# The dark side of TCP

## understanding TCP on very high-speed networks

**ACOMP 2008**

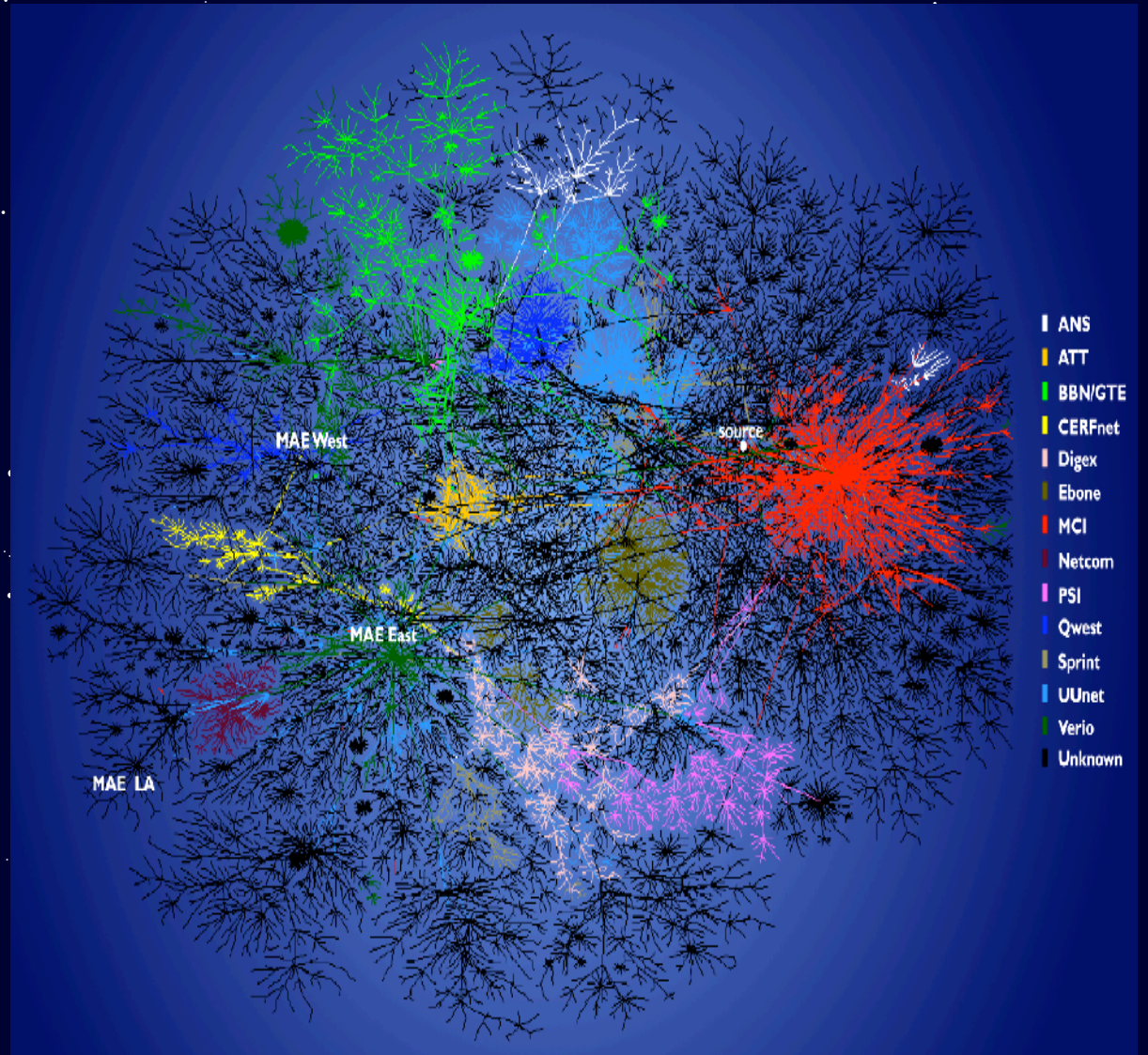HCMC, Bach Khoa University

March 11th, 2008
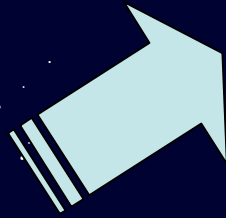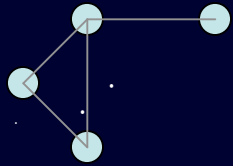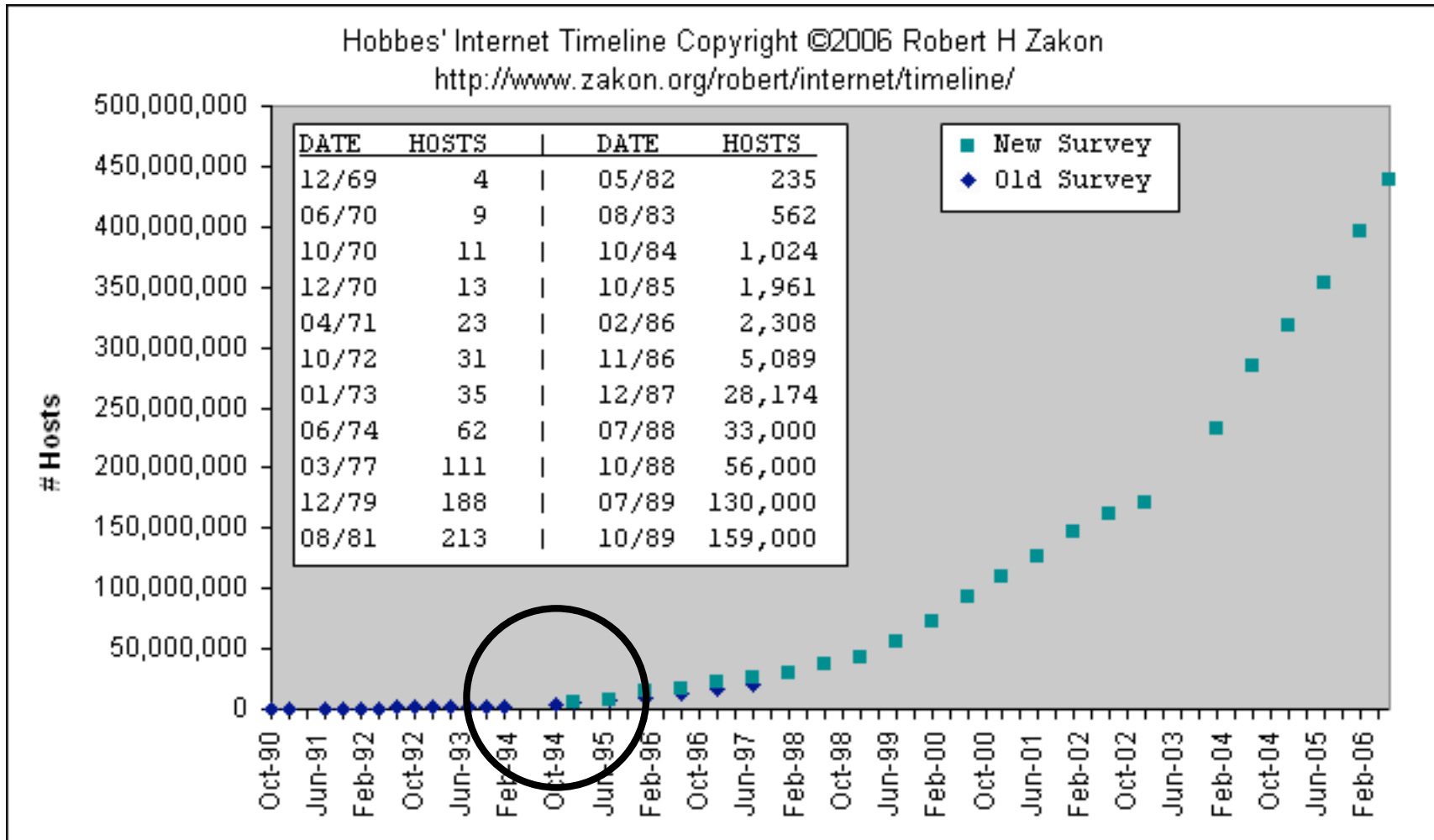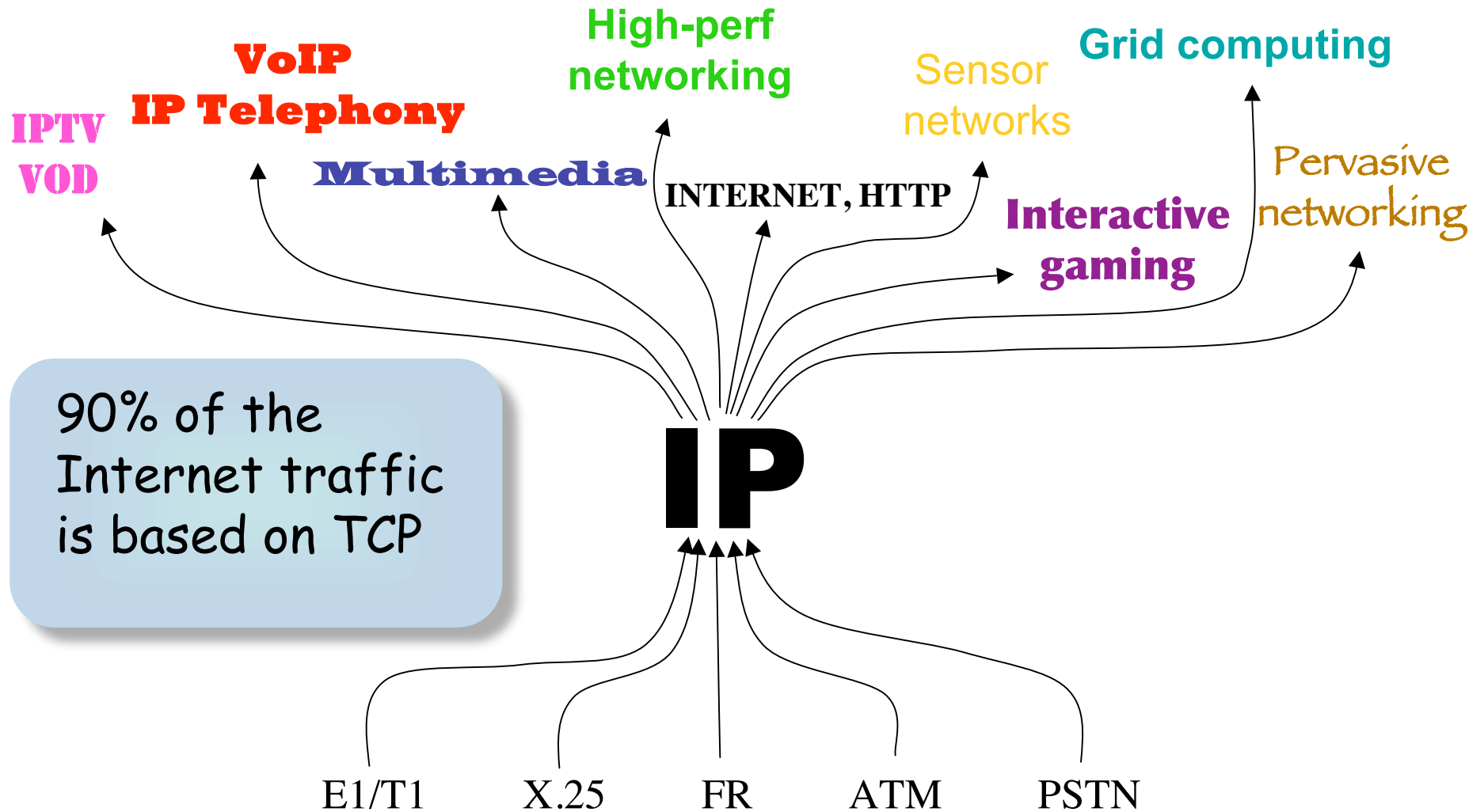
C. Pham

http://www.univ-pau.fr/~cpham

University of Pau, France

LIUPPA laboratory

# The big-bang of the Internet

# # Internet host



Hobbes' Internet Timeline Copyright ©2006 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

| DATE | HOSTS | | DATE | HOSTS |
|---|---|---|---|---|
| 12/69 | 4 | | 05/82 | 235 |
| 06/70 | 9 | | 08/83 | 562 |
| 10/70 | 11 | | 10/84 | 1,024 |
| 12/70 | 13 | | 10/85 | 1,961 |
| 04/71 | 23 | | 02/86 | 2,308 |
| 10/72 | 31 | | 11/86 | 5,089 |
| 01/73 | 35 | | 12/87 | 28,174 |
| 06/74 | 62 | | 07/88 | 33,000 |
| 03/77 | 111 | | 10/88 | 56,000 |
| 12/79 | 188 | | 07/89 | 130,000 |
| 08/81 | 213 | | 10/89 | 159,000 |

# Towards all IP & TCP!

High-perf networking

VoIP
IP Telephony

Grid computing

IPTV
VOD

Multimedia

Sensor networks

INTERNET, HTTP

Pervasive networking

Interactive gaming

90% of the Internet traffic is based on TCP

**IP**

E1/T1     X.25     FR     ATM     PSTN

# What TCP brings

- [ ] stream-based
  - [ ] segments a ... umbers
  - [ ] only consec...
- [ ] reliability
  - [ ] missing seg... g) and retransmitted
- [ ] flow control
  - [ ] receiver is ... y based)
- [ ] congestion co...
  - [ ] network is ... based)
  - [ ] protocol tr...
- [ ] connection ha...
  - [ ] explicit est...
- [ ] full-duplex co...
  - [ ] an ACK can be a data segment at the same time (piggybacking)

# A brief history of TCP

1975
**Three-way handshake**
*Raymond Tomlinson*
In SIGCOMM 75

1974
**TCP** described by
*Vint Cerf* and *Bob Kahn*
In IEEE Trans Comm

1982
**TCP & IP**
RFC 793 & 791

1983
**BSD Unix 4.2**
supports TCP/IP

1984
**Nagel's algorithm**
to reduce overhead
of small packets;
predicts congestion
collapse

1986
**Congestion
collapse**
observed

1987
**Karn's algorithm**
to better estimate
round-trip time

1988
**Van Jacobson's
algorithms**
congestion avoidance
and congestion control
(*most* implemented in
**4.3BSD Tahoe**)

1990
**4.3BSD Reno**
fast retransmit
delayed ACK's

1975        1980                    1985                              1990

# ...in the nineties

1994
**T/TCP**
(Braden)
Transaction
TCP

1996
**SACK TCP**
(Floyd et al)
Selective
Acknowledgement

1993
**TCP Vegas**
(Brakmo et al)
real congestion
*avoidance*

1994
**ECN**
(Floyd)
Explicit
Congestion
Notification

1996
**Hoe**
Improving TCP
startup

1996
**FACK TCP**
(Mathis et al)
extension to SACK

?

1993    1994    1996

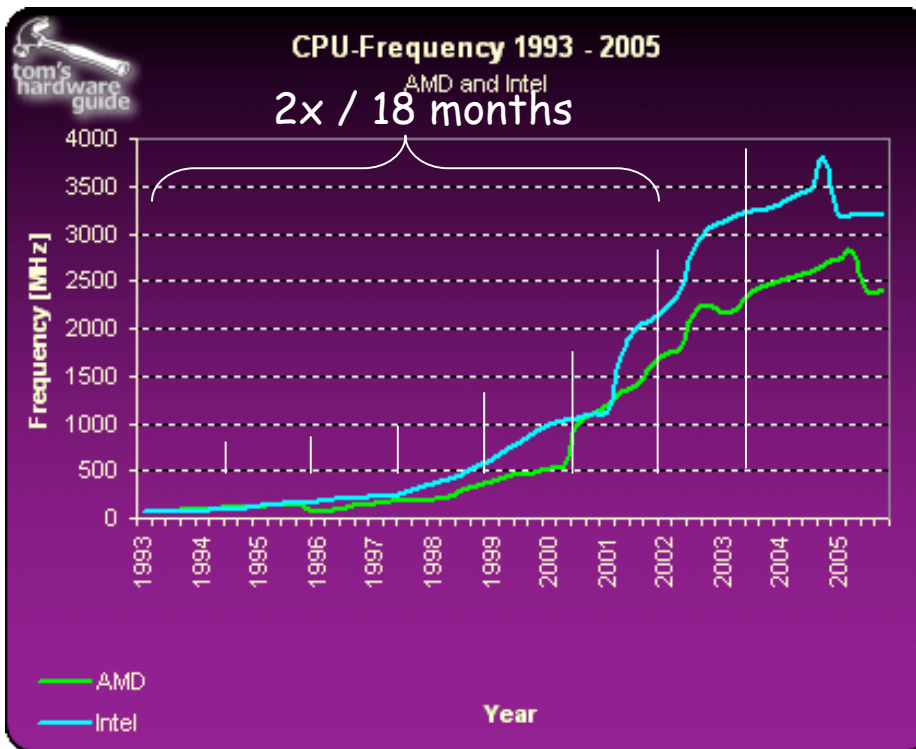# huge variety of communicating devices!

Internet

Wireless sensor nodes

# 1st revolution: Wireless Networks

- ❑ WiFi, WiMax
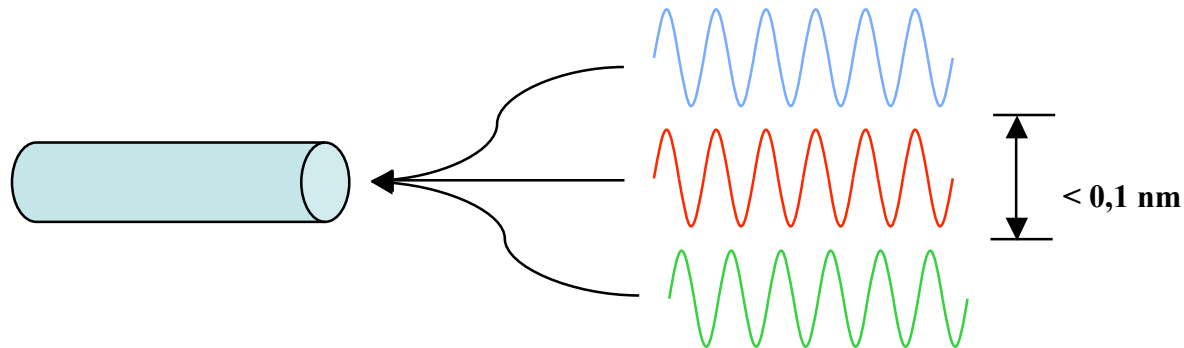- ❑ BlueTooth, ZigBee, IrDA...
- ❑ GSM, GPRS, EDGE, UMTS, 4G,...

