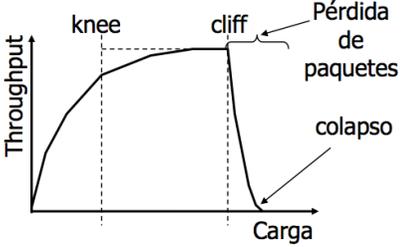


Control de Congestión

Controlar vs. Evitar

- Controlar
 - Permanecer a la izquierda del "cliff"
- Evitar
 - Permanecer a la izquierda del "knee"



The graph plots Throughput on the y-axis and Carga (Load) on the x-axis. The curve starts at the origin, rises steeply, then levels off at a point labeled 'knee'. After the knee, the curve continues to rise slightly to a peak and then drops sharply to zero at a point labeled 'cliff'. The region between the knee and the cliff is labeled 'Pérdida de paquetes' (Packet loss). The region after the cliff is labeled 'colapso' (collapse).

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 6 / 61

Control de Congestión

Modelos

- Modelo end-to-end
 - Los extremos son la fuente de la demanda
 - Los extremos deben estimar los tiempos y grado de congestión y reducir la demanda
 - Los nodos intermedios deben monitorear el estado de la red.
- Modelo basado en la red
 - Los extremos no son confiables
 - El nodo de la red tiene control sobre el tráfico
 - Acciones más rápidas

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 7 / 61

Notas:

Control de Congestión

TCP – Modelo basado en la Red

- Utiliza tres variables:
 - cwnd: ventana de congestión
 - rcv_win: ventana del receptor. Publicada en el segmento.
 - ssthresh: valor del umbral de "slow start". Actualiza cwnd.
- Para el envío
 - ventana = $\min(\text{rcv_win}, \text{cwnd})$

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 8 / 61

Control de Congestión

Slow Start

- ✓ Inicializa el sistema y descubre la congestión rápidamente.
- ✓ Incrementa cwnd hasta la congestión
- ✓ Estima el óptimo cwnd
- ✓ Detecta congestión por pérdida de segmentos
- ✓ Desventajas
 - Detección tardía
 - Enlaces de alta velocidad -- > ventanas mayores -- > mayor pérdida
 - Interacción con el algoritmo de retransmisión

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 9 / 61

Notas:

Control de Congestión

Slow Start ...

- ✓ En el comienzo o después de congestión
 $cwnd = 1$
- ✓ Después de cada segmento validado
 $cwnd < - - cwnd + 1$
- ✓ TCP detiene el crecimiento de cwnd cuando
 $cwnd \geq ssthresh$
- ✓ Pese al incremento unitario resulta exponencial

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 10 / 61

Control de Congestión

Slow Start – Ejemplo

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 11 / 61

Notas:

Estrategia del control de Congestión

- ✓ Incremento aditivo, comenzar con *ssthresh*, incrementar *cwnd* lentamente.
- ✓ Disminución multiplicativa
Cortar la ventana de congestión drásticamente si se detecta una pérdida.
- ✓ Si $cwnd > ssthresh$ entonces, por cada ACK

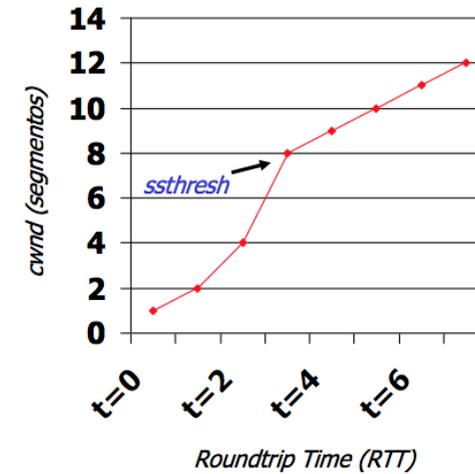
$$cwnd \leftarrow cwnd + \frac{1}{cwnd}$$

- ✓ *cwnd* se incrementa en 1 si todos los segmentos recibieron su validación.



