



哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY

立足航天，服务国防，面向国民经济主战场



《计算机网络》

第3章 传输层

聂兰顺



主要内容

本章学习目标

- ❖ 理解传输层服务
- ❖ 理解端到端原则
- ❖ 掌握传输层复用/分解方法
- ❖ 掌握UDP协议
- ❖ 掌握TCP协议
 - TCP协议特点
 - TCP段结构
 - TCP可靠数据传输
 - TCP流量控制
 - TCP连接控制
 - TCP拥塞控制
 - TCP公平性

主要内容

- ❖ 3.1 传输层服务
- ❖ 3.2 传输层多路复用/分用
- ❖ 3.3 UDP协议
- ❖ 3.4 可靠数据传输原理
- ❖ 3.5 TCP协议
 - 3.5.1 TCP段结构
 - 3.5.2 TCP可靠数据传输
 - 3.5.3 TCP流量控制
 - 3.5.4 TCP连接控制
 - 3.5.5 TCP拥塞控制
 - 3.5.6 TCP性能





哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY

立足航天，服务国防，面向国民经济主战场



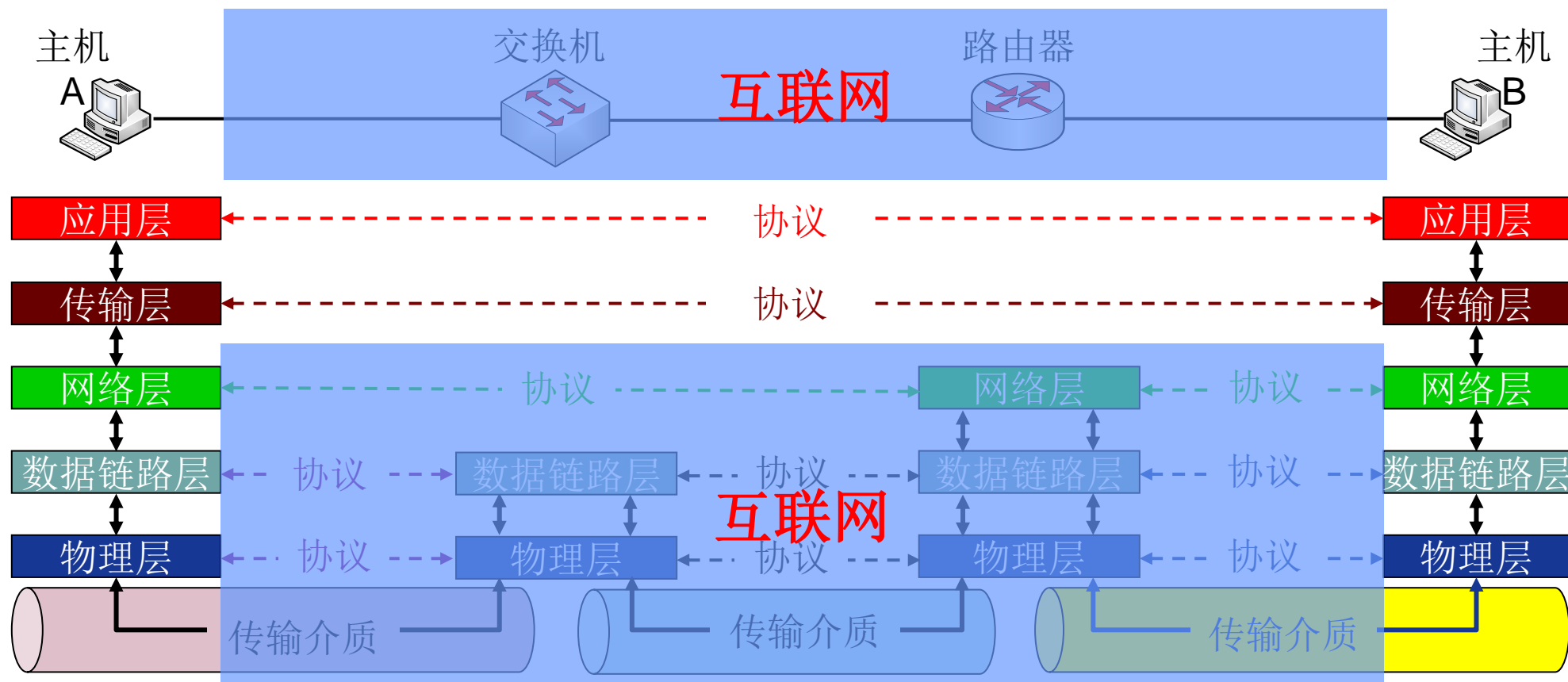
3.1 传输层服务

聂兰顺



传输层？

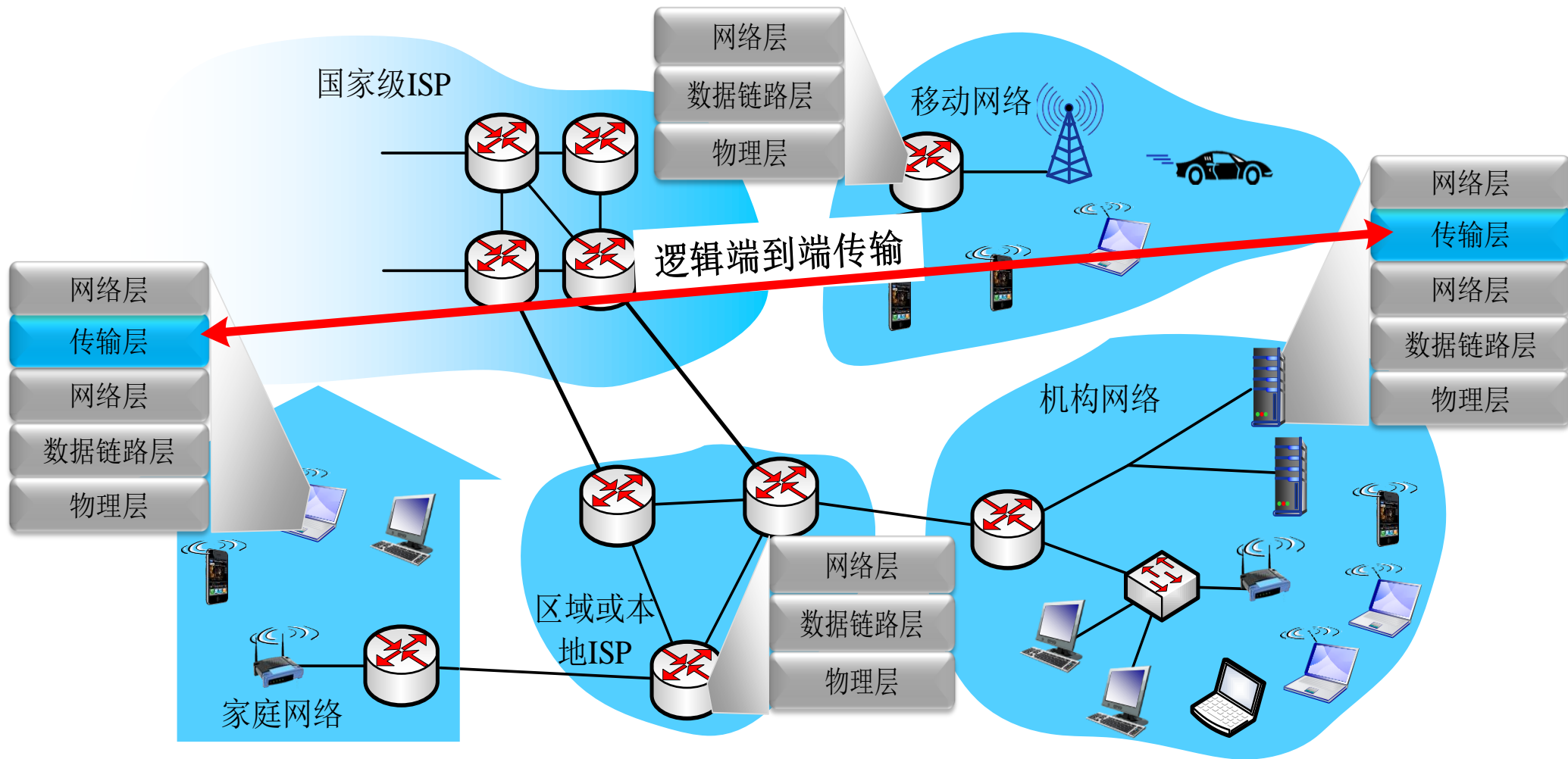
3.1 传输层服务





传输层？

3.1 传输层服务

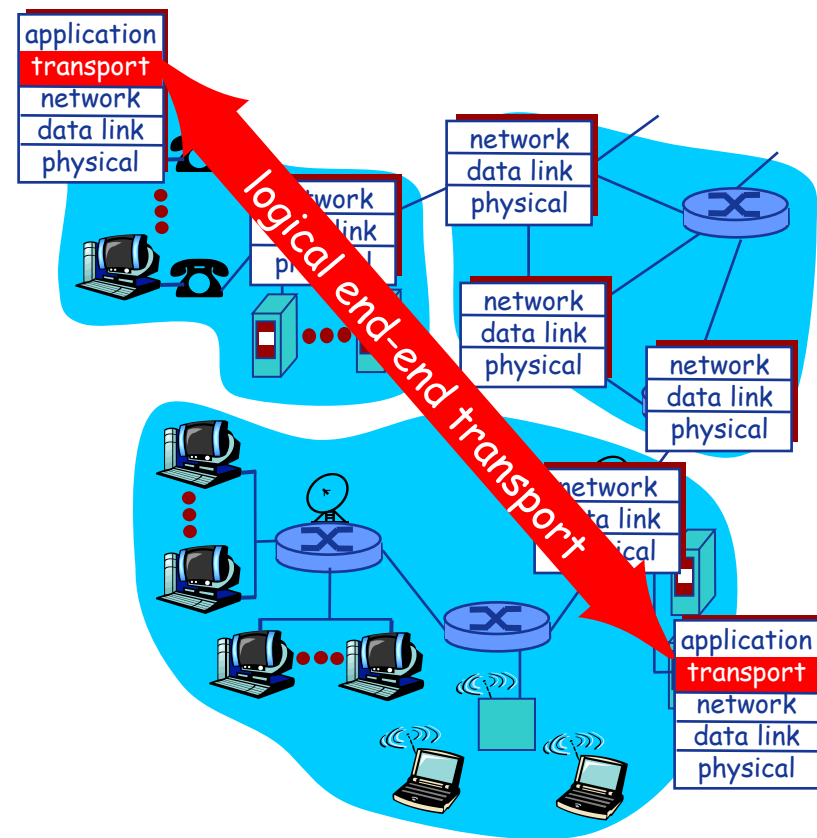




传输层服务和协议

3.1 传输层服务

- ❖ 传输层协议为运行在不同Host上的进程提供了一种**逻辑通信机制**
- ❖ 端系统运行传输层协议
 - **发送方**：将应用递交的消息分成一个或多个的Segment，并向下传给网络层。
 - **接收方**：将接收到的segment组装成消息，并向上交给应用层。
- ❖ 传输层可以为应用提供多种协议
 - Internet的TCP
 - Internet的UDP





传输层 VS. 网络层

3.1 传输层服务

- ❖ 网络层：提供**主机**之间的逻辑通信机制
- ❖ 传输层：提供**应用进程**之间的逻辑通信机制
 - 位于网络层之上
 - 依赖于网络层服务
 - 对网络层服务进行（可能的）增强

家庭类比:

12个孩子给12个孩子发信

- ❖ 应用进程 = 孩子
- ❖ 应用消息 = 信封里的信
- ❖ 主机 = 房子
- ❖ 传输层协议 = 李雷和韩梅梅
- ❖ 网络层协议 = 邮政服务





Internet传输层协议

3.1 传输层服务

❖ 可靠、按序的交付服务(TCP)

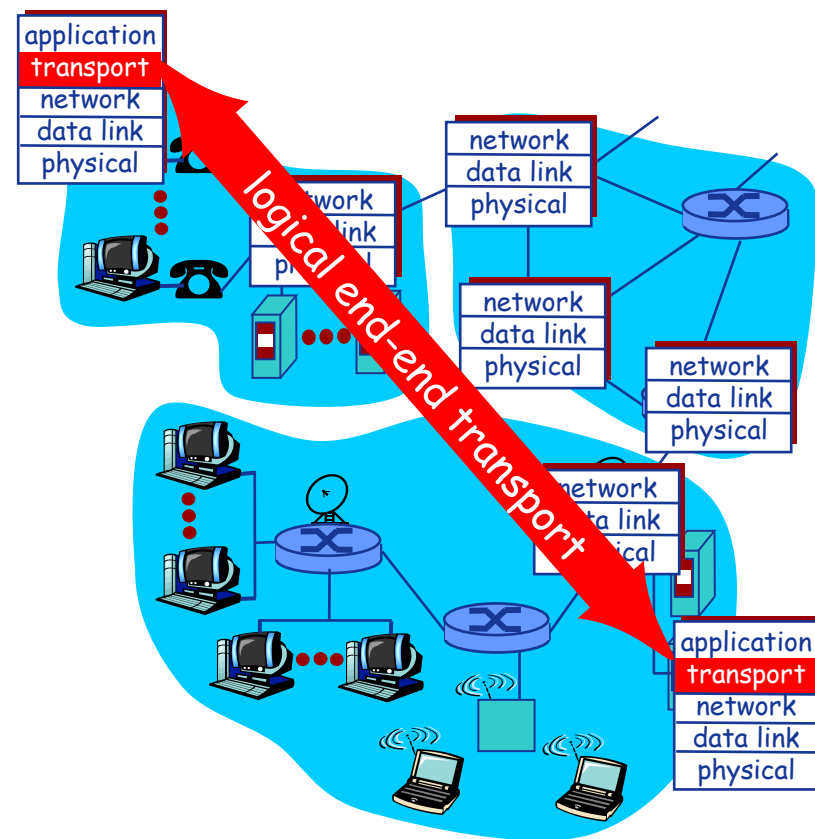
- 拥塞控制
- 流量控制
- 连接建立

❖ 不可靠的交付服务(UDP)

- 基于“**尽力而为 (Best-effort)**”的网络层，没有做（可靠性方面的）扩展

❖ 两种服务均不保证

- 延迟
- 带宽





哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY

立足航天，服务国防，面向国民经济主战场



3.2 传输层多路复用/分用

聂兰顺



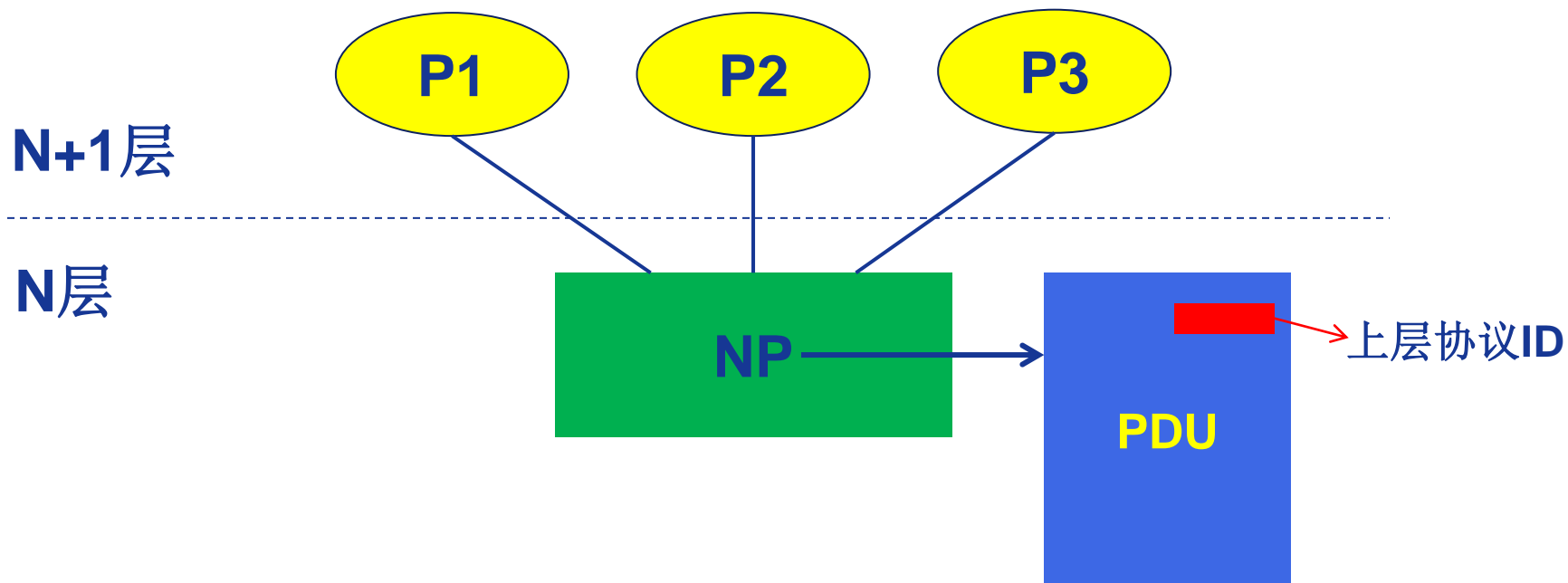
为什么需要多路复用/分用？

3.1 传输层服务

3.2 传输层多路复用/分解

Q: 为什么需要实现复用与分解？如何实现复用与分解？

A: 如果某层的一个协议/实体直接为上层的多个协议/实体提供服务，则需要复用/分用



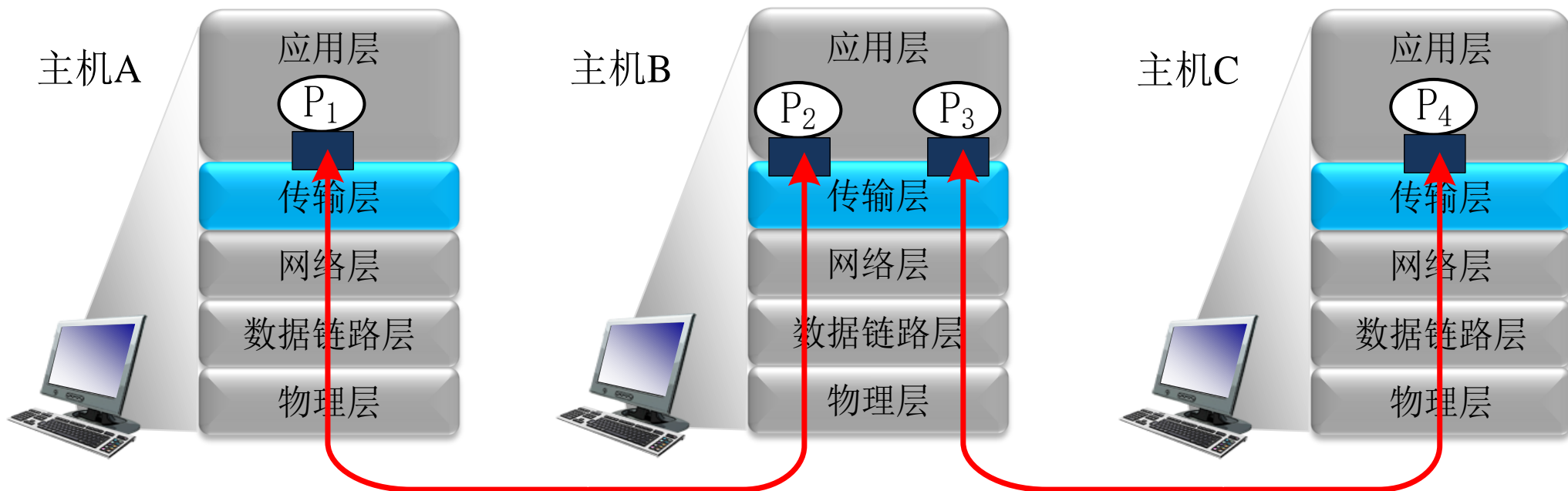
➤ 复用与分解只在传输层进行吗？



传输层多路复用/分用？

3.1 传输层服务

3.2 传输层多路复用/分解



图例：○ 进程 ■ 套接字

- 传输层如何实现复用与分解功能？
- 可能通过其他方式实现复用与分解吗？





传输层多路复用/分用

3.1 传输层服务

3.2 传输层多路复用/分解

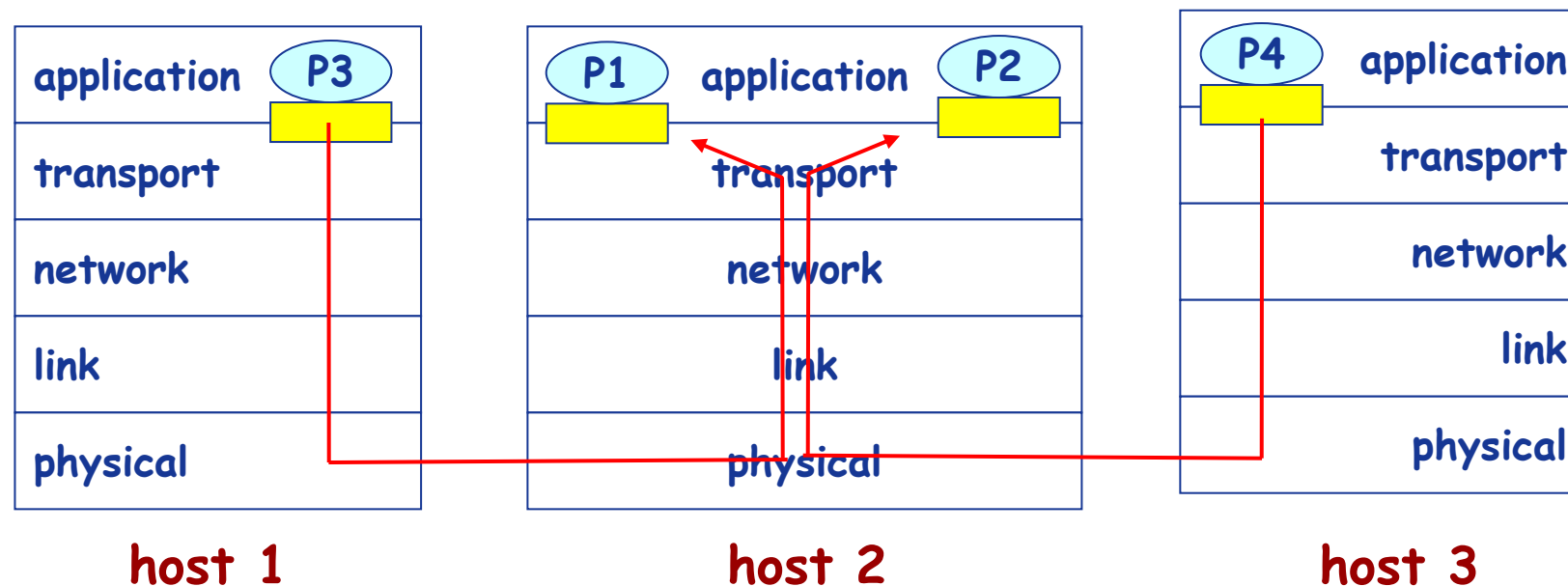
发送端进行多路复用:

从多个**Socket**接收数据, 为每块数据封装上头部信息, 生成**Segment**, 交给网络层

接收端进行多路分用:

传输层依据头部信息将收到的**Segment**交给正确的**Socket**, 即不同的进程

■ = socket ○ = process





传输层分用如何工作？

3.1 传输层服务

3.2 传输层多路复用/分解

❖ 主机接收到IP数据报(datagram)

- 每个数据报携带源IP地址、目的IP地址。
- 每个数据报携带一个传输层的段(Segment)。
- 每个段携带源端口号和目的端口号

❖ 主机收到Segment之后，传输层协议提取IP地址和端口号信息，将Segment导向相应的Socket

- TCP做更多处理



TCP/UDP 段格式



传输层无连接分用

3.1 传输层服务

3.2 传输层多路复用/分解

❖ 创建Socket并绑定端口号

```
DatagramSocket mySocket1 = new  
    DatagramSocket(9911);
```

```
DatagramSocket mySocket2 = new  
    DatagramSocket(9922);
```

❖ UDP的Socket用二元组标识

- (目的IP地址, 目的端口号)

❖ 主机收到UDP段后

- 检查段中的目的端口号
- 将UDP段导向绑定在该端口号的Socket

❖ 来自不同源IP地址和/或源端口号的IP数据包被导向同一个Socket

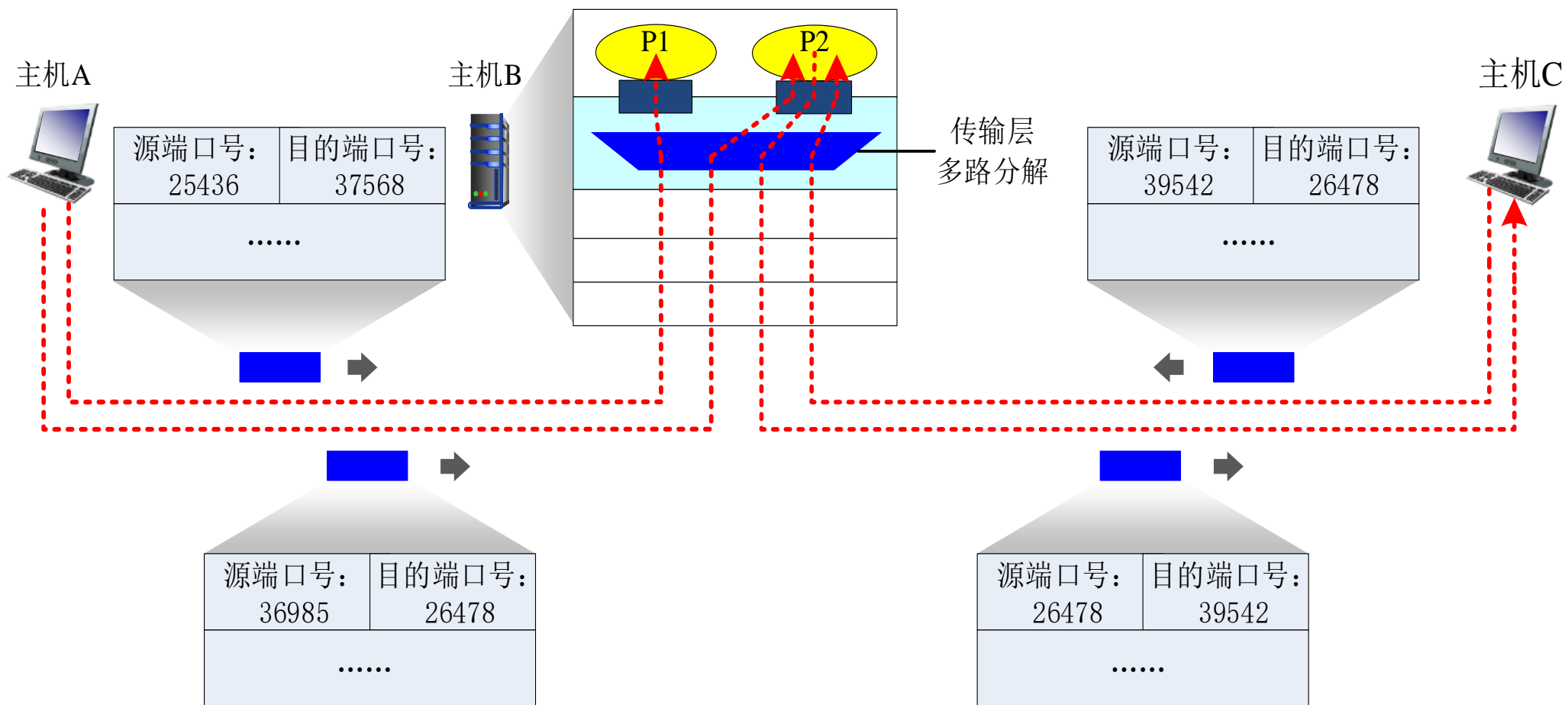




传输层无连接分用

3.1 传输层服务

3.2 传输层多路复用/分解





传输层面向连接的分用

3.1 传输层服务

3.2 传输层多路复用/分解

- ❖ TCP的Socket用四元组标识
 - 源IP地址
 - 源端口号
 - 目的IP地址
 - 目的端口号
- ❖ 接收端利用所有的四个值将Segment导向正确的Socket
- ❖ 服务器可能同时支持多个TCP Socket
 - 每个Socket用自己的四元组标识
- ❖ Web服务器为每个客户端创建不同的Socket

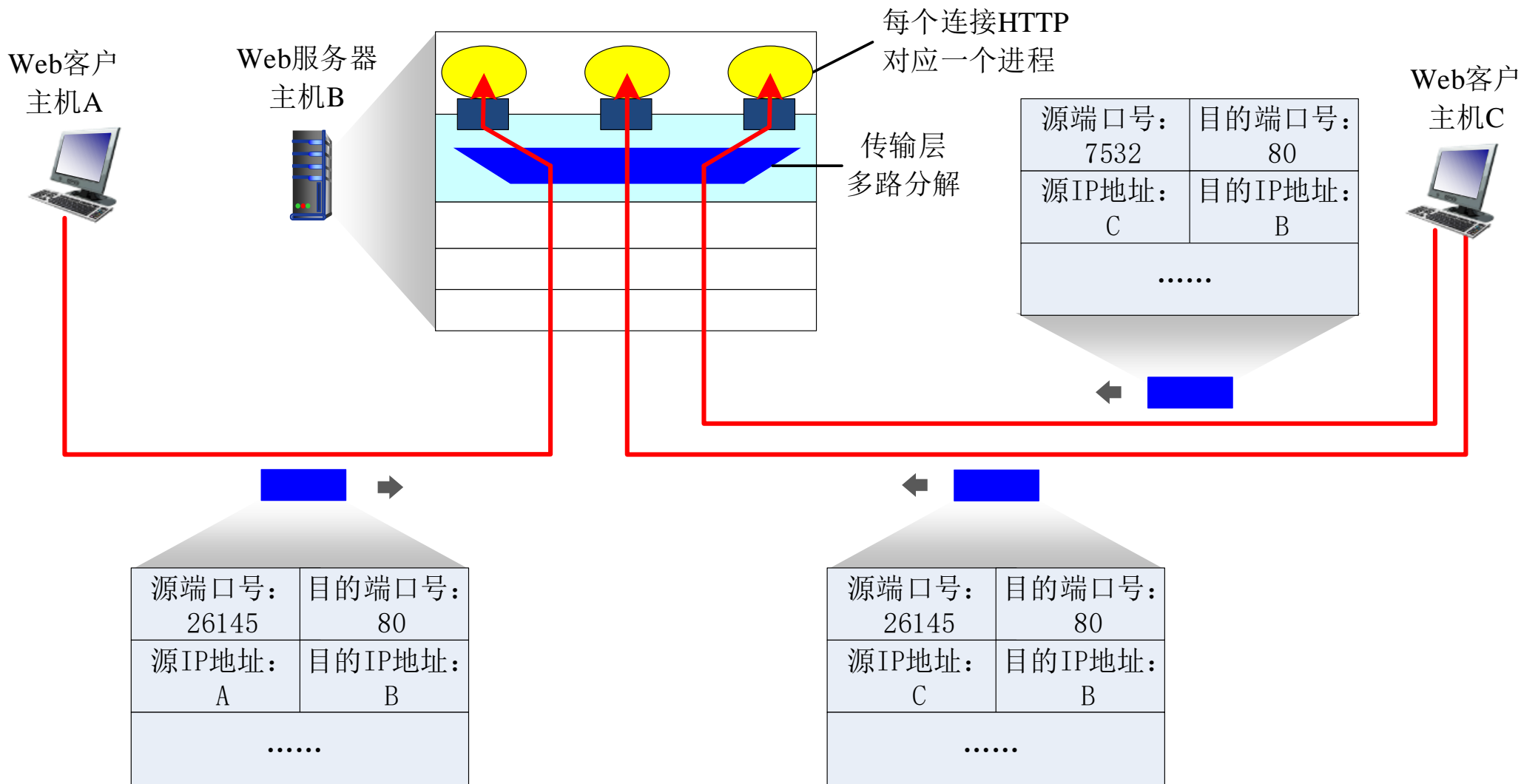




传输层面向连接的分用

3.1 传输层服务

3.2 传输层多路复用/分解





哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY

立足航天，服务国防，面向国民经济主战场



3.3 UDP协议

聂兰顺



UDP: 用户数据报协议[RFC 768]

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

❖ 基于Internet IP协议

- 复用/分用
- 简单的错误校验

❖ “Best effort” 服务，UDP段可能

- 丢失
- 非按序到达

❖ 无连接

- UDP发送方和接收方之间不需要握手
- 每个UDP段的处理独立于其他段

为什么需要UDP?