Access Point

# 2nd revolution: going optical



CPU-Frequency 1993 - 2005
AMD and Intel
2x / 18 months

2x / 7 months

Source « Optical fibers for Ultra-Large Capacity
Transmission » by J. Grochocinski

# DWDM, bandwidth for free?

DWDM: Dense Wavelength Division Multiplexing

< 0,1 nm

2Gbps

10Gbps

10, 40, 160 Gbps are available!

From Computer Desktop Encyclopedia
Reproduced with permission.
© 2001 Metromedia Fiber Network

# Fibers everywhere?

**NEWS of Dec 15th, 2004**

Verizon and SBC are deploying large optical fiber infrastructures in the US using FTTC or FTTP scenario

**NEWS from Japan and South Korea**

the first echnology users at the % of high-rket). In TH users 1 %.

**NEWS of May 31st, 2005**

US Fiber-to-the-home (FTTH) installations have grown 83% since October 2004, now reaching 398 n 43 states

rack to pass homes with nd of 2005

**NEWS of July, 2006**

France Telecom will deploy an FTTH test-bed infrastructure in Paris. 2.5 Gbps in download and 1.2Gbps in upload!

Gbps

2.5Gbp

campus

GigaEth

ore

160 Gbps

# SONET/SDH in the core
## 95% of exploited OF use SONET/SDH



**Digital switch**

**Digital switch**

**n*30*64 Kb/s**

**n*2048 Kb/s**

*Optical Fiber or Microwave Link*

**MUX PDH/SDH**

**MUX PDH/SDH**

SDH :
| | | |
|---|---|---|
| STM-1 | : | 155.520 Mb/s |
| STM-4 | : | 622.080 Mb/s |
| STM-16 | : | 2488.320 Mb/s |

# SONET/SDH transport network infrastructure

Add Drop Multiplexer

DCS or ADM

DCS or ADM

**rings**

DCS or ADM

DCS or ADM

**rings**

DCS or ADM

DCS or ADM

**SONET/SDH now offers**
Native Ethernet interface
Generic Framing Procedure
Virtual Concatenation

# The new networks

- vBNS
- Abilene
- SUPERNET
- DREN
- CA*NET
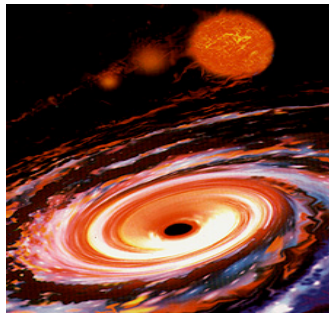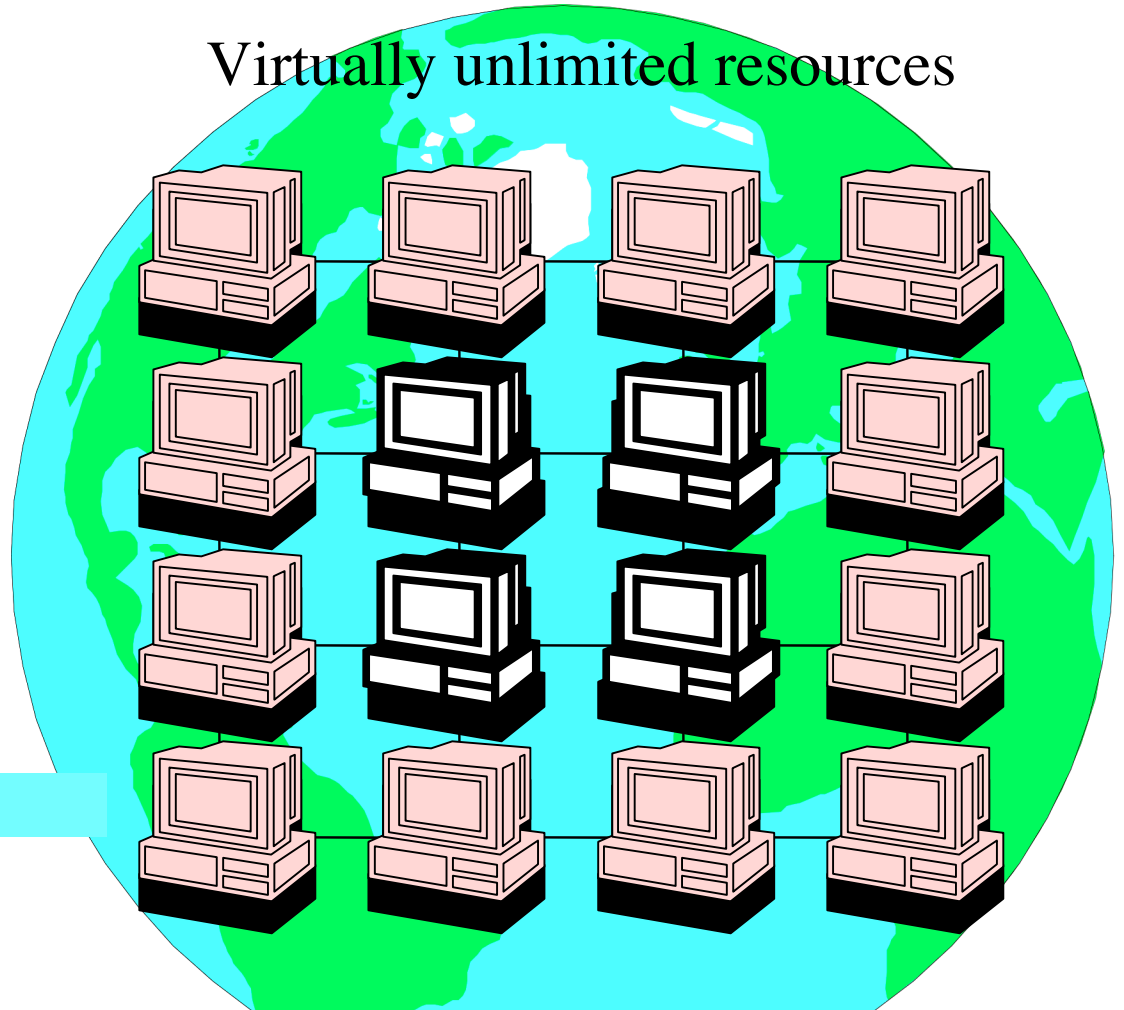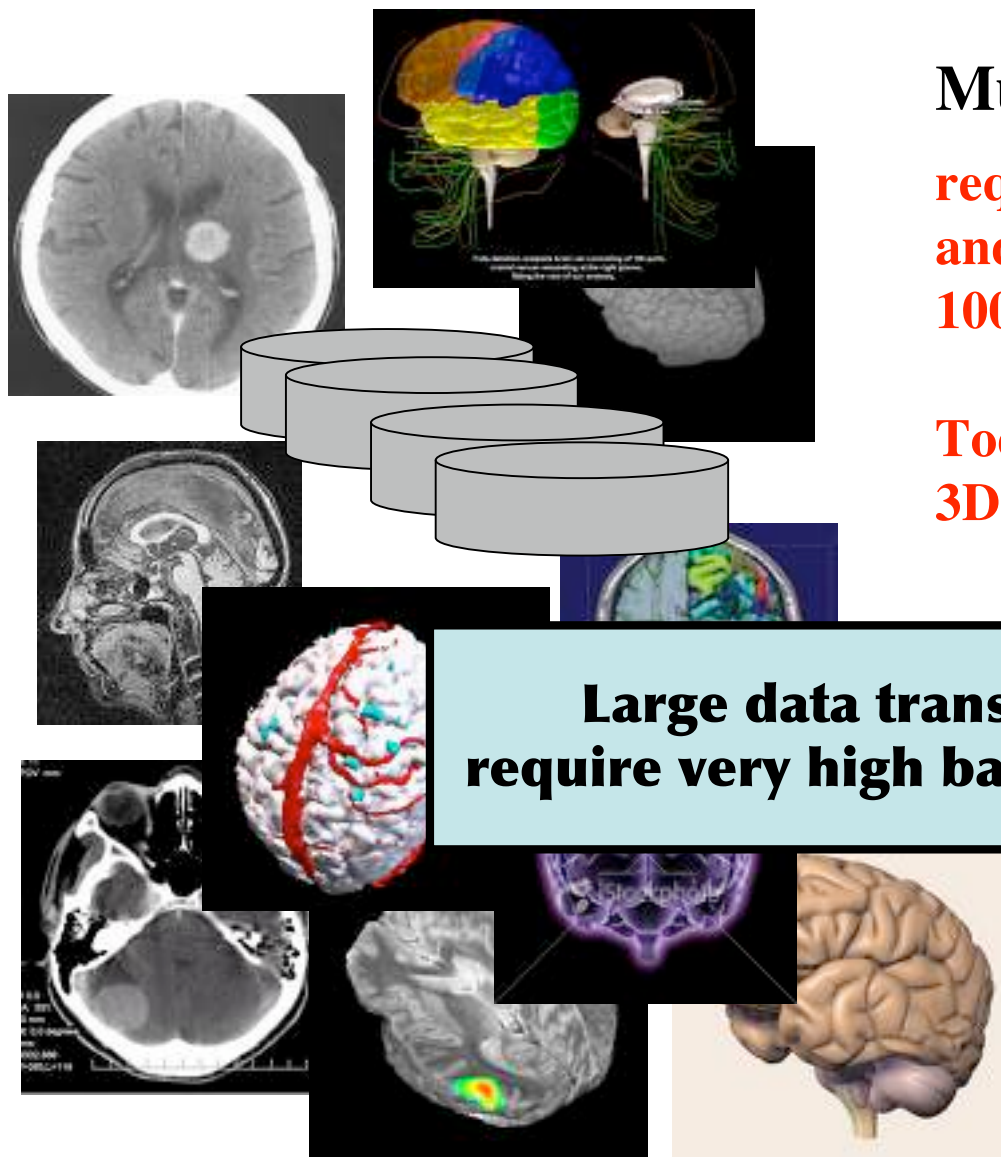- GEANT
- DATATAG
- …much more to come!

GEANT

Backbone Topology April 2004

# Computational grids



user application

1PFlops

Virtually unlimited resources

from Dorian Arnold: Netsolve Happenings

# Real-time interactive large-scale scientific collaborations

**Multimodality brain mapping**

require the ability to process, share, and interactively visualize multiple 100Gbytes datasets!

Today, to visualize and explore eight 3D images require 64Gb/s !

**Large data transfers require very high bandwidth**

# Very High-Speed Networks

**Optical fiber
40 Gbps**

200000km/s, delay of 5ms every 1000km
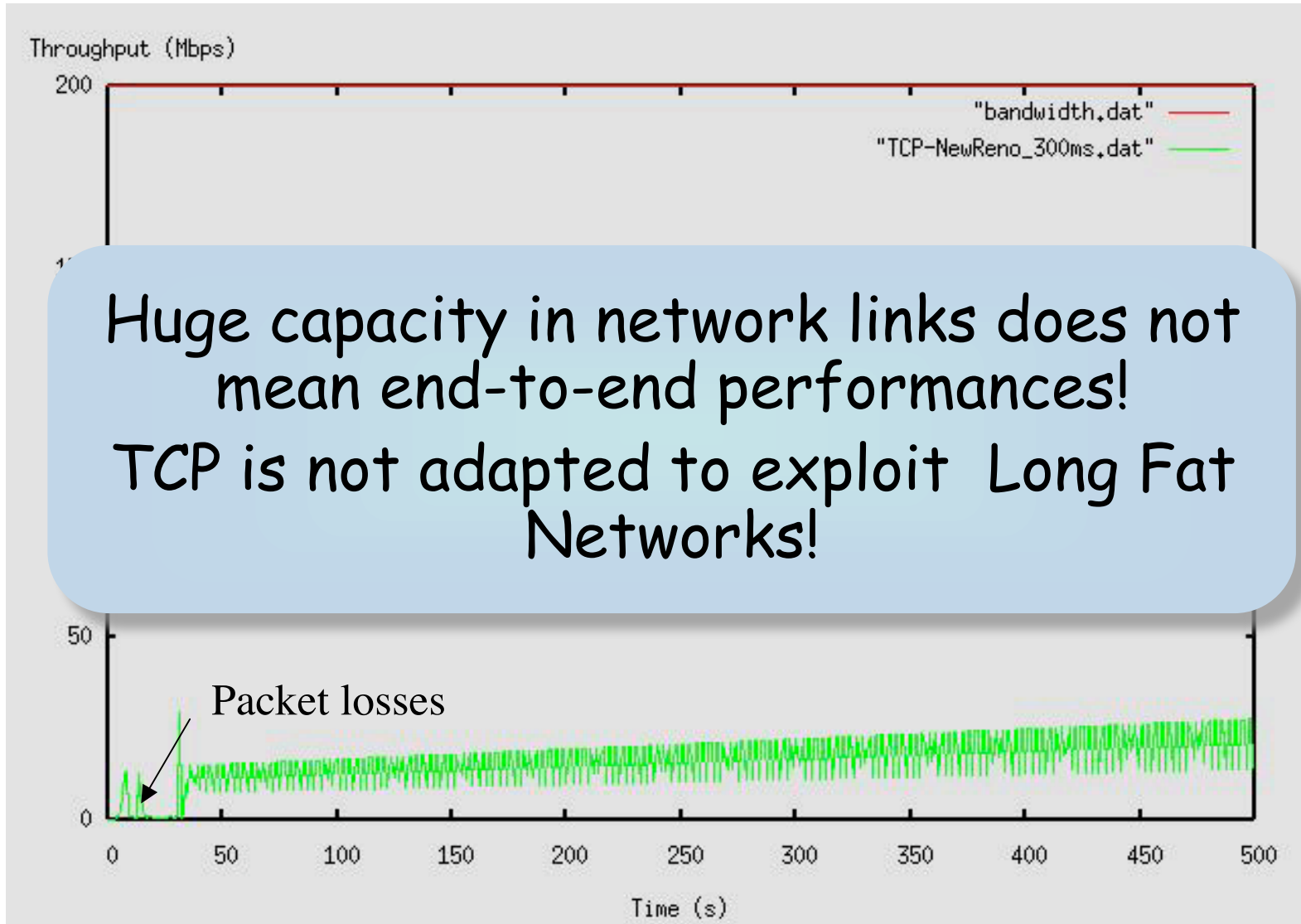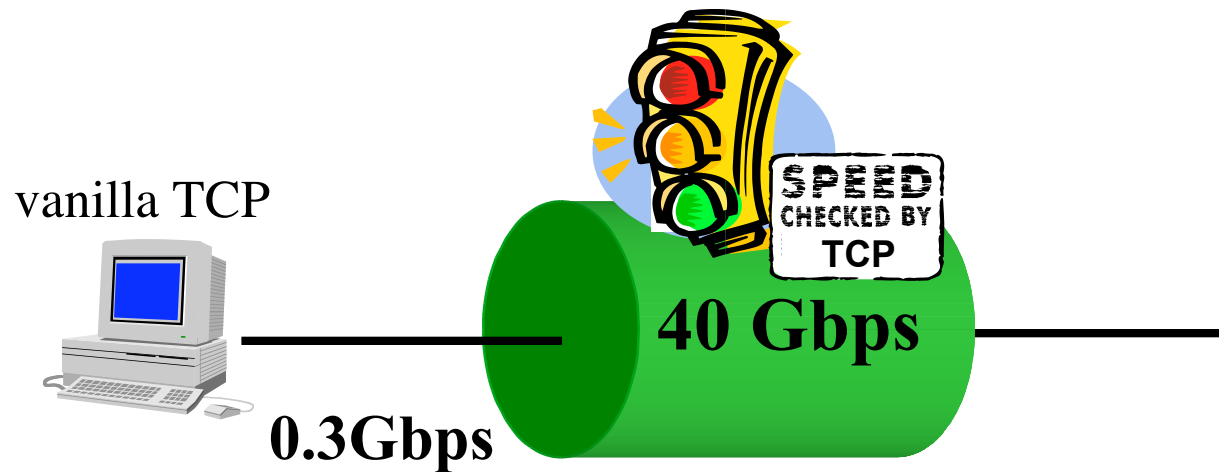
❑ Today's backbone links are optical, DWDM-based, and offer gigabit rates

❑ Transmission time <<< propagation time

❑ Duplicating a 10GB database should not be a problem anymore
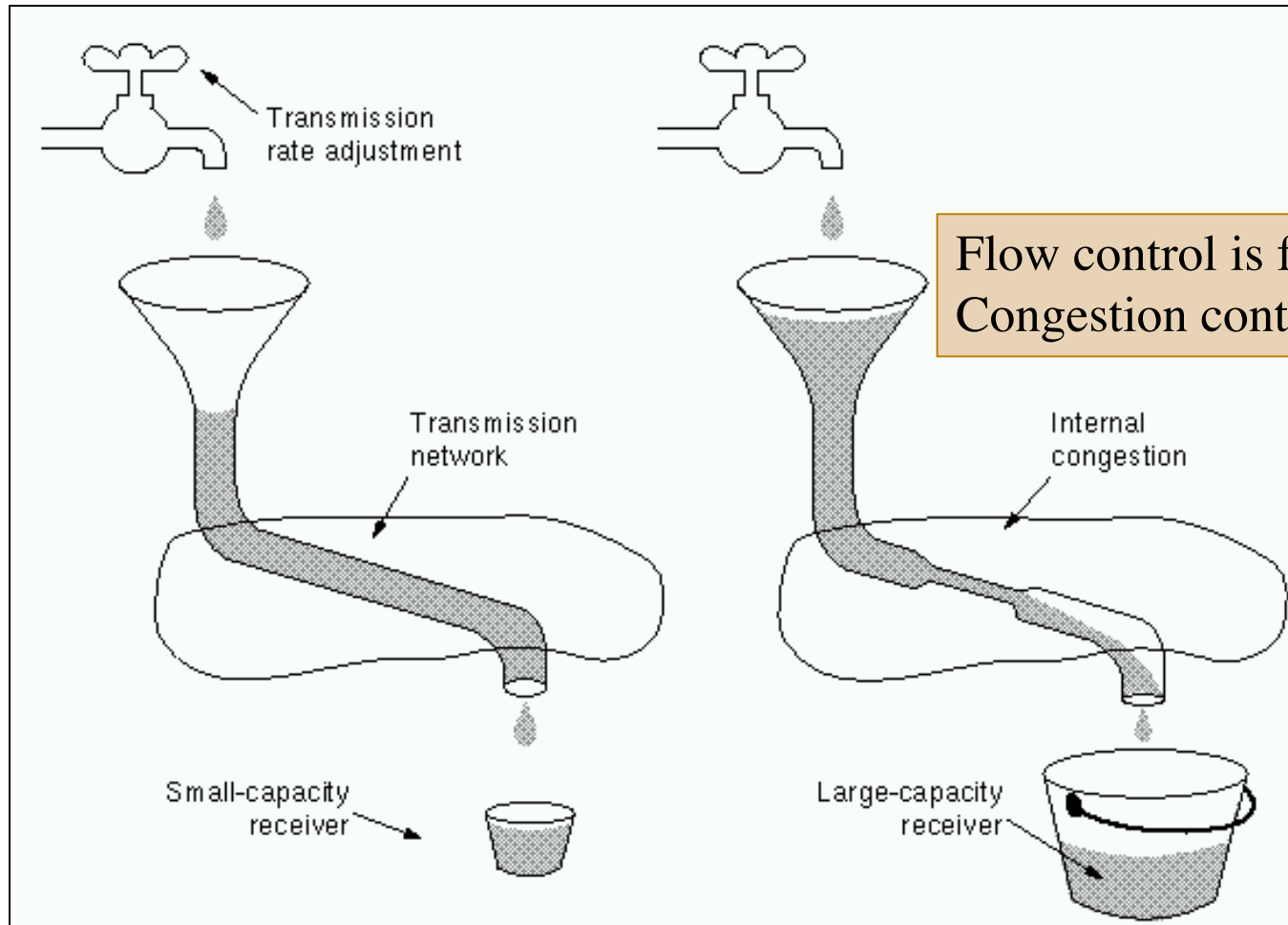
# The reality check: TCP on a 200Mbps link



Throughput (Mbps)

200

"bandwidth.dat"
"TCP-NewReno_300ms.dat"

Huge capacity in network links does not mean end-to-end performances!

