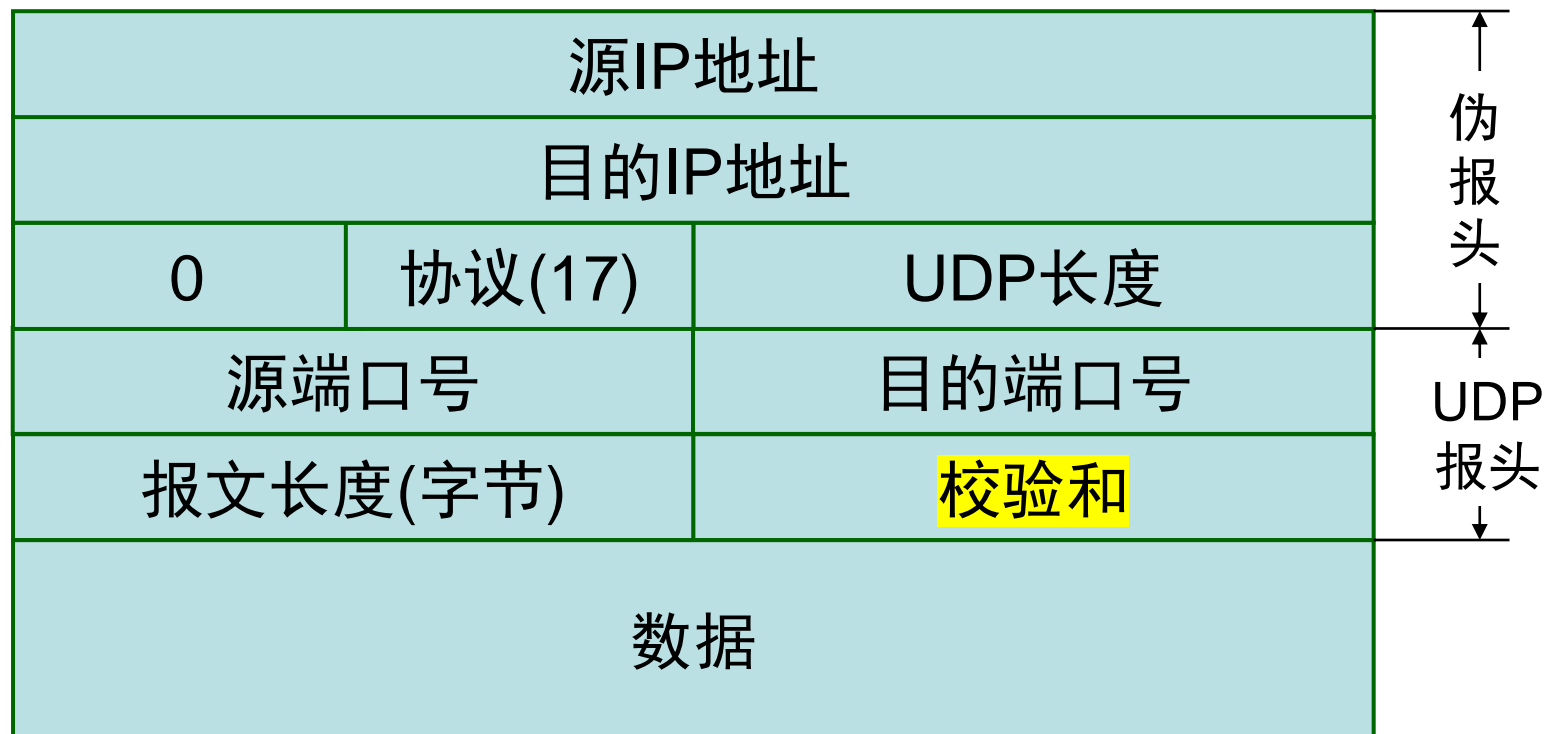
# UDP报头格式(2)

0	15	16	31
源端口号		目的端口号	
报文长度(字节)		校验和	

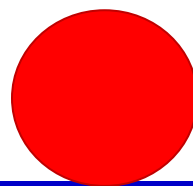
- **UDP数据报的长度(Length):** UDP长度字段指的是UDP首部和UDP数据的字节长度。UDP数据报长度是IP全长减去IP首部的长度。最小值是8。
- **UDP校验和(Checksum):** 可选字段，用来检验传输过程中是否出现了错误。UDP校验和覆盖伪首部、UDP首部和UDP数据。

## 6.2.2 UDP校验和

- UDP在报中包含一个12字节的伪报头以计算校验和。该伪报头包含IP报头的某些域，目的是让UDP检测数据确已到达正确的目的端。



# 计算 UDP 校验和的例子



12 字节 伪首部	153.19.8.104			
	171.3.14.11			
8 字节 UDP 首部	全 0	17	15	
	1087		13	
7 字节 数据	15		全 0	
	数据	数据	数据	数据
	数据	数据	数据	全 0

填充

10011001 00010011 → 153.19  
 00001000 01101000 → 8.104  
 10101011 00000011 → 171.3  
 00001110 00001011 → 14.11  
 00000000 00010001 → 0 和 17  
 00000000 00001111 → 15  
 00000100 00111111 → 1087  
 00000000 00001101 → 13  
 00000000 00001111 → 15  
 00000000 00000000 → 0 (校验和)  
 01010100 01000101 → 数据  
 01010011 01010100 → 数据  
 01001001 01001110 → 数据  
 01000111 00000000 → 数据和 0 (填充)

按二进制反码运算求和 10010110 11101101 → 求和得出的结果  
 将得出的结果求反码 01101001 00010010 → 校验和

## 6.3 传输控制协议TCP

- TCP提供了一种可靠的面向连接的字节流传输层服务，TCP提供端到端的流量控制，并计算和验证一个强制性的端到端检查和。



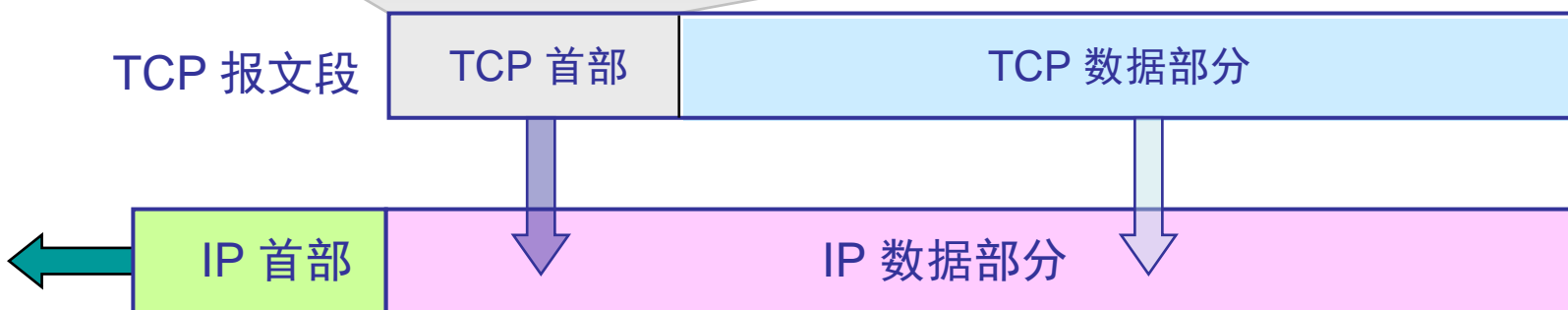
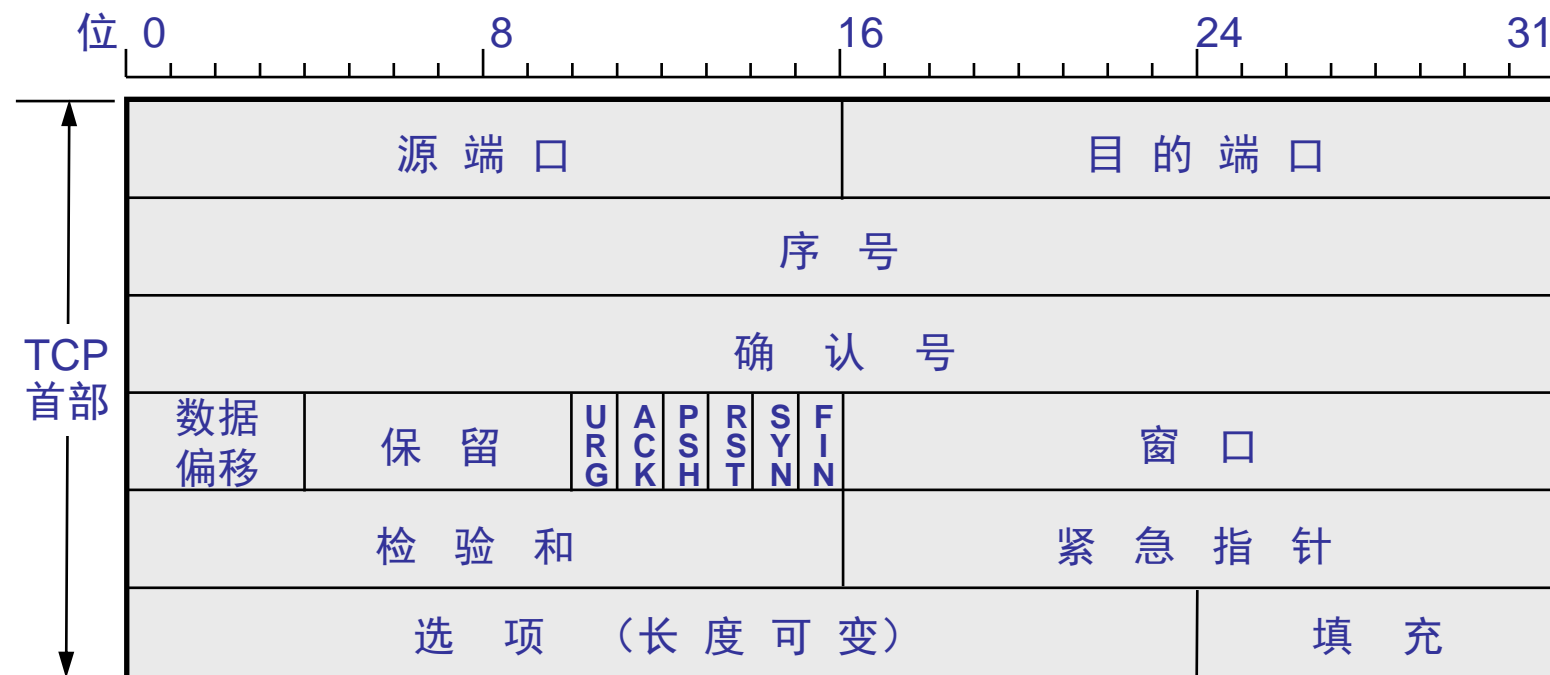
# TCP的主要特点

