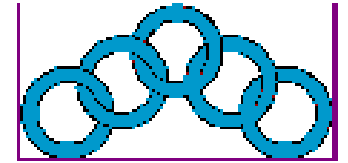




SLAC



Stanford Linear Accelerator Center



# Passive and Active Monitoring on a High-performance Network

*Les Cottrell, Warren Matthews, Davide Salomoni,  
Connie Logg – SLAC*

[www.slac.stanford.edu/grp/scs/net/talk/pam-apr01/](http://www.slac.stanford.edu/grp/scs/net/talk/pam-apr01/)

Presented at PAM-2001, Amsterdam April 23-24, 2001

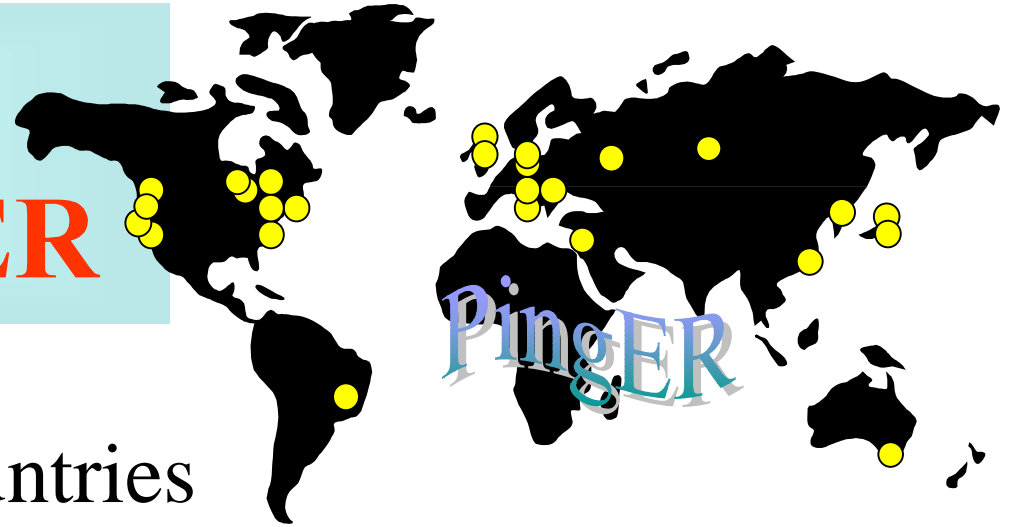
Partially funded by DOE/MICS Field Work Proposal on Internet End-to-end Performance Monitoring (IEPM), also supported by IUPAP



# Outline

- Results from active monitoring with PingER:
  - RTT, Loss, “jitter”
- Passive border monitoring results
- High perf throughput
  - achieving, measuring and impact
- Simulation of high perf throughput

# Active WAN Monitoring/PingER

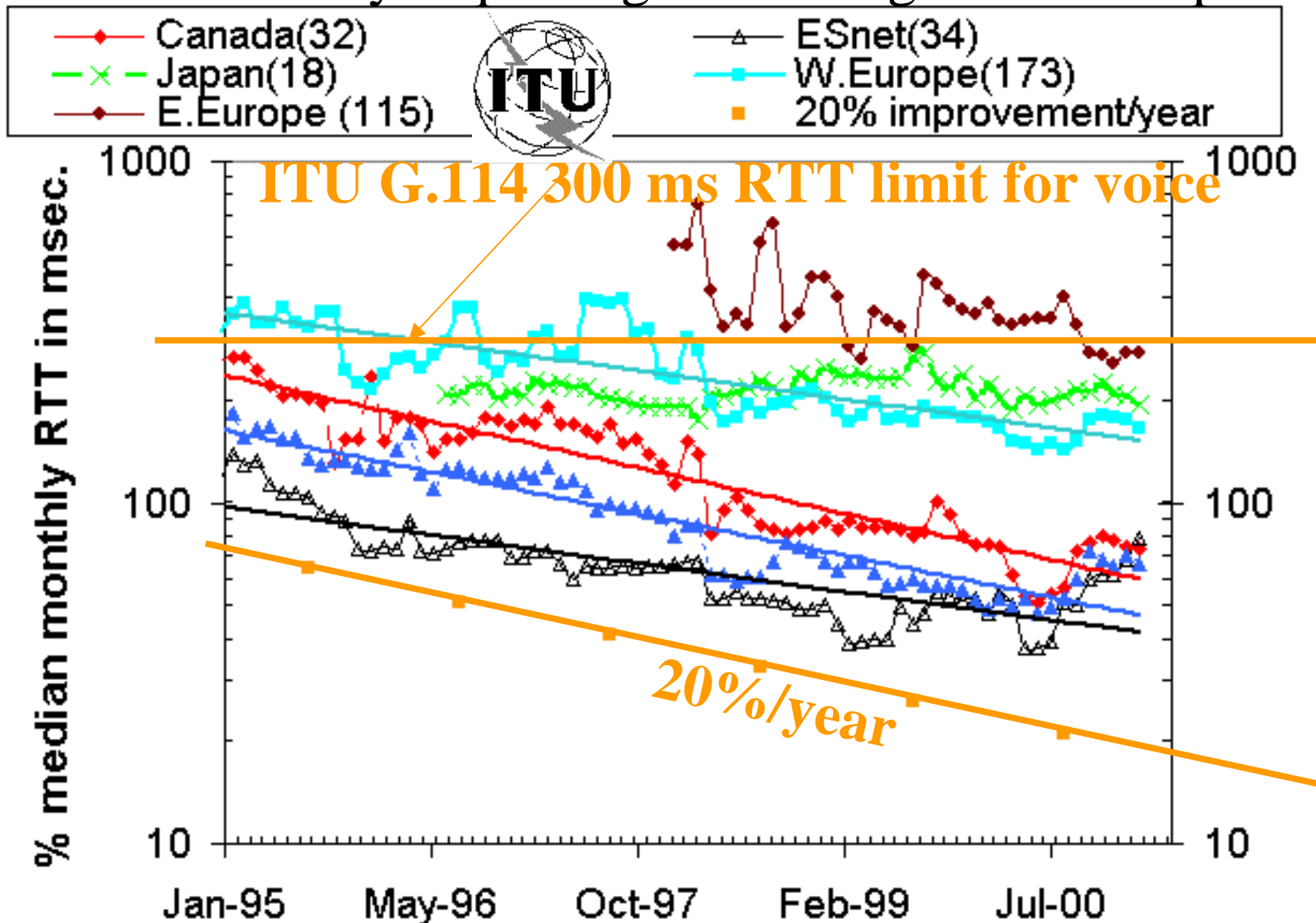


- Measurements from
  - 32 monitors in 14 countries
  - Over 600 remote hosts in over 72 countries
  - Over 3300 monitor-remote site pairs
  - Measurements go back to Jan-95
  - Reports on RTT, loss, reachability, IPDV, throughput, reordering, duplicates, looking at CLP (for bursty losses)...
- Uses ubiquitous “ping” facility of TCP/IP
- Countries monitored
  - Contain 78% of world population
  - 99% of online users of Internet

# RTT from ESnet to Groups of Sites

$RTT \sim \text{distance}/(0.6 * c) + \text{hops} * \text{router delay}$

Router delay = queuing + clocking in & out + processing



# RTT Region to Region

OK

White 0-64ms

Green 64-128ms

Yellow 128-256ms

NOT OK

Pink 256-512ms

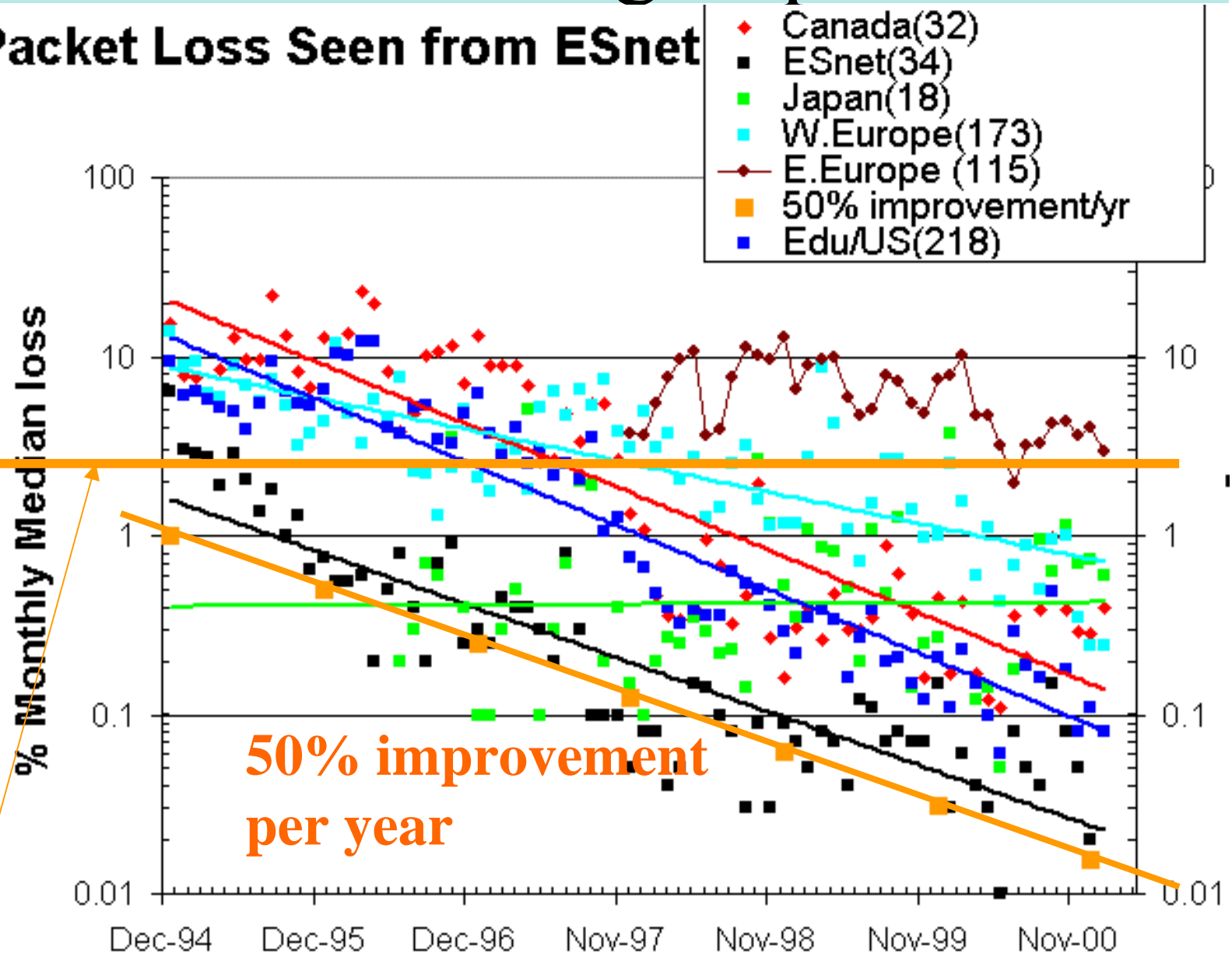
Red > 512ms

<u>WORLD</u>	<u>Australasia</u>	<u>East Europe</u>	<u>North America</u>	<u>West Europe</u>	<u>South America</u>	<u>Asia</u>
<u>Australasia</u>	3.95	714.74	300.68	454.68	389.69	373.59
<u>East Europe</u>		359.03	235.66	87.37	278.01	319.64
<u>North America</u>		244.24	69.44	153.23	223.06	203.83
<u>West Europe</u>		385.14	163.32	42.97	260.47	290.68
<u>South America</u>		626.39	421.45	590.69	18.93	780.00
<u>Asia</u>		472.57	327.85	321.99	447.02	24.00
<u>Africa</u>		770.90	804.67	804.15		
<u>Aria</u>			772.00	416.88		
<u>null</u>			501.30			
<u>Middle East</u>			1108.97			
<u>Central America</u>			436.00			