TCP is not adapted to exploit  Long Fat Networks!

50

Packet losses

0

0   50   100   150   200   250   300   350   400   450   500

Time (s)

# The things about TCP your mother never told you!

vanilla TCP

0.3Gbps

40 Gbps

SPEED CHECKED BY TCP

❑ If you want to transfer a 1Go file with a standard TCP stack, you will need minutes even with a 40Gbps (how much in $?) link!
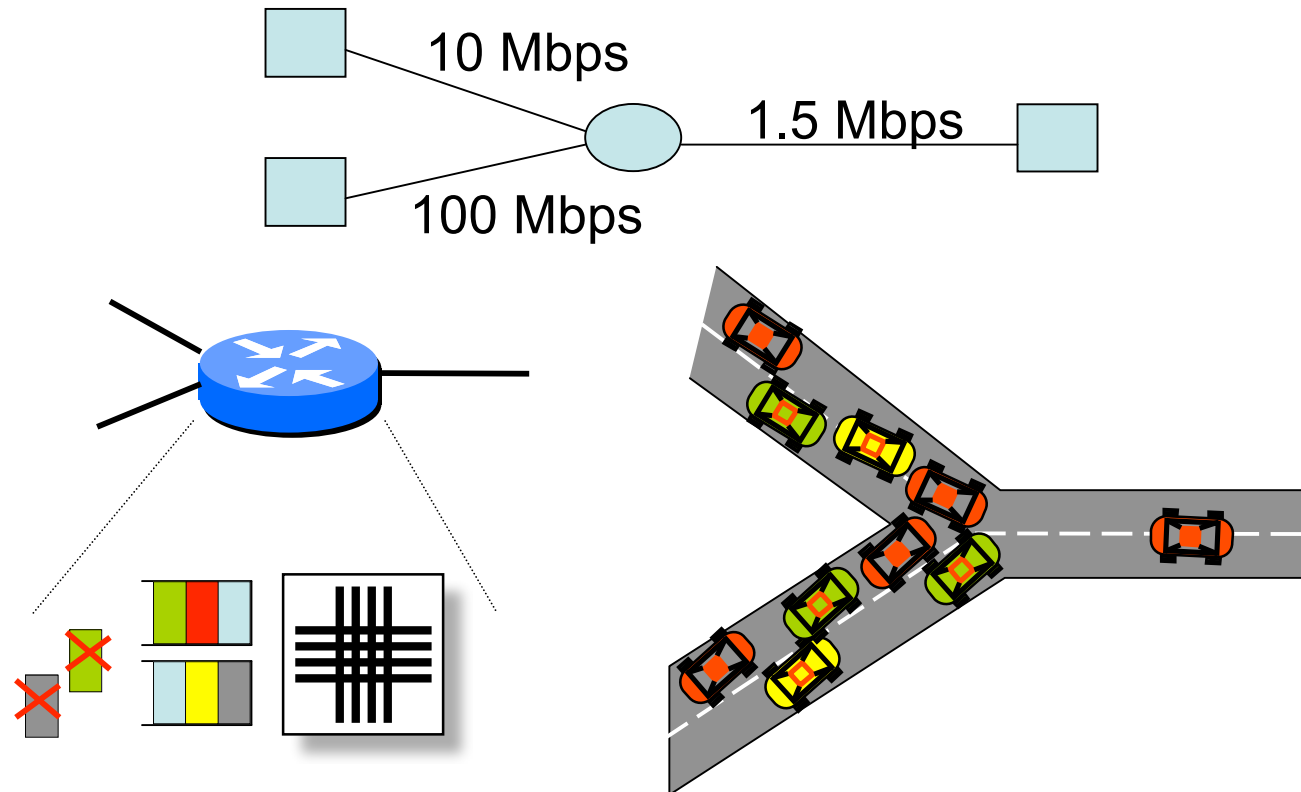
# Let's go back to the origin!



Flow control is for receivers
Congestion control is for the network

Congestion collapse was first observed in 1986 by V. Jacobson. Congestion control was added to TCP (TCP Reno) in 1988.

From Computer Networks, A. Tanenbaum
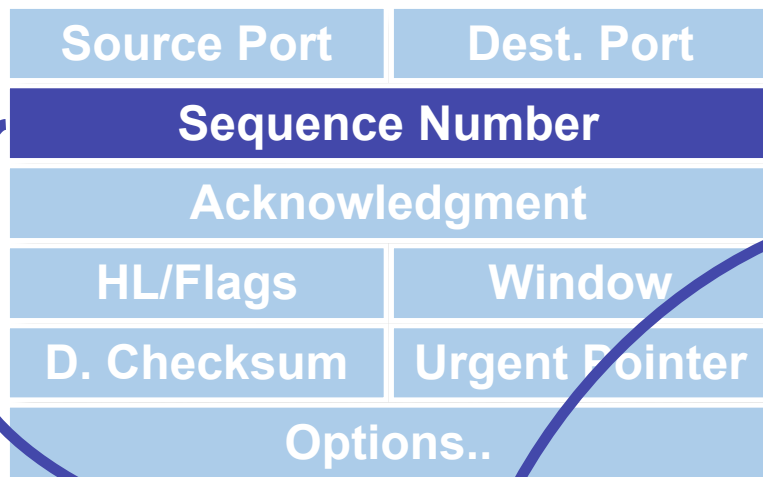
# The congestion phenomenon



☐ Too many packets sent to the same interface.

☐ Difference bandwidth from one network to another
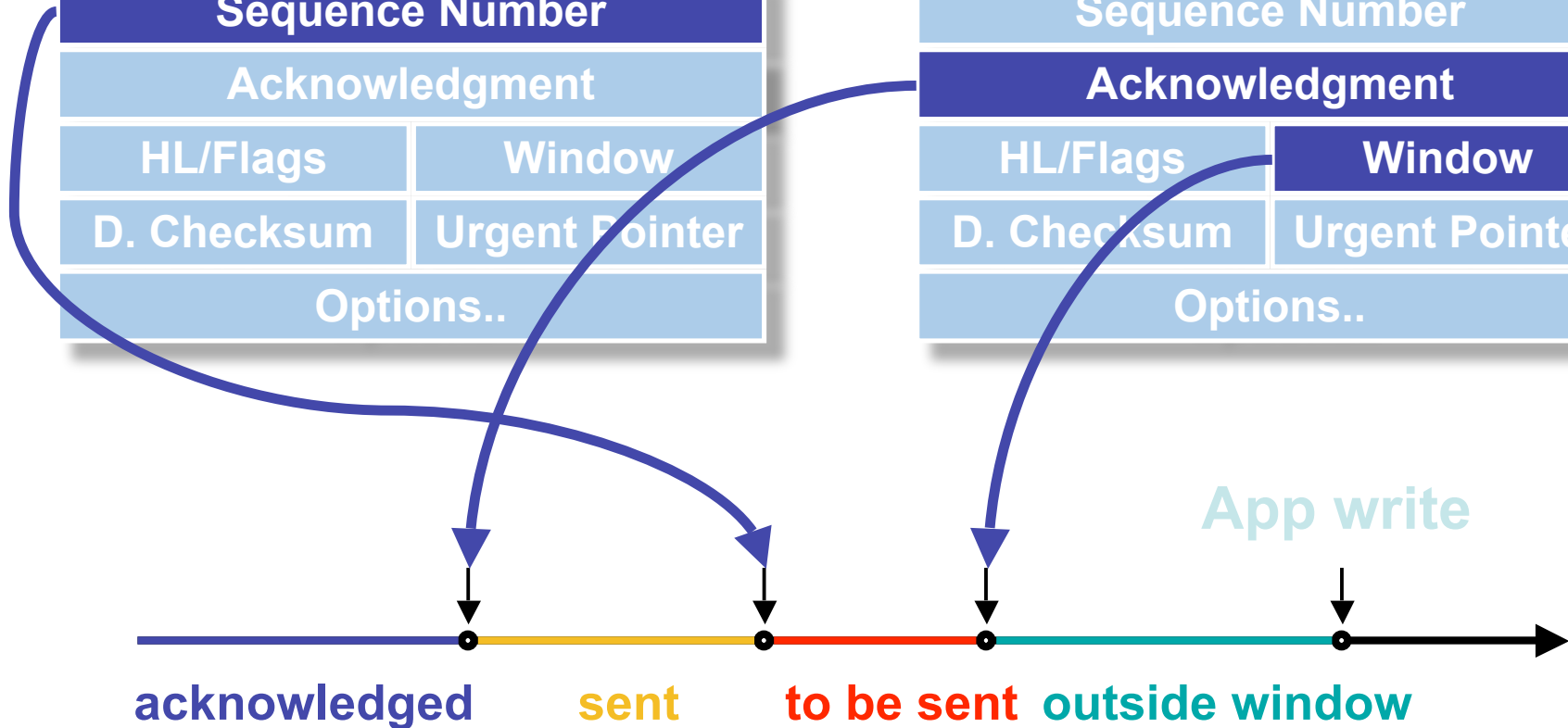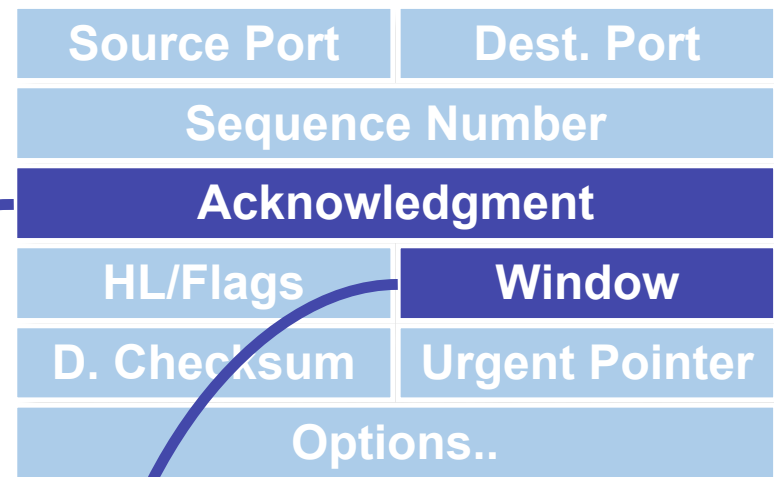
## Main consequence: packet losses in routers
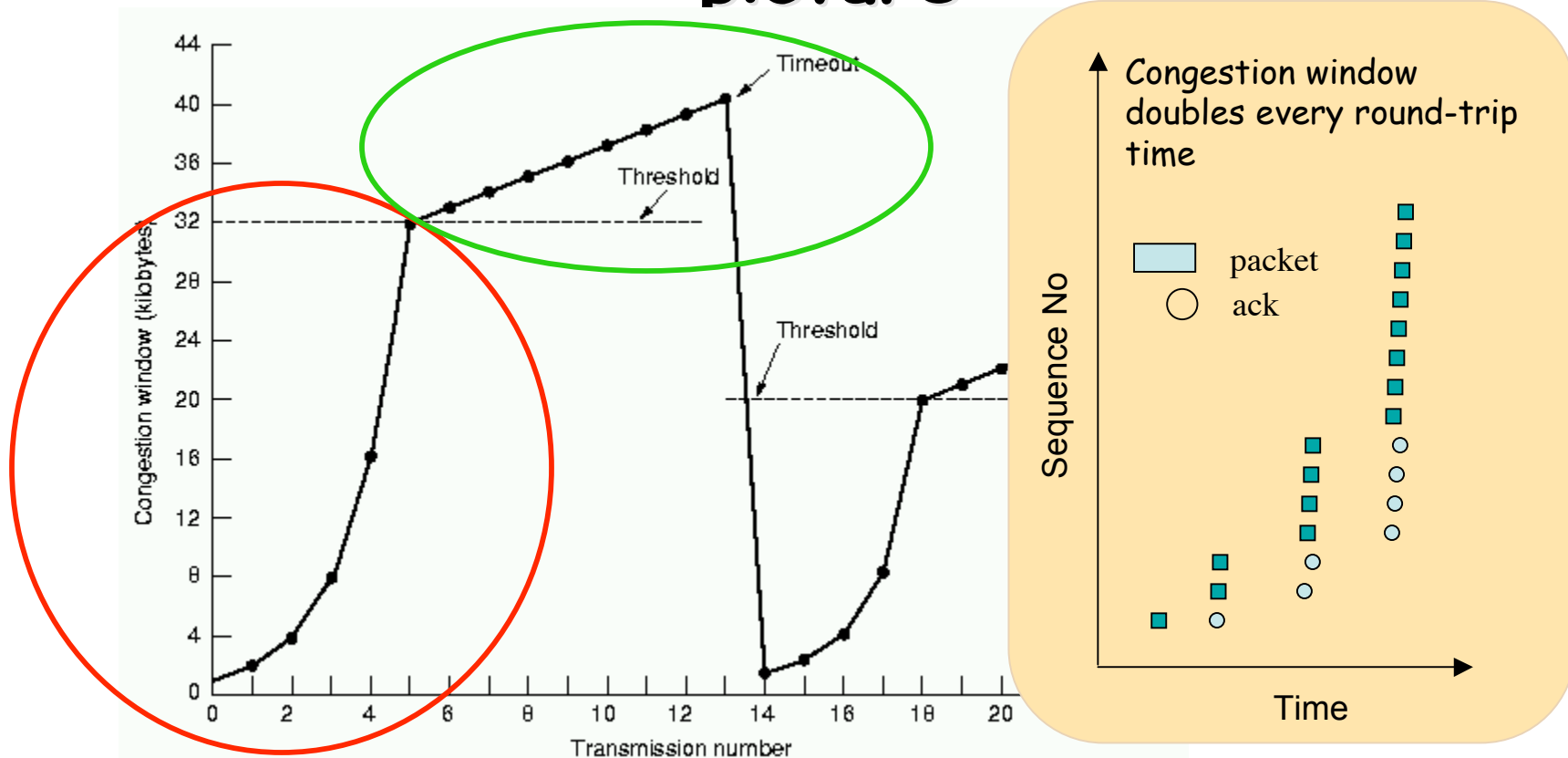
# Flow control
## prevents receiver's buffer overfow

**Packet Sent**

| Source Port | Dest. Port |
|---|---|
| Sequence Number | |
| Acknowledgment | |
| HL/Flags | Window |
| D. Checksum | Urgent Pointer |
| Options.. | |

**Packet Received**

| Source Port | Dest. Port |
|---|---|
| Sequence Number | |
| Acknowledgment | |
| HL/Flags | Window |
| D. Checksum | Urgent Pointer |
| Options.. | |

App write

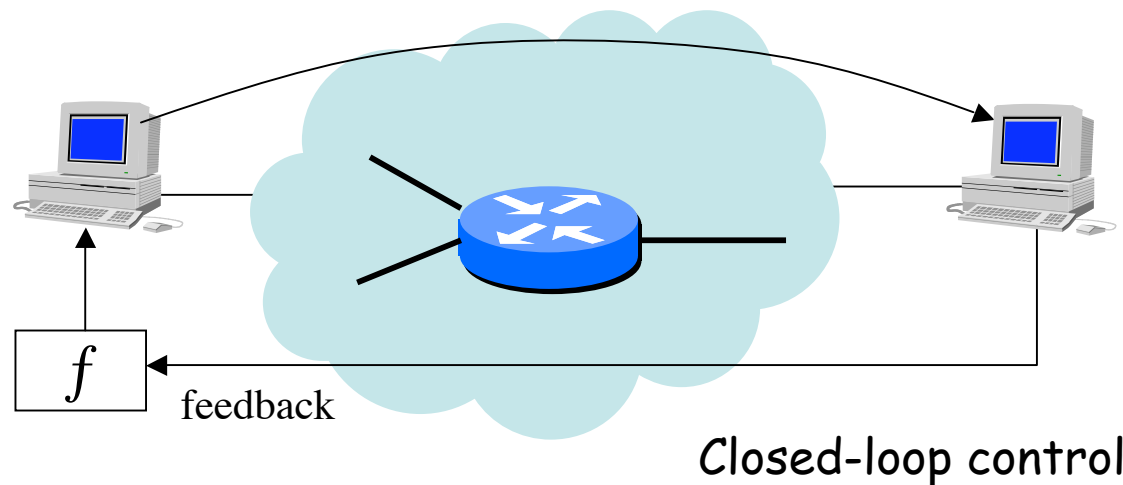**acknowledged**     **sent**     **to be sent**   **outside window**