- TCP提供的服务有以下主要特点：
  - 面向连接
  - 点对点->不能广播、不能组播、只能单播
  - 完全的可靠性
  - 全双工通信
  - 流接口
  - 可靠的连接建立
  - 友好的连接关闭

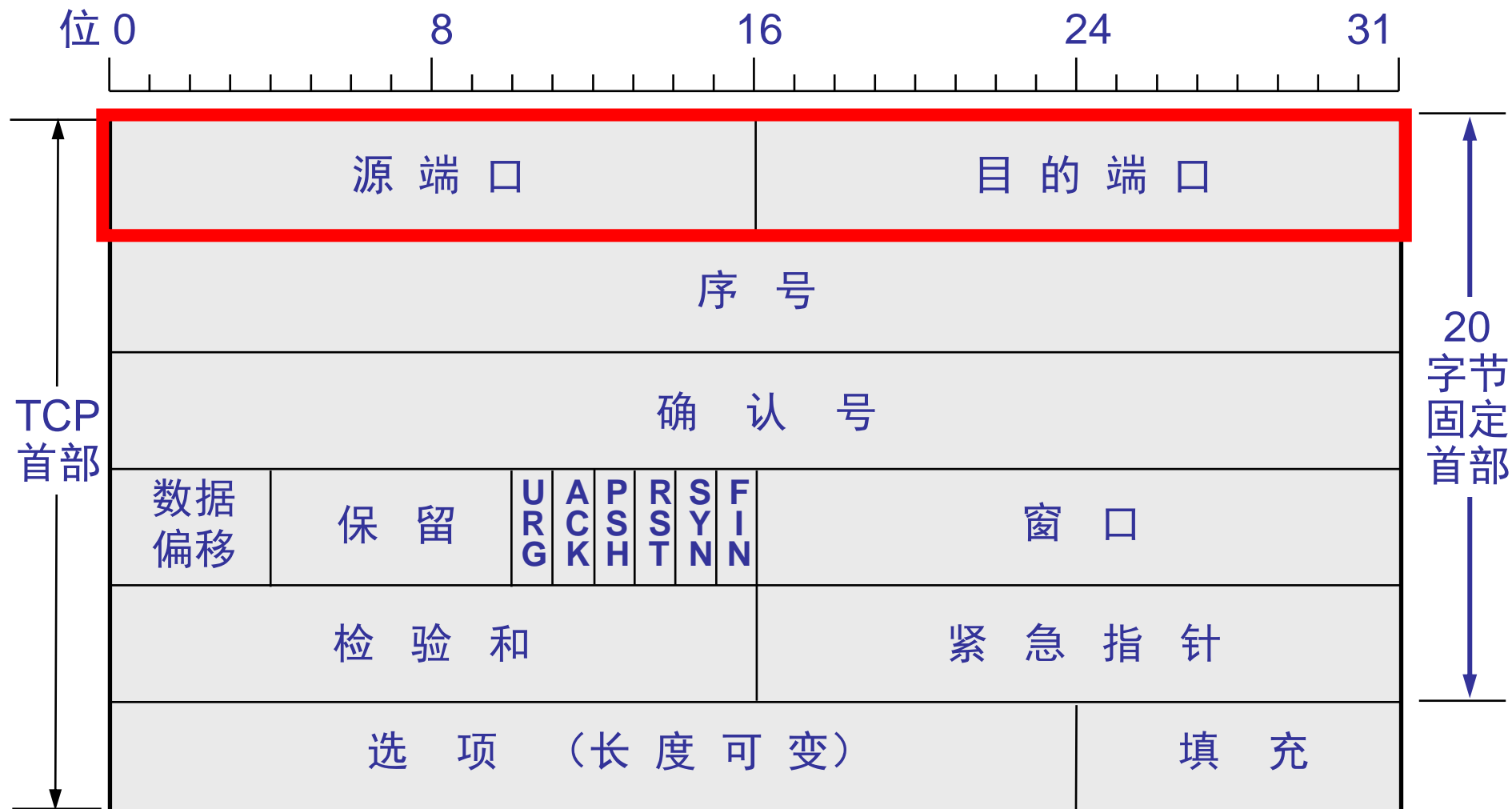
## 6.3.1 TCP提供的服务

- TCP在IP数据报中的封装
- 尽管TCP和UDP都使用相同的网络层（IP），TCP却向应用层提供与UDP完全不同的服务。TCP提供一种面向连接的、可靠的字节流服务。
- TCP协议可以表述为一个没有选择确认或否认的滑动窗口协议，它的窗口大小是可变的。

## 6.3.2 TCP的报头格式



# TCP的报头字段(1)



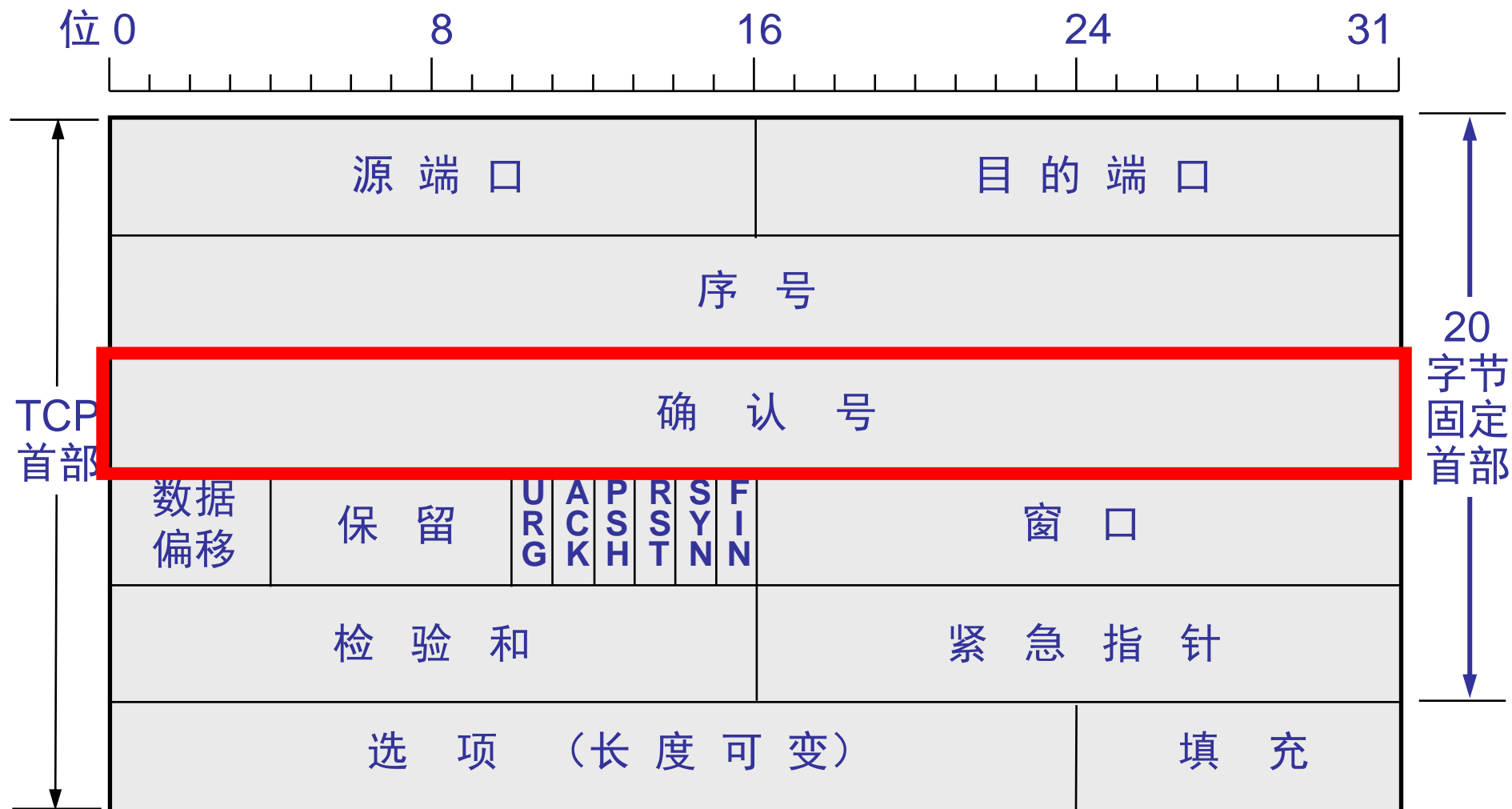
源端口和目的端口字段各占 2 字节。端口是传输层与应用层的服务接口。传输层的复用和分用功能都要通过端口才能实现。

# TCP的报头字段(2)



序号字段：占4字节。TCP 连接中传送的数据流中的每一个字节都编上一个序号。序号字段的值则指的是本报文段所发送的数据的第一个字节的序号。

# TCP的报头字段(3)



确认号字段：占4字节，是期望收到对方的下一个报文段的数据的第一个字节的序号。

# TCP的报头字段(4)



数据偏移(即首部长度的): 占4位, 它指出 TCP 报文段的数据起始处距离 TCP 报文段的起始处有多远。“数据偏移”的单位是 32 位字(以 4 字节为计算单位)。