OK within regions, N. America OK with Europe, Japan

# Loss seen from US to groups of Sites

## Packet Loss Seen from ESnet



50% improvement  
per year

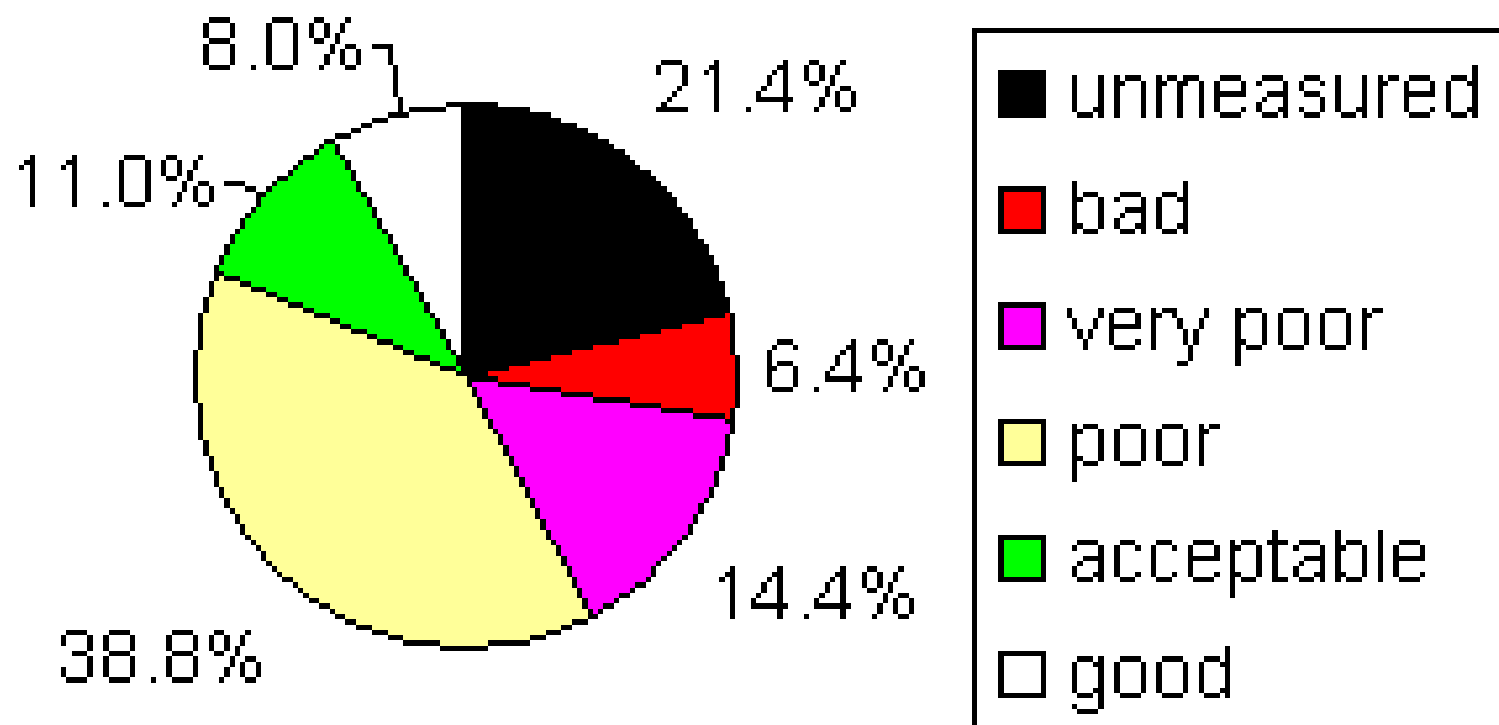
ETSI limit for loss (assumes random losses)



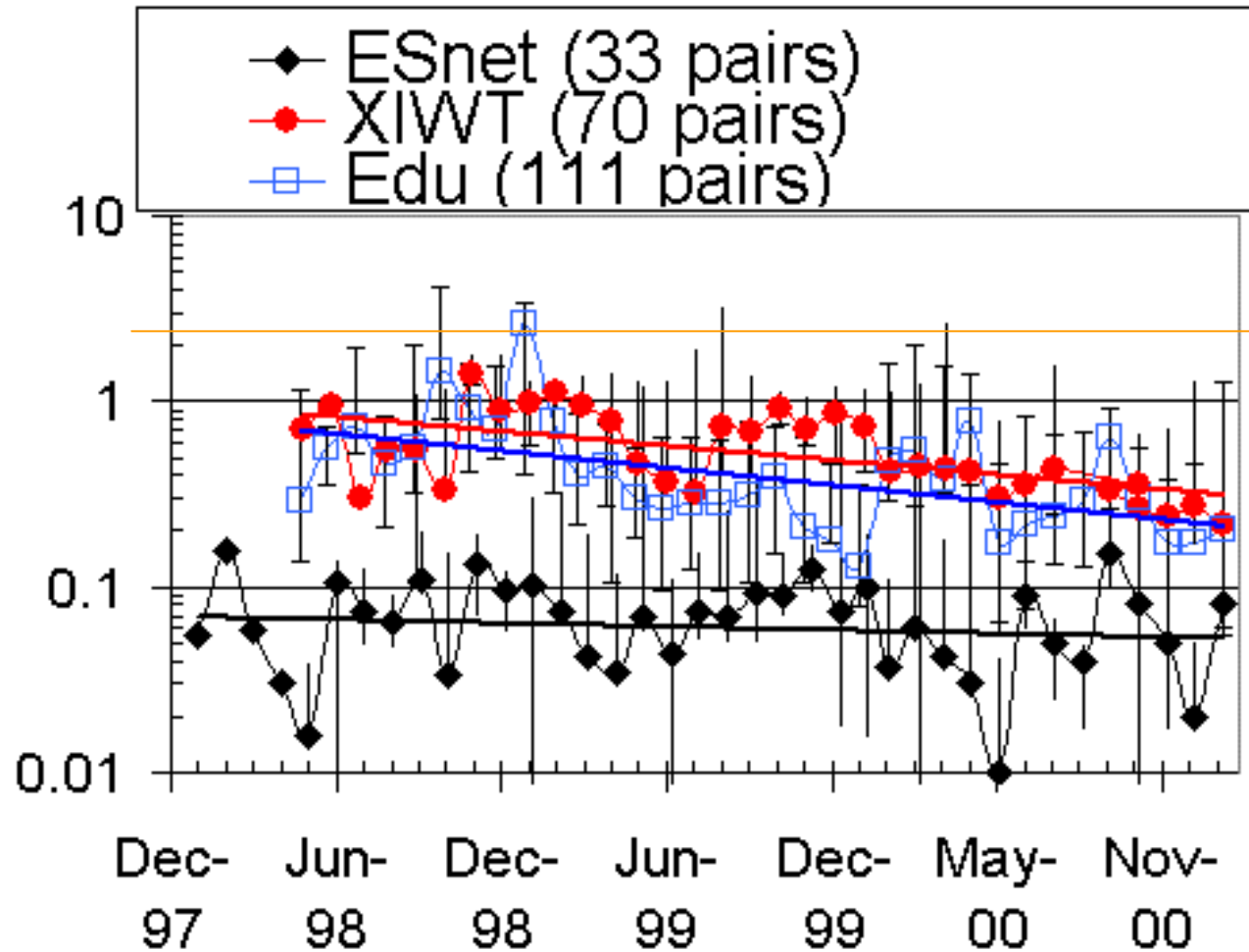
# Loss to world from US

Using year 2000, fraction of world's population/country from [www.nua.ie/surveys/how\\_many\\_online/](http://www.nua.ie/surveys/how_many_online/)

Fraction of world's population with measured loss performance, seen from US



# Losses within US for various Nets



**ESnet vs. Edu vs XIWT loss**



In general performance is good (i.e.  $\leq 1\%$ )

ESnet holding steady

Edu (vBNS/Abilene) & XIWT (70% .com) improving,

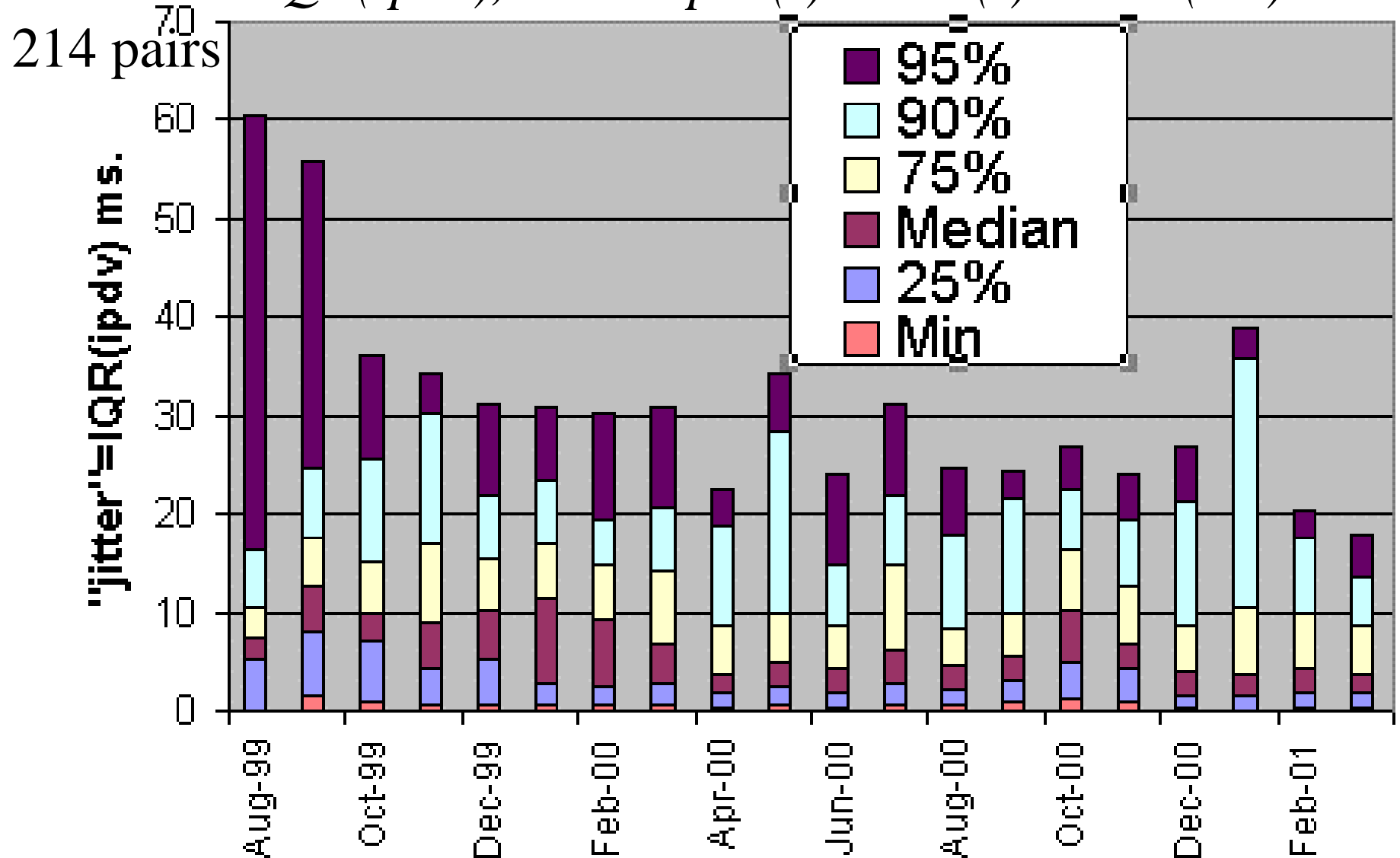


# Losses between Regions

<u>WORLD</u>	<u>Australasia</u>	<u>East Europe</u>	<u>North America</u>	<u>West Europe</u>	<u>South America</u>	<u>Asia</u>
<u>Australasia</u>	0.10	<u>6.91</u>	<u>0.81</u>	<u>1.12</u>	<u>0.31</u>	<u>1.65</u>
<u>East Europe</u>		5.69	<u>2.66</u>	<u>1.34</u>	<u>1.91</u>	<u>2.63</u>
<u>North America</u>		<u>6.39</u>	0.14	<u>0.23</u>	<u>0.02</u>	<u>0.36</u>
<u>West Europe</u>		<u>4.64</u>	<u>0.28</u>	0.30	<u>-0.20</u>	<u>0.43</u>
<u>South America</u>		<u>15.18</u>	<u>4.13</u>	<u>3.11</u>	-1.37	<u>5.02</u>
<u>Asia</u>		<u>5.93</u>	<u>1.68</u>	<u>1.30</u>	<u>0.00</u>	0.32
<u>Africa</u>		<u>32.73</u>	<u>2.10</u>	<u>5.58</u>		
<u>Aria</u>			<u>2.50</u>	<u>5.00</u>		
<u>null</u>			<u>1.12</u>			
<u>Middle East</u>			<u>3.71</u>			
<u>Central America</u>			<u>1.64</u>			