# TCP congestion control: the big picture



- cwnd grows exponentially (slow start), then linearly (*congestion avoidance*) with 1 more segment per RTT
- If loss, divides threshold by 2 (multiplicative decrease) and restart with cwnd=1 packet

# From the control theory point of view



$f$

feedback

Closed-loop control

❑ Feedback should be frequent, but not too much otherwise there will be oscillations

❑ Can not control the behavior with a time granularity less than the feedback period

# Congestion: A Close-up View
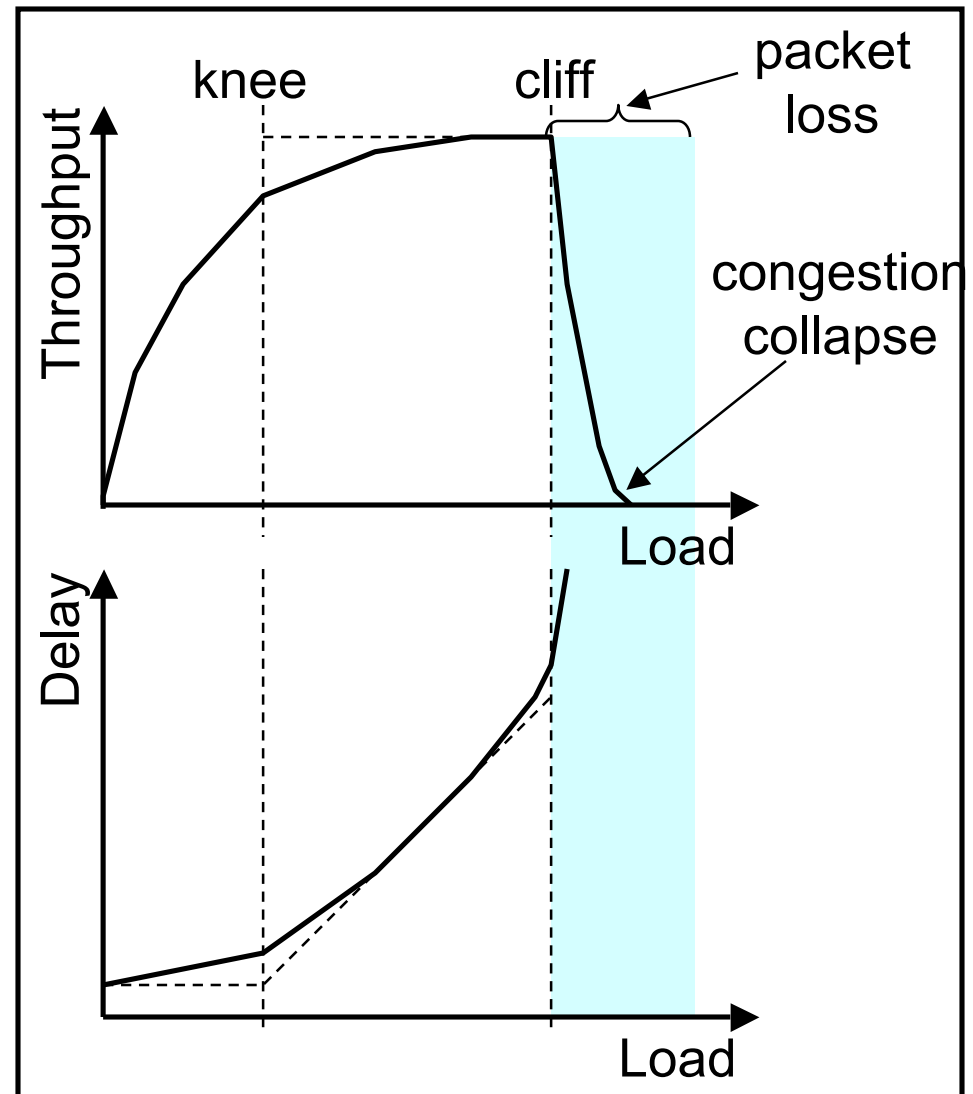
- ❏ knee – point after which
    - ❏ throughput increases very slowly
    - ❏ delay increases fast
- ❏ cliff – point after which
    - ❏ throughput starts to decrease very fast to zero (congestion collapse)
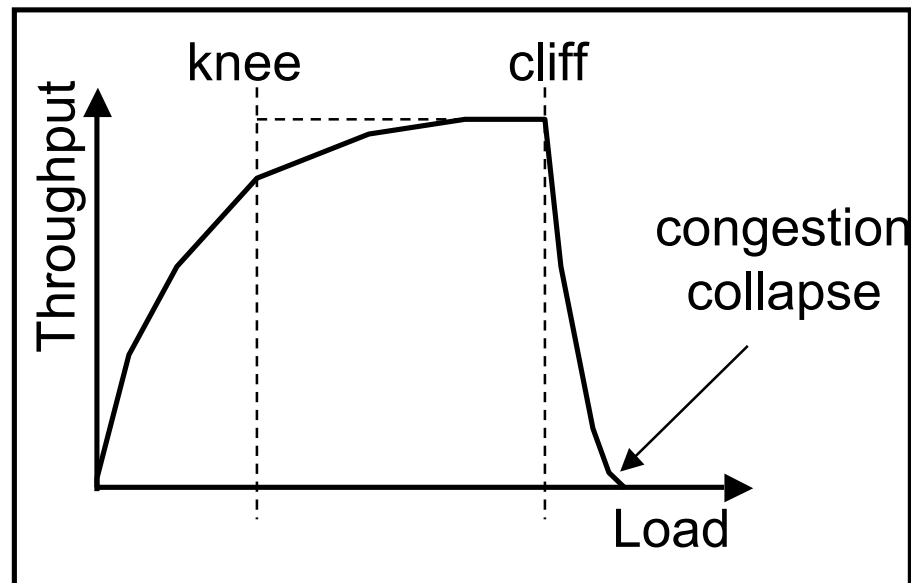    - ❏ delay approaches infinity
- ❏ Note (in an M/M/1 queue)
    - ❏ delay = $1/(1 - \text{utilization})$

# Congestion Control vs. Congestion Avoidance

❑ Congestion control goal
  ❑ stay left of cliff
❑ Congestion avoidance goal
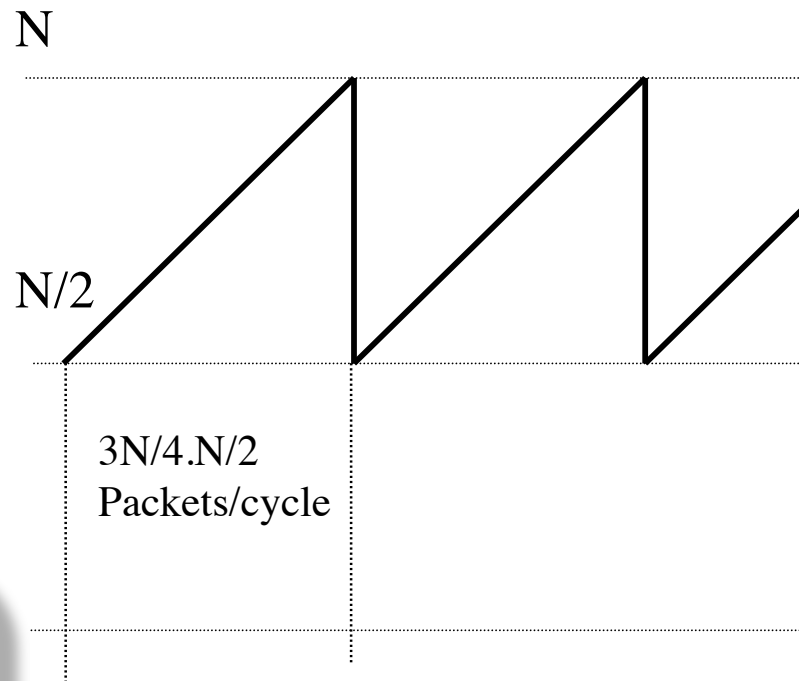  ❑ stay left of knee
❑ Right of cliff:
  ❑ Congestion collapse

# The TCP saw-tooth curve

**TCP behavior in steady state**

Isolated packet losses trigger the fast recovery procedure instead of the slow-start.

N

N/2

3N/4.N/2
Packets/cycle

☐ The TCP steady-state behavior is referred to as the Additive Increase-Multiplicative Decrease process
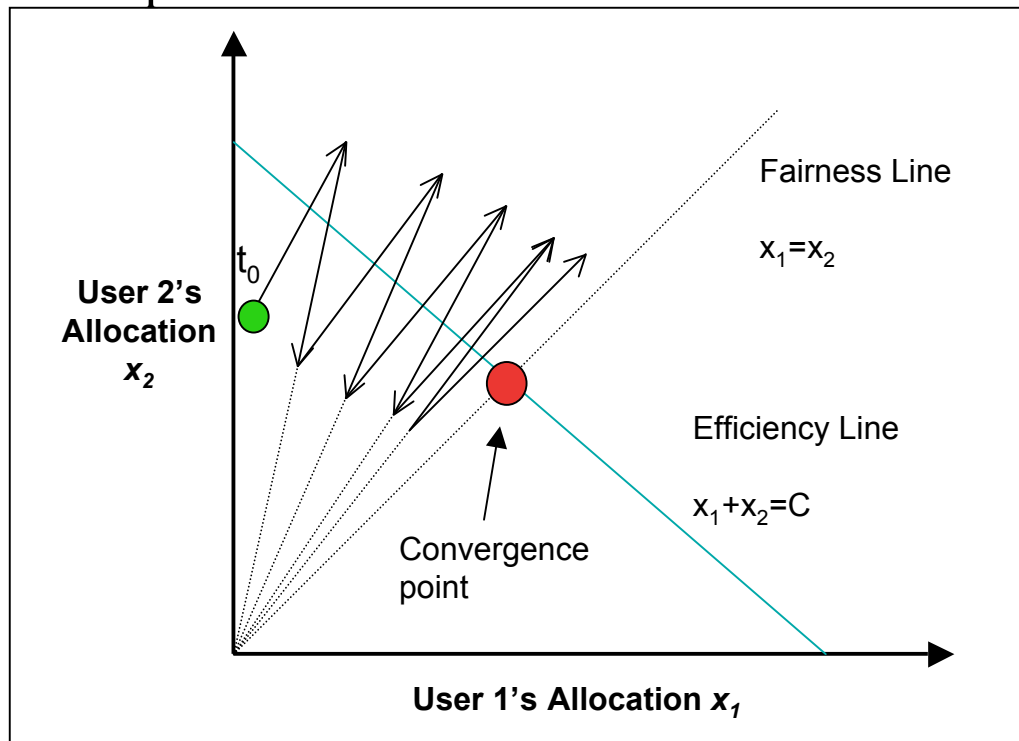
no loss:
cwnd = cwnd + 1
loss:
cwnd = cwnd*0.5
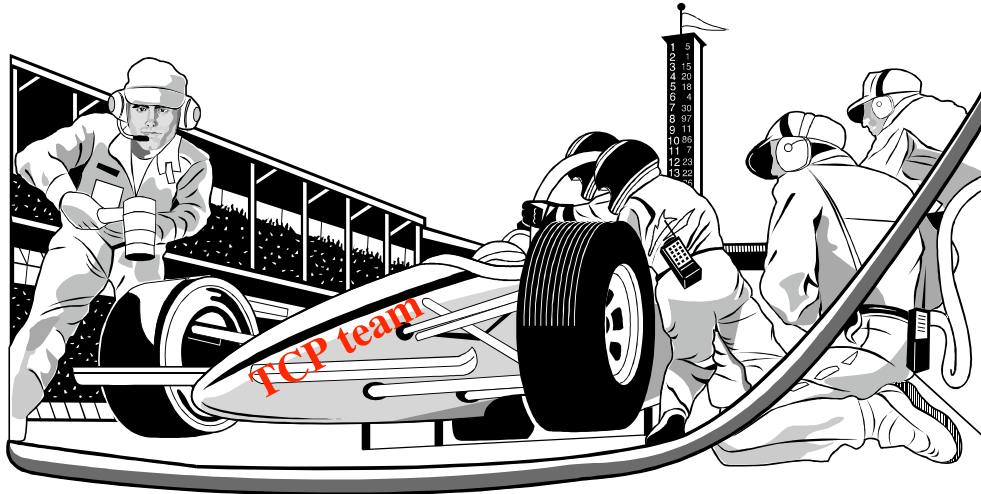
# AIMD

Phase plot



Fairness is preserved under Multiplicative Decrease since the user's allocation ratio remains the same

Ex:

$$\frac{x_2}{x_1} = \frac{x_2\, b}{x_1\, b}$$

- ❑ Assumption: decrease policy must (at minimum) reverse the load increase over-and-above efficiency line
- ❑ Implication: decrease factor should be conservatively set to account for any congestion detection lags etc

# Tuning stand for TCP
## the dark side of speed!

**TCP performances depend on**

☐ **TCP & network parameters**
- Congestion window size, *ssthresh* (threshold)
- RTO timeout settings
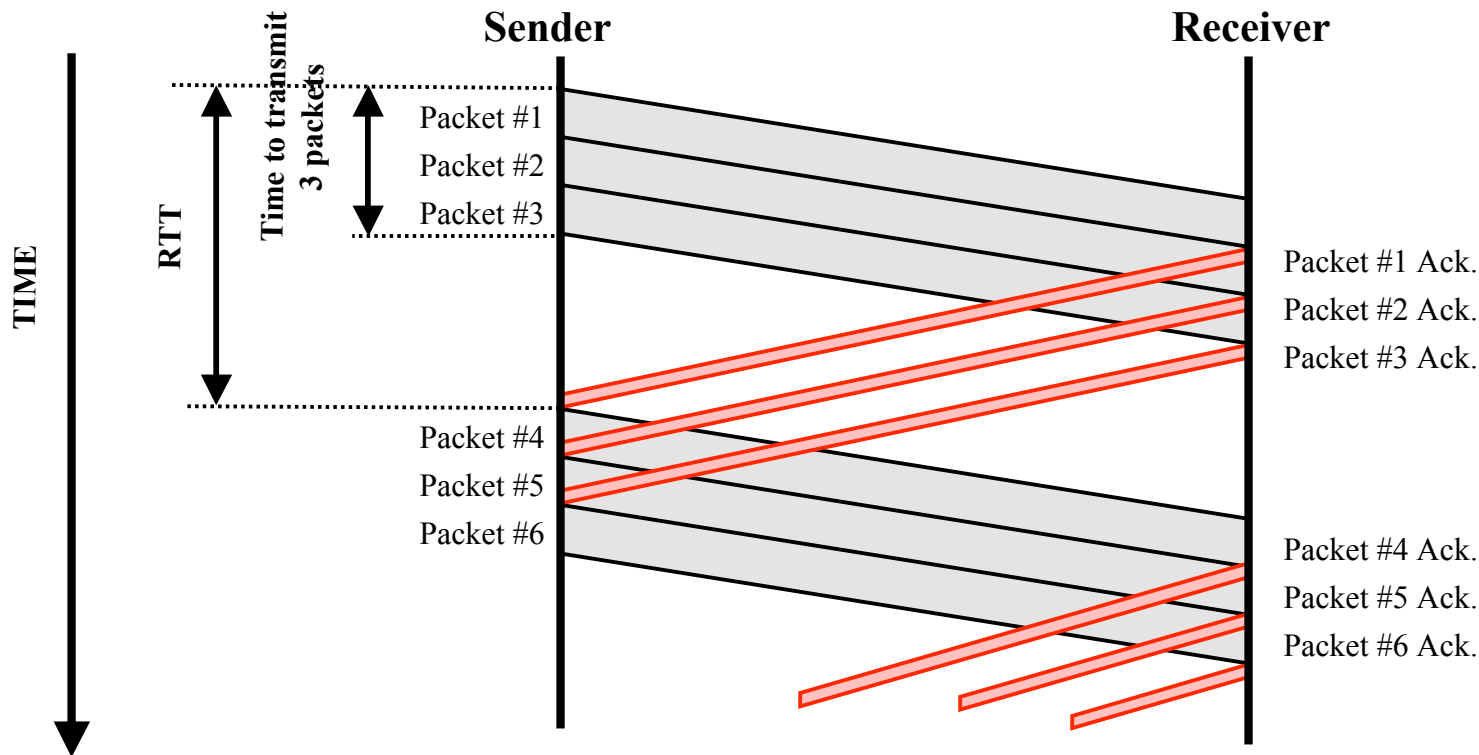- SACKs
- Packet size

☐ **System parameters**
- TCP and OS buffer size (in comm. subsys., drivers…)

**NEED A SPECIALIST!**

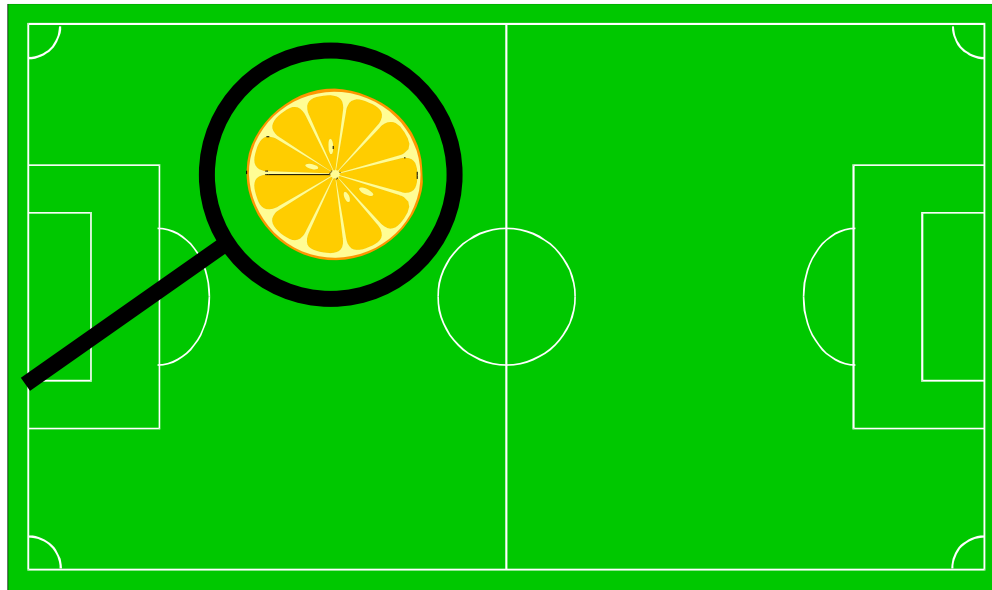# First problem: window size

❑ The default maximum window size is 64Kbytes. Then the sender has to wait for acks.