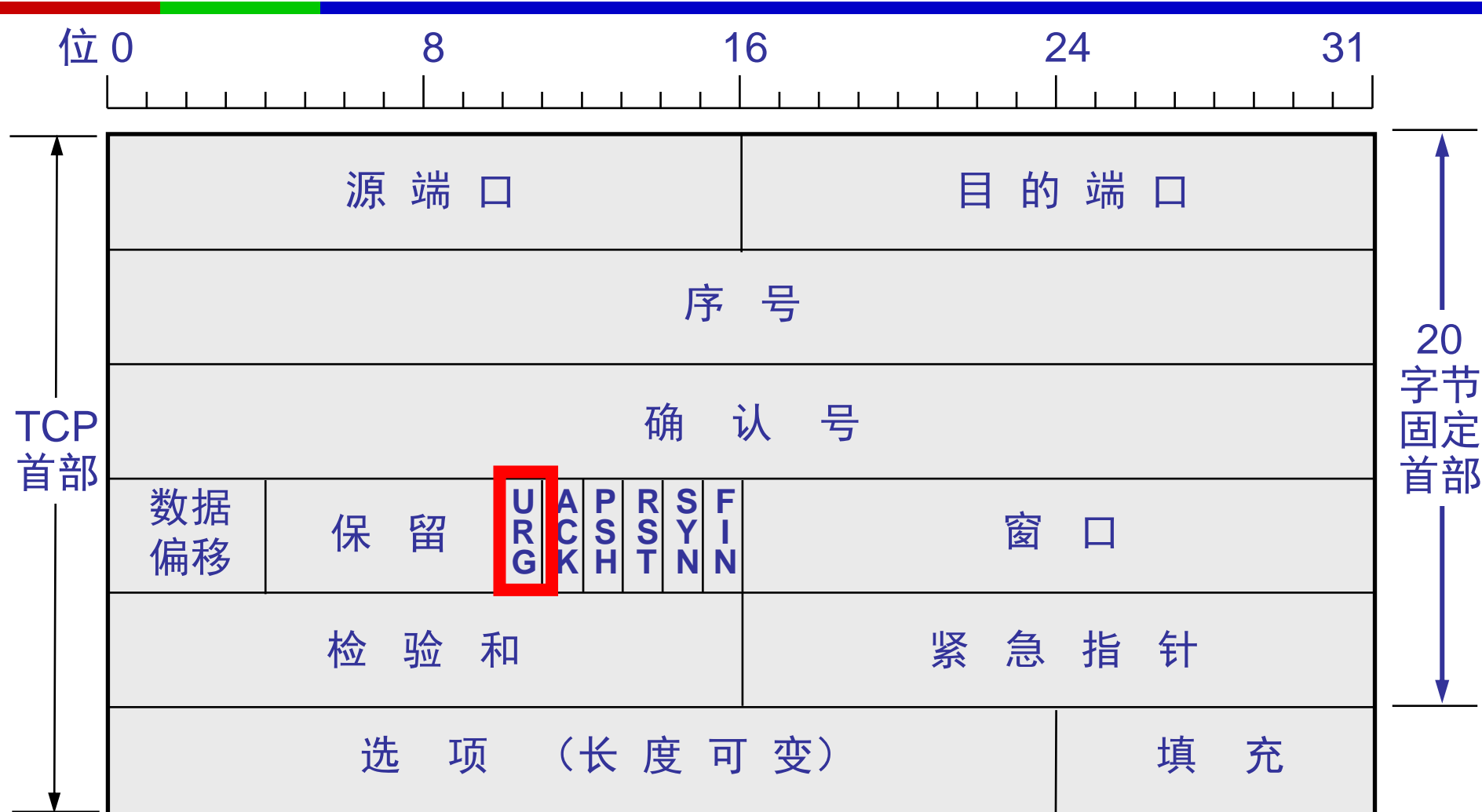
# TCP的报头字段(5)