- ❖ 无需建立连接
(减少延迟)
- ❖ 实现简单: 无需维护连接状态
- ❖ 头部开销少
- ❖ 没有拥塞控制: 应用可更好地控制发送时间和速率





UDP: 用户数据报协议[RFC 768]

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

❖ 常用于流媒体应用

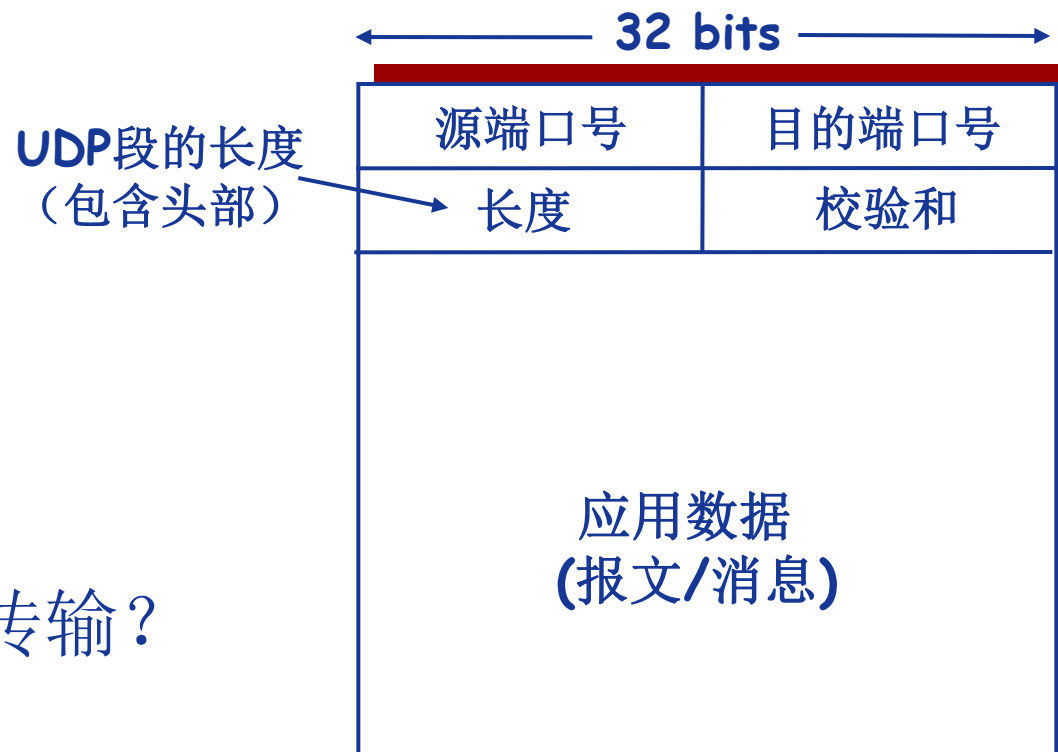
- 容忍丢失
- 速率敏感

❖ UDP还用于

- DNS
- SNMP

❖ 在UDP上实现可靠数据传输？

- 在应用层增加可靠性机制
- 应用特定的错误恢复机制
 - 例如：停等协议、滑动窗口协议



UDP报文段格式



UDP校验和(checksum)

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

目的：检测**UDP**段在传输中是否发生错误（如位翻转）

❖ 发送方

- 将参与校验和计算的所有内容视为**16-bit**整数序列
- 校验和计算：
 - 计算整数序列的**和 (sum)**
 - 进位也要加在和的后面
 - 将和按位求反（即反码），得到**校验和 (checksum)**
- 发送方将校验和放入校验和字段

❖ 接收方

- 针对收到的**UDP**报文段，按发送方同样的方法构建**16位**整数序列
- 按相同算法计算整数序列的**和 (sum)**
- 若**sum=1111111111111111**，则无错；否则，有错



UDP校验和(checksum)

3.1 传输层服务

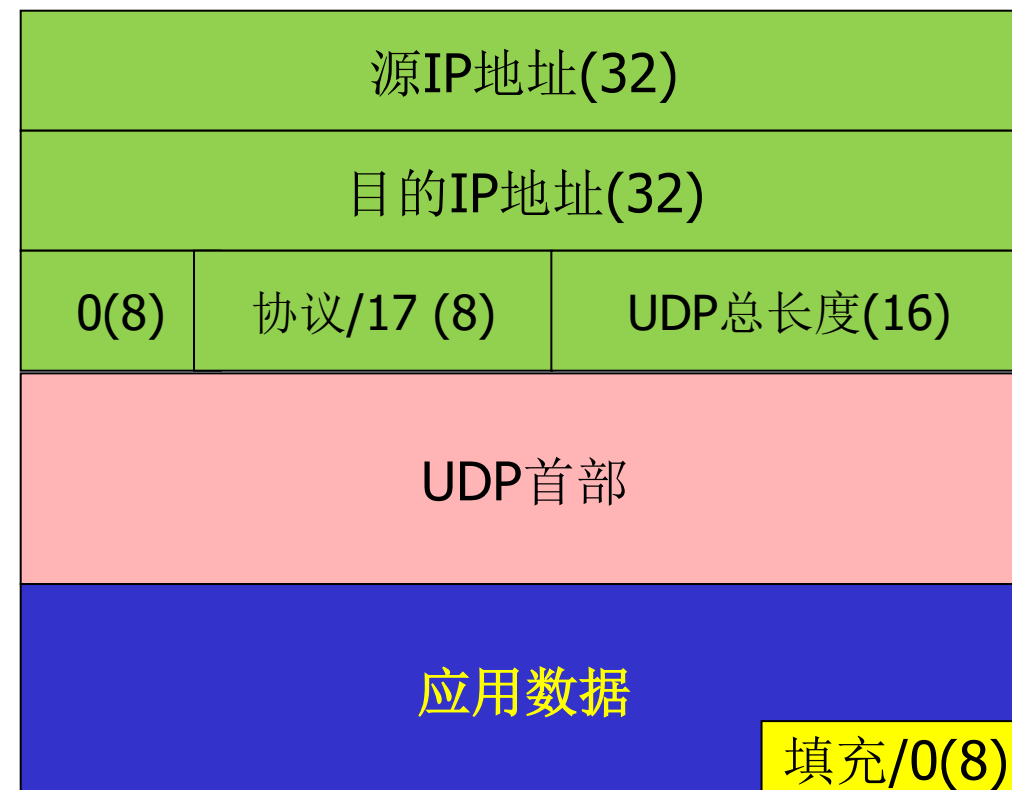
3.2 传输层多路复用/分解

3.3 UDP协议

❖ 3 部分:

- 伪首部
(Pseudo head)
- UDP首部
- 应用数据

伪首部





校验和计算示例

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

❖ 示例:

	1	1	1	0	0	1	1	0	0	1	1	0	0	1	1	0
	1	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
<hr/>																
回卷	1	1	0	1	1	1	0	1	1	1	0	1	1	1	0	1
	<hr/>															1
sum	1	0	1	1	1	0	1	1	1	0	1	1	1	1	0	0
checksum	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	1

❖ 注意:

- 最高位进位必须被加进去



哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY

立足航天，服务国防，面向国民经济主战场



3.4 可靠数据传输原理

聂兰顺

可靠数据传输原理

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 什么是可靠？

- 不错、不丢
不乱、不多

❖ 可靠数据传输协议

- 可靠数据传输对应应用层、传输层、链路层都很重要
- 网络Top-10问题
- 信道的不可靠特性决定了可靠数据传输协议(rdt)的复杂性

application layer
transport layer



(a) provided service



可靠数据传输原理

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 什么是可靠？

- 不错、不丢
不乱、不多

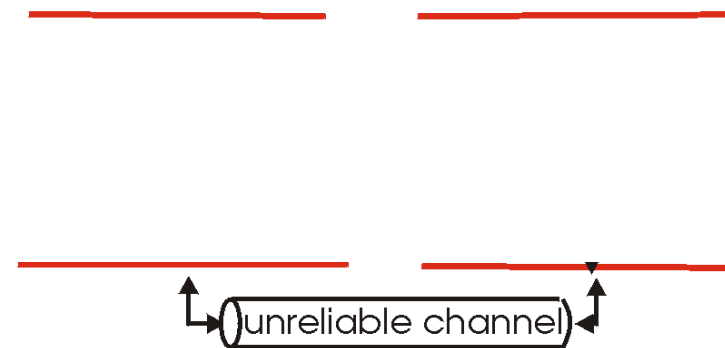
❖ 可靠数据传输协议

- 可靠数据传输对应用层、传输层、链路层都很重要
- 网络Top-10问题
- 信道的不可靠特性决定了可靠数据传输协议(rdt)的复杂性

application layer
transport layer



(a) provided service



(b) service implementation



可靠数据传输原理

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

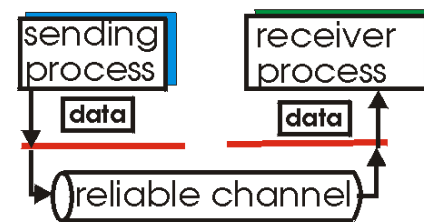
❖ 什么是可靠？

- 不错、不丢
不乱、不多

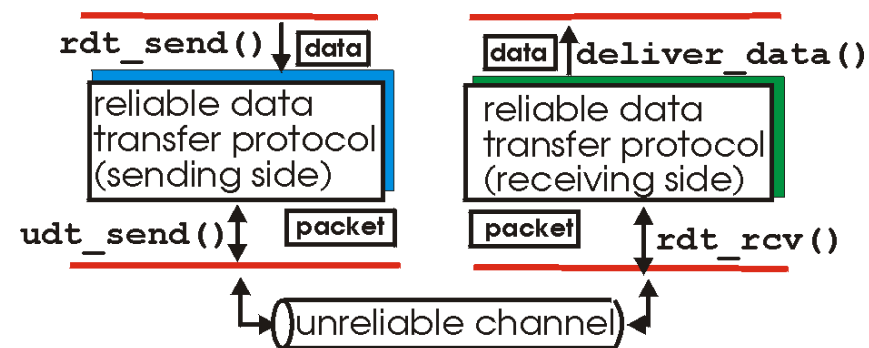
❖ 可靠数据传输协议

- 可靠数据传输对应用层、传输层、链路层都很重要
- 网络Top-10问题
- 信道的不可靠特性决定了可靠数据传输协议(rdt)的复杂性

application layer
transport layer



(a) provided service



(b) service implementation





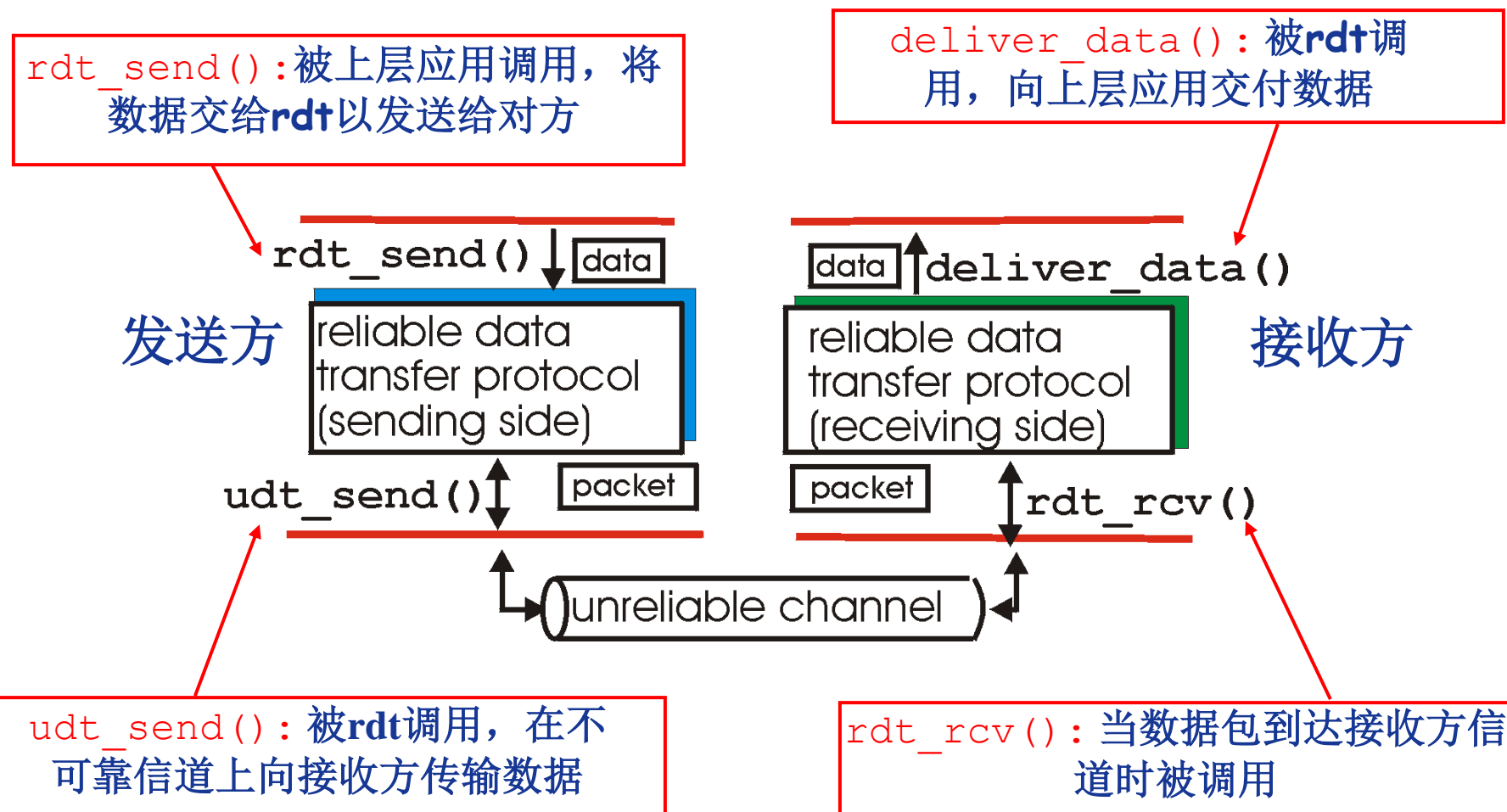
可靠数据传输协议基本结构:接口

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



可靠数据传输协议

3.1 传输层服务

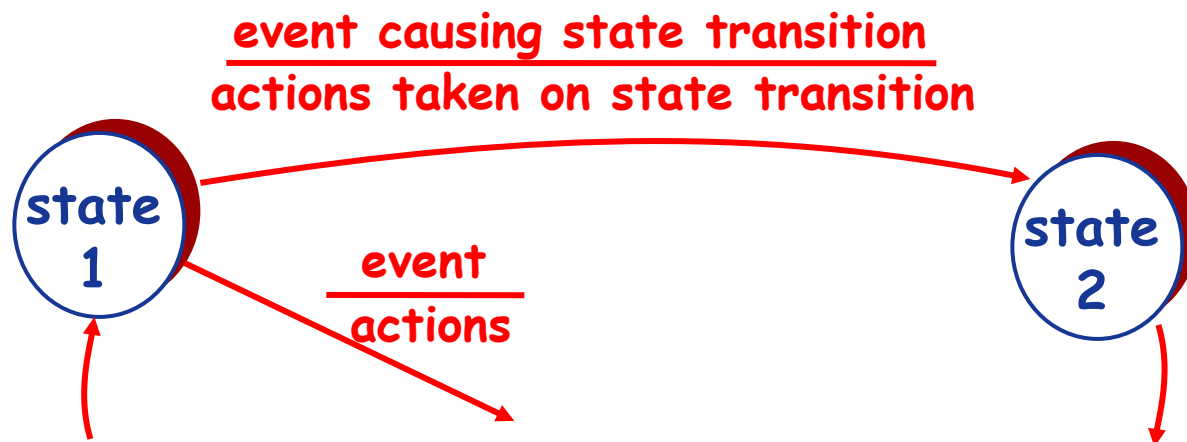
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

- ❖ 渐进地设计可靠数据传输协议的发送方和接收方
- ❖ 只考虑单向数据传输
 - 但控制信息双向流动
- ❖ 利用状态机(Finite State Machine, FSM)刻画传输协议

state: when in this "state" next state uniquely determined by next event





Rdt 1.0: 可靠信道上的可靠数据传输

3.1 传输层服务

3.2 传输层多路复用/分解

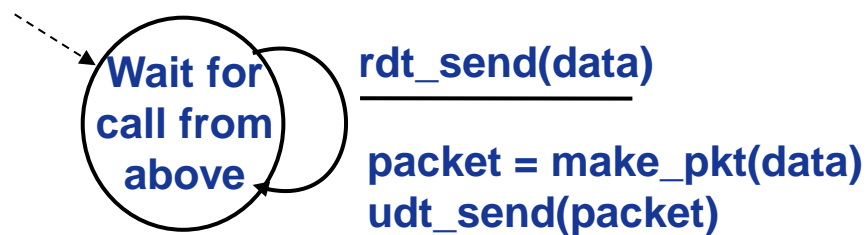
3.3 UDP协议

3.4 可靠数据传输原理

❖ 底层信道完全可靠

- 不会发生错误(bit error)
- 不会丢弃分组

❖ 发送方和接收方的FSM独立



sender



receiver





Rdt 2.0: 产生位错误的信道

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 底层信道可能翻转分组中的位(bit)

- 利用**校验和**检测位错误

❖ 如何从错误中恢复?

- **确认机制(Acknowledgements, ACK)**: 接收方显式地告知发送方分组已正确接收
- **NAK**: 接收方显式地告知发送方分组有错误
- 发送方收到**NAK**后, **重传**分组

❖ 基于这种重传机制的rdt协议称为**ARQ**(Automatic Repeat reQuest)协议

❖ Rdt 2.0中引入的新机制

- 差错检测
- 接收方反馈控制消息: **ACK/NAK**
- 重传





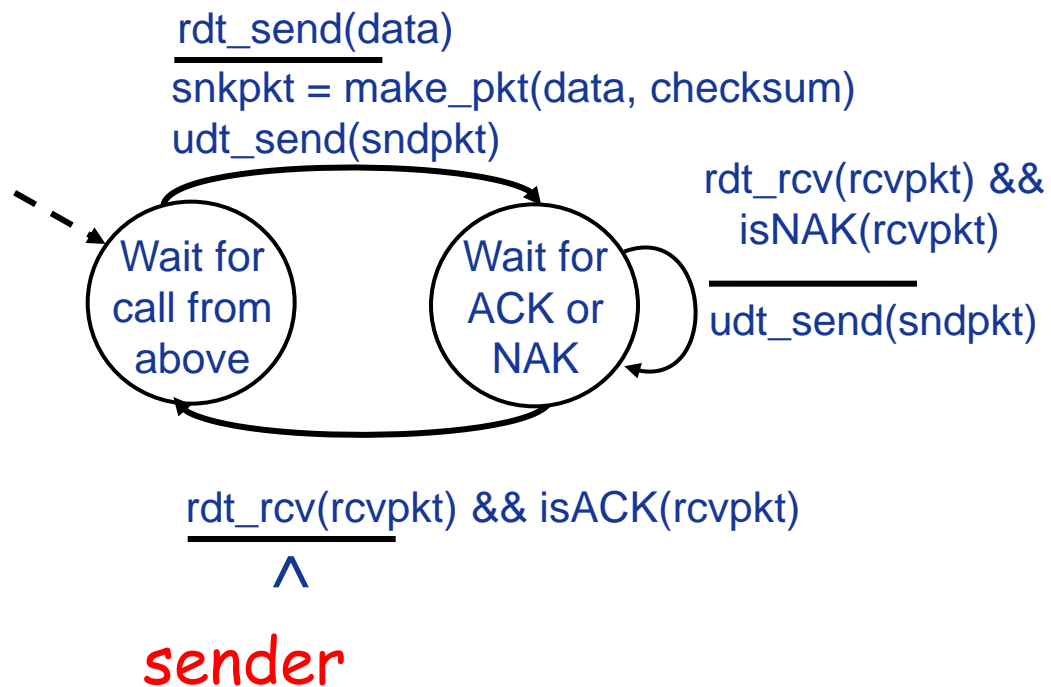
Rdt 2.0: FSM规约

3.1 传输层服务

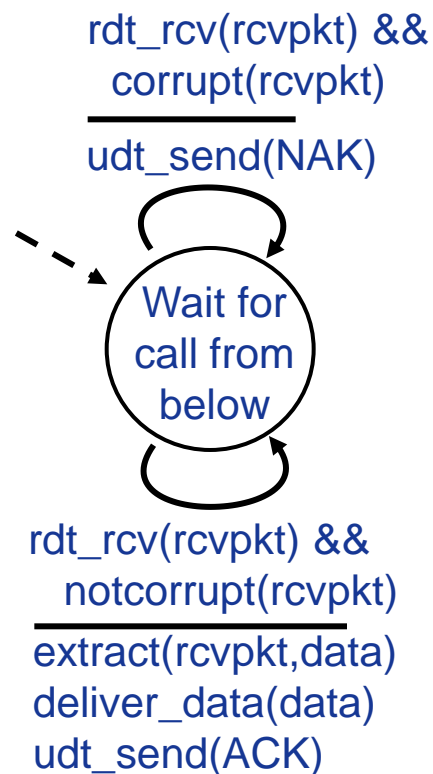
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



receiver





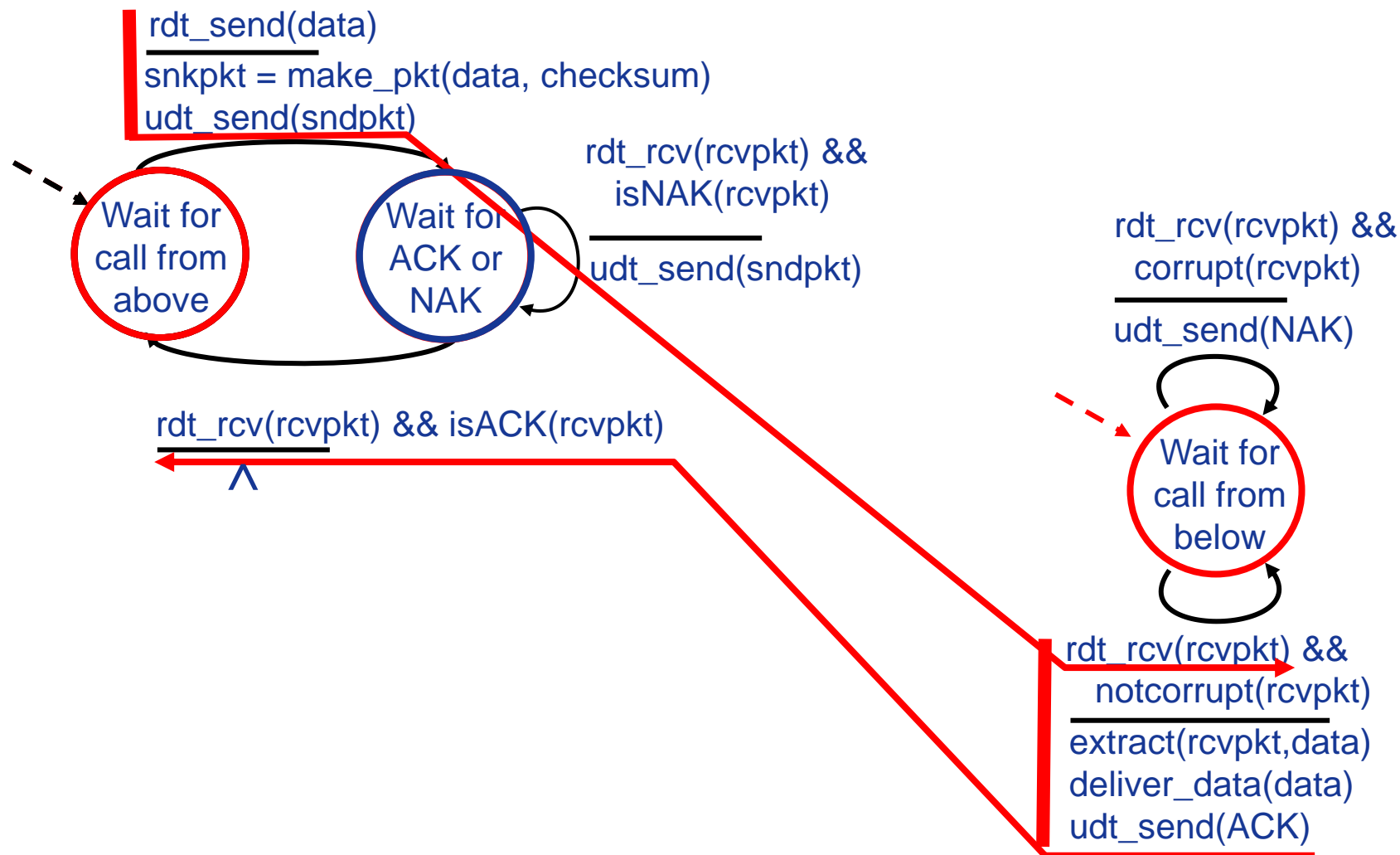
Rdt 2.0: 无错误场景

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



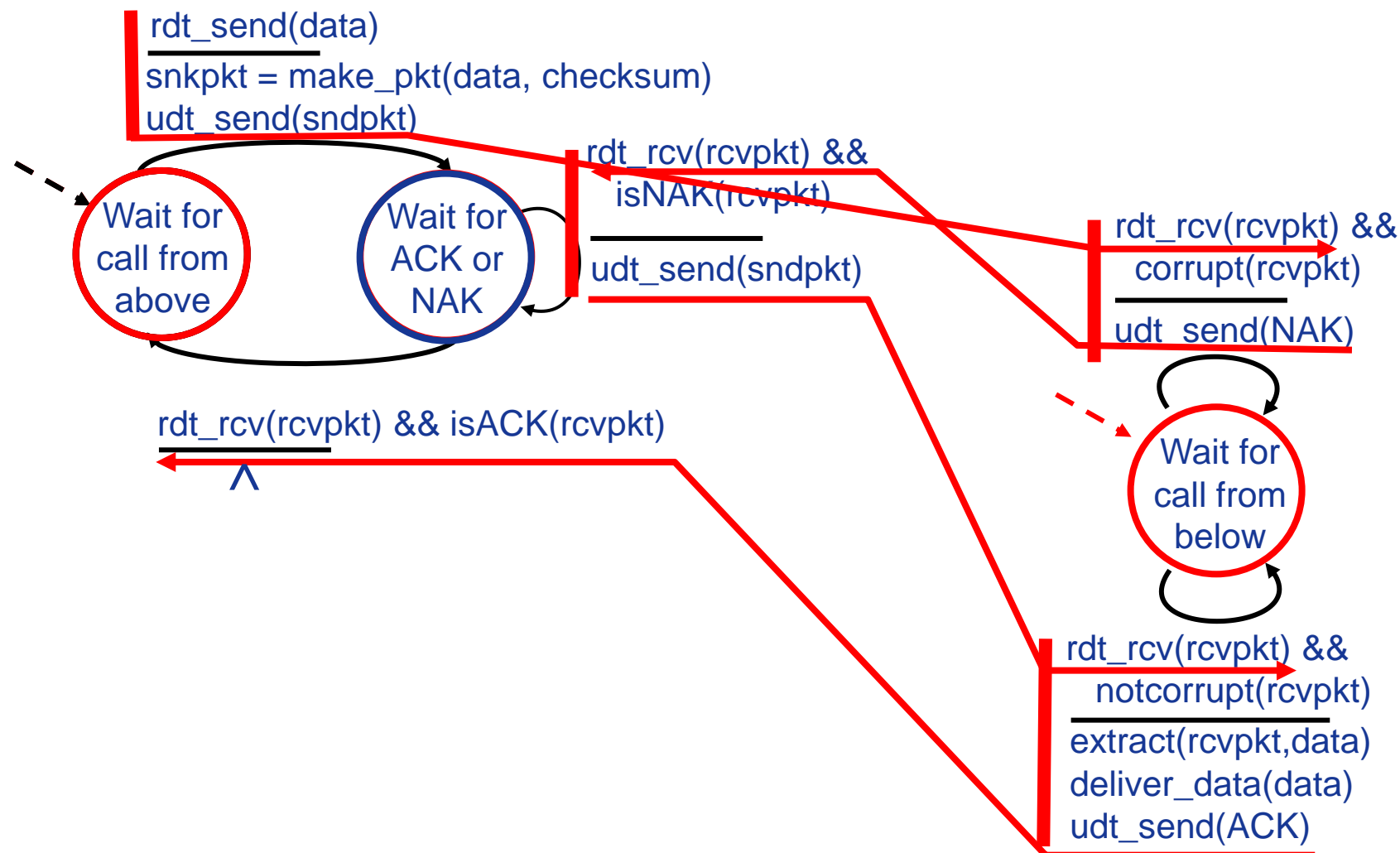


Rdt 2.0: 有错误场景

3.1 数据链路层基本服务

3.2 差错检测与纠正

3.3 可靠数据传输原理





Rdt 2.0有什么缺陷？



3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 如果ACK/NAK消息发生错误/被破坏(corrupted)会怎么样？

- 为ACK/NAK增加校验和，检错并纠错
- 发送方收到被破坏ACK/NAK时不知道接收方发生了什么
- 如果ACK/NAK坏掉，发送方重传
- 不能简单的重传：产生重复分组

❖ 如何解决重复分组问题？

- 序列号(Sequence number): 发送方给每个分组增加序列号
- 接收方丢弃重复分组

stop and wait

Sender sends one packet,
then waits for receiver
response





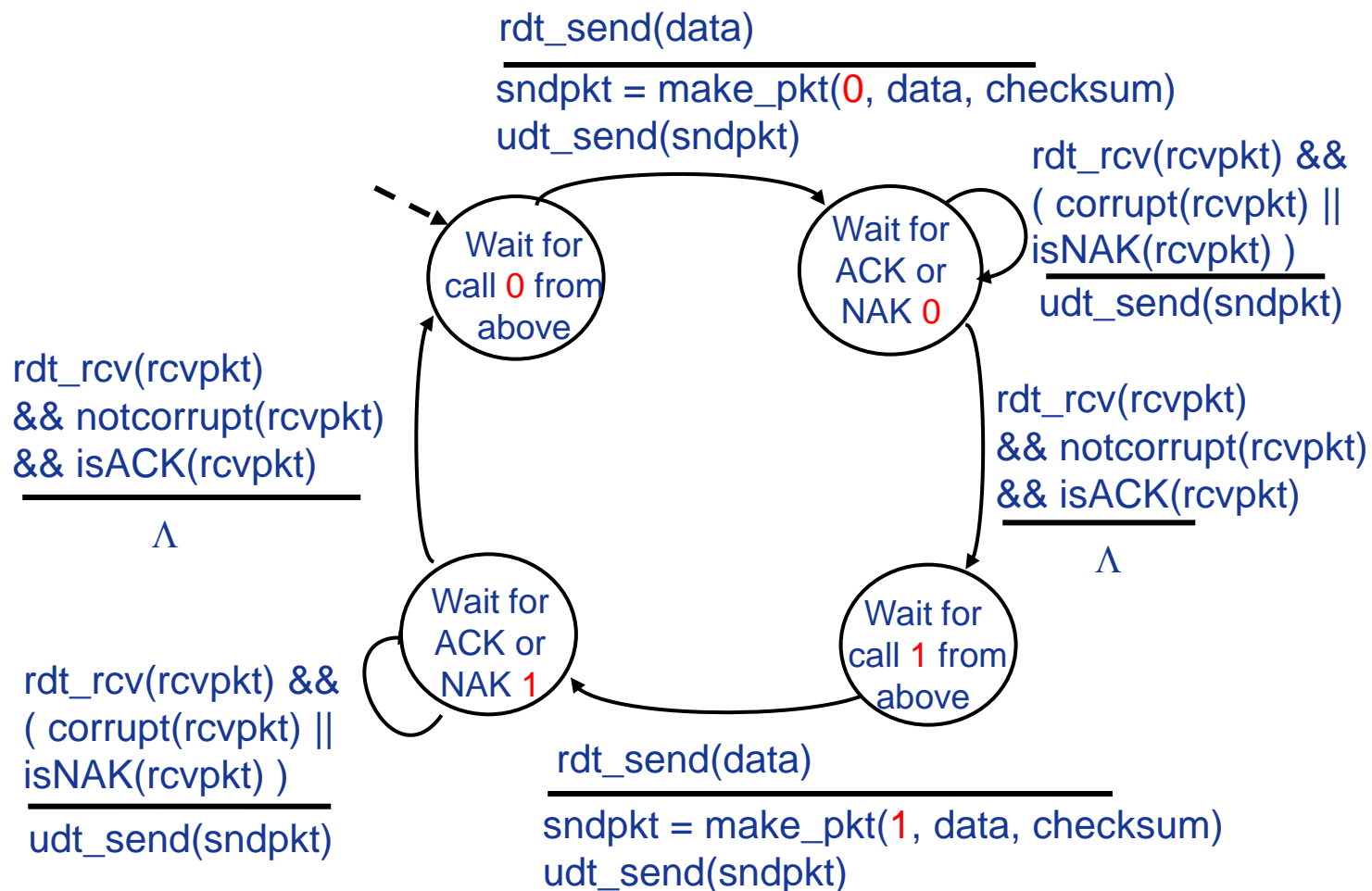
Rdt 2.1: 应对ACK/NAK破坏

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理





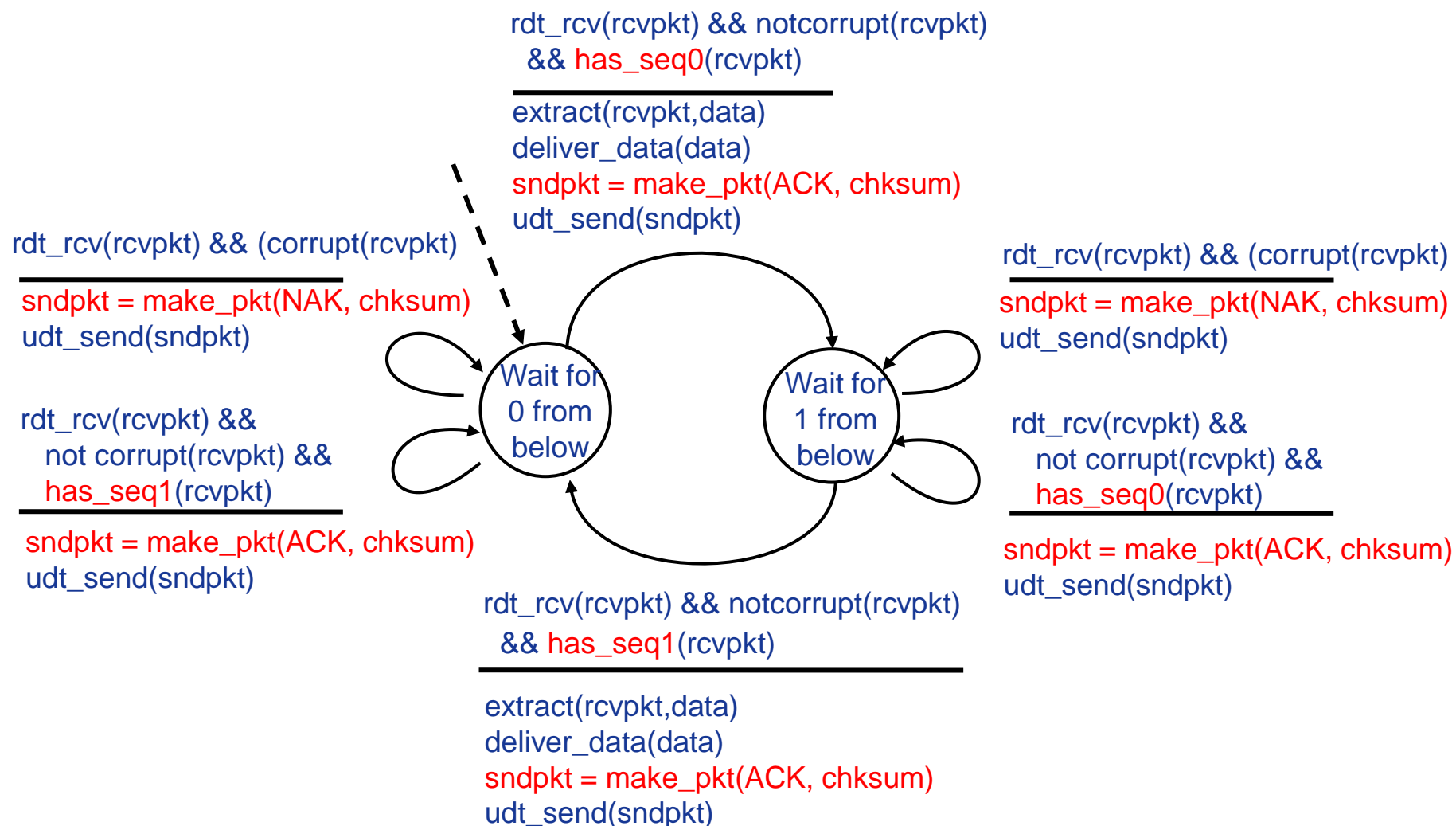
Rdt 2.1: 应对ACK/NAK破坏

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理





Rdt 2.1 vs. Rdt 2.0

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 发送方:

- ❑ 为每个分组增加了序列号
- ❑ 两个序列号(0, 1)就够用, 为什么?
- ❑ 需校验ACK/NAK消息是否发生错误
- ❑ 状态数量翻倍
 - ❑ 状态必须“记住”“当前”的分组序列号

❖ 接收方

- ❑ 需判断分组是否是重复
 - ❑ 当前所处状态提供了期望收到分组的序列号
- ❑ 注意: 接收方无法知道ACK/NAK是否被发送方正确收到





Rdt 2.2: 无NAK消息协议

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 我们真的需要两种确认消息(ACK + NAK)吗?

