



EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

ACH2026

Redes de Computadores

Capítulo 3 - Camada de Transporte

Profa. Dra. Cíntia B. Margi
Setembro/2009



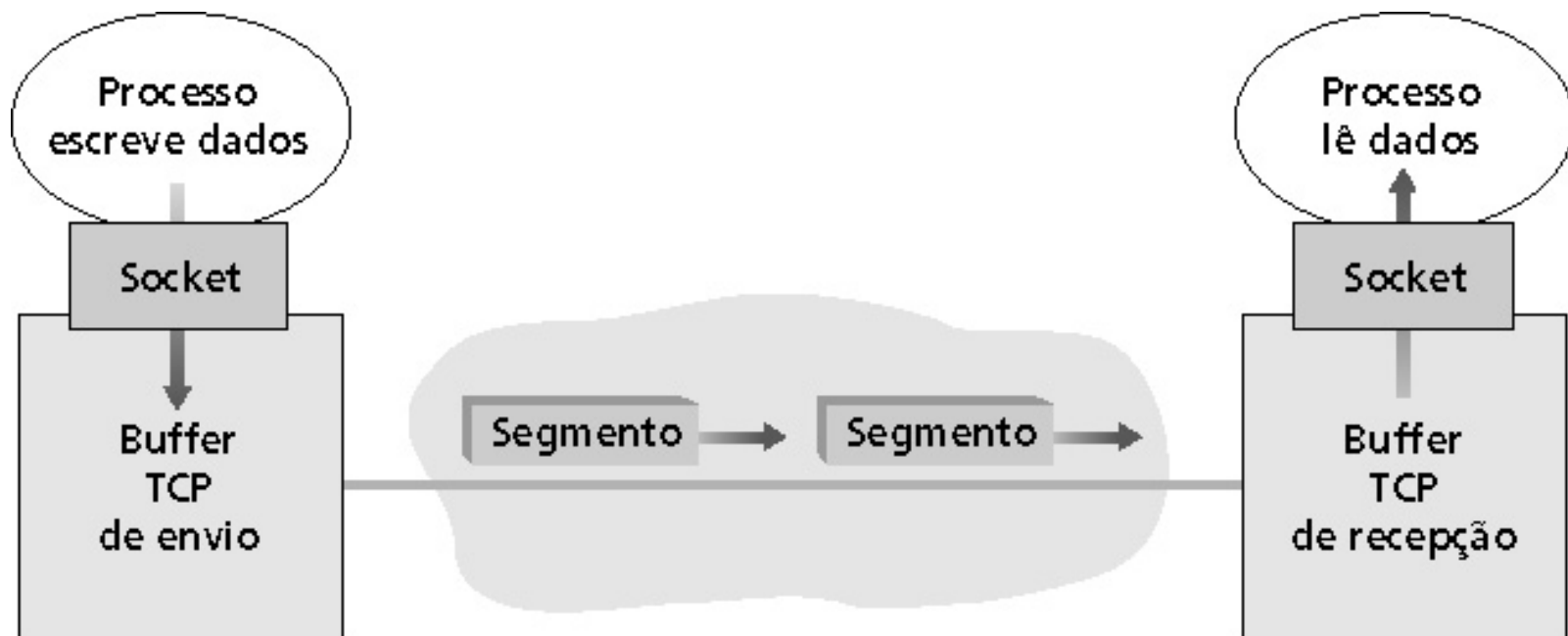
Transferência de Confiável de Dados

- Mecanismos:
 - Soma de verificação
 - Temporizador
 - Número de Seqüência
 - Mensagem de Reconhecimento
 - Mensagem de Reconhecimento Negativo
 - Janela



TCP

- Transmission Control Protocol.
- RFCs: 793, 1122, 1323, 2018, 2581.





TCP: Visão Geral

- ∇ Ponto-a-ponto:
 - Um transmissor, um receptor.
- ∇ Orientado à conexão:
 - apresentação: troca de mensagens de controle;
 - inicia o estado do transmissor e do receptor antes da troca de dados.
- ∇ Confiável;
- ∇ Full-duplex:
 - Transmissão bidirecional na mesma conexão.



TCP: Visão Geral

- ∇ Buffers de transmissão e de recepção:
 - Stream seqüencial de bytes;
 - MSS (maximum segment size): quantidade máxima de dados da camada de aplicação.
 - MTU (maximum transmission unit): tamanho máximo do quadro (camada de enlace) que remetente pode enviar.
 - Exemplos: 1460 bytes, 536 bytes.
 - Controle de fluxo:
 - Transmissor não esgota a capacidade do receptor.

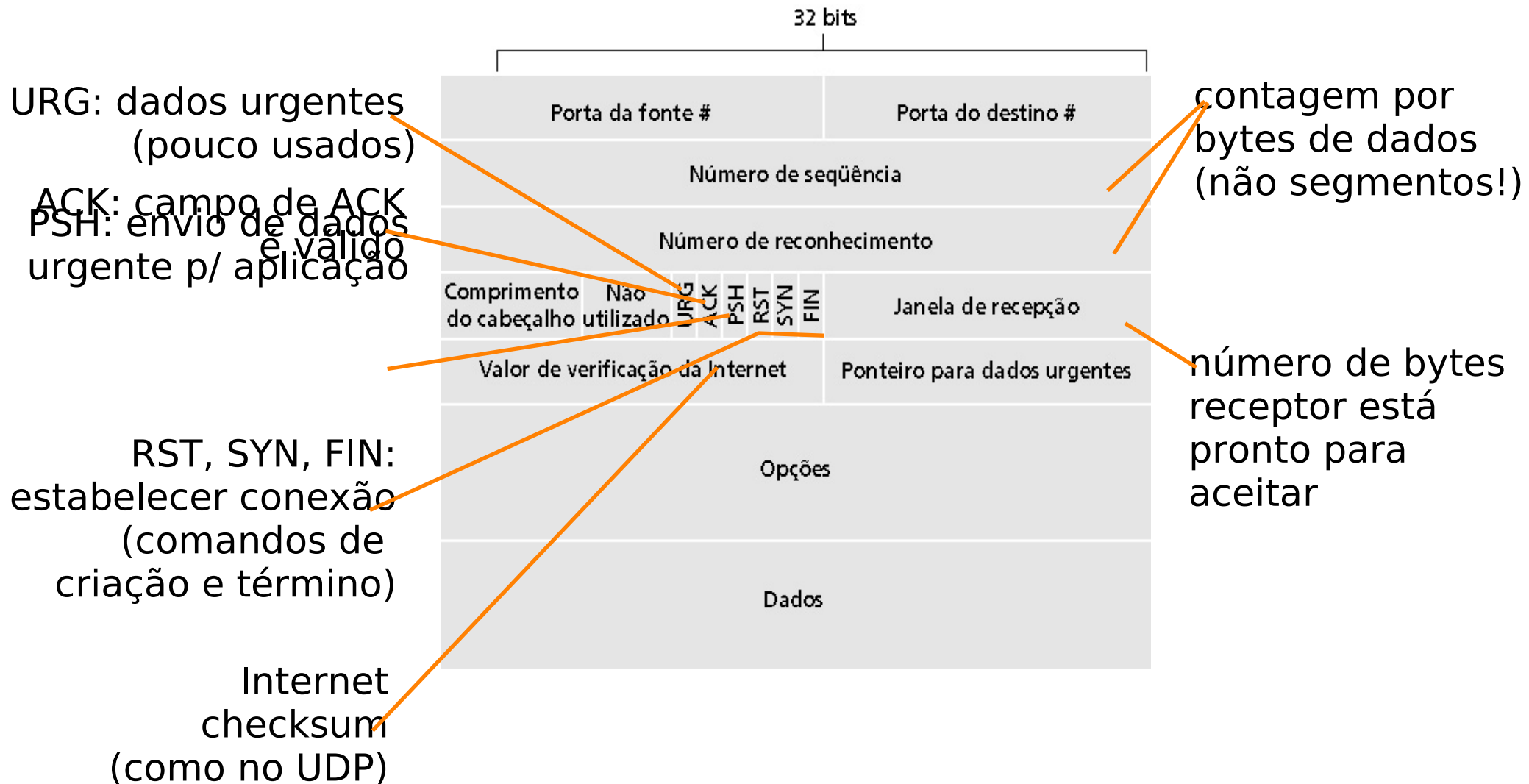


TCP: Visão Geral

- ∇ Suporte Paralelismo:
 - Múltiplos segmentos sem confirmação;
 - Controle de congestionamento e de fluxo definem tamanho da janela.
- ∇ Applet ilustrando controle de fluxo:
http://media.pearsoncmg.com/aw/aw_kurose_network_2/applets/flow/flowcontrol.html



Estrutura do Segmento TCP





Conexão TCP

- Transmissor estabelece conexão com o receptor antes de trocar segmentos de dados .
- Inicializar variáveis:
 - Números de sequência
 - Buffers, controle de fluxo (ex.: **RcvWindow**)
- ∇ **Cliente:** iniciador da conexão
`Socket clientSocket = new Socket("hostname", "port number");`
- ∇ **Servidor:** chamado pelo cliente
`Socket connectionSocket = welcomeSocket.accept();`

TCP: Estabelecimento de Conexão

Three way handshake:

1: sistema final cliente envia TCP SYN ao servidor

- especifica número de seqüência inicial.

2: sistema final servidor que recebe o SYN, responde com segmento SYNACK

- reconhece o SYN recebido;
- aloca buffers;
- especifica o número de seqüência inicial do servidor.

3: sistema final cliente reconhece o ACK.