# “Jitter” from N. America to W. Europe

“Jitter” =  $IQR(ipdv)$ , where  $ipdv(i) = RTT(i) - RTT(i-1)$



ETSI: DTR/TIPHON-05001 V1.2.5 (1998-09) good speech < 75ms jitter<sup>10</sup>

# “Jitter” between regions

Median	US	W. Europe	E. Europe
US	2	4	15
W. Europe	4	5	22
E. Europe	18	11	29
Asia	18	14	46
S. America	20	17	38
Australasia	4	4	28
Africa	28		

90th%	US	W. Europe	E. Europe
US	9	14	25
W. Europe	13	16	39
E. Europe	134	38	130
Asia	54	292	56
S. America	81	32	53
Australasia	50	15	
Africa	171	178	265

Median	ESnet	Edu	Com
ESnet	2	6	
Edu	2	1	
Com			5

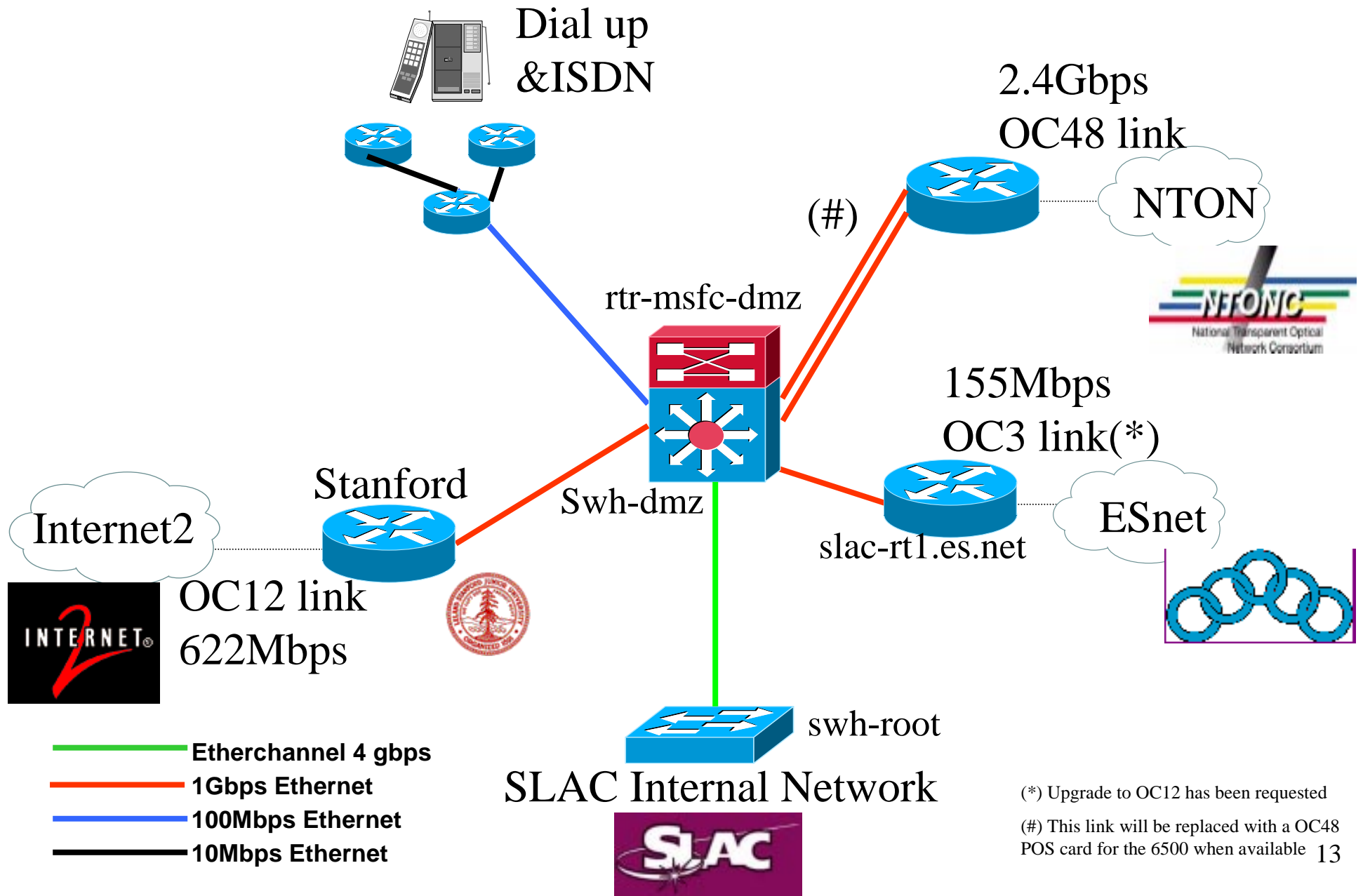
ETSI: DTR/TIPHON-05001 V1.2.5 (1998-09)

75ms=Good    125ms=Med    225ms=Poor

# Passive site border monitoring

- Use SNMP to get utilization etc.
- Used to use OC3Mon with CoralReef for flows etc. but now have GigE interfaces
- Use **Cisco Netflow** in Catalyst 6509 with MSFC, only on border at the moment
- Gather about 200MBytes/day of flow data
- Data recorded in binary every 10 minutes into RRD
- The raw data records include source and destination addresses and ports, the protocol, packet, octet and flow counts, and start and end times of the flows
  - Much less detailed than OC3Mon, but good compromise
  - Top talkers history and daily (from & to), tlds, vlans, protocol and application utilization, flow times, time series, distributions
- Use for network & security

# Simplified SLAC DMZ Network, 2001



# SLAC Traffic profile

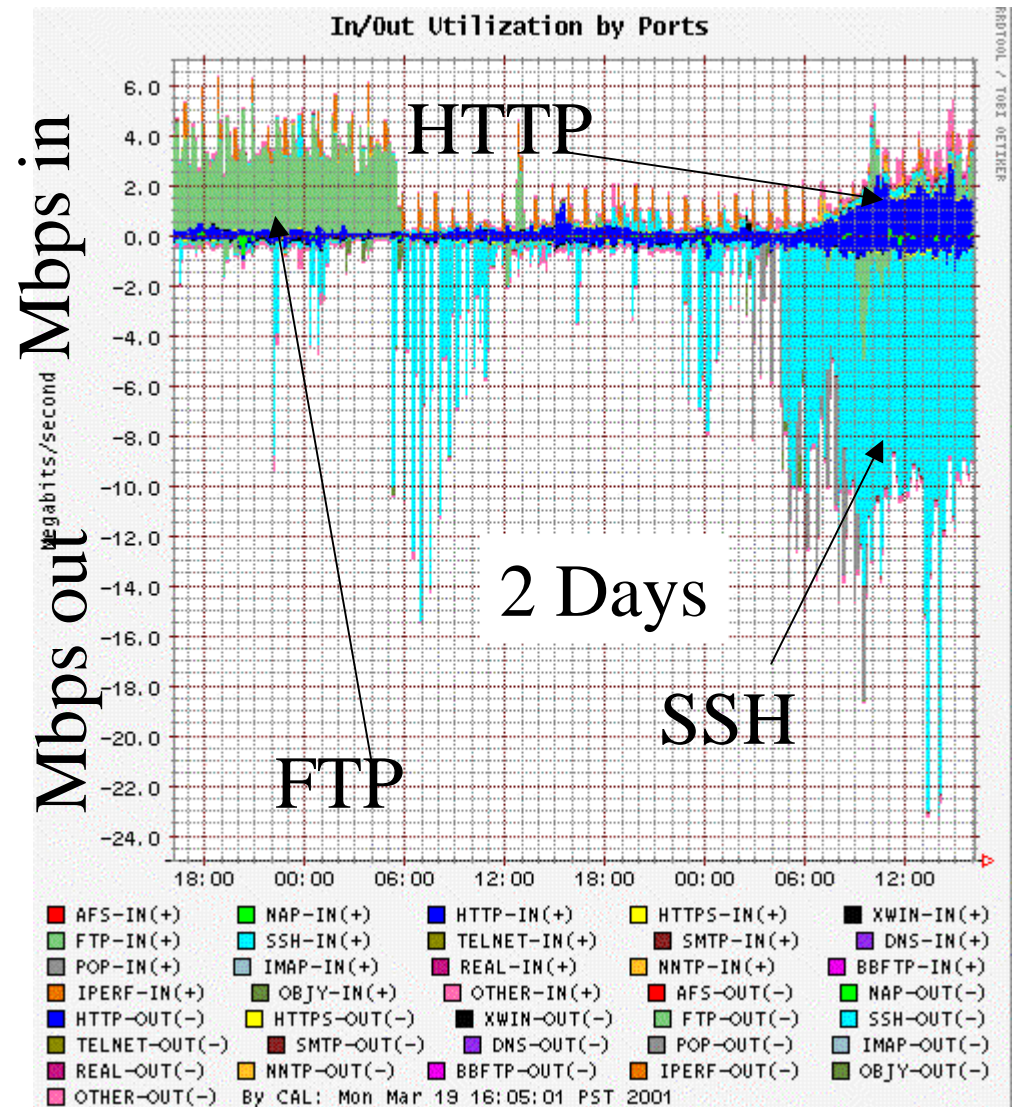
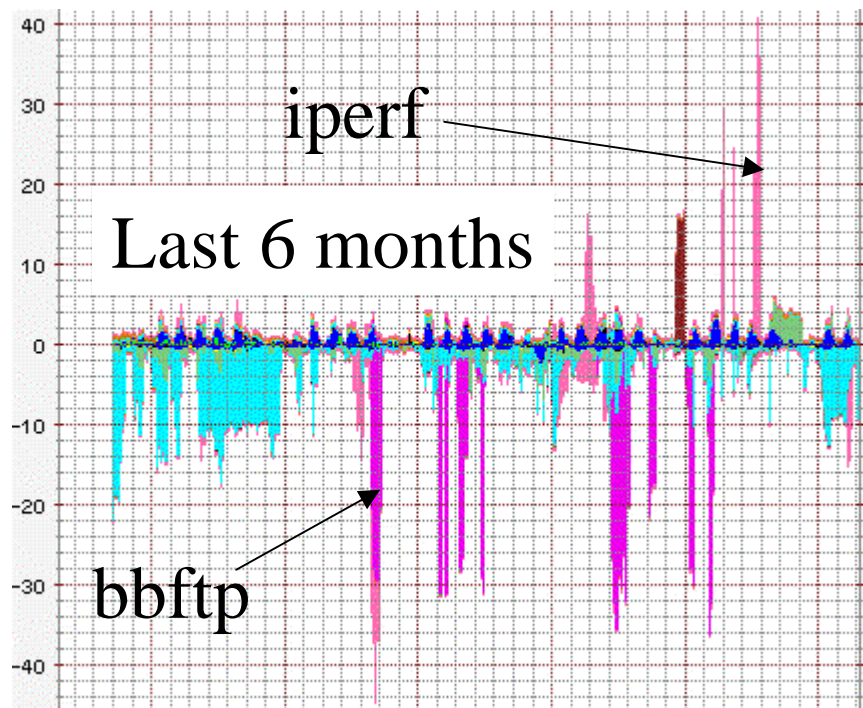
SLAC offsite links:

OC3 to ESnet, 1Gbps to Stanford U & thence OC12 to I2

OC48 to NTON

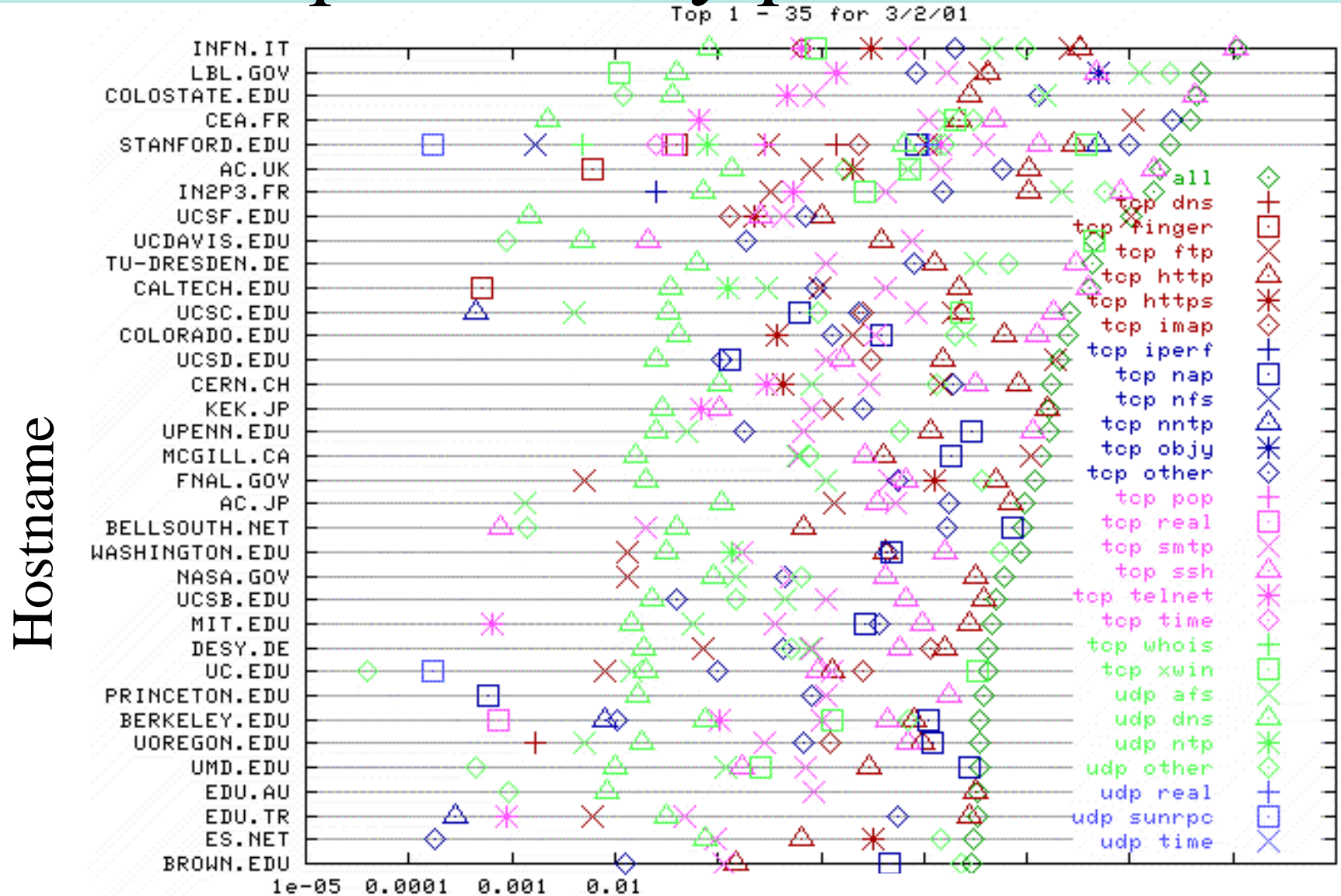
Profile

bulk-data xfer dominates





# Top talkers by protocol



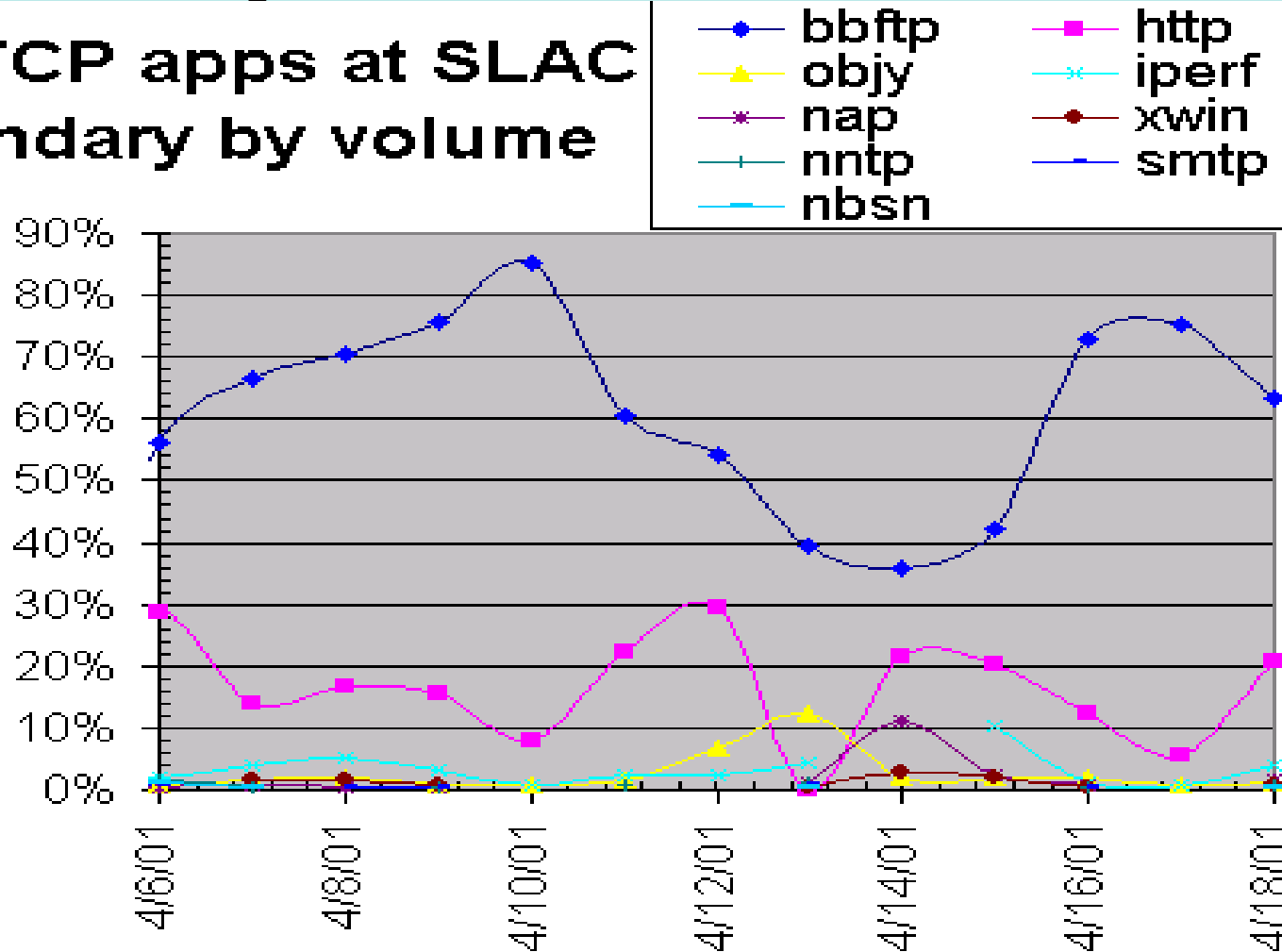
Volume dominated by single  
Application - bbftp

1 100 10000  
MBytes/day (log scale)