❖ 与rdt 2.1功能相同，但是只使用ACK

❖ 如何实现?

- 接收方通过ACK告知最后一个被正确接收的分组
- 在ACK消息中显式地加入被确认分组的序列号

❖ 发送方收到重复ACK之后，采取与收到NAK消息相同的动作

- 重传当前分组





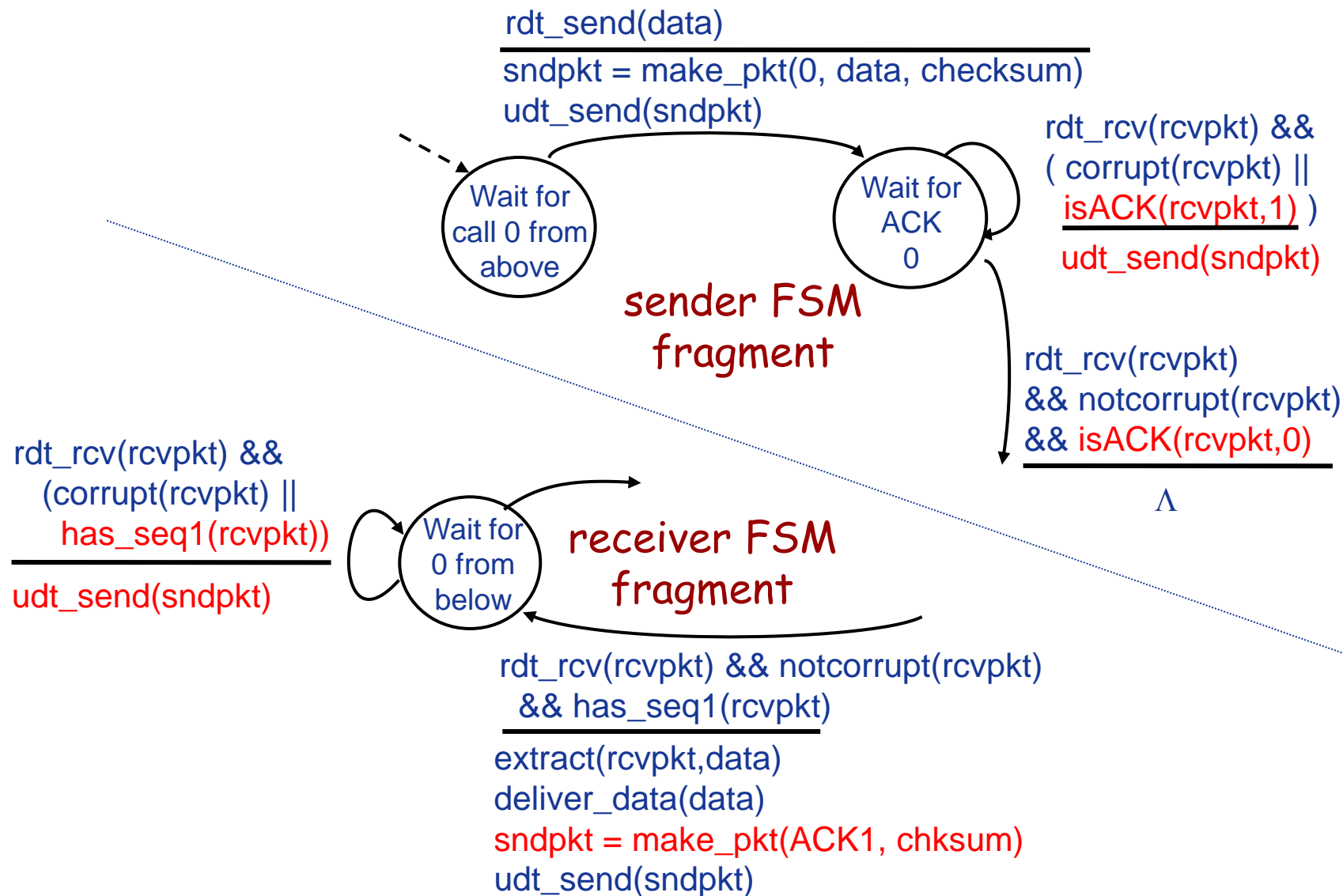
Rdt 2.2 FSM片段

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理





Rdt 3.0

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 如果信道既可能发生错误，也可能丢失分组，怎么办？

■ “校验和 + 序列号 + ACK + 重传”够用吗？

❖ 方法：发送方等待“合理”时间

- 如果没收到ACK，重传
- 如果分组或ACK只是延迟而不是丢了
 - 重传会产生重复，序列号机制能够处理
 - 接收方需在ACK中显式告知所确认的分组
- 需要定时器

❖ Rdt3.0---停-等协议





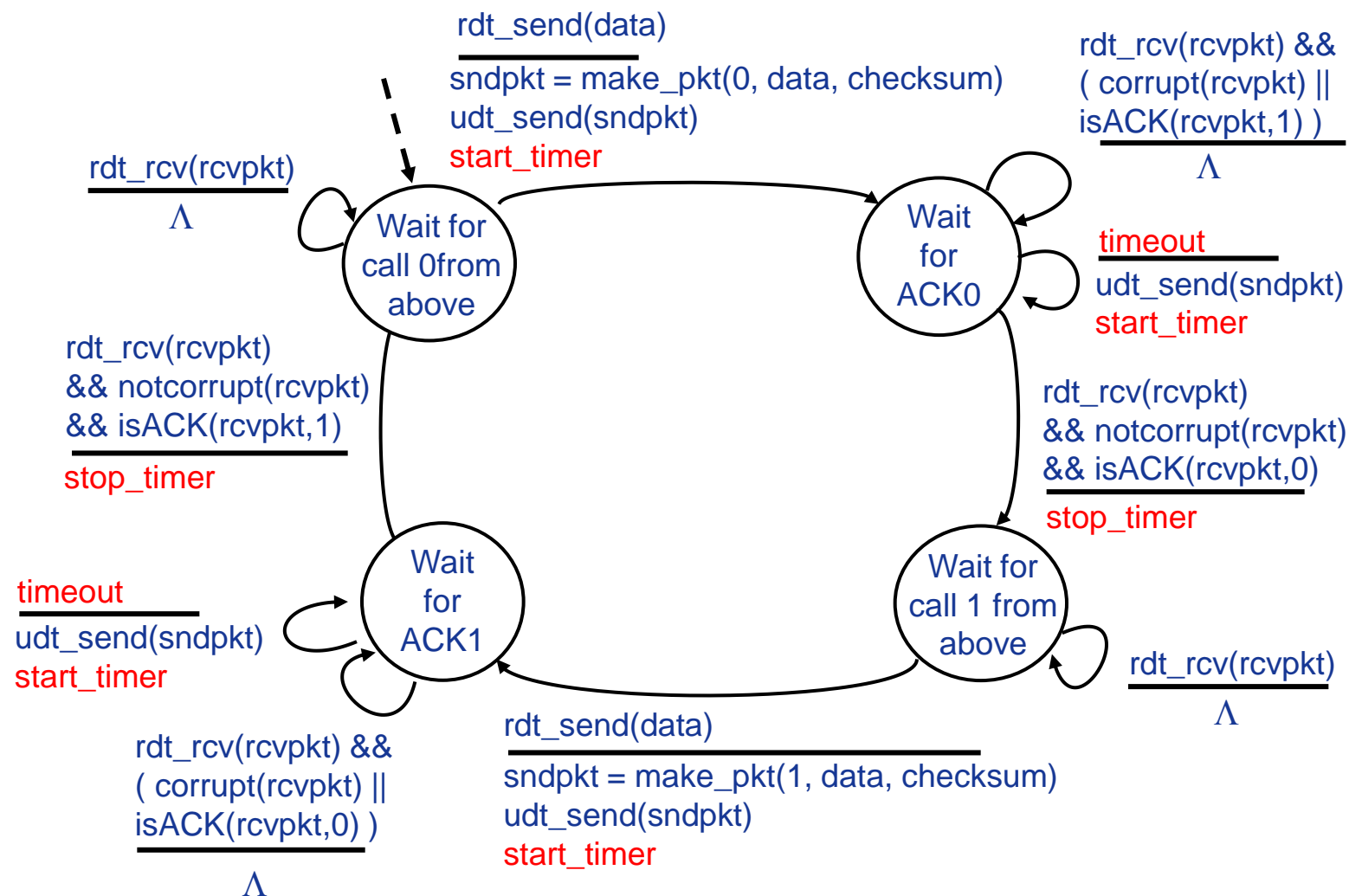
Rdt 3.0发送方FSM

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



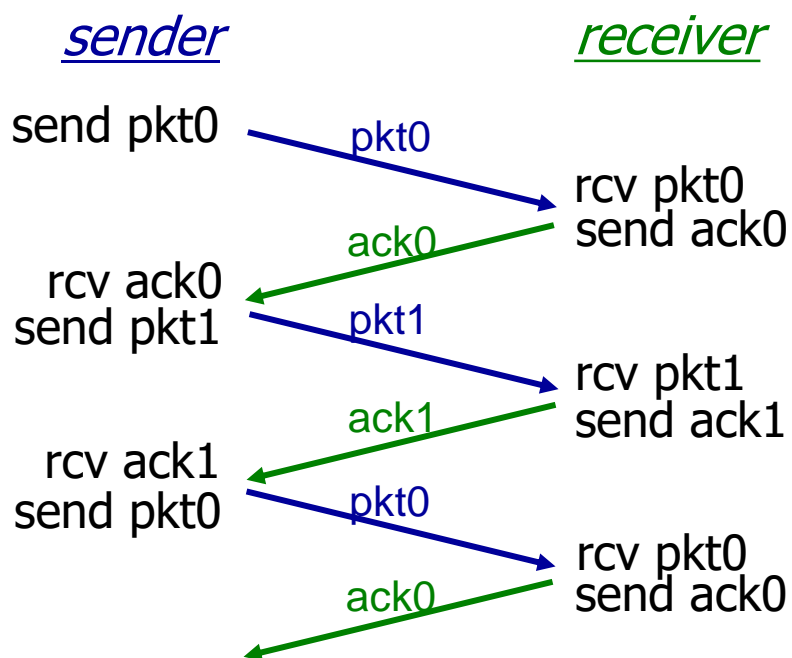
Rdt 3.0示例(1)

3.1 传输层服务

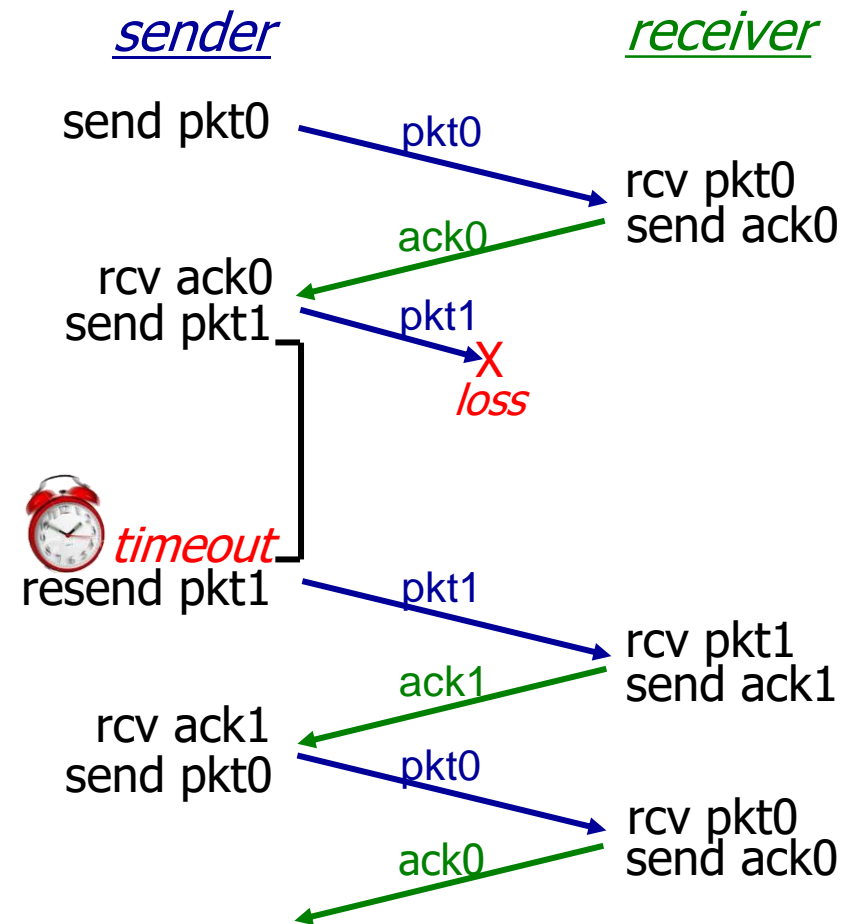
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



(a) no loss



(b) packet loss





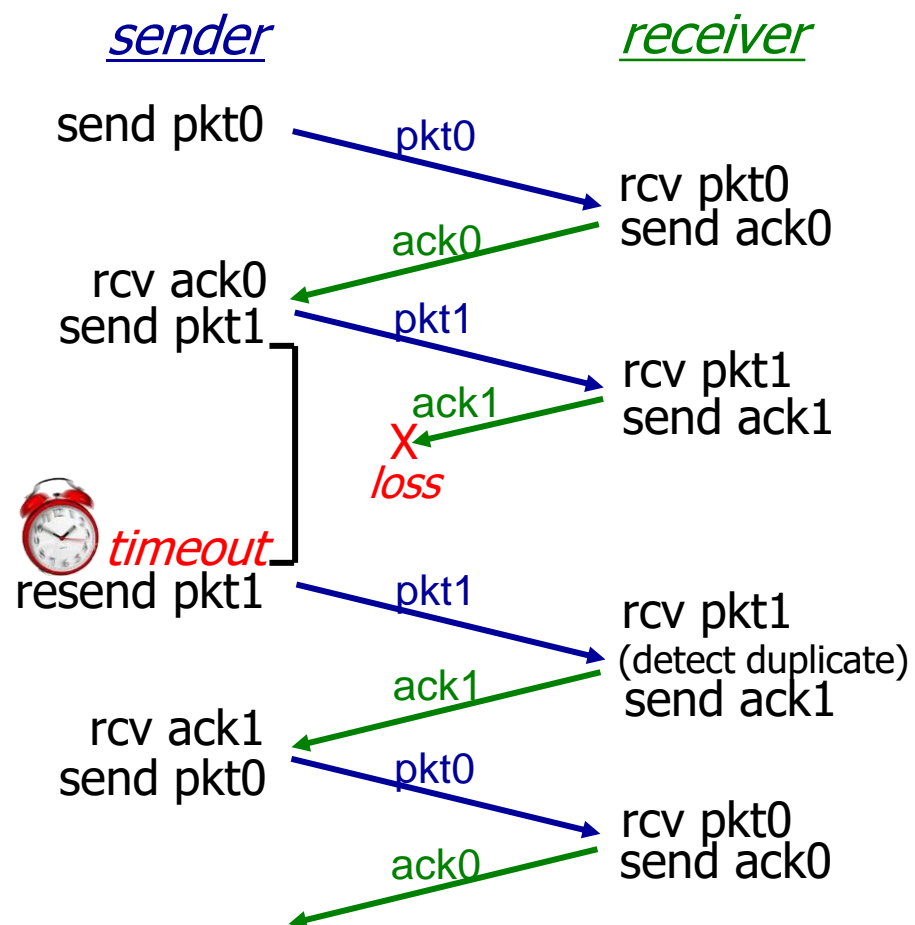
Rdt 3.0 示例(2)

3.1 传输层服务

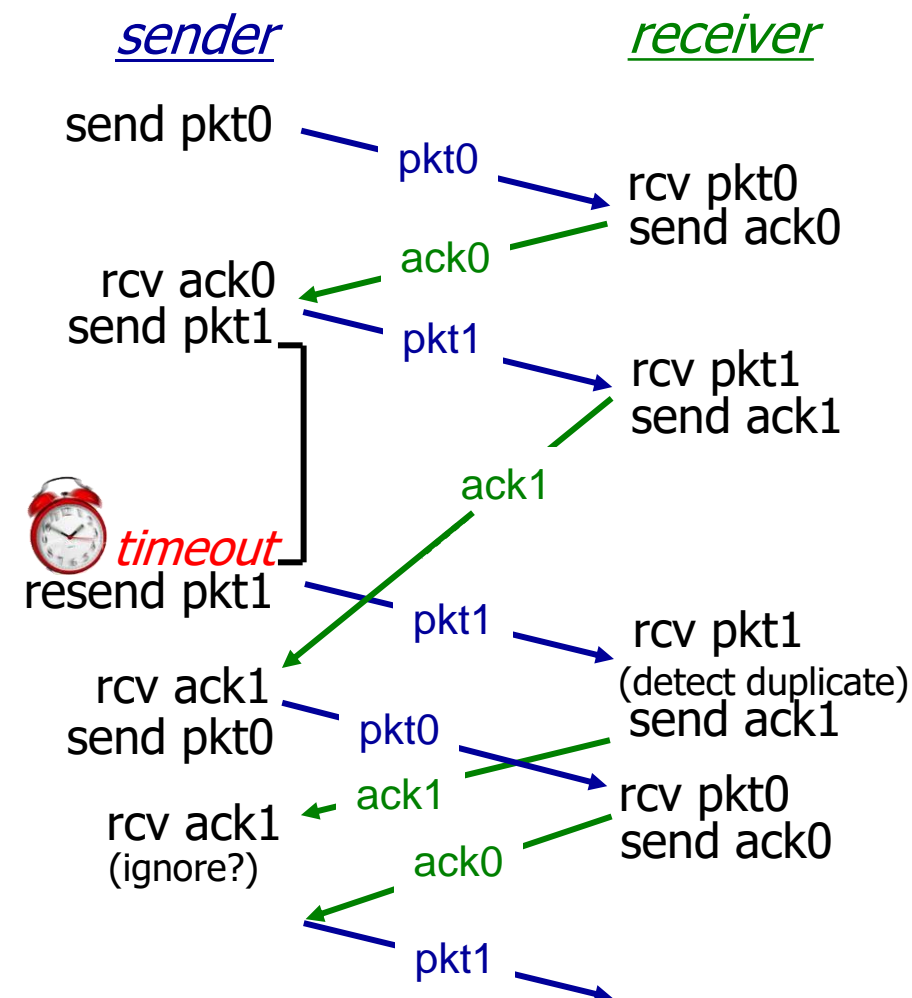
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



(c) ACK loss



(d) premature timeout/ delayed ACK





Rdt 3.0性能分析

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ Rdt 3.0能够正确工作，但性能很差

❖ 示例：1Gbps链路，15ms端到端传播延迟，1KB分组

$$T_{\text{transmit}} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8\text{kb/pkt}}{10^9 \text{ b/sec}} = 8 \text{ microsec}$$

■ 发送方利用率：发送方发送时间百分比

$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

■ 在1Gbps链路上每30毫秒才发送一个分组 → 33KB/sec

■ 网络协议限制了物理资源的利用





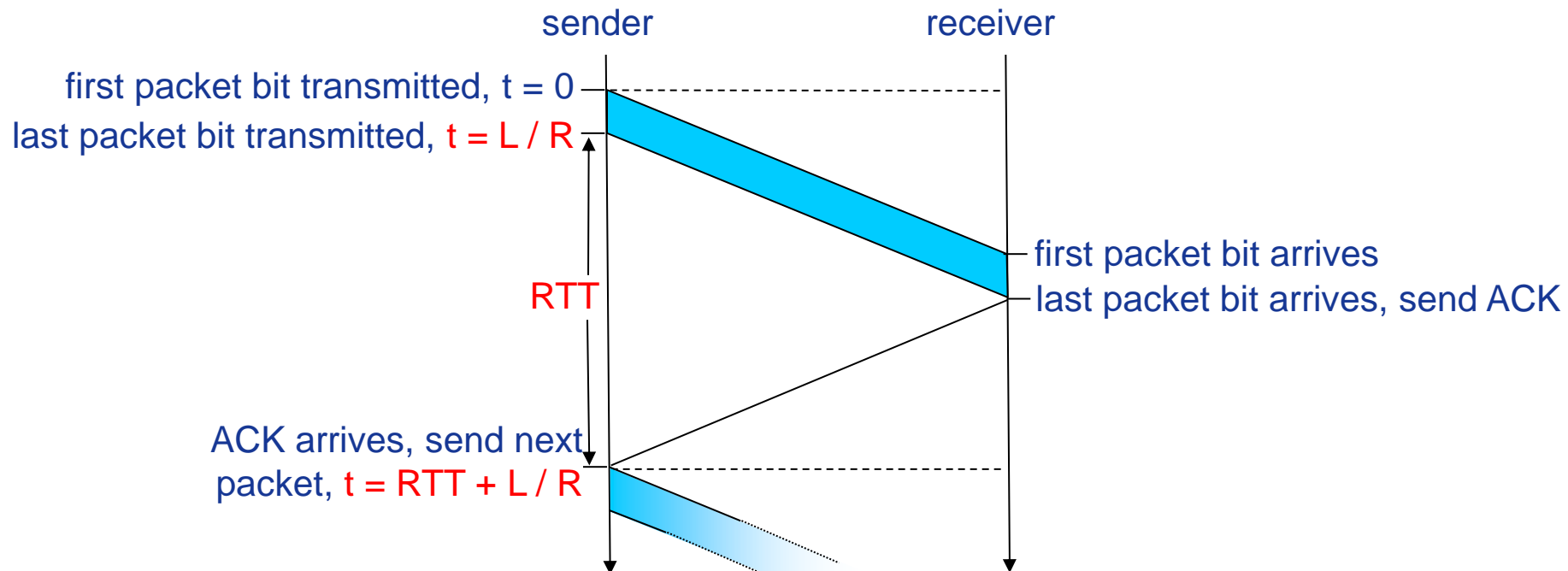
Rdt 3.0: 停等操作

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



$$U_{\text{sender}} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$



停等协议的信道利用率

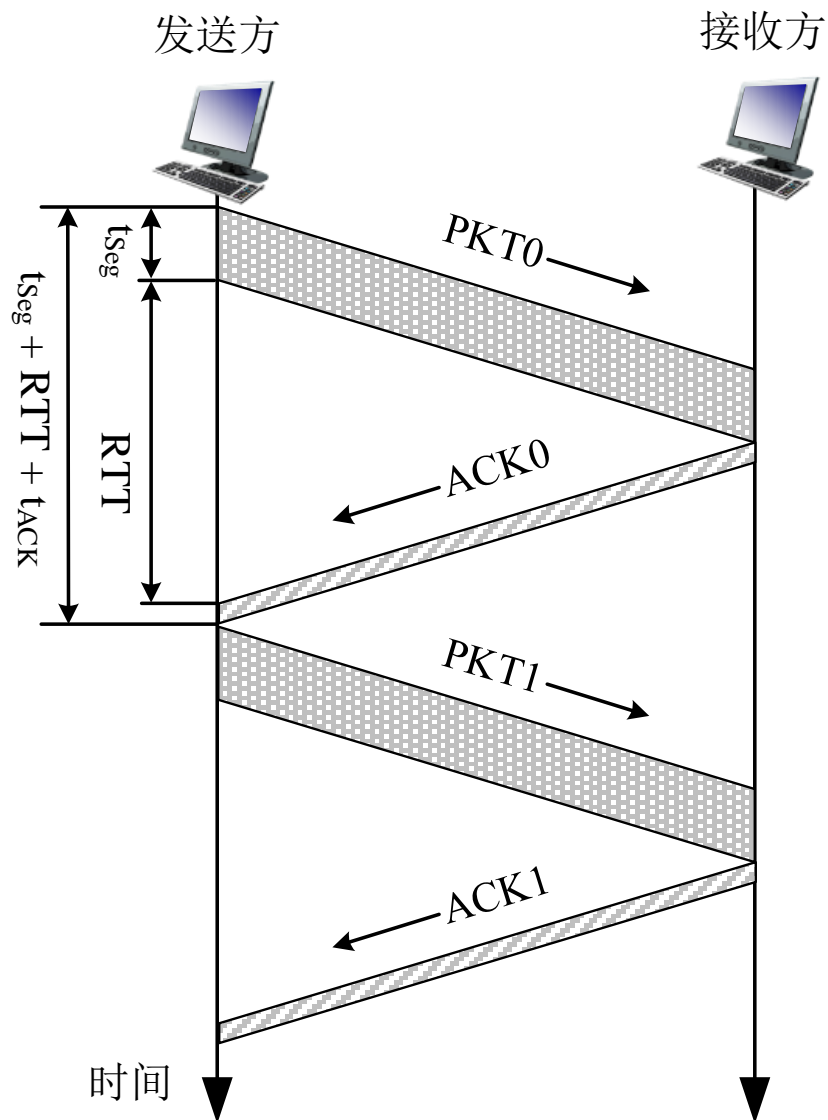
3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

$$U_{\text{Stop-Wait}} = \frac{t_{\text{seg}}}{t_{\text{seg}} + RTT + t_{\text{ACK}}}$$





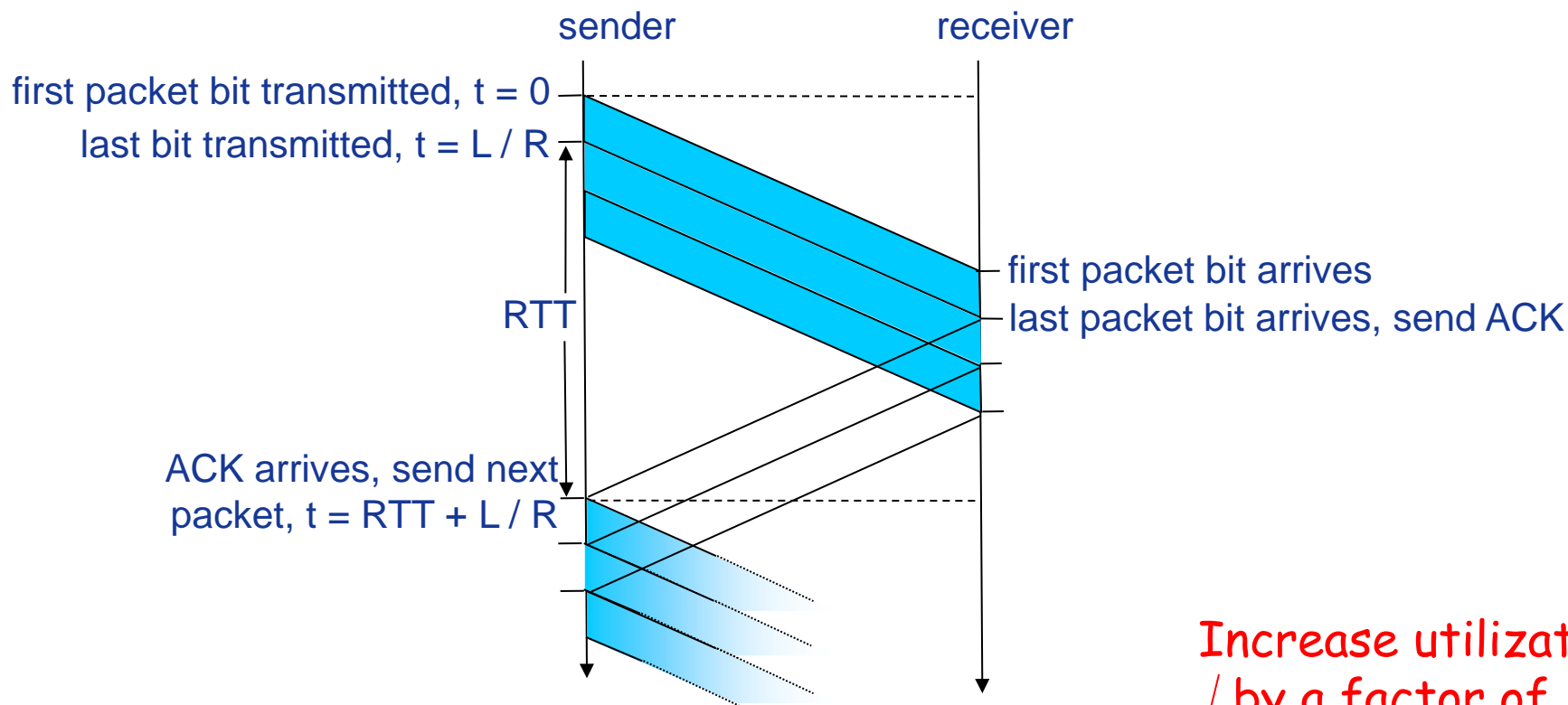
流水线机制：提高资源利用率

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



Increase utilization
by a factor of 3!