# First problem: window size

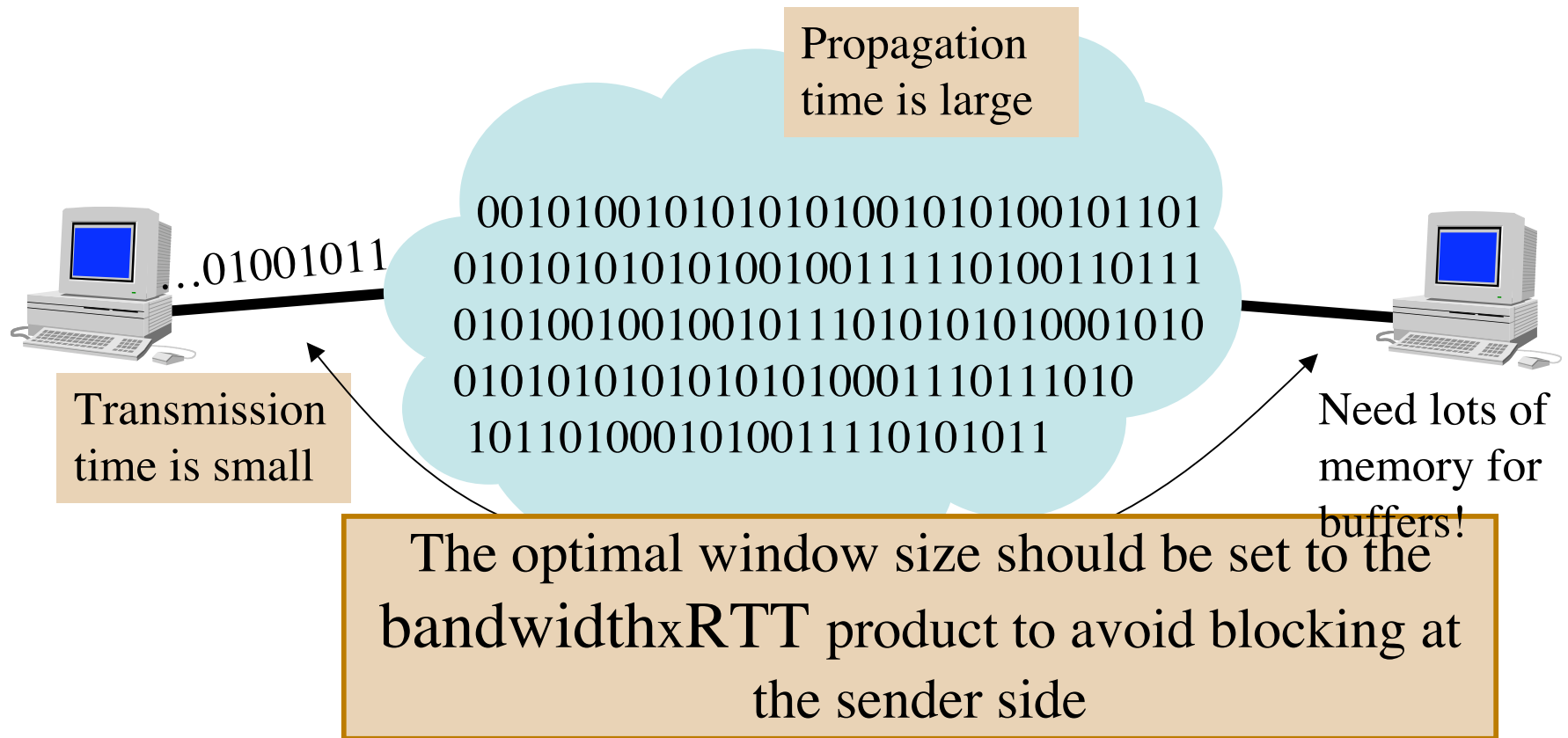❑ The default maximum window size is 64Kbytes. Then the sender has to wait for acks.

RTT=200ms Link is 0C-48 = 2.5 Gbps

Waiting time

# Rule of thumb on Long Fat Networks

capacity

❏High-~~speed~~ network

Propagation time is large

00101001010101010010101001011010101010101001001111101001101110101001001001011101010101010001010010101010101010101000111011101010110100010100111101010111

...01001011

Transmission time is small

Need lots of memory for buffers!

The optimal window size should be set to the bandwidthxRTT product to avoid blocking at the sender side
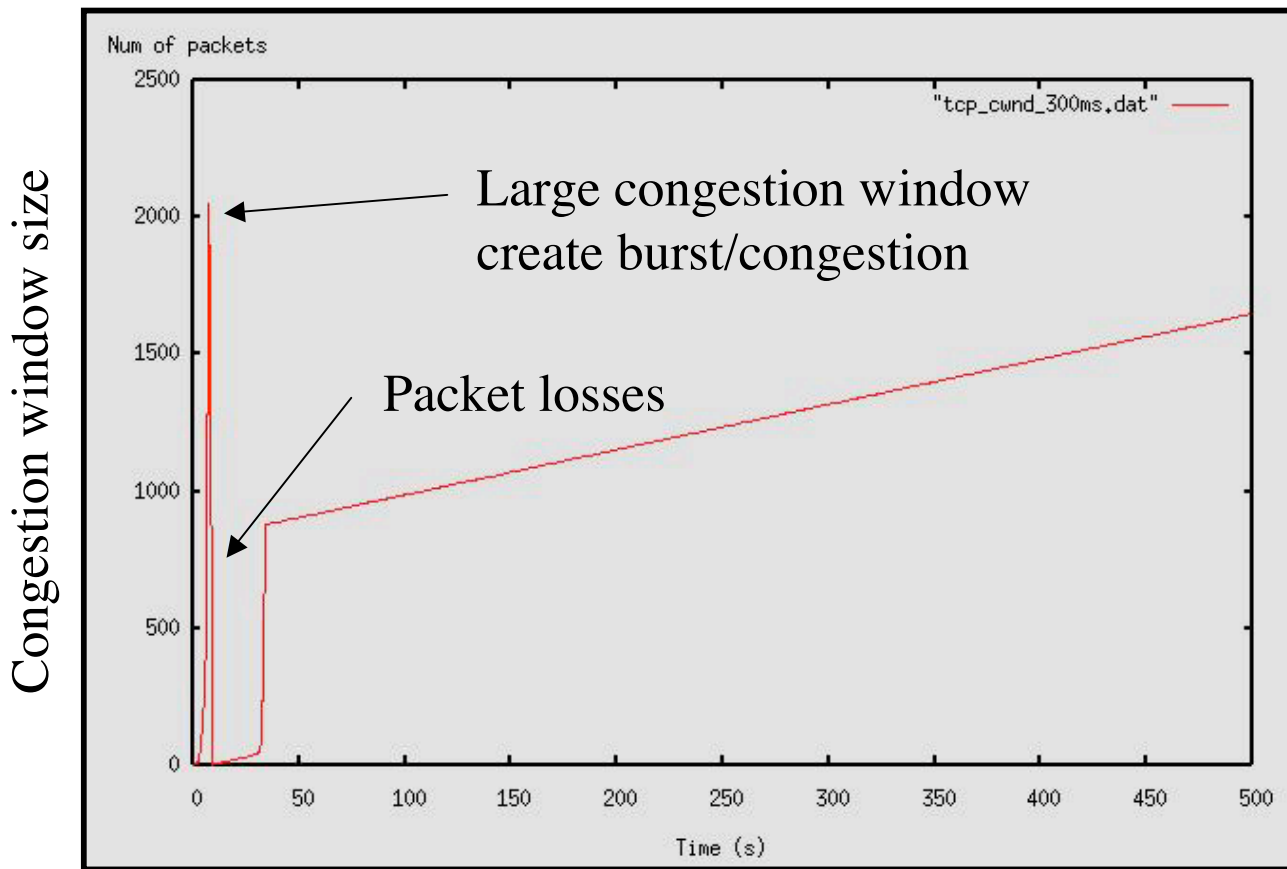
# Side effect of large windows
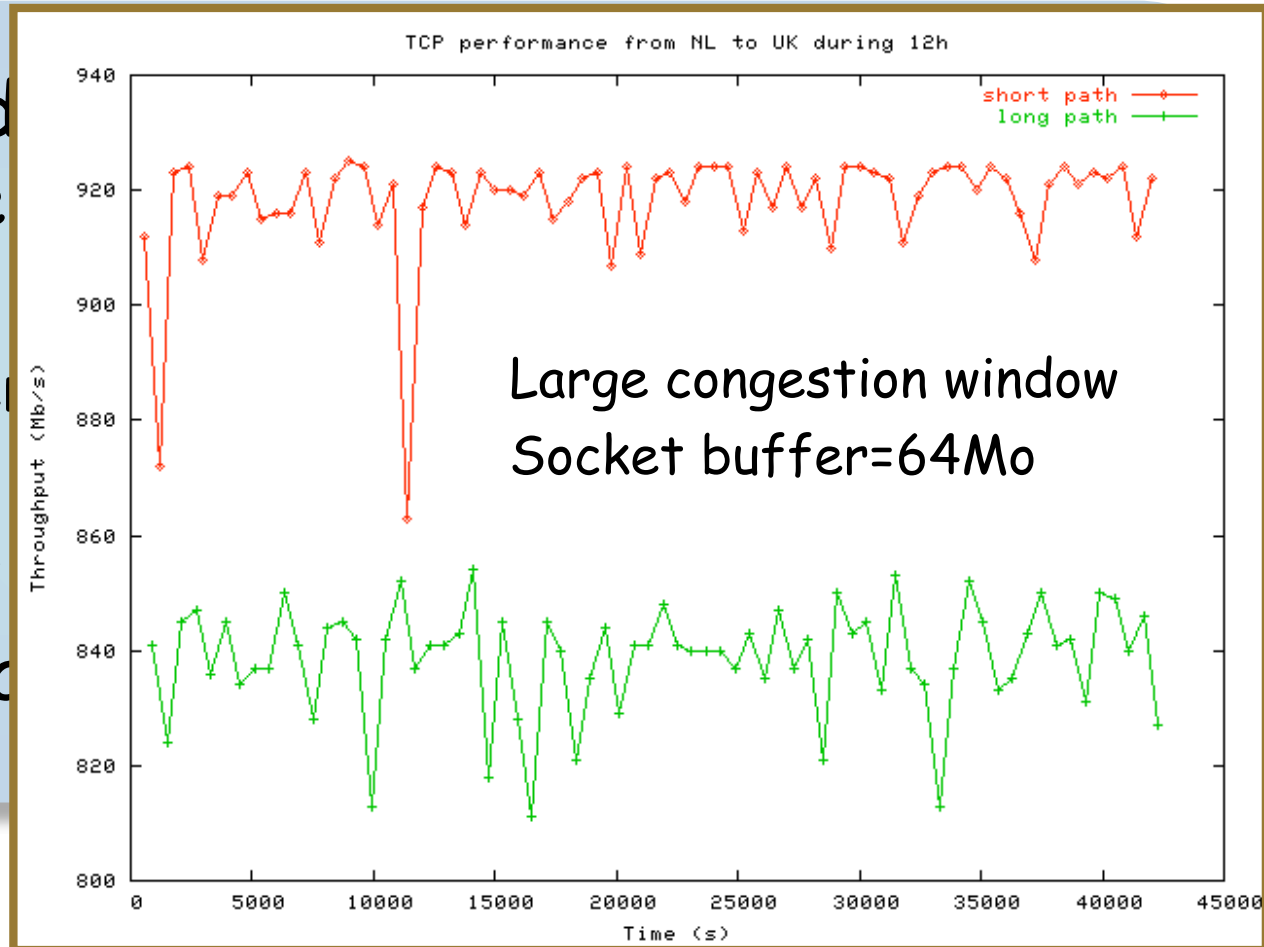
TCP becomes very sensitive to packet losses on LFN

# Pushing the limits of TCP

- Standard configuration (vanilla TCP) is not adequate on many OS, everything is under-sized
  - Receiver buffer
  - System buffer
  - Default block size
- Will manage to get near 1Gbps if well-tuned

# Pushing the limits of TCP

- ❑ Standard
  adequate
  sized
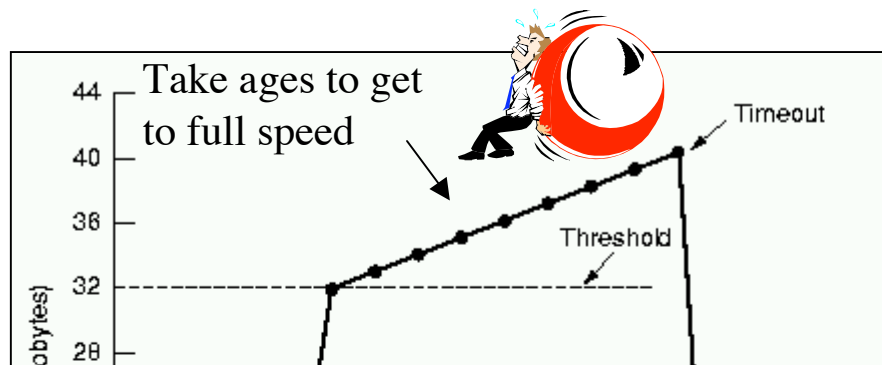  - ❑ Receiver
  - ❑ System
  - ❑ Default
- ❑ Will mand



TCP performance from NL to UK during 12h

Large congestion window
Socket buffer=64Mo

Source: M. Goutelle, GEANT test campaign

# Some TCP tuning guides

- http://www.psc.edu/networking/projects/tcptune/
- http://www.web100.org/
- http://rdweb.cns.vt.edu/public/notes/win2k-tcpip.htm
- http://www.sean.de/Solaris/soltune.html
- http://datatag.web.cern.ch/datatag/howto/tcp.html

# Problem on high capacity link?
# Additive increase is still too slow!

Take ages to get to full speed



With 100ms of round trip time, a connection needs 203 minutes (3h23) to send at 10Gbps starting from 1Mbps!

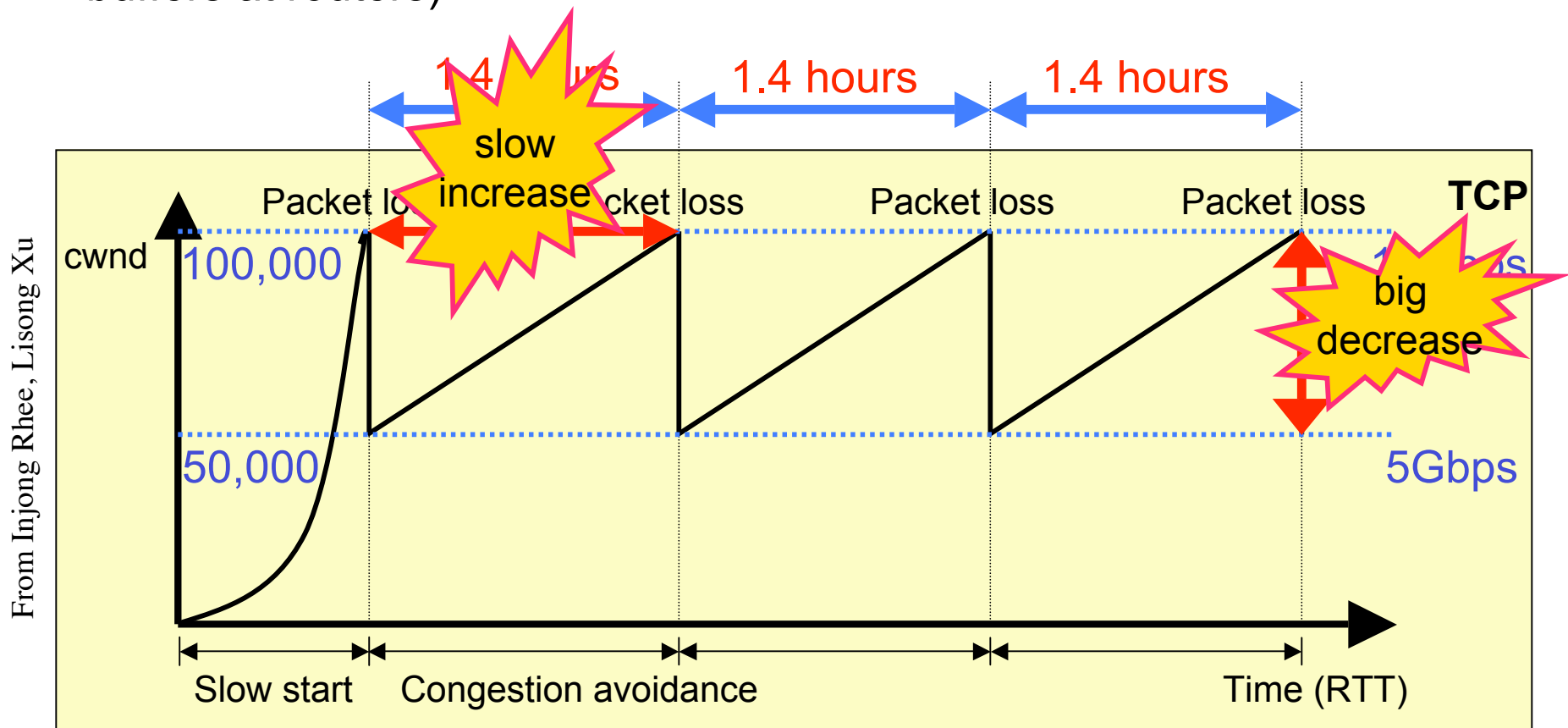**Once you get high throughput, maintaining it is difficult too!**