# Not your normal Internet site

Top TCP apps at SLAC  
boundary by volume

Daily percentage of  
byte volume



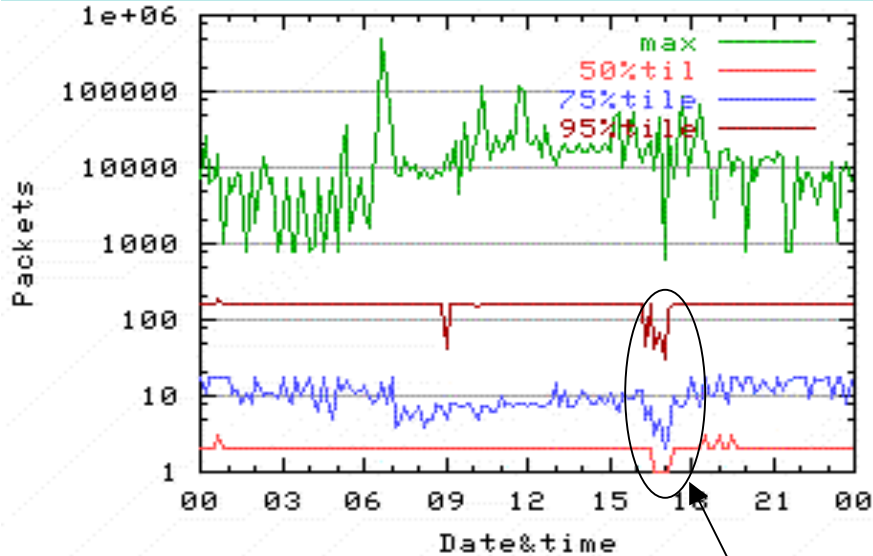
50-300Gbytes/day

Ames IXP: approximately 60-65% was HTTP, about 13% was NNTP

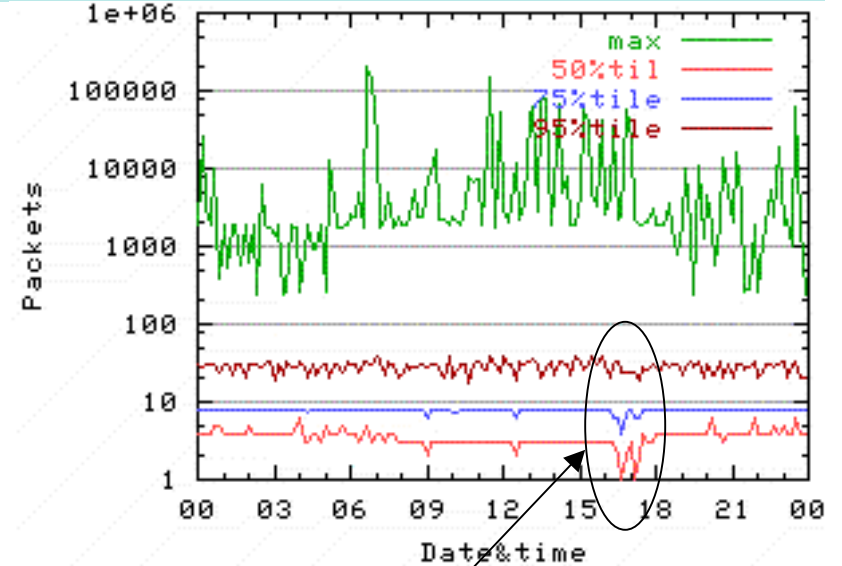
Uwisc: 34% HTTP, 24% FTP, 13% Napster



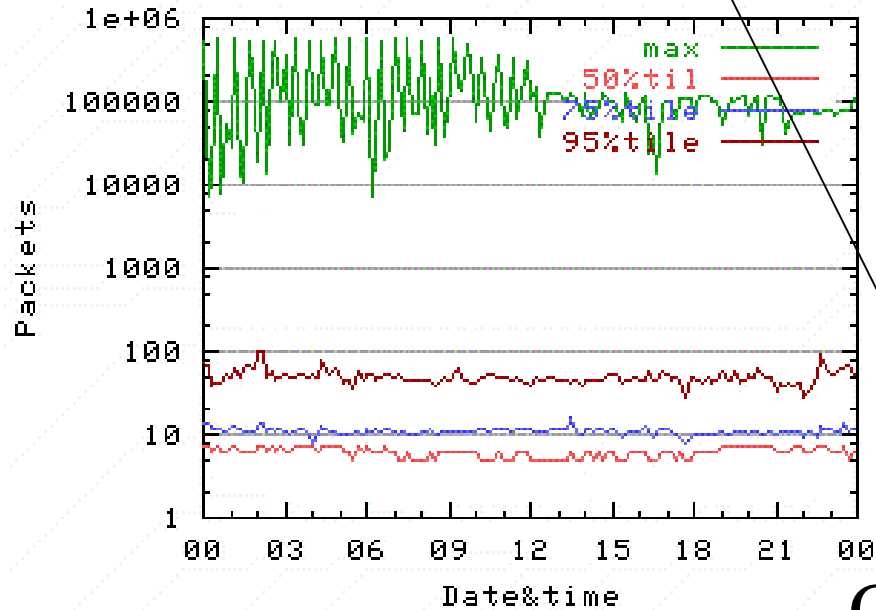
# Time series



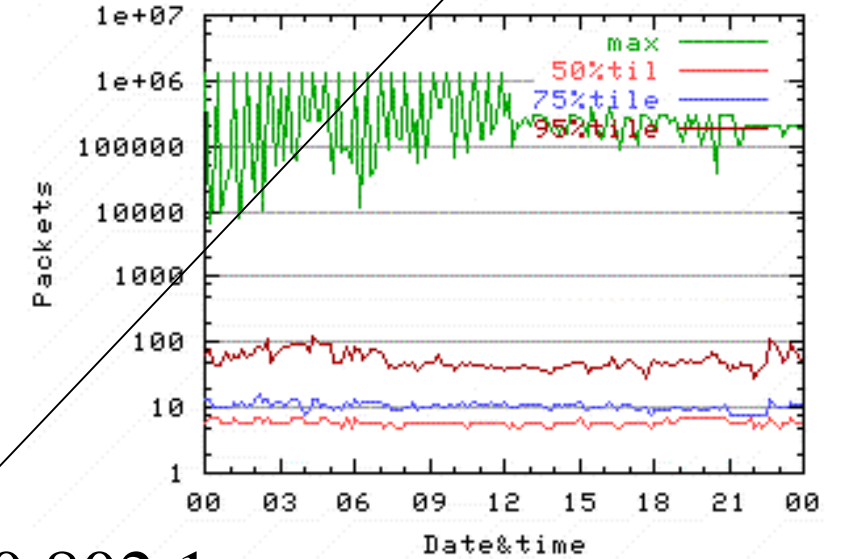
UDP



TCP



Outgoing

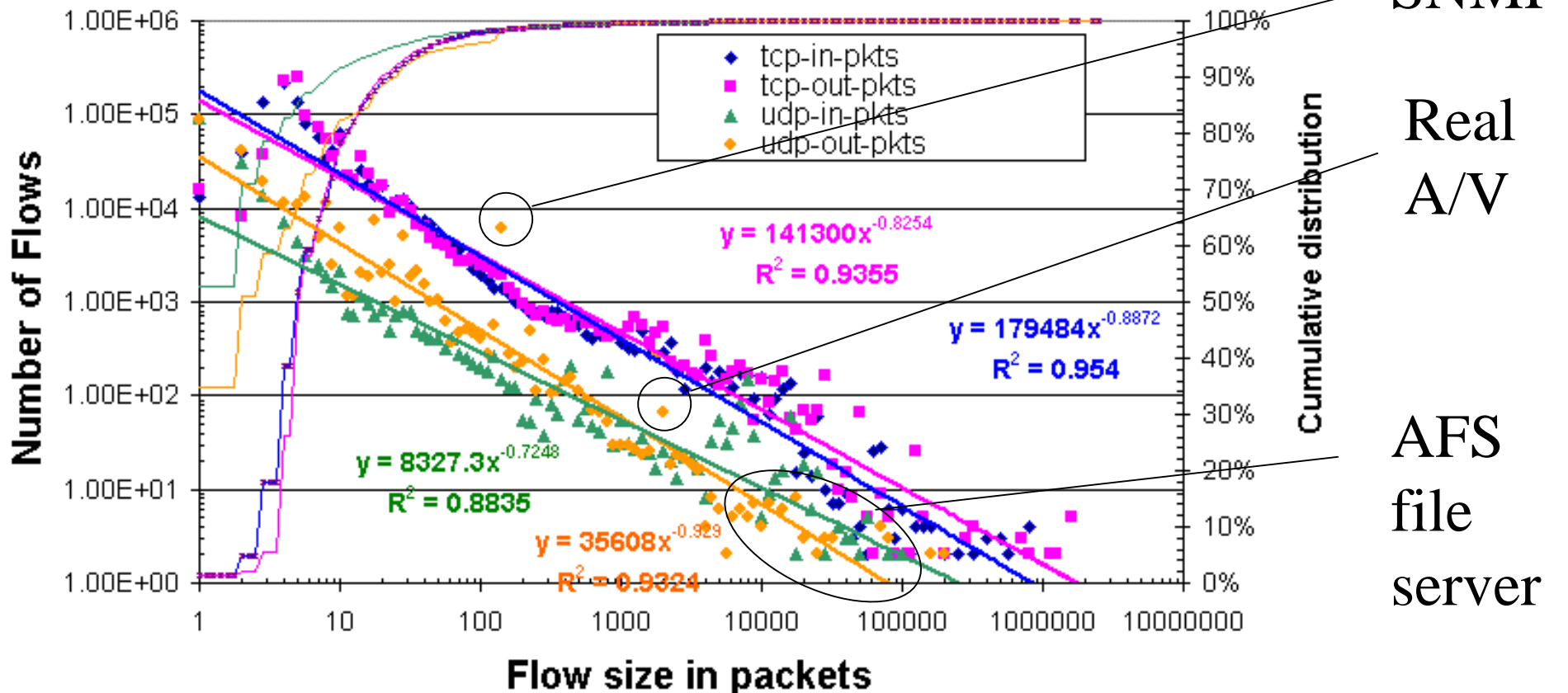


Incoming

Cat 4000 802.1q  
vs. ISL

# Flow sizes

Flow size distribution at SLAC border April 9, 2001



Confirms Nevil Brownlee's data measured at SDSC:

Heavy tailed, in ~ out, UDP flows shorter than TCP, packet~bytes

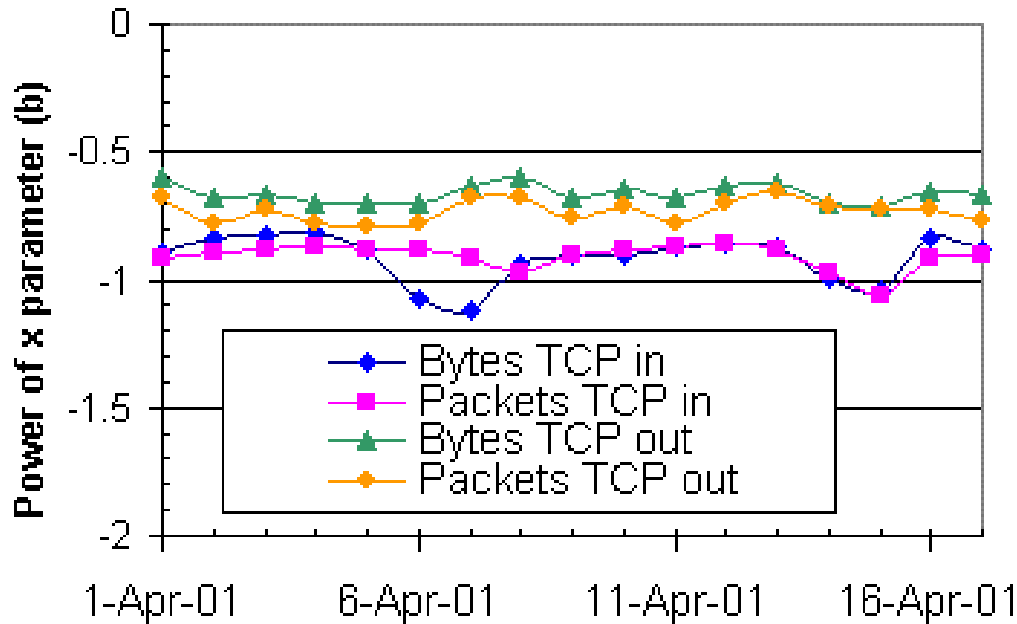
75% TCP-in < 5kBytes, 75% TCP-out < 1.5kBytes (<10pkts)

UDP 80% < 600Bytes (75% < 3 pkts), ~10 \* more TCP than UDP

Top UDP = AFS (>55%), Real(~25%), SNMP(~1.4%)

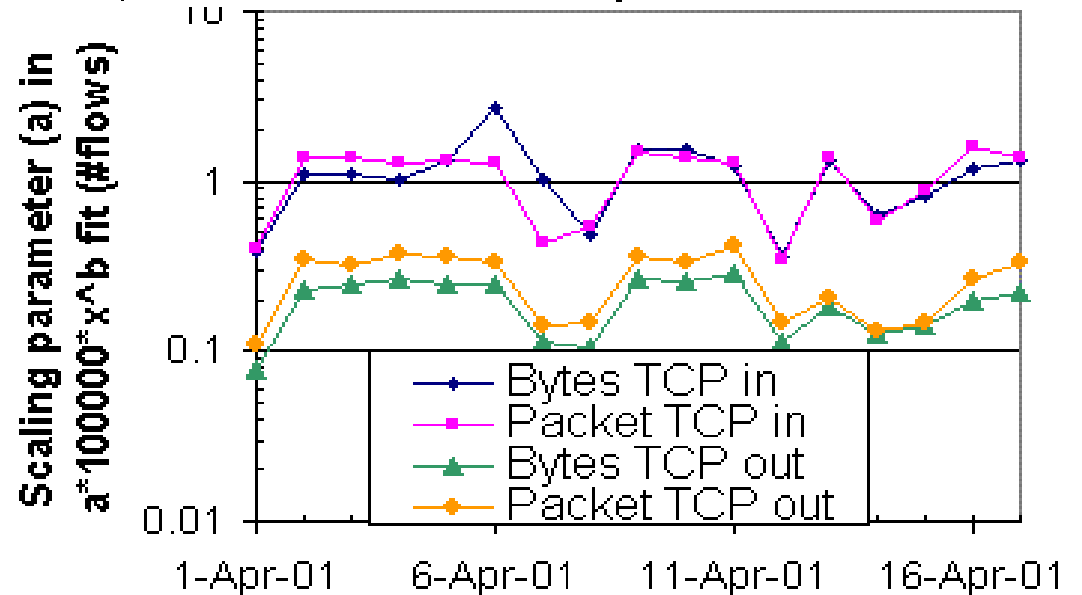
# Power law fit parameters by time

Slope of power law fit to Flow frequencies



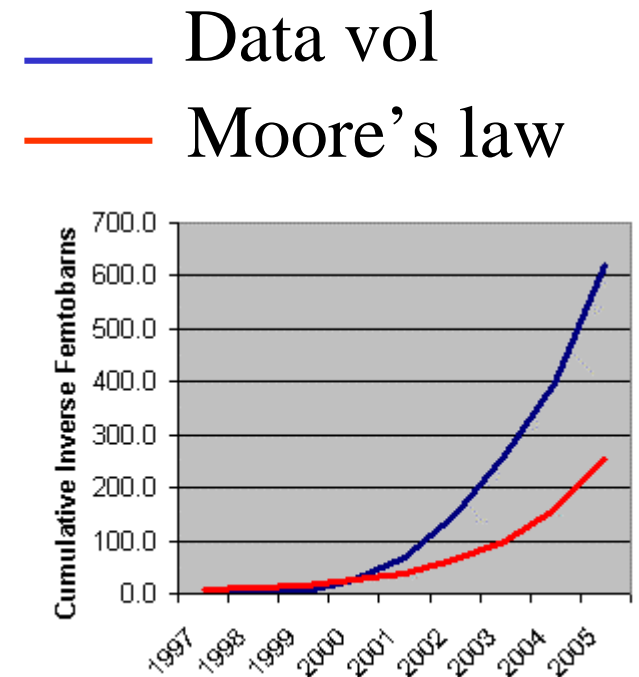
Just 2 parameters provide a reasonable description of the flow size distributions

Scaling parameter for power law fit to Flow frequencies



# App: High Speed Bulk Throughput

- **Driven by:**
  - Data intensive science, e.g. data grids
  - HENP data rates, e.g. BaBar 300TB/year, collection doubling yearly, i.e. PBytes in couple of years
  - Data rate from experiment  $\sim 20\text{MBytes/s} \sim 200\text{GBytes/d}$
  - Multiple regional computer centers (e.g. Lyon-FR, RAL-UK, INFN-IT, LBNL-CA, LLNL-CA, Caltech-CA) need copies of data
  - Boeing 747 high throughput, BUT poor latency ( $\sim 2$  weeks) & very people intensive
- So need high-speed networks and ability to utilize
  - High speed today = few hundred GBytes/day