Combinando...

ssthresh = 8



Notas:

Control de Congestión

Detección de paquetes perdidos

- ✓ Esperar RTO
- ✓ RTO es usualmente dos veces RTT
- ✓ Degradación de performance
- ✓ No esperar RTO
 - Utilizar mecanismos alternativos
 - Utilizar RTO si no actúan los anteriores

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 16 / 61

Control de Congestión

Fast Retransmit y Fast Recovery

- ✓ Frente a un segmento fuera de orden se debe enviar un ACK.
- ✓ Provoca duplicación de ACKs.
- ✓ Esta duplicación se ve como debida a:
 - Paquetes perdidos
 - Reordenamiento de paquetes
- ✓ No se puede discriminar
- ✓ Si se reciben 3 ACKs duplicados se considera que se debe a un paquete perdido
- ✓ Al recibir el tercer ACK repetido se retransmite sin esperar el RTO.
- ✓ Eso es Fast Retransmit
- ✓ Luego se ejecuta "congestion avoidanc"?, no slow start.
- ✓ Eso es Fast Recovery.

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 17 / 61

Notas:

Control de Congestión

Integrando...

- ✓ Slow Start.
- ✓ Congestion Avoidance.
- ✓ Si aparecen ACKs duplicados
Fast Retransmit y Fast Recovery.
Congestion Avoidance.
- ✓ Si RTO
Slow Start
- ✓ Resumiendo, TCP Reno.

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 18 / 61

Control de Congestión

TCP Vegas

- ✓ 1994
- ✓ Crecimiento más lento que el slow start
- ✓ Nuevo mecanismo de retransmisión
- ✓ Se chequea el TO al recibir el primer ACK duplicado
- ✓ Nuevo algoritmo de congestion avoidance
- ✓ Evita las oscilaciones de Reno
- ✓ Monitorea la diferencia entre el throughput estimado y el real
- ✓ Trata de reducir a cero los paquetes almacenados en los buffers de los routers

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 19 / 61

Notas:

Control de Congestión

Ventana TCP Vegas

$$w_s(t + 1) = \begin{cases} w_s(t) + \frac{1}{D_s(t)}, & \text{Si } \frac{w_s(t)}{d_s} - \frac{w_s(t)}{D_s(t)} < \alpha_s \\ w_s(t) - \frac{1}{D_s(t)}, & \text{Si } \frac{w_s(t)}{d_s} - \frac{w_s(t)}{D_s(t)} > \beta_s \\ w_s(t), & \alpha_s < \frac{w_s(t)}{d_s} - \frac{w_s(t)}{D_s(t)} < \beta_s \end{cases}$$

$D_s(t) = \text{RTT de la fuente } s$
 $d_s = \text{mínimo RTT de la fuente } s$
 α y β , *parámetros*

Marrone (LINTI-UNLP)
CI
8 de octubre de 2021 20 / 61

Control de Congestión

Throughput TCP Vegas

Marrone (LINTI-UNLP)
CI
8 de octubre de 2021 21 / 61

Notas:

Performance TCP

TCP Persist Timer

- ✓ Es necesario hacer una suerte de polling
- ✓ Pueden producirse "deadlocks"
- ✓ Al recibirse una ventana de 0 se activa el persist timer
- ✓ Normalmente 5 segundos
- ✓ Cuando expira se envía un segmento de 1 byte para verificar el estado del receptor
- ✓ El receptor le contesta acorde con el estado en que se encuentra
- ✓ Si el receptor vuelve a constestar con ACK , wind=0, entonces el persist timer hace un back-off binario (10, 20, 40 seg...)
- ✓ Al contestar el ACK no validando el byte recibido el emisor continúa enviando este byte de prueba
- ✓ La diferencia con el RTO es que en este caso el emisor envía permanentemente esta prueba hasta que la ventana se incremente o se haga un reset de la conexión

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 28 / 61

Performance TCP

Keepalive Timer

- ✓ En TCP si no hay intercambio de datos no hay tráfico alguno, pero la conexión persiste
- ✓ Hasta que haya una caída en algún extremo o reboot de alguno de los hosts
- ✓ Se define para interrogar al otro extremo por su estado
- ✓ No es parte de la especificación de TCP
- ✓ No se recomienda porque puede provocar caídas en caso que la falla sea transitoria
- ✓ También incrementa el uso del ancho de banda disponible como cualquier otra acción de control
- ✓ En algunos casos es necesario. Ej, Telnet

