# The dark side of TCP

## and moving forwards...

C. Pham

www.univ-pau.fr/~cpham

Présentation tiré d'un tutorial
effectué le 5/7/07 pour IEEE DFMA

Mise à jour: sept 2011

LIUPPA
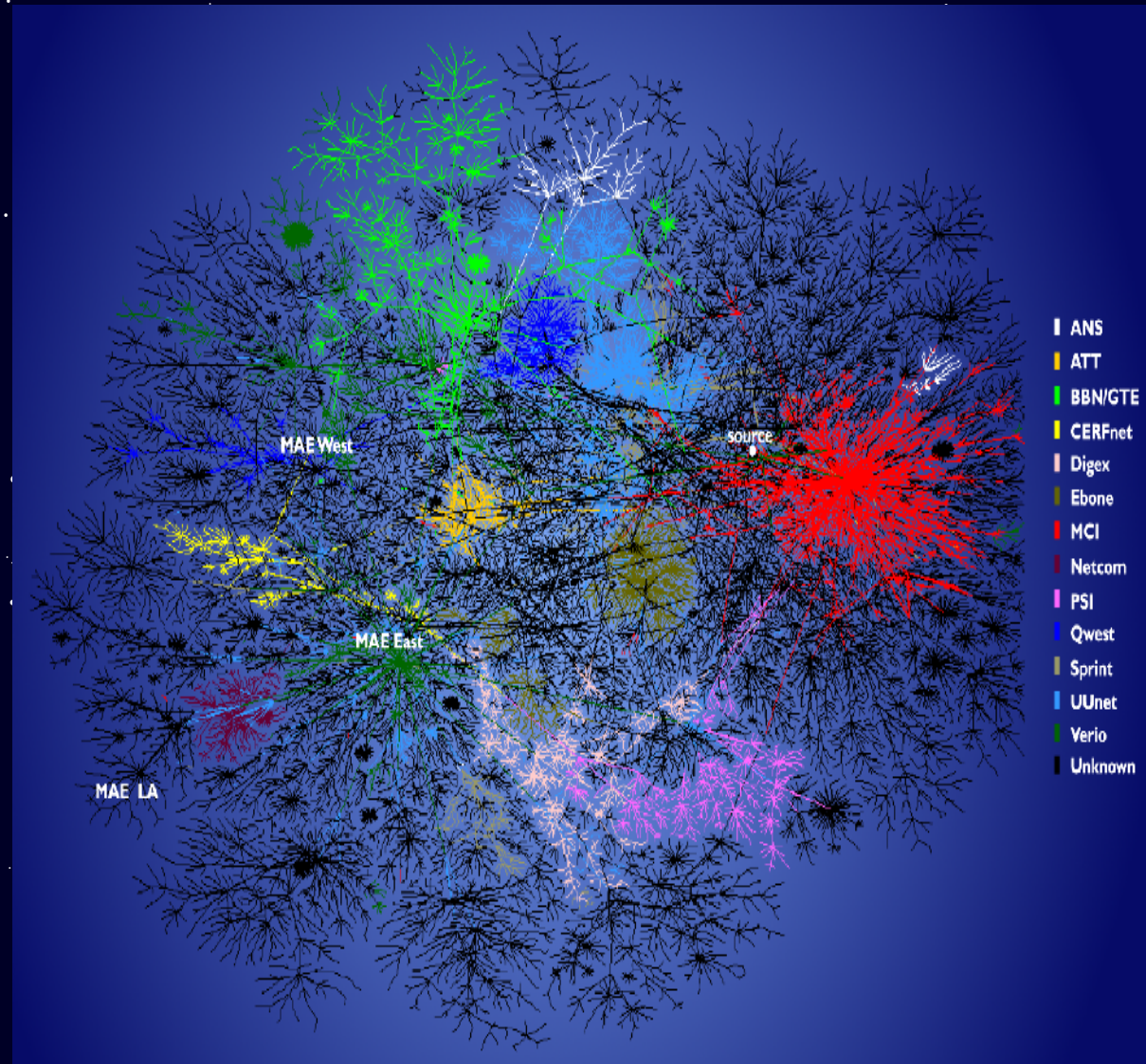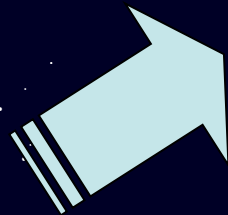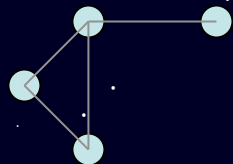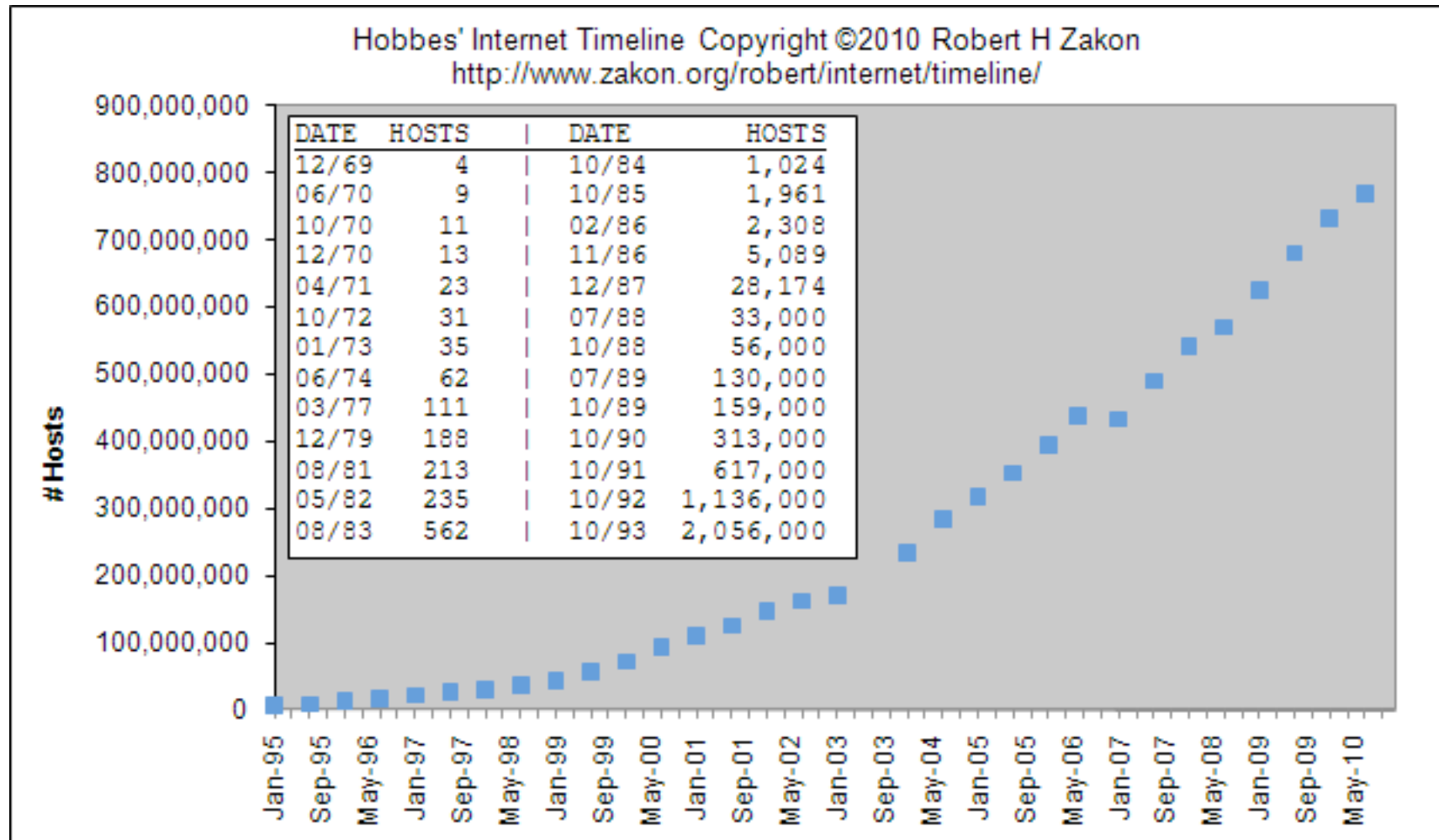
# The dark side of TCP

## and moving forwards...

**DFMA 07**

ENST, Paris, France

July 5th, 2007

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 ©2010 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

| DATE | HOSTS | | DATE | HOSTS |
|---|---|---|---|---|
| 12/69 | 4 | \| | 10/84 | 1,024 |
| 06/70 | 9 | \| | 10/85 | 1,961 |
| 10/70 | 11 | \| | 02/86 | 2,308 |
| 12/70 | 13 | \| | 11/86 | 5,089 |
| 04/71 | 23 | \| | 12/87 | 28,174 |
| 10/72 | 31 | \| | 07/88 | 33,000 |
| 01/73 | 35 | \| | 10/88 | 56,000 |
| 06/74 | 62 | \| | 07/89 | 130,000 |
| 03/77 | 111 | \| | 10/89 | 159,000 |
| 12/79 | 188 | \| | 10/90 | 313,000 |
| 08/81 | 213 | \| | 10/91 | 617,000 |
| 05/82 | 235 | \| | 10/92 | 1,136,000 |
| 08/83 | 562 | \| | 10/93 | 2,056,000 |

# # of www sites



Hobbes' Internet Timeline Copyright ©2010 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

| DATE | SITES | | DATE | SITES |
|---|---|---|---|---|
| 12/90 | 1 | | 06/95 | 23,500 |
| 12/91 | 10 | | 01/96 | 100,000 |
| 12/92 | 50 | | 06/96 | 252,000 |
| 06/93 | 130 | | 01/97 | 646,162 |
| 09/93 | 204 | | 06/97 | 1,117,259 |
| 10/93 | 228 | | 01/98 | 1,834,710 |
| 12/93 | 623 | | 06/98 | 2,410,067 |
| 06/94 | 2,738 | | 01/99 | 4,062,280 |
| 12/94 | 10,022 | | 07/99 | 6,598,697 |

# # of facebook account



Hobbes' Internet Timeline Copyright ©2010 Robert H Zakon
http://www.zakon.org/robert/internet/timeline/

# Towards all IP

**IPTV VOD**

**VoIP IP Telephony**

**Multimedia**

**High-perf networking**

INTERNET, HTTP

**Sensor networks**

**Grid computing**

**Interactive gaming**

Pervasive networking

# IP

E1/T1    X.25    FR    ATM    PSTN

# What's wrong?

The Internet has evolved from a <span style="color:red">wired network</span> for FTP, HTTP and e-mail...

" …the world has changed, the use of the Internet has changed and, fundamentally, the architecture has not evolved to take account of that. " (P. Howell, BT)

Internet

... to a fantastic infrastructure with a large variety of <span style="color:red">communicating devices</span> and high diversity of <span style="color:red">access</span> and traffic <span style="color:red">characteristics</span>

**Ubiqu**

**Mobil**
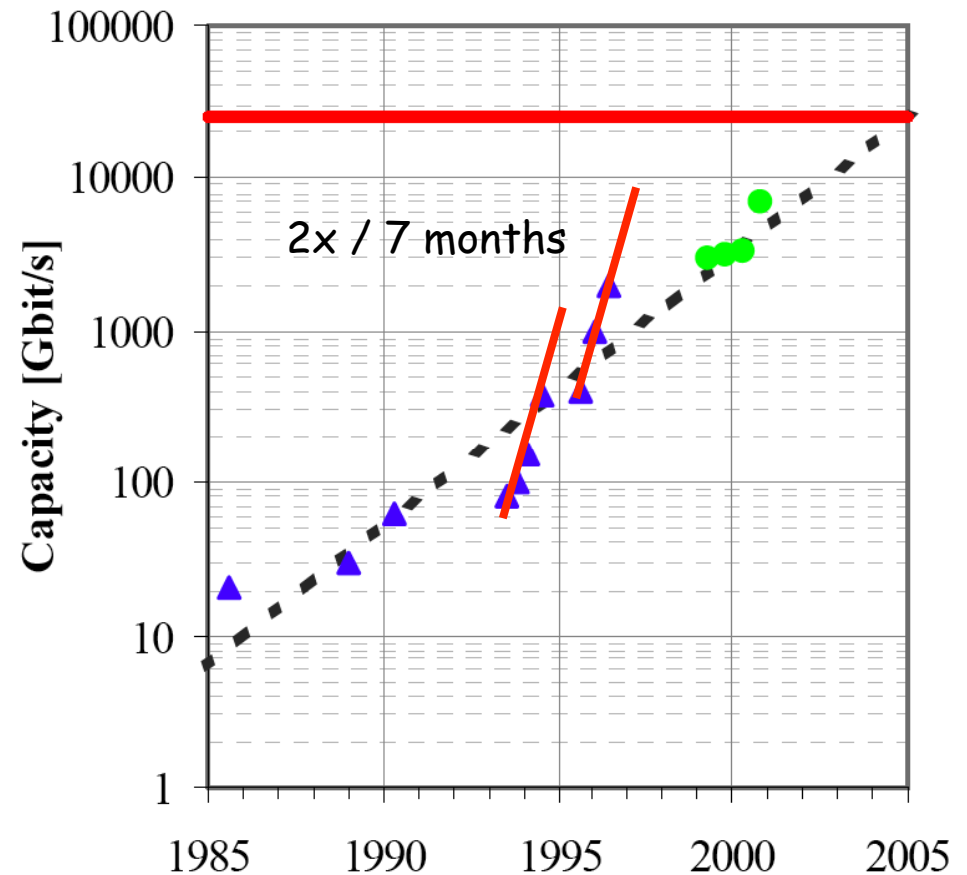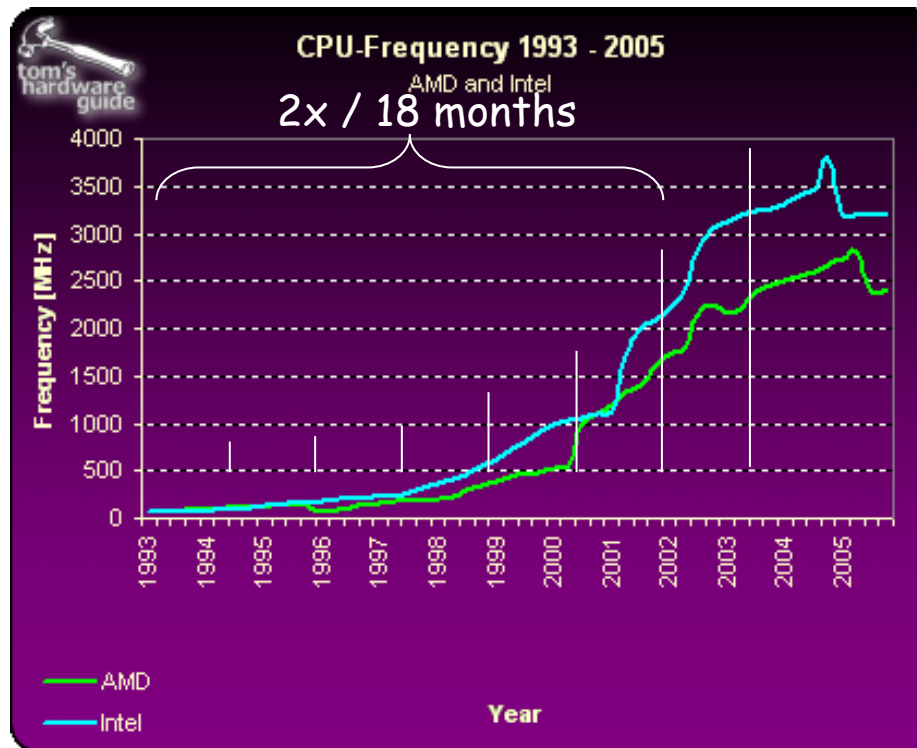
**Ad-Hoc**

**Telephony**

**MULTIMEDIA**

**Streaming**

# 1st revolution: Wireless Networks

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

Access Point

# 2nd revolution: going optical



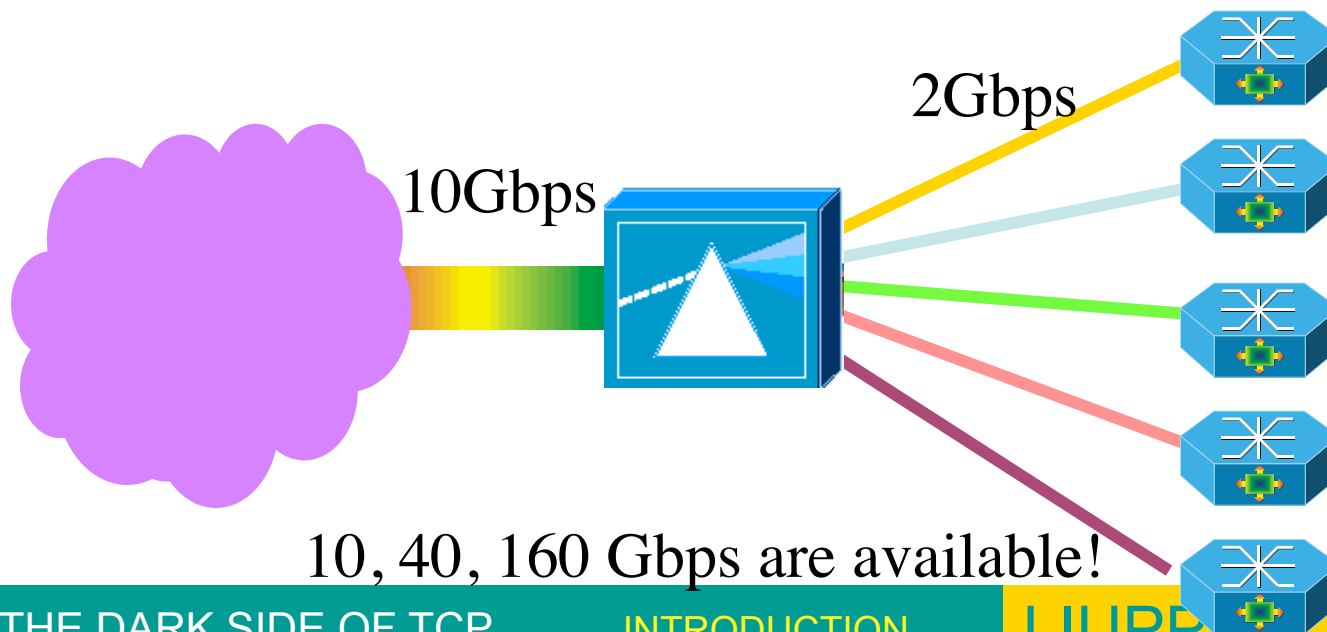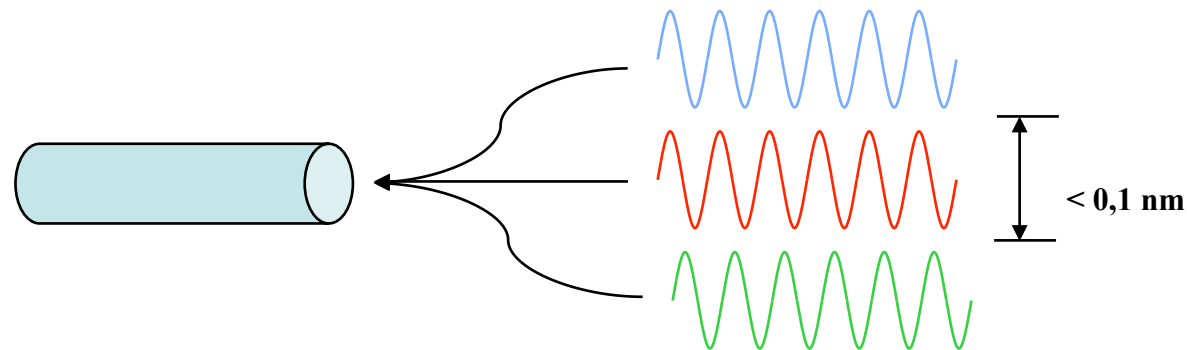CPU-Frequency 1993 - 2005
AMD and Intel

2x / 18 months

2x / 7 months

Capacity [Gbit/s]

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 ... er at the ...gh- ... ...n ...rs

**NEWS of May 31st, 2005**

US Fiber-to-the-home (FTTH) installations have grow ...

**NEWS of July,**

France Telecom wi ... an FTTH test- ... infrastructure in P... Gbps in downloa... 1.2Gbps in upl...

**NEWS of July, 2011**

France Telecom-Orange and Free will deploy FTTH in 5 millions home distributed in 1300 cities

2.5Gbp

campus

GigaEth

...ore

...160 Gbps

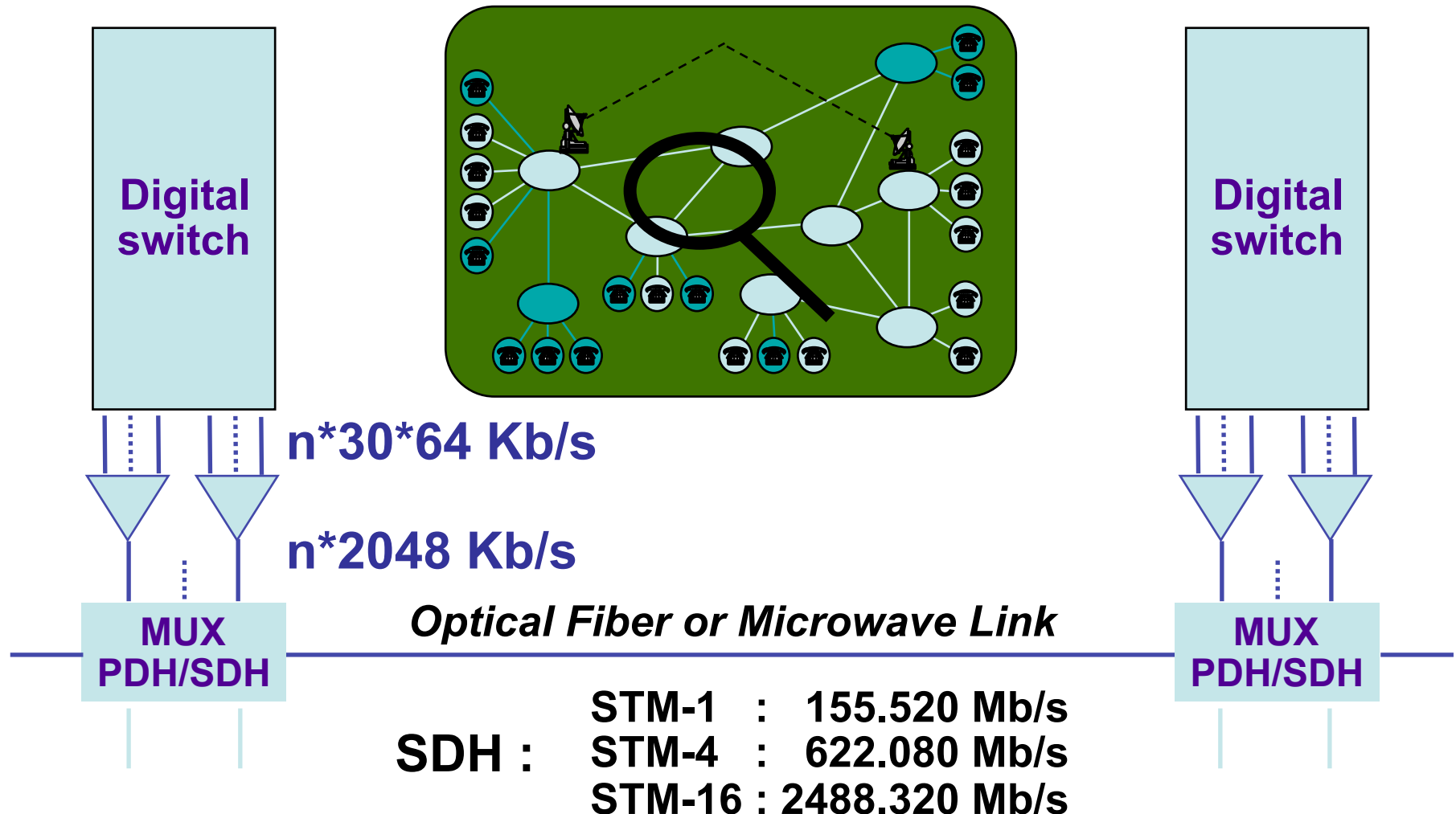# Latest news



source: FTTH Council

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



**Digital switch**

$n*30*64$ Kb/s

$n*2048$ Kb/s

**MUX PDH/SDH**

*Optical Fiber or Microwave Link*

**Digital switch**

**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

DCS or ADM

DCS or ADM

DCS or ADM

DCS or ADM

**rings**

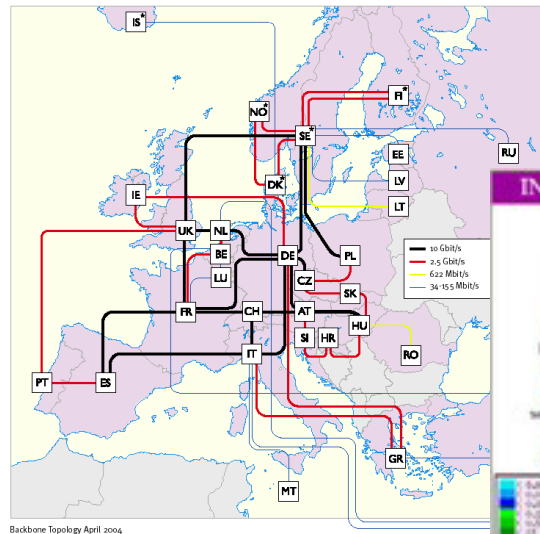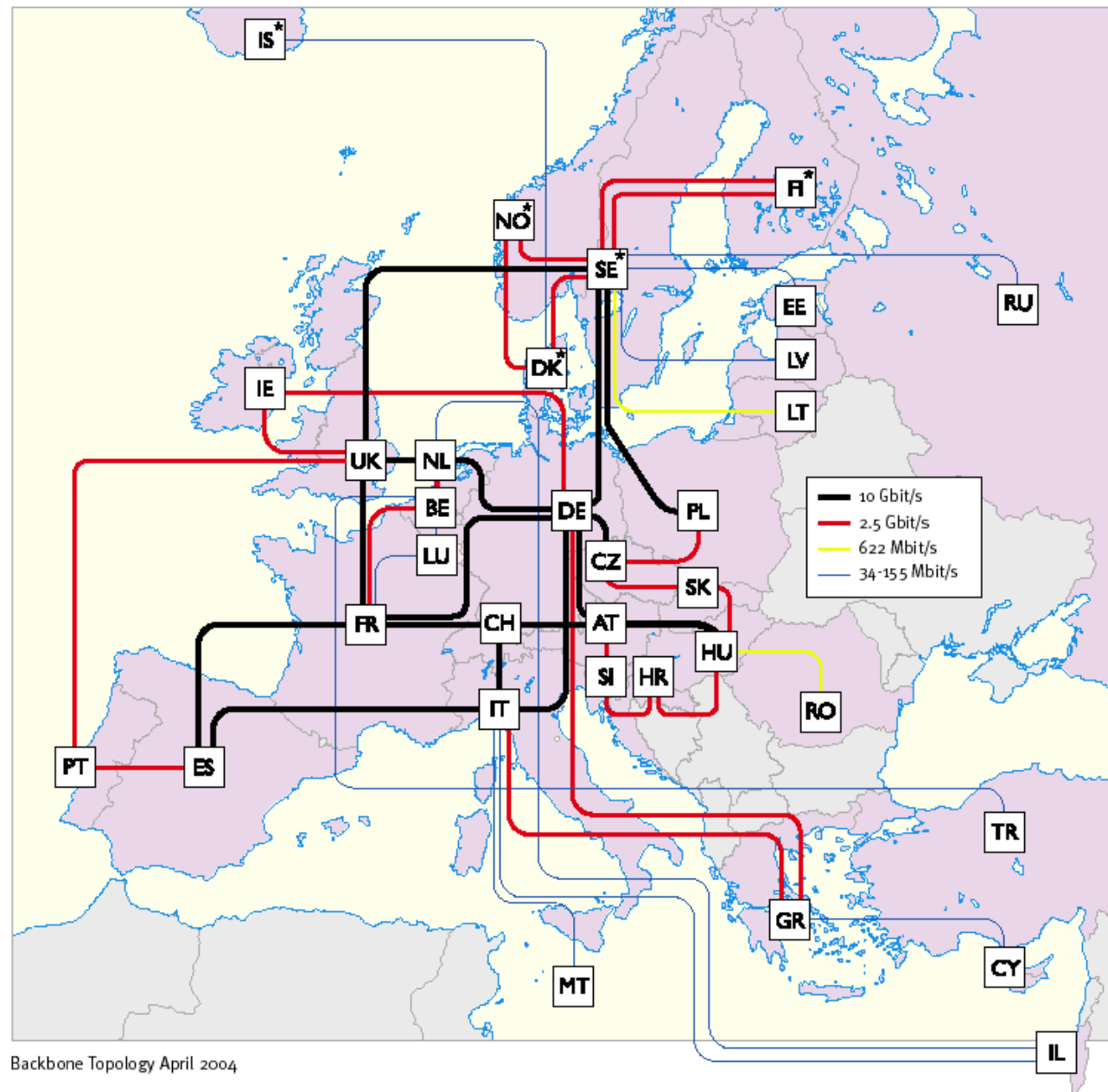**rings**

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

LIUPPA

# The new networks

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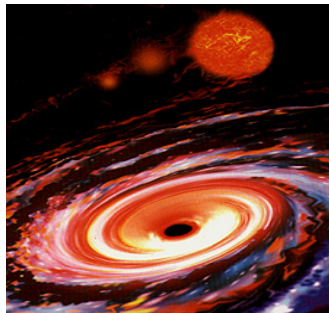- DREN
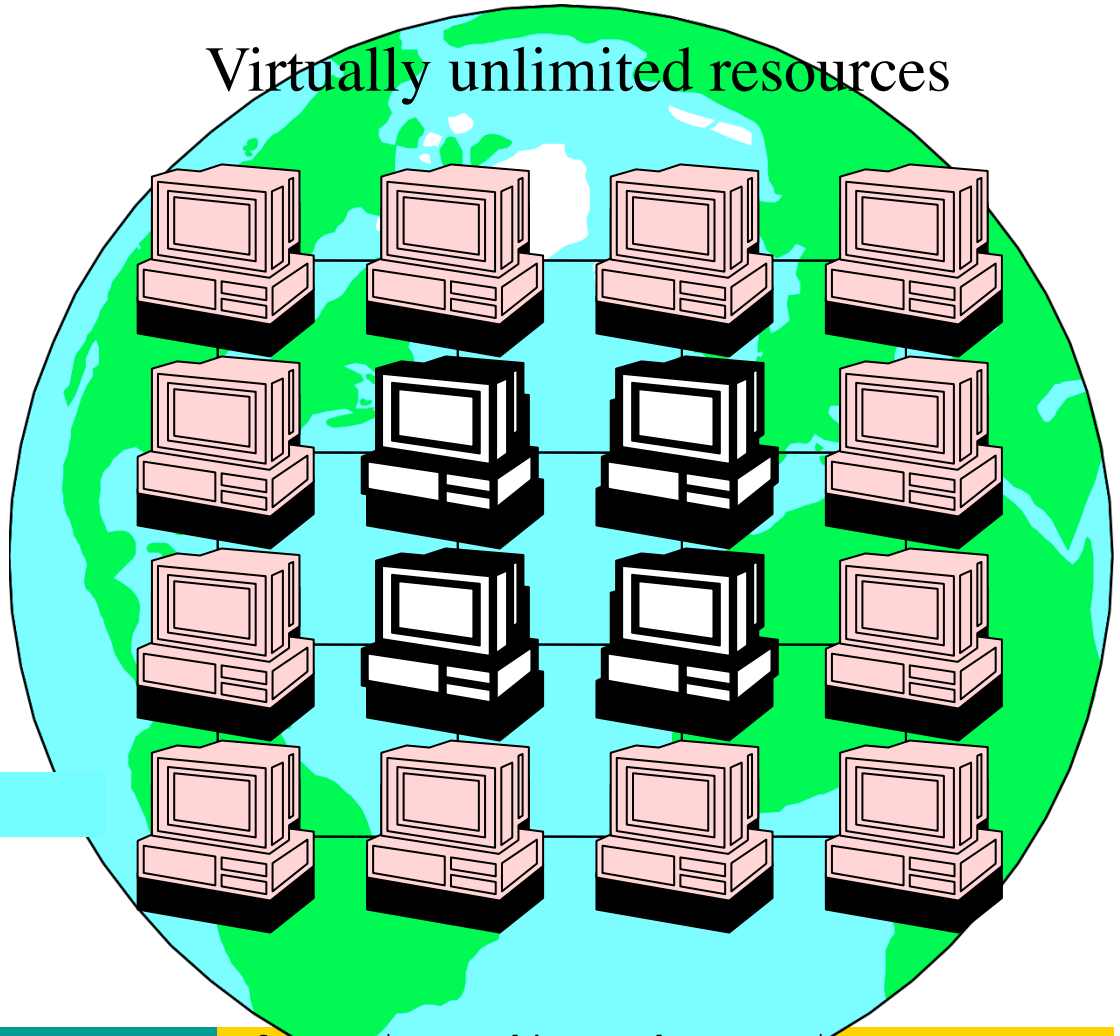- CA*NET
- GEANT
- DATATAG
- ...much more to come!

# GEANT



Backbone Topology April 2004

Legend:
- 10 Gbit/s
- 2.5 Gbit/s
- 622 Mbit/s
- 34-155 Mbit/s

# Computational grids

user application

1PFlops

Virtually unlimited resources

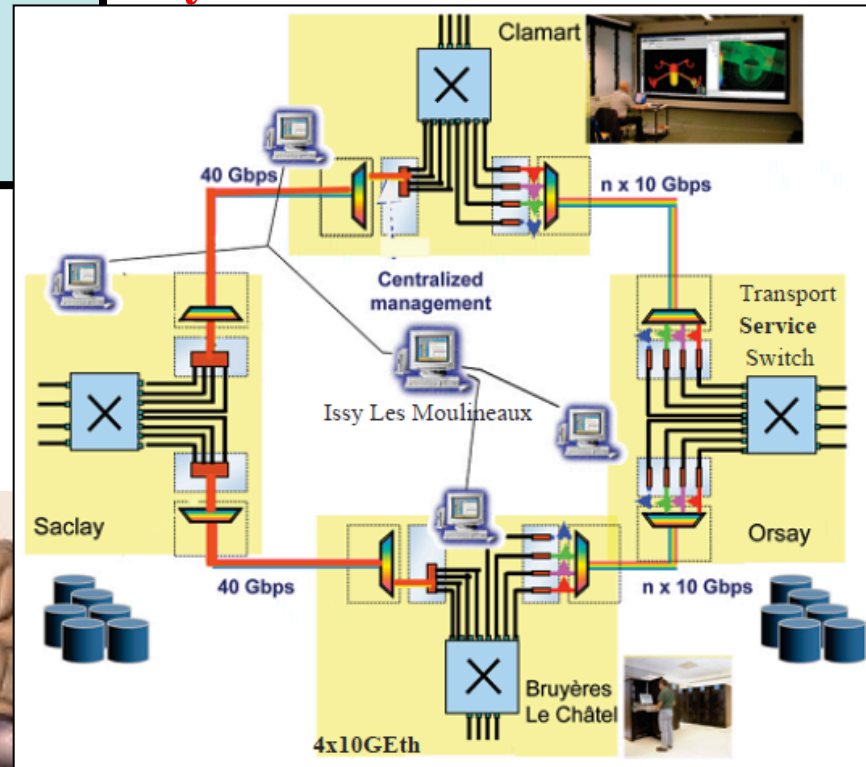from Dorian Arnold: Netsolve Happenings

# Real-time interactive large-scale scientific collaborations



**Large data transfers require very high bandwidth**

**Carriocas project (2006-2009) 40 Gb/wavelength**

**Multimodality brain mapping**

require the ability to process, share, and interactively visualize multiple 0Gbytes datasets!
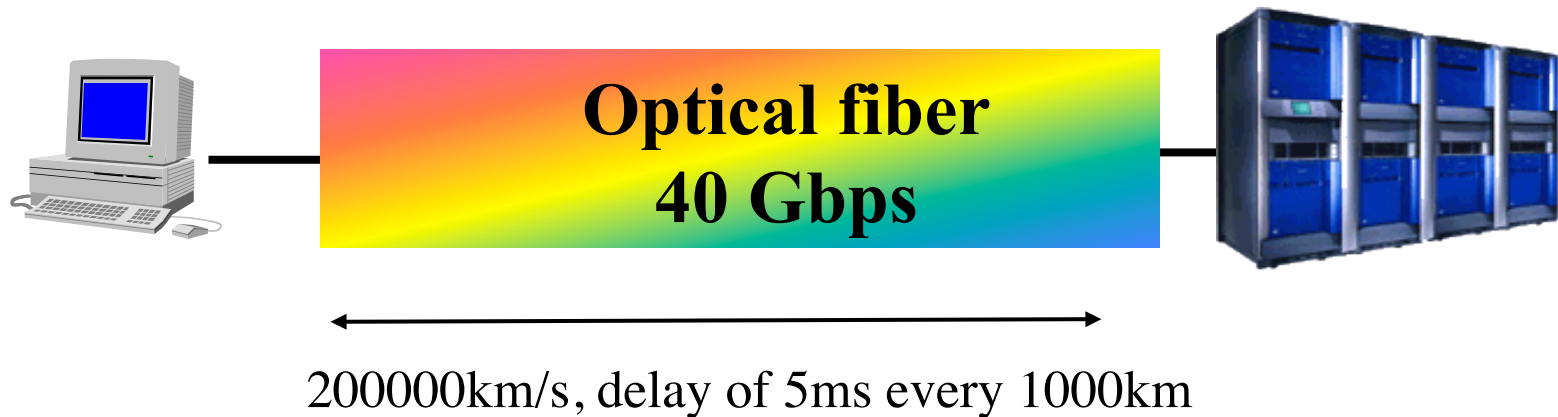
# Very High-Speed Networks



**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)

"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!

Packet losses

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

# Flow control
## prevents receiver's buffer overfow

**Packet Sent**

**Packet Received**

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

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

**App write**

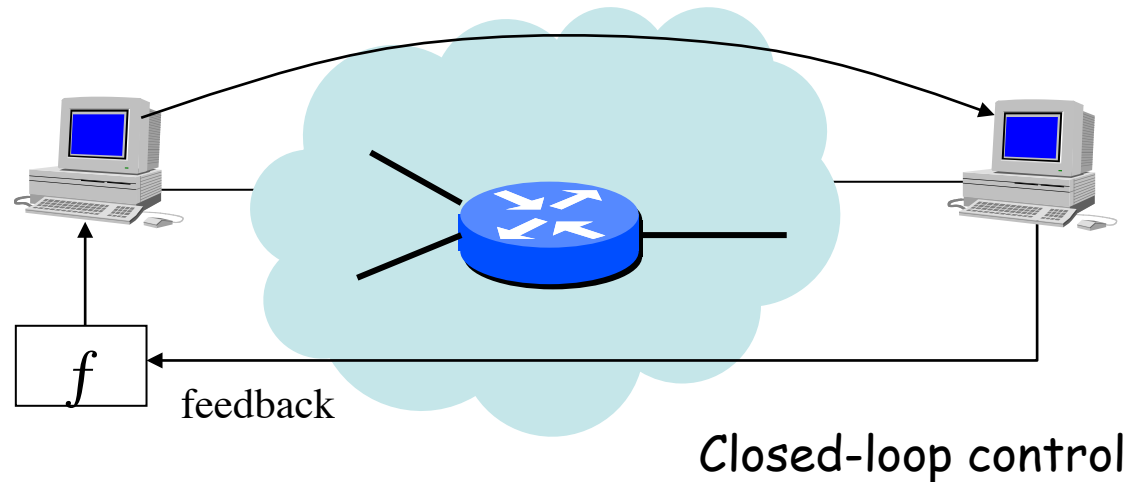**acknowledged**     **sent**     **to be sent**  **outside window**

# 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

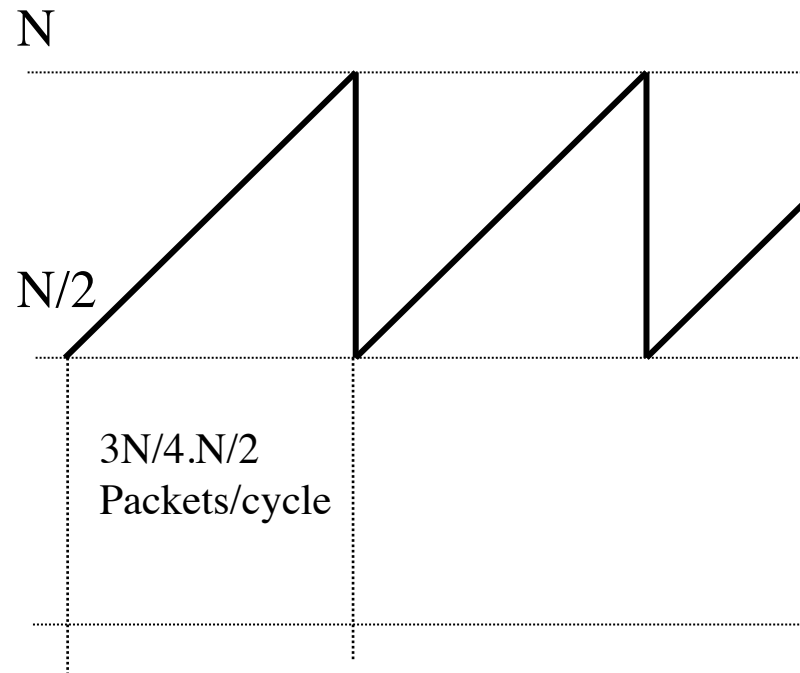# From the control theory point of view

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

# 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: $\dfrac{x_2}{x_1} = \dfrac{x_2\, b}{x_1\, b}$

- ❑ 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.

**Sender** **Receiver**

TIME

RTT

Time to transmit 3 packets

Packet #1
Packet #2
Packet #3

Packet #1 Ack.
Packet #2 Ack.
Packet #3 Ack.

Packet #4
Packet #5
Packet #6

Packet #4 Ack.
Packet #5 Ack.
Packet #6 Ack.

# 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

...01001011

0010100101010101010010101001011 01
010101010101001001111110100110111
0101001001001011101010101010001010
01010101010101010100011101110100
10110100010100111101010111

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
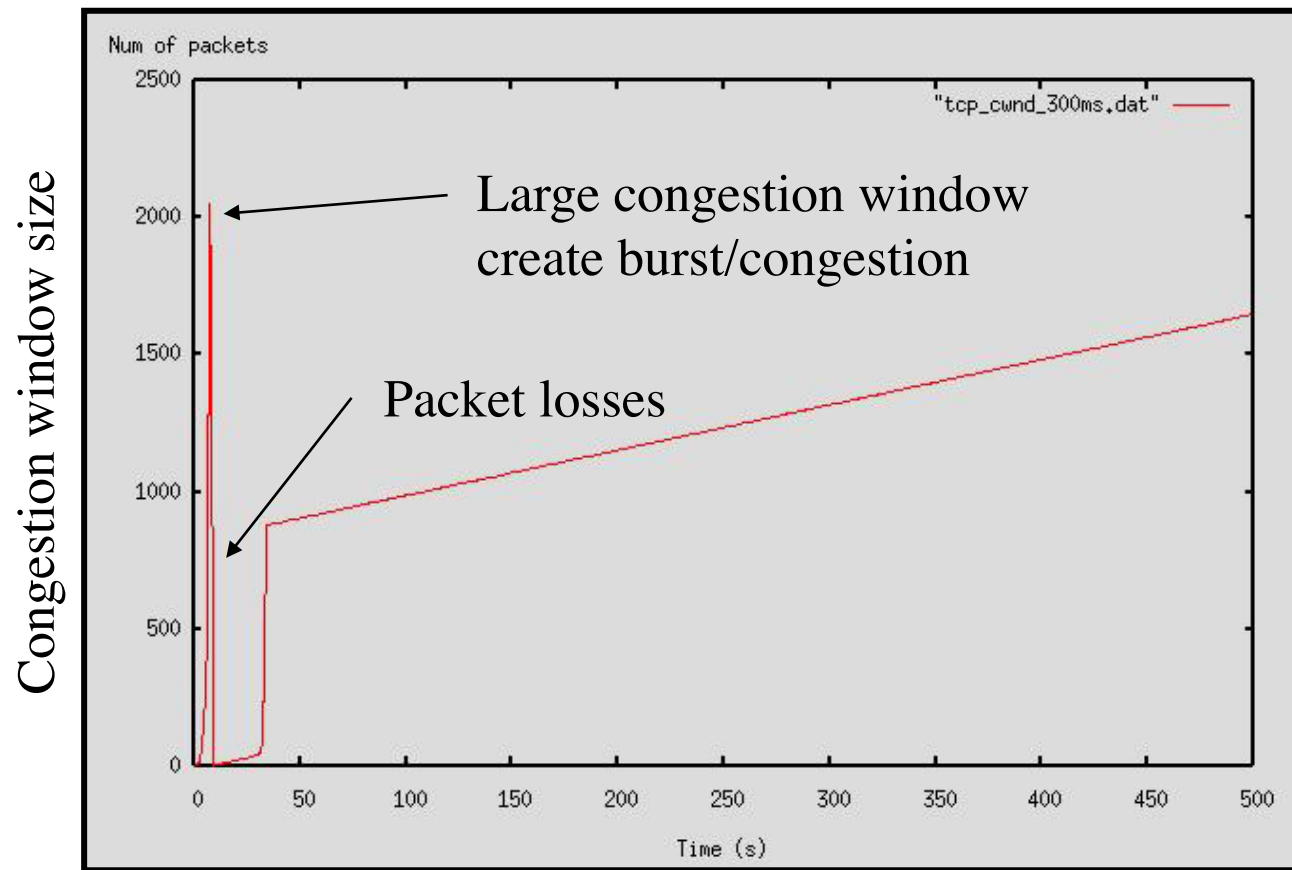
    ❑ Default block size

❑ 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.psc.edu/networking/projects/tcptune/)
- [http://www.web100.org/](http://www.web100.org/)
- [http://rdweb.cns.vt.edu/public/notes/win2k-tcpip.htm](http://rdweb.cns.vt.edu/public/notes/win2k-tcpip.htm)
- [http://www.sean.de/Solaris/soltune.html](http://www.sean.de/Solaris/soltune.html)
- [http://datatag.web.cern.ch/datatag/howto/tcp.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!**
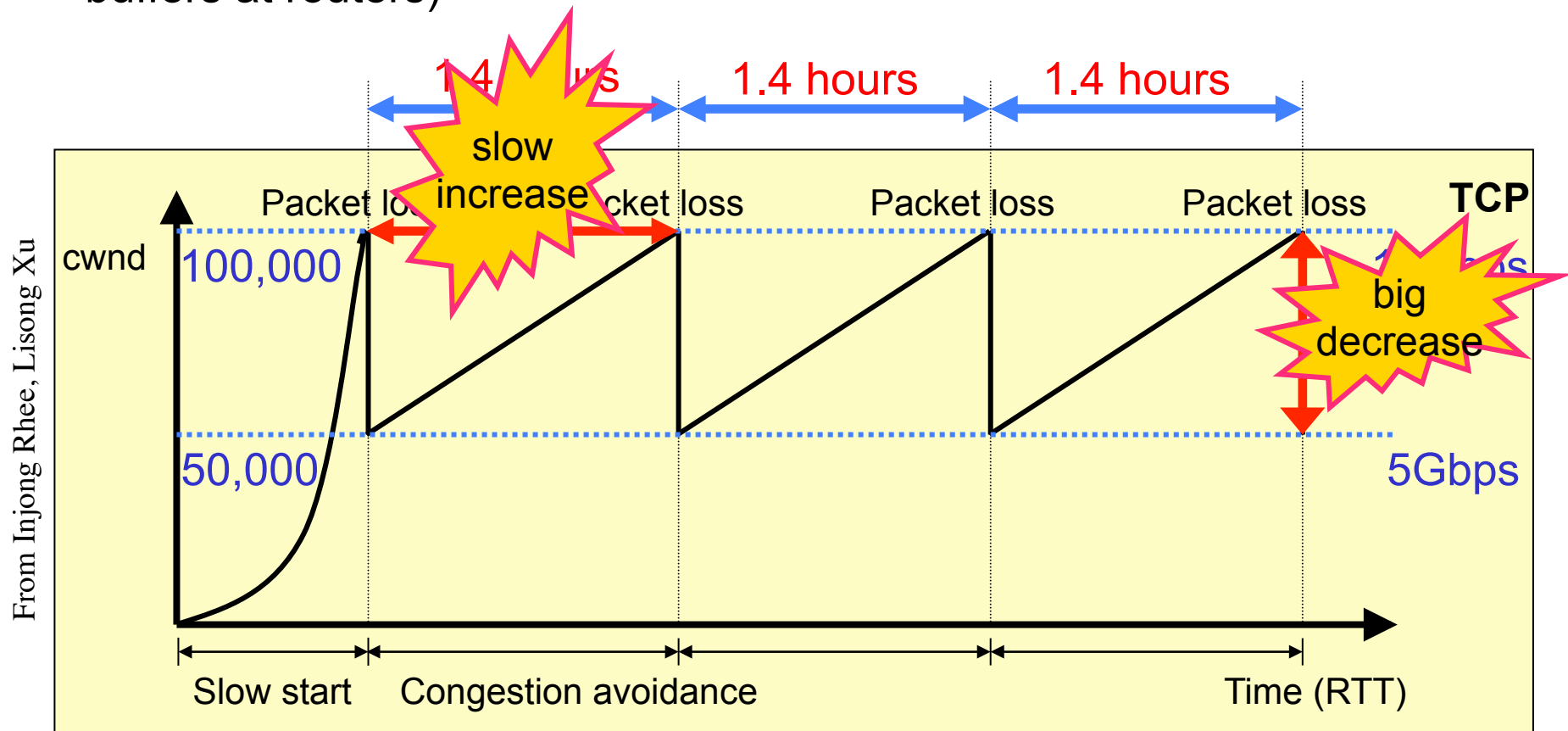
- Sustaining high congestion windows:
A Standard TCP connection with:
    - 1500-byte packets;
    - a 100 ms round-trip time;
    - a steady-state throughput of 10 Gbps;
would require:
    - an average congestion window of 83,333 segments;
    - and at most one drop (or mark) every 5,000,000,000 packets (or equivalently, at most one drop every 1 2/3 hours).

This is not realistic.

From S. Floyd

LIUPPA

# TCP rules:
# slow increase, big decrease

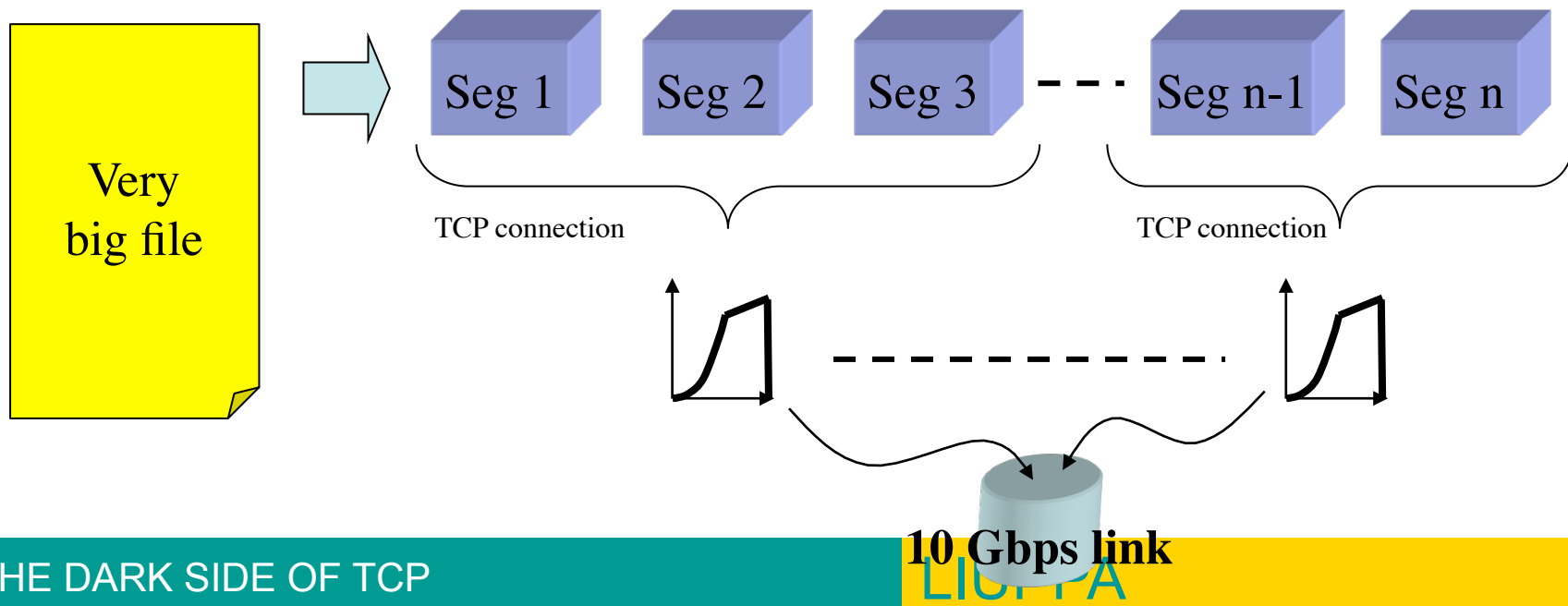A TCP connection with 1250-Byte packet size and 100ms RTT is running over a 10Gbps link (assuming no other connections, and no buffers at routers)

LIUPPA

# Going faster (cheating?)
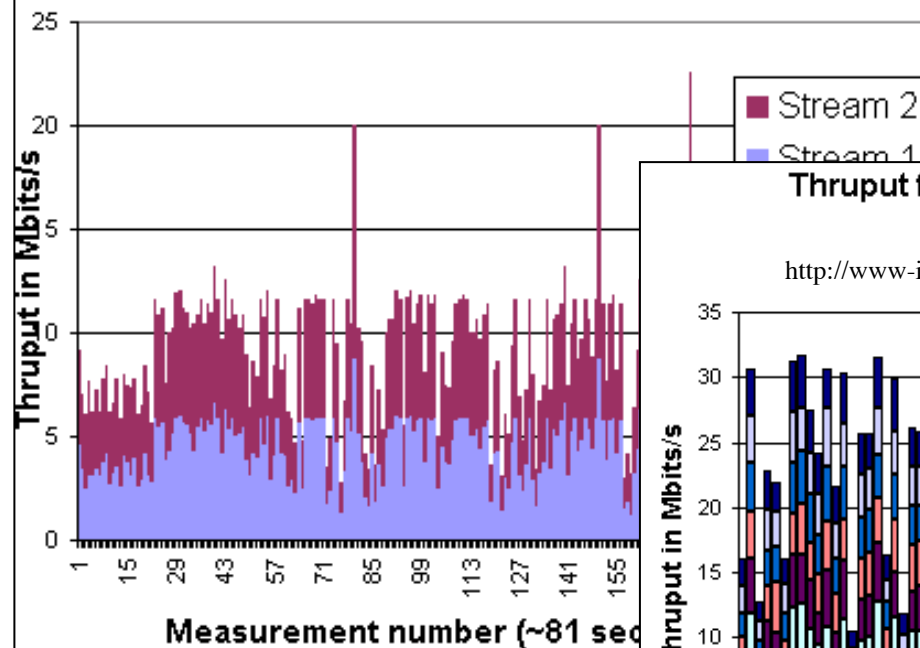## *n* flows is better than 1

❑ 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

# Some results from IEPM/ SLAC



Thruput SLAC to CERN with 256kByte window & 2 streams

Stream 2
Stream 1

Thruput in Mbits/s

Measurement number (~81 sec

**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

Thruput in Mbits/s

Stream 1   Stream 2   Stream 3   Stream 4
Stream 5   Stream 6   Stream 7   Stream 8

Measurement number (each separated by ~ 162 seconds)

# Multiple streams

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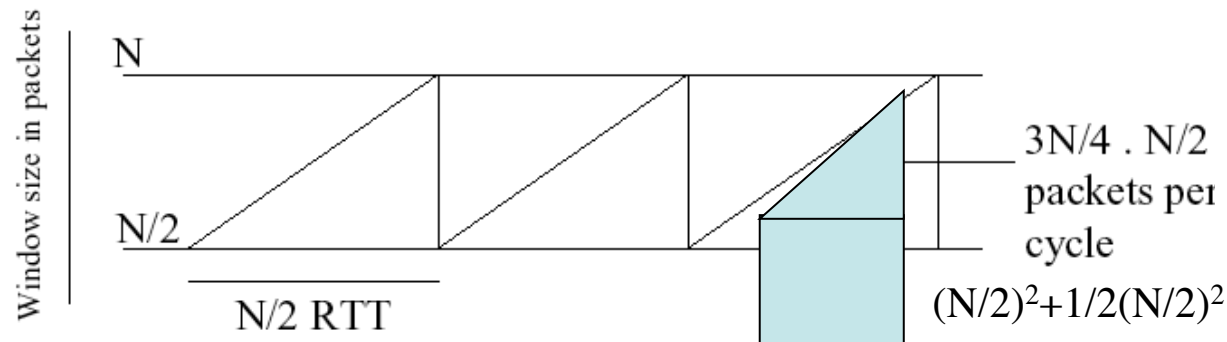  - bbFTP (http://doc.in2p3.fr/bbftp/)

# 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



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 = 3N²/8 = 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

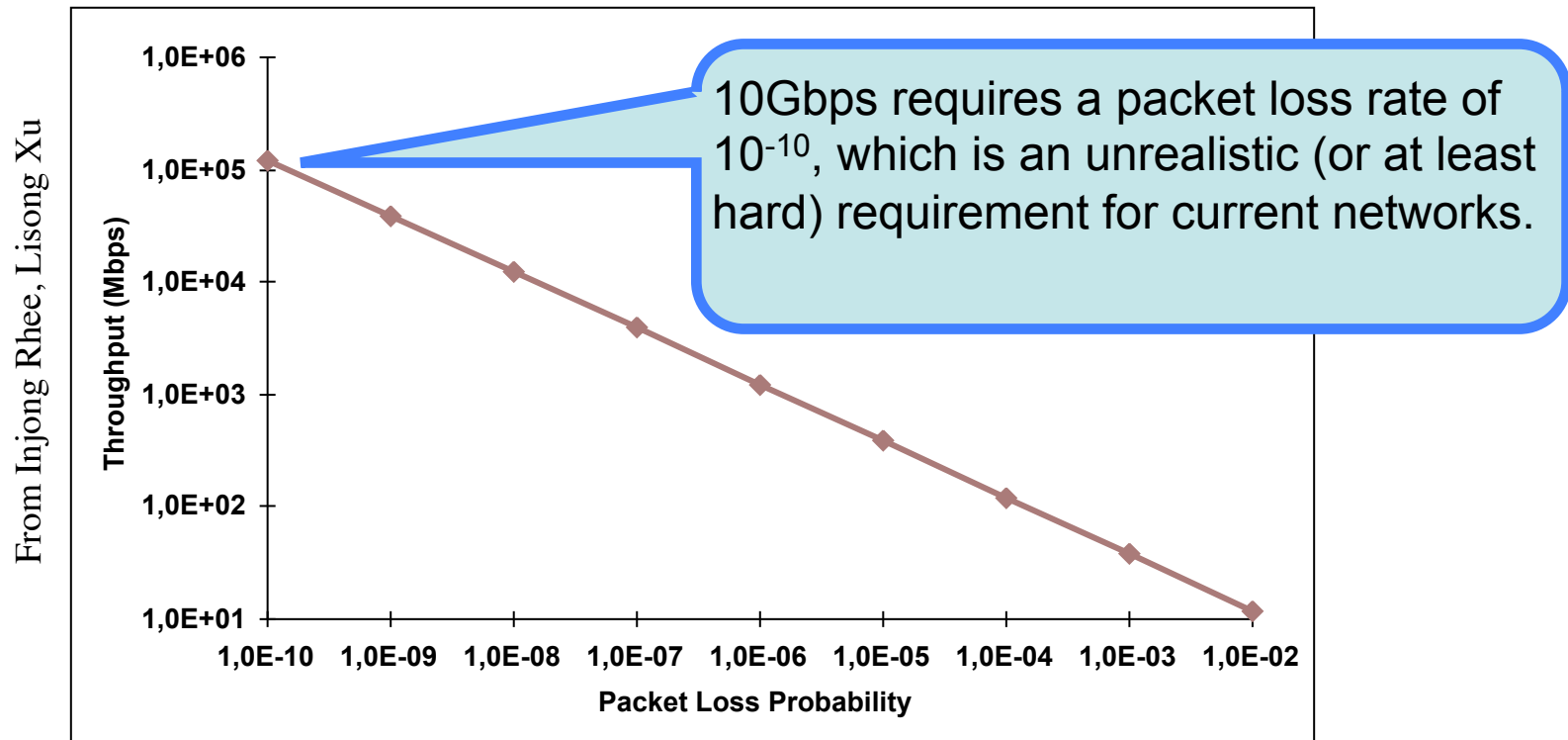$$Throughput = \frac{W}{RTT} = \sqrt{\frac{3}{2}}\frac{MTU}{RTT\sqrt{p}} = \sqrt{\frac{3}{2}}\frac{1}{RTT\sqrt{p}}$$

# 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

From Injong Rhee, Lisong Xu

**Throughput (Mbps)**

- 1,0E+06
- 1,0E+05
- 1,0E+04
- 1,0E+03
- 1,0E+02
- 1,0E+01

**Packet Loss Probability**

1,0E-10  1,0E-09  1,0E-08  1,0E-07  1,0E-06  1,0E-05  1,0E-04  1,0E-03  1,0E-02

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

# AIMD, general case

cwnd = cwnd + 1
⬇
cwnd = cwnd + 32

cwnd = cwnd * (1-1/2)
⬇
cwnd = cwnd * (1-1/8)

The throughput of AIMD is always about 13 times larger than that of TCP

NOT TCP Friendly!!!

**What's wrong?**

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

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

Throughput (Mbps)

1,0E+05
1,0E+04
1,0E+03
1,0E+02
1,0E+01

1,0E-07    1,0E-06    1,0E-05    1,0E-04    1,0E-03    1,0E-02

**Packet Loss Probability**

TCP
AIMD

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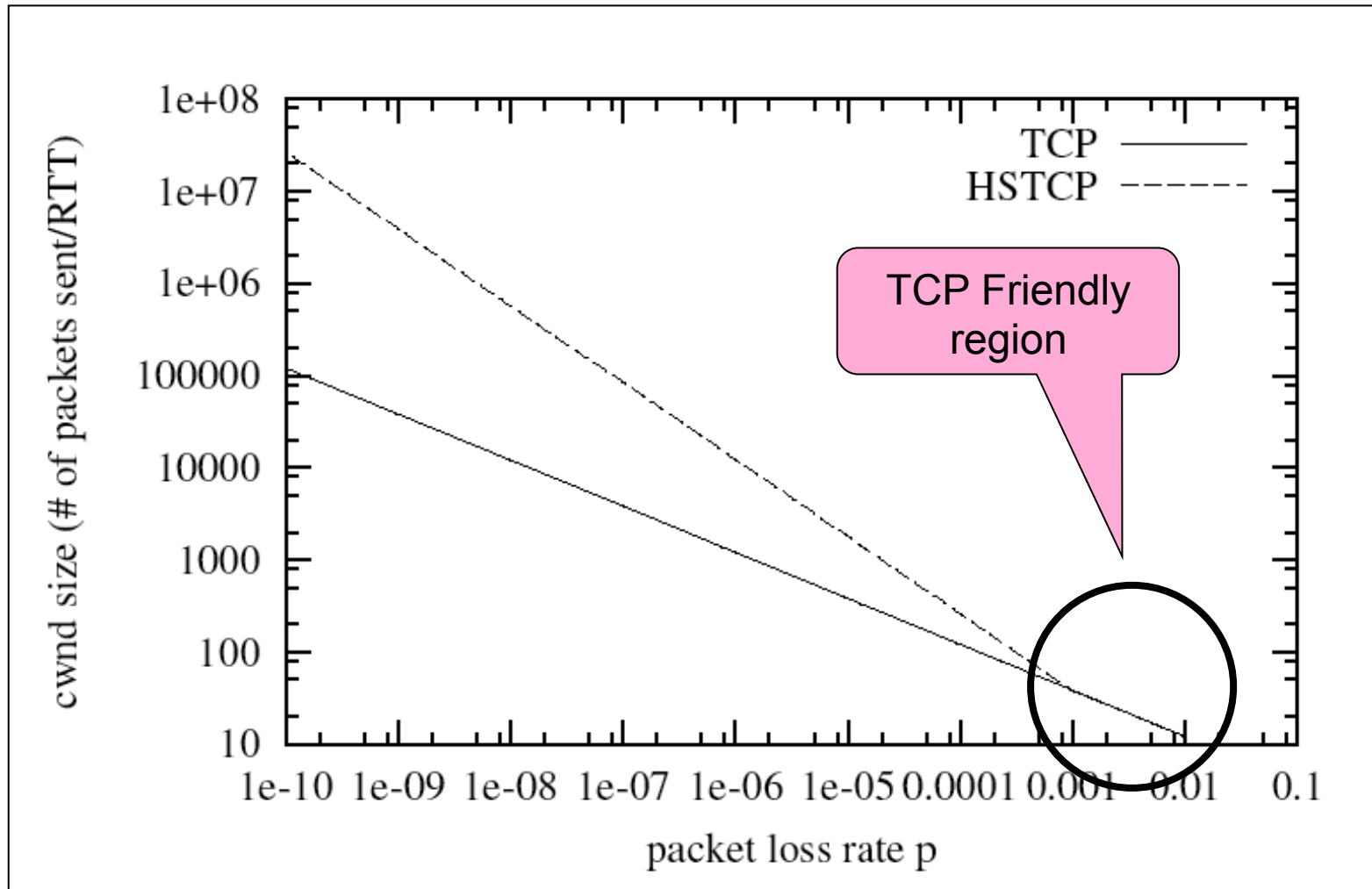
From draft-ietf-tsvwg-highspeed-01.txt

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



THE DARK SIDE OF TCP     BEYOND TCP     LIUPPA

# Relation with AIMD

□ **TCP-AIMD**
- □ Additive increase: a=1
- □ Multiplicative decrease: b=1/2

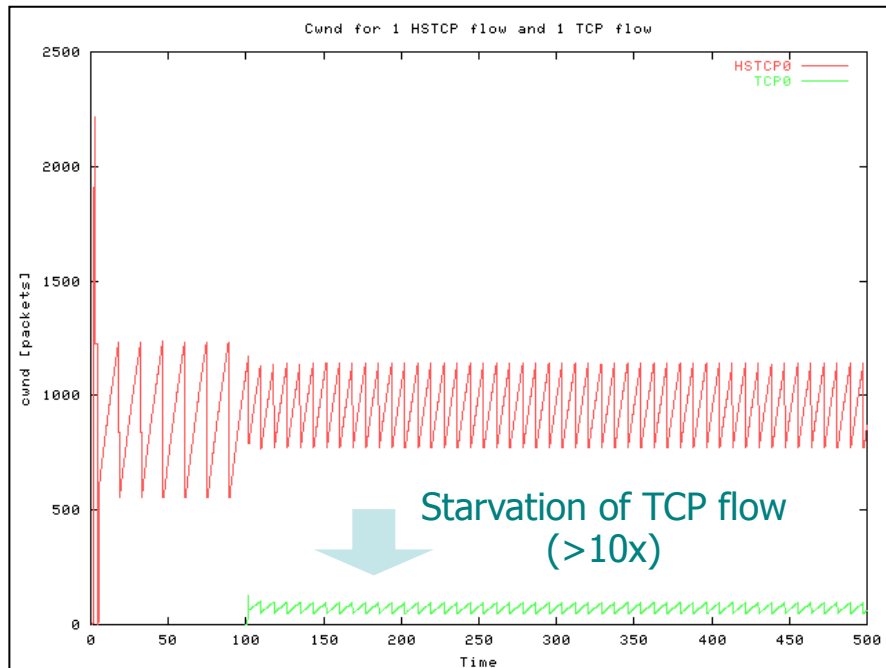□ **HSTCP-AIMD**
- □ Link a & b to congestion window size
- □ a = a(cwnd), b=b(cwnd)
- □ General rules
  - the larger cwnd, the larger the increment
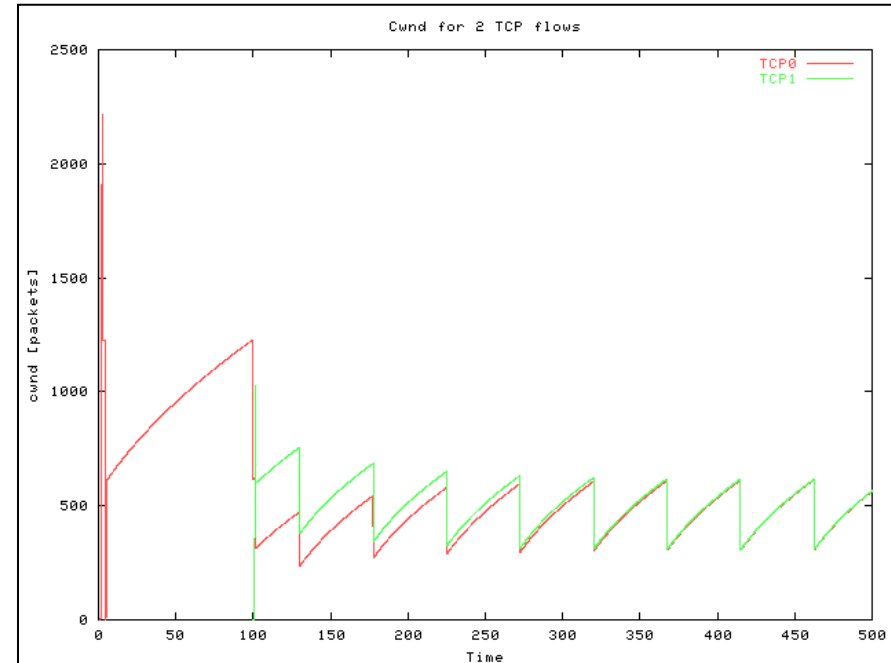  - The larger cwnd, the smaller the decrement

no loss:
cwnd = cwnd + 1

loss:
cwnd = cwnd*0.5

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

# Talking about dark side...



Cwnd for 1 HSTCP flow and 1 TCP flow

Starvation of TCP flow (>10x)



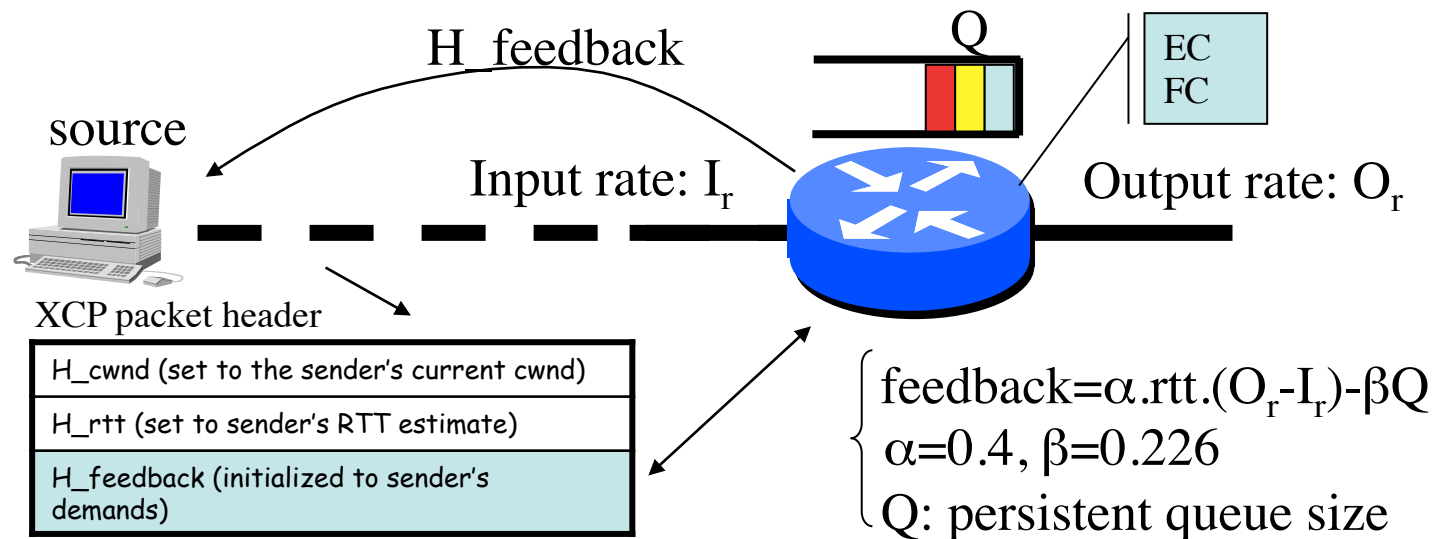Cwnd for 2 TCP flows

**1 HSTCP and 1 TCP flow**

**2 TCP flows**

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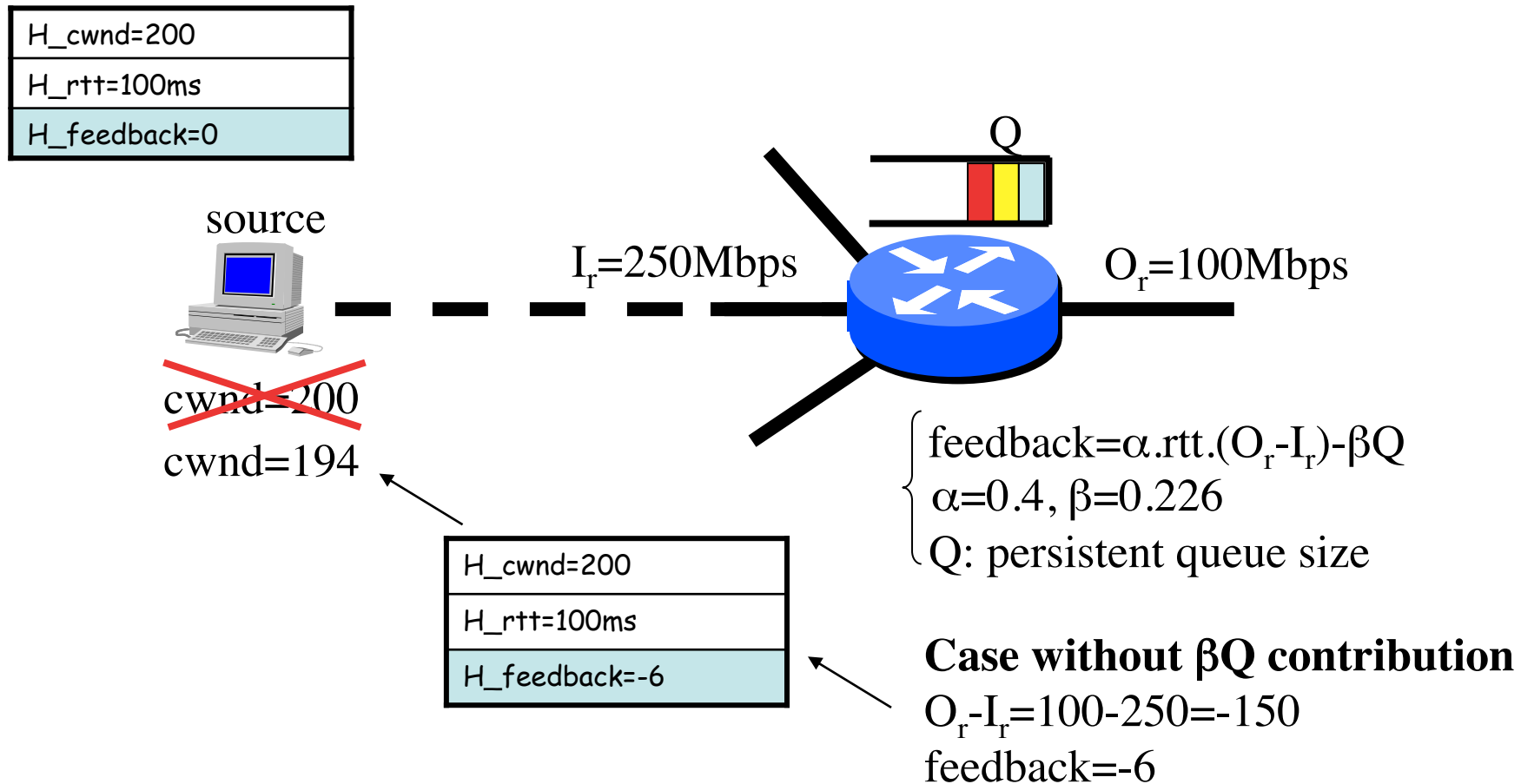Bottleneck BW=50Mbps
Qsize=BW*RTT
Qtype=DropTail

# 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

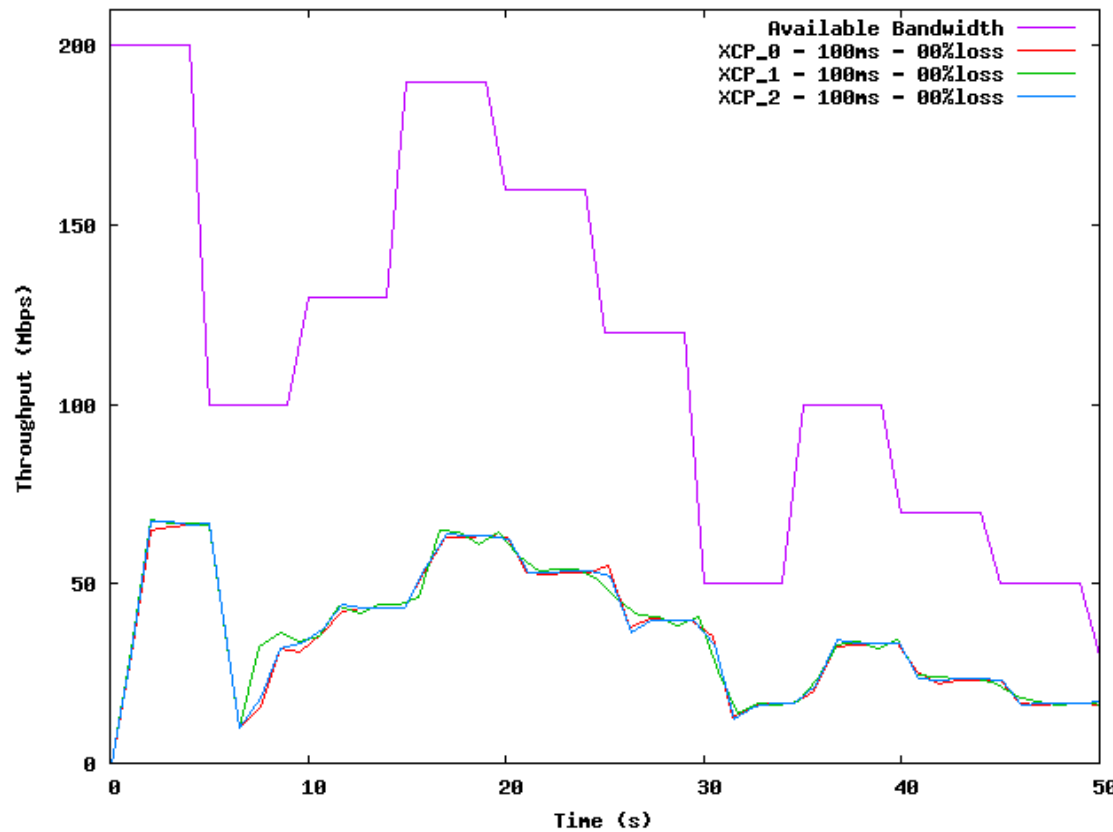❑ Feedback is sent back to the source in special fields of the packet header

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 |

source

$I_r$=250Mbps

$Q$

$O_r$=100Mbps

~~cwnd=200~~
cwnd=194

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

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

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

# XCP
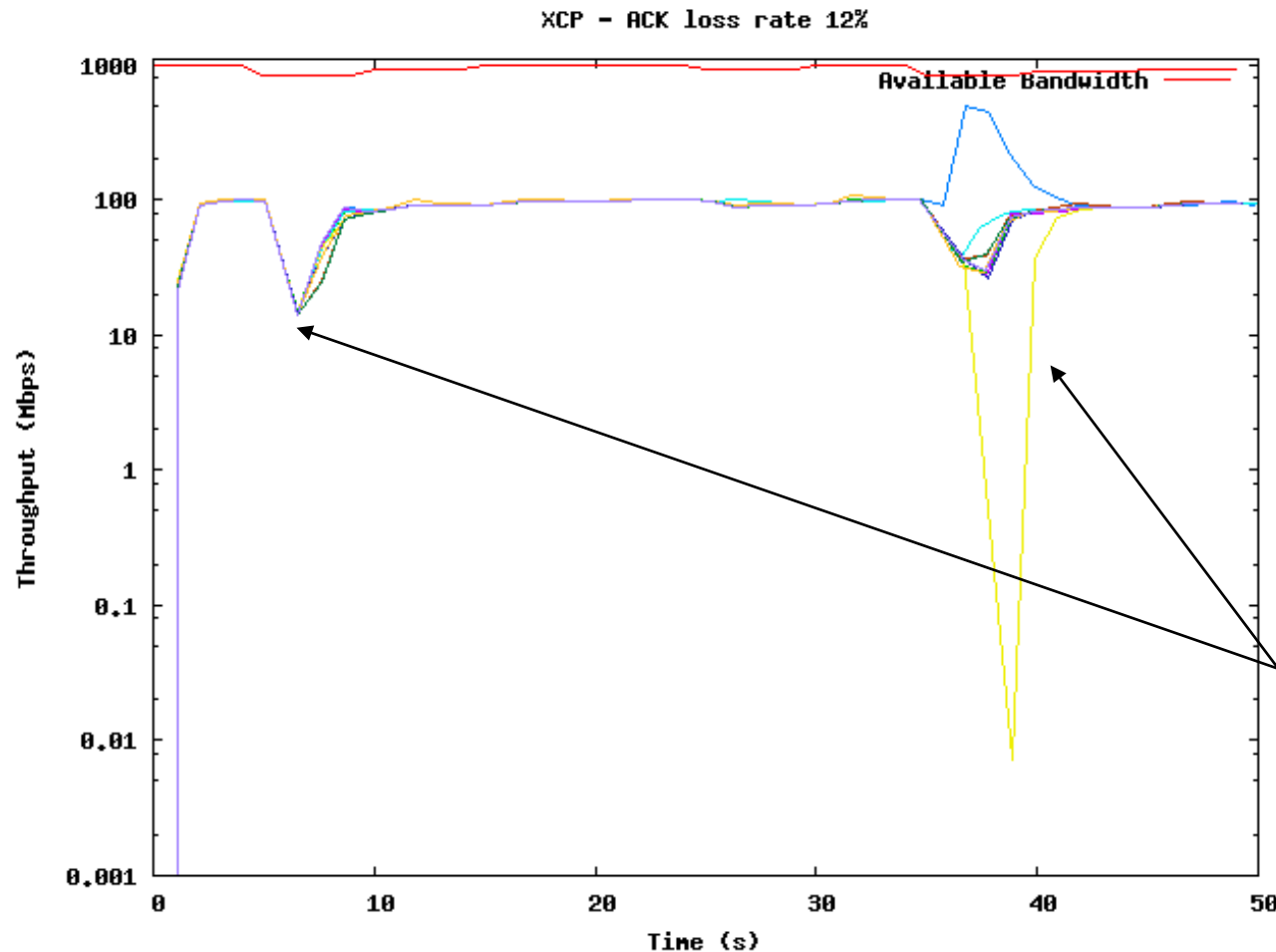## Variable bandwidth environments



Good fairness and stability even in variable bandwidth environments

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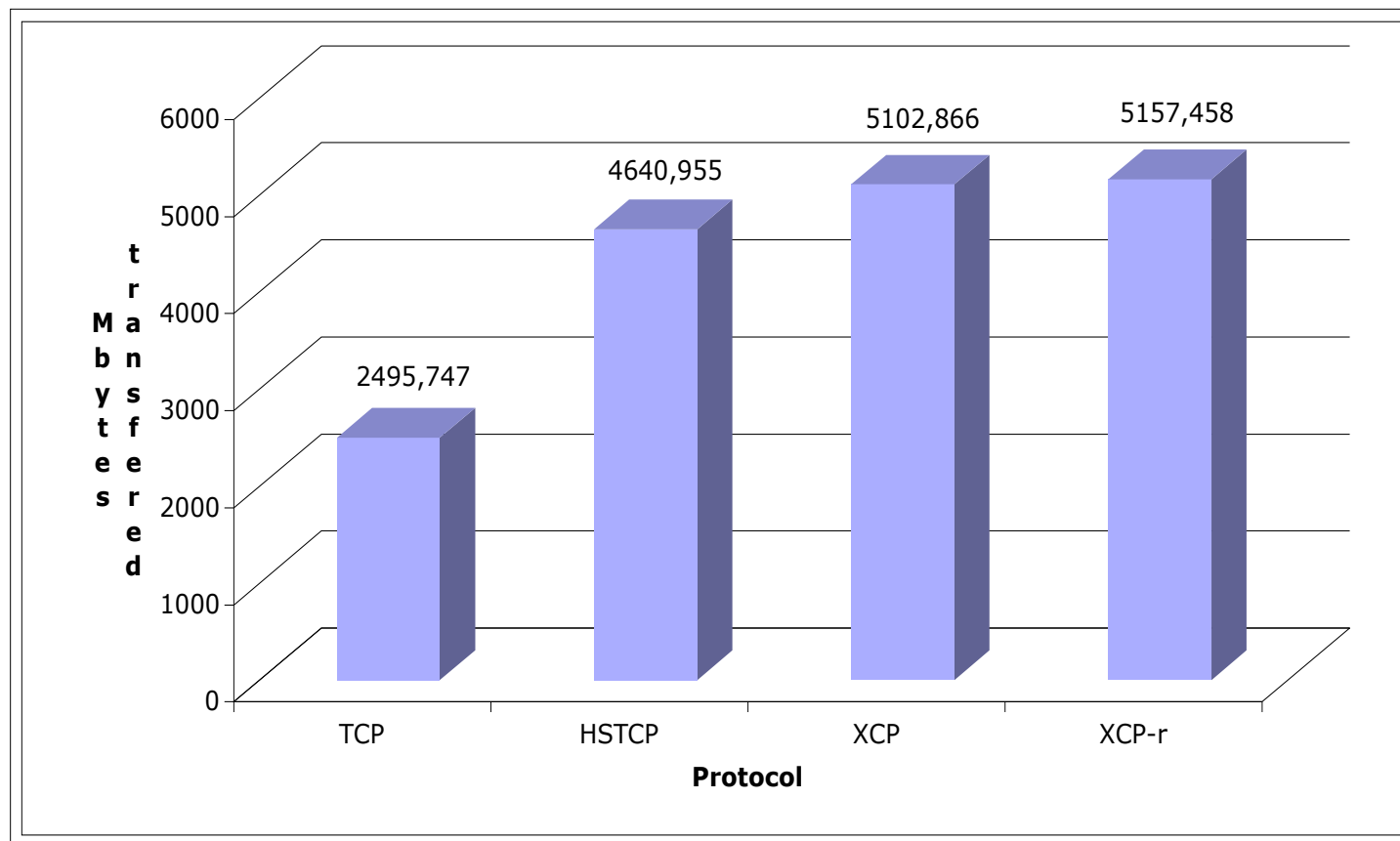## A more robust version of XCP



10 flows sharing a 1Gbps link

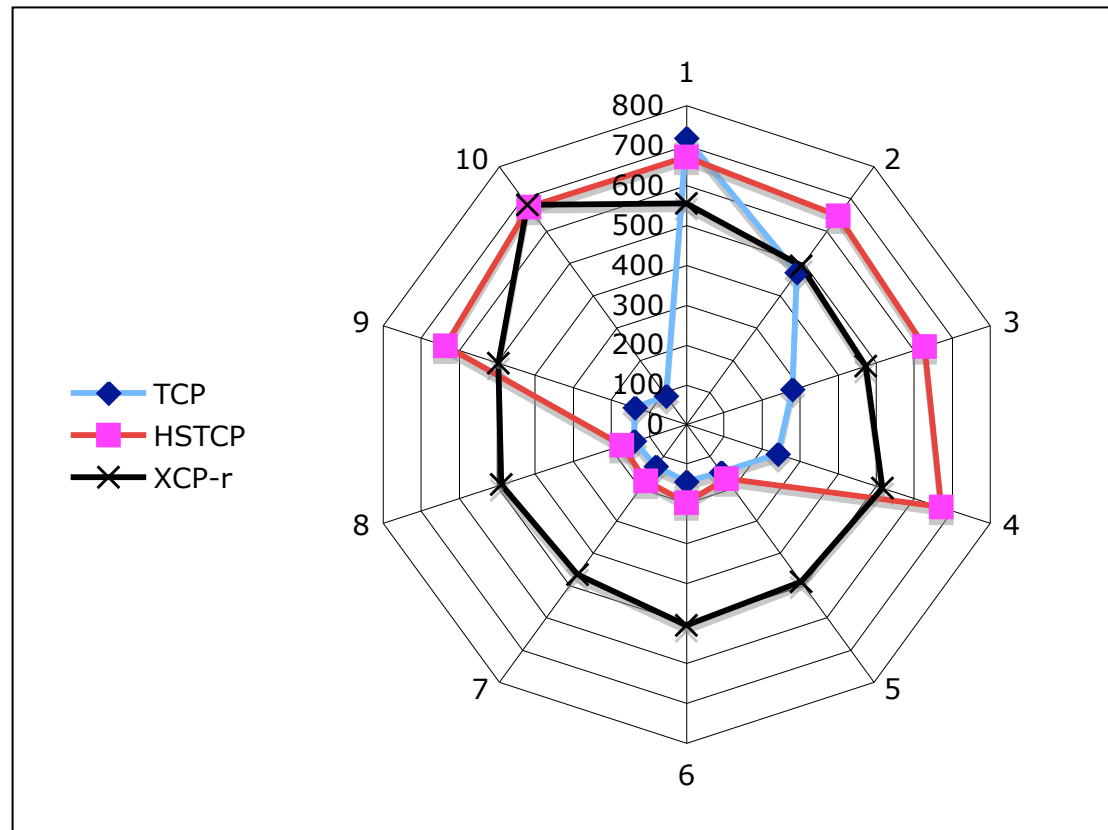Fast recovery after the timeouts and better fairness level

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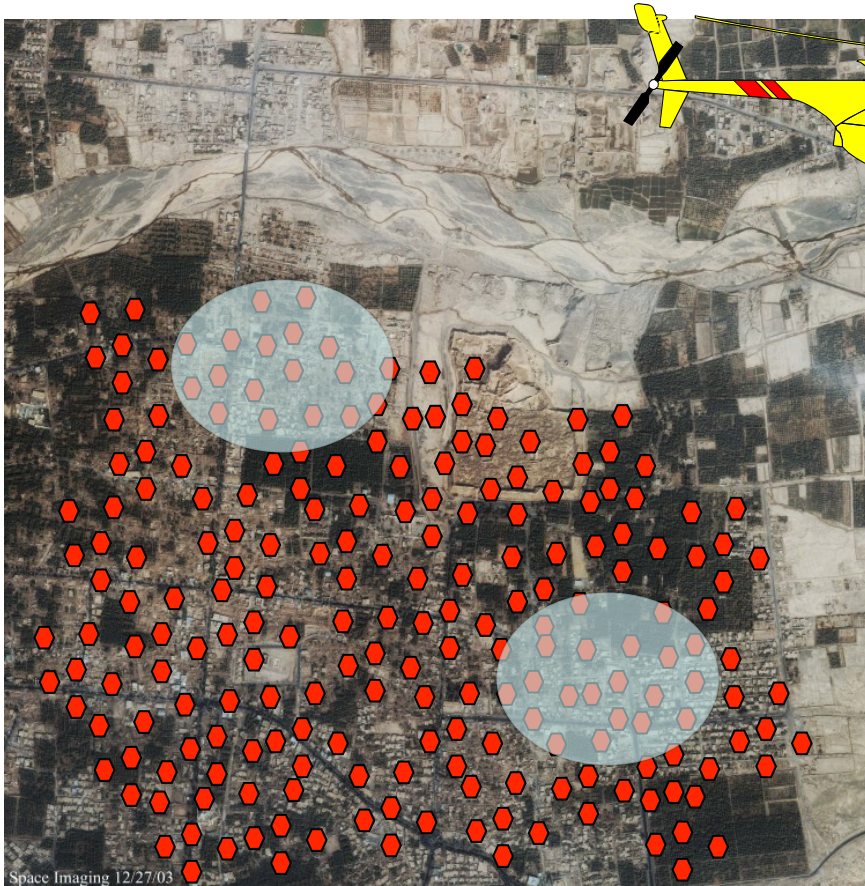- « The Web100 project will provide the software and tools necessary for end-hosts to automatically and transparently achieve high bandwidth data rates (100 Mbps) over the high performance research networks »
- Actually it's not limited to 100Mbps!
- Recommended solution for end-users to deploy and test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



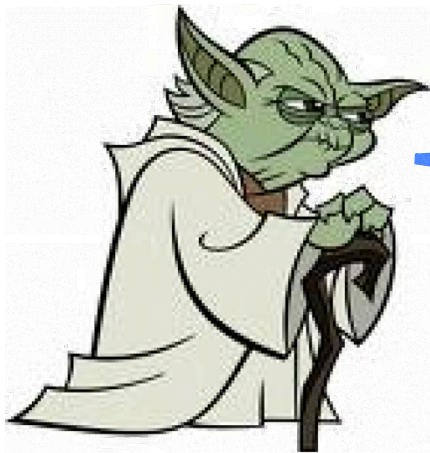Rapid deployment of fire detection systems in high-risk places

LIUPPA

# Conclusions

❑Understanding the dark side allows to move forwards!

❑However…

vanilla TCP

10GB file

40 Gbps

**MAY THE FORCE BE WITH YOU!**