$$U_{\text{sender}} = \frac{3 * L / R}{RTT + L / R} = \frac{.024}{30.008} = 0.0008$$





流水线协议

3.1 传输层服务

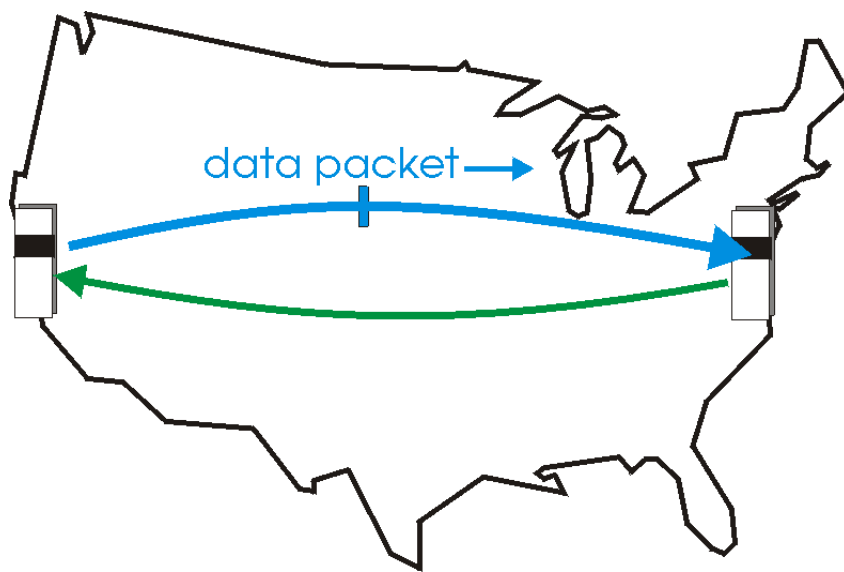
3.2 传输层多路复用/分解

3.3 UDP协议

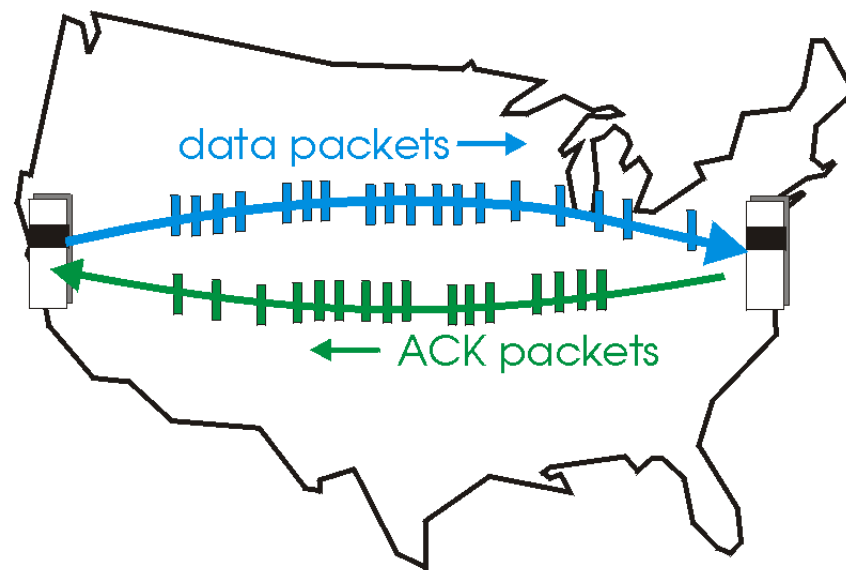
3.4 可靠数据传输原理

❖ 允许发送方在收到ACK之前连续发送多个分组

- 更大的**序列号范围**
- 发送方和/或接收方需要更大的存储空间以**缓存分组**



(a) a stop-and-wait protocol in operation



(b) a pipelined protocol in operation





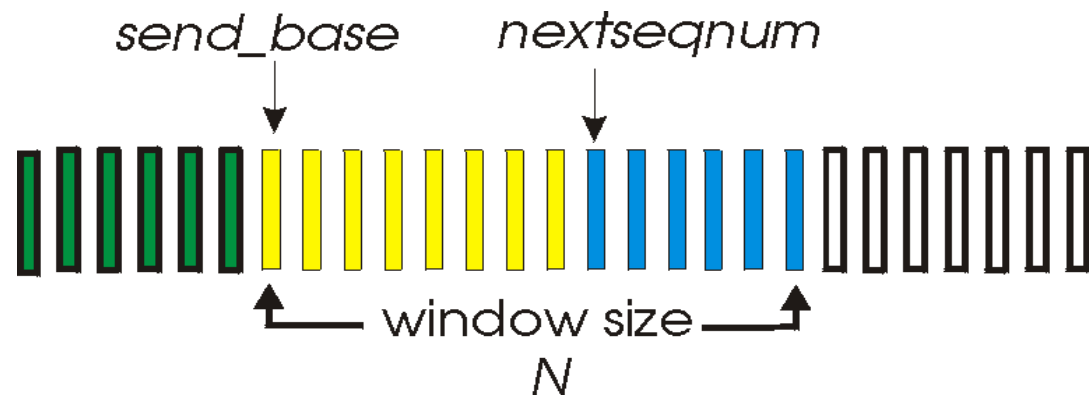
滑动窗口协议

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



❖ 滑动窗口协议: Sliding-window protocol

❖ 窗口

- 允许使用的序列号范围
- 窗口尺寸为 N : 最多有 N 个等待确认的消息

❖ 滑动窗口

- 随着协议的运行, 窗口在序列号空间内向前滑动

❖ 滑动窗口协议: GBN, SR





滑动窗口协议

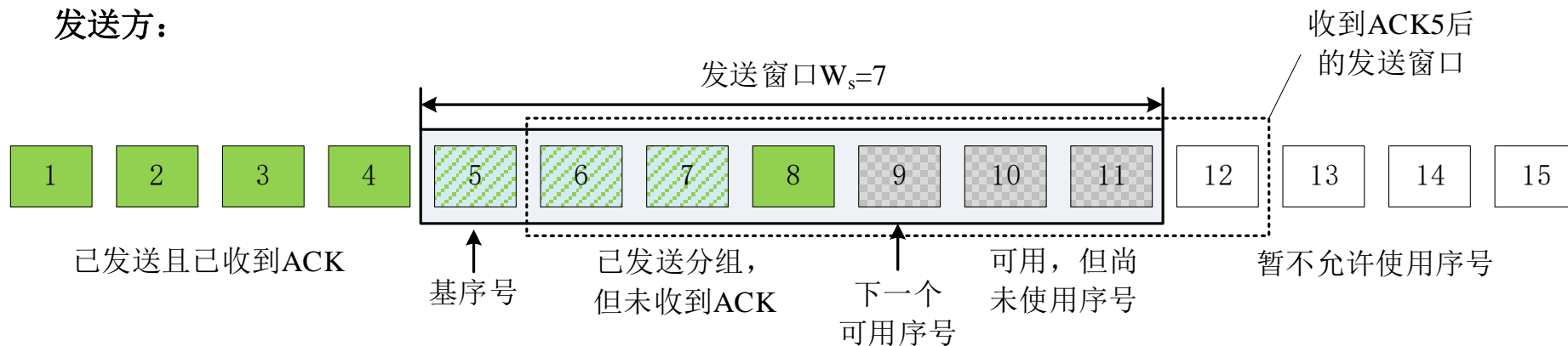
3.1 传输层服务

3.2 传输层多路复用/分解

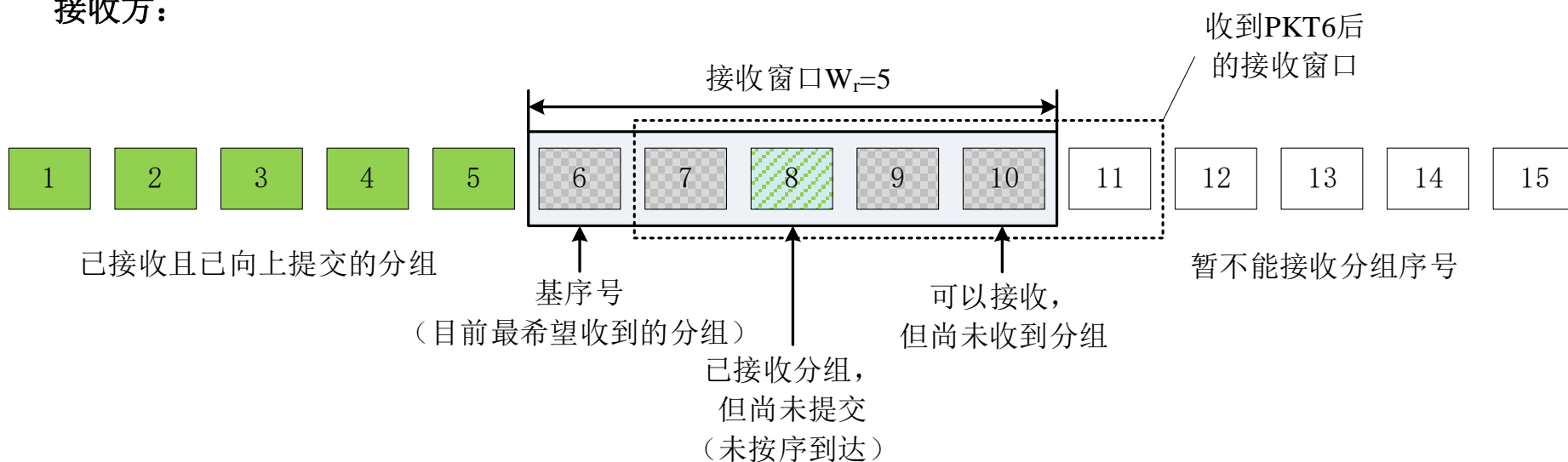
3.3 UDP协议

3.4 可靠数据传输原理

发送方:



接收方:





滑动窗口协议的信道利用率

3.1 传输层服务

3.2 传输层多路复用/分解

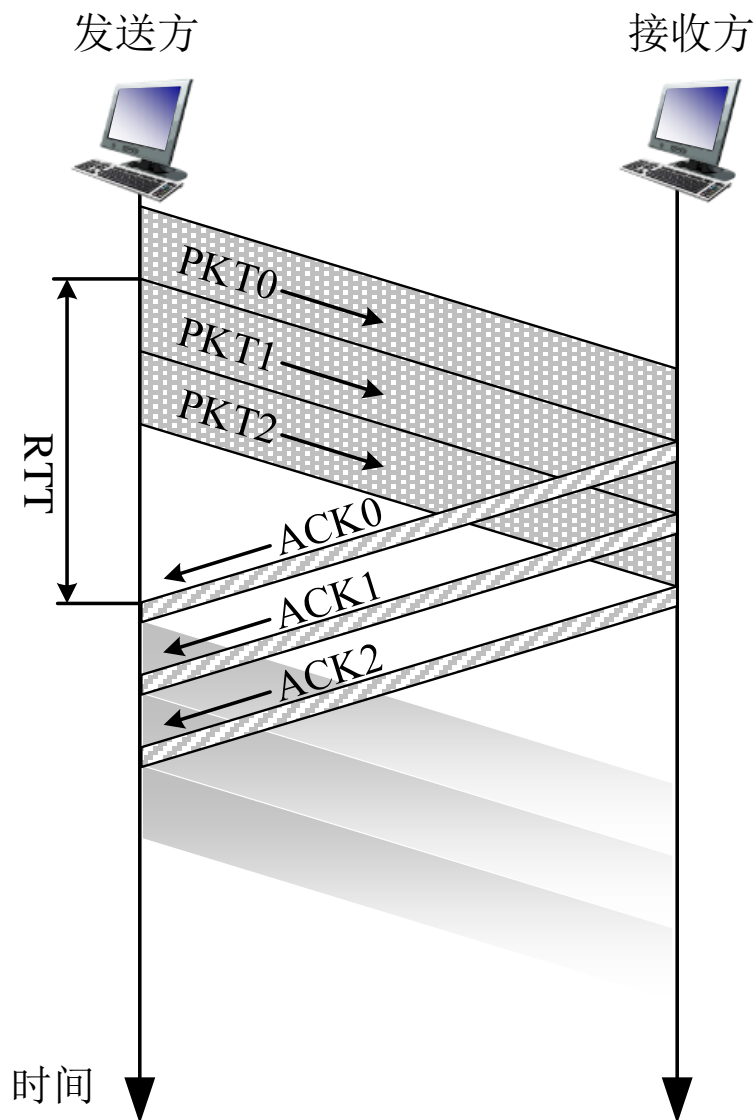
3.3 UDP协议

3.4 可靠数据传输原理

$$U = \frac{W_s \times t_{seg}}{t_{seg} + RTT + t_{ACK}}$$

$$U = \frac{W_s \times L/R}{L/R + 2dp + L'/R}$$

$$U = \frac{W_s \times L}{L + 2dpR + L'}$$





Go-Back-N(GBN)协议：发送方

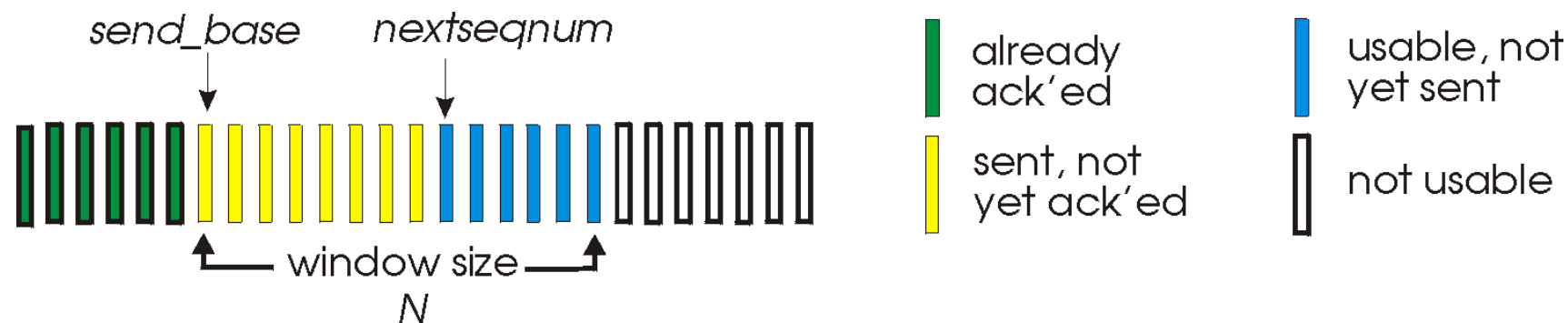
3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

- ❖ 分组头部包含k-bit序列号
- ❖ 窗口尺寸为N，最多允许N个分组未确认



- ❖ ACK(n): 确认到序列号n(包含n)的分组均已被正确接收
 - 累积确认
 - 可能收到重复ACK
- ❖ 为“空中”的分组设置计时器(timer)
- ❖ 超时Timeout(n)事件: 重传序列号大于等于n, 还未收到ACK的所有分组





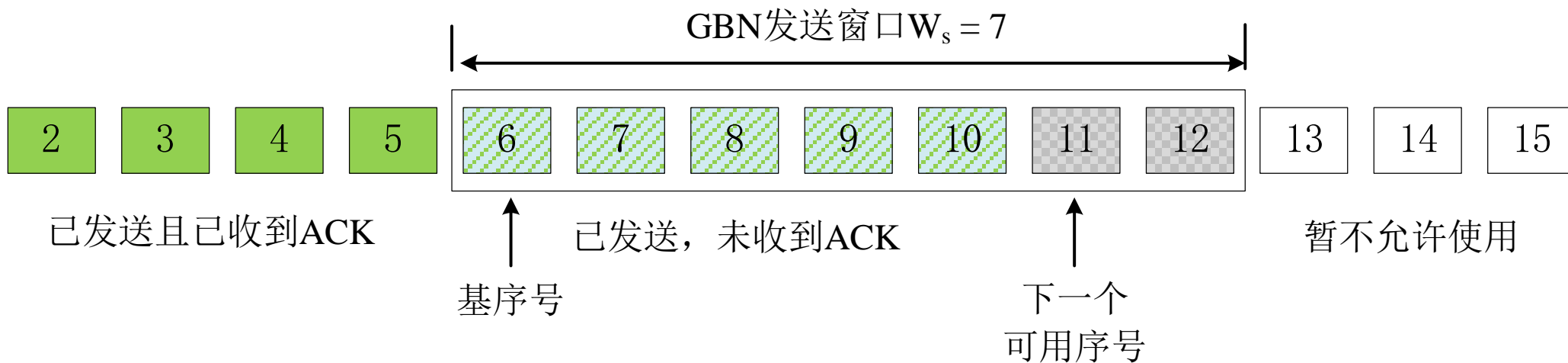
Go-Back-N(GBN)协议：发送方

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理





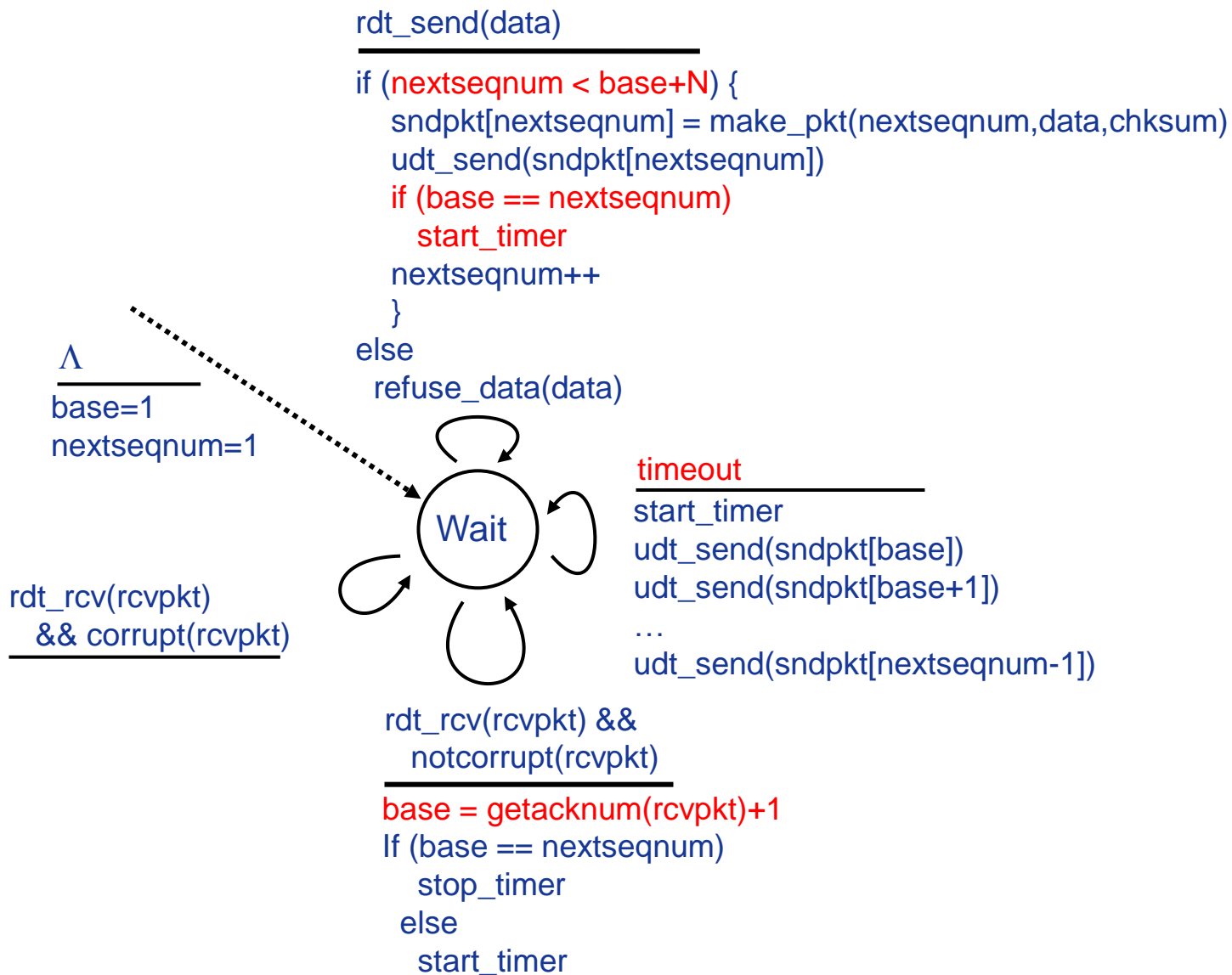
GBN: 发送方扩展FSM

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理





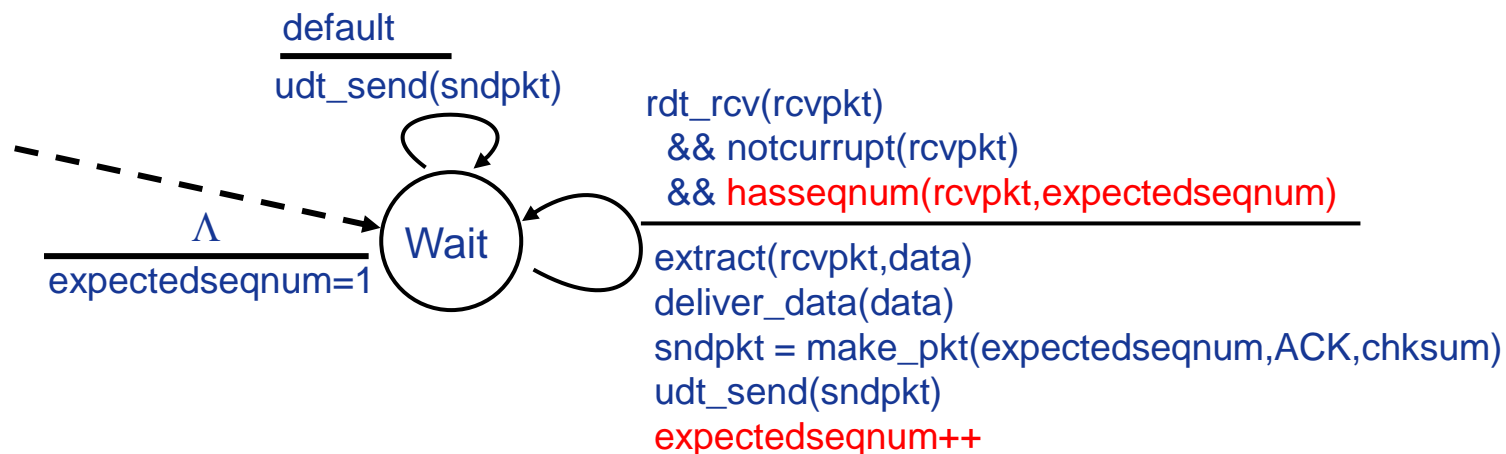
GBN: 接收方扩展FSM

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



❖ ACK机制: 发送拥有最高序列号的、已被正确接收的分组的ACK

- 可能产生重复ACK
- 只需要记住唯一的**expectedseqnum**, 即 **$W_R=1$**

❖ 乱序到达的分组:

- 直接丢弃→接收方没有缓存
- 重新确认序列号最大的、按序到达的分组





GBN示例

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

ignore duplicate ACK

*pkt 2 timeout*

send pkt2

send pkt3

send pkt4

send pkt5

receiver

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, discard,
(re)send ack1receive pkt4, discard,
(re)send ack1receive pkt5, discard,
(re)send ack1

rcv pkt2, deliver, send ack2

rcv pkt3, deliver, send ack3

rcv pkt4, deliver, send ack4

rcv pkt5, deliver, send ack5





例题

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

【例1】数据链路层采用后退N帧（GBN）协议，发送方已经发送了编号为0~7的帧。当计时器超时时，若发送方只收到0、2、3号帧的确认，则发送方需要重发的帧数是多少？分别是那几个帧？

【解】根据GBN协议工作原理，GBN协议的确认是累积确认，所以此时发送端需要重发的帧数是4个，依次分别是4、5、6、7号帧。





Selective Repeat协议

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ **GBN有什么缺陷？**

❖ **SR协议：**

❖ 接收方对每个分组单独进行确认

- 设置缓存机制，缓存乱序到达的分组

❖ 发送方只重传那些没收到ACK的分组

- 为每个分组设置定时器

❖ 发送方窗口

- N个连续的序列号
- 限制已发送且未确认的分组数

❖ 接收窗口

- 可以接收的无差错到达的分组序号





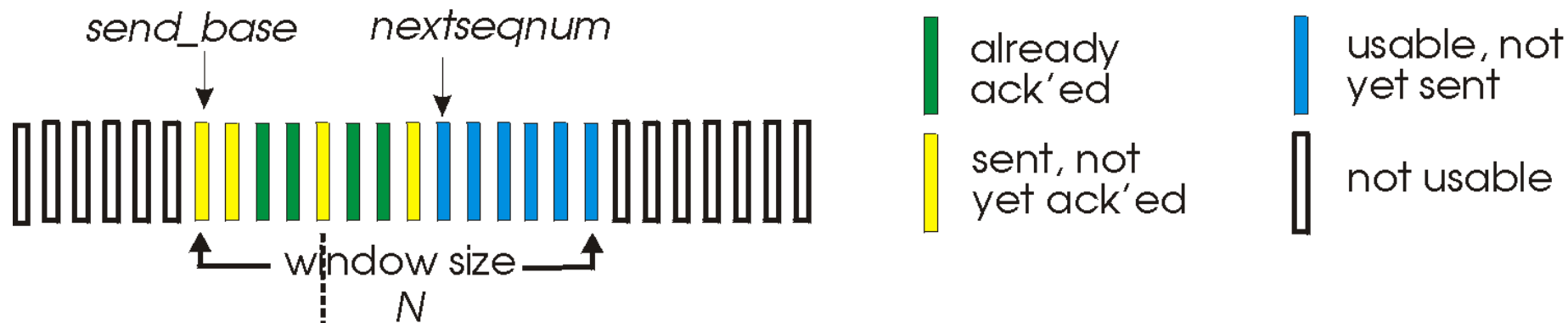
SR协议：发送方/接收方窗口

3.1 传输层服务

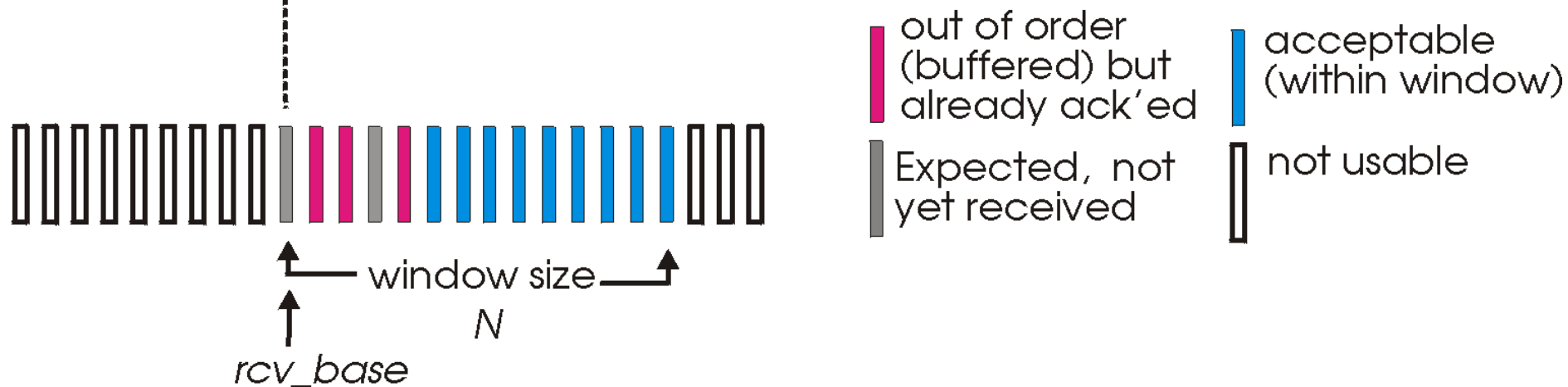
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



(a) sender view of sequence numbers



(b) receiver view of sequence numbers






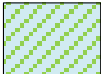
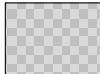

SR协议：发送方/接收方窗口

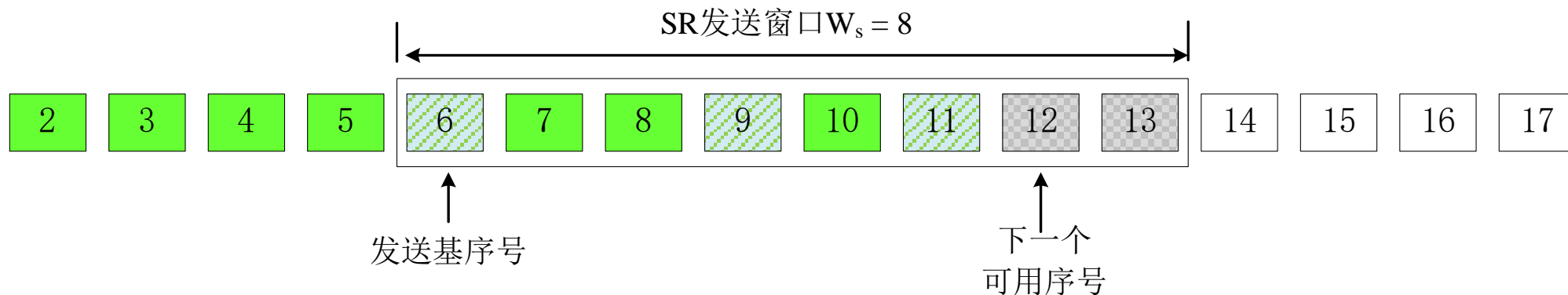
3.1 传输层服务





3.2 传输层多路复用/分解

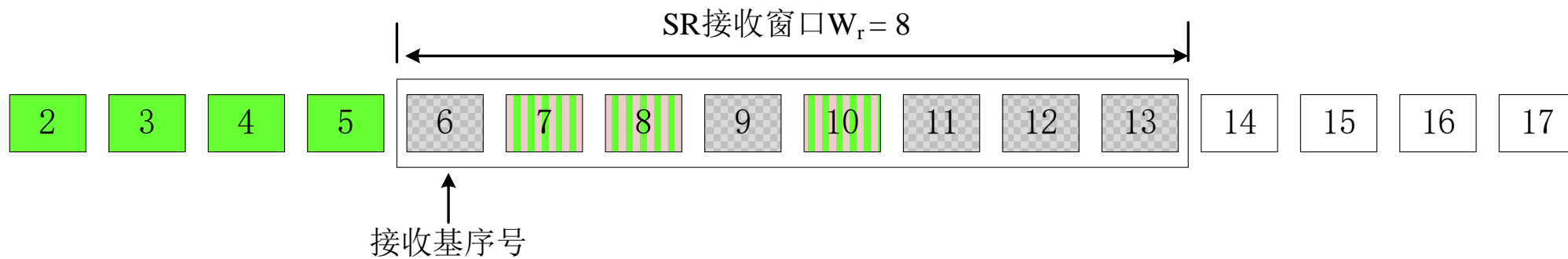
3.3 UDP协议

3.4 可靠数据传输原理

发送方：  已发送且已收到ACK  已发送，未收到ACK  可用序号  不可用序号



接收方：  已经接收，并已提交  已接收，并已确认，但暂不能提交  可以接收  暂不能接收





SR协议

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

—sender—

data from above :

- ❖ if next available seq # in window, send pkt

timeout(n):

- ❖ resend pkt n, restart timer

ACK(n) in [sendbase,sendbase+N]:

- ❖ mark pkt n as received
- ❖ if n is smallest unACKed pkt, advance window base to next unACKed seq #

—receiver—

pkt n in [rcvbase, rcvbase+N-1]

- ❑ send ACK(n)
- ❑ out-of-order: buffer
- ❑ in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

pkt n in [rcvbase-N,rcvbase-1]

- ❑ ACK(n)

otherwise:

- ❑ ignore

思考：SR协议还可以有其他设计吗？





SR协议示例

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

record ack3 arrived



pkt 2 timeout

send pkt2

(but not 3,4,5)

receiver

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, buffer,
send ack3

receive pkt4, buffer,
send ack4

receive pkt5, buffer,
send ack5

rcv pkt2; deliver pkt2,
pkt3, pkt4, pkt5; send ack2

Q: 当ack2到达发送方会怎么样?