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 29 / 61

Notas:

Performance TCP

Keepalive Timer... I

- ✓ En el extremo en que esté habilitado si no hay actividad en 2 horas envía un segmento de prueba
- ✓ El segmento de prueba es similar al de persist timer
- ✓ Puede también enviarse vacío como un ACK pero con un número de secuencia inesperado
- ✓ El receptor se puede encontrar en alguno de estos estados:
 - Activo:** Responderá al segmento de prueba. El emisor resetea el timer por otras 2 horas. Si en ese intervalo aparece tráfico entonces se vuelve a resetear
 - Caído:** o en proceso de reboot. El emisor no recibirá respuesta y se genera un timeout a los 75 segundos. Repite 10 veces en intervalos de 75 seg. Si no recibe respuesta finaliza la conexión
 - Fin-reboot:** el emisor recibirá una respuesta que será un reset de la conexión

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 30 / 61

Performance TCP

Keepalive Timer... II

Receptor activo pero inalcanzable por el emisor: el emisor no recibe respuesta y genera un timeout de 75 seg al cabo de los cuales retransmite el byte de prueba. Es similar al caso 2

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 31 / 61

Notas:

Performance TCP

"Wrap around"

- ✓ El número de secuencia tiene 32 bits
- ✓ En velocidades altas el espacio de 32 bits puede reciclar dentro del intervalo de tiempo que un segmento es retardado en la red
- ✓ Para conseguir una operación libre de este error:

$$2^{31} / B > MSL(seg)$$

donde B es la velocidad efectiva del enlace en bytes/seg

$$T_{wrap} = 2^{31} / B$$

Red	B (Mbps)	Twrap (seg)
Ethernet	1.25	1700
FDDI	12.5	170
Gigabit	125	17



Marrone (LINTI-UNLP) CI 8 de octubre de 2021 32 / 61

Performance TCP

PAWS

- ✓ Definido para rechazar segmentos duplicados y reencarnaciones
- ✓ Utiliza la opción de timestamp
- ✓ Asume que los segmentos de datos y ACK recibidos contienen un valor de TS monótono no-decreciente
- ✓ Un segmento puede ser descartado como duplicado si se recibe con un valor de TS menor que uno anterior
- ✓ "Menor que" significa que si s y t son valores de TS, entonces

$$s < t \text{ si } 0 < (t - s) < 2^{31}$$

- ✓ Los valores de TS enviados en <SYN> y/o <SYN,ACK> inicializan PAWS
- ✓ No requiere sincronización de relojes entre emisor y receptor

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 33 / 61

Notas:

TCP SACK

SACK

- ✓ Produce la retransmisión selectiva
- ✓ RFC 2018
- ✓ Comprende dos opciones de TCP
 - ✓ "Sack-Permitted Option"
 - ✓ "Sack Option"

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 35 / 61

TCP SACK

Sack Permitted Option

- ✓ Enviada con el SYN
- ✓ Habilita el uso de SACK

```
+-----+-----+  
| Kind=4 | Length=2 |  
+-----+-----+
```

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 36 / 61

Notas:

TCP SACK

Sack Option – Estructura

```
+-----+-----+-----+-----+
|  NOP  |  NOP  | Kind=5 | Length |
+-----+-----+-----+-----+
| Left Edge of 1st Block |
+-----+-----+-----+-----+
| Right Edge of 1st Block |
+-----+-----+-----+-----+
|                               |
/                               /
|                               |
+-----+-----+-----+-----+
| Left Edge of nth Block |
+-----+-----+-----+-----+
| Right Edge of nth Block |
+-----+-----+-----+-----+
```

FACULTAD DE INGENIERIA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 37 / 61

Los NOP se agregan para alinear la opción a palabras de 32 bits (4 bytes).

Left Edge of Block

Es el primer número de secuencia del block.