保留字段——占 6 位，保留为今后使用，但目前应置为 0。

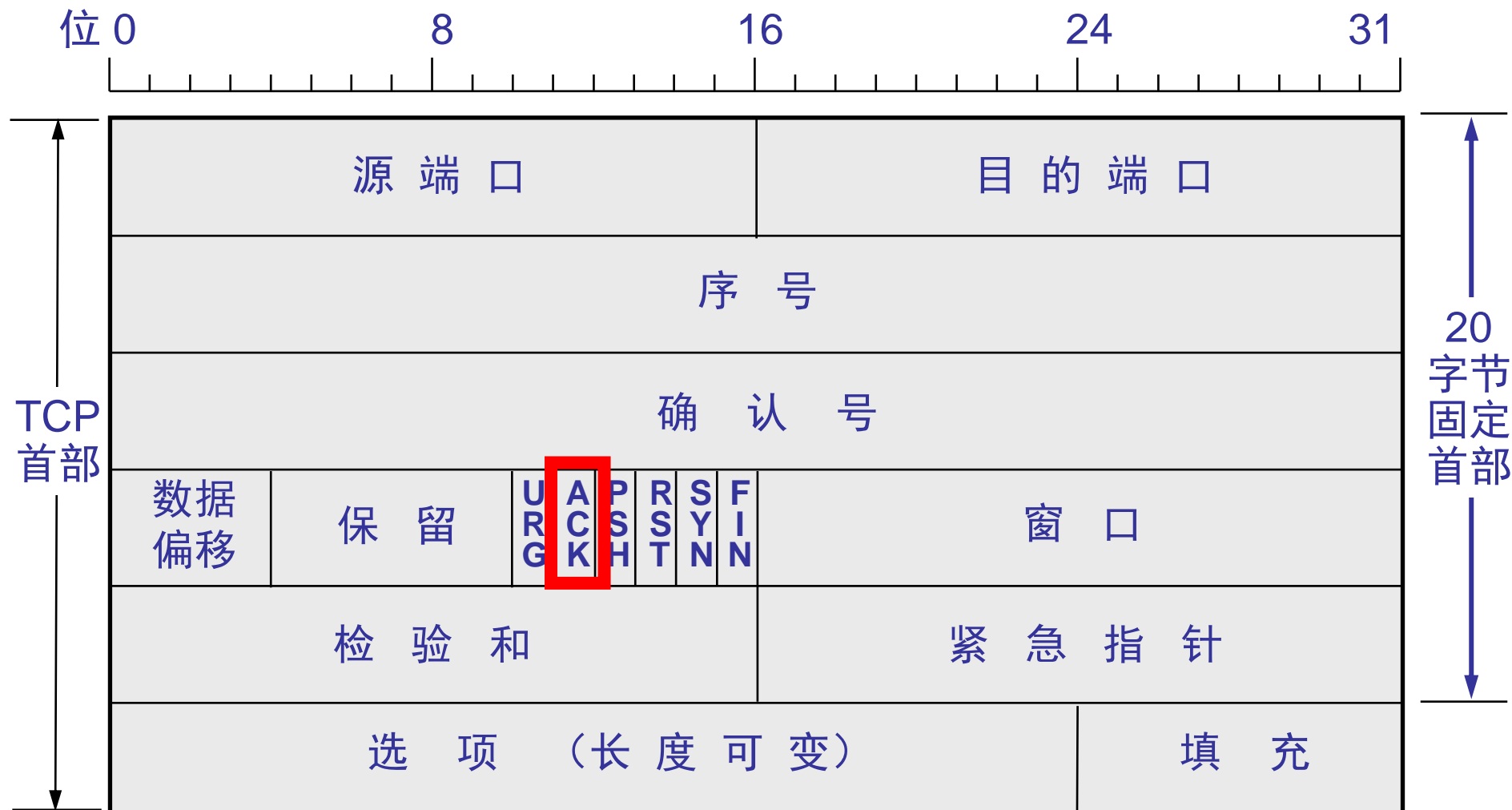


# TCP的报头字段(6) 6个标志位—URG



**紧急URG:** 当URG = 1时, 表明紧急指针字段有效。此报文段中有紧急数据, 系统应尽快传送(相当于高优先级的数据)。

# TCP的报头字段(7) 6个标志位—ACK



确认ACK：只有当 ACK=1时确认号字段才有效。当ACK=0 时，确认号无效。

# TCP的报头字段(8) 6个标志位—PSH



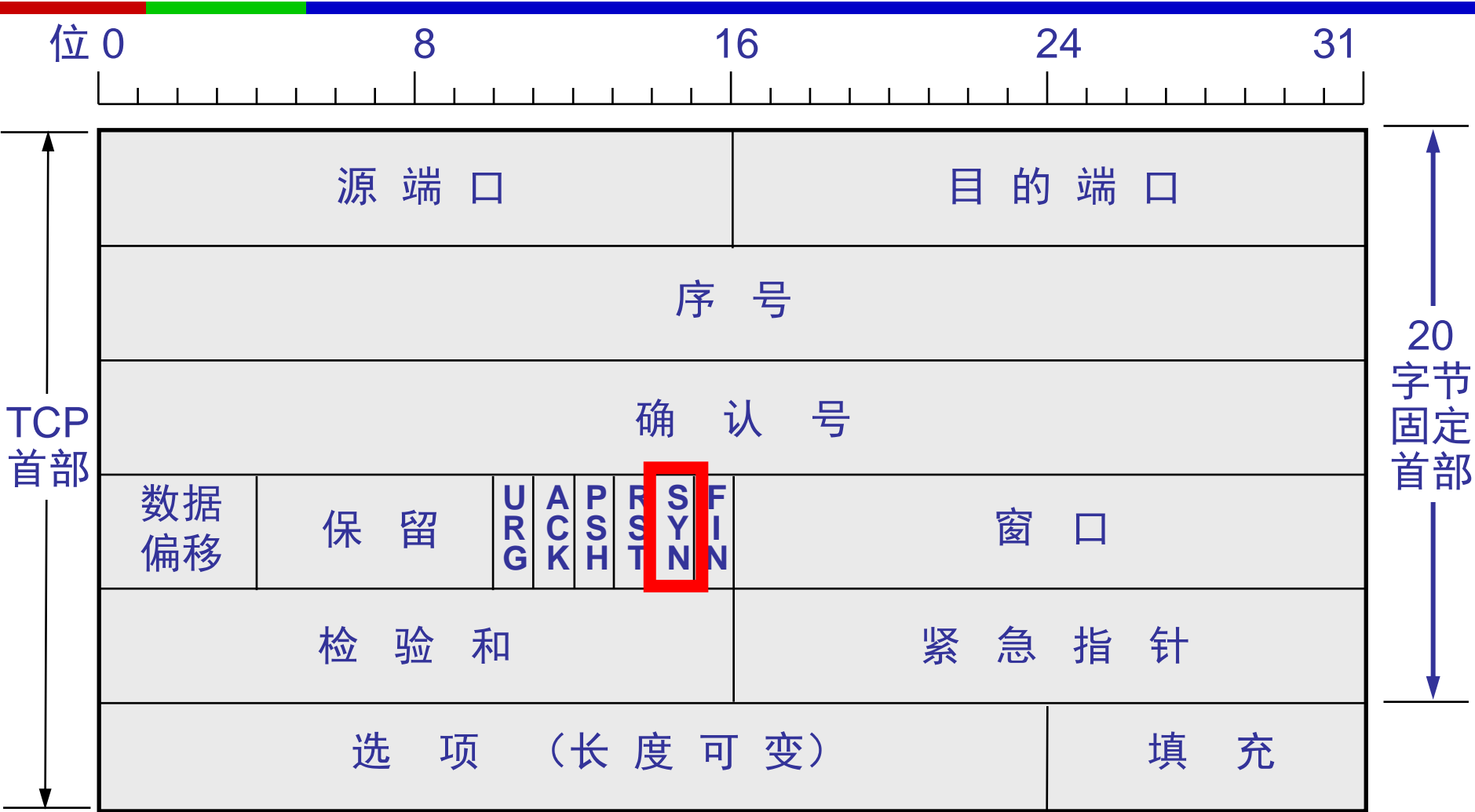
**推送 PSH：** 接收TCP收到PSH=1的报文段，就尽快地交付接收应用进程，而不再等到整个缓存都填满了后再向上交付。

# TCP的报头字段(9) 6个标志位—RST



**复位 RST**：当RST=1时，表明TCP连接中出现严重差错(如由于主机崩溃或其他原因)，必须释放连接，然后再重新建立运输连接。

# TCP的报头字段(10) 6个标志位—SYN



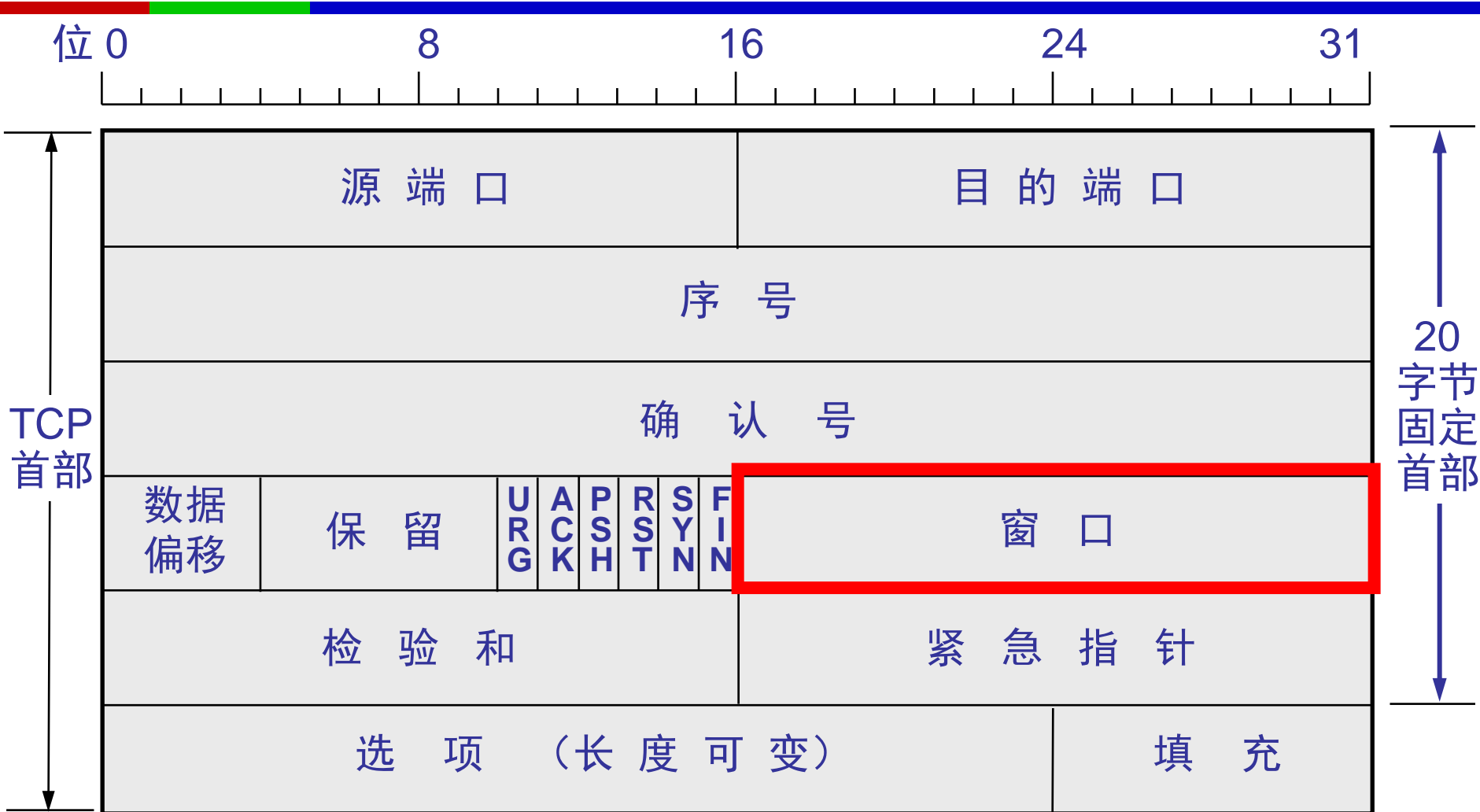
同步 SYN：同步SYN=1表示这是一个连接请求或连接接受报文。

# TCP的报头字段(11) 6个标志位—FIN



**终止FIN:** 用来释放一个连接。FIN=1表明此报文段的发送端的数据已发送完毕，并要求释放运输连接。

# TCP的报头字段(12)



窗口字段占2字节，用来让对方设置发送窗口的依据，单位为字节。

# TCP的报头字段(13)



检验和：占2字节。检验和字段检验的范围包括首部和数据这两部分。在计算检验和时，要在TCP报文段的前面加上12字节的伪首部。



# TCP的报头字段(14)



紧急指针字段：占16位，是一个正的偏移量，指出在本报文段中紧急数据共有多少个字节(紧急数据放在本报文段数据的最前面)，和序号字段中的值相加表示紧急数据最后一个字节的序号。

# TCP的报头字段(15)

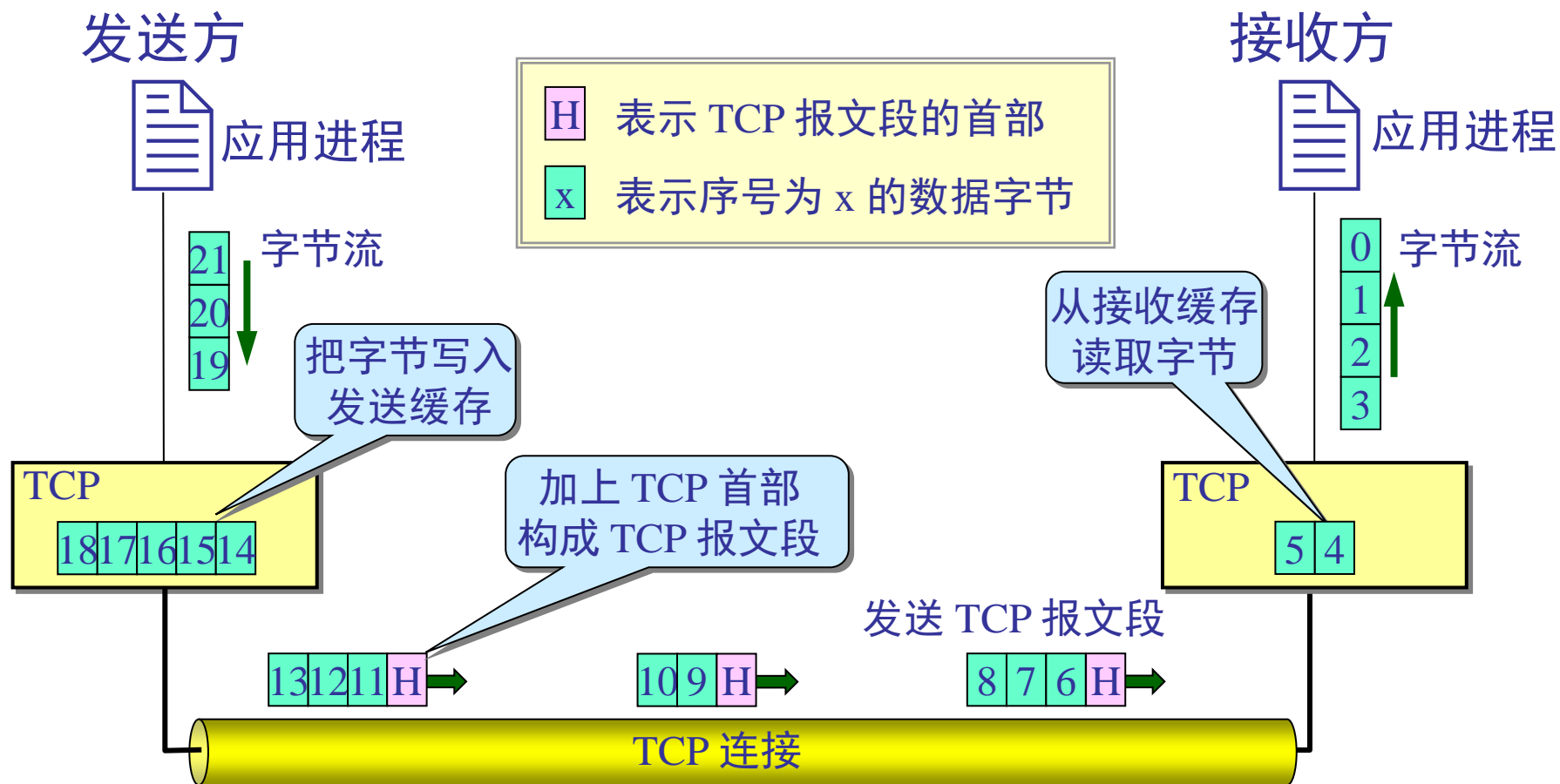


选项字段：长度可变。TCP最初只规定了一种选项，即最大报文段长度 MSS。MSS 告诉对方TCP：“我的缓存所能接收的报文段的数据字段的最大长度是MSS个字节。”

# TCP特性(1)

- 在TCP连接中每个传输的字节都被计数，**确认序号包含**发送确认的一端**所期望收到的下一个序号**。因此，确认序号应当是上次已经成功收到数据字节序号加1。只有**ACK**标志为1时，确认序号字段才有效。
- 发送**ACK**无需任何代价，因为32bit的确认序号字段和**ACK**标志一样，总是**TCP**首部的一部分。因此，一旦一个**TCP**连接建立起来，**ACK**标志总是被设置1。

# TCP面向流的概念



# TCP特性(2)

- TCP 可以描述为一个没有选择确认或否定的滑动窗口协议。TCP 缺少选择确认是因为 TCP 首部中的确认序号表示发方已成功收到字节，但还不包含确认序号所指的字节。当前还无法对数据流中选定的部分进行确认。
- 例如，如果1~1024字节已经成功收到，下一个报文段中包含序号从2049 ~ 3072的字节，收端并不能确认这个新的报文段，它所能做的就是发回一个确认序号为1025的ACK。

# TCP所采用的技术

- 对付分组重复和乱序传递的排序技术
  - 发送端为每个分组附加一个序号
- 对付分组丢失的重传技术
  - 使用带重传的正向确认机制
- 避免分组重复的技术
- 防止数据过荷的流量控制技术

# 超时重传时间的选择

- 重传机制是 TCP 中最重要和最复杂的问题之一。
- TCP 每发送一个报文段，就对这个报文段设置一次计时器。只要计时器设置的重传时间到但还没有收到确认，就要重传这一报文段。
- 往返时延的方差很大
  - 由于 TCP 的下层是一个互联网环境，IP 数据报所选择的路由变化很大。因而传输层的往返时间的方差也很大。

# 加权平均往返时间

- TCP 保留了 RTT 的一个加权平均往返时间 RTTS（又称为平滑的往返时间）
- 第一次测量到 RTT 样本时， $RTT_s$  值就取为所测量到的 RTT 样本值。以后每测量到一个新的 RTT 样本，按下式重新计算一次  $RTT_s$ ：

$$\begin{aligned} \text{新的 } RTT_s &= (1-\alpha) \times (\text{旧的 } RTT_s) \\ &\quad + \alpha \times (\text{新的 RTT 样本}) \end{aligned}$$

式中， $0 \leq \alpha < 1$ 。若  $\alpha$  很接近于 0，表示 RTT 值更新较慢。若选择  $\alpha$  接近于 1，则表示 RTT 值更新较快。

- RFC 2988 推荐的  $\alpha$  值为 1/8，即 0.125。



# 超时重传时间 RTO

- RTO 应略大于上面得出的加权平均往返时间  $RTT_S$ 。
- RFC 2988 建议使用下式计算 RTO:

$$RTO = RTT_S + 4 \times RTT_D$$

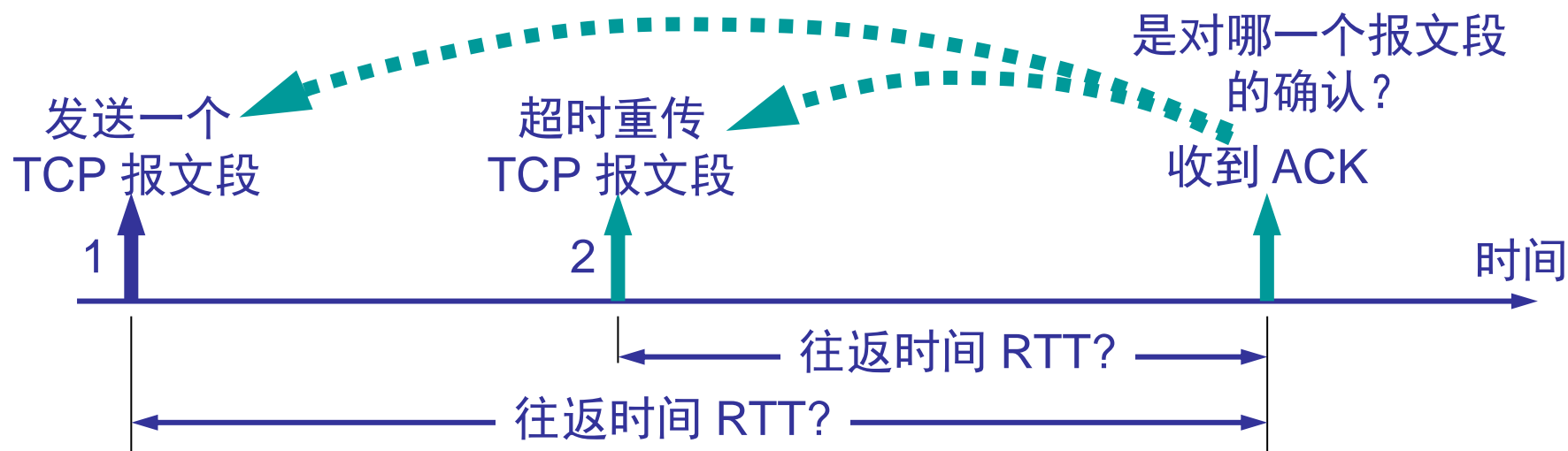
$RTT_D$  是 RTT 的偏差的加权平均值:

- RFC 2988 建议这样计算  $RTT_D$ 。第一次测量时,  $RTT_D$  值取为测量到的 RTT 样本值的一半。在以后的测量中, 则使用下式计算加权平均的  $RTT_D$ :

$$\begin{aligned} \text{新的 } RTT_D = & (1 - \beta) \times (\text{旧的 } RTT_D) \\ & + \beta \times |RTT_S - \text{新的 } RTT \text{ 样本}| \end{aligned}$$

$\beta$  是个小于 1 的系数, 其推荐值是 1/4, 即 0.25。

# Karn算法及其修正



- 报文段每重传一次，就把 RTO 增大一些：

$$\text{新的 RTO} = \gamma \times (\text{旧的 RTO})$$

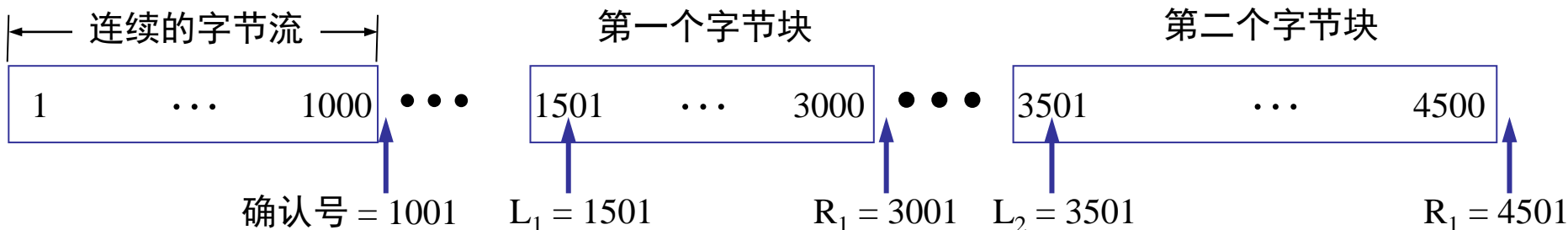
系数 $\gamma$ 的典型值是 2。

# TCP SACK

- TCP无法对一个报文段进行否认。

例如，如果收到包含1025~2048字节的报文段，但它的检查和有错误，TCP接收端所能做的就是发回一个确认序号为1025的ACK。

- 选择确认SACK (Selective ACK, RFC2018)



# RFC 2018的一些规定

- 如果要使用选择确认，那么在建立 **TCP** 连接时，就要在 **TCP** 首部的选项中加入“允许 **SACK**”的选项，而双方必须都事先商定好。原来首部中的“确认号字段”的用法仍然不变。
- 如果使用选择确认，那么需要2个字节用来指明该选项：
  - 1个字节用于**SACK**选项
  - 1个字节用于指明选项的长度
- 由于首部选项的长度最多只有 40 字节，而指明一个边界就要用掉 4 字节，因此在选项中最多只能指明 4 个字节块的边界信息。