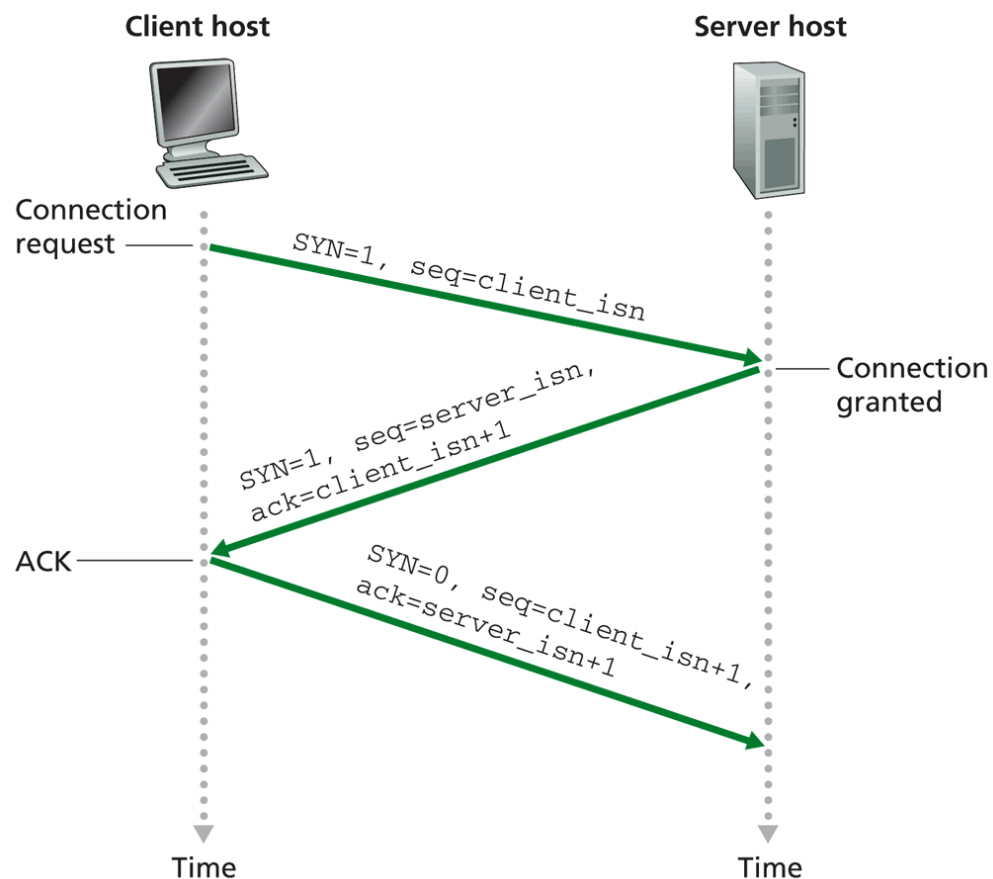


Figure 3.38 ♦ TCP three-way handshake: segment exchange

TCP: Fechando a Conexão

cliente fecha o socket:
clientSocket.close();

1: o cliente envia o segmento TCP FIN ao servidor._

2: servidor recebe FIN, responde com ACK. Fecha a conexão, envia FIN.

3: cliente recebe FIN, responde com ACK.

∇ Entra “espera temporizada”
- vai responder com ACK a FINs recebidos.

4: servidor, recebe ACK
Conexão fechada.

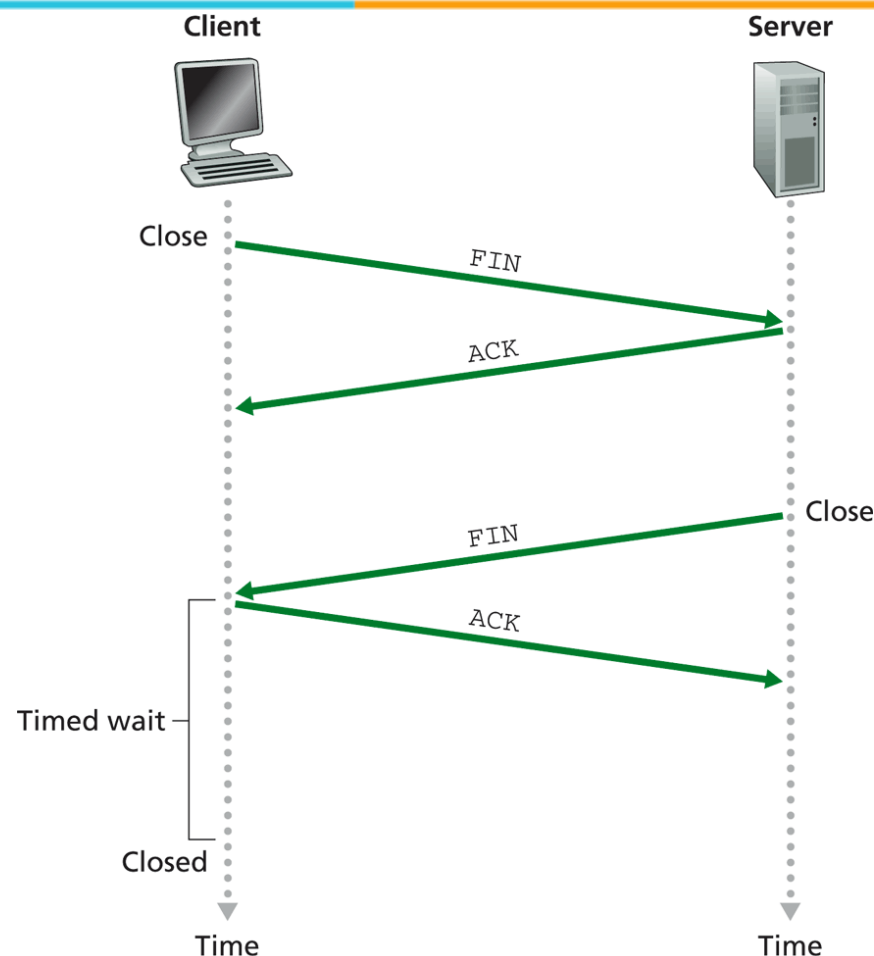
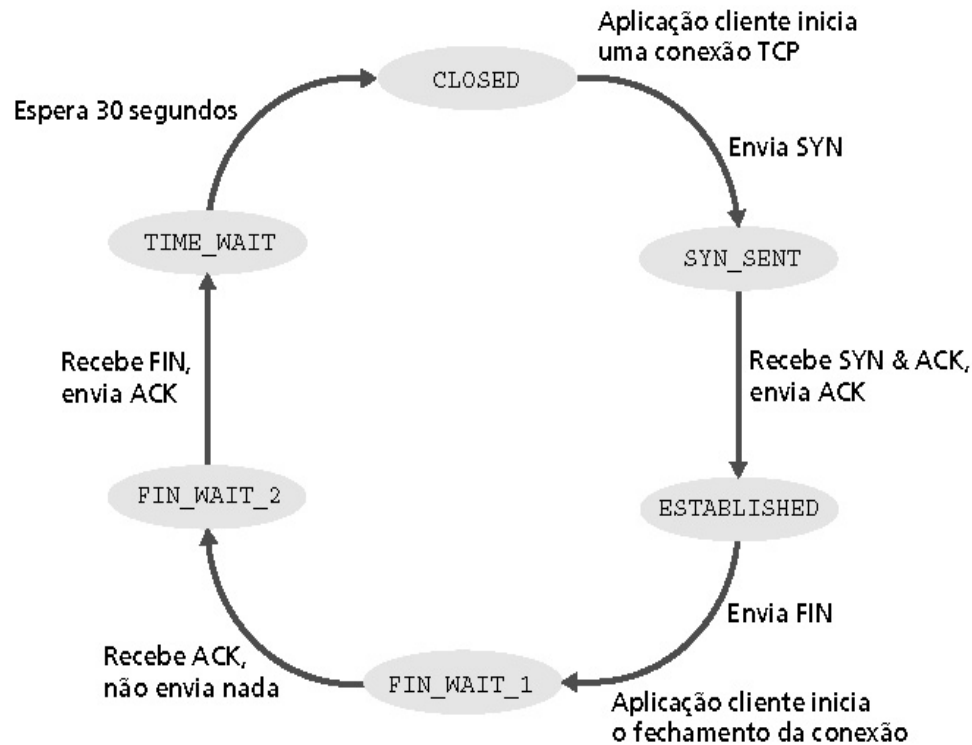


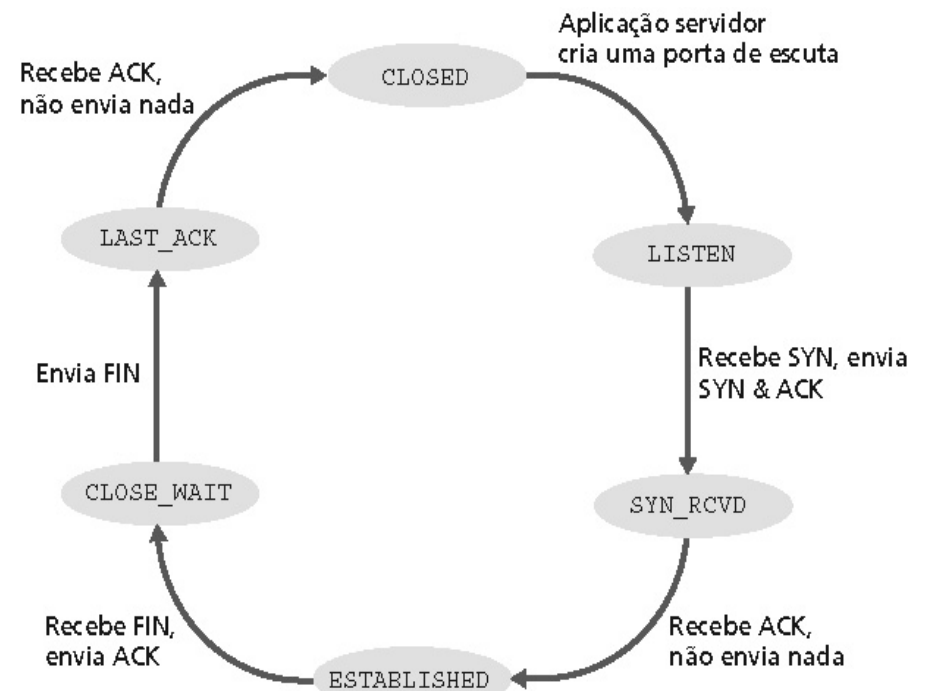
Figure 3.39 ♦ Closing a TCP connection



TCP: Máquina de Estados



Estados do cliente



Estados do servidor

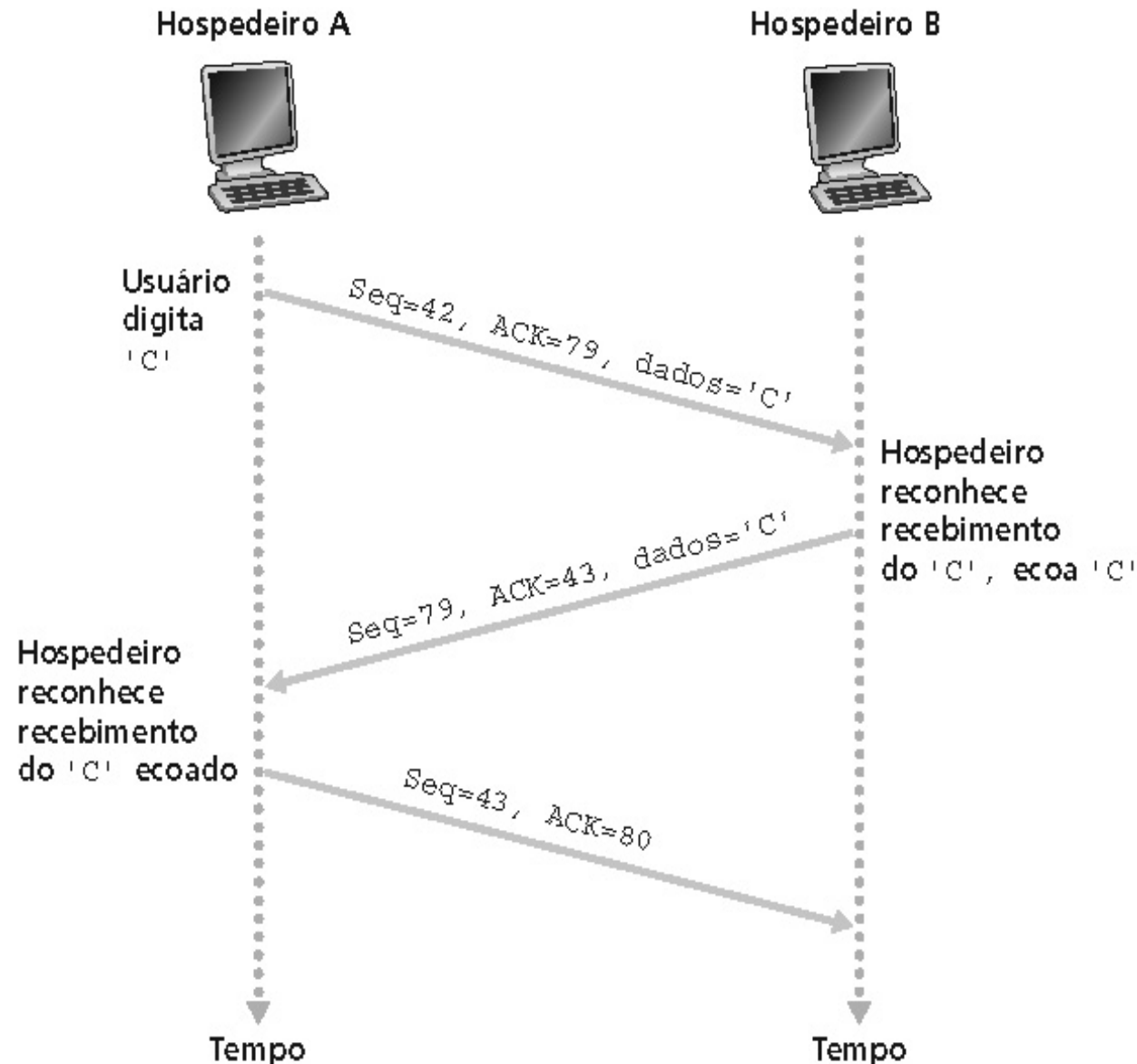


Números de seqüência e ACK

- ∇ Números de seqüência:
 - Número do primeiro byte nos segmentos de dados.
- ∇ ACKs:
 - Número do próximo byte esperado do outro lado;
 - ACK cumulativo.
- ∇ Como o receptor trata segmentos fora de ordem?
 - A especificação do TCP não define, fica a critério do implementador.



Estudo de Caso: Telnet





EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

Como determinar o valor dos temporizadores?

Como determinar o valor dos temporizadores?

- ∇ Maior que o RTT (Round Trip Time)!
 - Porém RTT varia!
- ∇ Muito curto: temporização prematura.
 - Retransmissões desnecessárias.
- ∇ Muito longo: a reação à perda de segmento fica lenta.
- ∇ Como estimar o RTT?

Amostras de RTT

- ∇ **SampleRTT**: tempo medido da transmissão de um segmento até a respectiva confirmação.
 - Ignora retransmissões e segmentos reconhecidos de forma cumulativa.
- ∇ SampleRTT varia de forma rápida, é desejável um amortecedor para a estimativa do RTT.
 - Usar várias medidas recentes, não apenas o último **SampleRTT** obtido.



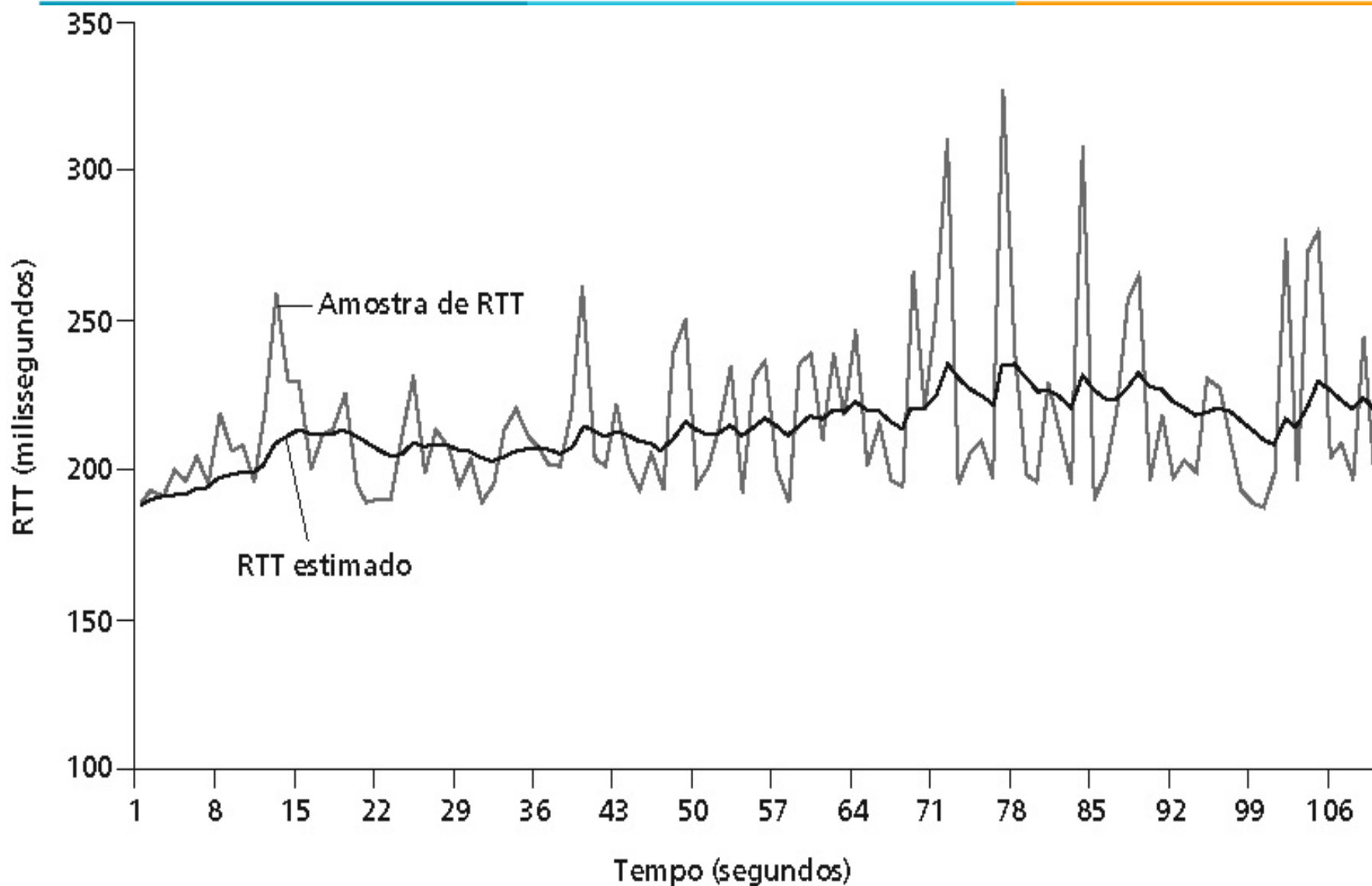
Estimando o RTT

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

- ∇ Média móvel com peso exponencial
- ∇ Influência de uma dada amostra decresce de forma exponencial
- ∇ Valor típico: $\alpha = 0,125$



Exemplo de Estimativa do RTT



Definindo a temporização

- ∇ **EstimatedRTT** mais “margem de segurança”.
- ∇ Primeiro estimar o quanto o SampleRTT se desvia do EstimatedRTT:

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

(tipicamente $\beta = 0.25$)



Definindo a temporização

∇ Então ajustar o intervalo de
temporização:

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



EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

TCP: Transferência de Dados Confiável



TCP: Transferência de Dados Confiável

∀ TCP cria serviços de rdt em cima do serviço não-confiável do IP.

- segmentos em paralelo;
- ACKs cumulativos;
- TCP usa tempo de retransmissão simples.

∀ Retransmissões são disparadas por:

- eventos de tempo de confirmação;
- ACKs duplicados.



Eventos do Transmissor TCP

- ∇ Considere um transmissor TCP simples:
 - ignore ACKs duplicados, controle de fluxo e controle de congestionamento.
- ∇ Dado recebido da aplicação:
 - cria segmento com número de sequência;
 - Número de sequência é o número do byte-stream do 1º byte de dados no segmento;
 - inicia o temporizador se ele ainda não estiver em execução (considera o segmento não-confirmado mais antigo);
 - tempo de expiração: TimeoutInterval.



Eventos do Transmissor TCP

- ∇ Tempo de confirmação:
 - retransmite o segmento que provocou o tempo de confirmação;
 - reinicia o temporizador.
- ∇ ACK recebido:
 - Quando houver o ACK de segmentos anteriormente não confirmados:
 - atualizar o que foi confirmado;
 - iniciar o temporizador se houver segmentos pendentes.



```
/* Assume sender is not constrained by TCP flow or congestion control, that data from above is less than MSS in size, and that data transfer is in one direction only. */
```

```
NextSeqNum=InitialSeqNumber
```

```
SendBase=InitialSeqNumber
```

```
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)  
        break;
```

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

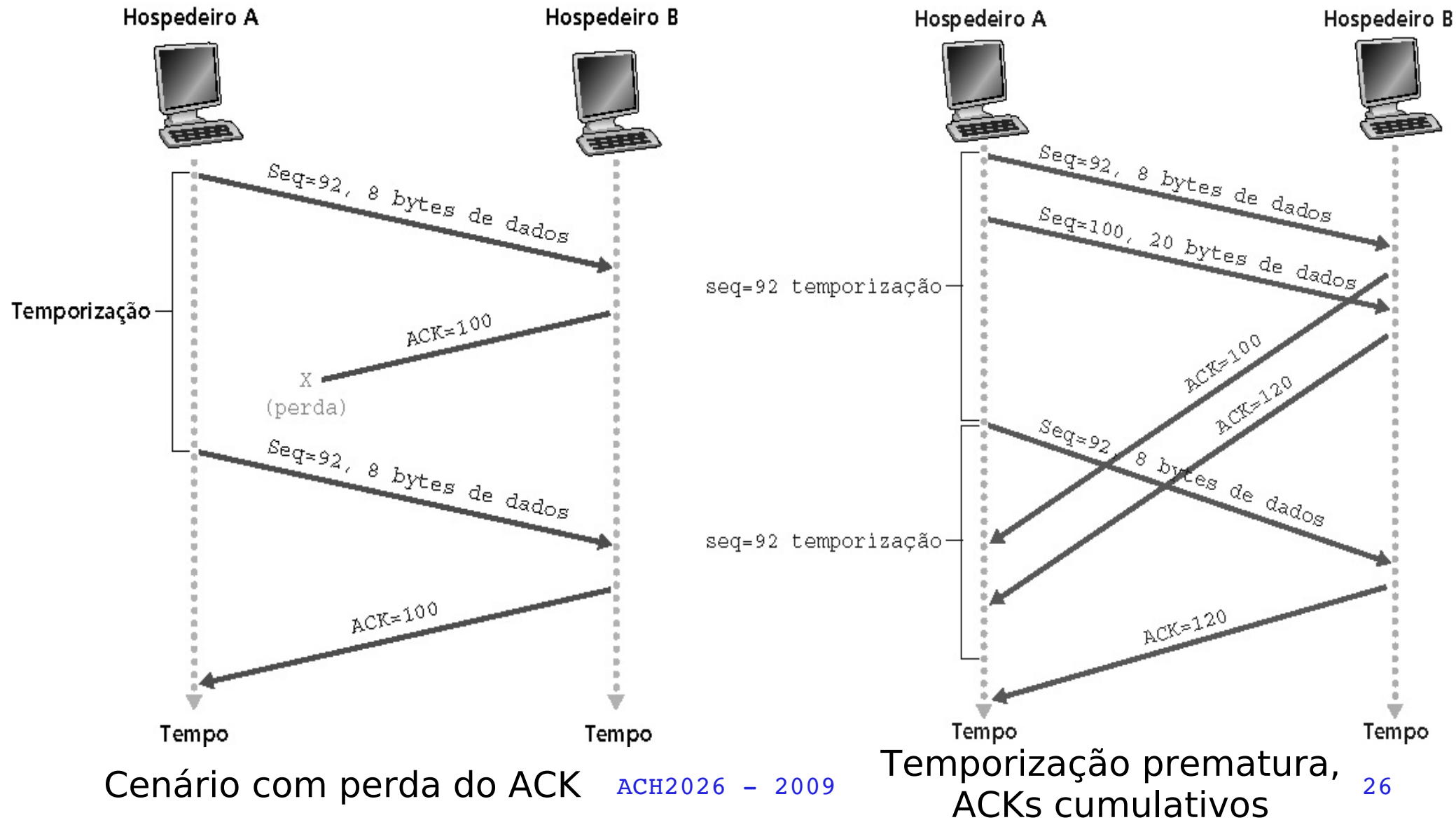
```
    event: ACK received, with ACK field value of y  
        if (y > SendBase) {  
            SendBase=y  
            if (there are currently any not-yet-acknowledged segments)  
                start timer  
        }  
        break;
```

```
} /* end of loop forever */
```

Figure 3.33 ♦ Simplified TCP sender

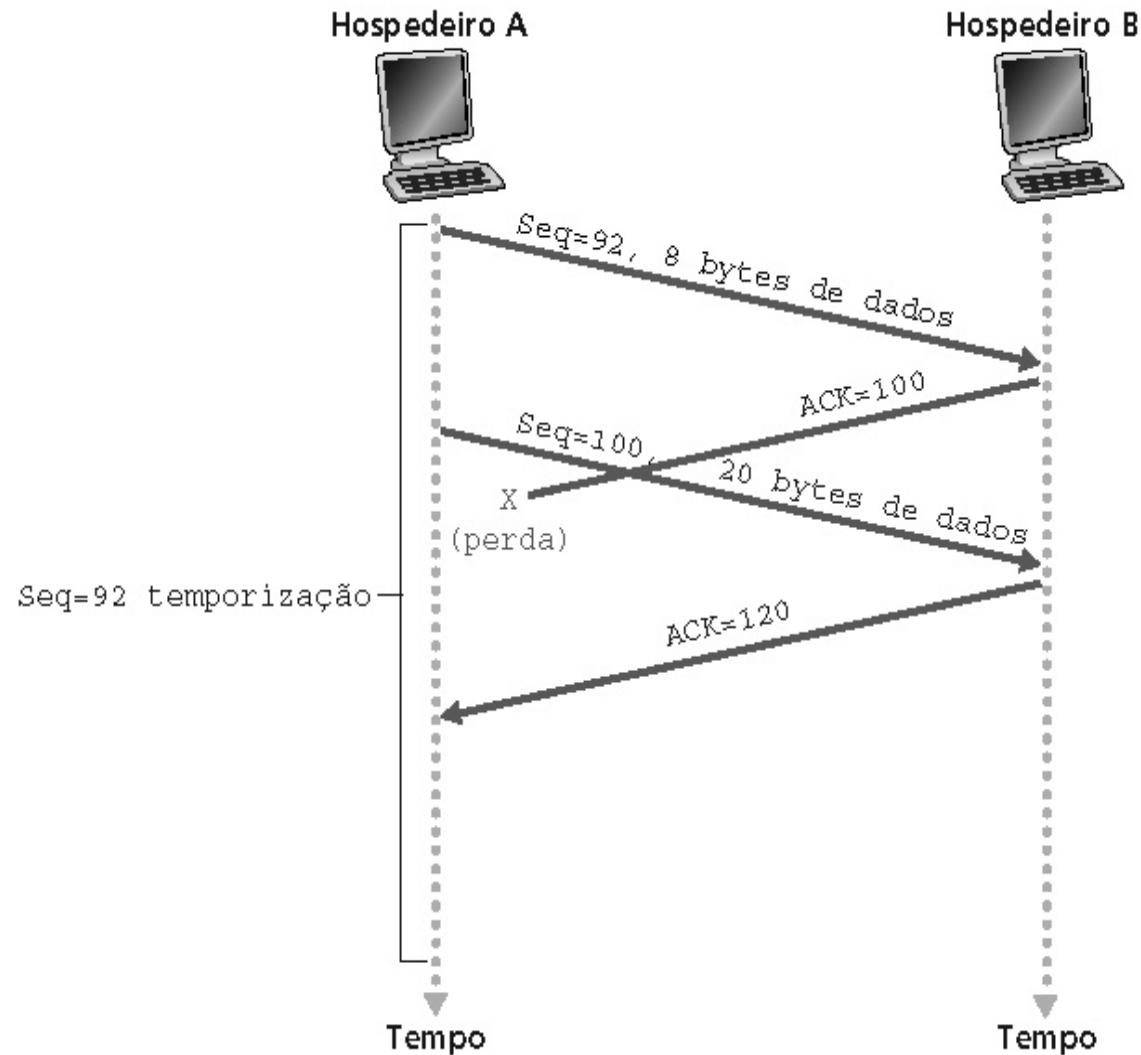


TCP: Retransmissão





TCP: Retransmissão



Cenário de ACK cumulativo



Geração de ACK [RFC 1122,2581]

Evento no receptor	Ação do receptor TCP
Segmento chega em ordem, não há lacunas, segmentos anteriores já aceitos	ACK retardado. Espera até 500 ms pelo próximo segmento. Se não chegar, envia ACK.
Segmento chega em ordem, não há lacunas, um ACK atrasado pendente.	Imediatamente envia um ACK cumulativo.
Segmento chega fora de ordem, número de seqüência chegou maior: gap detectado.	Envia ACK duplicado, indicando número de seqüência do próximo byte esperado.
Chegada de segmento que parcial ou completamente preenche o gap.	Reconhece imediatamente se o segmento começa na borda inferior do gap.



Retransmissão Rápida

- Com frequência, o tempo de expiração é relativamente longo...
- Detecta segmentos perdidos por meio de ACKs duplicados:
 - transmissor frequentemente envia muitos segmentos *back-to-back*, e se um segmento é perdido, haverá muitos ACKs duplicados.
- Se o transmissor recebe 3 ACKs para o mesmo dado, supõe que o segmento após foi perdido:
 - Retransmissão rápida: reenvia o segmento antes de o temporizador expirar.



EACH

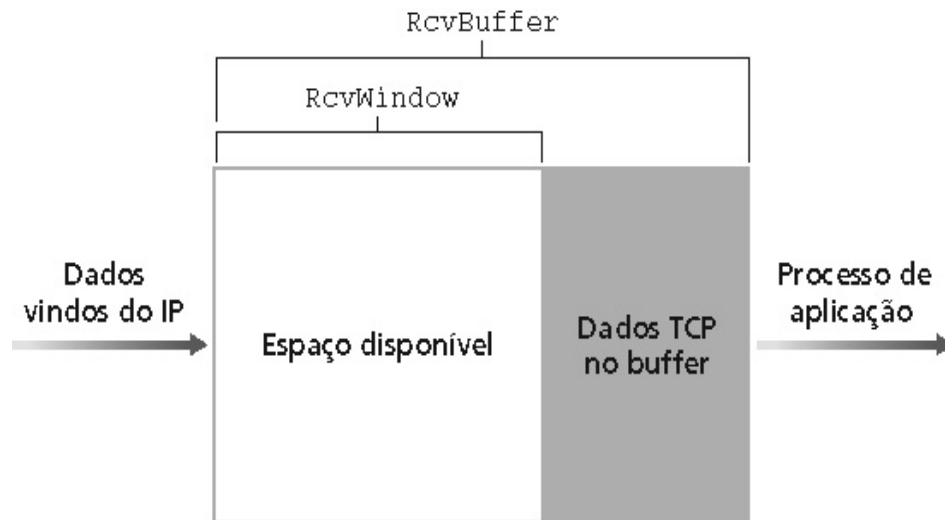
Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

TCP: Controle de fluxo



Controle de Fluxo do TCP

- Lado receptor da conexão TCP possui um buffer de recepção:



- Processos de aplicação podem ser lentos para ler o buffer

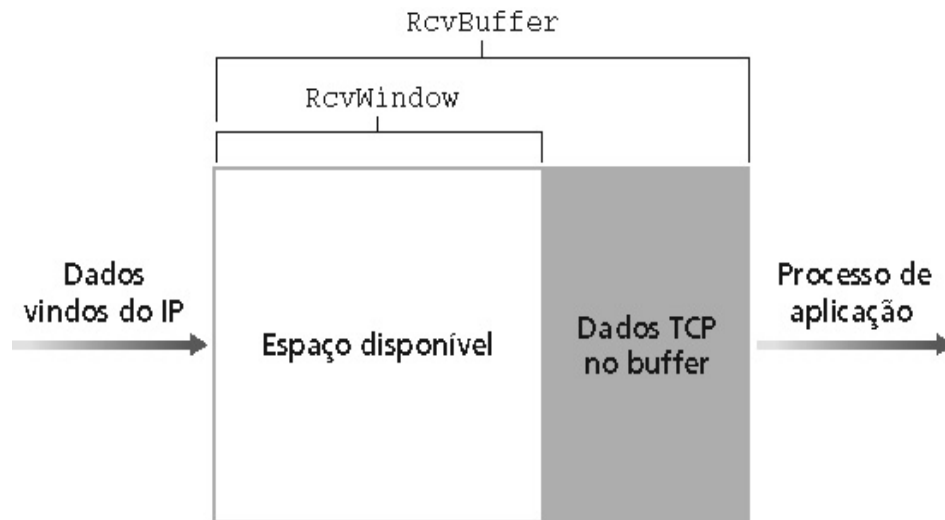
Controle de fluxo

Transmissor não deve esgotar os buffers de recepção enviando dados rápido demais.

- Serviço de **speed-matching**: encontra a taxa de envio adequada à taxa de vazão da aplicação receptora.



Controle de Fluxo do TCP: Como funciona?



- Receptor informa a área disponível incluindo valor **RcvWindow** nos segmentos
- Transmissor limita os dados não confinados ao **RcvWindow**
 - Garantia contra overflow no buffer do receptor

(suponha que o receptor TCP descarte segmentos fora de ordem)

- Espaço disponível no buffer
- = **RcvWindow**
= **RcvBuffer - [LastByteRcvd - LastByteRead]**



EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

TCP: Controle de congestionamento



Congestionamento

- ∇ Informalmente: “muitas fontes enviando dados acima da capacidade da **rede** de tratá-los”.
- ∇ Diferente de controle de fluxo!
- ∇ Sintomas:
 - perda de pacotes (saturação de buffer nos roteadores);
 - atrasos grandes (filas nos buffers dos roteadores).



Mecanismos de Controle de Congestionamento

Fim-a-fim:

- Não usa realimentação explícita da rede.
- Congestionamento é inferido a partir das perdas e dos atrasos observados nos sistemas finais.
- Abordagem usada pelo TCP.



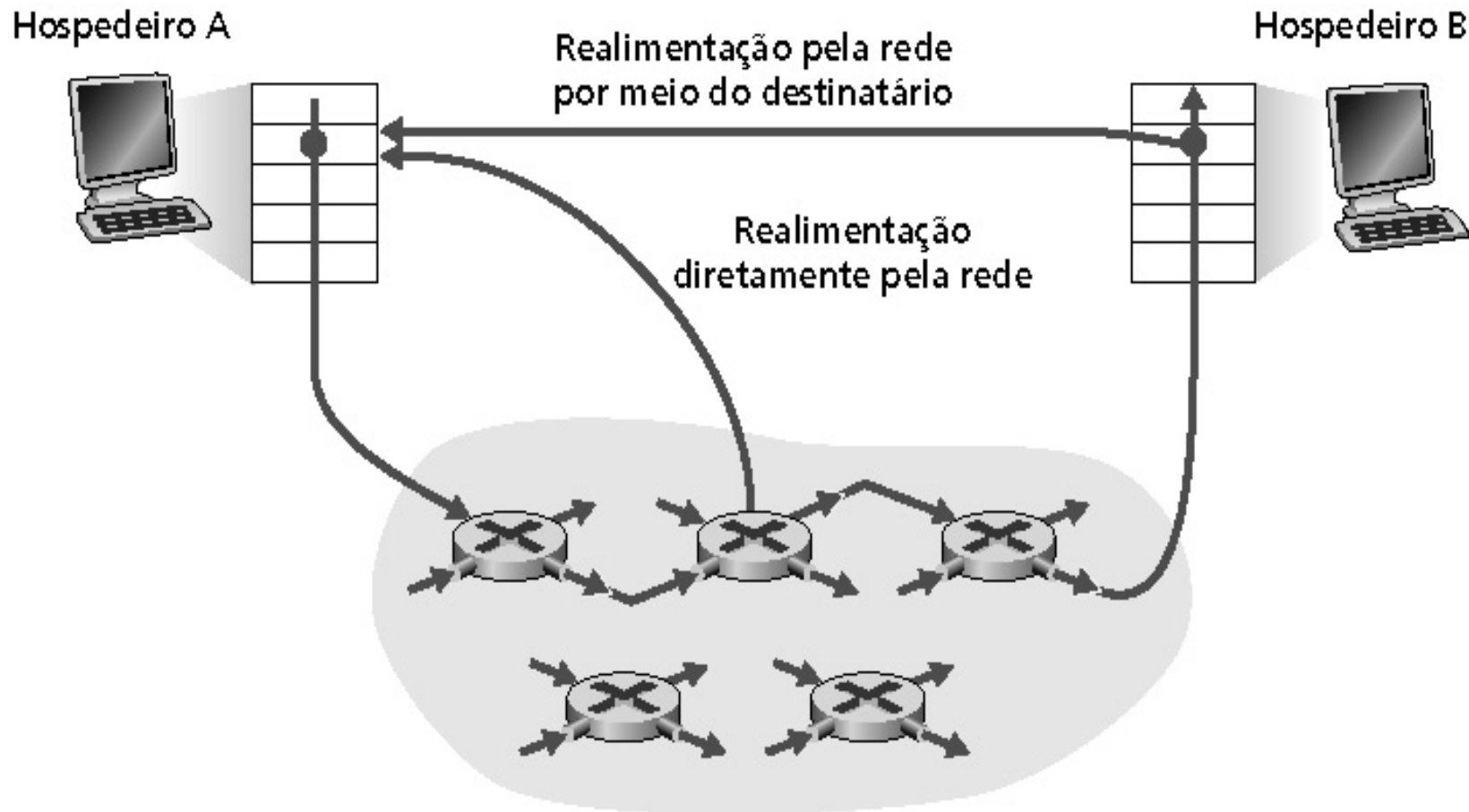
Mecanismos de Controle de Congestionamento

Assistido pela rede:

- Roteadores enviam informações para os sistemas finais.
- Bit único indicando o congestionamento (SNA, DECbit, TCP/IP ECN, ATM).
- Taxa explícita do transmissor poderia ser enviada.



Mecanismos de Controle de Congestionamento





EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

TCP: Controle de Congestionamento



TCP: Controle de Congestionamento

- ∇ Controle fim-a-fim (sem assistência da rede).
- ∇ Transmissor limita a transmissão:
 $\text{LastByteSent} - \text{LastByteAcked} \leq \text{CongWin}$
- ∇ Aproximadamente,
 $\text{rate} = \text{CongWin} / \text{RTT Bytes/seg}$
- ∇ **CongWin** é dinâmico, função de congestionamento das redes detectadas.



Como o transmissor detecta congestionamento?

- ∇ Através da perda de segmentos!
- ∇ A perda é detectada por:
 - tempo de confirmação;
 - *ou* 3 ACKs duplicados.
- ∇ Consequentemente, transmissor TCP reduz a taxa (**CongWin**) após o evento de perda.



TCP: Controle de Congestionamento

- ∇ Algoritmo possui três mecanismos:
- Aumento Aditivo, Diminuição Multiplicativa (AIMD);
 - Partida lenta;
 - Reação a eventos de esgotamento de temporização.



TCP: AIMD

-
- Gráfico de Janela de congestionamento vs Tempo. O eixo Y é rotulado 'Janela de congestionamento' com marcas em 8K, 16K e 24K. O eixo X é rotulado 'Tempo'. A linha mostra um padrão de crescimento linear seguido por quedas abruptas, formando uma série de picos decrescentes.



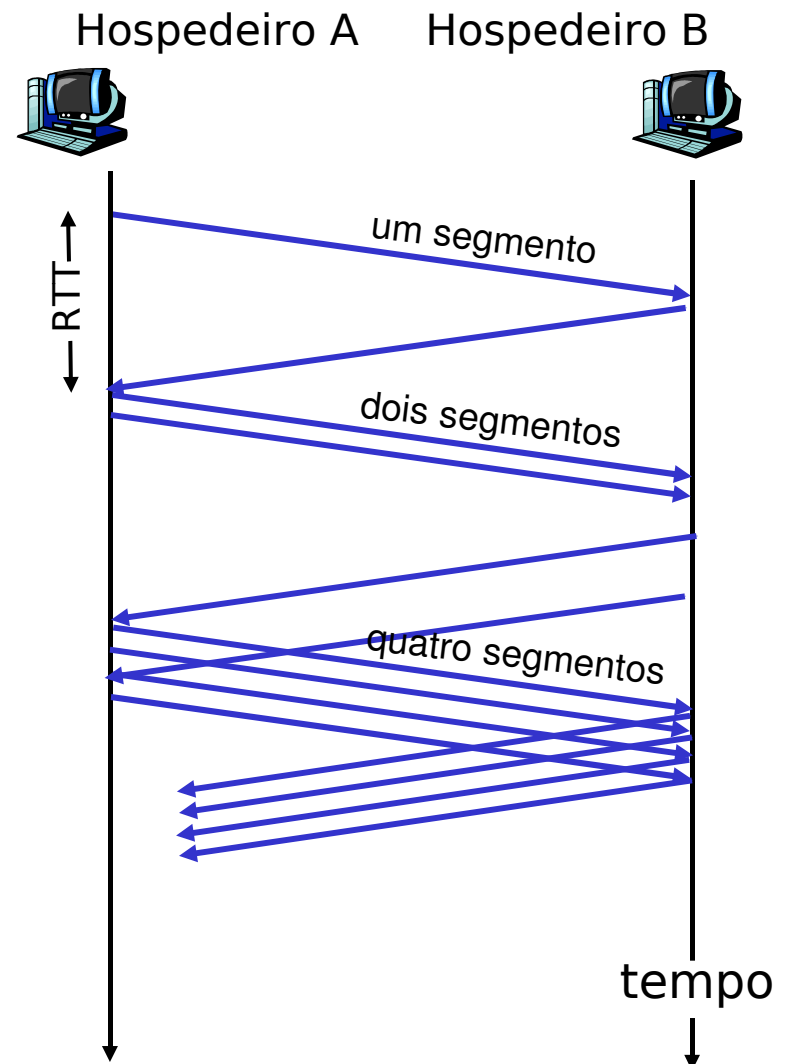
TCP: Partida Lenta

- ∇ Quando a conexão começa, **CongWin** = 1 MSS
 - Exemplo: MSS = 500 bytes e RTT = 200 ms.
 - Taxa inicial = 20 kbps.
- ∇ Largura de banda disponível pode ser $\gg \text{MSS/RTT}$
 - Desejável aumentar rapidamente até a taxa respeitável.

TCP: Partida Lenta

- Quando a conexão começa, a taxa aumenta de modo exponencial até a ocorrência do primeiro evento de perda.
- Dobra o **CongWin** a cada RTT.
- Faz-se incrementando o **CongWin** para cada ACK recebido.

Resumindo: taxa inicial é lenta mas aumenta de modo exponencialmente rápido.





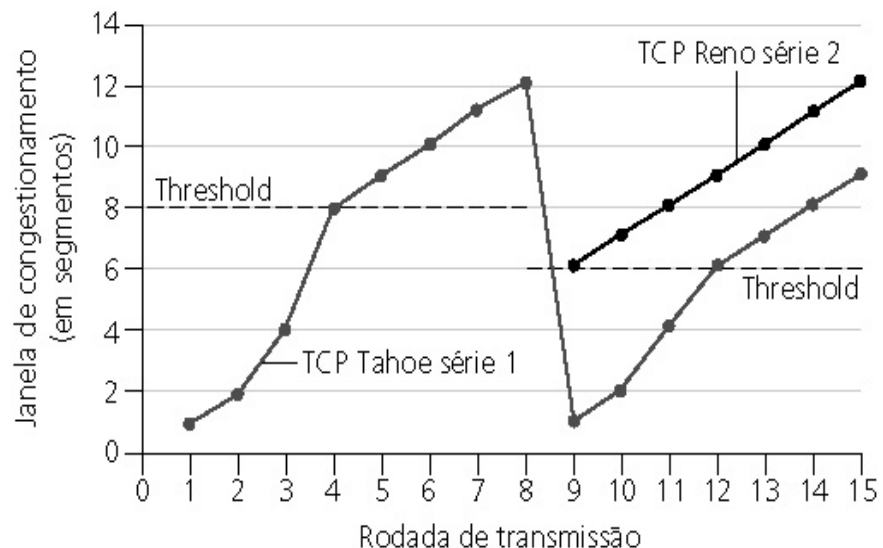
TCP Reno: Reação a *timeouts*

- ∇ Após 3 ACKs duplicados:
 - **CongWin** é cortado pela metade;
 - janela então cresce linearmente.
- ∇ Após *timeout* de confirmação:
 - **CongWin** é ajustado para 1 MSS;
 - janela então cresce exponencialmente até um limite, então cresce linearmente.
- **Filosofia:**
 - se ocorreu *timeout*, dados não estão sendo entregues!



TCP: Reação a *timeouts*

- ▽ Quando o aumento exponencial deve tornar-se linear?
- ▽ Quando **CongWin** obtiver 1/2 do seu valor antes do tempo de confirmação.
- ▽ Implementação:
 - No evento de perda, o limiar (variável) é ajustado para 1/2 do CongWin logo antes do evento de perda.





TCP: Controle de Congestionamento

- ∇ Quando **CongWin** está abaixo do limite (**Threshold**), o transmissor em fase de slow-start, a janela cresce exponencialmente.
- ∇ Quando **CongWin** está acima do limite (**Threshold**), o transmissor em fase de congestion-avoidance, a janela cresce linearmente.
- ∇ Quando ocorrem três ACK duplicados, o limiar (**Threshold**) é ajustado em **CongWin/2** e **CongWin** é ajustado para **Threshold**.
- ∇ Quando ocorre tempo de confirmação, o **Threshold** é ajustado para **CongWin/2** e o **CongWin** é ajustado para 1 MSS.

Resumindo...

Evento	Estado	Ação do transmissor TCP	Comentário
ACK recebido para dado previamente não confirmado	partida lenta (SS)	$\text{CongWin} = \text{CongWin} + \text{MSS}$, If ($\text{CongWin} > \text{Threshold}$) ajusta estado para “prevenção de congestionamento”	Resulta em dobrar o CongWin a cada RTT
ACK recebido para dado previamente não confirmado	prevenção de congestionamento (CA)	$\text{CongWin} = \text{CongWin} + \text{MSS} * (\text{MSS}/\text{CongWin})$	Aumento aditivo, resulta no aumento do CongWin em 1 MSS a cada RTT
Evento de perda detectado por três ACKs duplicados	SS or CA	$\text{Threshold} = \text{CongWin}/2$, $\text{CongWin} = \text{Threshold}$, ajusta estado para “prevenção de congestionamento”	Recuperação rápida, implementando redução multiplicativa o CongWin não cairá abaixo de 1 MSS.
Tempo de confirmação	SS or CA	$\text{Threshold} = \text{CongWin}/2$, $\text{CongWin} = 1 \text{ MSS}$, ajusta estado para “partida lenta”	Entra em partida lenta
ACK duplicado	SS or CA	Incrementa o contador de ACK duplicado para o segmento que está sendo confirmado	CongWin e Threshold não mudam

Vazão do TCP (simples)

- ∇ Como obter a vazão (*throughput*) média do TCP como função do tamanho da janela e do RTT?
 - Ignore a partida lenta.
- ∇ Seja W o tamanho da janela quando ocorre perda.
- ∇ Quando a janela é **W** , a vazão é W/RTT .
- ∇ Após a perda, a janela cai para $W/2$, e a vazão para $W/2\text{RTT}$.

Vazão média: $0,75 W/\text{RTT}$

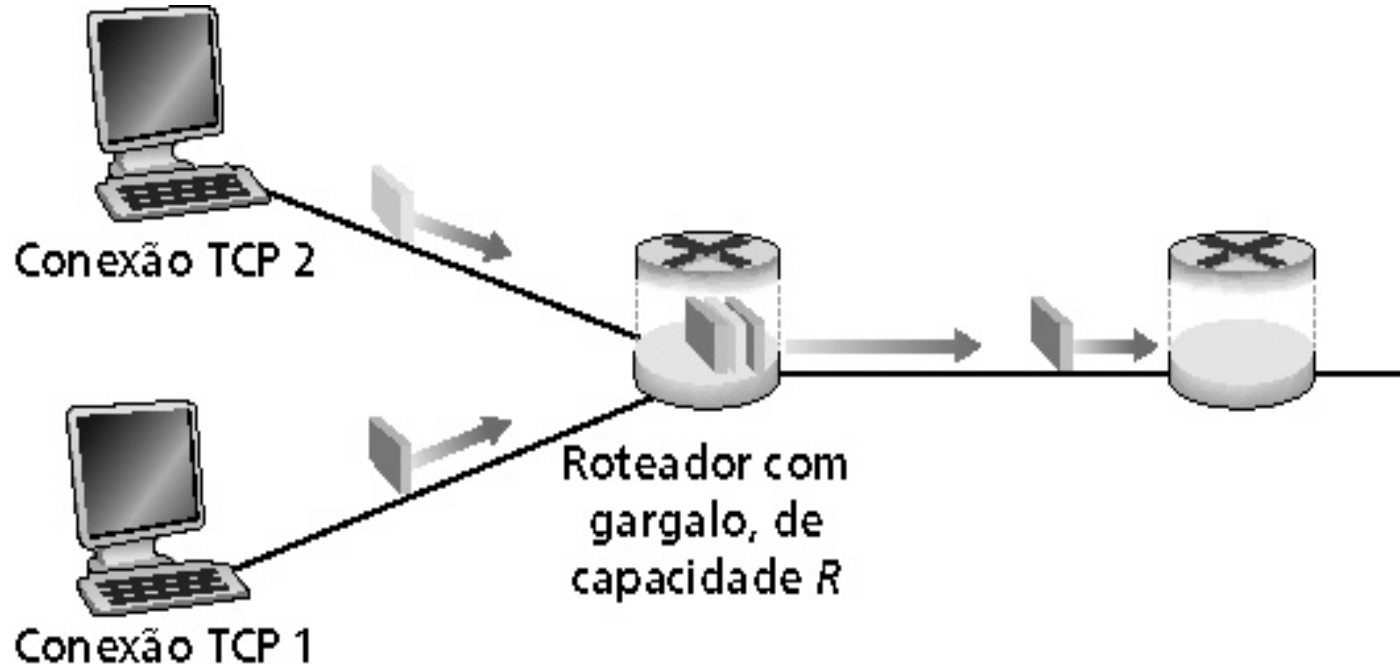
TCP do “futuro”

- ∇ Exemplo: segmento de 1500 bytes, RTT de 100 ms.
- ∇ Para obter *throughput* de 10 Gbps,
 - janela $W = 83.333$ segmentos!
- ∇ Throughput em termos da taxa de perda:
$$\frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$
- ∇ $L = 2 \cdot 10^{-10}$, onde L é a probabilidade de perda. Muito baixa! **Irreal!**
- ∇ São necessárias novas versões de TCP para alta velocidade!



Justiça (Eqüidade) do TCP

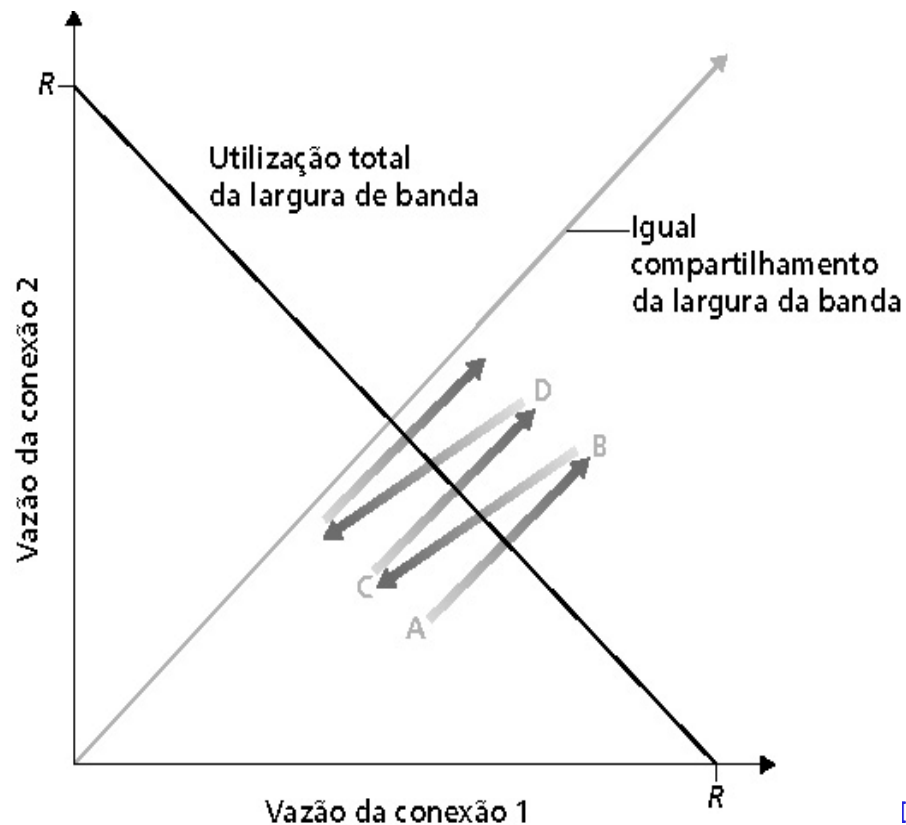
- Objetivo de eqüidade: se K sessões TCP compartilham o mesmo enlace do gargalo com largura de banda R , cada uma deve ter taxa média de R/K .