- Sustaining high congestion windows:
A Standard TCP connection with:
  – 1500-byte packets;
  – a 100 ms round-trip time;
  – a steady-state throughput of 10 Gbps;
would require:
  – an average congestion window of 83,333 segments;
  – and at most one drop (or mark) every 5,000,000,000 packets (or equivalently, at most one drop every 1 2/3 hours).
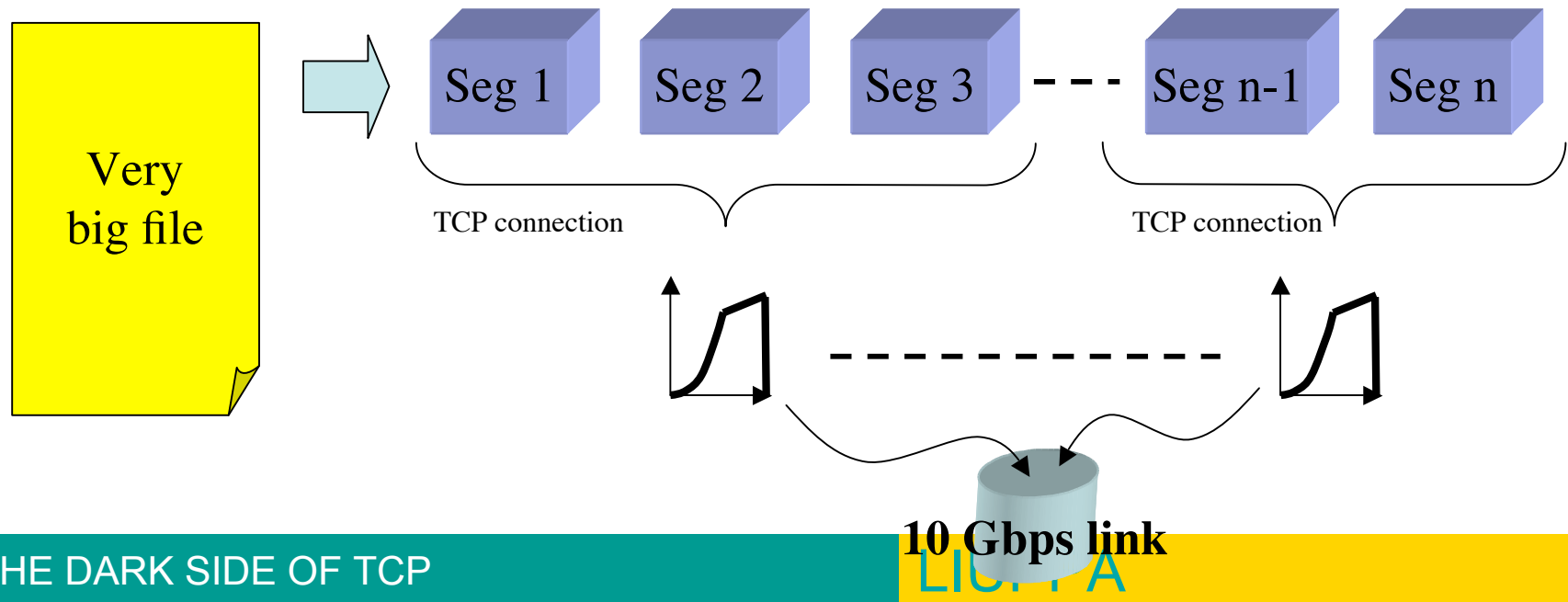This is not realistic.

From S. Floyd

# TCP rules:
# slow increase, big decrease

A TCP connection with 1250-Byte packet size and 100ms RTT is running over a 10Gbps link (assuming no other connections, and no buffers at routers)
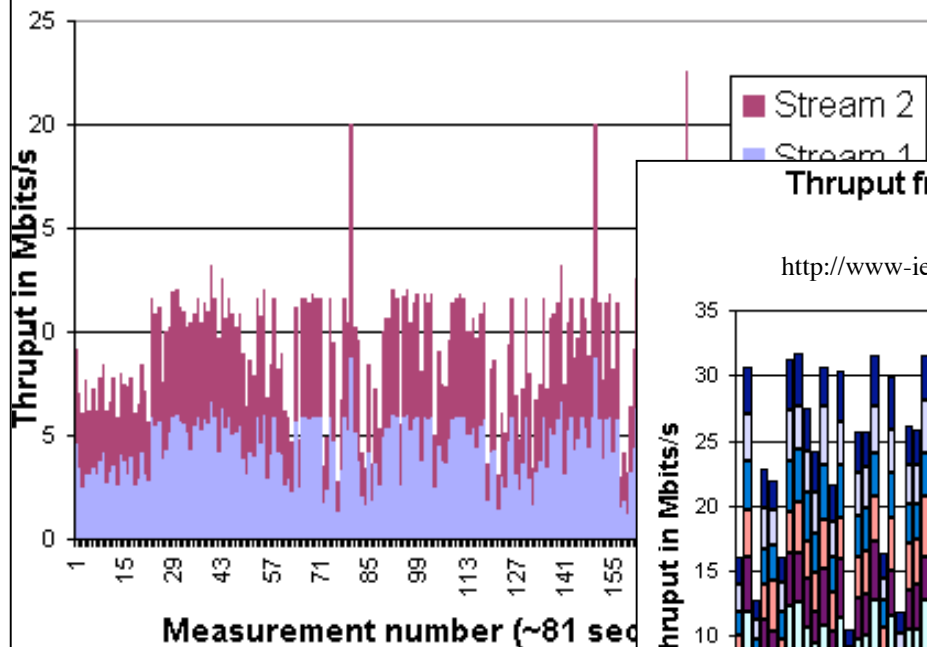
# Going faster (cheating?)
## *n* flows is better than 1

❑ The CC limits the throughput of a TCP connection: so why not use more than 1 connection for the same file?
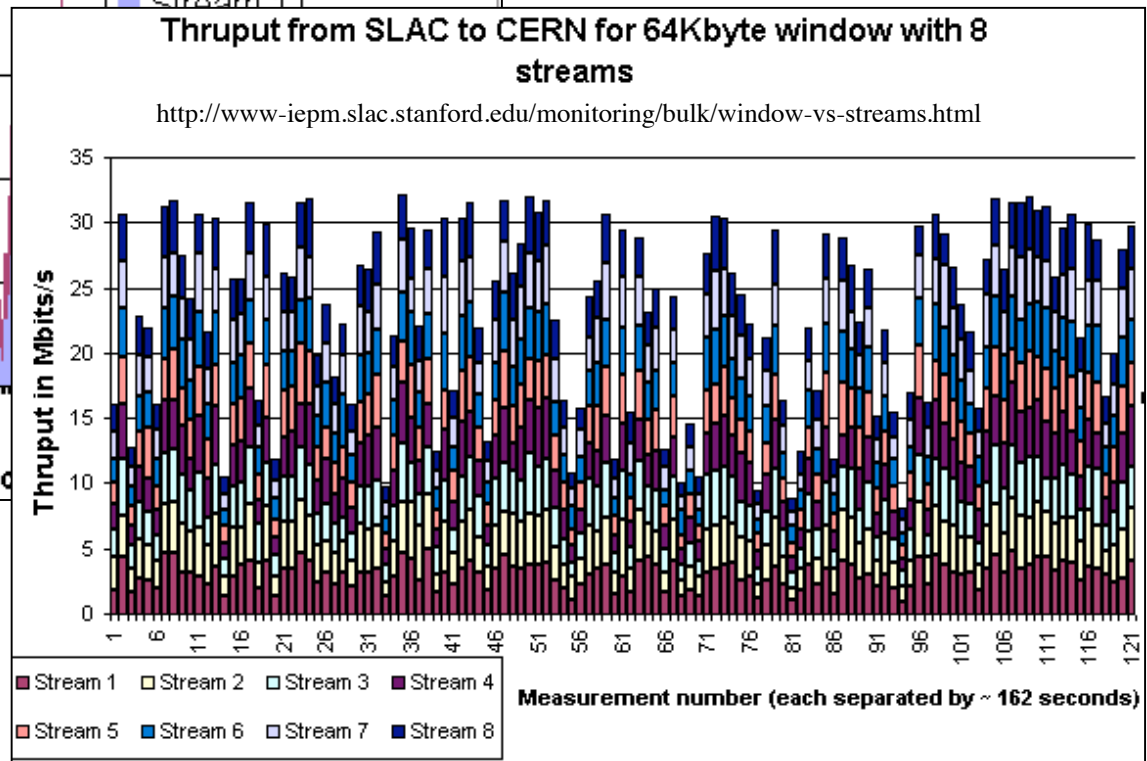
Very big file

Seg 1    Seg 2    Seg 3    - - -    Seg n-1    Seg n

TCP connection    TCP connection

10 Gbps link

# Some results from IEPM/SLAC



Thruput SLAC to CERN with 256kByte window & 2 streams

**More streams is better than larger congestion windows**

Thruput from SLAC to CERN for 64Kbyte window with 8 streams

http://www-iepm.slac.stanford.edu/monitoring/bulk/window-vs-streams.html

# Multiple streams

- No/few modifications to transport protocols (i.e. TCP)
  - Parallel socket libraries
  - GridFTP (http://www.globus.org/datagrid/gridftp.html)
  - bbFTP (http://doc.in2p3.fr/bbftp/)

# New transport protocols

- New transport protocols are those that are not only optimizations of TCP
- New behaviors, new rules, new requirements! Everything is possible!
- New protocols are then not necessarily TCP compatible!
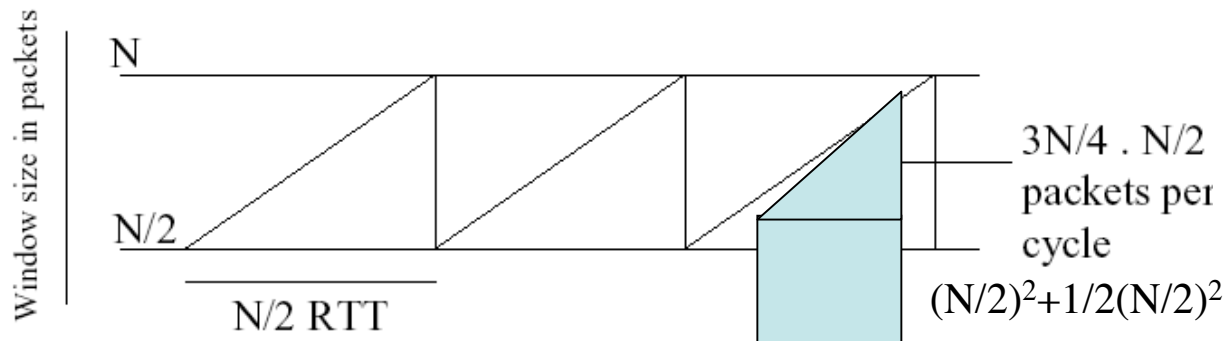
# The new transport protocol strip

# Response function

- Throughput = f(p, RTT)
- TCP's response function



Window size in packets

N

N/2

N/2 RTT

3N/4 . N/2 packets per cycle

$(N/2)^2 + 1/2(N/2)^2$

Average window size (in packets) = W = 3N/4 , from (N+N/2)/2

Number of packets per cycle = 3N/4 . N/2 = $3N^2/8$ = 1/ p

- Where p is the packet loss ratio (which should remain small enough)
- So $N = \sqrt{\dfrac{8}{3p}}$

Average throughput (in packets/sec) = B = W / RTT = 3N / 4 RTT

$$Throughput = \frac{W}{RTT} = \sqrt{\frac{3}{2}}\frac{MTU}{RTT\sqrt{p}} = \sqrt{\frac{3}{2}}\frac{1}{RTT\sqrt{p}}$$

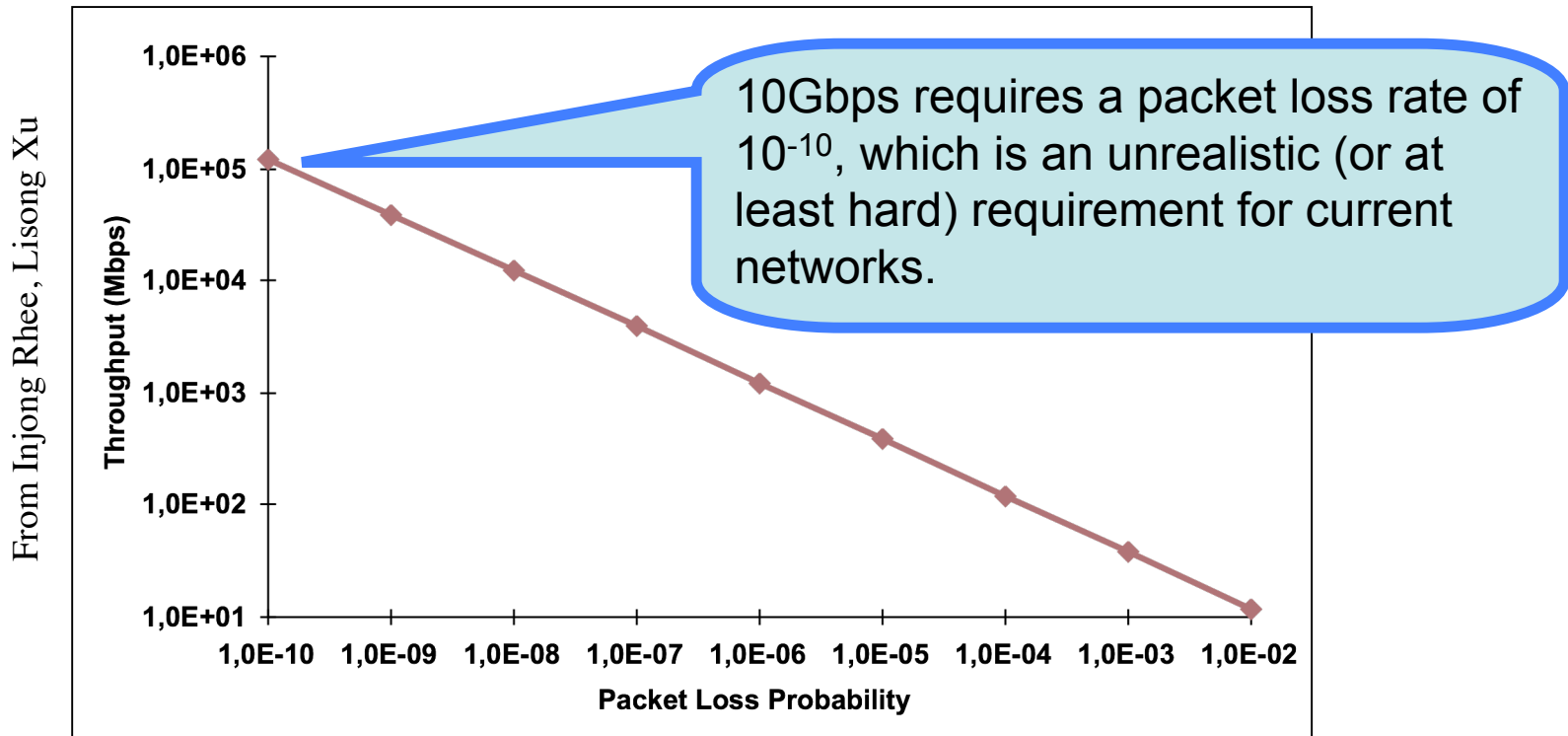# TCP's response function in image

$$Throughput = \frac{W}{RTT} = \sqrt{\frac{3}{2}} \frac{MTU}{RTT\sqrt{p}}$$
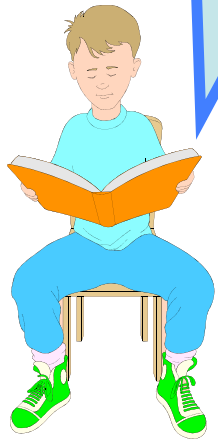
$MTU$ : Packet Size
$RTT$ : Round-Trip Time
$P$    : Packet Loss Probability



10Gbps requires a packet loss rate of $10^{-10}$, which is an unrealistic (or at least hard) requirement for current networks.

From Injong Rhee, Lisong Xu

# AIMD, general case



**TCP:** $R = \dfrac{MSS}{RTT}\dfrac{1.2}{p^{0.5}}$

**AIMD:** $R = \dfrac{MSS}{RTT}\dfrac{15.5}{p^{0.5}}$

Inspired from Injong Rhee, Lisong Xu

# High Speed TCP [Floyd]

❏ Modifies the response function to allow for more link utilization in current high-speed networks where the loss rate is smaller than that of the networks TCP was designed for (at most $10^{-2}$)

| TCP Throughput (Mbps) | RTTs Between Losses | W | P |
|---:|---:|---:|---|
| 1 | 5.5 | 8.3 | 0.02 |
| 10 | 55.5 | 83.3 | 0.0002 |
| 100 | 555.5 | 833.3 | 0.000002 |
| 1000 | 5555.5 | 8333.3 | 0.00000002 |
| 10000 | 55555.5 | 83333.3 | 0.0000000002 |

Table 1: RTTs Between Congestion Events for Standard TCP, for 1500-Byte Packets and a Round-Trip Time of 0.1 Seconds.