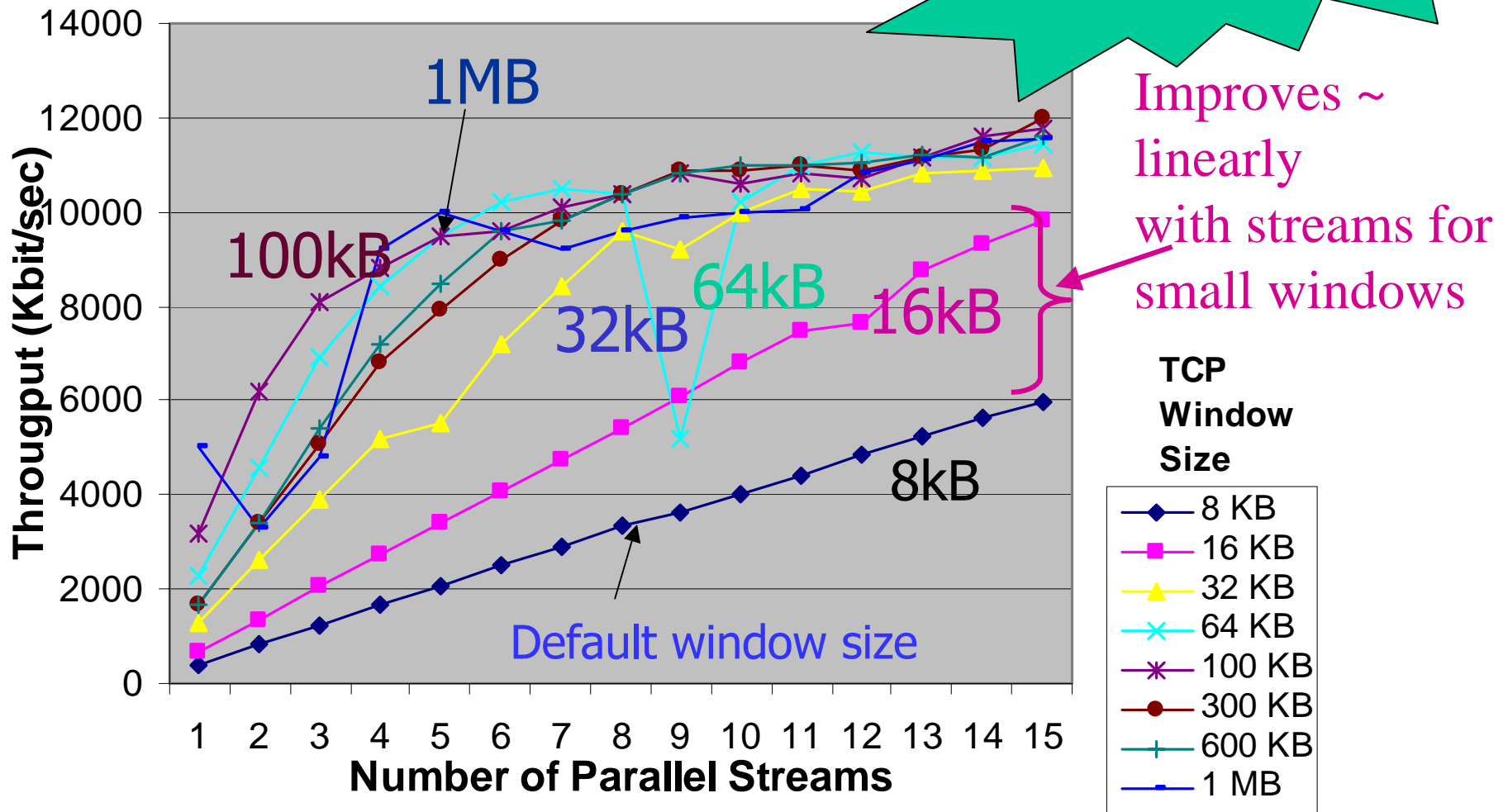
# Measuring TCP throughput

- Selected about a dozen major collaborator sites in CA, CO, IL, FR, UK over last 9 months
  - Of interest to SLAC
  - Can get logon accounts
- Use **iperf**
  - Choose window size and # parallel streams
  - Run for 10 seconds together with ping (loaded)
  - Stop iperf, run ping (unloaded) for 10 seconds
  - Change window or number of streams & repeat
- Record streams, window, throughput (Mbits/s), loaded & unloaded ping responses

# SLAC to CERN thruput vs windows & streams

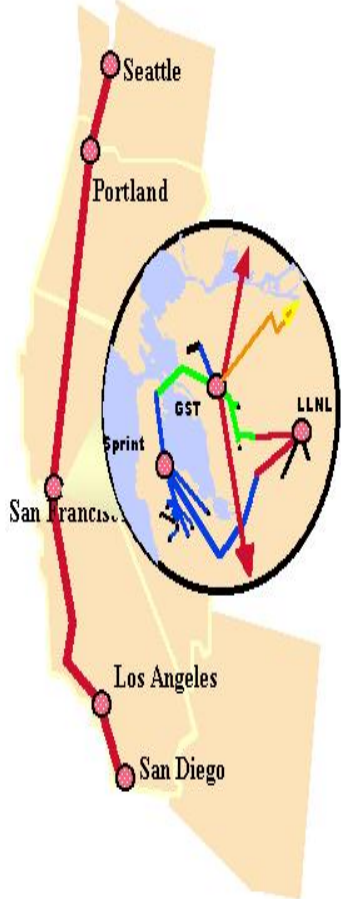
iperf file transfer (2MB) between SLAC  
and CERN  
25 Feb 2000

Hi-perf = big windows  
& multiple streams

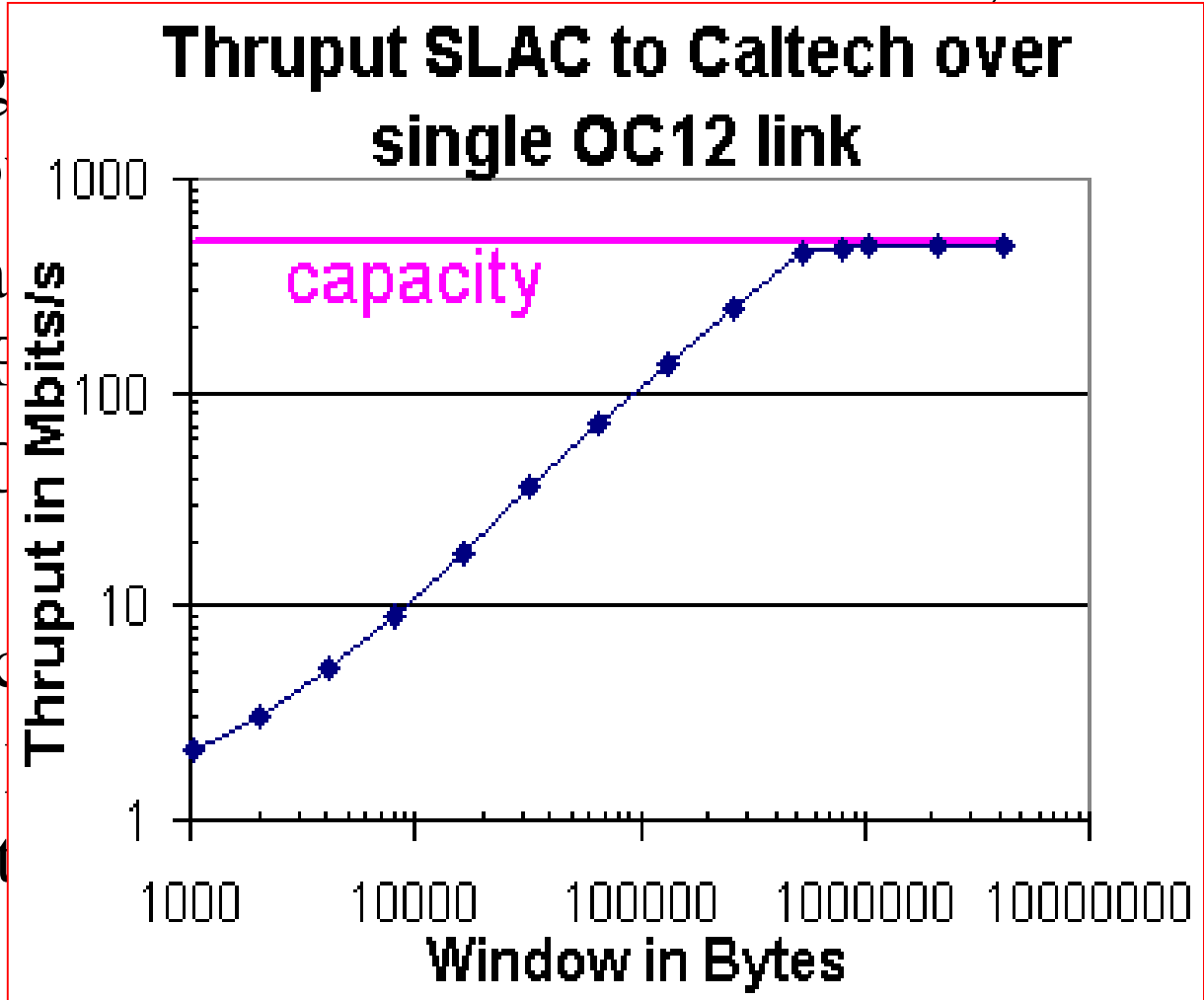




# Progress towards goal: 100 Mbytes/s Site-to-Site



- Focus on SLAC – Caltech over NTON;
- Using ... & do ...
- Repla ... with ...
- SLAC ... and 2 ...
- Caltec ...
- ~500 ... recent

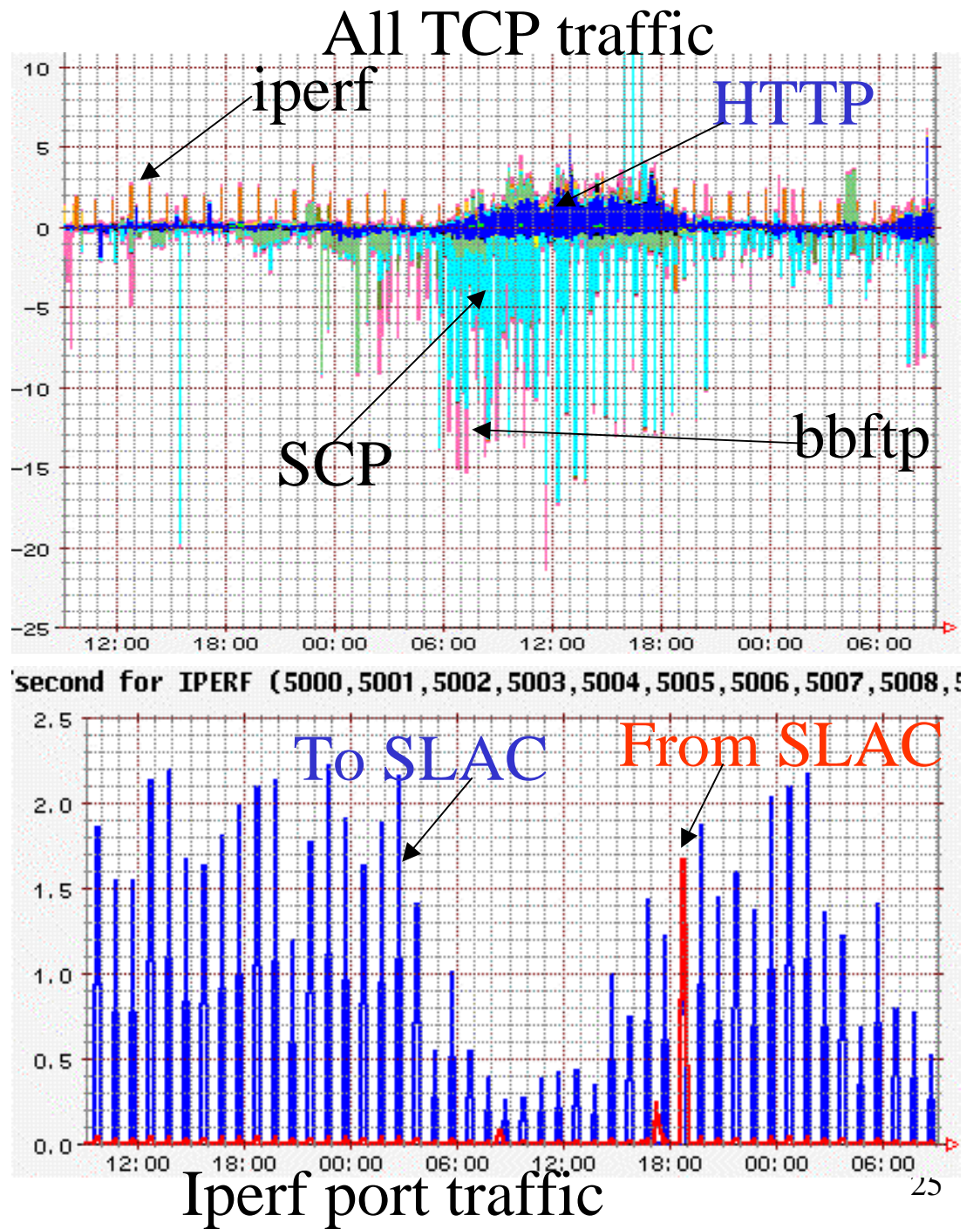


# SC2000 WAN Challenge

- SC2000, Dallas to SLAC RTT ~ 48msec
  - SLAC/FNAL booth: Dell PowerEdge PIII 2 \* 550MHz with 64bit PCI + Dell 850MHz both running Linux, each with GigE, connected to Cat 6009 with 2GigE bonded to Extreme SC2000 floor switch
  - NTON: OC48 to GSR to Cat 5500 Gig E to Sun E4500 4\*460MHz and Sun E4500 6\*336MHz
- Internet 2: 300 Mbits/s
- NTON 960Mbits/s Dallas to SLAC mem-to-mem
- Details:
  - [www-iepm.slac.stanford.edu/monitoring/bulk/sc2k.html](http://www-iepm.slac.stanford.edu/monitoring/bulk/sc2k.html)