SR协议：困境

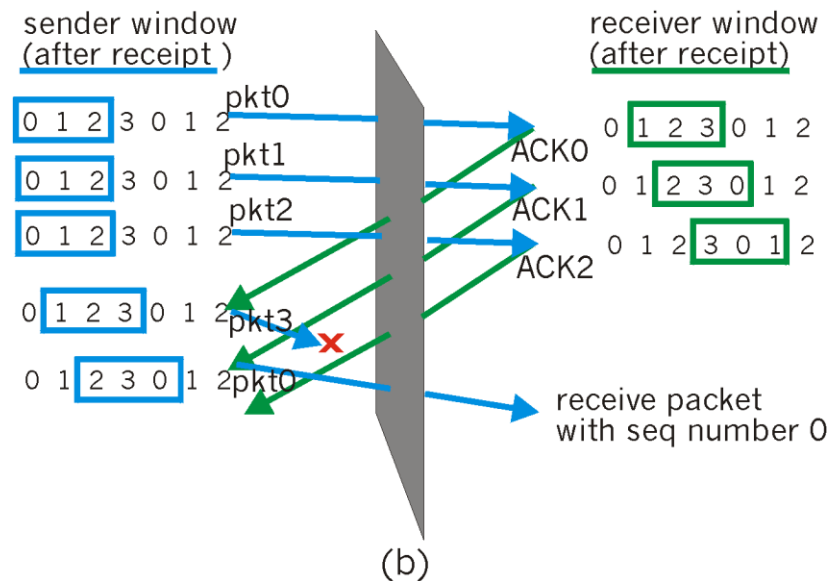
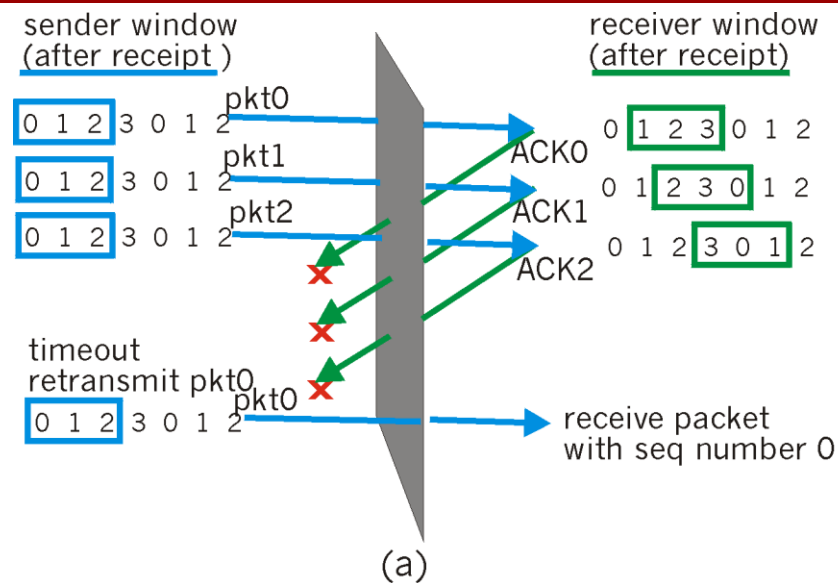
3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

- ❖ 序列号: 0, 1, 2, 3
- ❖ 窗口尺寸: 3
- ❖ 接收方能区分开右侧两种不同的场景吗?
- ❖ (a)中, 发送方重发0号分组, 接收方收到后会如何处理?



窗口大小与序号空间的约束条件？

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

❖ 问题：序列号空间大小与窗口尺寸需满足什么关系？

$$W_s + W_r \leq 2^k$$

- W_s 发送窗口, W_r 接收窗口, k 序号位数

❖ 对于GBN协议: $W_r=1$

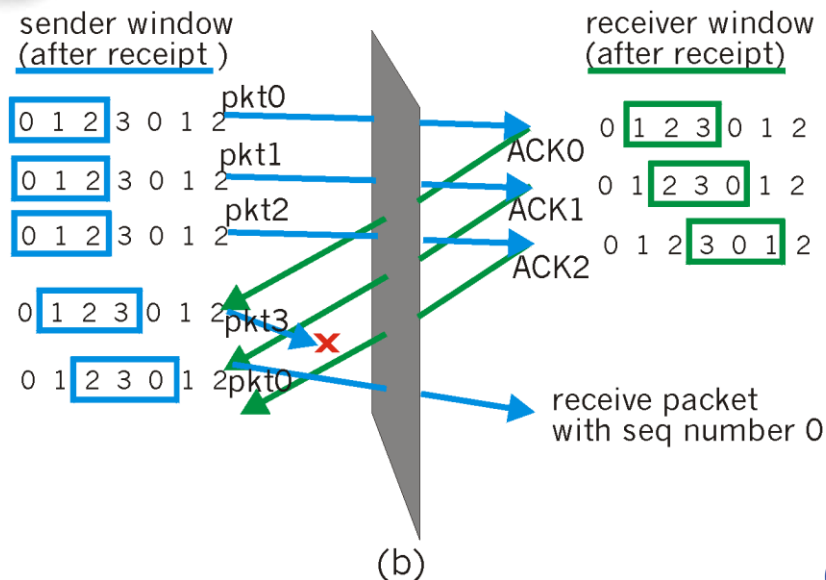
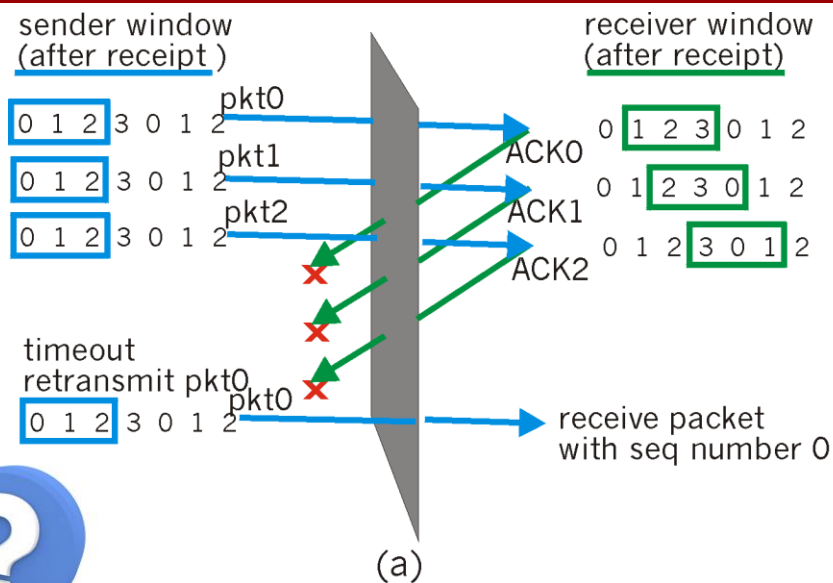
$$W_S \leq 2^k - 1$$

❖ 对于典型的 $W_s=W_r=W$ 的SR协议

$$W_S \leq 2^{(k-1)}$$

❖ 对于停-等协议, 即 $W_s=W_r=1$

$$k \geq 1$$





滑动窗口协议的窗口大小

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理



❖ 讨论:

- 1. 滑动窗口协议的窗口大小影响协议哪些性能?
- 2. 哪些因素会影响滑动窗口大小的确定?

❖ 性能:

- 信道利用率
- 吞吐率
-

❖ 因素:

- 序号空间
- 缓存大小
- 流量控制
- 拥塞控制
-



单选题 1分



3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

主机甲通过128 kbps卫星链路，采用滑动窗口协议向主机乙发送数据，链路单向传播延迟为250 ms，帧长为1000字节。不考虑确认帧的开销，为使链路利用率不小于80%，帧序号的比特数至少是_____。

A

2

B

3

C

4

D

5



提交



哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY

立足航天，服务国防，面向国民经济主战场



3.5 TCP协议

聂兰顺



TCP概述: RFCs-793, 1122, 1323, 2018, 2581

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

❖ 点对点

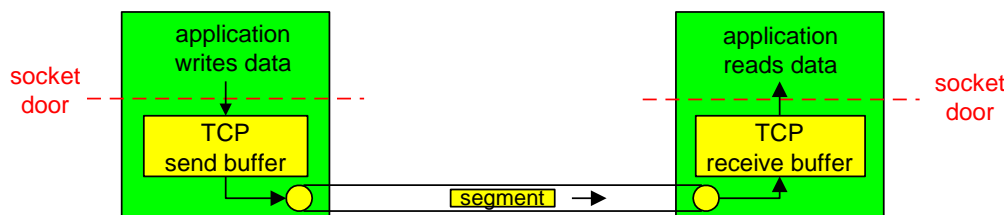
- 一个发送方，一个接收方

❖ 可靠的、按序字节流

❖ 流水线机制

- TCP拥塞控制和流量控制机制设置窗口尺寸

❖ 发送方/接收方缓存



❖ 全双工(full-duplex)

- 同一连接中能够传输双向数据流

❖ 面向连接

- 通信双方在发送数据之前必须建立连接。
- 连接状态只在连接的两端中维护，在沿途节点中并不维护状态。
- TCP连接包括：两台主机上的缓存、连接状态变量、socket等

❖ 流量控制机制

❖ 拥塞控制



TCP段结构

3.1 传输层服务

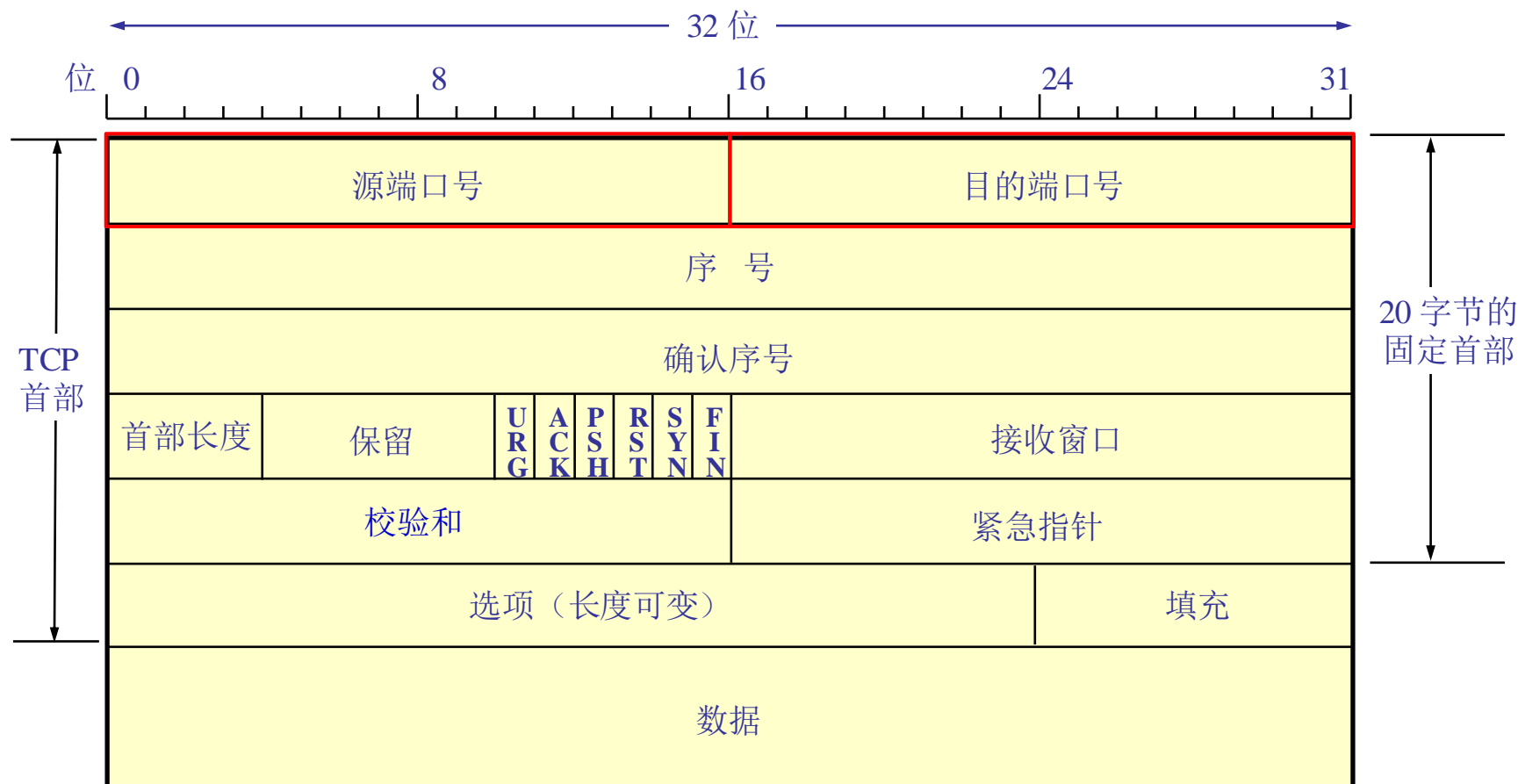
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 源端口号与目的端口号字段分别占16位

■ 多路复用/分解



TCP段结构

3.1 传输层服务

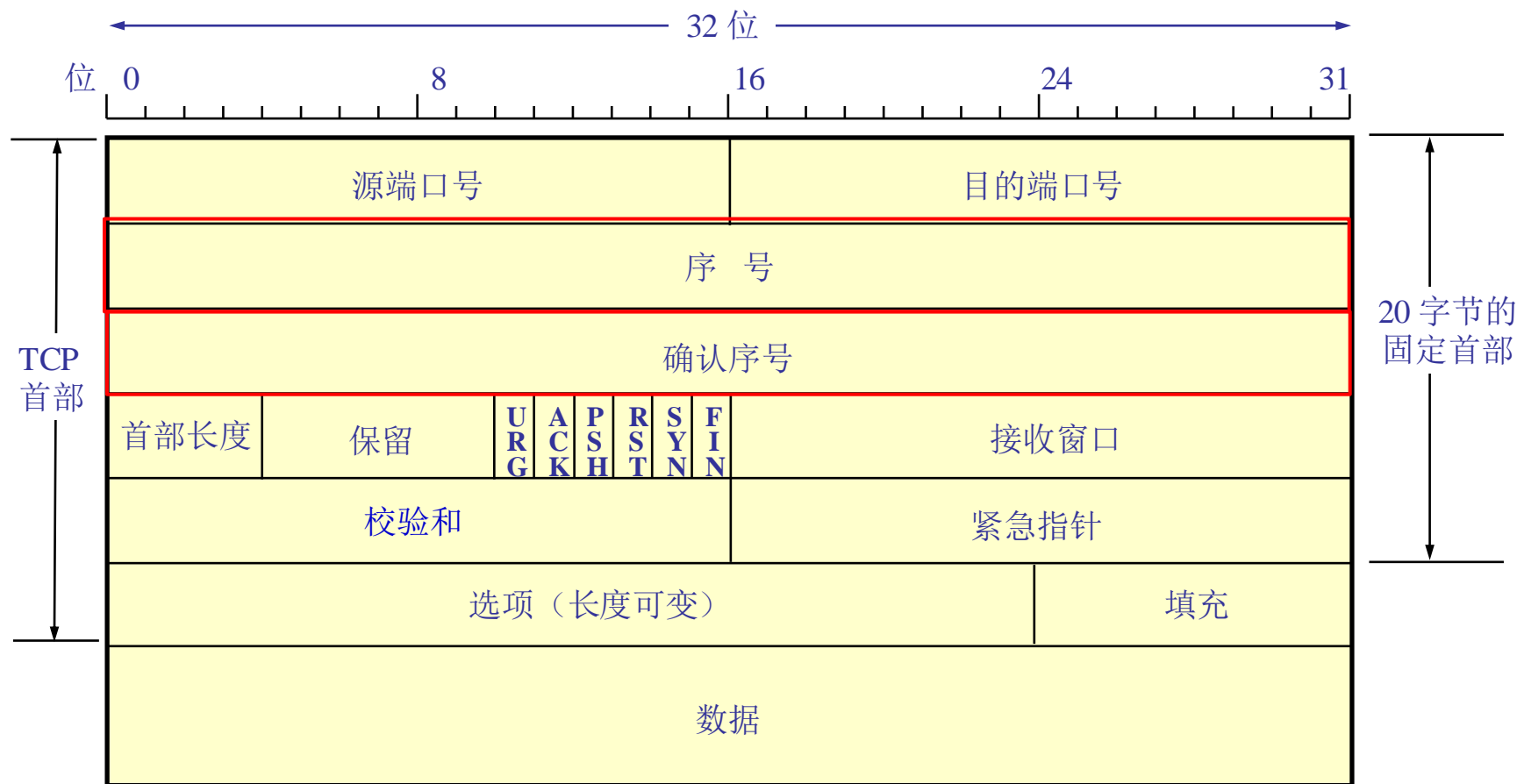
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 序号字段与确认序号字段分别占32位

- 对每个应用层数据的每个字节进行编号
- 确认序号是期望从对方接收数据的字节序号，累计确认



TCP段结构

3.1 传输层服务

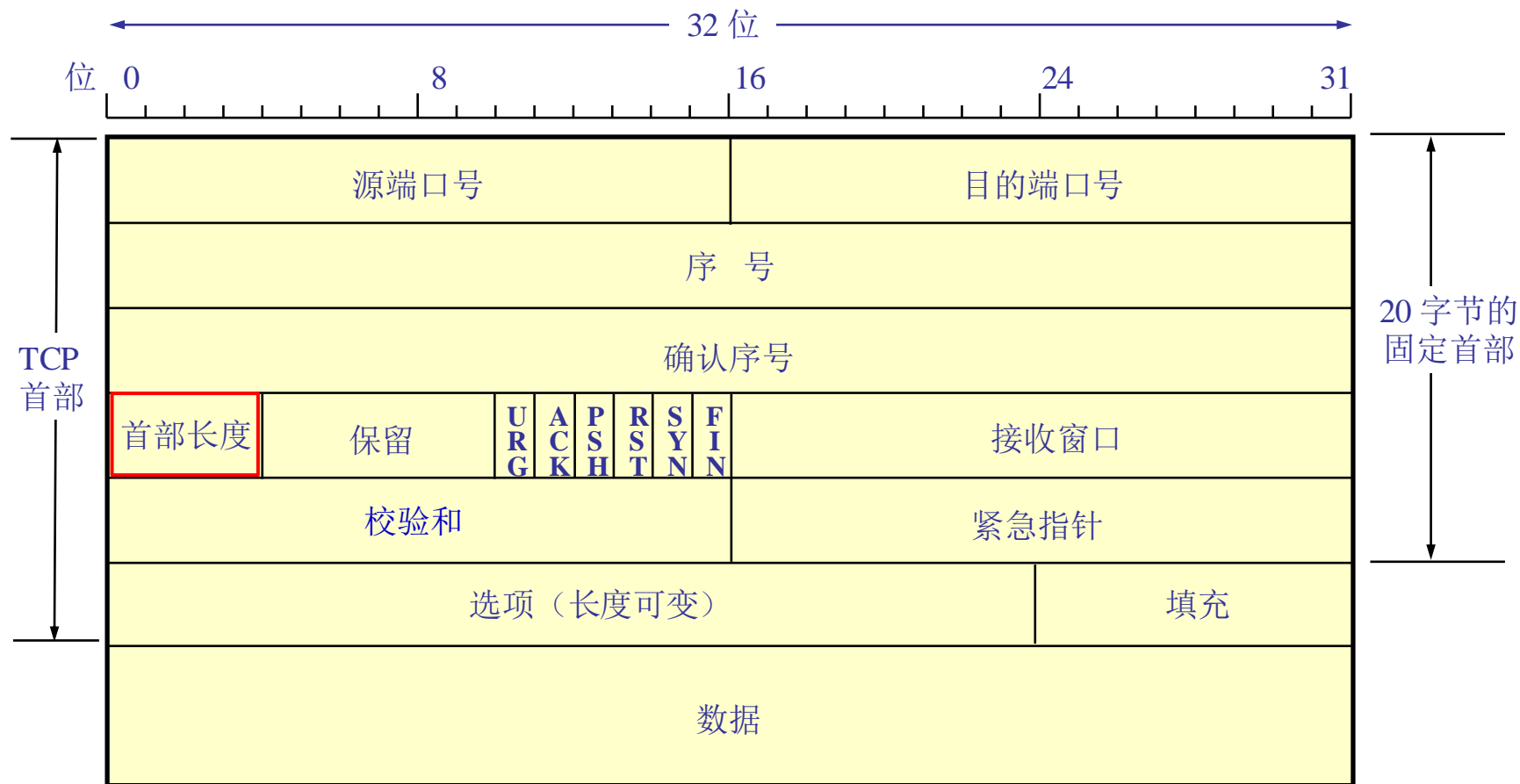
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 首部长度的字段占4位

■ 4字节为计算单位





TCP段结构

3.1 传输层服务

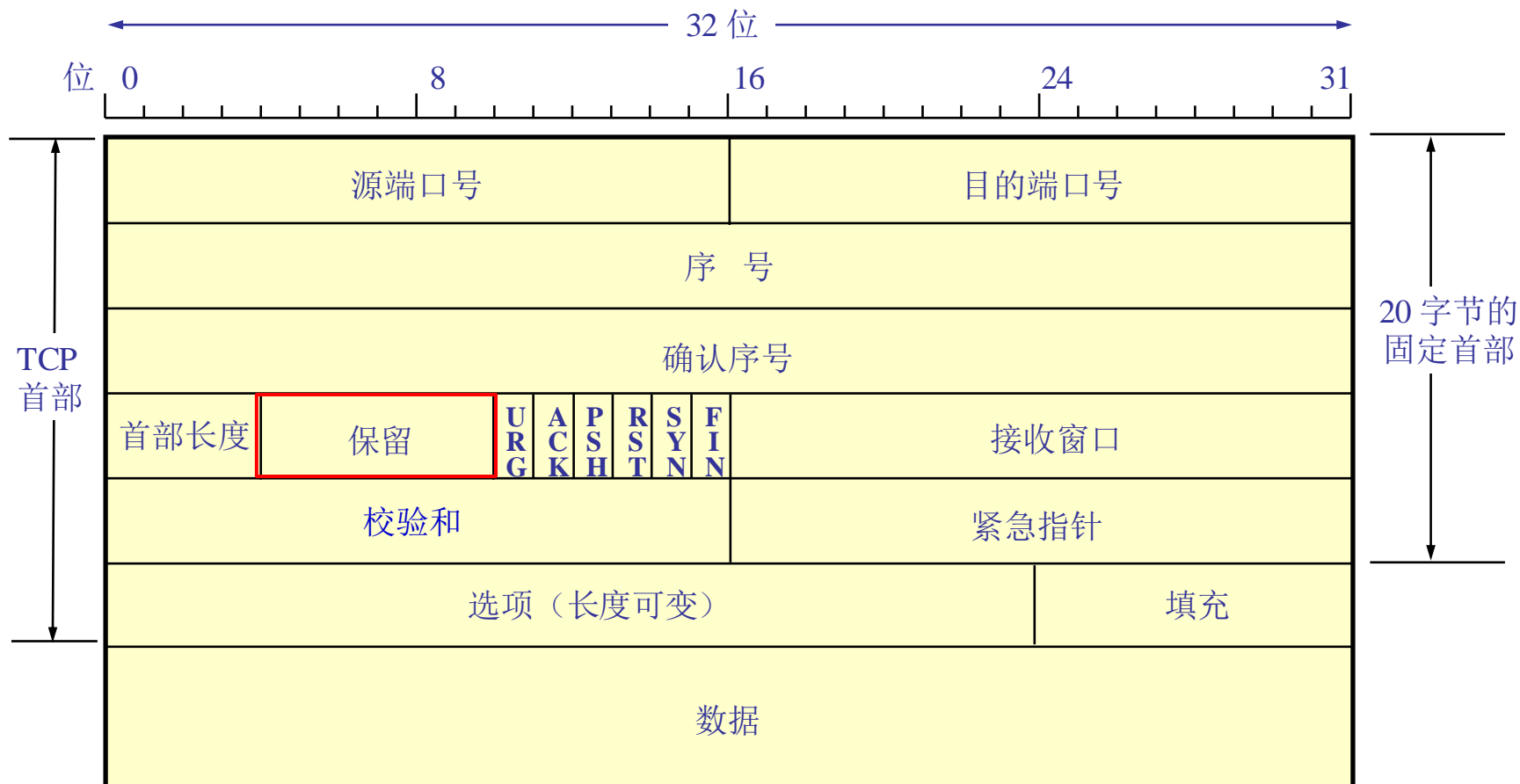
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 保留字段占6位

■ 目前值为0





TCP段结构

3.1 传输层服务

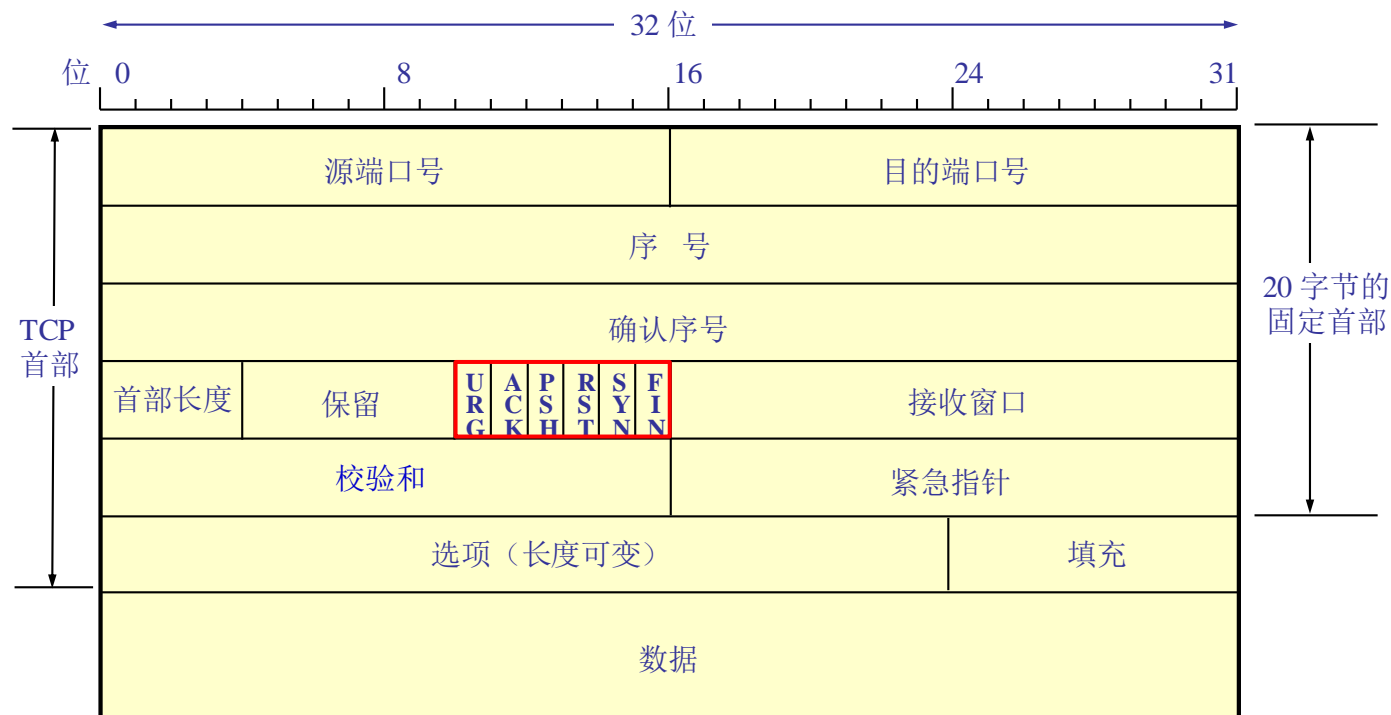
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 6位标志位（字段）

- URG=1时，表明紧急指针字段有效
- ACK=1时，标识确认序号字段有效
- PSH=1时，尽快将段中数据交付接收应用进程
- RST=1时，重新建立TCP连接
- SYN=1时，表示该TCP段是一个建立新连接请求控制段
- FIN=1时，表明请求释放TCP连接



TCP段结构

3.1 传输层服务

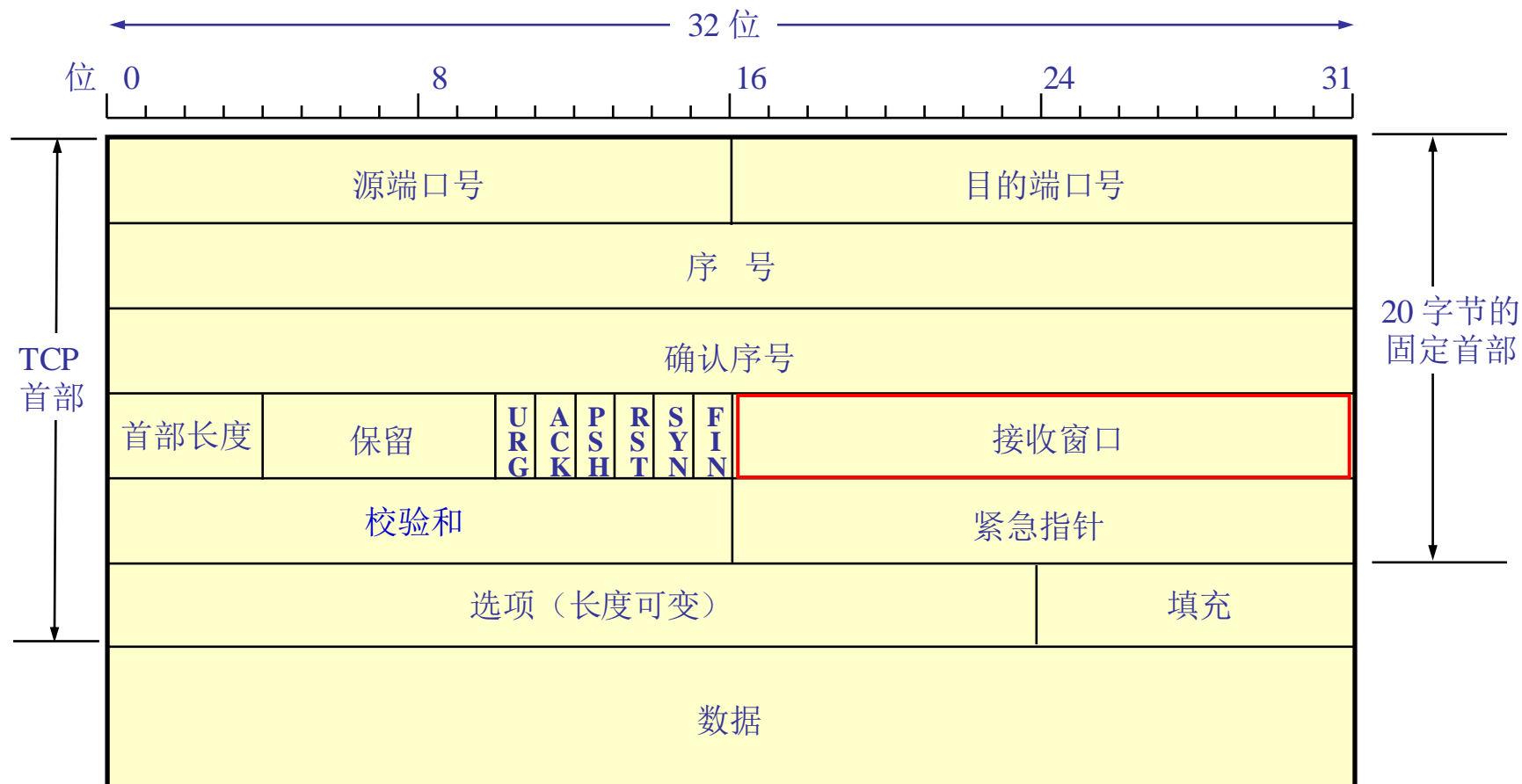
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 接收窗口字段占16位

■ 流量控制



TCP段结构

3.1 传输层服务

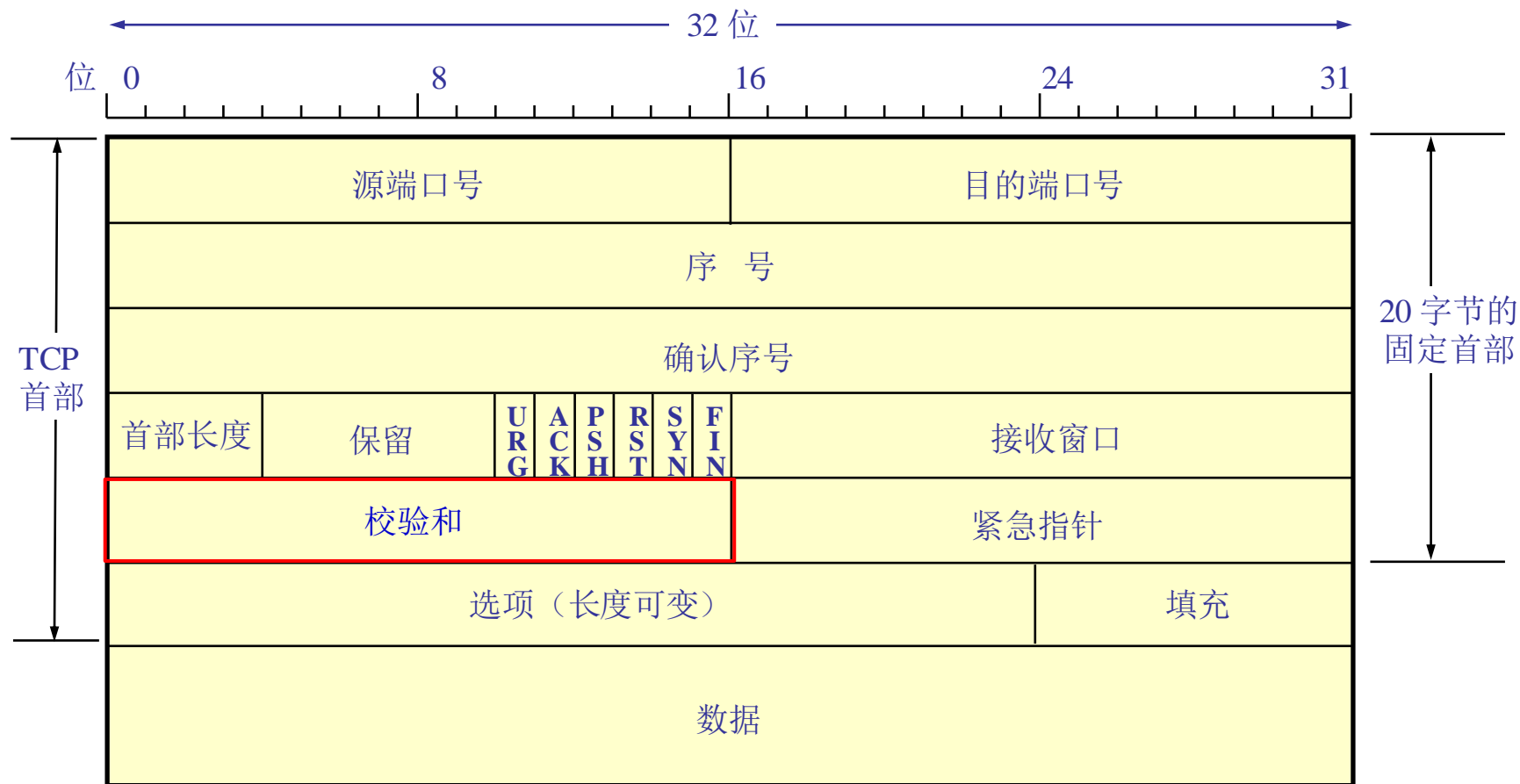
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 校验和字段占16位

- 包括TCP伪首部、TCP首部和应用层数据三部分



TCP段结构

3.1 传输层服务

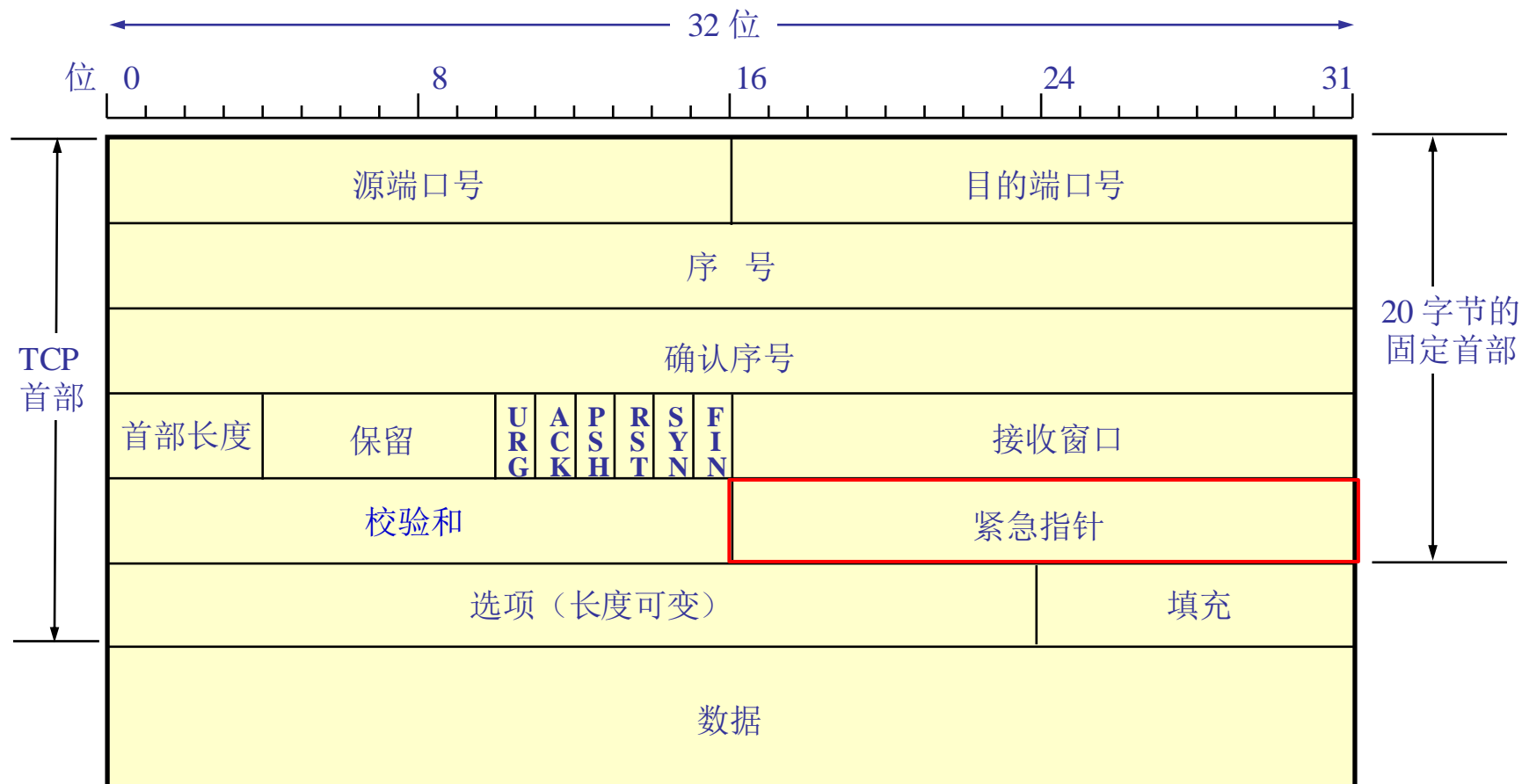
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 紧急指针字段占16位

- URG=1时才有效
- 指出紧急数据最后一个字节在数据中的位置



TCP段结构

3.1 传输层服务

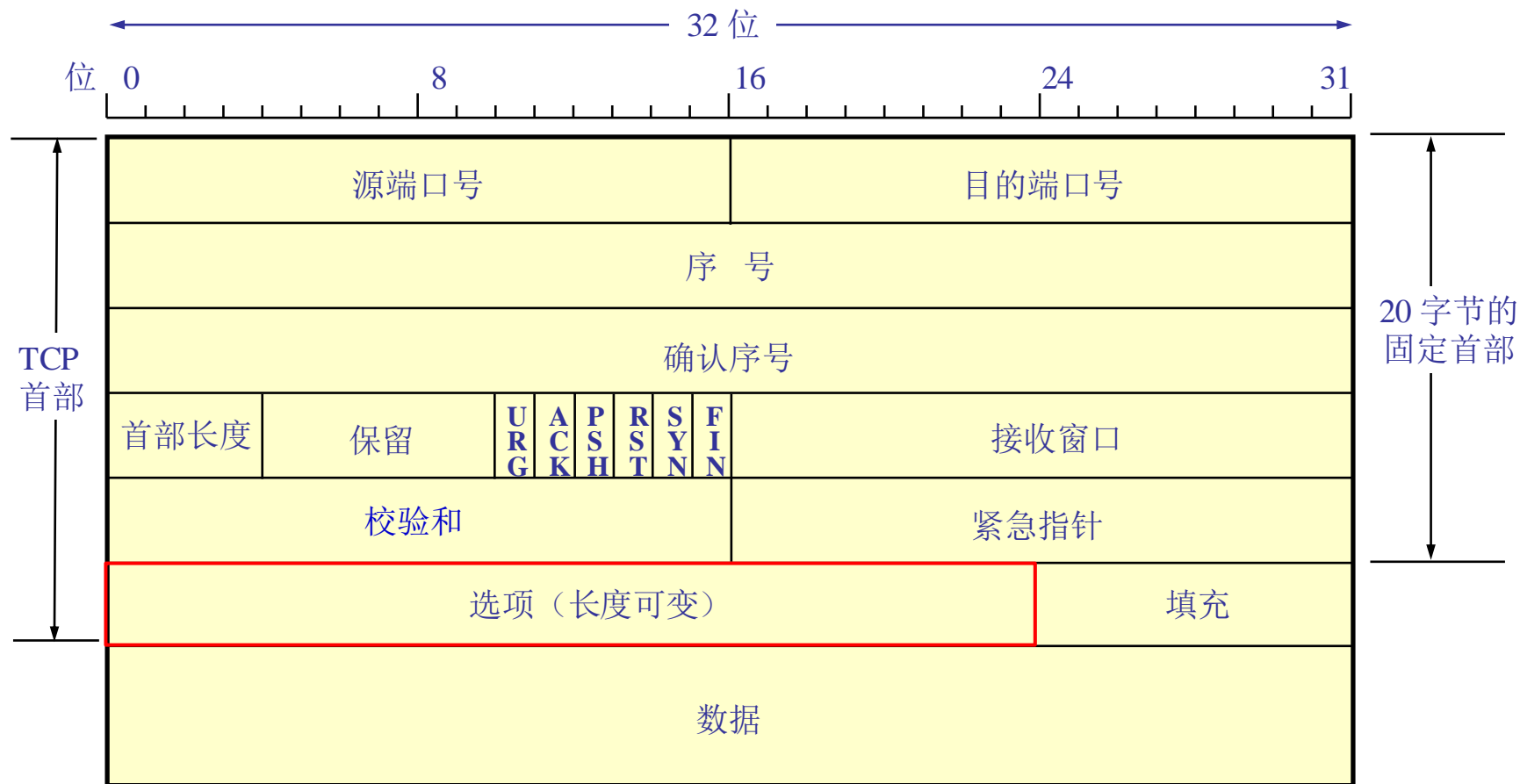
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 选项字段的长度可变

- 最大段长度MSS
- 时间戳
- SACK



TCP段结构

3.1 传输层服务

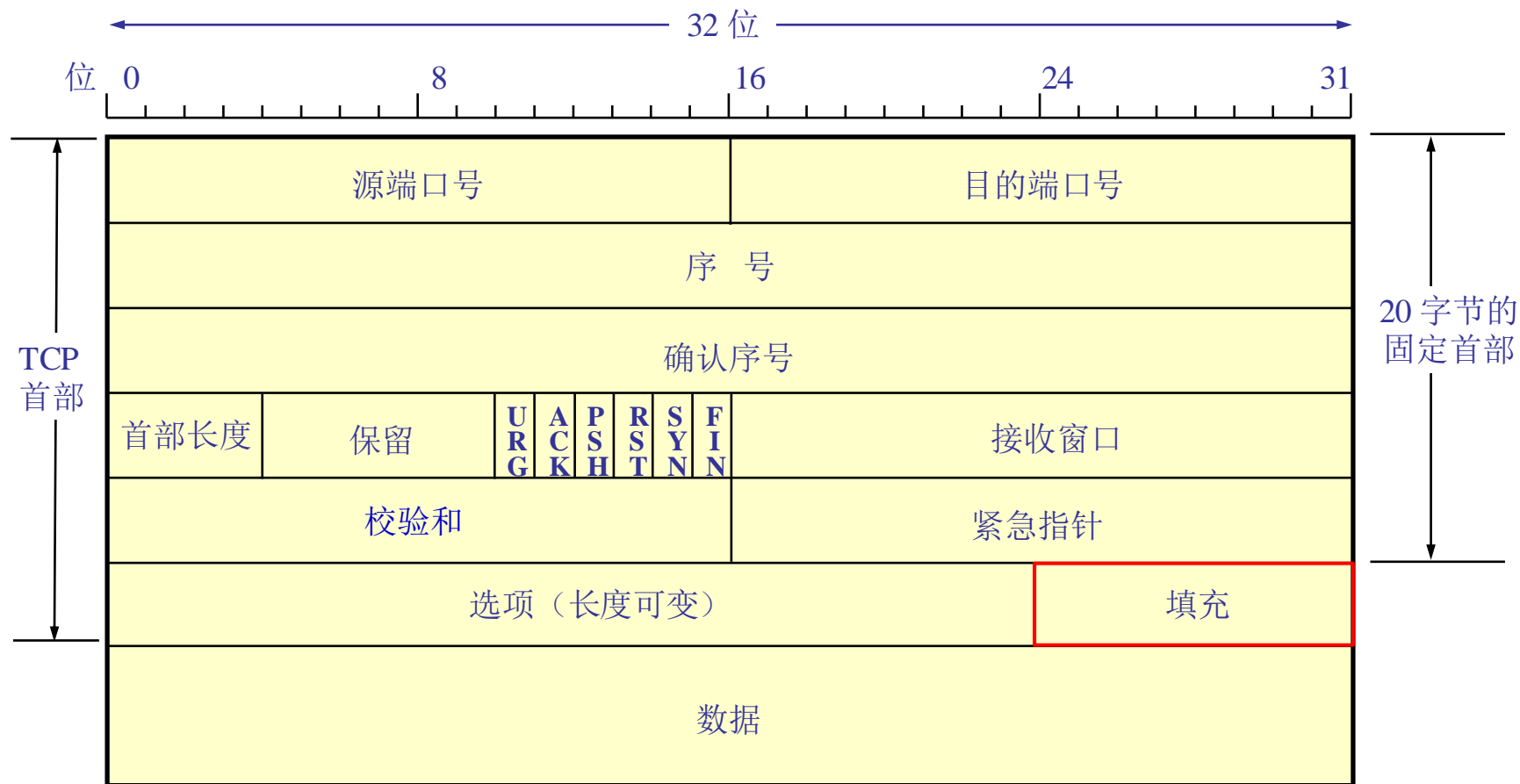
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



❖ 填充字段，长度为0~3个字节

- 取值全0



TCP: 序列号和ACK

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP段结构



序列号:

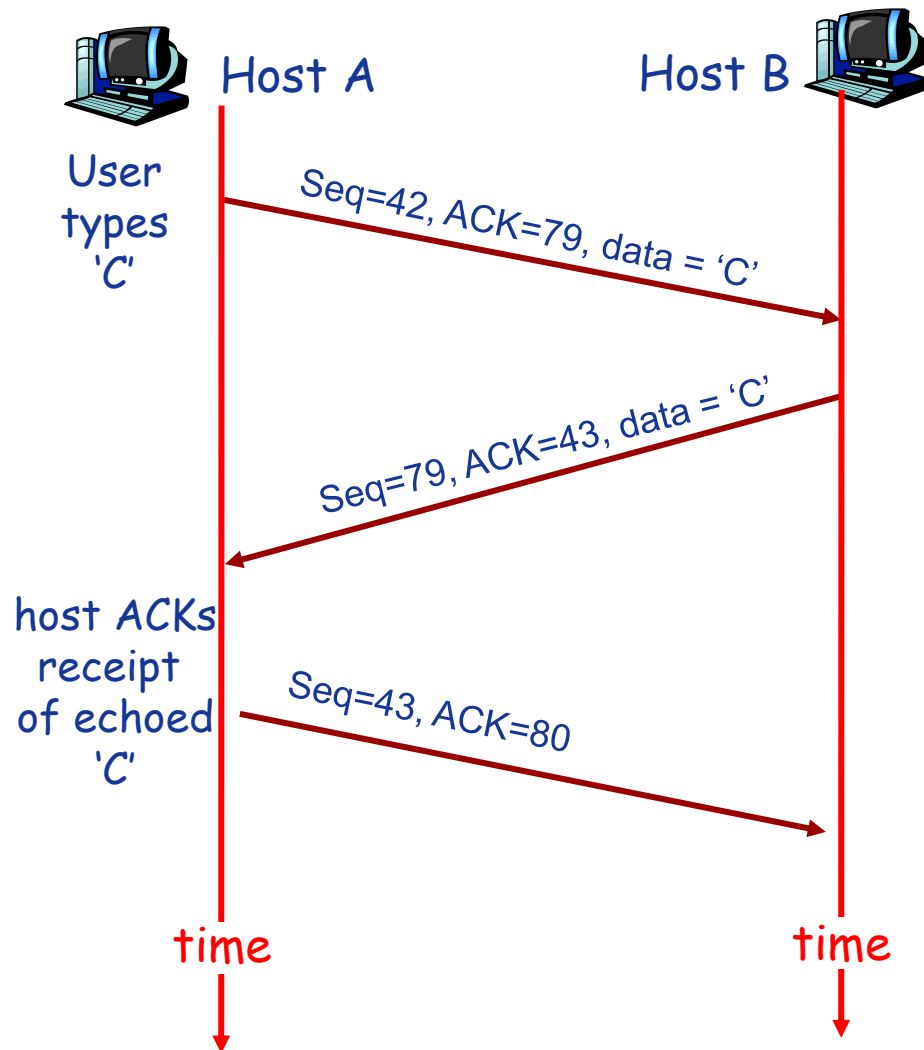
- 序号是段（ **Segment** ）中第1个字节的编号，而不是段的“连续”编号
- 建立**TCP**连接时，双方随机选择序列号

ACKs:

- **期望**接收到的下一个字节的序列号
- **累计确认**: 该序列号之前的所有字节均已被正确接收到

Q: 接收方如何处理乱序到达的段?

- **A:** **TCP**规范中没有规定，由**TCP**的实现者做出决策



简单telnet



TCP可靠数据传输概述

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输

- ❖ TCP在IP层的不可靠服务基础上实现可靠数据传输服务
- ❖ 流水线机制
- ❖ 累积确认
- ❖ TCP使用单一重传定时器
- ❖ 触发重传的事件
 - 超时
 - 收到重复ACK
- ❖ 渐进式





TCP RTT和超时

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输

❖ **问题：** 如何设置定时器的超时时间？

❖ 大于RTT

- 但是RTT是变化的

❖ 过短：

- 不必要的重传

❖ 过长：

- 对段丢失时间反应慢

❖ **问题：** 如何估计RTT？

❖ SampleRTT: 测量从段发出去到收到ACK的时间

- 忽略重传

❖ SampleRTT变化

- 测量多个SampleRTT，求平均值，形成RTT的估计值
EstimatedRTT

$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

指数加权移动平均

α 典型值：0.125



TCP RTT和超时

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输

定时器超时时间的设置?

- **EstimatedRTT + “安全边界”**
- **EstimatedRTT变化大→较大的边界**

测量RTT的变化值: SampleRTT与EstimatedRTT的差值

$$\text{DevRTT} = (1 - \beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically, $\beta = 0.25$)

定时器超时时间的设置:

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$





TCP发送方事件

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输



❖ 从应用层收到数据

- 创建Segment
- 序列号是Segment第一个字节的编号
- 开启计时器
- 设置超时时间: TimeoutInterval

❖ 超时

- 重传引起超时的段
- 重启定时器

❖ 收到ACK

- 如果确认此前未确认的段
 - 更新SendBase
 - 如果窗口中还有未被确认的分组, 重新启动定时器



TCP发送端程序

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输



```
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum
loop (forever) {
  switch(event)
  event: data received from application above
    create TCP segment with sequence number NextSeqNum
    if (timer currently not running)
      start timer
    pass segment to IP
    NextSeqNum = NextSeqNum + length(data)

  event: timer timeout
    retransmit not-yet-acknowledged segment with
      smallest sequence number
    start timer

  event: ACK received, with ACK field value of y
    if (y > SendBase) {
      SendBase = y
      if (there are currently not-yet-acknowledged segments)
        start timer
    }
} /* end of loop forever */
```



TCP重传示例

3.1 传输层服务

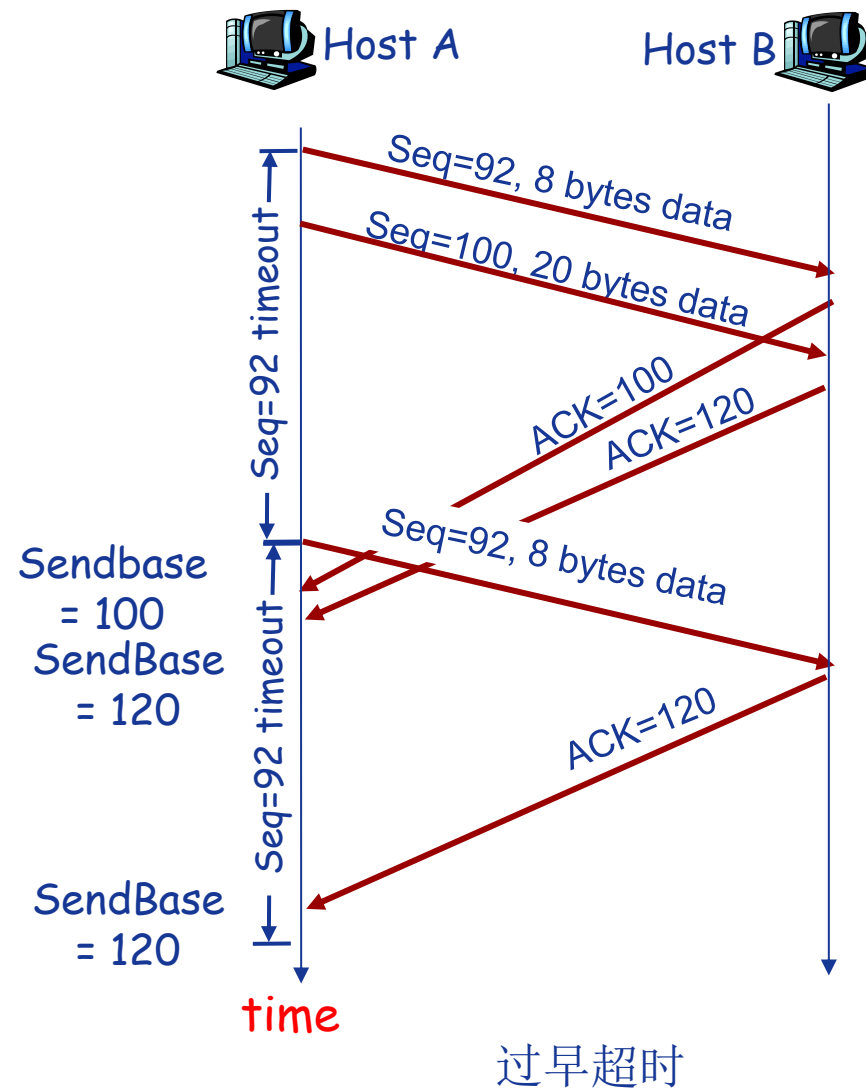
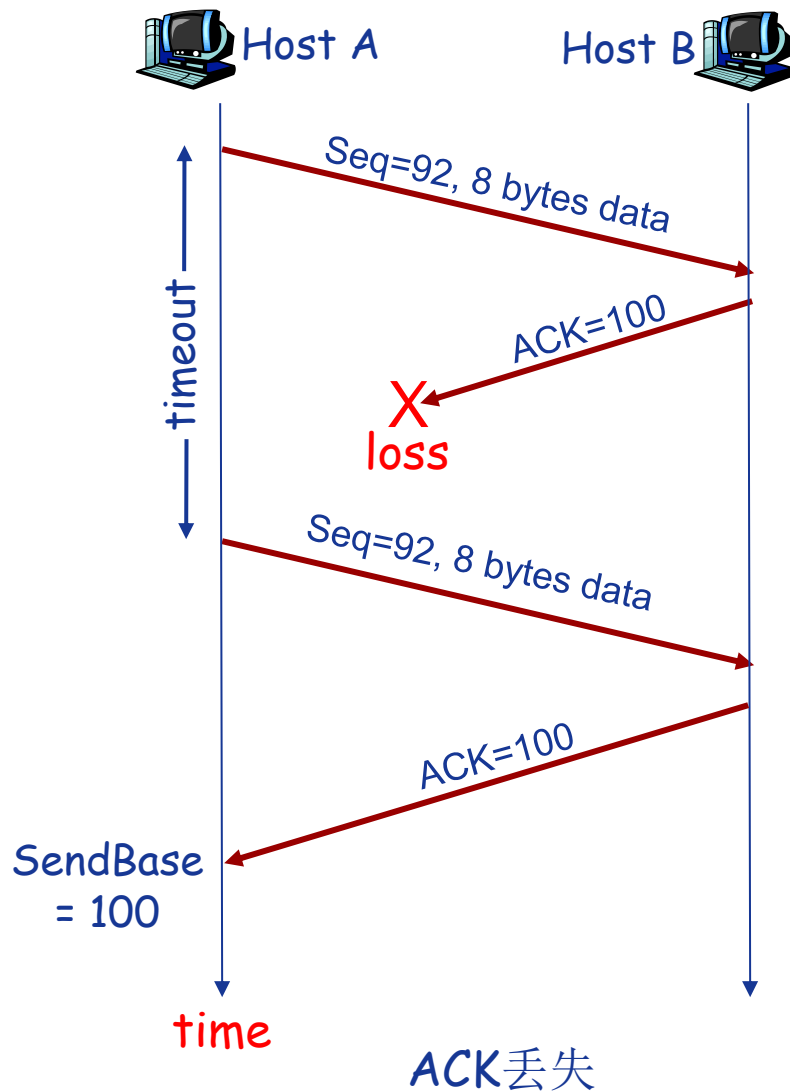
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输





TCP重传示例

3.1 传输层服务

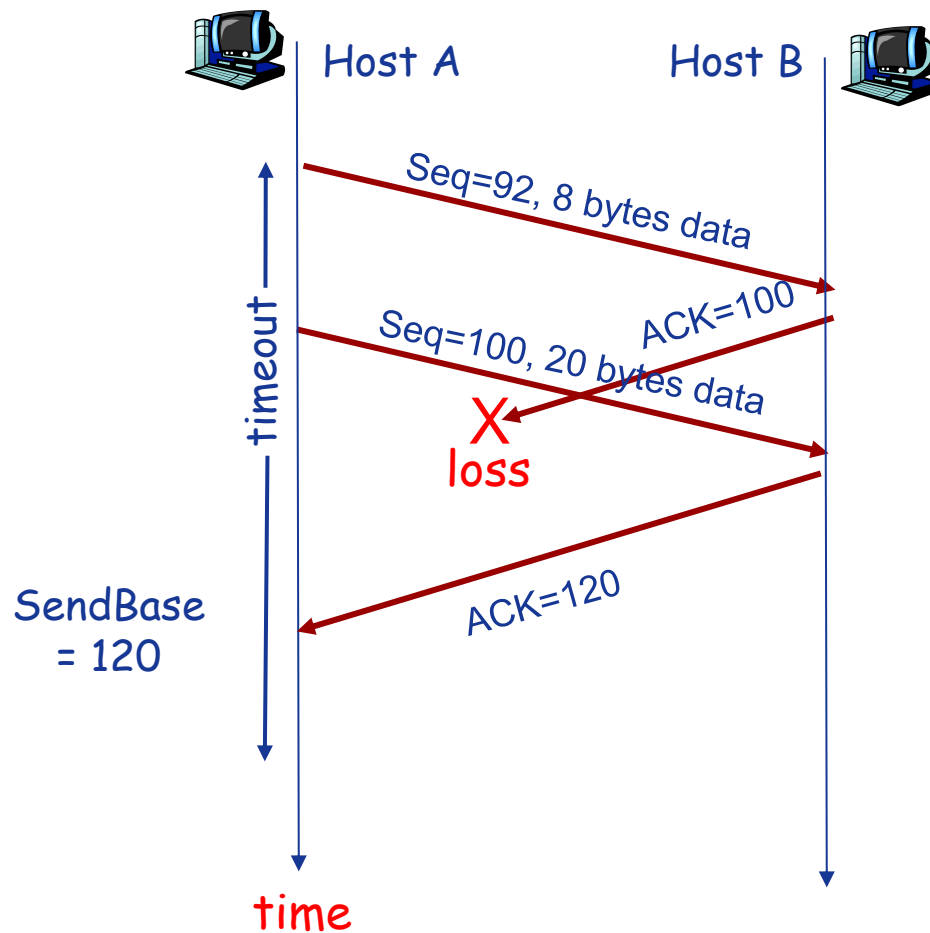
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输



累积确认



TCP ACK生成: RFC 1122, RFC 2581

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输



Event at Receiver	TCP Receiver action
Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
Arrival of in-order segment with expected seq #. One other segment has ACK pending	Immediately send single cumulative ACK, ACKing both in-order segments
Arrival of out-of-order segment higher-than-expect seq. # . Gap detected	Immediately send <i>duplicate ACK</i> , indicating seq. # of next expected byte
Arrival of segment that partially or completely fills gap	Immediate send ACK, provided that segment startsat lower end of gap

快速重传机制

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输



❖ TCP的实现中，如果发生超时，超时时间间隔将重新设置，即将超时时间间隔加倍，导致其**很大**

- 重发丢失的分组之前要等待很长时间

❖ 通过重复ACK检测分组丢失

- Sender会背靠背地发送多个分组
- 如果某个分组丢失，可能会引发多个重复的ACK

❖ 如果sender收到对同一数据的3个额外ACK，则假定该数据之后的段已经丢失

- **快速重传**：在定时器超时之前即进行重传

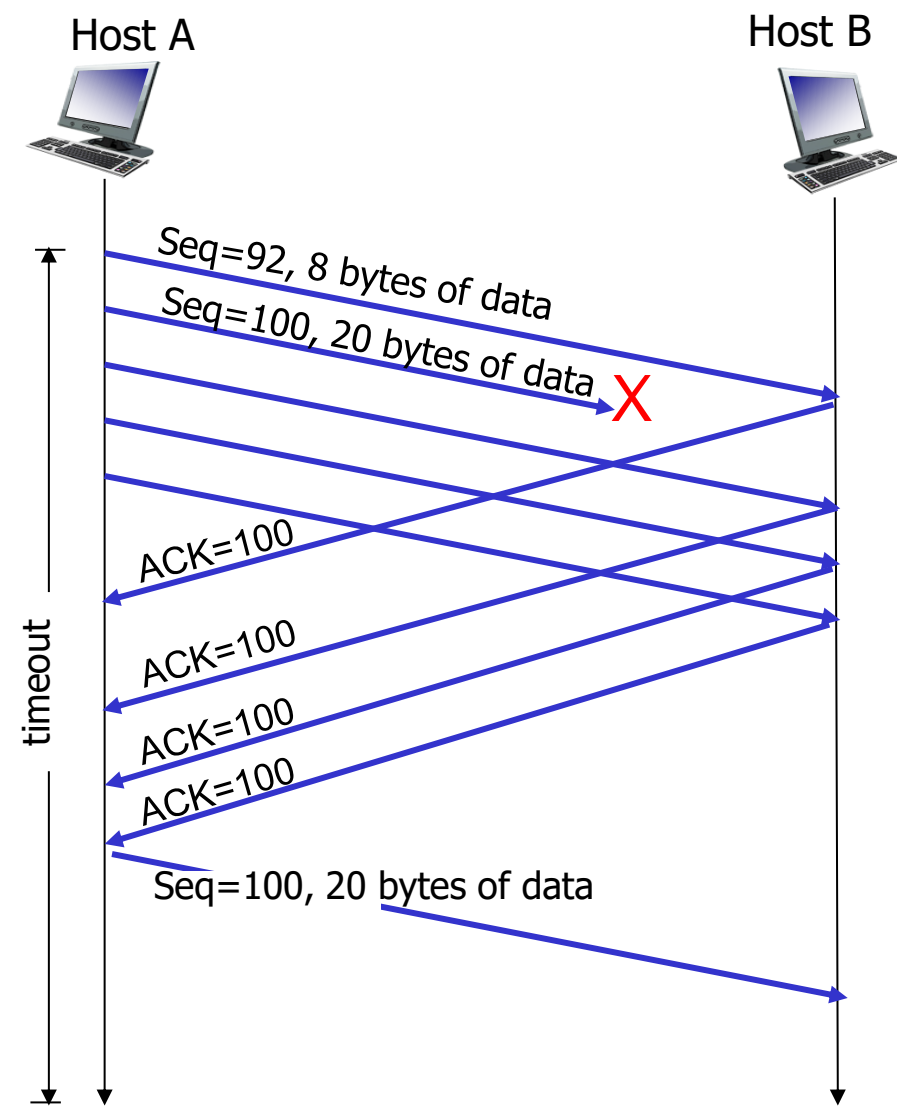


快速重传机制

- 3.1 传输层服务
- 3.2 传输层多路复用/分解
- 3.3 UDP协议
- 3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输