# 提高性能的有关措施

## ■ Nagle算法：

- 若发送进程把要发送的数据逐个字节地送到TCP的发送缓存，则发送方把第1个字节发送出去，发后续的字节缓存起来。
  - 当收到对第1个字节的确认后，在将缓存中的数据组装成一个报文段发送出去，同时继续对后续数据进行缓存
  - 当到达的数据已达到发送窗口大小的一半或者已达到最大报文段时，立即发送一个报文段。
- 糊涂窗口综合症：TCP接收缓存已满，而应用进程一次只从缓存中读取一个字节。

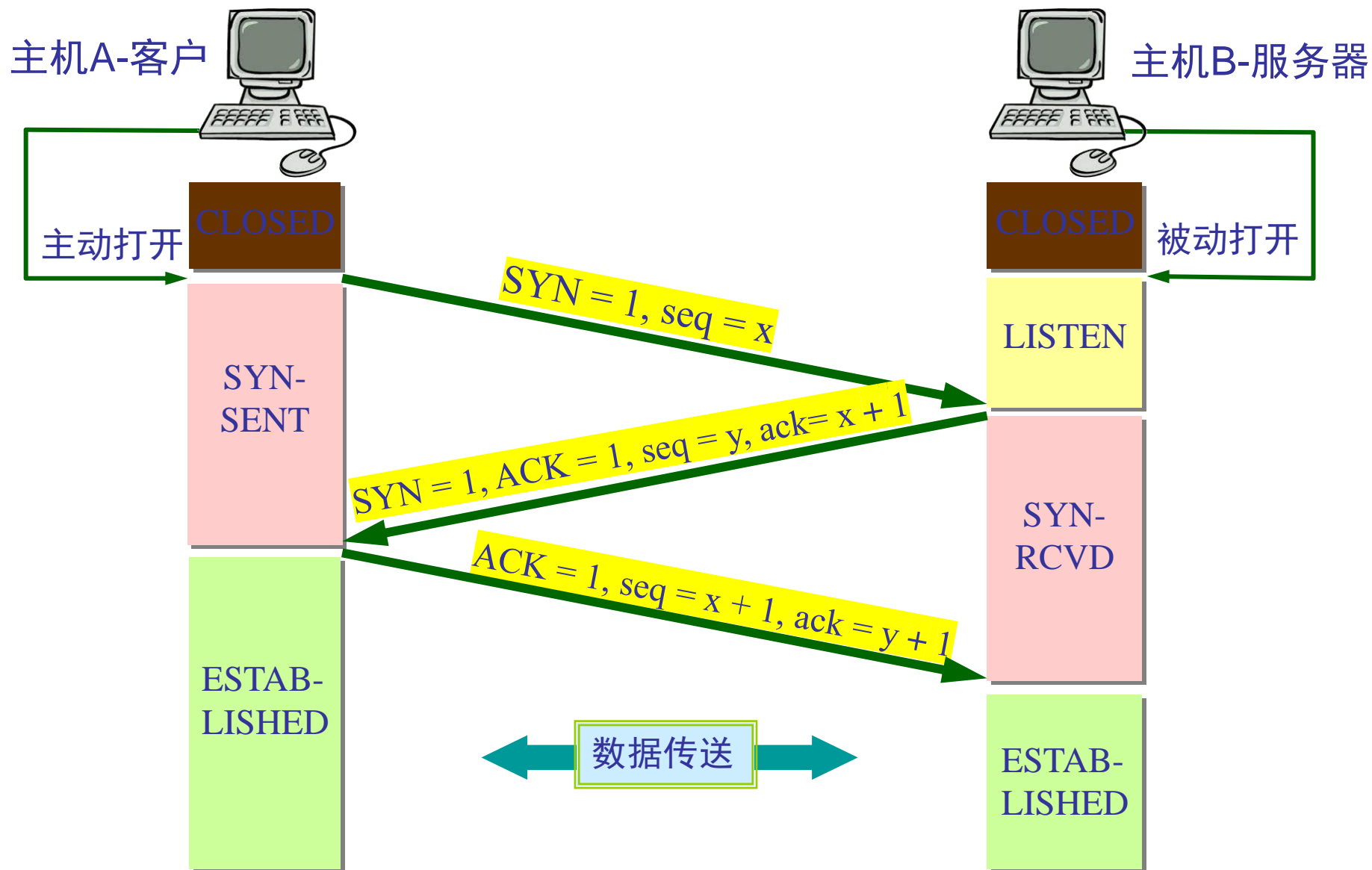
## 6.3.3 TCP连接的建立和释放

- TCP是一个面向连接的协议，通信双方在发送数据之前都必须建立一个TCP连接。
- TCP连接的建立过程称为三次握手(3-way handshake)。
- 建立连接的三次握手过程的一个关键是涉及序号的选择问题。

# 三次握手

- 第一次握手：建立连接时，源主机发送SYN包(本次连接的seq=x)到目的主机，并进入SYN\_SEND状态，等待目的主机的确认
- 第二次握手：目的主机收到SYN请求包后，如果同意连接，则发送SYN的确认(ack=x+1)给源主机，同时也发送一个自己的SYN包(seq=y)，即SYN+ACK包，此时目的主机进入SYN\_RECV状态
- 第三次握手：源主机收到目的主机的SYN+ACK包后，向源主机发送确认包ACK(ack=y+1)，此包发送后，源端和目的端进入ESTABLISHED状态，完成三次握手

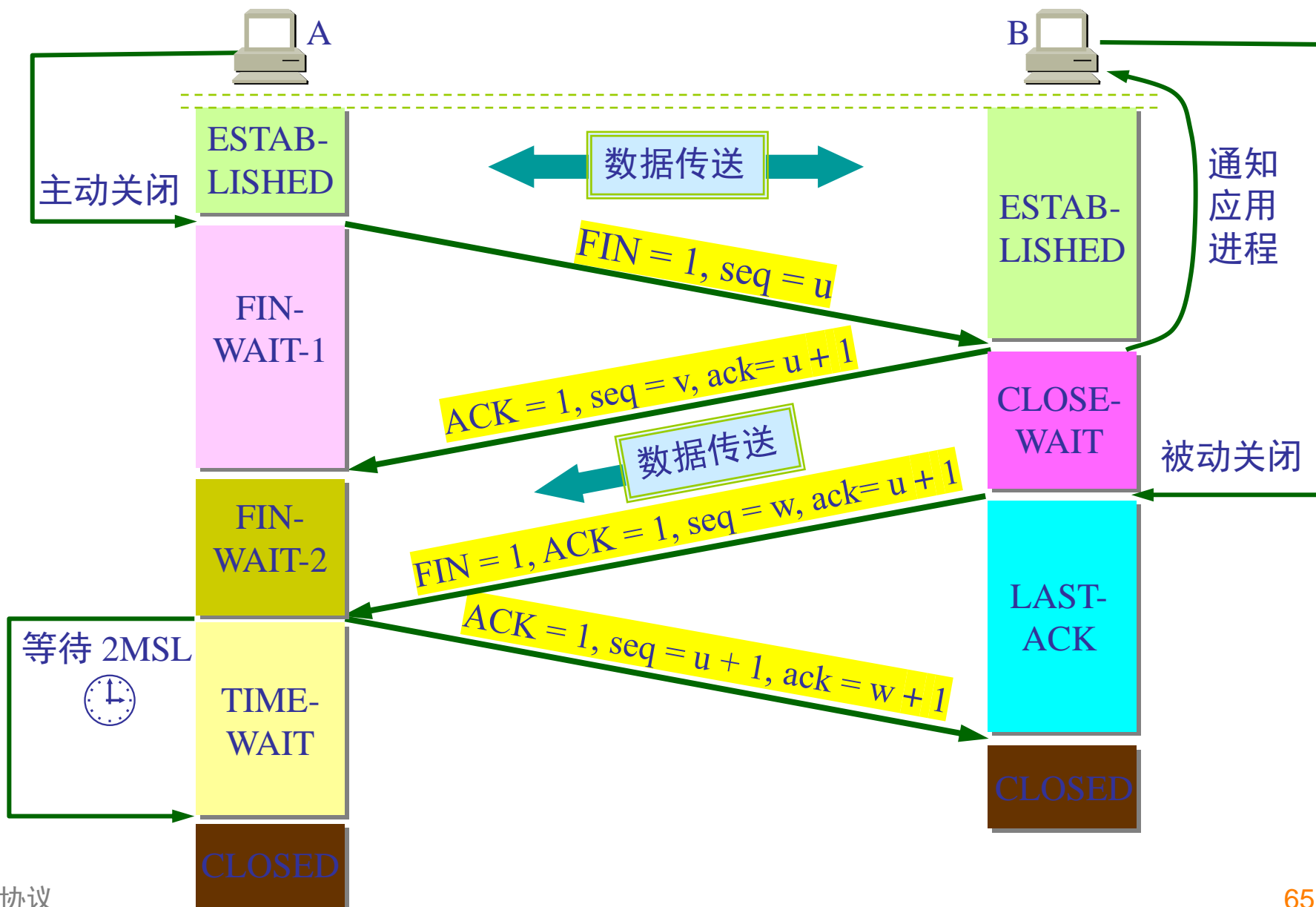
# 连接建立报文序列 背!!