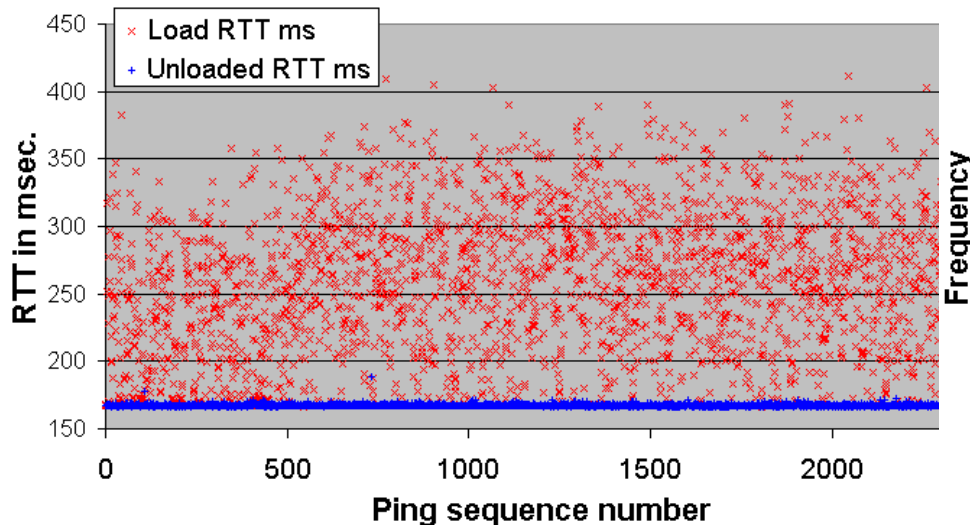
# Impact of cross-traffic on Iperf between SLAC & W. Europe



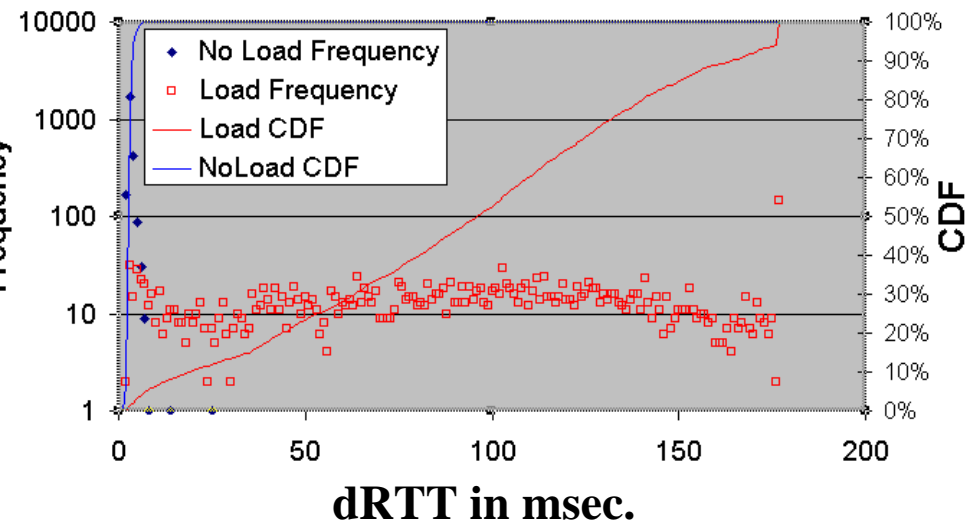
# Impact on Others

- Make ping measurements with & without iperf loading
  - Loss loaded(unloaded)
  - RTT

Ping RTT between SLAC & CERN for a loaded and unloaded link



Ping RTTs for loaded and unloaded link between SLAc & CERN



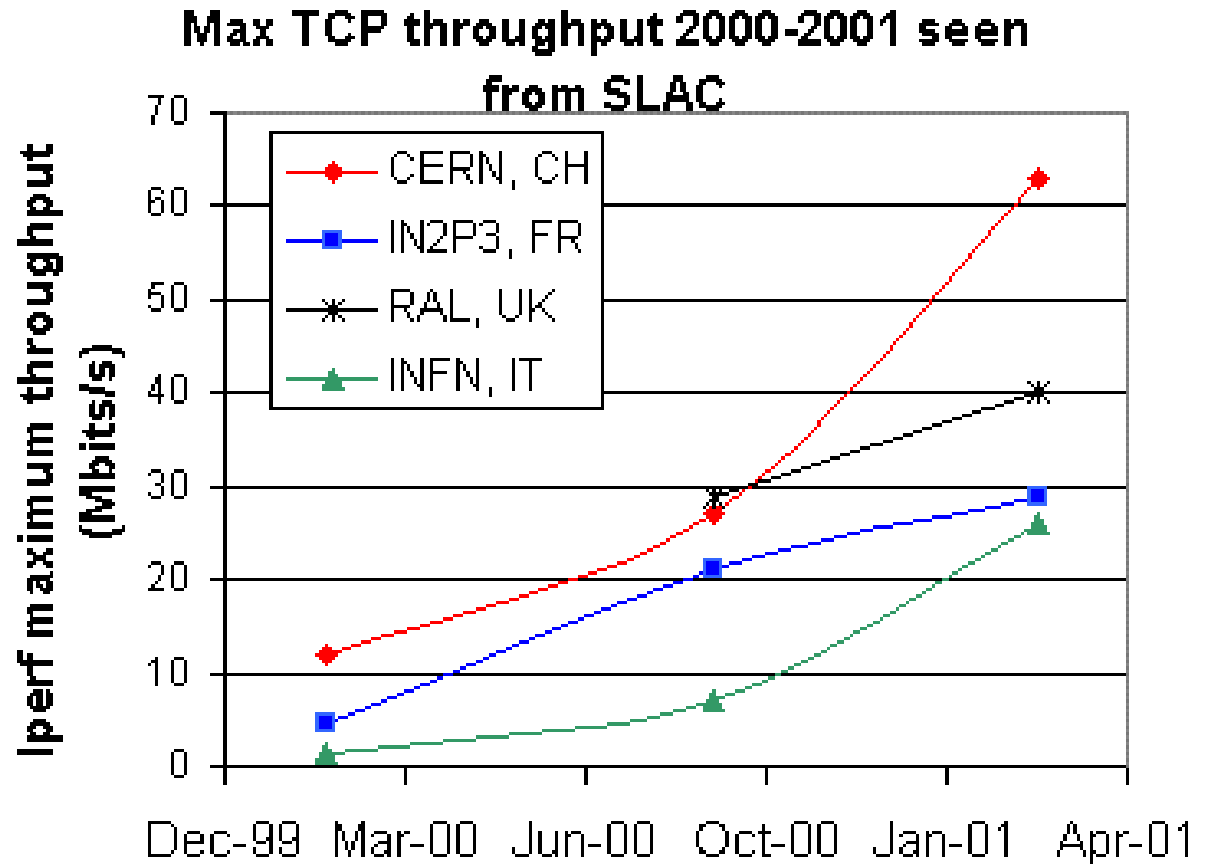
# Improvements for major International BaBar sites

Links are being improved: ESnet, PHYnet, GARR, Janet, TEN-155

Improvements to come:

IN2P3 => 155Mbps

RAL => 622Mbps



Throughput improvements of 2 to 16 times in a year

# Iperf throughput conclusions 1/2

- Can saturate bottleneck links
- For a given **iperf** measurement, streams share throughput equally.
- For small window sizes throughput increases linearly with number of streams
- Predicted optimum window sizes can be large ( $>$  Mbyte)
- Need  $>$  1 stream to get optimum performance
- Can get close to max thruptut with small ( $\leq 32$ Mbyte) with sufficient (5-10) streams
- Improvements of 5 to 60 in thruptut by using multiple streams & larger windows
- Loss not sensitive to throughput

## Iperf thruput conclusions 2/2

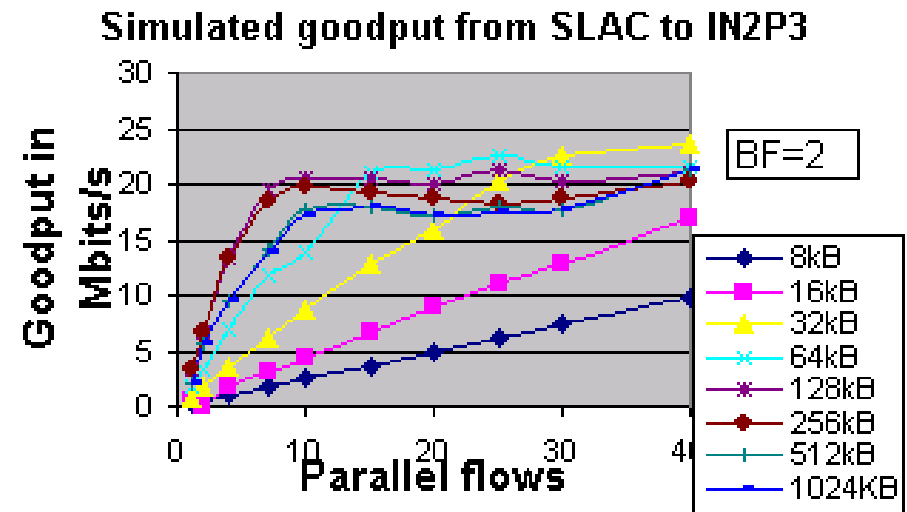
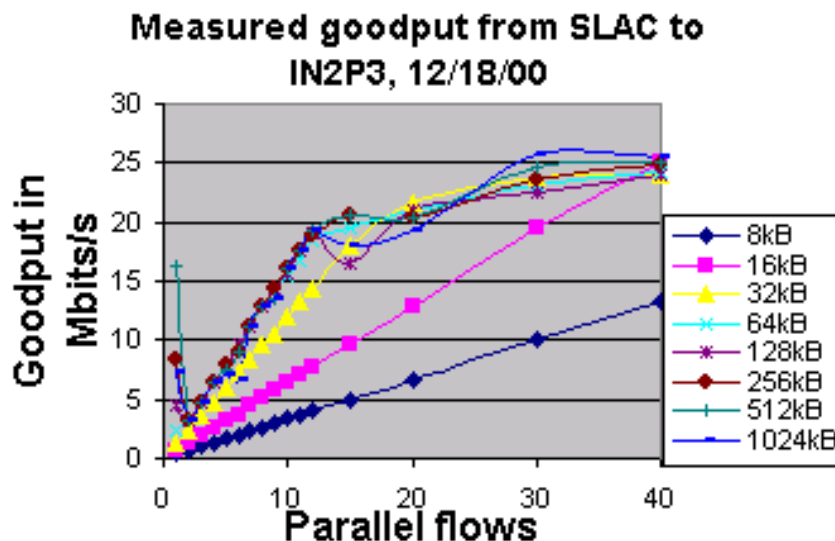
- For fixed *streams\*window* product, streams are more effective than window size:

Site	Window	Streams	Throughput
CERN	256kB	2	9.45Mbits/
CERN	64kB	8	26.8Mbits/s
Caltech	256kB	2	1.7Mbits/s
Caltech	64kB	8	4.6Mbits/s