3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输

主机甲与主机乙间已建立一个TCP连接，主机甲向主机乙发送了两个连续的TCP段，分别包含300字节和500字节的有效载荷，第一个段的序列号为200，主机乙正确接收到两个段后，发送给主机甲的确认序列号是

- ☐ A 500
- ☐ B 700
- ☐ C 800
- ☒ D 1000

提交



单选题 1分

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

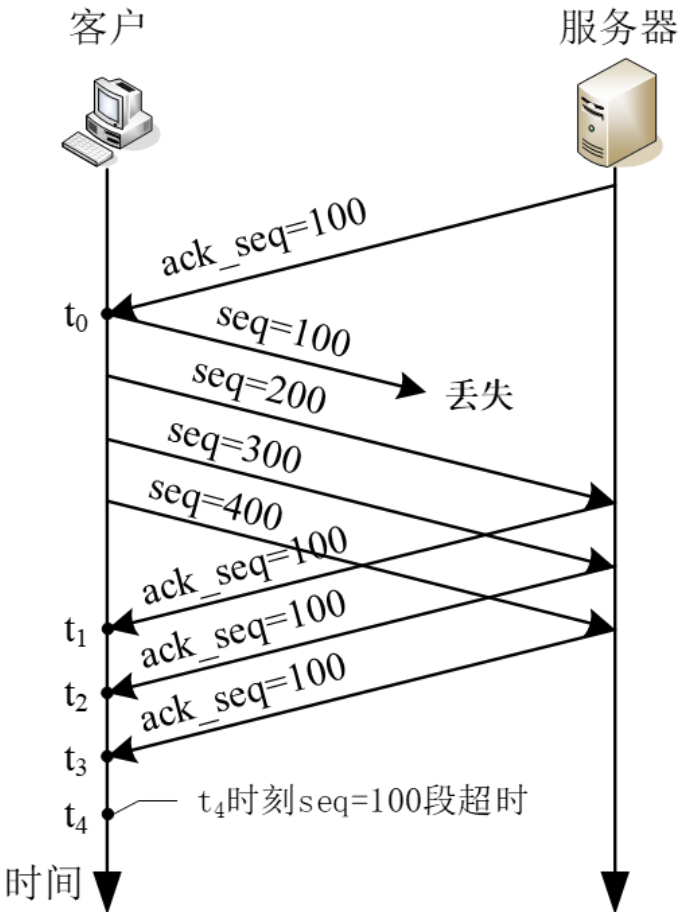
3.4 可靠数据传输原理

3.5 TCP协议

TCP可靠数据传输

某客户通过一个TCP连接向服务器发送数据的部分过程如下图所示。客户在 t_0 时刻第一次收到确认序列号 $ack_seq=100$ 的段，并发送序列号 $seq=100$ 的段，但发生丢失。若TCP支持快速重传，则客户重新发送 $seq=100$ 段的时刻是

- A t_1
- B t_2
- C t_3
- D t_4



提交



TCP流量控制

3.1 传输层服务

3.2 传输层多路复用/分解

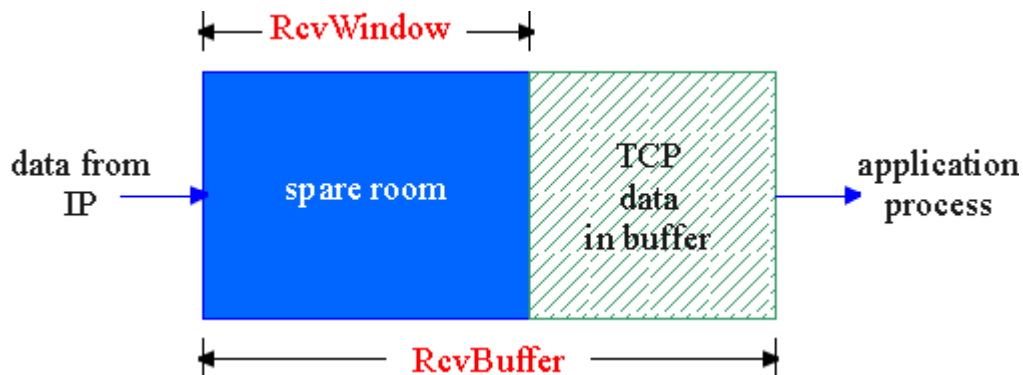
3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP流量控制

❖ 接收方为TCP连接分配buffer



flow control

发送方不会传输的太多、
太快以至于淹没接收方
(buffer溢出)

❖ 速度匹配机制

❑ 上层应用可能处理
buffer中数据的速度较慢



TCP流量控制

3.1 传输层服务

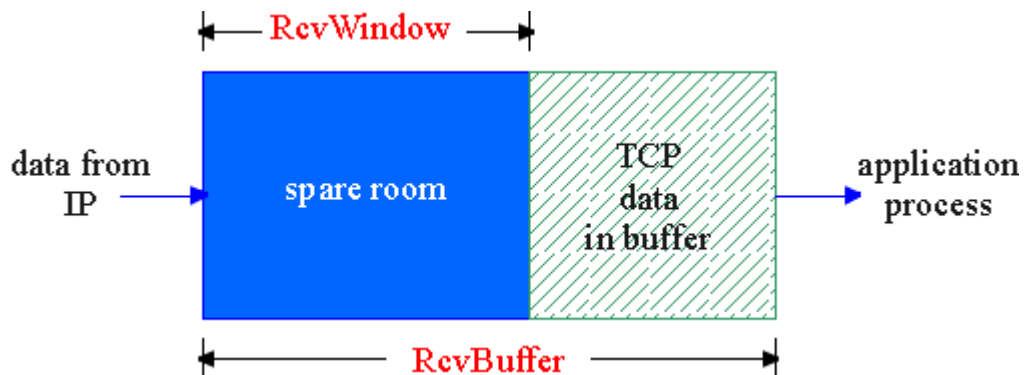
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP流量控制



(假定TCP receiver丢弃乱序的 segments)

❖ Buffer中的可用空间(spare room)

= RcvWindow

= RcvBuffer - [LastByteRcvd - LastByteRead]

❖ Receiver通过在Segment的头部字段将RcvWindow告诉Sender

❖ Sender限制自己已经发送的但还未收到ACK的数据不超过接收方的空闲RcvWindow尺寸

❖ Receiver告知Sender RcvWindow=0, 会出现什么情况?





TCP连接管理

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP连接管理

❖ TCP sender和receiver在传输数据前需要建立连接

❖ 初始化TCP变量

- Seq. #
- Buffer和流量控制信息

❖ Client: 连接发起者

```
Socket clientSocket = new  
Socket("hostname", "port  
number");
```

❖ Server: 等待客户连接请求

```
Socket connectionSocket =  
welcomeSocket.accept();
```

Three way handshake:

Step 1: client host sends TCP SYN segment to server

- specifies initial seq #
- no data

Step 2: server host receives SYN, replies with SYNACK segment

- server allocates buffers
- specifies server initial seq. #

Step 3: client receives SYNACK, replies with ACK segment, which may contain data



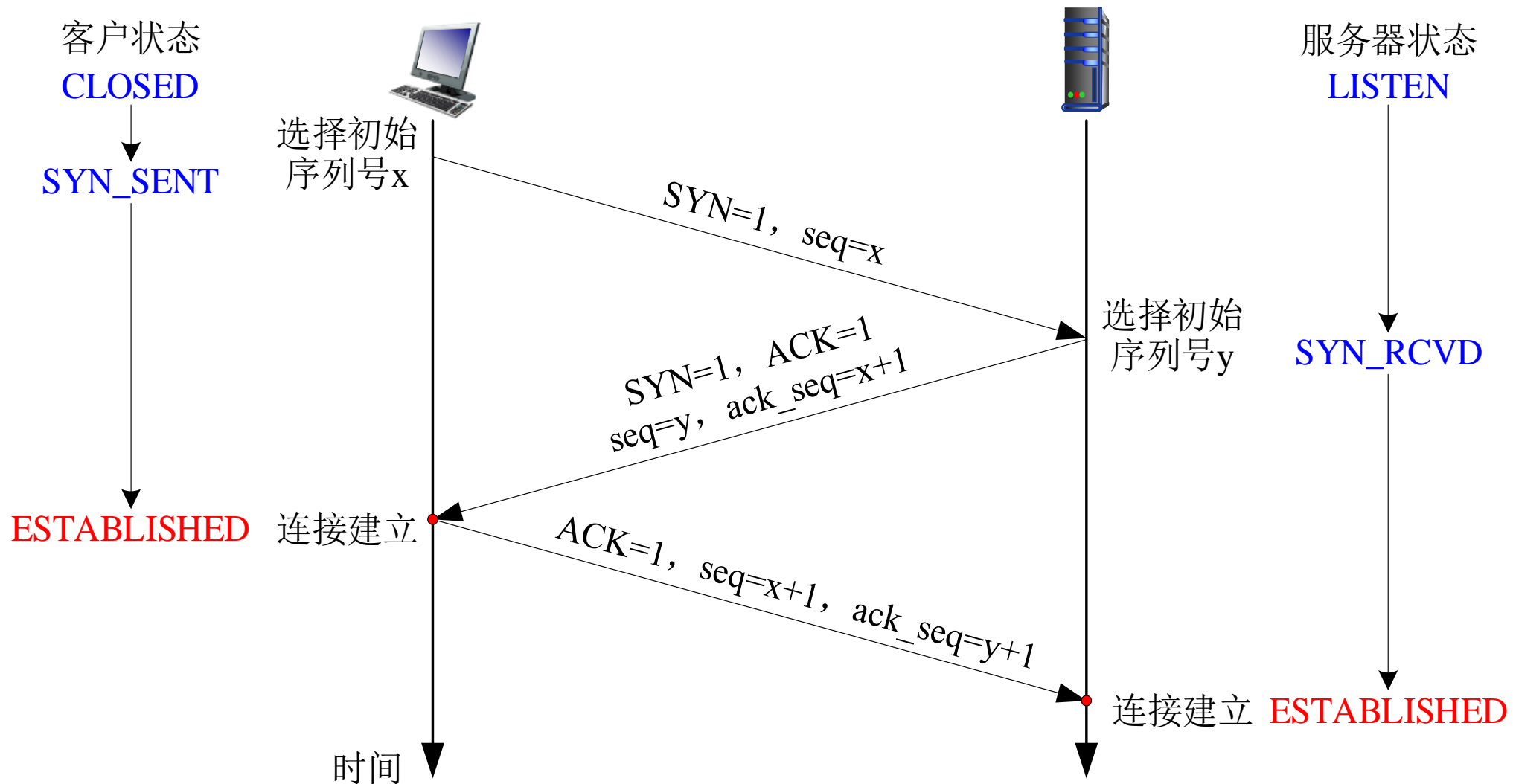


TCP连接管理：建立

- 3.1 传输层服务
- 3.2 传输层多路复用/分解
- 3.3 UDP协议
- 3.4 可靠数据传输原理

3.5 TCP协议

TCP连接管理



TCP连接管理：关闭

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP连接管理



Closing a connection:

client closes socket: `clientSocket.close()`;

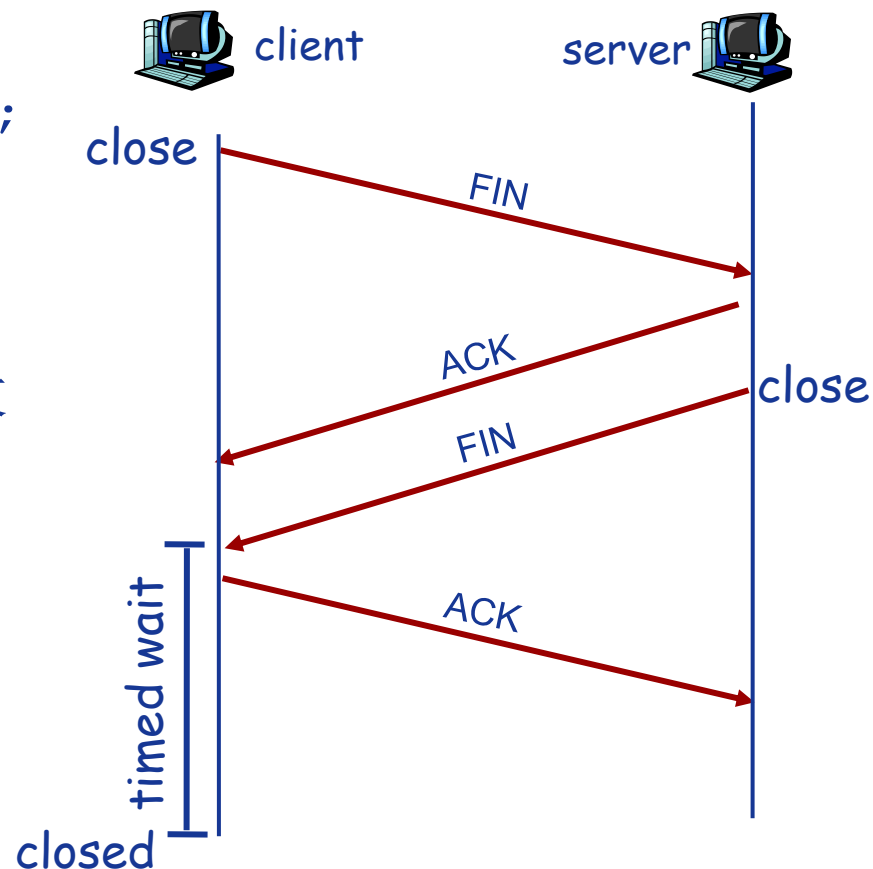
Step 1: client向server发送TCP FIN 控制
segment

Step 2: server 收到FIN, 回复ACK. 关闭连接, 发
送FIN.

Step 3: client 收到FIN, 回复ACK.

- 进入“等待”–如果收到FIN, 会重新发送
ACK

Step 4: server收到ACK. 连接关闭.





TCP连接管理：断连过程

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP连接管理

客户状态
ESTABLISHED

FIN_WAIT_1

FIN_WAIT_2

TIME_WAIT

CLOSED

能收不能
发数据能收不能
发数据延时等待
2MSL

时间

服务器状态
ESTABLISHED

CLOSE_WAIT

LAST_ACK

CLOSED

仍能发
数据不再发
数据

关闭连接

FIN=1, seq=u

ACK=1, seq=v, ack_seq=u+1

FIN=1, ACK=1
seq=w, ack_seq=u+1

ACK=1, seq=u+1, ack_seq=w+1

关闭连接



TCP连接管理

3.1 传输层服务

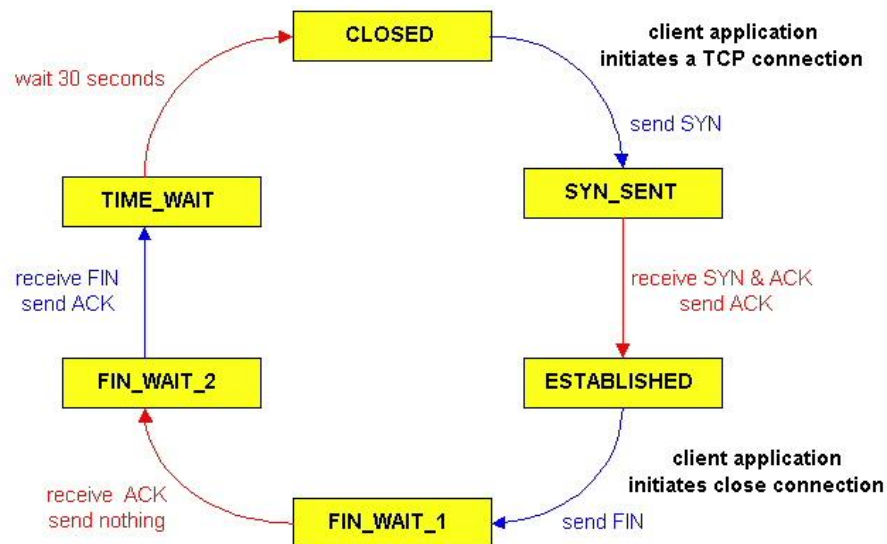
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

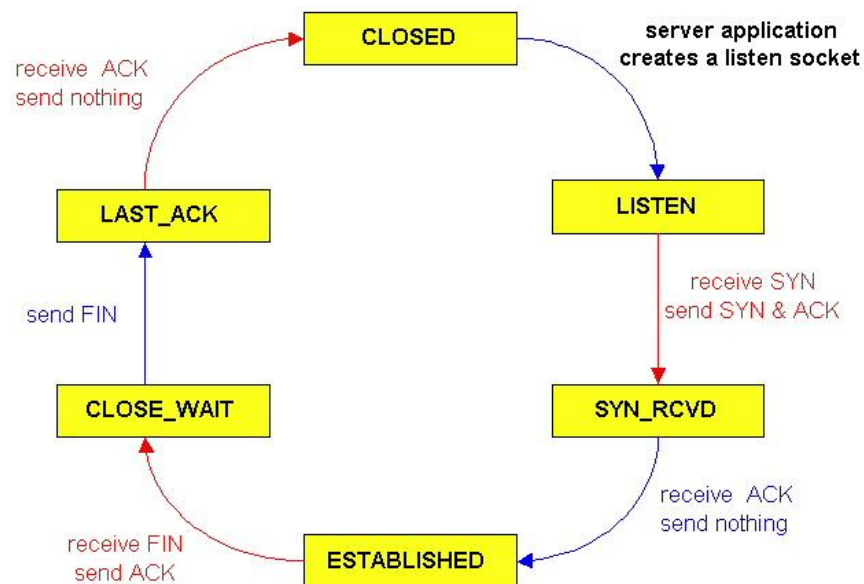
3.5 TCP协议

TCP连接管理



TCP client lifecycle

TCP server lifecycle



3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP连接管理



若主机甲主动发起一个与主机乙的**TCP**连接，甲、乙选择的初始序列号分别为**2018**和**2020**，则第三次握手**TCP**段的确认序列号是

- ☐ A 2019
- ☐ B 2020
- ☒ C 2021
- ☐ D 2022

提交



TCP拥塞控制的基本原理

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制

❖ Sender限制发送速率

$\text{LastByteSent} - \text{LastByteAcked}$

$\leq \text{CongWin}$

$$\text{rate} \approx \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$

❖ CongWin:

- 动态调整以改变发送速率
- 反映所感知到的网络拥塞

问题：如何感知网络拥塞？

❖ Loss事件=timeout或3个重复ACK

❖ 发生loss事件后，发送方降低速率

如何合理地调整发送速率？

❖ 加性增—乘性减: AIMD

❖ 慢启动: SS





加性增一乘性减: AIMD

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

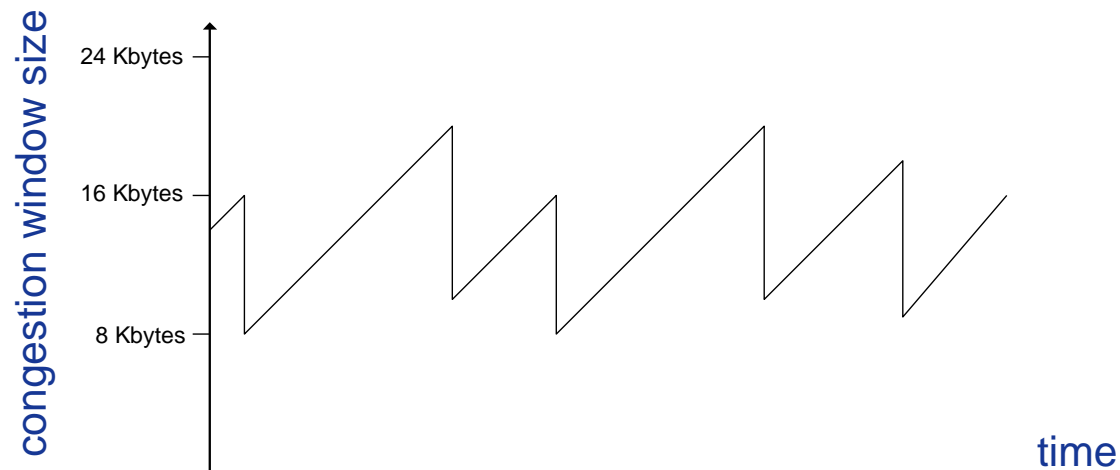
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制

锯齿行为: 探测
可用带宽

- ❖ **原理**: 逐渐增加发送速率, 谨慎探测可用带宽, 直到发生丢包
- ❖ **方法**: AIMD
 - **Additive Increase**: 每个RTT将CongWin增大1个MSS-拥塞避免
 - **Multiplicative Decrease**: 发生丢包后将CongWin减半



TCP慢启动: SS

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制



❖ TCP连接建立时, CongWin=1

- 例: MSS=500 byte,
RTT=200msec
- 初始速率=20k bps

❖ 可用带宽可能远远高于 初始速率:

- 希望快速增长

❖ 原理:

- 当连接开始时, 指数性增长

Slowstart algorithm

```
initialize: Congwin = 1
for (each segment ACKed)
    Congwin++
until (loss event OR
      CongWin > threshold)
```



TCP慢启动: SS

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

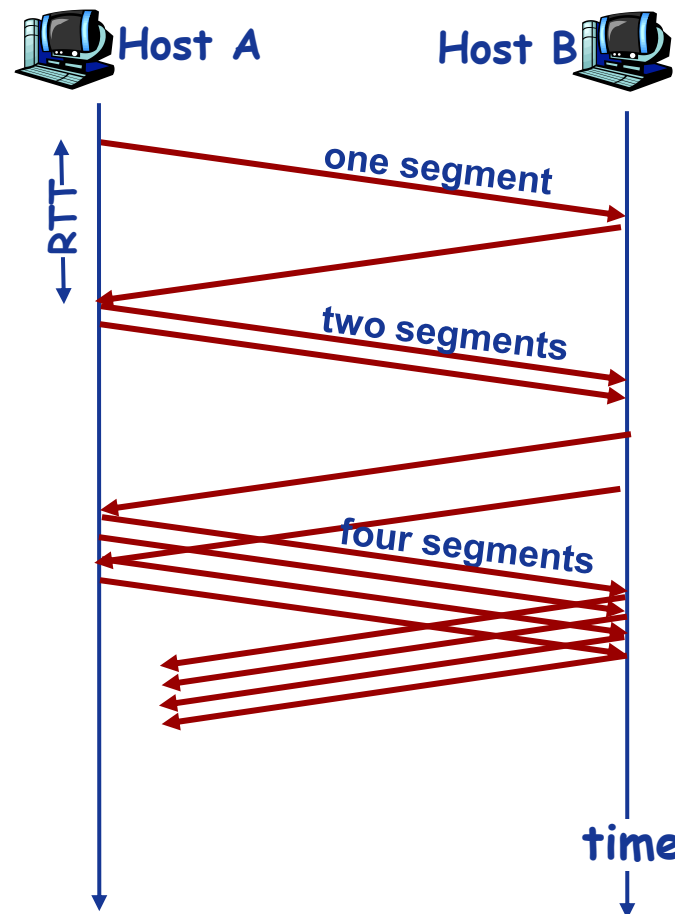
TCP拥塞控制



❖ 指数性增长

- 每个RTT将CongWin翻倍
- 收到每个ACK进行CongWin++操作

❖ 初始速率很慢，但是快速攀升





TCP慢启动: SS

3.1 传输层服务

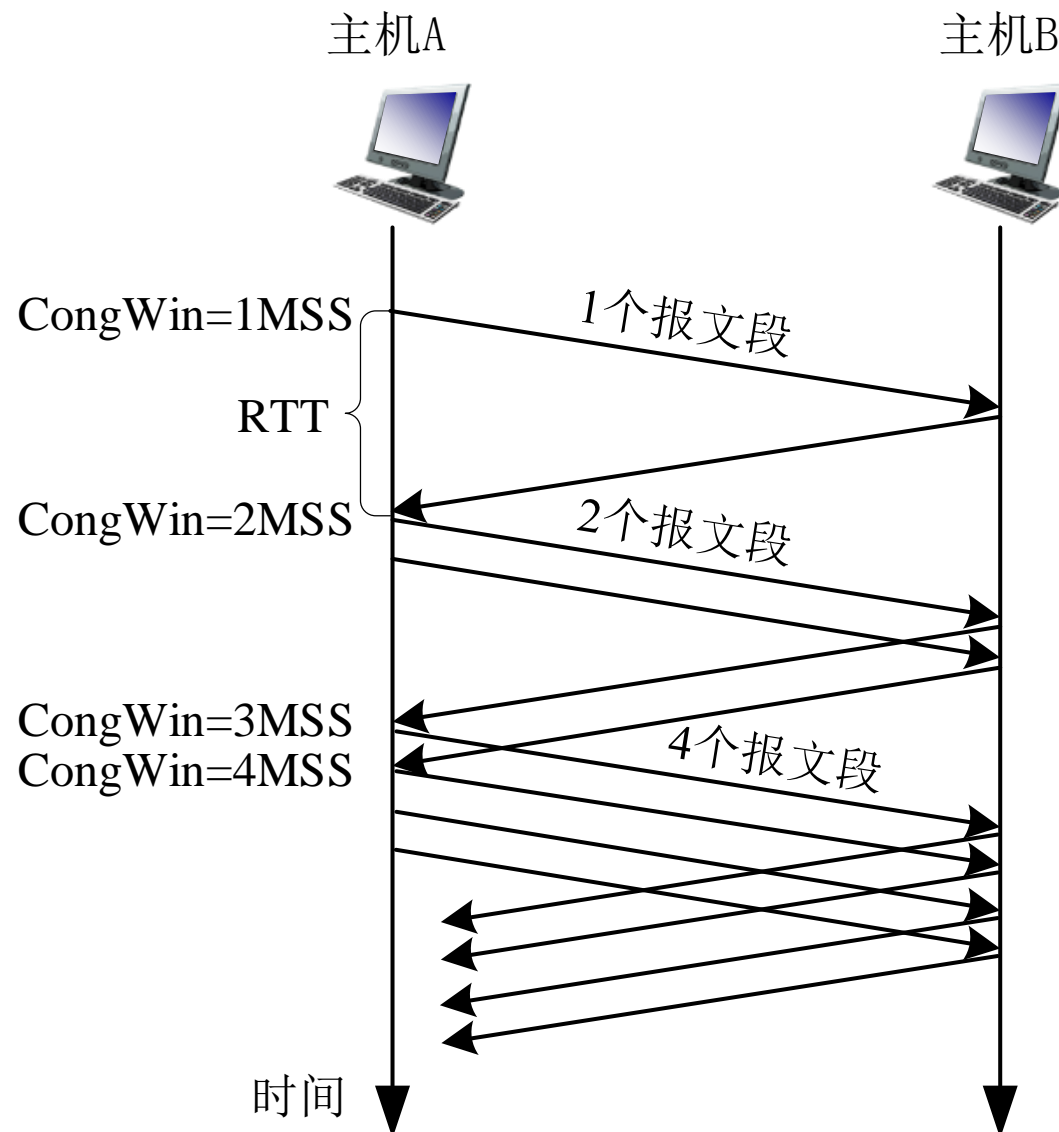
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制





Threshold变量

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制



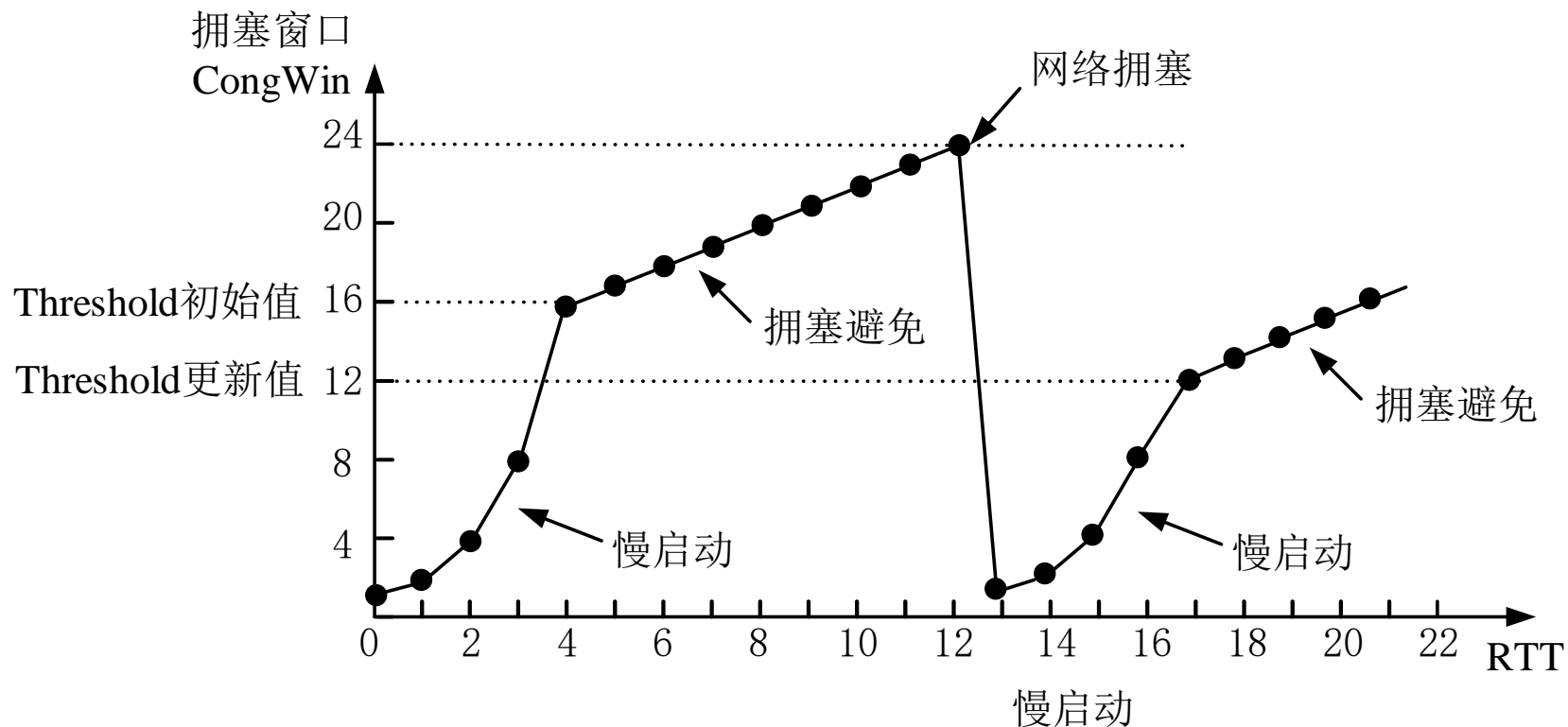
Q:何时应该指数性增长切换为线性增长(拥塞避免)?

A: 当CongWin达到Loss事件前值的1/2时.

实现方法:

❖ 变量 **Threshold**

❖ Loss事件发生时, **Threshold** 被设为Loss事件前**CongWin** 值的1/2。





Loss事件的处理

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制

❖ 3个重复ACKs:

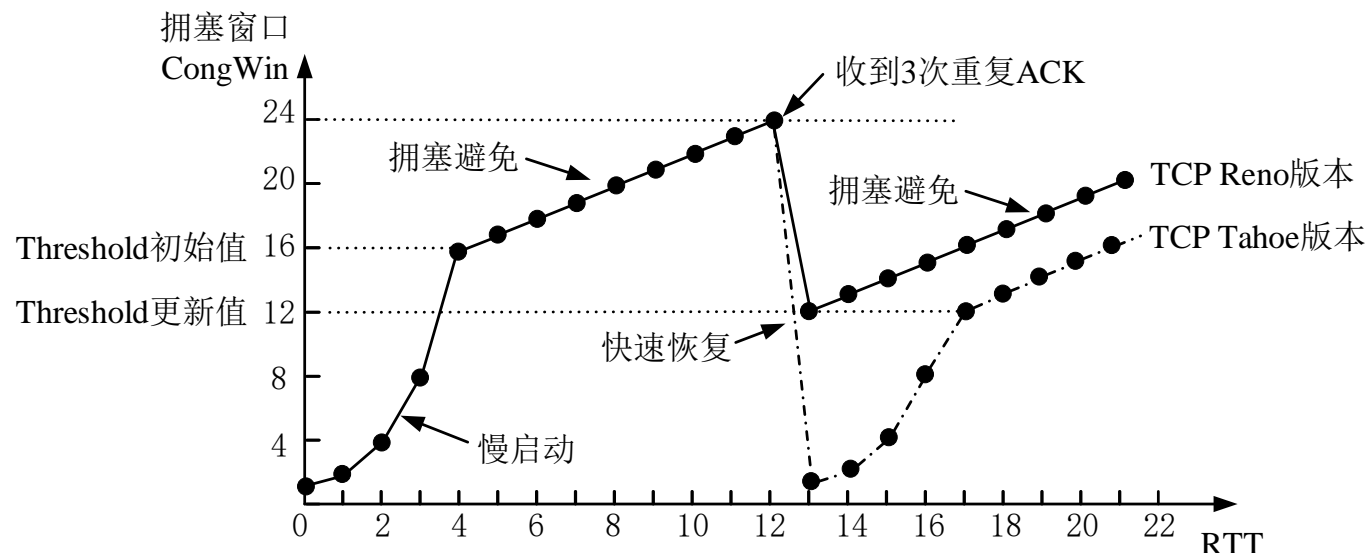
- CongWin切换到一半
- 然后线性增长

❖ Timeout事件:

- CongWin直接设为1个MSS
- 然后指数增长
- 达到threshold后, 再线性增长

Philosophy:

- ❑ 3个重复ACKs表示网络还能够传输一些 segments
- ❑ timeout事件表明拥塞更为严重





TCP拥塞控制：总结

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制

- ❖ When CongWin is below Threshold, sender in **slow-start** phase, window grows exponentially.
- ❖ When CongWin is above Threshold, sender is in **congestion-avoidance** phase, window grows linearly.
- ❖ When a **triple duplicate ACK** occurs, Threshold set to $\text{CongWin}/2$ and CongWin set to Threshold.
- ❖ When **timeout** occurs, Threshold set to $\text{CongWin}/2$ and CongWin is set to 1 MSS.





TCP拥塞控制

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

State	Event	TCP Sender Action	Commentary
Slow Start (SS)	ACK receipt for previously unacked data	$CongWin = CongWin + MSS$, If $(CongWin > Threshold)$ set state to "Congestion Avoidance"	Resulting in a doubling of CongWin every RTT
Congestion Avoidance (CA)	ACK receipt for previously unacked data	$CongWin = CongWin + MSS * (MSS / CongWin)$	Additive increase, resulting in increase of CongWin by 1 MSS every RTT
SS or CA	Loss event detected by triple duplicate ACK	$Threshold = CongWin / 2$, $CongWin = Threshold$, Set state to "Congestion Avoidance"	Fast recovery, implementing multiplicative decrease. CongWin will not drop below 1 MSS.
SS or CA	Timeout	$Threshold = CongWin / 2$, $CongWin = 1 MSS$, Set state to "Slow Start"	Enter slow start
SS or CA	Duplicate ACK	Increment duplicate ACK count for segment being acked	CongWin and Threshold not changed

TCP拥塞控制





TCP拥塞控制算法

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制



Th = ?

CongWin = 1 MSS

/* slow start or exponential increase */

While (No Packet Loss and CongWin < Th) {

 send CongWin TCP segments

 for each ACK increase CongWin by 1

}

/* congestion avoidance or linear increase */

While (No Packet Loss) {

 send CongWin TCP segments

 for CongWin ACKs, increase CongWin by 1

}

Th = CongWin/2

If (3 Dup ACKs) CongWin = Th;

If (timeout) CongWin=1;



例题

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制



【例1】一个TCP连接总是以1 KB的最大段长发送TCP段，发送方有足够多的数据要发送。当拥塞窗口为16 KB时发生了超时，如果接下来的4个RTT（往返时间）时间内的TCP段的传输都是成功的，那么当第4个RTT时间内发送的所有TCP段都得到肯定应答时，拥塞窗口大小是多少？

【解】 $\text{threshold} = 16/2 = 8 \text{ KB}$, $\text{CongWin} = 1 \text{ KB}$, 1个RTT后, $\text{CongWin} = 2 \text{ KB}$, 2个RTT后, $\text{CongWin} = 4 \text{ KB}$, 3个RTT后, $\text{CongWin} = 8 \text{ KB}$, Slowstart is over; 4个RTT后, $\text{CongWin} = 9 \text{ KB}$

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制



主机甲和主机乙之间已建立了一个TCP连接，TCP最大段长度为1000字节。若主机甲的当前拥塞窗口为4 000字节，在主机甲向主机乙连续发送2个最大段后，成功收到主机乙发送的对第一个段的确认段，确认段中通告的接收窗口大小为2000字节，则此时主机甲还可以向主机乙发送的最大字节数是

- ☒ A 1000
- ☐ B 2000
- ☐ C 3000
- ☐ D 4000

提交



TCP的吞吐率

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP吞吐率分析



- ❖ 给定拥塞窗口大小和RTT，TCP的平均吞吐率是多少？
 - 忽略掉Slow start
- ❖ 假定发生超时时CongWin的大小为W，吞吐率是 W/RTT
- ❖ 超时后， $CongWin=W/2$ ，吞吐率是 $W/2RTT$
- ❖ 平均吞吐率为： $0.75W/RTT$

TCP的吞吐率

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP吞吐率分析

❖ 举例：每个Segment有1500个byte, RTT是100ms, 希望获得10Gbps的吞吐率

- $\text{throughput} = W * \text{MSS} * 8 / \text{RTT}$, 则
- $W = \text{throughput} * \text{RTT} / (\text{MSS} * 8)$
- $\text{throughput} = 10\text{Gbps}$, 则 $W = 83,333$

❖ 窗口大小为83,333



TCP的吞吐率

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP吞吐率分析



❖ 吞吐率与丢包率(loss rate, L)的关系

- CongWin从 $W/2$ 增加至 W 时出现第一个丢包，那么一共发送的分组数为

$$W/2 + (W/2 + 1) + (W/2 + 2) + \dots + W = 3W^2/8 + 3W/4$$

- W 很大时， $3W^2/8 \gg 3W/4$ ，因此 $L \approx 8/(3W^2)$

$$W = \sqrt{\frac{8}{3L}} \quad \text{Throughput} = \frac{0.75 \cdot MSS \cdot \sqrt{\frac{8}{3L}}}{RTT} \approx \frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

❖ $L = 2 \cdot 10^{-10}$ **Wow!!!**

❖ 高速网络下需要设计新的TCP



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

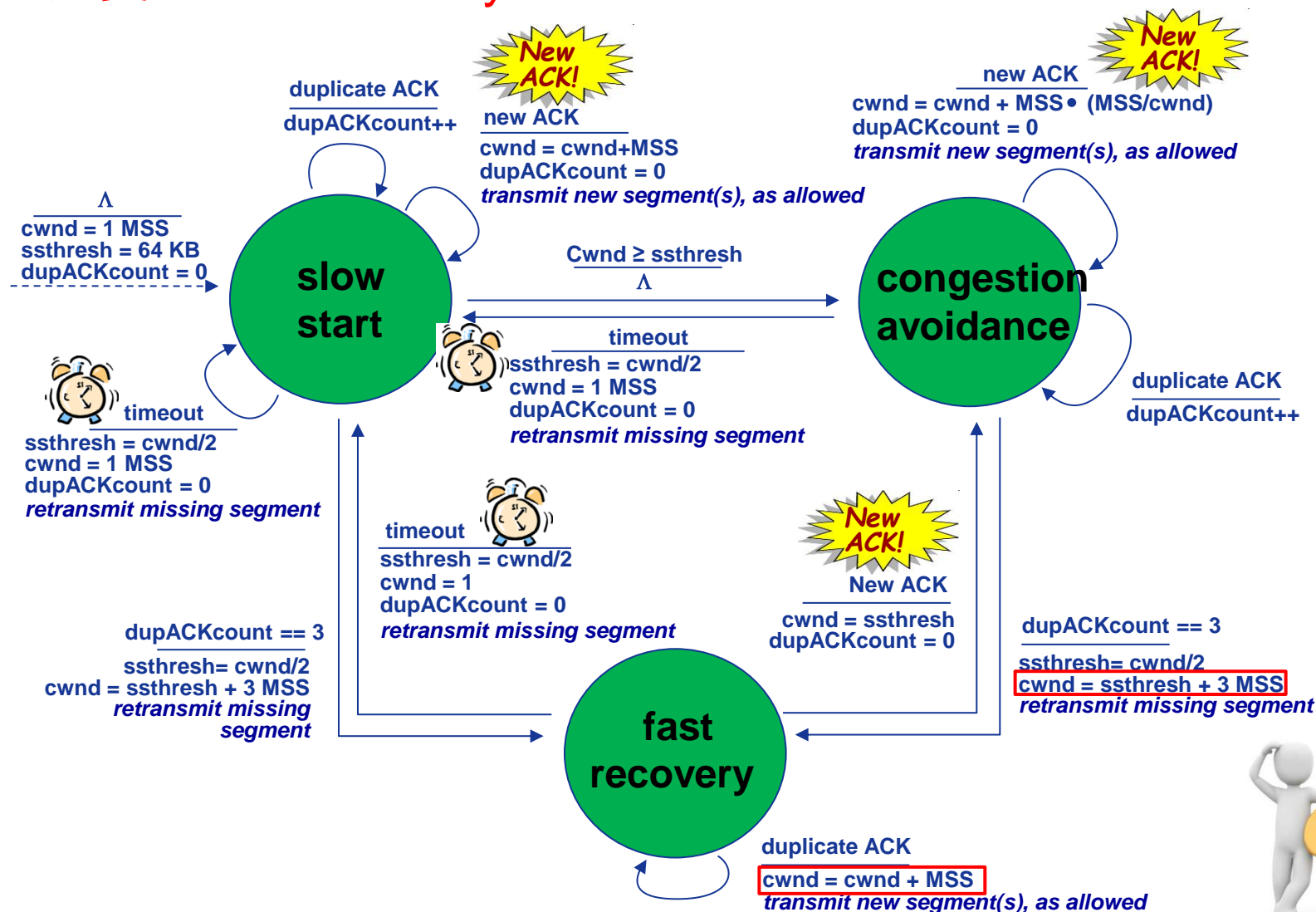
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进



❖ 快速回复 (fast recovery)



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

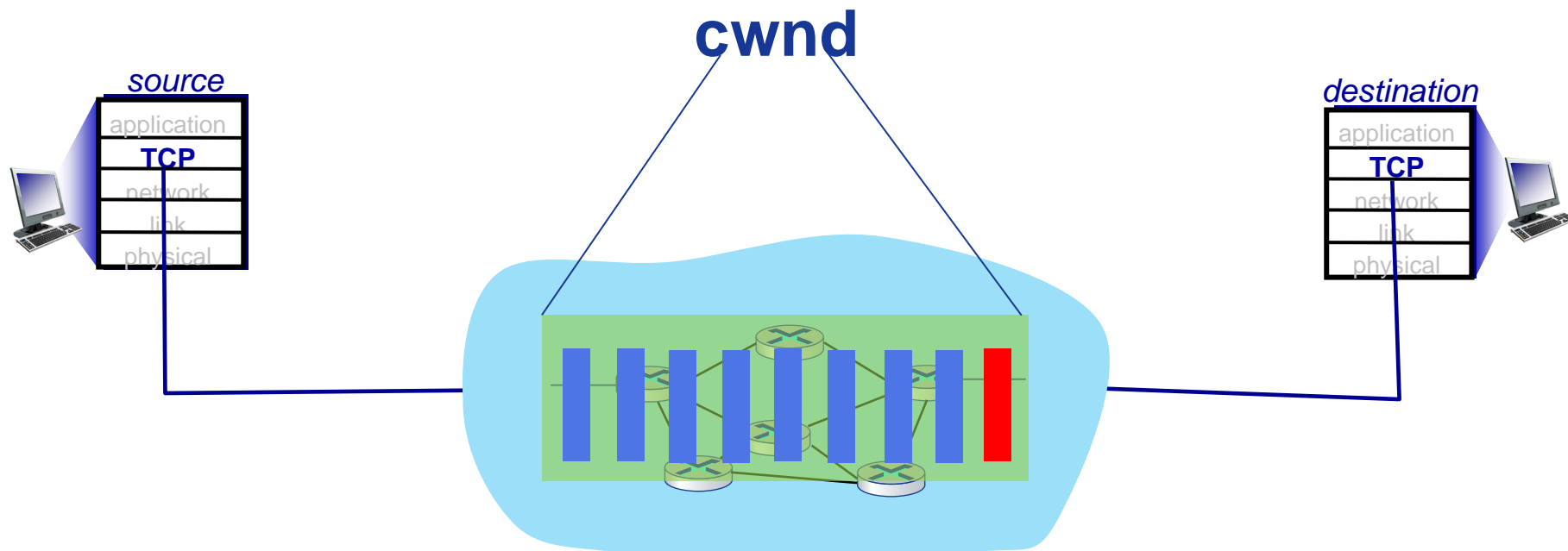
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进

❖ 快速回复 (fast recovery)

- 为什么窗口要膨胀？
- 为什么会出现3次重复确认？





TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

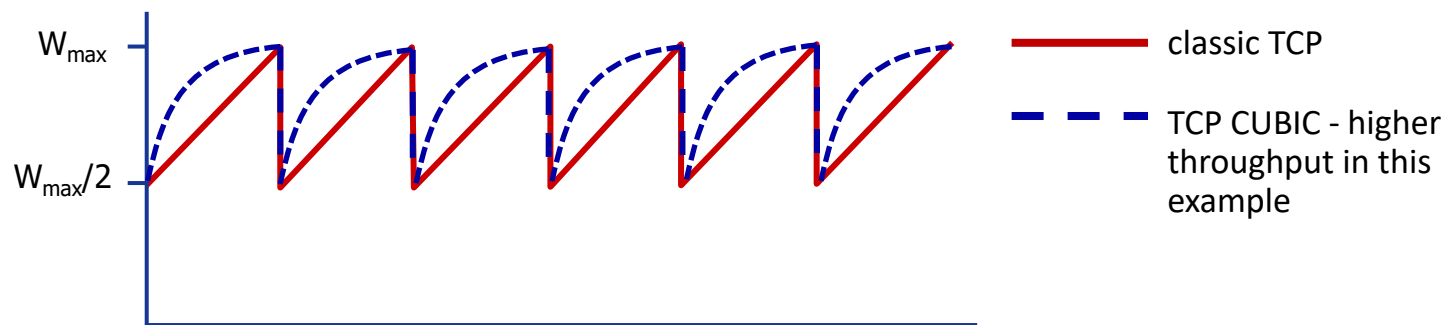
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进

❖ TCP CUBIC

- Is there a better way than AIMD to “probe” for usable bandwidth?
- Insight/intuition:
 - W_{\max} : sending rate at which congestion loss was detected
 - congestion state of bottleneck link probably (?) hasn't changed much
- after cutting rate/window in half on loss, initially ramp to W_{\max} *faster*, but then approach W_{\max} more *slowly*



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

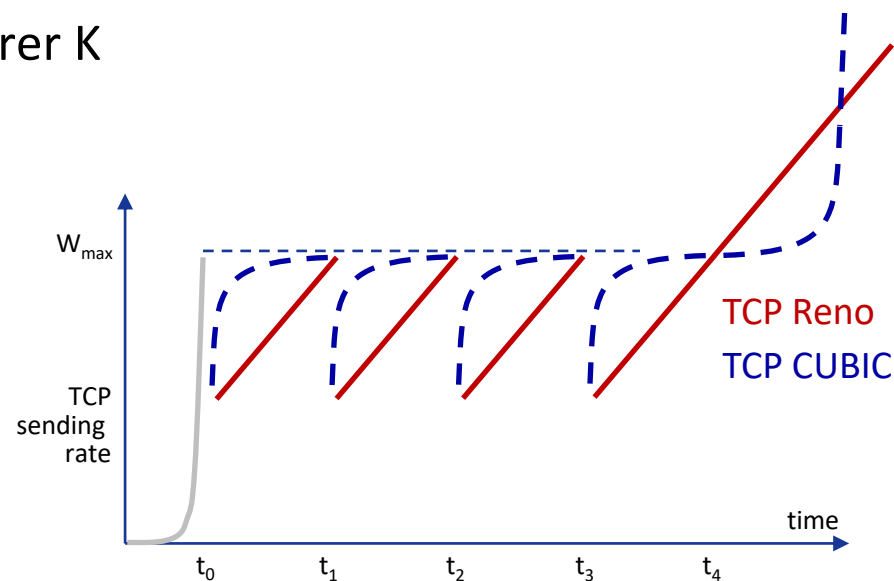
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进

❖ TCP CUBIC

- K: point in time when TCP window size will reach W_{\max}
 - K itself is tuneable
- increase W as a function of the *cube* of the distance between current time and K
 - larger increases when further away from K
 - smaller increases (cautious) when nearer K
- TCP CUBIC default in Linux, most popular TCP for popular Web servers



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

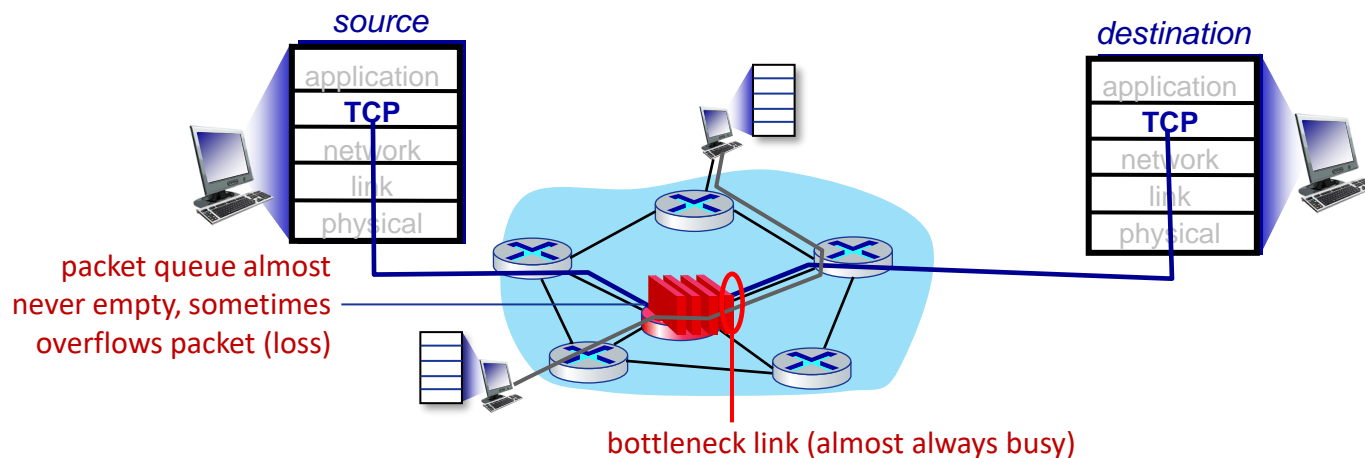
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进

❖ Delay-based TCP congestion control

- TCP (classic, CUBIC) increase TCP's sending rate until packet loss occurs at some router's output: the *bottleneck link*



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

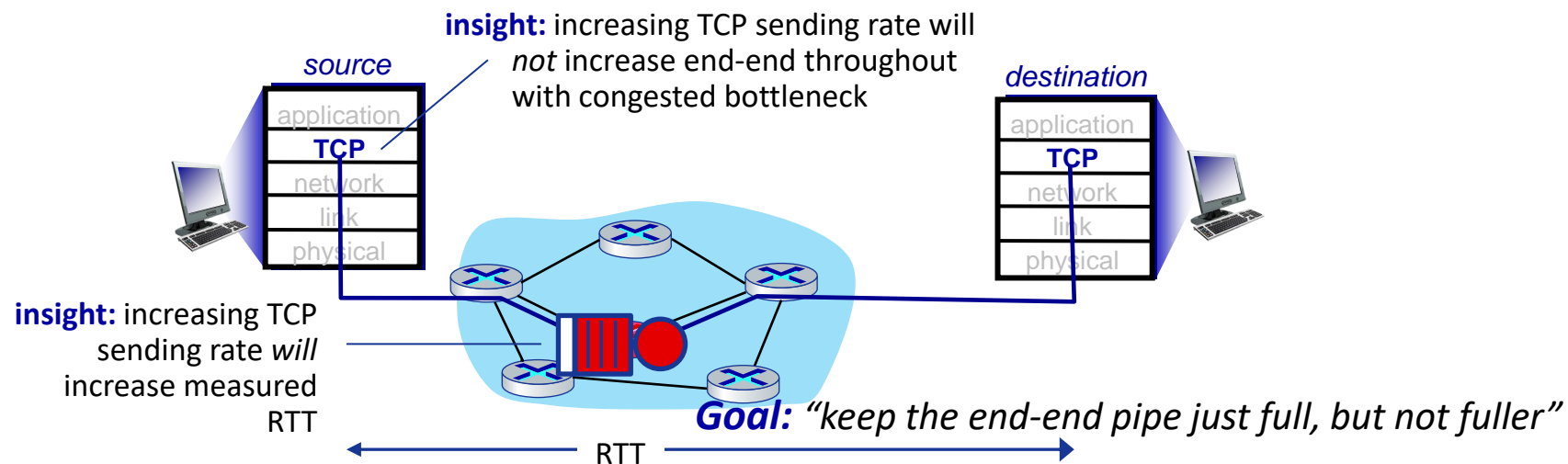
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进

❖ Delay-based TCP congestion control

- TCP (classic, CUBIC) increase TCP's sending rate until packet loss occurs at some router's output: the *bottleneck link*
- understanding congestion: useful to focus on congested bottleneck link



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

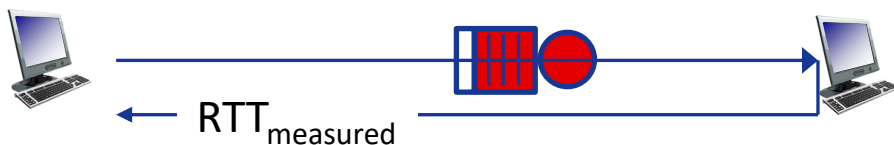
3.4 可靠数据传输原理

3.5 TCP协议

TCP拥塞控制改进

❖ Delay-based TCP congestion control

- Keeping sender-to-receiver pipe “just full enough, but no fuller”: keep bottleneck link busy transmitting, but avoid high delays/buffering



$$\text{measured throughput} = \frac{\text{\# bytes sent in last RTT interval}}{\text{RTT}_{\text{measured}}}$$

Delay-based approach:

- RTT_{\min} - minimum observed RTT (uncongested path)
- uncongested throughput with congestion window cwnd is $\text{cwnd}/\text{RTT}_{\min}$

if measured throughput “very close” to uncongested throughput
increase cwnd linearly /* since path not congested */
else if measured throughput “far below” uncongested throughput
decrease cwnd linearly /* since path is congested */



TCP拥塞控制的改进

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

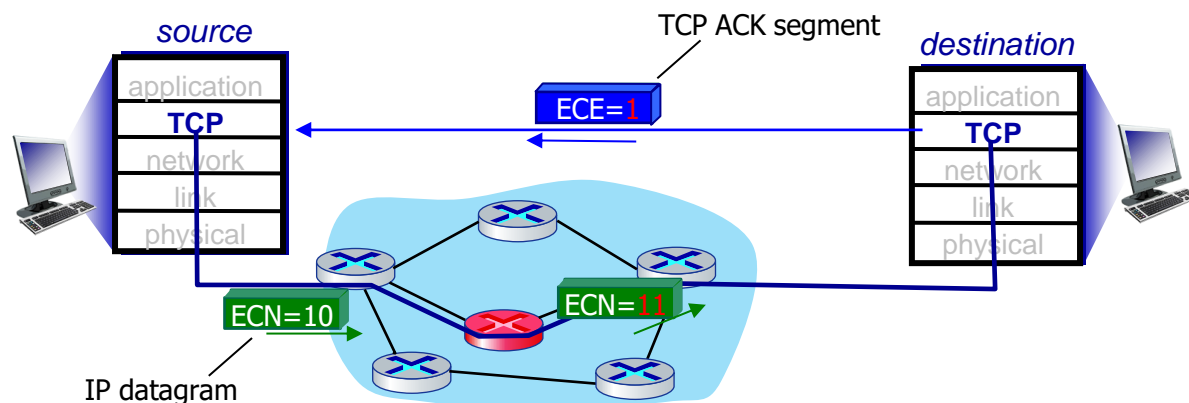
3.5 TCP协议

TCP拥塞控制改进

❖ Explicit congestion notification (ECN)

TCP deployments often implement *network-assisted* congestion control:

- two bits in IP header (ToS field) marked *by network router* to indicate congestion
 - *policy* to determine marking chosen by network operator
- congestion indication carried to destination
- destination sets ECE bit on ACK segment to notify sender of congestion
- involves both IP (IP header ECN bit marking) and TCP (TCP header C,E bit marking)



TCP的公平性

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

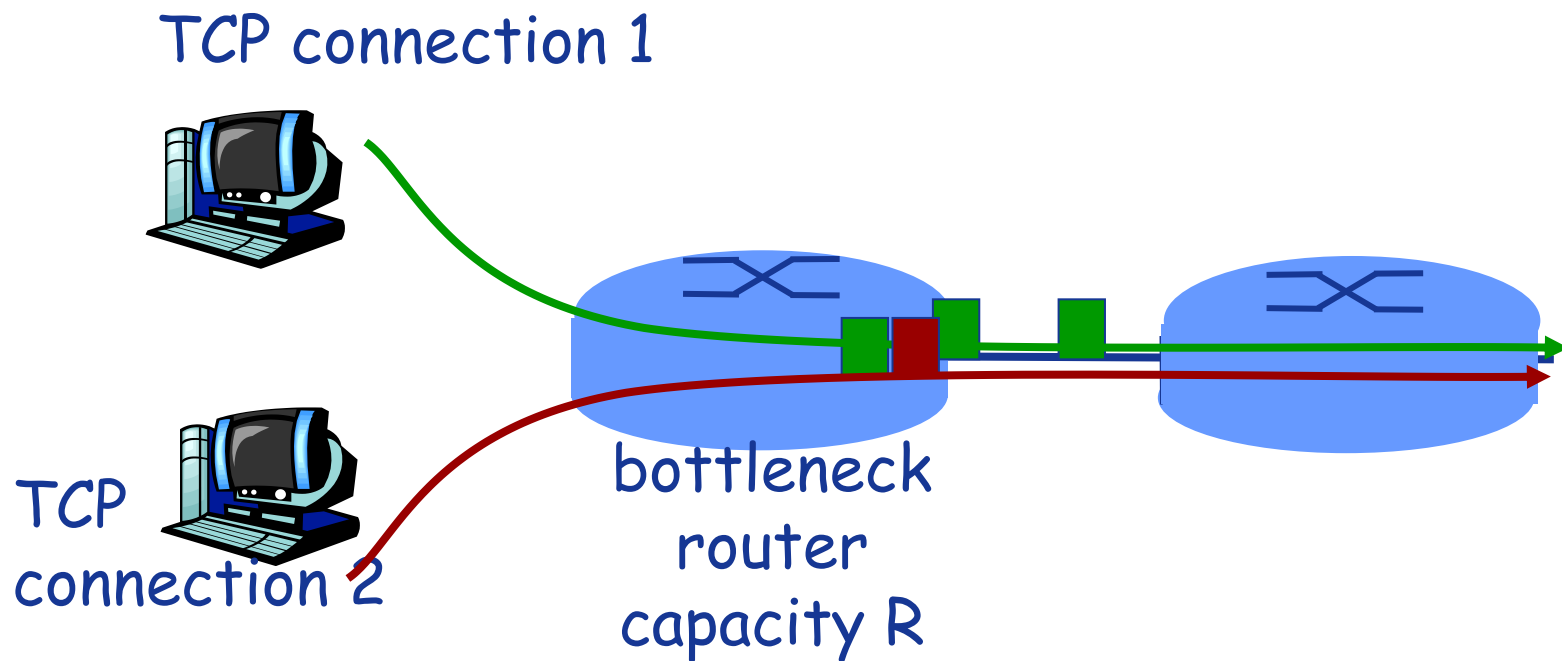
3.4 可靠数据传输原理

3.5 TCP协议

TCP性能分析

❖ 公平性？

- 如果K个TCP Session共享相同的瓶颈带宽R，那么每个Session的平均速率为 R/K





TCP具有公平性吗？

❖ 是的！

3.1 传输层服务

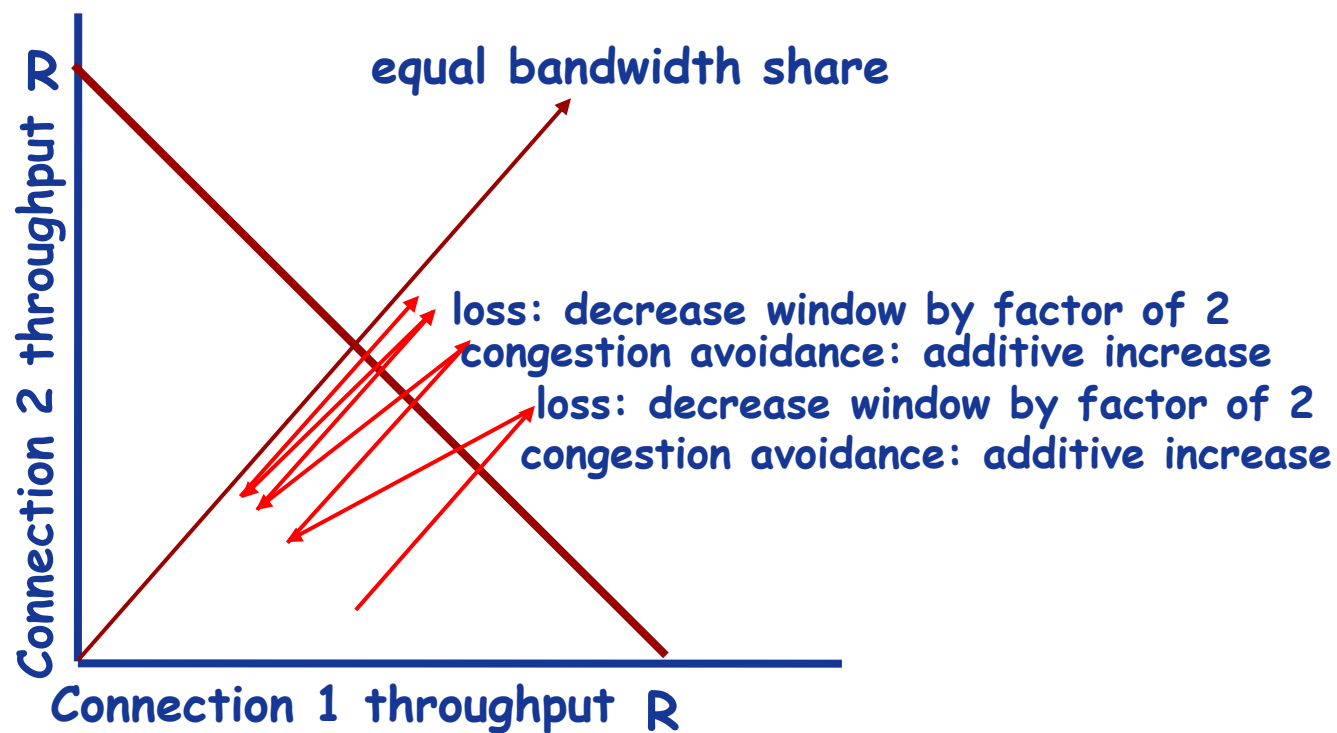
3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP性能分析





TCP的公平性

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP性能分析



❖ 公平性与UDP

- 多媒体应用通常不使用TCP，以免被拥塞控制机制限制速率
- 使用UDP：以恒定速率发送，能够容忍丢失
- 产生了不公平

❖ 研究：TCP friendly

❖ 公平性与并行TCP连接

- 某些应用会打开多个并行连接
- Web浏览器
- 产生公平性问题

❖ 例子：链路速率为 R ，已有9个连接

- 若新的应用请求建立1个TCP连接，则获得 $R/10$ 的速率
- 若新的应用请求建立11个TCP连接，则获得 $R/2$ 的速率



TCP的公平性

3.1 传输层服务

3.2 传输层多路复用/分解

3.3 UDP协议

3.4 可靠数据传输原理

3.5 TCP协议

TCP性能分析



❖ 公平性与UDP

- 多媒体应用通常不使用TCP，以免被拥塞控制机制限制速率
- 使用UDP：以恒定速率发送，能够容忍丢失
- 产生了不公平

❖ 研究：TCP friendly

❖ 公平性与并行TCP连接

- 某些应用会打开多个并行连接
- Web浏览器
- 产生公平性问题

❖ 例子：链路速率为 R ，已有9个连接

- 若新的应用请求建立1个TCP连接，则获得 $R/10$ 的速率
- 若新的应用请求建立11个TCP连接，则获得 $R/2$ 的速率



哈爾濱工業大學
HARBIN INSTITUTE OF TECHNOLOGY



立足航天，服务国防，面向国民经济主战场



谢谢!