# TCP连接的释放 背

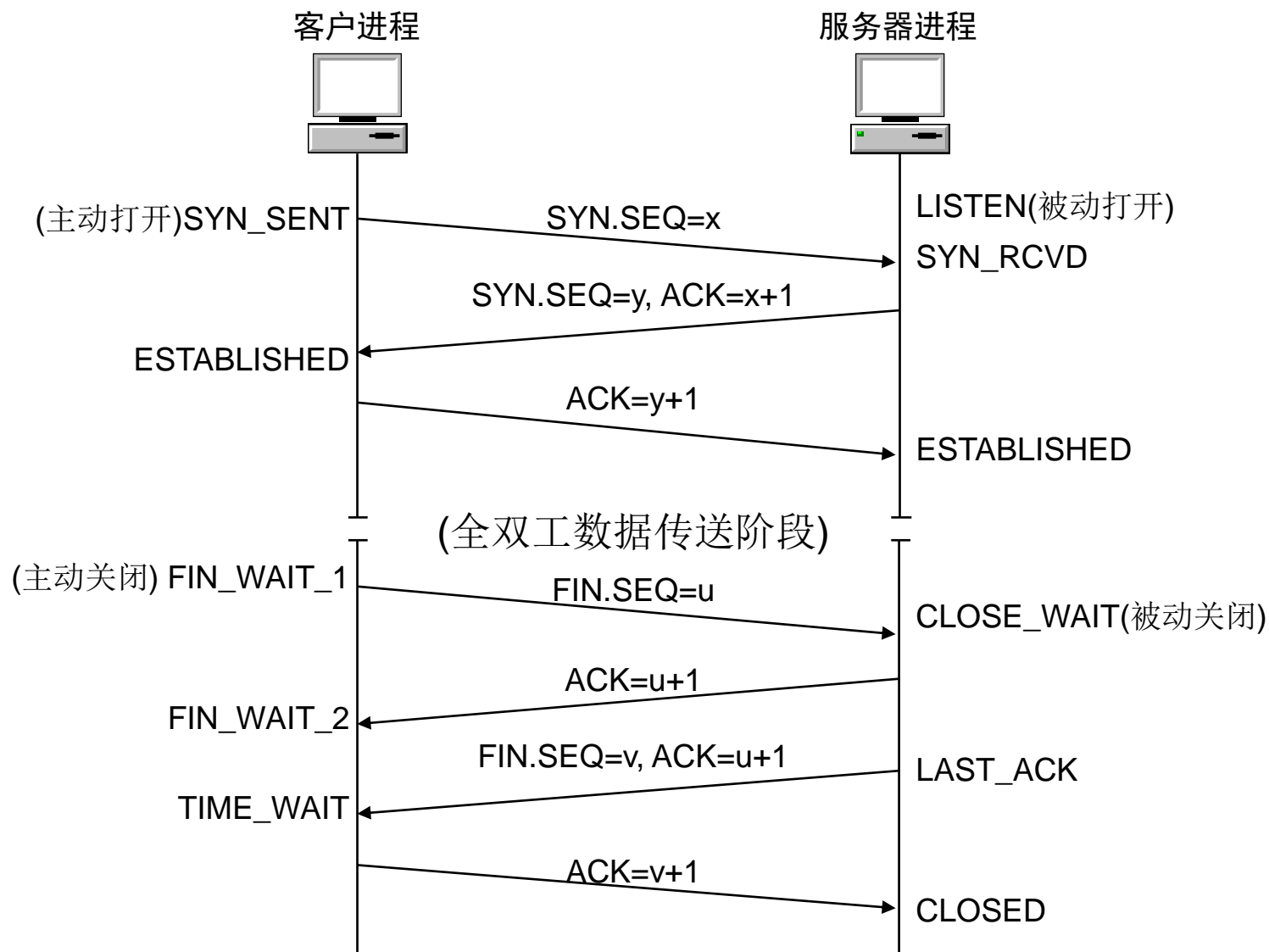
- 建立一个TCP连接需要三次握手，而释放一个TCP连接需要经过4次握手。



# A必须等待2MSL的时间

- MSL: 最长报文寿命Maximum Segment Lifetime
- 第一, 为了保证 A 发送的最后一个 ACK 报文段能够到达 B。
- 第二, 防止 “已失效的连接请求报文段” 出现在本连接中。A 在发送完最后一个 ACK 报文段后, 再经过时间 2MSL, 就可以使本连接持续的时间内所产生的所有报文段, 都从网络中消失。这样就可以使下一个新的连接中不会出现这种旧的连接请求报文段。

# TCP正常的连接建立和关闭



## 6.3.4 传输层拥塞控制

### ■ TCP拥塞控制有4种算法：

- 慢启动
- 拥塞避免
- 快速重传
- 快速恢复

### ■ 几个术语

- 拥塞窗口cwnd
- 接收端窗口rwnd
- 发送端最大数据段尺寸SMSS
- 慢启动阈值sssthresh

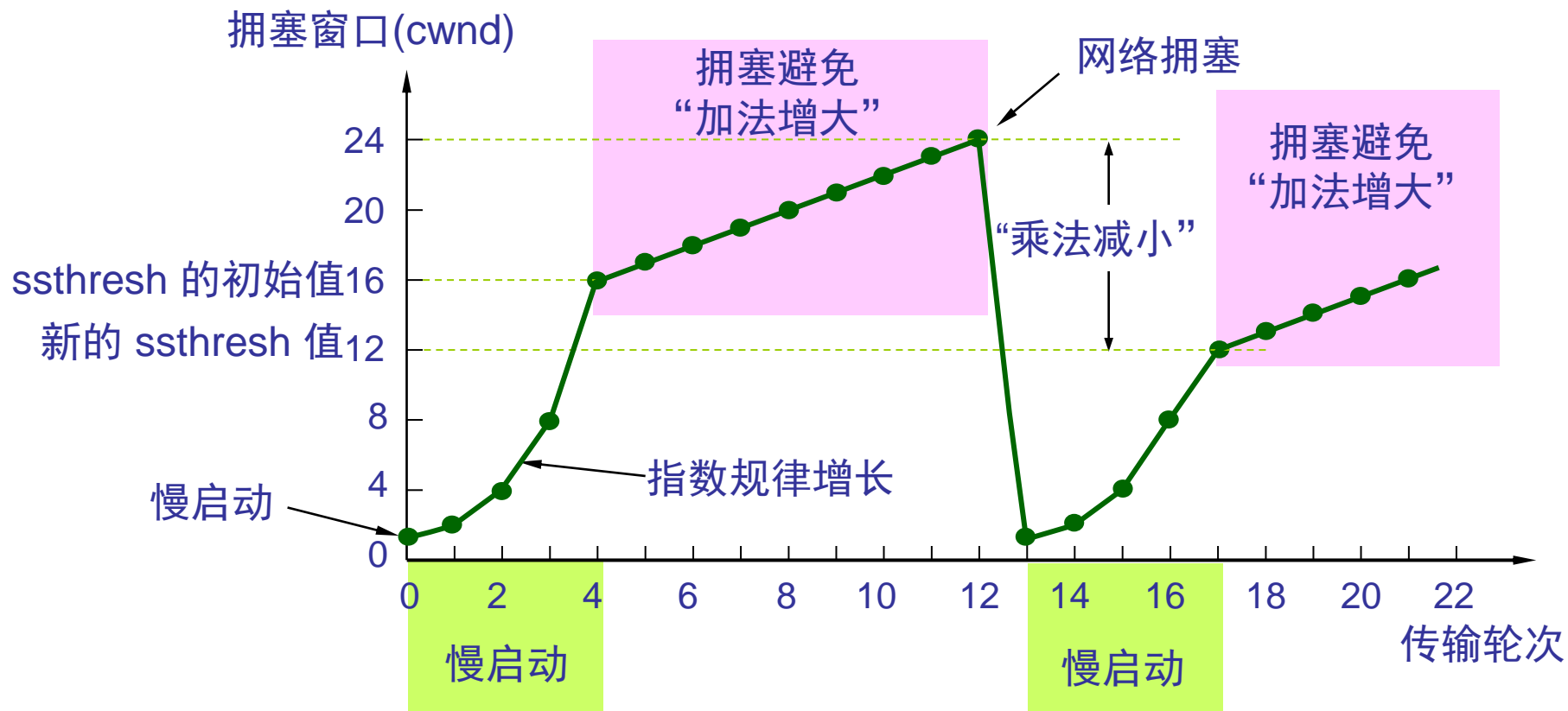
# 慢启动

- 在慢启动期间，发送方将初始的拥塞窗口(cwnd)设置为1个发送端最大数据段尺寸(SMSS)字节；在第1个超时周期内没有丢失报文的情况下，拥塞窗口(cwnd)设置为2个发送端最大数据段尺寸(SMSS)字节；在第2个超时周期内没有丢失报文的情况下，拥塞窗口(cwnd)设置为4个发送端最大数据段尺寸(SMSS)字节；在第3个超时周期内没有丢失报文的情况下，拥塞窗口(cwnd)设置为8个发送端最大数据段尺寸(SMSS)字节；依此方式，拥塞窗口(cwnd)按指数方式增长，直到拥塞窗口(cwnd)超过慢启动阈值(ssthresh)。
- 如果发生超时事件，则会重新进行慢启动阶段，初始窗口仍是一个SMSS(最大数据段长度)

# 拥塞避免

- 当拥塞窗口(cwnd)超过慢启动阈值(ssthresh)或者当拥塞窗口(cwnd)大小达到慢启动阈值(ssthresh)的大小，进入拥塞避免期间。
- 在拥塞避免期间，在没有丢失报文的情况下，拥塞窗口(cwnd)按线性方式增长，即每收到一个确认回答(ACK)，拥塞窗口(cwnd)的大小增加1个发送端最大数据段尺寸(SMSS)字节。
- 当检测到数据段丢失时，则将慢启动阈值(ssthresh)设置为当前拥塞窗口(cwnd)的一半，并重新开始慢启动算法。

# 慢启动和拥塞避免工作过程



# 乘法减小(multiplicative decrease)

- “乘法减小”是指不论在慢开始阶段还是拥塞避免阶段，只要出现一次超时（即出现一次网络拥塞），就把慢开始门限值 **ssthresh** 设置为当前的拥塞窗口值 **cwnd** 乘以 0.5。
- 当网络频繁出现拥塞时，**ssthresh** 值就下降得很快，以大大减少注入到网络中的分组数。