- There is an optimum number of streams above which performance flattens out
- See [www-iepm.slac.stanford.edu/monitoring/bulk/](http://www-iepm.slac.stanford.edu/monitoring/bulk/)

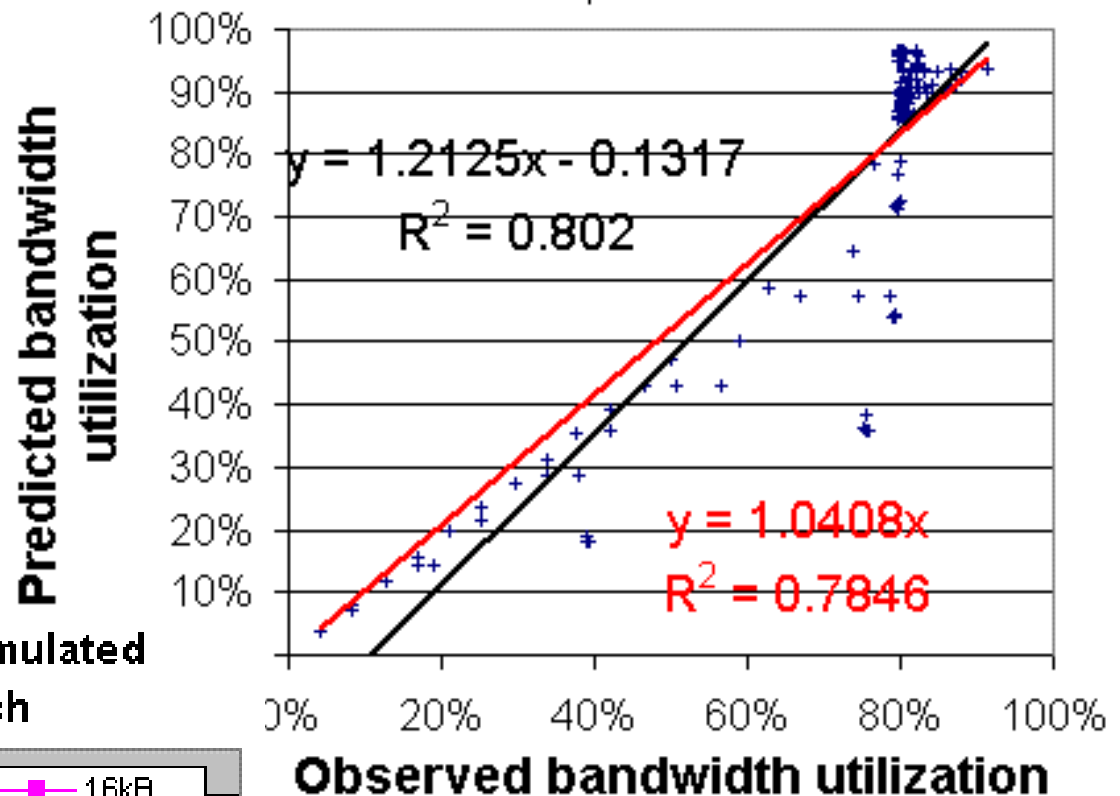
# Network Simulator (ns-2)

- From UCB, simulates network
  - Choice of stack (Reno, Tahoe, Vegas, SACK...)
  - RTT, bandwidth, flows, windows, queue lengths ...
- Compare with measured results
  - Agrees well
  - Confirms observations (e.g. linear growth in throughput for small window sizes as increase number of flows)

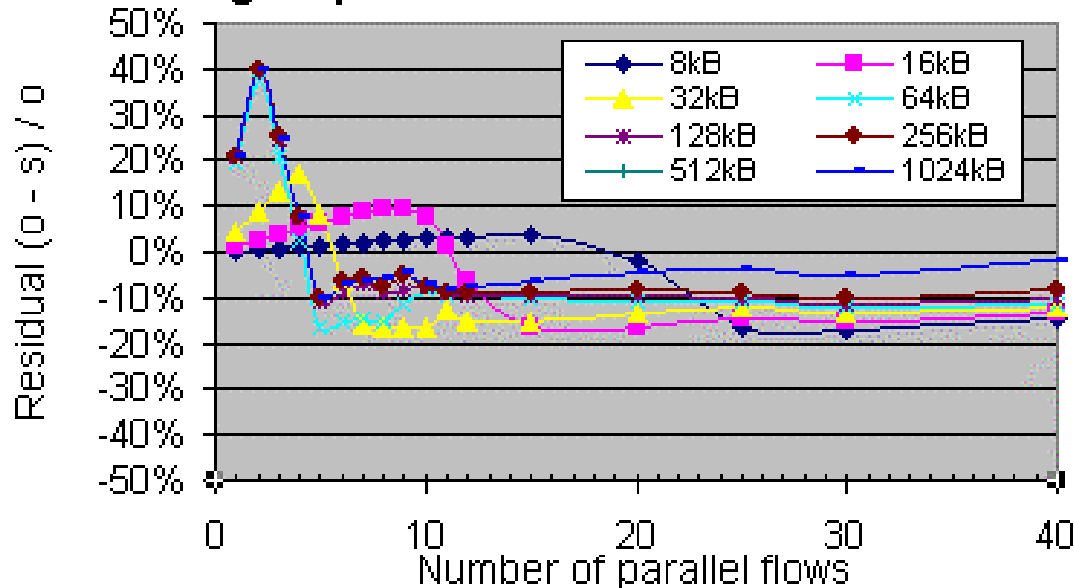


# Agreement of ns2 with observed

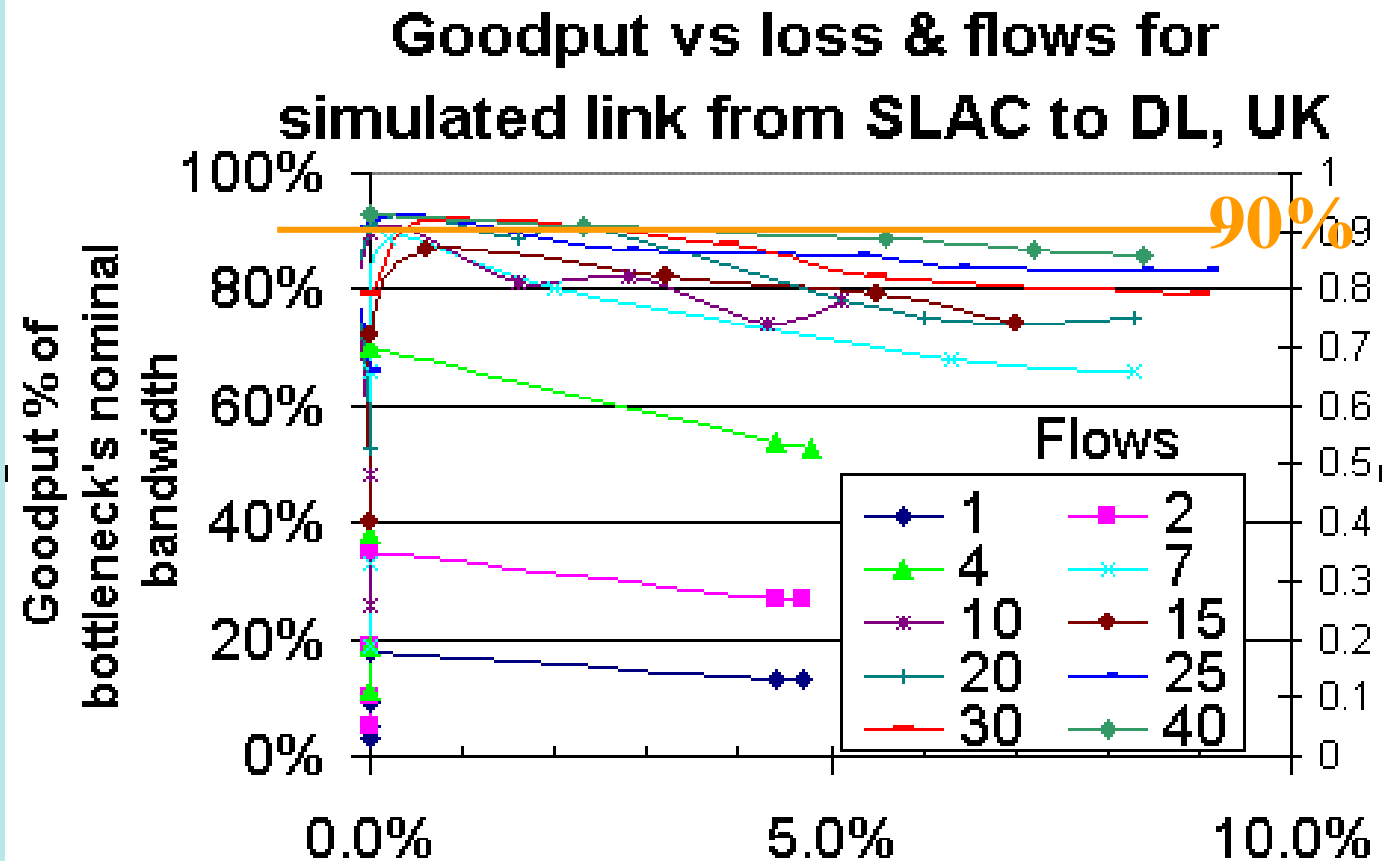
**Predicted vs observed utilization for SLAC to Caltech, Dec 25, 2001**



**Residuals between observed & simulated goodput from SLAC to Caltech**



# Ns-2 thruput & loss predict



BW=10Mbps, RTT=162ms, Q=80, BF=2 **Packet loss**

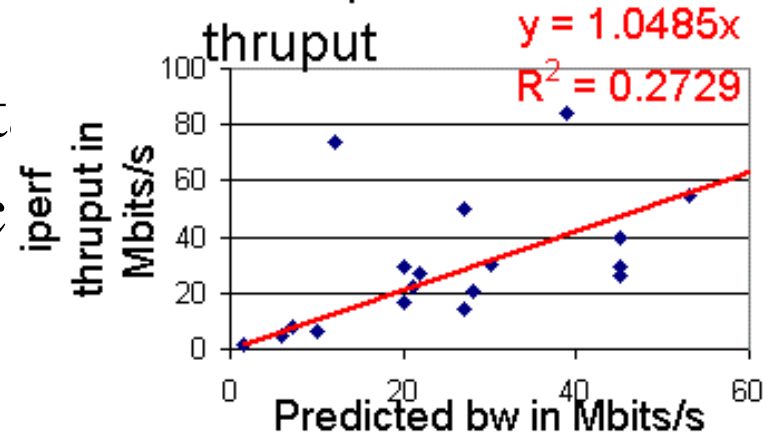
- Indicates on unloaded link can get 70% of available bandwidth without causing noticeable packet loss
- Can get over 80-90% of available bandwidth
- Can overdrive: no extra throughput BUT extra loss



# Simulator benefits

- No traffic on network (nb throughput can use 90%)
- Can do what if experiments
- No need to install iperf servers or have accounts
- No need to configure host to allow large windows
- BUT
  - Need to estimate simulator parameters, e.g.
    - RTT use ping or synack
    - Bandwidth, use pchar, pipechar etc., moderately accurate
- **AND *its not the real thing***
  - Need to validate vs. observed data
  - Need to simulate cross-traffic etc

Measured vs. predicted



# WAN thruput conclusions

- High FTP performance across WAN links is possible
  - Even with 20-30Mbps bottleneck can do > 100Gbytes/day
- OS must support **big windows** selectable by application
- Need **multiple parallel streams**
- **Loss is important** in particular interval between losses
- Compression looks promising, but needs cpu power
- Can get close to max thruput with small (<=32Mbyte) with sufficient (5-10) streams
- **Improvements of 5 to 60 in thruput** by using multiple streams & larger windows
- Impacts others users, need **Less than Best Effort** QoS service

# More Information

- This talk:
  - [www.slac.stanford.edu/grp/scs/net/talk/slac-wan-perf-apr01.htm](http://www.slac.stanford.edu/grp/scs/net/talk/slac-wan-perf-apr01.htm)
- IEPM/PingER home site
  - [www-iepm.slac.stanford.edu/](http://www-iepm.slac.stanford.edu/)
- Transfer tools:
  - <http://hepwww.rl.ac.uk/Adye/talks/010402-ftp/html/sld015.htm>
- TCP Tuning:
  - [www.ncne.nlanr.net/training/presentations/tcp-tutorial.ppt](http://www.ncne.nlanr.net/training/presentations/tcp-tutorial.ppt)