From draft-ietf-tsvwg-highspeed-01.txt

# Modifying the response

| Packet Drop Rate P | Congestion Window W | RTTs Between Losses |
|---|---|---|
| 10^-2 | 12 | 8 |
| 10^-3 | 38 | 25 |
| 10^-4 | 120 | 80 |
| 10^-5 | 379 | 252 |
| 10^-6 | 1200 | 800 |
| 10^-7 | 3795 | 2530 |
| 10^-8 | 12000 | 8000 |
| 10^-9 | 37948 | 25298 |
| 10^-10 | 120000 | 80000 |

Table 2: TCP Response Function for Standard TCP.  The average congestion window W in MSS-sized segments is given as a function of the packet drop rate P.
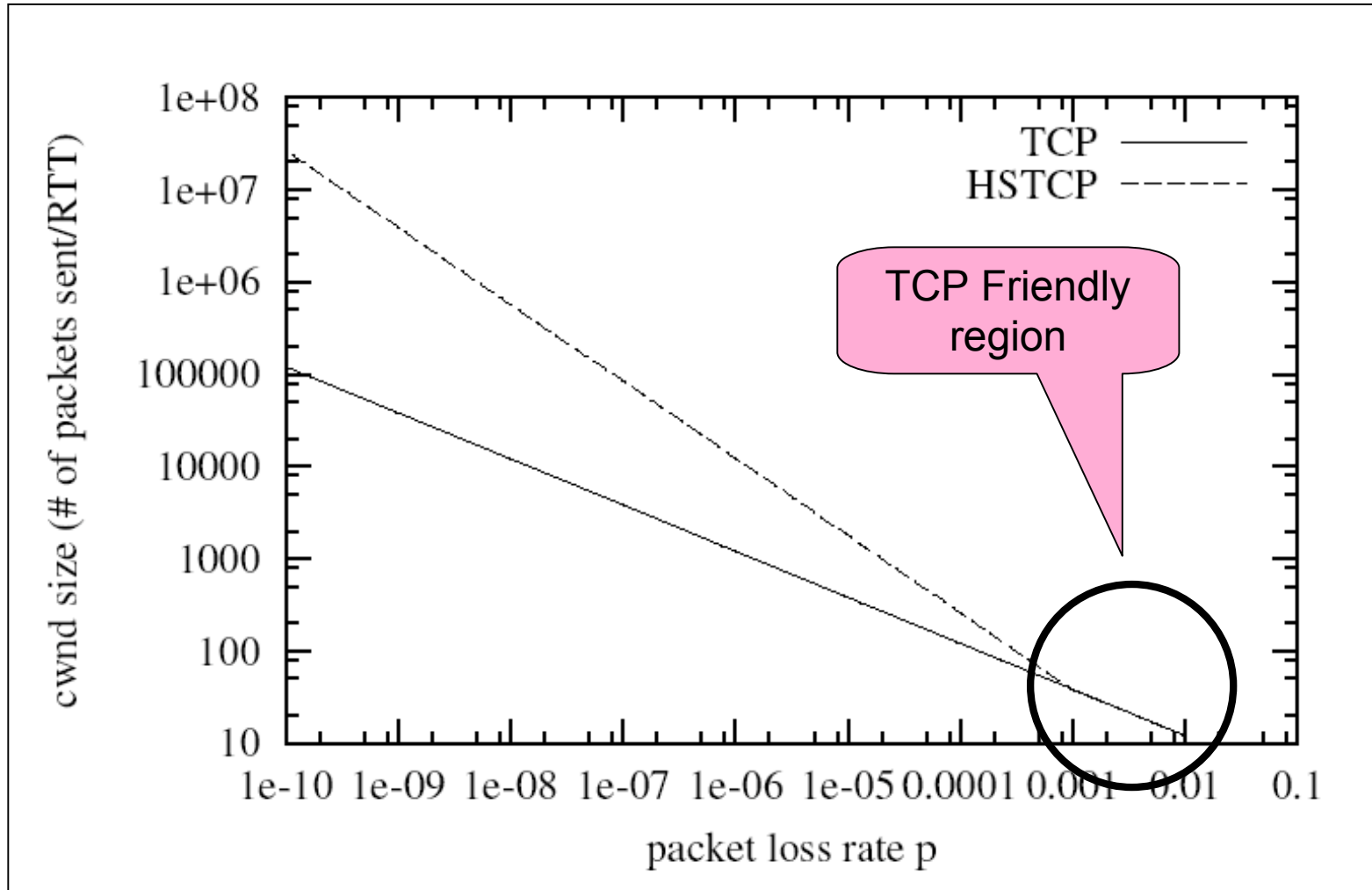
From draft-ietf-tsvwg-highspeed-01.txt

To specify a modified response function for HighSpeed TCP, we use three parameters, Low_Window, High_Window, and High_P.  To Ensure TCP compatibility, the HighSpeed response function uses the same response function as Standard TCP when the current congestion window is at most Low_Window, and uses the HighSpeed response function when the current congestion window is greater than Low_Window.  In this document we set Low_Window to 38 MSS-sized segments, corresponding to a packet drop rate of 10^-3 for TCP.

| Packet Drop Rate P | Congestion Window W | RTTs Between Losses |
|---|---|---|
| 10^-2 | 12 | 8 |
| 10^-3 | 38 | 25 |
| 10^-4 | 263 | 38 |
| 10^-5 | 1795 | 57 |
| 10^-6 | 12279 | 83 |
| 10^-7 | 83981 | 123 |
| 10^-8 | 574356 | 180 |
| 10^-9 | 3928088 | 264 |
| 10^-10 | 26864653 | 388 |

Table 3: TCP Response Function for HighSpeed TCP.  The average congestion window W in MSS-sized segments is given as a function of the packet drop rate P.

# See it in image

# Relation with AIMD

□ TCP-AIMD
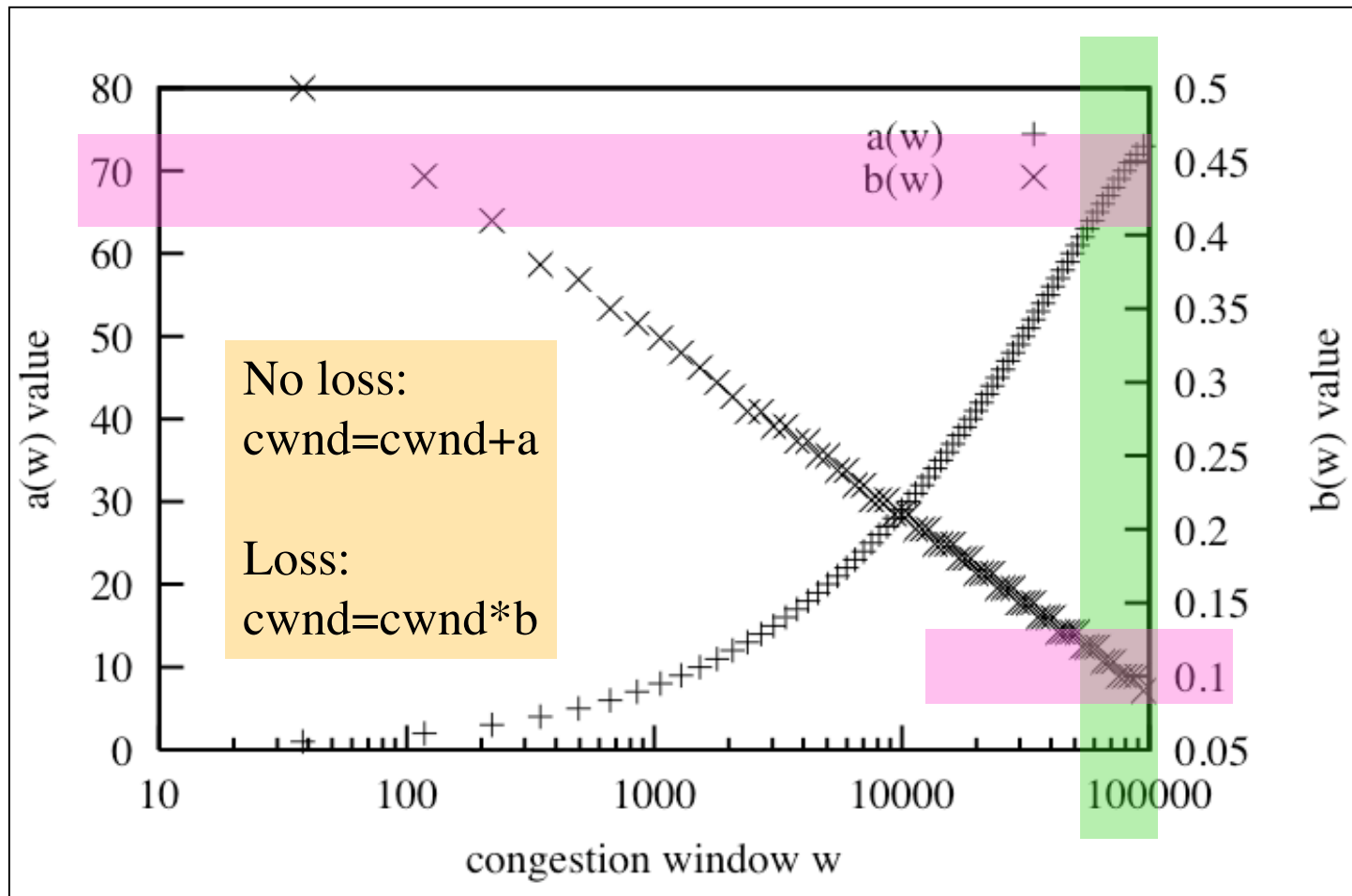  □ Additive increase: a=1
  □ Multiplicative decrease: b=1/2
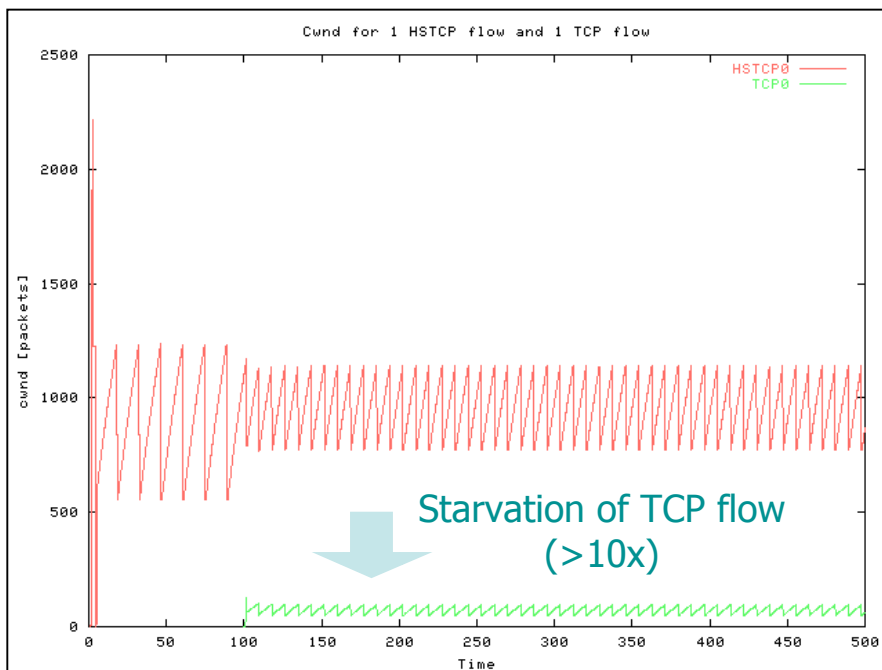
□ HSTCP-AIMD
  □ Link a & b to congestion window size
  □ a = a(cwnd), b=b(cwnd)
  □ General rules
    • the larger cwnd, the larger the increment
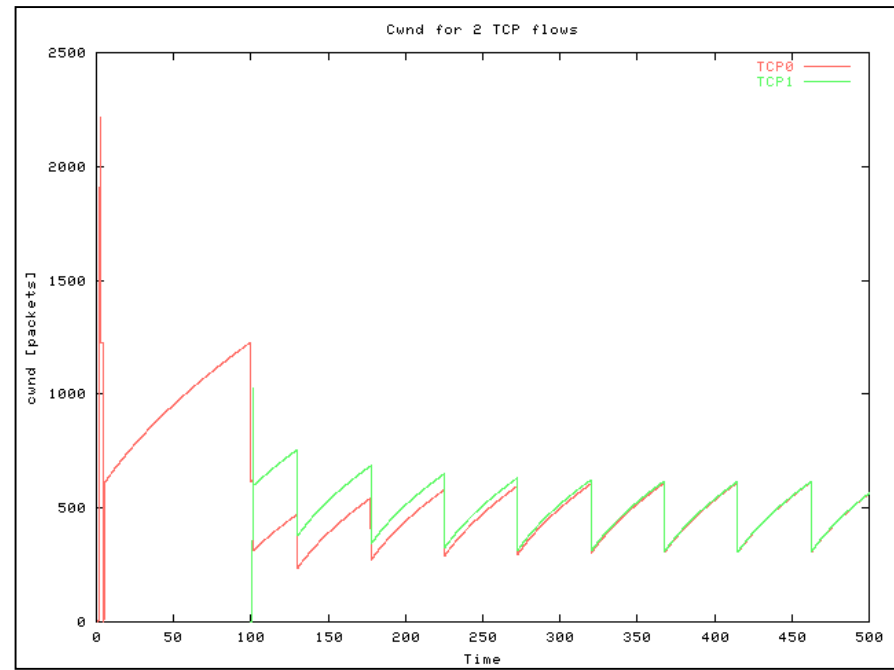    • The larger cwnd, the smaller the decrement

no loss:

cwnd = cwnd + 1

loss:

cwnd = cwnd*0.5

# Quick to grab bandwidth, slow to give some back!



No loss:
cwnd=cwnd+a

Loss:
cwnd=cwnd*b

# Talking about dark side...



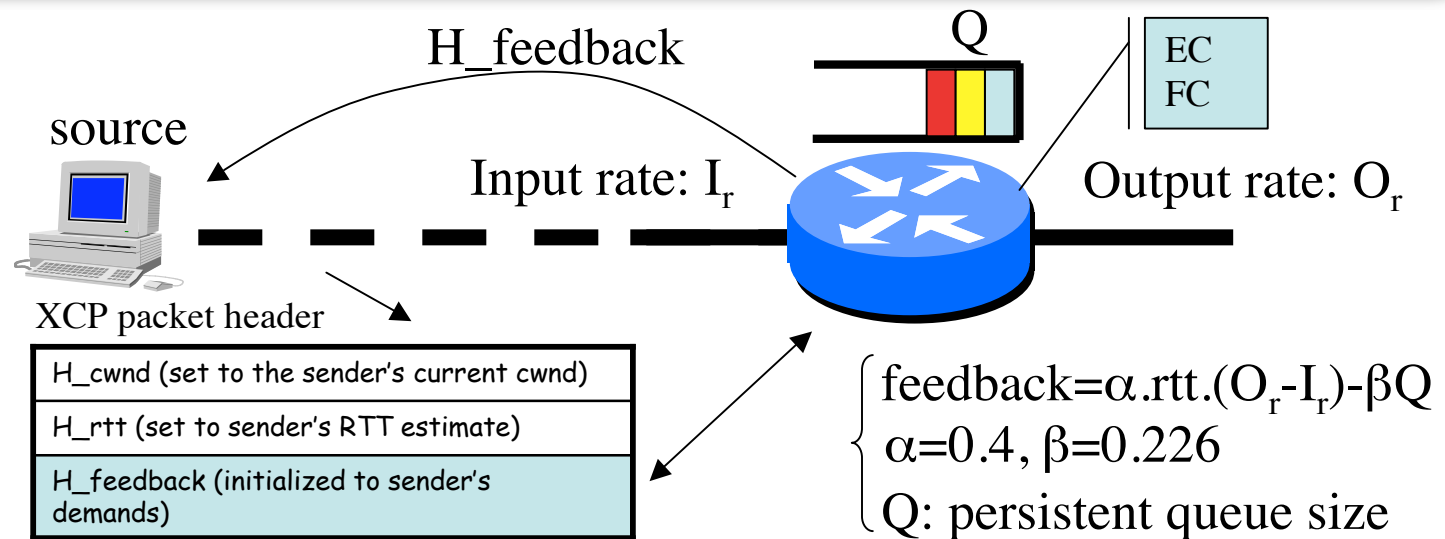**1 HSTCP and 1 TCP flow**

**2 TCP flows**

SETUP  RTT=100ms
Bottleneck BW=50Mbps
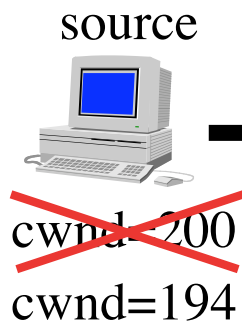Qsize=BW*RTT
Qtype=DropTail

# XCP [Katabi02]

- ❑ XCP is a router-assisted solution, generalized the ECN concepts (FR, TCP-ECN)
- ❑ XCP routers can compute the available bandwidth by monitoring the input rate and the output rate
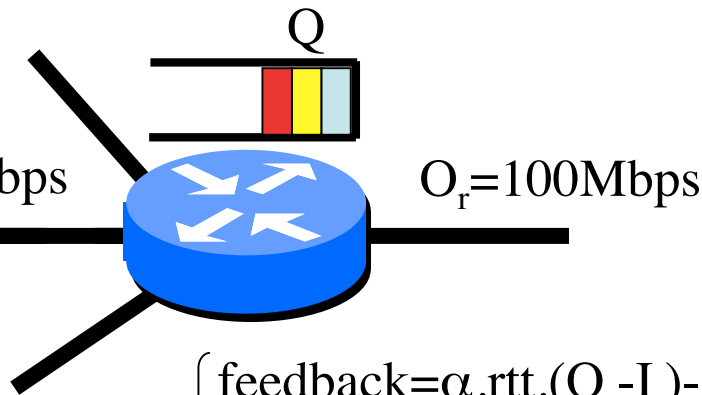- ❑ Feedback is sent back to the source in special fields of the packet header

H_feedback

Q

EC
FC

source

Input rate: $I_r$

Output rate: $O_r$

XCP packet header

| H_cwnd (set to the sender's current cwnd) |
| H_rtt (set to sender's RTT estimate) |
| H_feedback (initialized to sender's demands) |

$$\begin{cases} \text{feedback} = \alpha.\text{rtt}.(O_r - I_r) - \beta Q \\ \alpha = 0.4, \beta = 0.226 \\ Q: \text{persistent queue size} \end{cases}$$

# XCP in action

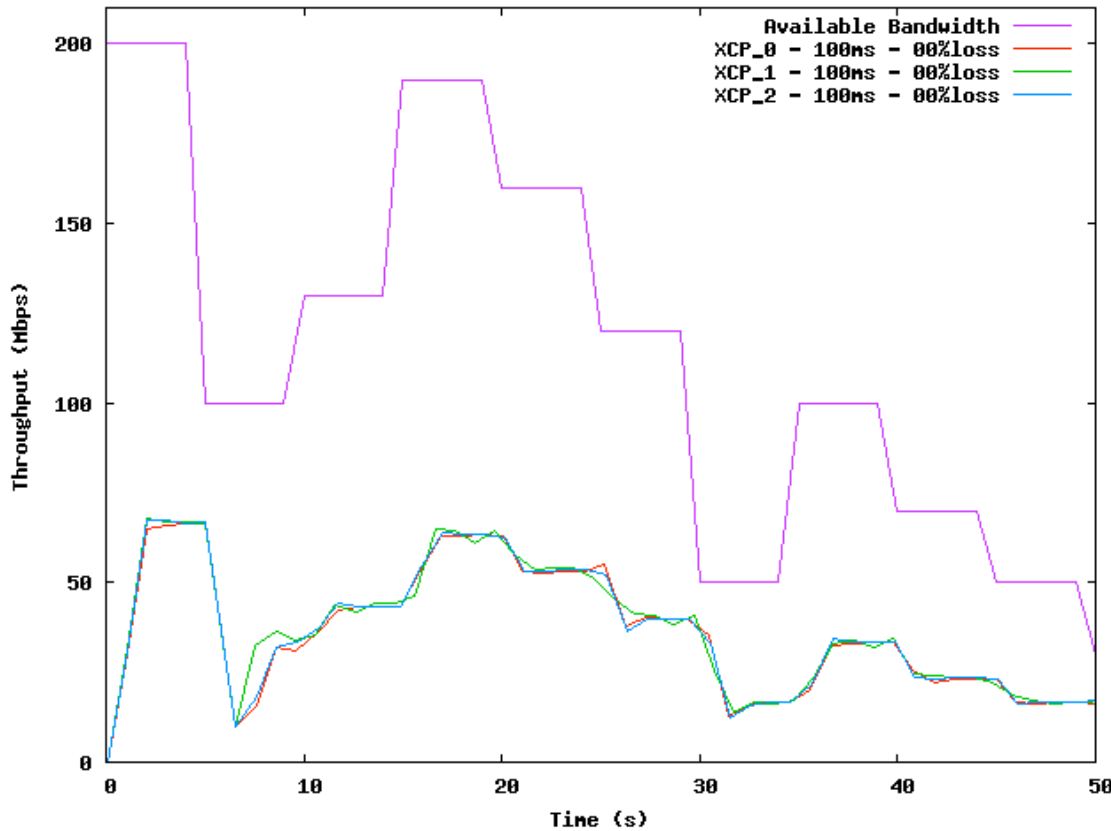Feedback value represents a window increment/decrement

| H_cwnd=200 |
|---|
| H_rtt=100ms |
| H_feedback=0 |

Q

source

$I_r$=250Mbps

$O_r$=100Mbps

cwnd=~~200~~
cwnd=194

$$feedback = \alpha . rtt . (O_r - I_r) - \beta Q$$
$$\alpha = 0.4, \beta = 0.226$$
Q: persistent queue size

| H_cwnd=200 |
|---|
| H_rtt=100ms |
| H_feedback=-6 |

**Case without βQ contribution**
$O_r - I_r$=100-250=-150
feedback=-6

# XCP
## Variable bandwidth environments

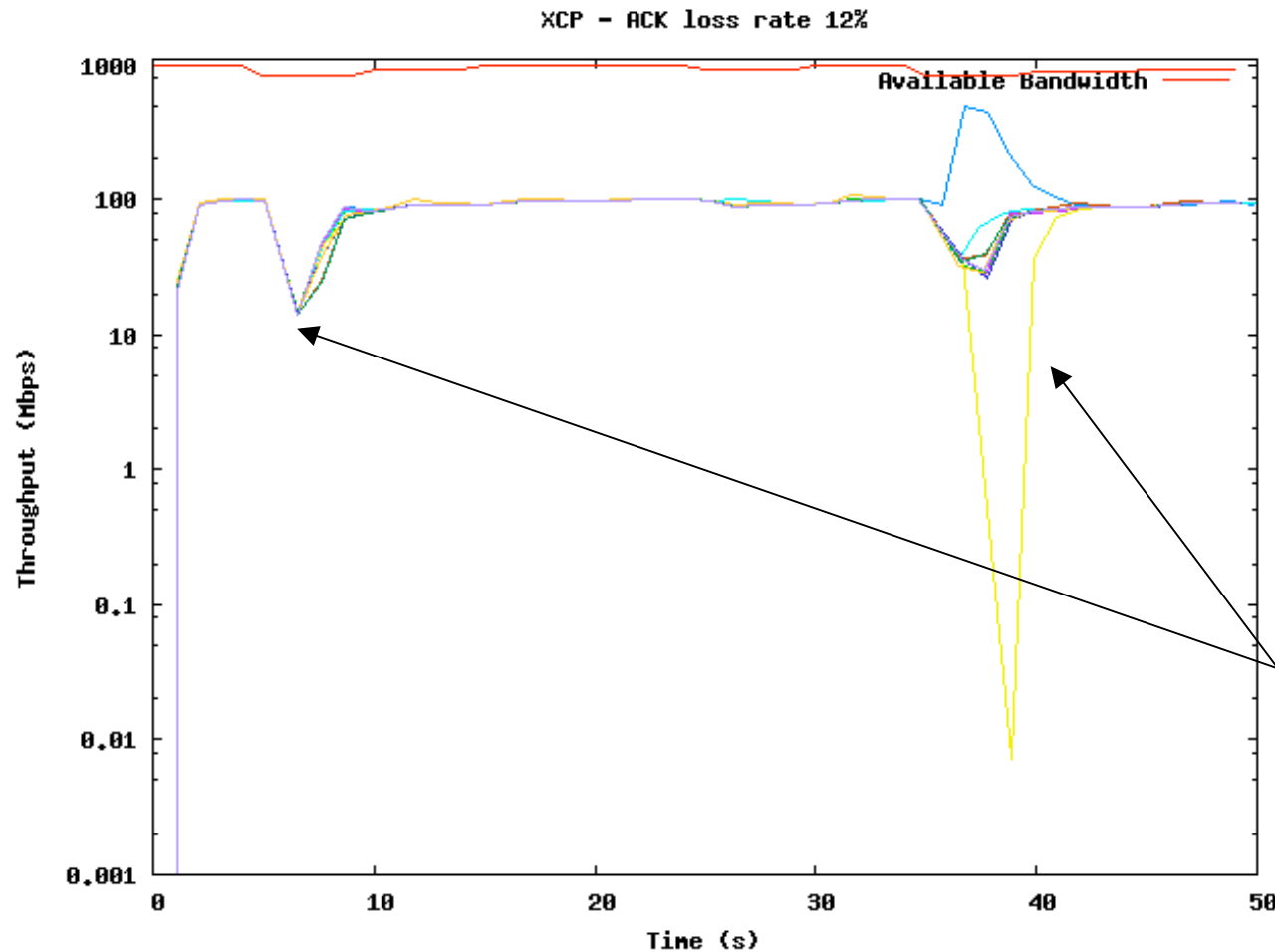

Good fairness and stability even in variable bandwidth environments

# XCP-r [Pacheco&Pham05]
## A more robust version of XCP

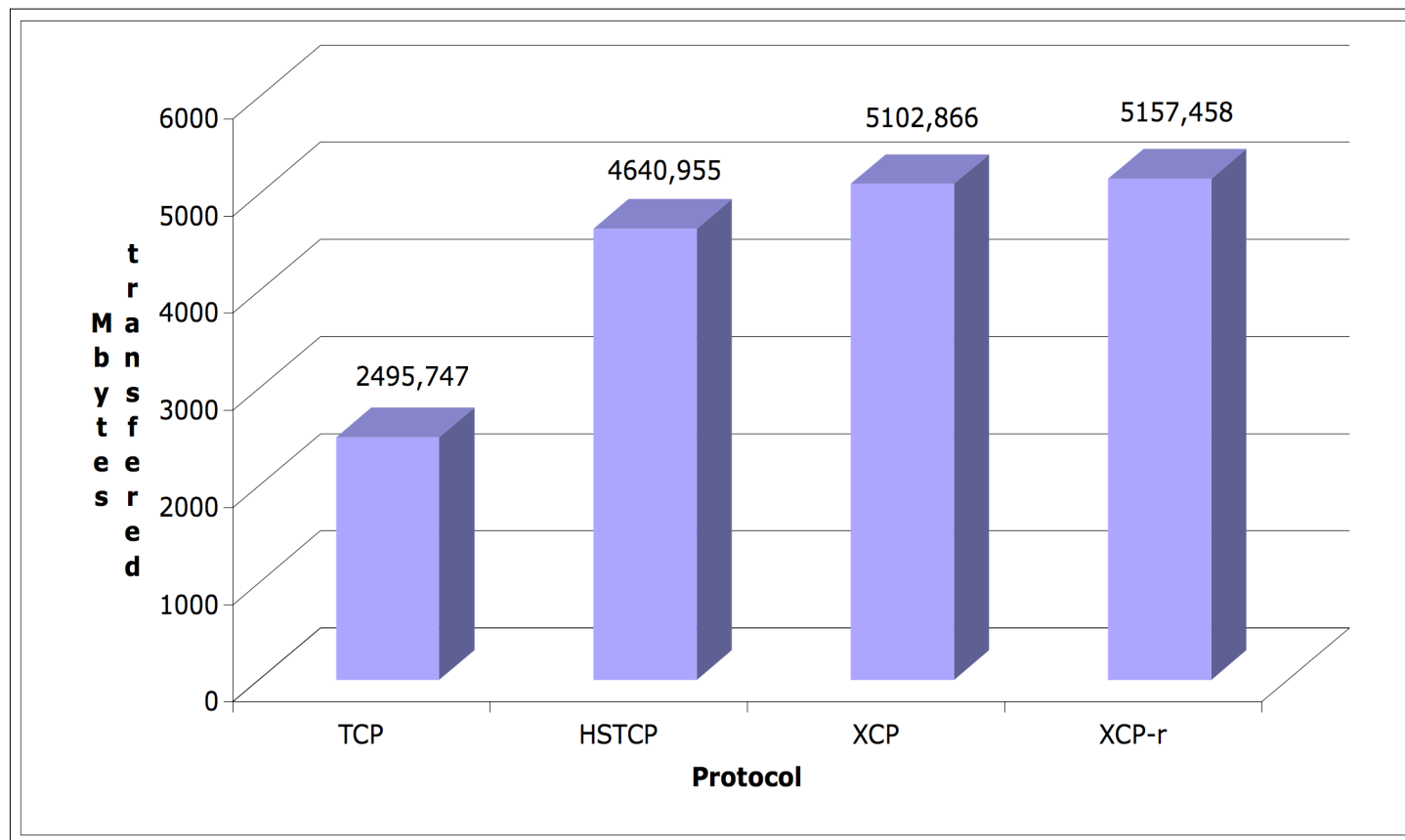

XCP - ACK loss rate 12%

10 flows sharing a 1Gbps link

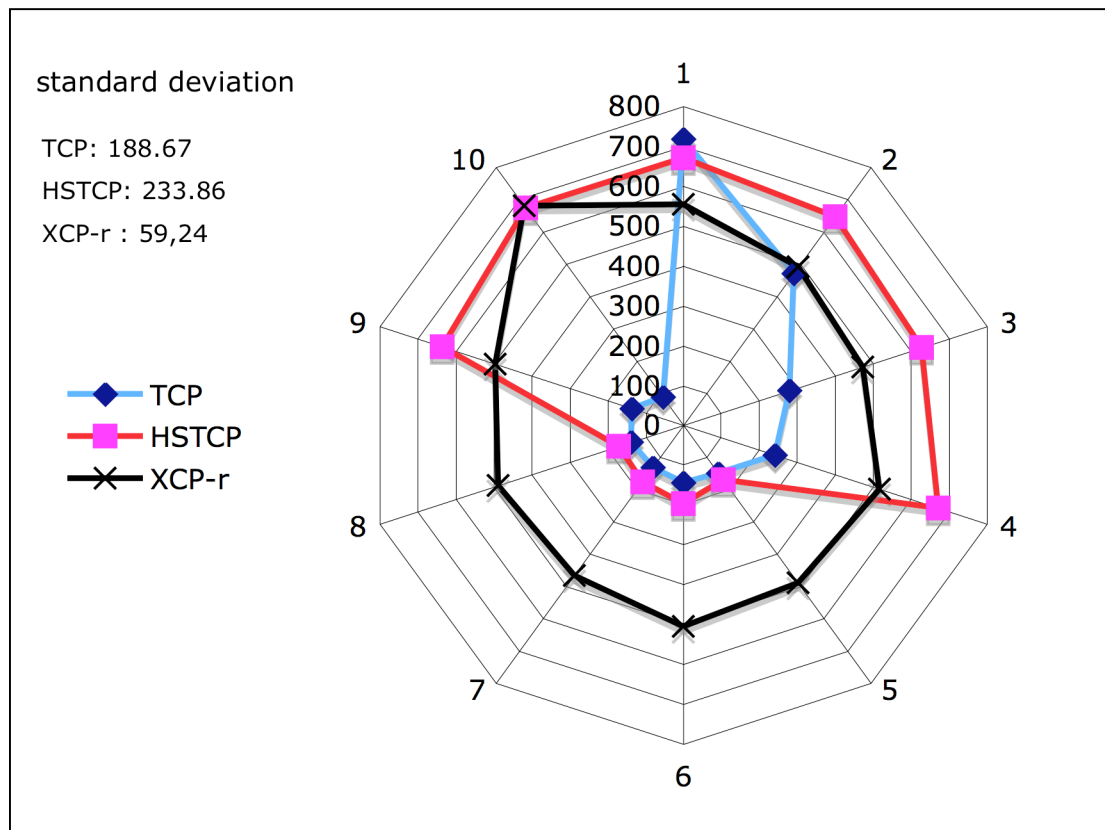Fast recovery after the timeouts and better fairness level

# XCP-r performance

Amount of data transfered in 50s, 10 flows, 1Gbps link, 200ms RTT

# XCP-r fairness

TCP and HSTCP are not really fair...

# Nothing is perfect :-(

- Multiple or parallel streams
  - How many streams?
  - Tradeoff between window size and number of streams
- New protocol
  - Fairness issues?
  - Deployment issues?
  - Still too early to know the side effects

# Where to find the new protocols?

- ❑ HSTCP
  - http://www.icir.org/floyd/hstcp.html
- ❑ STCP on Linux 2.4.19
  - http://www-lce.eng.cam.ac.uk/~ctk21/scalable/
- ❑ FAST
  - http://netlab.caltech.edu/FAST/
- ❑ XCP
  - http://www.ana.lcs.mit.edu/dina/XCP/
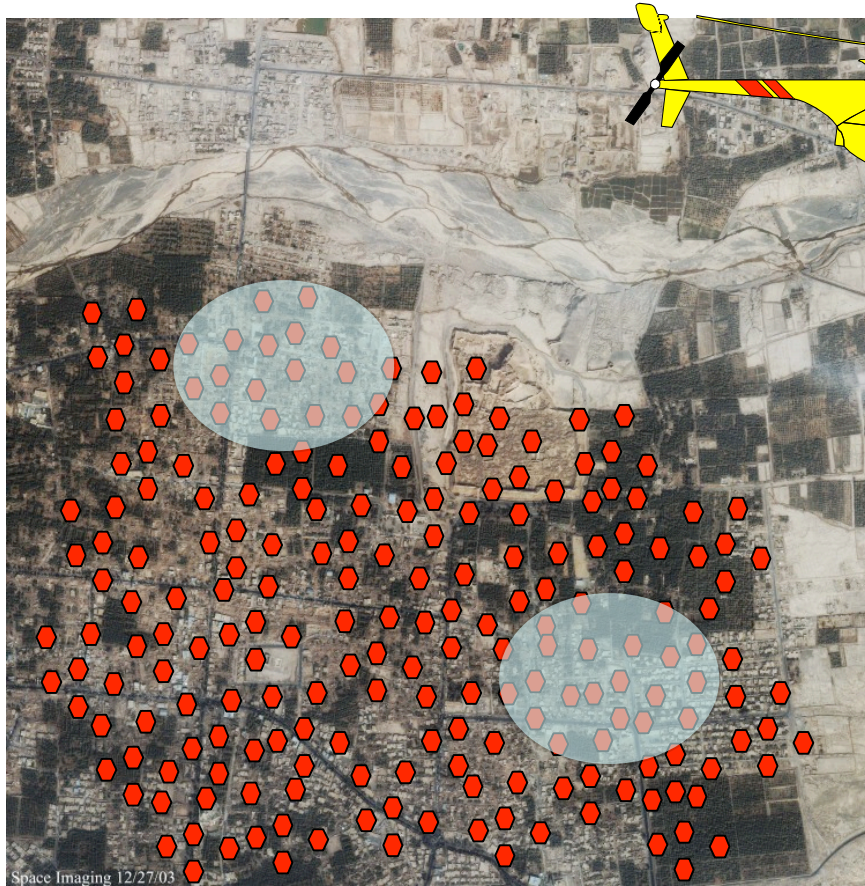  - http://www.isi.edu/isi-xcp/#software

# Web100 project

- [www.web100.org](www.web100.org)
- « The Web100 project will provide the software and tools necessary for end-hosts to automatically and transparently achieve high bandwidth data rates (100 Mbps) over the high performance research networks »
- Actually it's not limited to 100Mbps!
- Recommended solution for end-users to deploy and test high-speed transport solutions

# Hostile environments

- Asymetric networks
  - Satellite links & terrestrial links
- Wireless (WiFi, WiMax)
  - High loss probability
  - Losses ≠ congestions
- Ad-Hoc (PDA)
  - Small capacity
- Wireless Sensor Networks
  - **All of the above mentioned problems!**

# New sensor applications
## disaster relief - security



Real-time organization and optimization of rescue in large scale disasters

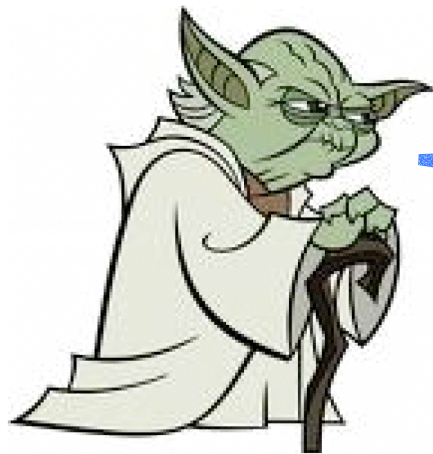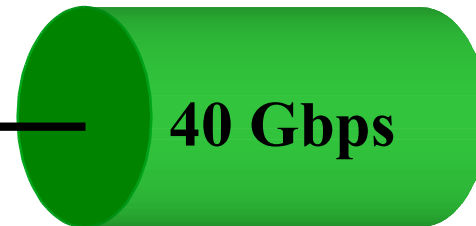Rapid deployment of fire detection systems in high-risk places

# Conclusions

❑ Understanding the dark side allows to move forwards!

❑ However...

vanilla TCP

10GB file

40 Gbps

**MAY THE FORCE BE WITH YOU!**