O TCP é justo?

- ∇ Duas sessões competindo pela banda:
- O aumento aditivo fornece uma inclinação de 1, quando a vazão aumenta.
 - Redução multiplicativa diminui a vazão proporcionalmente.



- perda: reduz janela por um fator de 2
- prevenção de congestionamento: aumento aditivo
- perda: reduz janela por um fator de 2
- prevenção de congestionamento: aumento aditivo
- Portanto, AIMD é justo!



Justiça e UDP

- ∇ Aplicações multimídia normalmente não usam TCP!
 - Não querem a taxa estrangulada pelo controle de congestionamento.
- ∇ Em vez disso, usam UDP:
 - trafega áudio/vídeo a taxas constantes, toleram perda de pacotes
- ∇ Área de pesquisa: TCP amigável.



Justiça e Conexões TCP Paralelas

- ∇ Nada previne as aplicações de abrirem conexões paralelas entre 2 hospedeiros!
- ∇ Web browsers fazem isso.
- ∇ Exemplo: enlace de taxa R suportando 9 conexões:
 - novas aplicações pedem 1 TCP, obtém taxa de $R/10$.
 - novas aplicações pedem 11 TCPs, obtém $R/2$!



Modelagem do atraso TCP

- ∇ Quanto tempo demora para receber um objeto de um servidor Web após enviar um pedido? Ou seja, qual a latência?
- ∇ Ignorando o congestionamento, o atraso é influenciado por:
 - estabelecimento de conexão TCP;
 - atraso de transferência de dados;
 - partida lenta.



Modelagem do atraso

TCP - Hipóteses

- ∇ Suponha um enlace entre o cliente e o servidor com taxa de dados R .
- ∇ Seja:
 - S : MSS (bits);
 - O : tamanho do objeto (bits).
- ∇ Não há retransmissões (sem perdas e corrupção de dados).
- ∇ Janela de congestionamento (W):
 - estática;
 - dinâmica.

Janela de Congestionamento Estática

Primeiro caso:

- $WS/R > RTT + S/R$
- ou seja, o ACK para o primeiro segmento na primeira janela retorna antes de enviar todos os dados.

$$\text{atraso} = 2RTT + O/R$$

ACH2

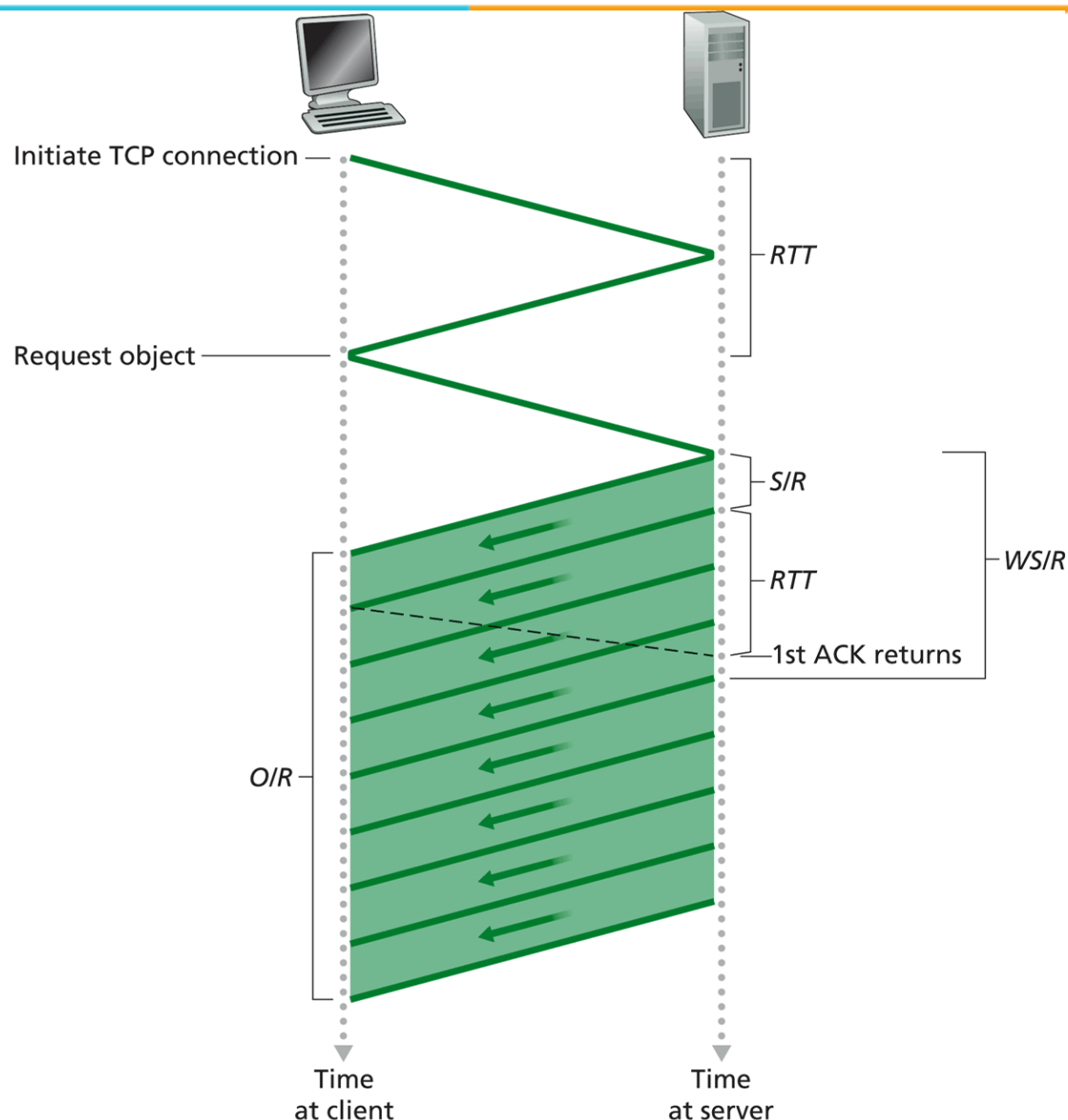


Figure 3.55 ♦ The case $WS/R > RTT + S/R$



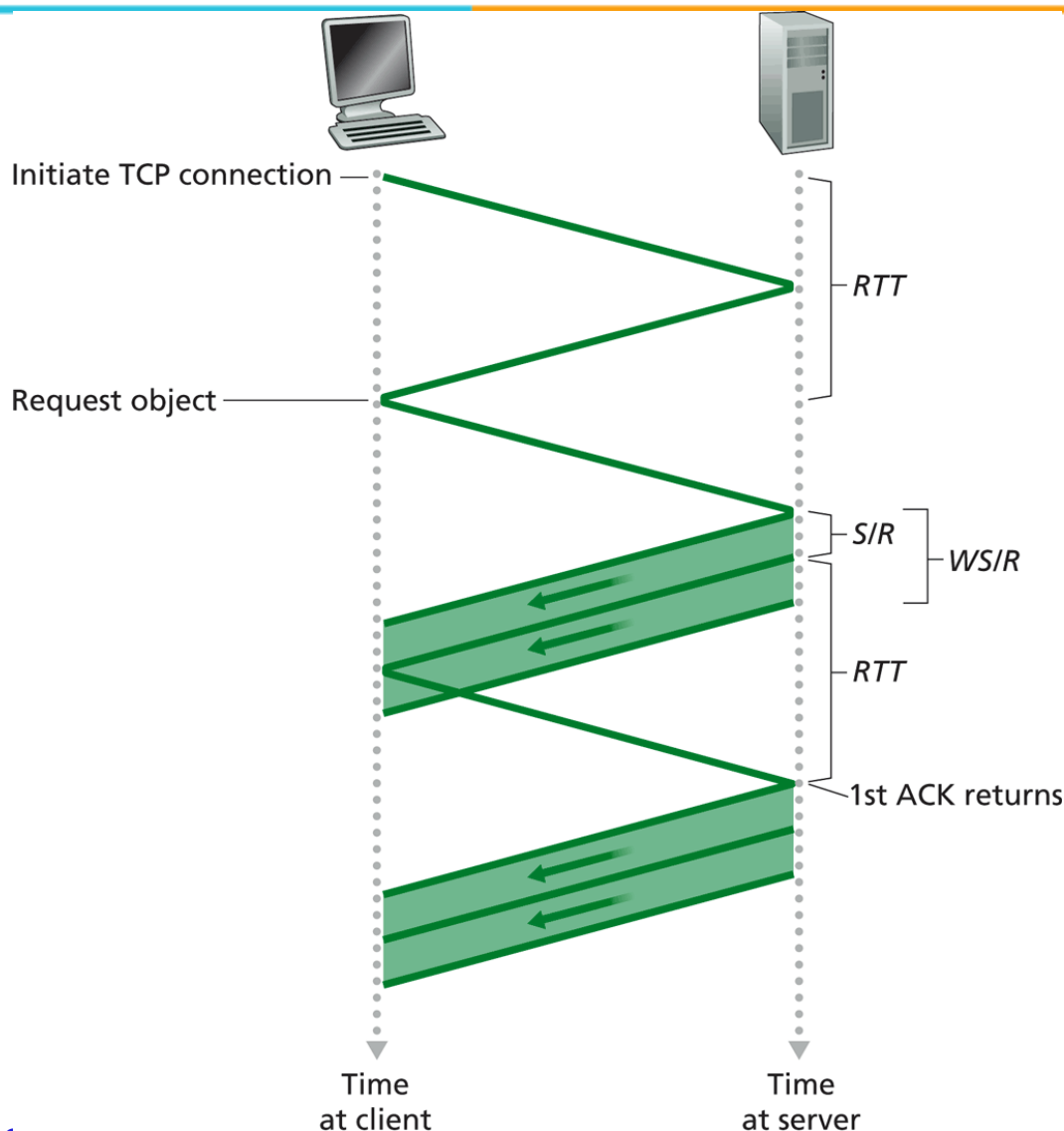
Janela de Congestionamento Estática

Segundo caso:

✓ $WS/R < RTT + S/R$
✓ ou seja, o ACK chega após enviar a janela de dados completa.

✓ $K = O/WS$

atraso = $2RTT + O/R + (K-1)[S/R + RTT - WS/R]$



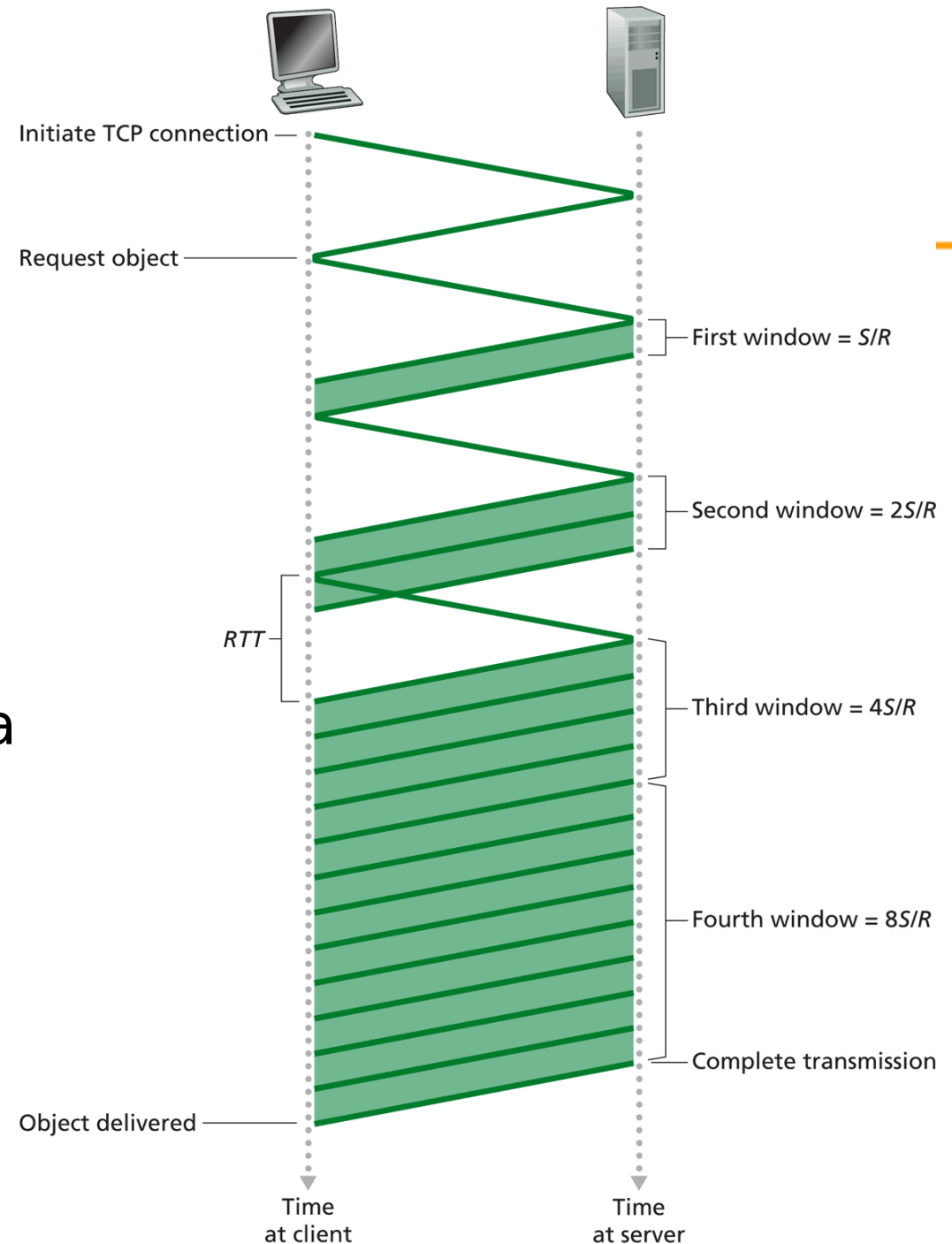


EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

Janela Dinâmica

Agora suponha que a janela cresça de acordo com os procedimentos da fase partida lenta.



ACH:

Figure 3.57 ♦ TCP timing during slow start



EACH

Escola de Artes, Ciências e Humanidades
da Universidade de São Paulo

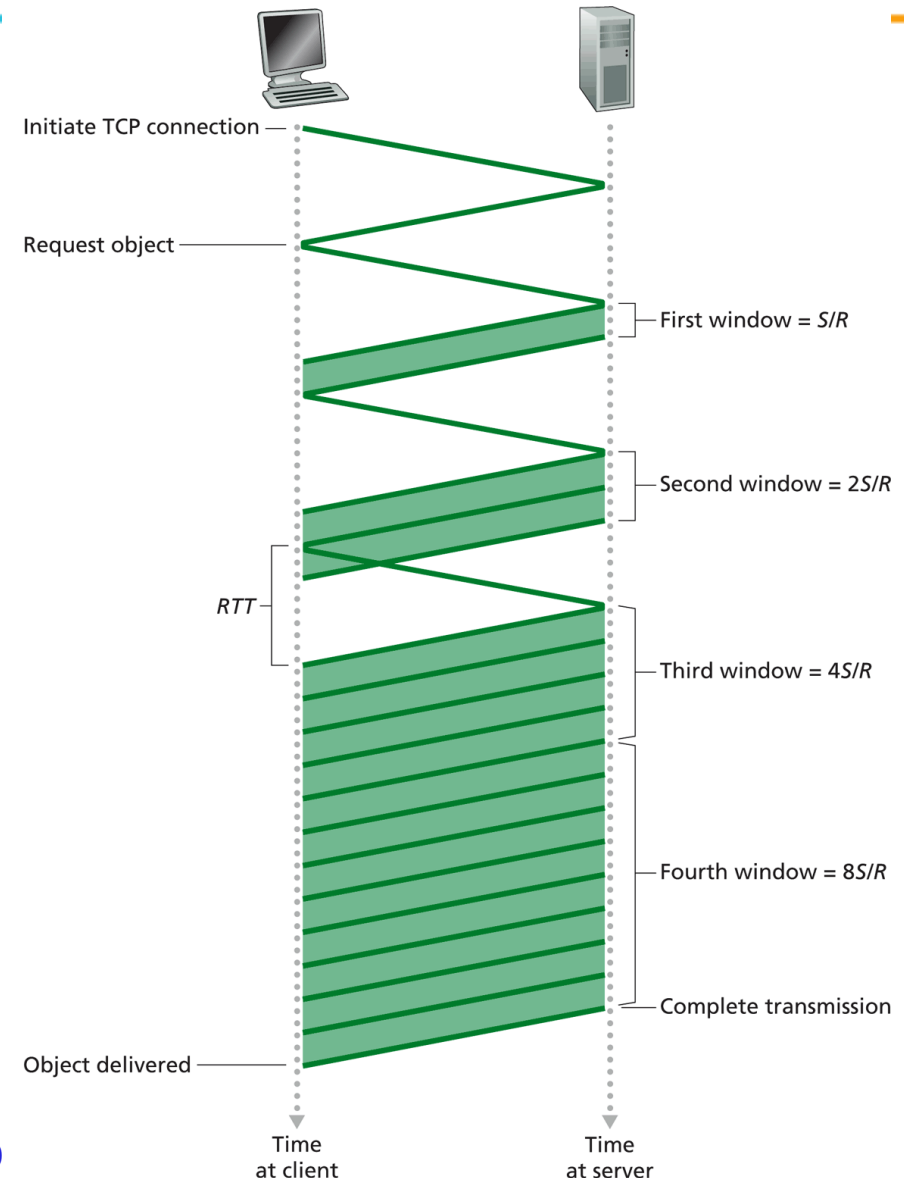
Janela de Congestionamento Dinâmica - Atraso

Componentes do atraso:

- 2 RTT para estabelecimento de conexão e requisição.
- O/R para transmitir um objeto.
- Servidor com Q períodos inativos devido à **partida lenta**.

Servidor inativo:

$P = \min\{K-1, Q\}$ vezes.



Janela de Congestionamento Dinâmica – Atraso

A latência (ou atraso) de um objeto de tamanho O é:

$$Latency = 2RTT + \frac{O}{R} + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R}$$

onde:

- P é o número de vezes em que o TCP fica bloqueado no servidor: $P = \min\{Q, K-1\}$
- Q é o número de vezes que o servidor ficaria bloqueado se o objeto fosse de tamanho infinito;
- K é o número de janelas que cobrem o objeto.



Modelagem de atraso: HTTP

Assuma que uma página Web consista em:

- 1 página HTML de base (de tamanho O bit)
- M imagens (cada uma de tamanho O bit)
- **HTTP não persistente:**
 - $M + 1$ conexões TCP nos servidores
 - **Tempo de resposta = $(M + 1)O/R + (M + 1)2RTT$ + soma dos períodos de inatividade**
- **HTTP persistente:**
 - **2 RTT** para requisitar e receber o arquivo HTML de base
 - **1 RTT** para requisitar e receber M imagens
 - **Tempo de resposta = $(M + 1)O/R + 3RTT$ + soma dos períodos de inatividade**
- **HTTP não persistente com X conexões paralelas**
 - Suponha o inteiro M/X
 - 1 conexão TCP para o arquivo de base
 - M/X ajusta as conexão paralelas para imagens
 - **Tempo de resposta = $(M + 1)O/R + (M/X + 1)2RTT$ + soma dos períodos de inatividade**



Exemplos de Atraso TCP

- $S=536$ bytes
- $RTT = 100$ ms (atraso intercontinental típico)
- $O = 100$ KB, número de janelas $K = 8$.

R	O/R	P	Latência Mínima	Latência com partida lenta
28Kbps	28,6s	1	28,8s	28,9s
100 Kbps	8s	2	8,2s	8,4s
1Mbps	800ms	5	1s	1,5s
10Mbps	80ms	7	0,28s	0,98s



A seguir...

- ∇ Saímos da “borda” da rede (camadas de aplicação e de transporte).
- ∇ Vamos para o “núcleo” da rede.

- ∇ Próxima aula (21/set):
 - Lab.
 - tragam as dúvidas!
- ∇ Prova 1: 25/set
 - Capítulos 1, 2 e 3.