# 加法增大(additive increase)

- “加法增大”是指执行拥塞避免算法后，在收到对所有报文段的确认后（即经过一个往返时间），就把拥塞窗口 **cwnd** 增加一个 **SMSS** 大小，使拥塞窗口缓慢增大，以防止网络过早出现拥塞。

# 注意！！！！

- “拥塞避免”并非指完全能够避免了拥塞。利用以上的措施要完全避免网络拥塞还是不可能的。
- “拥塞避免”在拥塞避免阶段把拥塞窗口控制为按线性规律增长，使网络比较不容易出现拥塞。

# 例题

- 假定最大报文段长度是1KB，TCP拥塞窗口是16KB，并发生了超时事件。如果接着4个轮次传输都是成功的，那么该窗口将是多大？

答：发生超时后，下一次传输的是1，接着是2、4、8个报文段。所以4个轮次后的拥塞窗口是8KB。

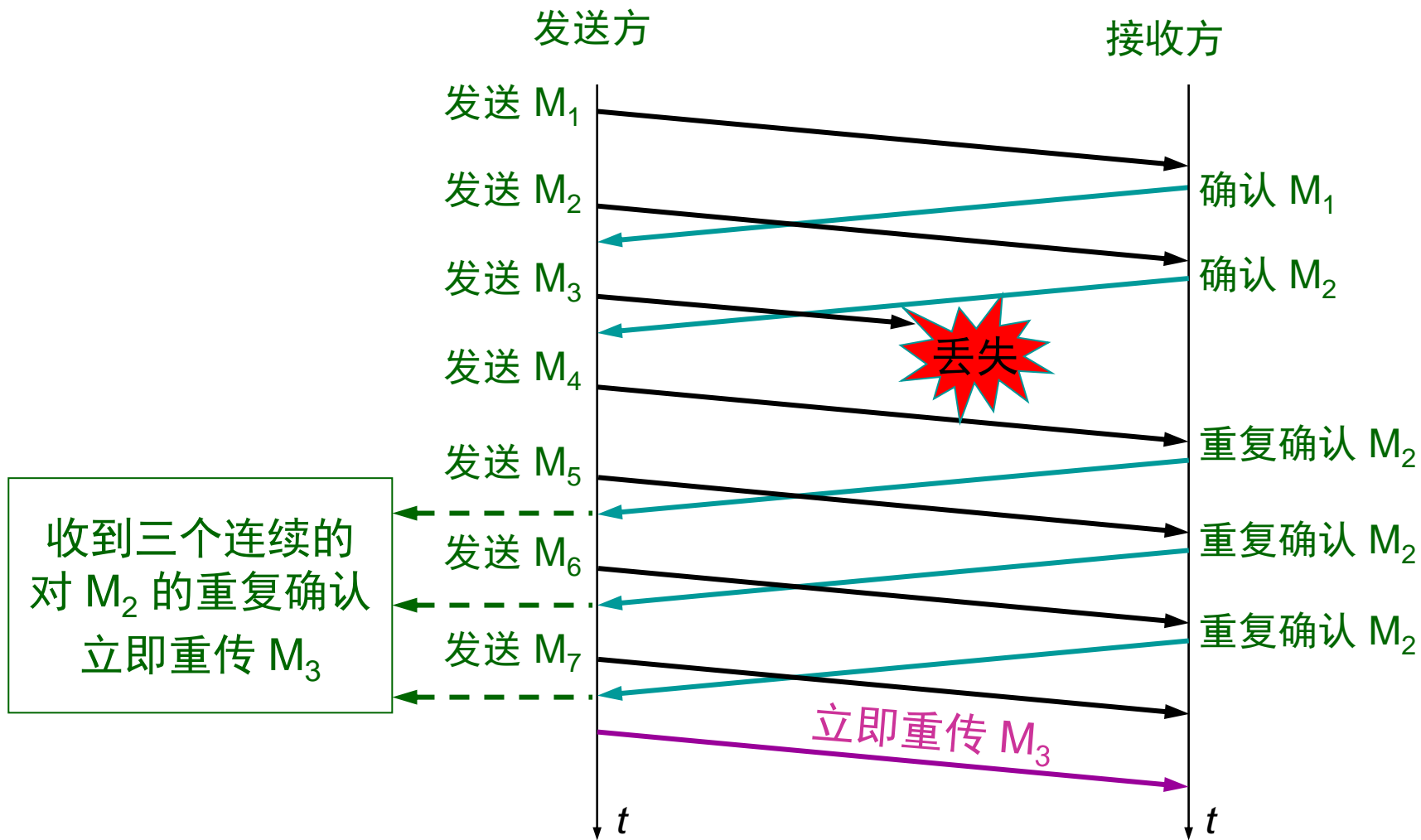


如果接着6个轮次都成功，  
CWND的值应该是多大？

# 快速重传/快速恢复

- 快重传算法首先要求接收方每收到一个失序的报文段后就立即发出重复确认。这样做可以让发送方及早知道有报文段没有到达接收方。
- 发送方只要一连收到三个重复确认就应当立即重传对方尚未收到的报文段。
- 不难看出，快重传并非取消重传计时器，而是在某些情况下可更早地重传丢失的报文段。

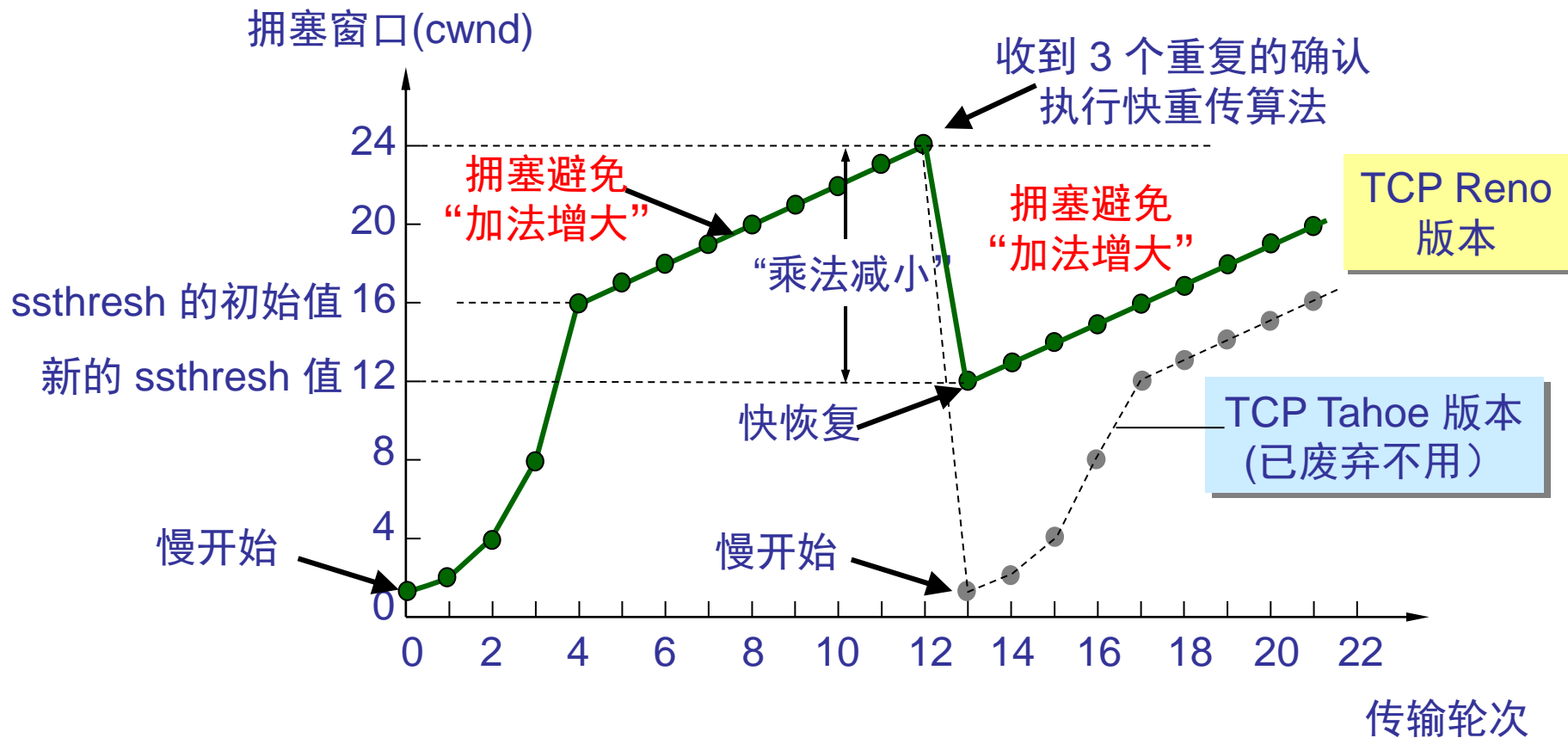
# 快速重传示例



# 快速恢复算法

- 当发送端收到连续三个重复的确认时，就执行“乘法减小”算法，把慢开始门限 **ssthresh** 减半。但接下去不执行慢开始算法。
- 由于发送方现在认为网络很可能没有发生拥塞，因此现在不执行慢开始算法，即拥塞窗口 **cwnd** 现在不设置为 1，而是设置为慢开始门限 **ssthresh** 减半后的数值，然后开始执行拥塞避免算法（“加法增大”），使拥塞窗口缓慢地线性增大。

# 连续收到三个重复的确认转入拥塞避免



# 发送窗口的上限值

- 发送方的发送窗口的上限值应当取为接收方窗口  $rwnd$  和拥塞窗口  $cwnd$  这两个变量中较小的一个，即应按以下公式确定：

发送窗口的上限值  $= \min(rwnd, cwnd)$

- 当  $rwnd < cwnd$  时，是接收方的接收能力限制发送窗口的最大值。
- 当  $cwnd < rwnd$  时，则是网络的拥塞限制发送窗口的最大值。