Right Edge of Block

Es el número de secuencia inmediato siguiente al último número de secuencia del block. Los bytes por debajo del block (Left Edge of Block - 1), y justo por encima del block (Right Edge of Block), no se reciben.

Dado que las opciones tienen una longitud máxima de 40 bytes esto hace que no se puedan indicar más de 4 huecos.

Al recibirse se produce la retransmisión de lo que indica la opción

Notas:

TCP SACK

Ejemplos SACK - rfc 2018

El transmisor envía 8 segmentos de 500 bytes de datos con SN:5000
 Caso 1:
 Se pierde el primer segmento y llegan OK los 7 restantes:

Triggering Segment	ACK	Left Edge	Right Edge
5000	(lost)		
5500	5000	5500	6000
6000	5000	5500	6500
6500	5000	5500	7000
7000	5000	5500	7500
7500	5000	5500	8000
8000	5000	5500	8500
8500	5000	5500	9000



FACULTAD DE INGENIERÍA Facultad de Ingeniería

Marrone (LINTI-UNLP)
CI
8 de octubre de 2021 38 / 61

TCP SACK

Ejemplos SACK - rfc 2018...

El transmisor envía 8 segmentos de 500 bytes de datos con SN:5000
 Caso 2:
 Se pierden los segmentos pares:

Triggering Segment	ACK	1er Bloque		2do Bloque		3er Bloque	
		Left Edge	Right Edge	Left Edge	Right Edge	Left Edge	Right Edge
5000	5500						
5500	(lost)						
6000	5500	6000	6500				
6500	(lost)						
7000	5500	7000	7500	6000	6500		
7500	(lost)						
8000	5500	8000	8500	7000	7500	6000	6500
8500	(lost)						



FACULTAD DE INGENIERÍA Facultad de Ingeniería

Marrone (LINTI-UNLP)
CI
8 de octubre de 2021 39 / 61

Notas:

TCP SACK

Ejemplos SACK - rfc 2018...

Continuando con el caso anterior, supongamos que se recibe el cuarto paquete a continuación. Entonces:

Triggering Segment	ACK	1er Bloque		2do Bloque		3er Bloque	
		Left	Right	Left	Right	Left	Right
6500	5500	Edge	Edge	Edge	Edge	Edge	Edge
		6000	7500	8000	8500		

Si luego llega el segundo segmento:

Triggering Segment	ACK	1er Bloque		2do Bloque		3er Bloque	
		Left	Right	Left	Right	Left	Right
5500	7500	Edge	Edge	Edge	Edge	Edge	Edge
		8000	8500				


 FACULTAD DE INGENIERÍA
 Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 40 / 61

ECN

Generalidades

- ✓ Retomar el modelo basado en la red para control de congestión (FR, ATM)
- ✓ Participan IP y TCP
- ✓ Colabora con RED
- ✓ RED procede al marcado de IP
- ✓ TCP interpreta y completa la acción de control


 FACULTAD DE INGENIERÍA
 Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 42 / 61

Notas:

ECN

ECN – Secuencia de Eventos

- ✓ Se activa ECT "codepoint" en paquetes transmitidos por el emisor indicando que se soporta ECN.
- ✓ Un router ECN detecta congestión y detecta que el ECT "codepoint" está activo en el paquete que está pronto a marcar (antes lo descartaba).
- ✓ El router activa el CE y forwarda el paquete.

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 47 / 61

ECN

ECN – Secuencia de Eventos

- ✓ El receptor recibe el paquete con CE activo y activa el flag ECN-E en el header del próximo paquete TCP que enviará al emisor.
- ✓ El emisor recibe el paquete TCP con el bit de Flag ECN-Echo activo y actúa como si hubiera detectado un paquete perdido. Acorde con su mecanismo de control de congestión.
- ✓ El emisor activa el Flag de CWR en el header TCP del próximo paquete que envía al receptor.

