

A Camada de Transporte

Objetivos :

- ❑ entender os princípios por trás dos serviços da camada de transporte:
 - Multiplexação / demultiplexação
 - transferência de dados confiável
 - controle de fluxo
 - controle de congestionamento
- ❑ instanciação e implementação na Internet

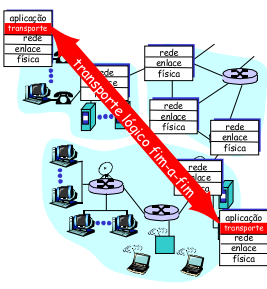
Resumo:

- ❑ serviços da camada de transporte
- ❑ Multiplexação / demultiplexação
- ❑ transporte sem conexão: UDP
 - princípios de transferência confiável de dados
- ❑ transporte orientado a conexão: TCP
 - transferência confiável
 - controle de fluxo
 - gerenciamento de conexão
- ❑ princípios de controle de congestionamento
- ❑ controle de congestionamento no TCP

1

Protocolos e Serviços de Transporte

- ❑ Fornecem **comunicação lógica** entre processos de aplicação rodando em hospedeiros diferentes
- ❑ Os protocolos de transporte são executados nos sistemas finais da rede



2

Serviço de Transporte vs Rede

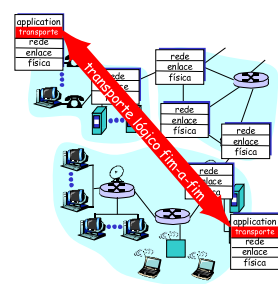
- ❑ **camada de rede:**
comunicação lógica entre hospedeiros
- ❑ **camada de transporte:**
comunicação lógica entre processos
 - utiliza e aprimora os serviços oferecidos pela camada de rede

3

Protocolos da Camada de Transporte

Serviços de Transporte da Internet:

- ❑ entrega *unicast*, sequencial e confiável: TCP
 - estabelecimento de conexão
 - controle de fluxo
 - controle de congestionamento
- ❑ entrega *unicast* ou *multicast*, não confiável (*best-effort*) e não sequencial: UDP
- ❑ não disponíveis:
 - garantia de retardo
 - garantia de banda
 - *multicast* confiável



4

Multiplexação/demultiplexação

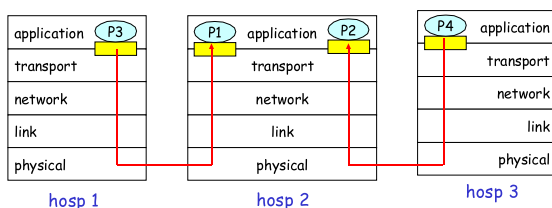
Demultiplex no receptor:

entrega dos segmentos recebidos ao soquete correto

Multiplex no emissor:

reunir dados de múltiplos soquetes; encapsular os dados com cabeçalho

■ = soquete ○ = processo



5

Demultiplexação sem conexão

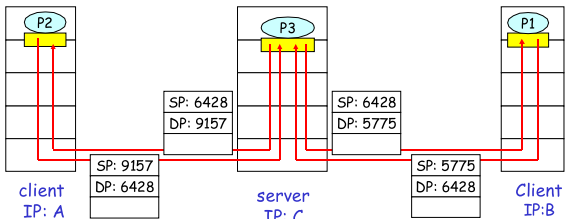
- ❑ Criar soquetes com números de portas:


```
DatagramSocket mySocket1 = new DatagramSocket(99111);
DatagramSocket mySocket2 = new DatagramSocket(99222);
```
- ❑ Soquete UDP identificado pela 2-tupla: (end IP dest, numero porta dest)
- ❑ Hospedeiro recebe um segmento UDP :
 - Verifica o número da porta de destino no segmento
 - Repassa o segmento para o soquete com esse número de porta
- ❑ Datagramas IP com endereços IP fonte e/ou números de portas fonte diferentes entregues ao mesmo soquete

6

Demultiplexação sem conexão (cont)

```
DatagramSocket serverSocket = new DatagramSocket(6428);
```



SP (Source Port) fornece o "endereço de retorno"

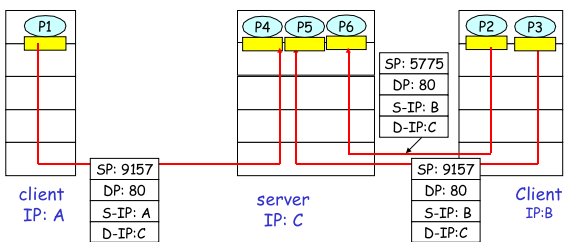
7

Demultiplexação orientada a conexão

- Soquete TCP identificado pela 4-tupla:
 - End. IP fonte
 - Num. porta fonte
 - End. IP destino
 - Num. porta destino
- Hospedeiro receptor usa os 4 valores para dirigir o segmento ao soquete apropriado
- Hospedeiros servidores podem ter vários soquetes TCP simultaneamente
 - Cada um identificado pela sua quadrupla
- Servidores Web usam soquetes diferentes para cada cliente conectado
 - HTTP não persistente terá um soquete diferente para cada requisição

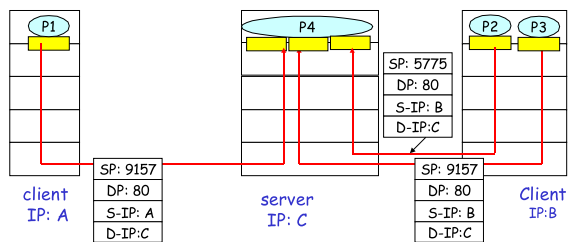
8

Demultiplexação orientada a conexão (cont)



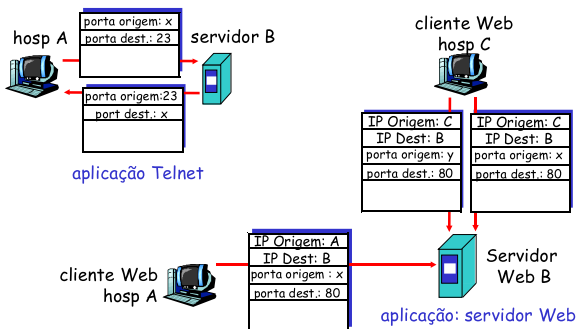
9

Demultiplexação orientada a conexão : *Servidor com Threads*



10

Multiplexação /Demultiplexação : exemplos



11

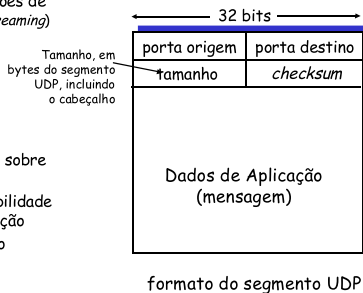
UDP: User Datagram Protocol [RFC 768]

- protocolo de transporte da Internet "sem gorduras", "sem frescuras"
- serviço "melhor esforço", segmentos UDP podem ser:
 - perdidos
 - entregues fora de ordem
- **sem conexão:**
 - não há apresentação entre o UDP transmissor e o receptor
 - cada segmento UDP é tratado de forma independente dos outros

12

UDP: mais

- muito usado por aplicações de multimídia contínua (*streaming*)
 - tolerantes a perdas
 - sensíveis à taxa
- outros usos do UDP:
 - DNS
 - SNMP
- transferência confiável sobre UDP:
 - acrescentar confiabilidade na camada de aplicação
 - recuperação de erro específica de cada aplicação



13

UDP : Checksum

Objetivo: detectar "erros" (ex., bits trocados) no segmento transmitido

Transmissor:

- trata o conteúdo do segmento como uma sequência de inteiros de 16 bits
- *checksum*: complemento de 1 da soma do conteúdo do segmento
 - colocado no campo de *checksum* do UDP

Receptor:

- computa o *checksum* do segmento recebido
- verifica se o *checksum* calculado é igual recebido:
 - NÃO - erro detectado
 - SIM - não há erros. Mas, talvez haja erros apesar disto

14

Checksum da Internet: Exemplo

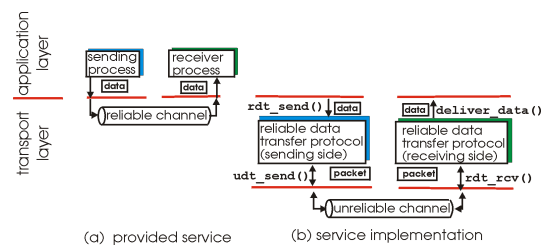
- Nota
 - Ao adicionar números, um vai-um do bit mais significativo deve ser somado ao resultado
- Exemplo: soma de dois inteiros de 16-bits

| | | | | | | | | | | | | | | | | |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 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 |
| Vai-um | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 1 |
| soma | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 |
| checksum | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |

15

Transferência Confiável de Dados

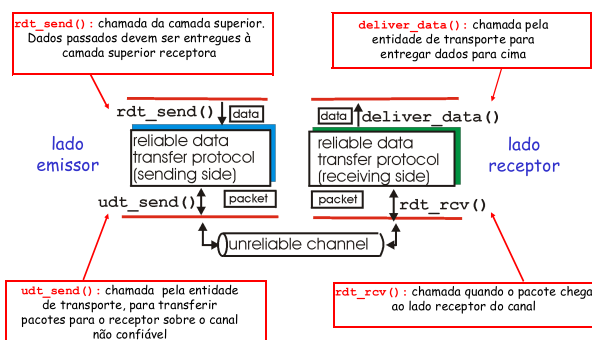
- importante nas camadas de aplicação, transporte e enlace
- um dos 10 tópicos mais importantes de redes!



as características dos canais não confiáveis determinarão a complexidade dos protocolos de transferência de dados confiável (rdt)

16

Transferência confiável: o ponto de partida



17

Transferência confiável: o ponto de partida

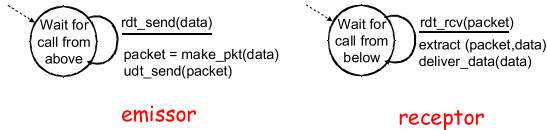
- desenvolvimento incremental dos lados transmissor e receptor de um protocolo de transferência de dados confiável (rdt)
- Considera-se apenas a transferência de dados unidirecional
 - entretanto informação de controle flui em ambas direções!
- uso de máquinas de estados finitos (MEF) para especificar o transmissor e o receptor



18

Rdt1.0: transferência confiável sobre canal confiável

- canal subjacente perfeitamente confiável
 - não há erros de bits
 - não há perdas de pacotes
- MEFs separadas para transmissor e receptor:
 - transmissor envia dados para o canal subjacente
 - receptor lê os dados do canal subjacente



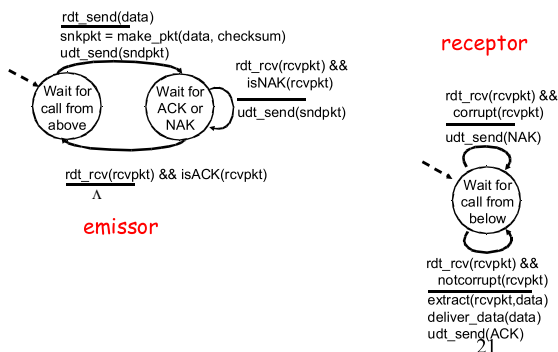
19

Rdt2.0: canal com erros de bit

- canal subjacente pode trocar valores de bits de um pacote
 - `checksum` detecta erros de bits
- a questão: como se recuperar desses erros:
 - reconhecimentos (ACKs)**: receptor avisa explicitamente ao transmissor que o pacote foi recebido corretamente
 - reconhecimentos negativos (NAKs)**: receptor avisa explicitamente ao transmissor que o pacote tem erros
 - transmissor reenvia o pacote quando da recepção de um NAK
- novos mecanismos no `rdt2.0`:
 - detecção de erros
 - retorno do receptor: mensagens de controle (`ACK`, `NAK`) `rcvr` -> `sender`

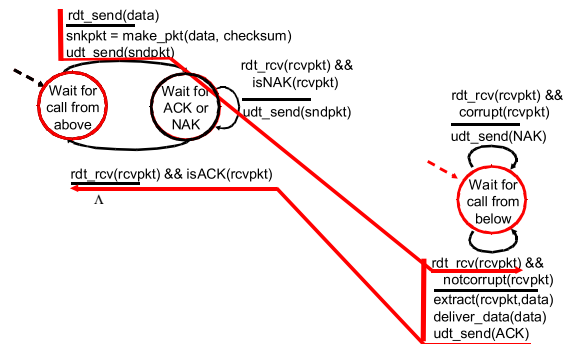
20

rdt2.0: especificação da MEF



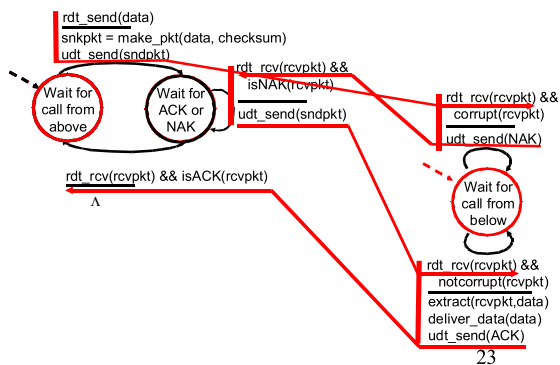
21

rdt2.0: operação sem de erros



22

rdt2.0: operação com erros



23

rdt2.0 tem um problema fatal!

O que acontece se o **ACK/NAK** é corrompido?

- emissor não sabe o que aconteceu no receptor!

O que fazer?

- retransmitir: pode gerar uma duplicata;

Tratamento de duplicatas:

- emissor acrescenta **número de sequência** em cada pacote
- emissor reenvia o pacote corrente se `ACK/NAK` for corrompido
- receptor descarta (não passa para a cima) pacotes duplicados

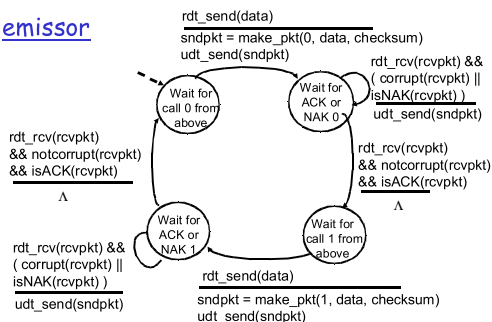
parar e esperar (stop and wait)

Emissor envia um pacote e então espera pela resposta do receptor

24

rdt2.1: trata ACK/NAKs corrompidos

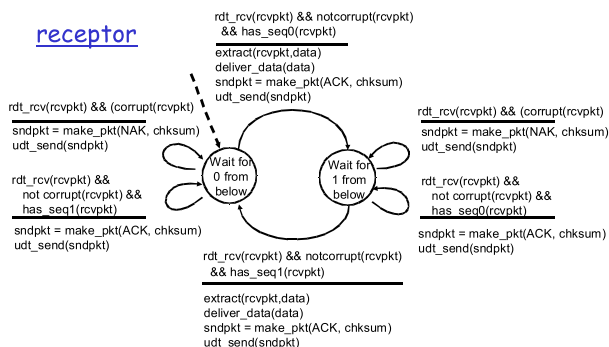
emissor



25

rdt2.1: trata ACK/NAKs corrompidos

receptor



26

rdt2.1: discussão

Emissor:

- adiciona número de sequência ao pacote
- dois números de sequência (0 e 1) bastam.
- deve verificar se os ACK/NAK recebidos estão corrompidos
- duplicação do número de estados (vs. rdt 2.0)
 - o estado deve "lembrar" se o pacote "corrente" tem número de sequência 0 ou 1

Receptor:

- deve verificar se o pacote recebido é duplicado
 - estado indica que pacote é esperado (0 ou 1)
- nota: receptor não tem como saber se seu último ACK/NAK foi recebido pelo emissor

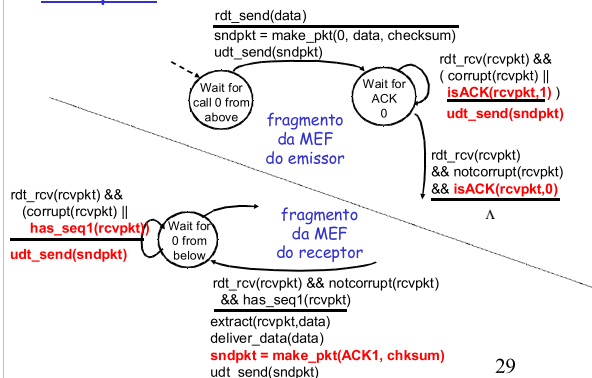
27

rdt2.2: um protocolo sem NAK

- mesma funcionalidade do rdt2.1, usando somente ACKs
- ao invés de enviar NAK, o receptor envia ACK para o último pacote recebido sem erro
 - receptor deve incluir explicitamente o número de sequência do pacote sendo reconhecido
- ACKs duplicados no transmissor resultam na mesma ação do NAK: *retransmissão do pacote corrente*

28

rdt2.2: trechos do emissor e do receptor



29

rdt3.0: canais com erros e perdas

Nova Hipótese: canal de transmissão pode também perder pacotes (dados ou ACKs)

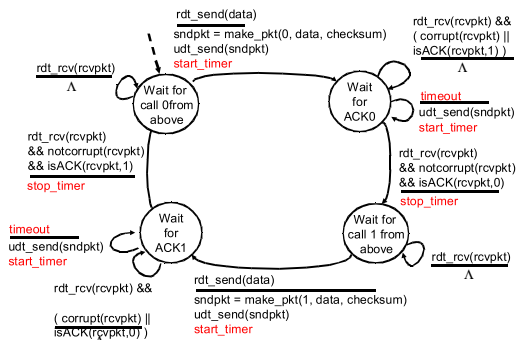
- checksum, números de sequência, ACKs, retransmissões serão úteis, mas não suficientes.

Abordagem: emissor espera um tempo "razoável" pelo ACK

- retransmite se nenhum ACK for recebido neste tempo
- se o pacote (ou ACK) estiver apenas atrasado (não perdido):
 - retransmissão será duplicata (números de sequência já tratam disso)
 - receptor deve especificar o número de sequência do pacote sendo reconhecido
- exige um temporizador

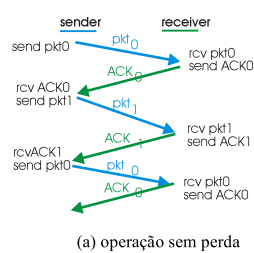
30

rdt3.0 emissor

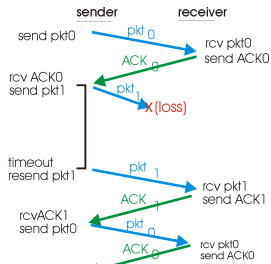


31

rdt3.0 em ação



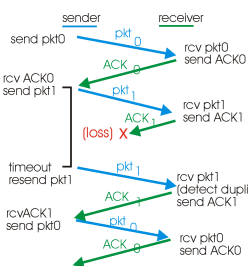
(a) operação sem perda



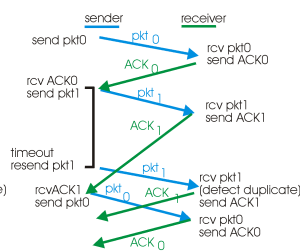
(b) pacote perdido

32

rdt3.0 em ação



(c) ACK perdido



(d) timeout prematuro

33

Desempenho do rdt3.0

- rdt3.0 funciona, mas o seu desempenho é sofrível
- exemplo: enlace de 1 Gbps, 15 ms de atraso de propagação, pacotes de 1000 bytes:

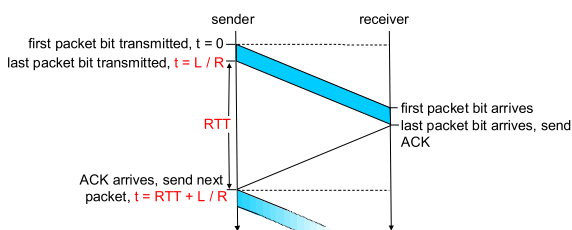
$$T_{\text{transmissão}} = \frac{8\text{kb}/\text{pkt}}{10^{**}9 \text{ b/s}} = 8 \mu\text{s}$$

$$U_{\text{emissor}} = \frac{L/R}{RTT+L/R} = \frac{0,008\text{ms}}{30,008\text{ms}} = 0,00027$$

- Um pacote de 1000 bytes a cada 30,008 ms -> 267 kbps de vazão sobre um canal de 1 Gbps
- o protocolo de rede limita o uso dos recursos físicos!

34

rdt3.0: funcionamento do "stop-and-wait"



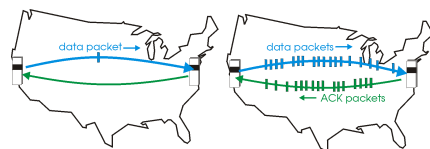
$$U_{\text{emissor}} = \frac{L/R}{RTT+L/R} = \frac{0,008\text{ms}}{30,008\text{ms}} = 0,00027$$

35

Protocolos com Pipelining/janela

Pipelining: o emissor permite vários pacotes "em trânsito" ainda não reconhecidos

- faixa de números de sequência deve ser aumentada
- armazenamento no transmissor e/ou no receptor



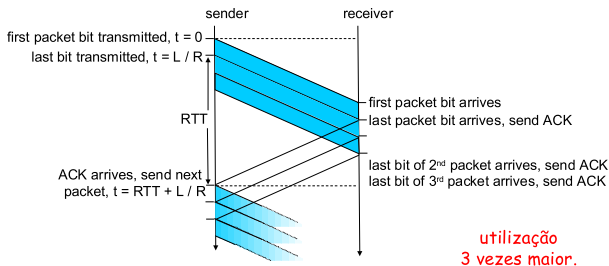
(a) operação do protocolo stop-and-wait

(b) operação do protocolo com pipelining

- Protocolos com pipelining: *go-back-N*, *retransmissão seletiva*

36

Pipelining: aumento da utilização



utilização
3 vezes maior.

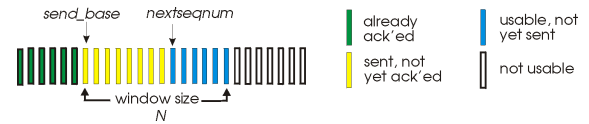
$$U_{\text{emissor}} = \frac{3 * L/R}{RTT + L/R} = \frac{3 * 0,008\text{ms}}{30,008\text{ms}} = 0,0008$$

37

Go-Back-N (GBN)

Emissor:

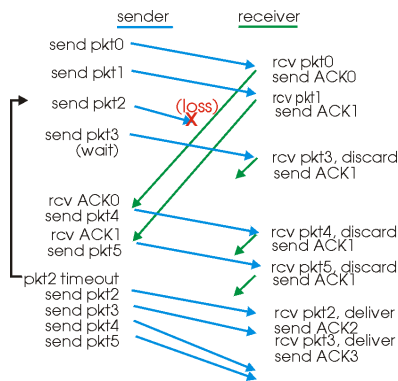
- Número de sequência com k bits no cabeçalho do pacote
- "janela" de tamanho N :
 - N pacotes não reconhecidos consecutivos permitidos



- ACK(n): reconhece todos os pacotes até o número de sequência n , inclusive. "ACK cumulativo"
- timeout(n): retransmite o pacote n e todos os pacotes subsequentes que estejam dentro da janela

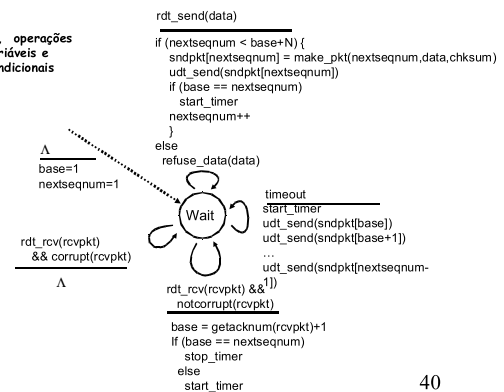
38

GBN em ação



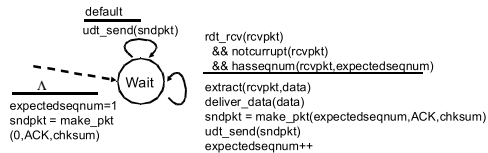
GBN: MEF estendida* do emissor

* variáveis, operações sobre variáveis e ações condicionais



40

GBN: MEF estendida do receptor



- sempre envia ACK para o pacote corretamente recebido com o maior número de sequência *em ordem*
 - pode gerar ACKs duplicados
 - precisa lembrar apenas do número de sequência esperado (*expectedseqnum*)
- pacotes fora de ordem:
 - descarta (não armazena) -> **não há buffer de recepção!**
 - Reconfirma pacote com o mais alto número de sequência em ordem

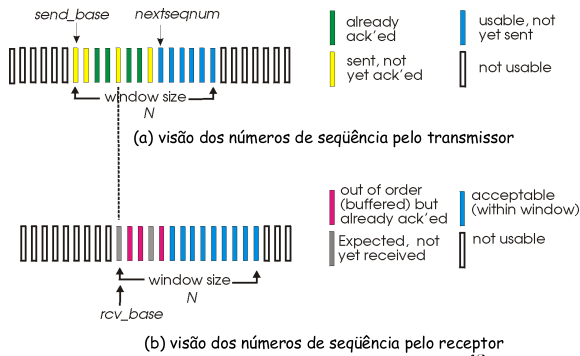
41

Retransmissão Seletiva (RS)

- receptor reconhece *individualmente* todos os pacotes recebidos corretamente
 - armazena pacotes, quando necessário, para eventual entrega em ordem para a camada superior
- transmissor somente reenvia os pacotes para os quais um ACK não foi recebido
 - transmissor temporiza cada pacote não reconhecido
- janela de transmissão
 - N números de sequência consecutivos
 - novamente limita a quantidade de pacotes enviados, mas não reconhecidos

42

Retransmissão seletiva: janelas do transmissor e do receptor



43

Retransmissão seletiva

emissor

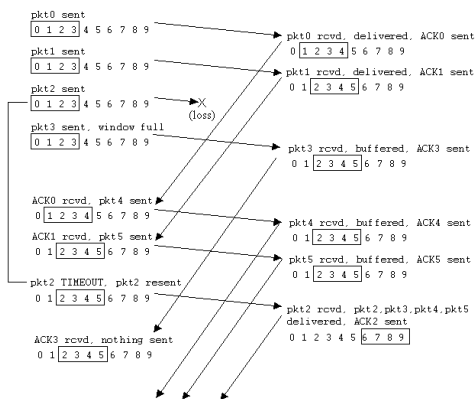
- dados da camada superior :**
- se** próximo número de sequência está na janela
então envia o pacote
senão armazena ou recusa os dados
 - timeout(n):**
 - reenvia pacote n
 - reinicia o temporizador.
 - ACK(n)** em [sendbase, sendbase+N]:
 - marca pacote n como recebido
 - se** n = sendbase
então avança a base da janela para o próximo número de sequência não reconhecido

receptor

- pacote n em** [rcvbase, rcvbase+N-1]
- envia ACK(n)
 - fora de ordem: armazena
 - em ordem:
 - entrega todos os pacotes armazenados em ordem
 - avança base da janela para o próximo pacote ainda não recebido
 - pkt n em** [rcvbase-N, rcvbase-1]
 - ACK(n)
 - caso contrário:**
 - ignora

44

Retransmissão seletiva em ação

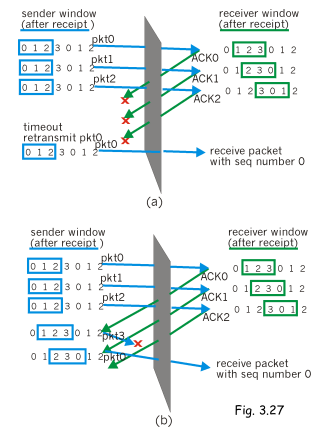


Retransmissão seletiva: dilema

- Exemplo:**
- #s de sequência: 0, 1, 2, 3
 - tamanho da janela=3

- receptor não vê diferença nos dois cenários!
- incorretamente passa dados duplicados como novos (figura a)

Q: relação entre o tamanho do espaço de números de sequência e o tamanho da janela?



46

Canal com memória

- Os protocolos apresentados pressupõem que o canal não reordena os pacotes, i.e., o canal não tem memória.
- Em canais com memória, um emissor não deve reutilizar um número de sequência até ter certeza de que ele não está mais presente na rede.

47

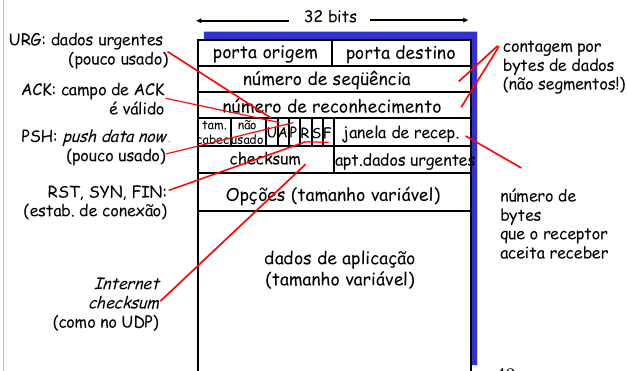
TCP: Visão Geral RFCs: 793, 1122, 1323, 2018, 2581

- ponto-a-ponto:**
 - um transmissor, um receptor
- fluxo de bytes, confiável, em ordem:**
 - não há limites de mensagens
- pipelined:**
 - Controle de congestão e de fluxo definem tamanho da janela
- buffers de transmissão e de recepção**
- dados full-duplex:**
 - transmissão bidirecional na mesma conexão
 - MSS: maximum segment size
 - quantidade máxima de dados da aplicação no segmento
- orientado a conexão:**
 - apresentação (handshaking) - troca de mensagens de controle: inicia o estado do transmissor e do receptor antes da troca de dados
- controle de fluxo:**
 - transmissor não esgota a capacidade do receptor



48

Estrutura do Segmento TCP



49

Números de Sequência e ACKs do TCP

Números de sequência:

- Número, no fluxo de dados, do primeiro byte no segmento de dados

ACKs:

- número do próximo byte esperado do outro lado
- ACK cumulativo

- Q: como o receptor trata segmentos fora de ordem?
- A especificação do TCP não define, fica a critério do implementador

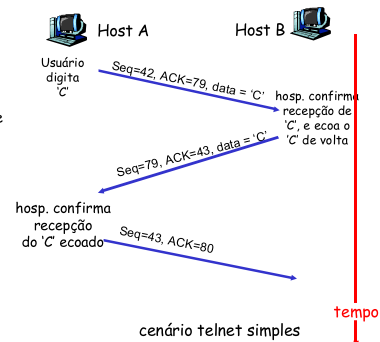


Fig. 3.30

TCP: RTT e Temporização

Q: como escolher o valor de temporização do TCP?

- maior que o RTT
- muito curto: temporização prematura
 - retransmissões desnecessárias
- muito longo: reação lenta à perda de segmentos

Q: Como estimar o RTT?

- SampleRTT: tempo medido desde a transmissão de um segmento até a recepção da 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

51

TCP: RTT e Temporização

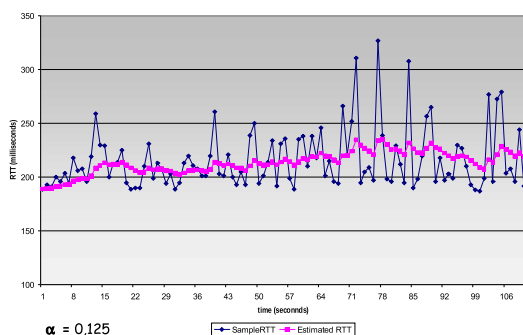
$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

- Média móvel com peso exponencial
- influência de uma amostra passada decresce de forma exponencial
- valor típico de $\alpha = 0,125$ (isto é $1/8$)

52

Amostras e Estimativas do RTT

RTT: galax.usmss.edu to fantasia.eurocom.fr



53

TCP: RTT e Temporização

Definindo o intervalo de temporização

- EstimatedRTT mais "margem de segurança"
- Primeiro estima-se de quanto SampleRTT se desvia de EstimatedRTT:

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

(tipicamente, $\beta = 0,25$)

- Então, define-se o intervalo de temporização:

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

54

```

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)
    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 not-yet-acknowledged segments)
        start_timer()
    }
    break;
} /* end of loop forever */

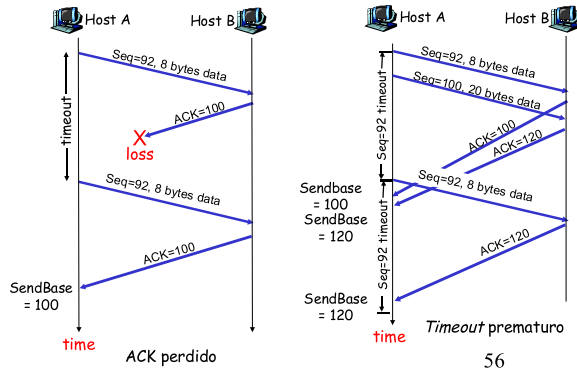
```

TCP:
Transf.
Confiável

Emissor
(simplificado)

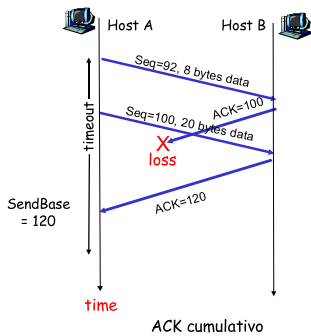
55

TCP: cenários de retransmissão



56

TCP: cenários de retransmissão (cont.)



57

```

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(to:=TimeoutInterval)
    pass segment to IP
    NextSeqNum = NextSeqNum + length(data)
    break;
  event: timer timeout
    retransmit not-yet-acknowledged segment with
    smallest sequence number
    // forma limitada de controle de congestionamento
    start_timer(to := 2 * to)
    break;
  event: ACK received, with ACK field value of y
    if (y > SendBase) {
      SendBase = y
      if (there are currently not-yet-acknowledged segments)
        start_timer(to := TimeoutInterval)
    }
    break;
} /* end of loop forever */

```

TCP:
Transf.
Confiável

Emissor
(simplificado)

58

TCP: Geração de ACKs [RFC 1122, RFC 2581]

| Evento no Receptor | Ação do TCP Receptor |
|--|---|
| segmento chega em ordem, não há lacunas, seg.s anteriores confirmados | ACK retardado. Espera até 500ms pelo próximo segmento. Se não chegar, envia ACK |
| segmento chega em ordem, não há lacunas, um ACK atrasado pendente | envia imediatamente um ACK cumulativo |
| segmento chega fora de ordem # de seq. maior que o esperado lacuna detectada | envia ACK repetido, indicando número de sequência do próximo byte esperado |
| chegada de segmento que preenche uma lacuna parcial ou completamente | reconhece imediatamente, se o segmento começa na borda inferior da lacuna |

59

Retransmissão Rápida [RFC 2581]

- Período de *timeout* em geral é relativamente longo:
 - Grande demora para o reenvio de pacotes perdidos
- Deteção de segmentos perdidos através de ACKs repetidos:
 - O emissor costuma mandar muitos segmentos um logo após o outro
 - Se um segmento é perdido, provavelmente cheguem muitos ACKs repetidos.
- SE o emissor receber 3 ACKs para o mesmo dado, ele supõe que o segmento após o dado confirmado foi perdido:
 - retransmissão rápida:** reenvia o segmento antes que seu temporizador expire

60

Algoritmo de Retransmissão Rápida

```

event: ACK received, with ACK field value of y
  if (y > SendBase) {
    SendBase = y
    if (there are currently not-yet-acknowledged segments)
      start timer
  }
  else {
    increment count of ACKs received for y
    if (count of ACKs received for y = 3) {
      resend segment with sequence number y
    }
  }

```

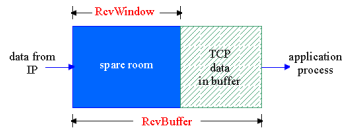
ACK repetido para segmento já confirmado

Retransmissão rápida

61

TCP: Controle de Fluxo

- Lado receptor do TCP tem um *buffer* de recepção:



controle de fluxo

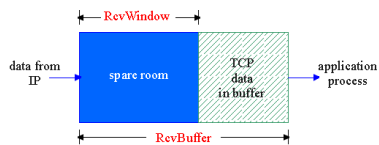
transmissor não deve fazer transbordar o buffer de recepção enviando dados rápido demais

- Processos de aplicação podem ser lentos na leitura do *buffer*

- Serviço de compatibilização das taxas de envio e recepção

62

TCP: Controle de Fluxo



$\text{LastByteRcvd} - \text{LastByteRead} \leq \text{RcvBuffer}$

$\text{RcvWindow} = \text{RcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$

receptor: informa dinamicamente o espaço livre incluindo *RcvWindow* nos segmentos

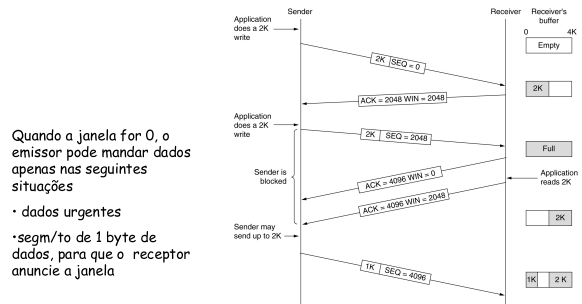
transmissor: limita a quantidade de dados transmitidos mas não reconhecidos a *RcvWindow*

LastByteRcvd: número do último byte na cadeia de dados colocado no buffer

LastByteRead: número do último byte na cadeia de dados lido do buffer

63

TCP: Controle de Fluxo



TANENBAUM, A. S., Computer Networks, 4rd. Ed., Prentice-Hall, 2003. Capítulo 6.

64

TCP: Gerência de Conexão

Lembre-se: O TCP transmissor estabelece uma "conexão" com o receptor antes de trocar segmentos de dados

- inicia variáveis:
 - o números de sequência
 - o buffers, info. de controle de fluxo (ex. *RcvWindow*)
- **cliente:** iniciador da conexão
- **servidor:** chamado pelo cliente

Apresentação em três vias

Passo 1: sistema final cliente envia um segmento TCP SYN ao servidor

- o especifica número de sequência inicial

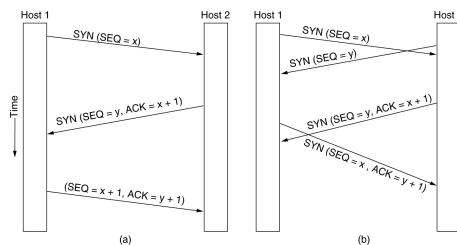
Passo 2: o sistema final servidor recebe o SYN, responde com segmento SYNACK

- o aloca buffers
- o especifica o número de sequência inicial do servidor

Passo 3: o cliente recebe SYNACK, responde com segmento ACK, que pode conter dados

65

TCP: Estabelecimento de Conexão



(a) Situação normal; (b) Colisão de chamada

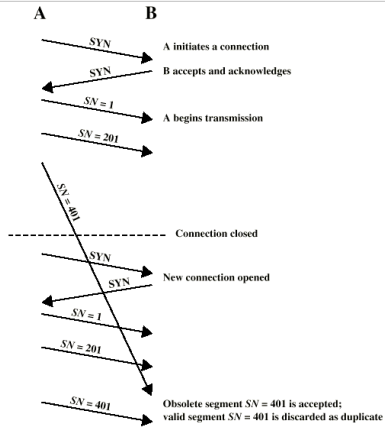
TANENBAUM, A. S., Computer Networks, 4rd. Ed., Prentice-Hall, 2003. Capítulo 6.

66

Apresentação em duas vias: Segmento de dados obsoleto

Obs. Numeração de segmentos não segue a convenção do TCP

STALLINGS, W. Data and Computer Communications. 6th ed. Prentice-Hall, 2000.

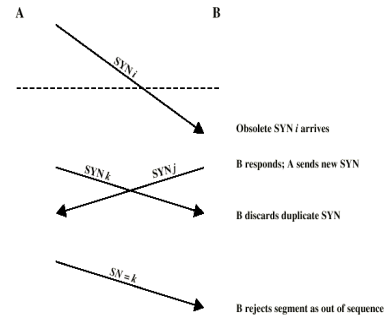


67

Apresentação em duas vias: Segmento SYN obsoleto

Obs. Numeração de segmentos não segue a convenção do TCP

STALLINGS, W. Data and Computer Communications. 6th ed. Prentice-Hall, 2000.

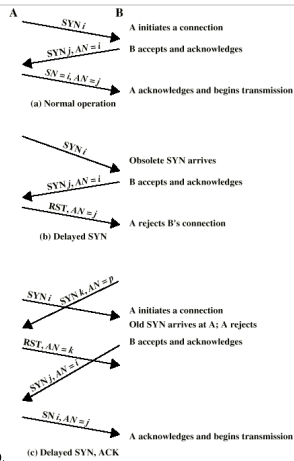


68

Apresentação em três vias: Exemplos

Obs. Numeração de segmentos não segue a convenção do TCP

STALLINGS, W. Data and Computer Communications. 6th ed. Prentice-Hall, 2000.



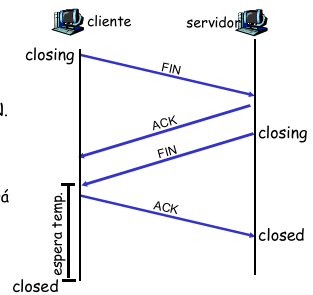
TCP Término de Conexão

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

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

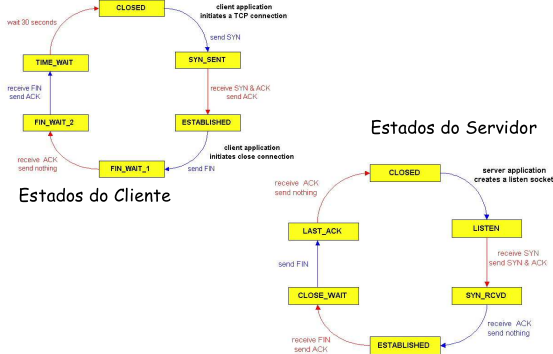
Passo 3: cliente recebe FIN, responde com ACK. Entra "espera temporizada" - responderá com ACK a FINs recebidos.

Passo 4: servidor recebe ACK. Conexão fechada.



70

TCP Controle de Conexão



71

TCP: Modelando a Gerência de Conexão

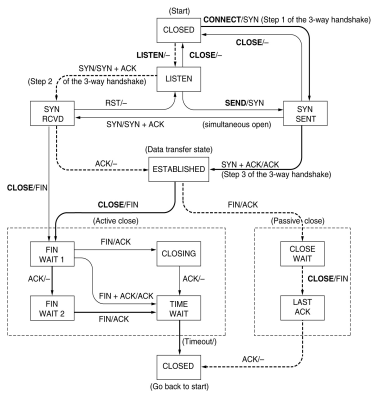
| State | Description |
|-------------|--|
| CLOSED | No connection is active or pending |
| LISTEN | The server is waiting for an incoming call |
| SYN RCVD | A connection request has arrived; wait for ACK |
| SYN SENT | The application has started to open a connection |
| ESTABLISHED | The normal data transfer state |
| FIN WAIT 1 | The application has said it is finished |
| FIN WAIT 2 | The other side has agreed to release |
| TIMED WAIT | Wait for all packets to die off |
| CLOSING | Both sides have tried to close simultaneously |
| CLOSE WAIT | The other side has initiated a release |
| LAST ACK | Wait for all packets to die off |

TANENBAUM, A. S., Computer Networks, 4nd. Ed., Prentice-Hall, 2003. Capítulo 6.

72

TCP: Modelando a Gerência de Conexão

- Linha grossa contínua: caminho normal do cliente
- Linha grossa tracejada: caminho normal do servidor
- Linhas finas: eventos não usuais
- Transições são rotuladas pelos eventos que as causam e pela ação resultante, separada por uma barra



TANENBAUM, A. S., Computer Networks, 4rd. Ed., Prentice-Hall, 2003. Capítulo 6.

73

Princípios de Controle de Congestionamento

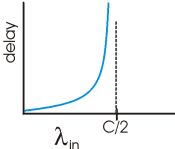
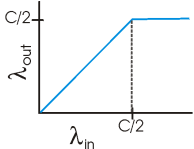
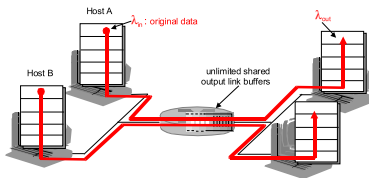
Congestionamento:

- informalmente: "muitas fontes enviando muitos dados muito rapidamente para a rede tratá-los"
- diferente de controle de fluxo!
- sintomas:
 - perda de pacotes (transbordamento de buffers nos roteadores)
 - atrasos elevados (filas nos buffers dos roteadores)
- um dos 10 problemas mais importantes em redes!

74

Causas/custos do congestionamento: cenário 1

- dois transmissores, dois receptores
- um roteador, buffers infinitos
- não há retransmissão

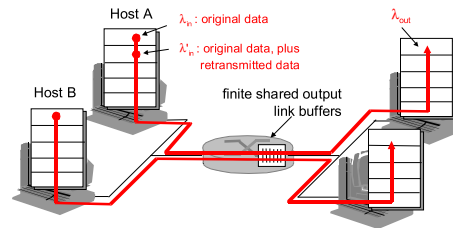


- Pode-se obter a vazão máxima
- Grandes atrasos quando congestionado

75

Causas/custos do congestionamento: cenário 2

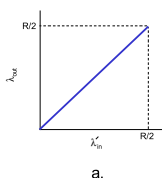
- um roteador, buffers finitos
- transmissor reenvia pacotes perdidos



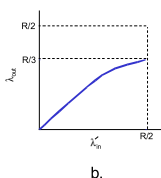
76

Causas/custos do congestionamento: cenário 2

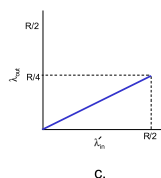
- (a) A envia um pacote apenas quando há um buffer disponível (A é adivinho!)
- (b) retransmissão "perfeita" somente quando há perdas (t.o. grande o bastante):
- (c) retransmissão de pacotes atrasados (não perdidos)



a.



b.



c.

"custos" do congestionamento:

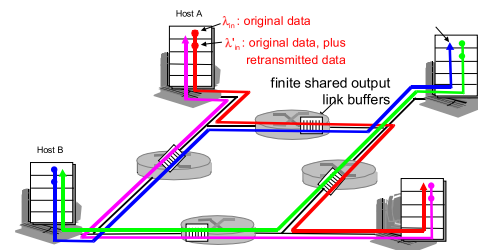
- mais trabalho (retransmissões) para um dado "tráfego bom"
- retransmissões desnecessárias: enlace transporta várias cópias do mesmo pacote

77

Causas/custos do congestionamento: cenário 3

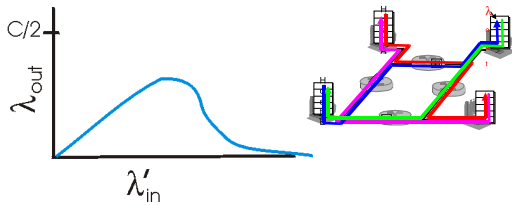
- quatro transmissores
- caminhos com múltiplos saltos
- temporizações/retransmissões

Q: o que acontece quando λ_{in} aumentam?



78

Causas/custos do congestionamento: cenário 3



Outro "custo" do congestionamento:

- quando um pacote é descartado, a capacidade de transmissão usada "corrente acima" para aquele pacote é desperdiçada!

79

Abordagens ao problema de controle de congestionamento

Duas abordagens gerais:

Controle de congestionamento fim-a-fim:

- não há realimentação explícita da rede
- congestionamento é inferido a partir das perdas e dos atrasos observados nos sistemas finais
- abordagem usada pelo TCP

Controle de congestionamento assistido pela rede:

- roteadores fornecem informações aos sistemas finais
 - bit indicando o congestionamento (SNA, DECbit, TCP/IP ECN, ATM)
 - taxa explícita a que o transmissor pode enviar

80

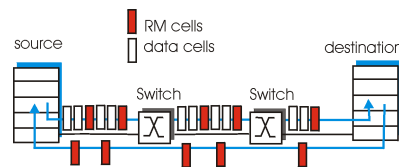
Estudo de caso: controle de congestionamento no serviço ATM ABR

ABR: available bit rate:

- "serviço elástico"
- se o caminho do transmissor está pouco usado:
 - transmissor pode usar a banda disponível
- se o caminho do transmissor está congestionado:
 - transmissor é limitado a uma taxa mínima garantida

81

Controle de congestionamento no ATM ABR



Células RM (resource management):

- bits nas células RM são modificados pelos comutadores
 - NI bit: não aumentar a taxa (congestionamento leve)
 - CI bit: indicação de congestionamento
- células RM são devolvidas pelo receptor, com os bits de indicação intactos
- comutadores tb. podem gerar e enviar células RM diretamente a uma fonte
- campo ER (explicit rate) de dois bytes nas células RM
 - comutador congestionado pode reduzir o valor de ER nas células
 - o transmissor envia dados de acordo com a vazão mínima no caminho
- bit EFCI nas células de dados: marcado pelos comutadores congestionados
 - se a célula de dados que precede uma célula RM tem o bit EFCI marcado, o destino marca o bit CI da célula RM devolvida à fonte

82

TCP: Controle de Congestionamento

- Controle fim-a-fim (sem assistência da rede)
- Emissor limita a transmissão:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \min\{\text{CongWin}, \text{RcvWindow}\}$$
- Aproximadamente,

- CongWin é dinamicamente ajustado pelo emissor

Como o emissor percebe o congestionamento?

- perda = timeout ou 3 acks duplicados (Reno)
- Emissor TCP reduz a taxa (CongWin) após uma perda

Três mecanismos:

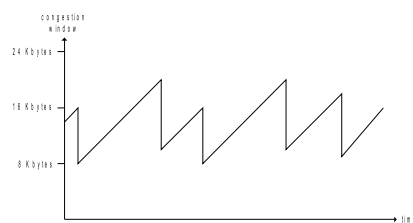
- AIMD
- Partida lenta
- Conservadorismo após timeouts

83

TCP: Controle de Congestionamento AIMD

Aumento aditivo: incrementar CongWin de 1 MSS a cada RTT na ausência de perdas: *testando*

Diminuição multiplicativa: cortar CongWin pela metade após uma perda

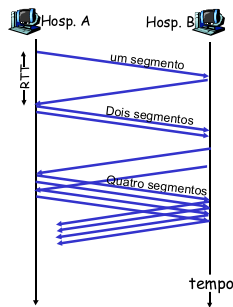


Conexão TCP ao longo do tempo

84

TCP: Controle de Congestionamento Partida Lenta (*Slow Start*)

- Quando uma conexão começa, $\text{CongWin} = 1 \text{ MSS}$
 - Exemplo: $\text{MSS} = 500 \text{ bytes}$, $\text{RTT} = 200 \text{ ms}$
 - Taxa inicial = 20 kbps
- Largura de banda disponível pode ser $\gg \text{MSS}/\text{RTT}$
 - Desejável subir rapidamente para uma taxa respeitável
- Quando a conexão começa, incrementar a taxa exponencialmente até o primeiro evento de perda:
 - Dobrar CongWin a cada RTT
 - Feito incrementando CongWin para cada ACK recebido
- **Resumo:** taxa inicial é baixa, mas cresce exponencialmente



85

TCP: Controle de Congestionamento Conservadorismo após *timeouts*

- Após 3 ACKs duplicados:
 - CongWin é cortada pela metade
 - janela cresce então linearmente
- Mas, após um *timeout*:
 - $\text{CongWin} = 1 \text{ MSS}$;
 - janela cresce então exponencialmente ...
 - ... até um patamar, então cresce linearmente

Motivos:

- 3 ACKs duplicados indicam que a rede é capaz de entregar alguns segmentos
- *timeout* antes de 3 ACKs duplicados é "mais alarmante"

86

TCP: Controle de Congestionamento Síntese

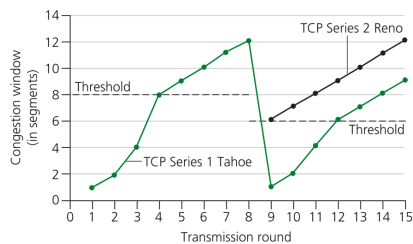


Figure 3.51 ♦ Evolution of TCP's congestion window (Tahoe and Reno)

87

Resumo: Controle de Congestionamentos no TCP

- Quando CongWin está abaixo de Threshold , o emissor está na fase de **partida lenta**, e a janela cresce exponencialmente.
- Quando CongWin está acima de Threshold , o emissor está na fase de **prevenção de congestionamento**, e a janela cresce linearmente.
- Quando ocorre um **triplo ACK duplicado**,
 $\text{Threshold} = \text{CongWin}/2$
 $\text{CongWin} = \text{Threshold}$.
- Quando ocorre um **timeout**,
 $\text{Threshold} = \text{CongWin}/2$
 $\text{CongWin} = 1 \text{ MSS}$.

88

TCP: controle de congestionamento no emissor

| Estado | Evento | Ação do TCP Emissor | Comentário |
|---------------------------|---------------------------------------|--|---|
| Slow Start (SS) | ACK para dados ainda não reconhecidos | $\text{CongWin} = \text{CongWin} + \text{MSS}$ If ($\text{CongWin} > \text{Threshold}$) State = CA | Resulta em dobrar CongWin a cada RTT |
| Congestion Avoidance (CA) | ACK para dados ainda não reconhecidos | $\text{CongWin} = \text{CongWin} + \text{MSS} * (\text{MSS}/\text{CongWin})$ | Aumento aditivo, resultando no aumento de CongWin de 1 MSS por RTT |
| SS or CA | ACK triplicado | $\text{Threshold} = \text{CongWin}/2$, $\text{CongWin} = \text{Threshold}$, State = CA | Recuperação rápida, implementação da diminuição multiplicativa. |
| SS or CA | Timeout | $\text{Threshold} = \text{CongWin}/2$, $\text{CongWin} = 1 \text{ MSS}$, State = SS | Entra na partida lenta |
| SS or CA | ACK duplicado | Incrementar o contador de duplicatas do segmento confirmado | CongWin e Threshold não são modificados |

89

TCP: vazão (*throughput*)

- Qual é a vazão média do TCP como função do tamanho da janela e do RTT?
 - Ignore a partida lenta
- Seja W o tamanho da janela, em bytes, quando ocorre uma perda
 - vazão é W/RTT
- Logo após a perda, a janela cai para $W/2$
 - vazão passa a $W/(2 \cdot \text{RTT})$
- Vazão aumenta linearmente entre $W/(2 \cdot \text{RTT})$ e W/RTT (o processo repete-se continuamente)
 - Vazão média de uma conexão = $0,75 W/\text{RTT}$

90

TCP: O Futuro

- Exemplo: segmentos de 1500 bytes, RTT de 100ms, deseja-se vazão de 10 Gbps

$$\text{Vazão média} = 0,75 W/\text{RTT}$$

$$10 \cdot 10^9 = 0,75 [(1500 \cdot 8) \cdot S]/100 \cdot 10^{-3}$$

$$S \sim 112.000$$

Precisamos de uma janela de ~112.000 segmentos (em trânsito)

- Vazão média em termos da taxa de perda:

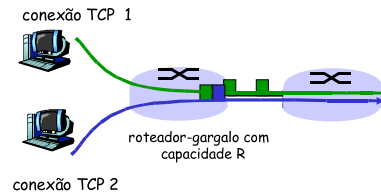
$$\text{Vazão média} = \frac{1,22 \cdot \text{MSS}}{\text{RTT} \sqrt{L}}$$

- Para alcançar a vazão de 10 Gbps precisamos de $L = 2 \cdot 10^9$ (um evento de perda a cada 5.000.000.000 de segmentos)
- Necessárias novas versões do TCP para redes de alta velocidade!

91

TCP: Equidade

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



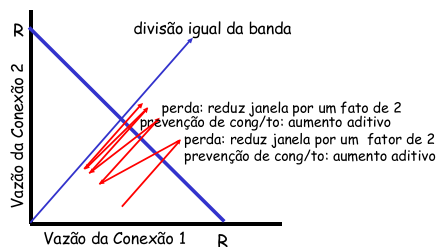
- Supomos que ambas conexões têm o mesmo MSS e RTT (mesma vazão para janelas iguais)
- Ignoramos a partida lenta.

92

Porque 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



93

TCP: Equidade (...)

Equidade e UDP

- Apl. multimídia geralmente não usam o TCP
 - não querem que a taxa seja reduzida pelo controle de congestionamento
- ... usam UDP:
 - Enviam áudio/vídeo a uma taxa constante
 - Toleram a perda de pacotes
- Área de pesquisa: controle de congestionamento X UDP

Equidade e Conexões TCP paralelas

- Nada impede as apl.s de abrirem conexões paralelas entre 2 hospedeiros
- Navegadores da Web fazem isso
- Exemplo: enlace de taxa R com 9 conexões;
 - Nova apl. abre 1 conexão TCP → recebe R/10
 - Nova apl. abre 9 conexões TCP → recebe R/2

94

A Camada de Transporte : Resumo

- princípios por trás dos serviços da camada de transporte:
 - Multiplexação / demultiplexação
 - transferência confiável de dados
 - controle de fluxo
 - controle de congestionamento
- instanciação e implementação na Internet
 - UDP
 - TCP

95