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 48 / 61

Notas:

AQM

Características

- RFC 2309 – Reemplazada por RFC 7567(2015)
- Modelo de control de congestión basado en la red
- Controlar la longitud de la cola de salida de los datagramas
- Monitorear el estado de la cola y calcular la longitud promedio
- Marcar los datagramas según esa longitud promedio
- Completa el mecanismo ECN.
- Propone diversos algoritmos
- Ejemplo: RED (Random Early Detection)

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 52 / 61

AQM

RED – Algoritmo

Dos algoritmos

- 1 Longitud promedio de la cola
 - Filtro pasabajos con promedio ponderado exponencial
 - Determina el índice de rafagosidad que se admitirá en el "buffer "
- 2 Marcado de paquetes
 - Determina la frecuencia de marcado de paquetes en función del nivel de congestión en la red
 - Se compara la longitud promedio con dos umbrales
 - Si está por debajo del menor no se marcan
 - Si está entre ambos se los marca con una probabilidad p_a
 - Si está por encima del mayor se lo marca
 - El marcado puede ser sólo eso o el descarte.

FACULTAD DE INFORMATICA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 53 / 61

Notas:

AQM

RED – Algoritmo – Longitud promedio

Initialization:

$$avg \leftarrow 0$$
$$count \leftarrow -1$$

for each packet arrival
calculate the new average queue size avg :
if the queue is nonempty

$$avg \leftarrow (1 - w_q) avg + w_q q$$

else

$$m \leftarrow f(time - q_{time})$$
$$avg \leftarrow (1 - w_q)^m avg$$


FACULTAD DE INGENIERIA
Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 54 / 61

AQM

RED – Algoritmo – Mercado de paquetes

if $min_{th} \leq avg < max_{th}$
increment count
calculate probability p_a :

$$p_b \leftarrow max_p (avg - min_{th}) / (max_{th} - min_{th})$$
$$p_a \leftarrow p_b / (1 - count \times p_b)$$

with probability p_a :
mark the arriving packet

$$count \leftarrow 0$$

else if $max_{th} \leq avg$
mark the arriving packet

$$count \leftarrow 0$$

else $count \leftarrow -1$
when queue becomes empty

$$q_{time} \leftarrow time$$


FACULTAD DE INGENIERIA
Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 55 / 61

Notas:

Variaciones del mercado

- Medir la longitud de la cola en bytes en vez de paquetes
- La longitud refleja el retardo promedio en el router
- Se deben hacer cambios:

$$p_b \leftarrow \max_p (avg - min_{th}) / (max_{th} - min_{th})$$

$$p_b \leftarrow p_b \text{ PacketSize} / \text{MaximumPacketSize}$$

$$p_a \leftarrow p_b / (1 - count \times p_b)$$

- Es más probable que se marque un paquete FTP que uno de Telenet



RED-Algoritmo-Parámetros

Variables:

avg : longitud promedio de la cola

q_{time} : Comienzo del tiempo ocioso de la cola

$count$: Paquetes desde la última marcación

Constantes:

w_q : peso de la cola

min_{th} : Umbral mínimo

max_{th} : Umbral máximo de la cola

max_p : valor máximo para p_b

Otros:

p_a : probabilidad de marcado

q : longitud actual de la cola

$time$: tiempo actual

$f(t)$: función lineal del tiempo t



Notas:

AQM

RED - Implementación

- ✓ Si *avg* está dentro de las cotas se debe marcar con probabilidad
- ✓ ¿Qué paquete se marca?
- ✓ $R = \text{Random}[0, 1]$
- ✓ Se marca si:
$$R < \frac{p_b}{1 - \text{count} \times p_b}$$

FACULTAD DE INGENIERÍA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 60 / 61

AQM

CC BY

Atribución-NoComercial-CompartirIgual
4.0 Internacional (CC BY-NC-SA 4.0)

Esta obra está sujeta a la licencia Atribución-NoComercial-CompartirIgual 4.0 Internacional (CC BY-NC-SA 4.0) de Creative Commons.

Para detalle de esta licencia visite
<https://creativecommons.org/licenses/by-nc-sa/4.0/>

FACULTAD DE INGENIERÍA Facultad de Ingeniería

Marrone (LINTI-UNLP) CI 8 de octubre de 2021 61 / 61

Notas:
