

What we have learned from developing and running ABwE

Jiri Navratil, Les R.Cottrell
(SLAC)

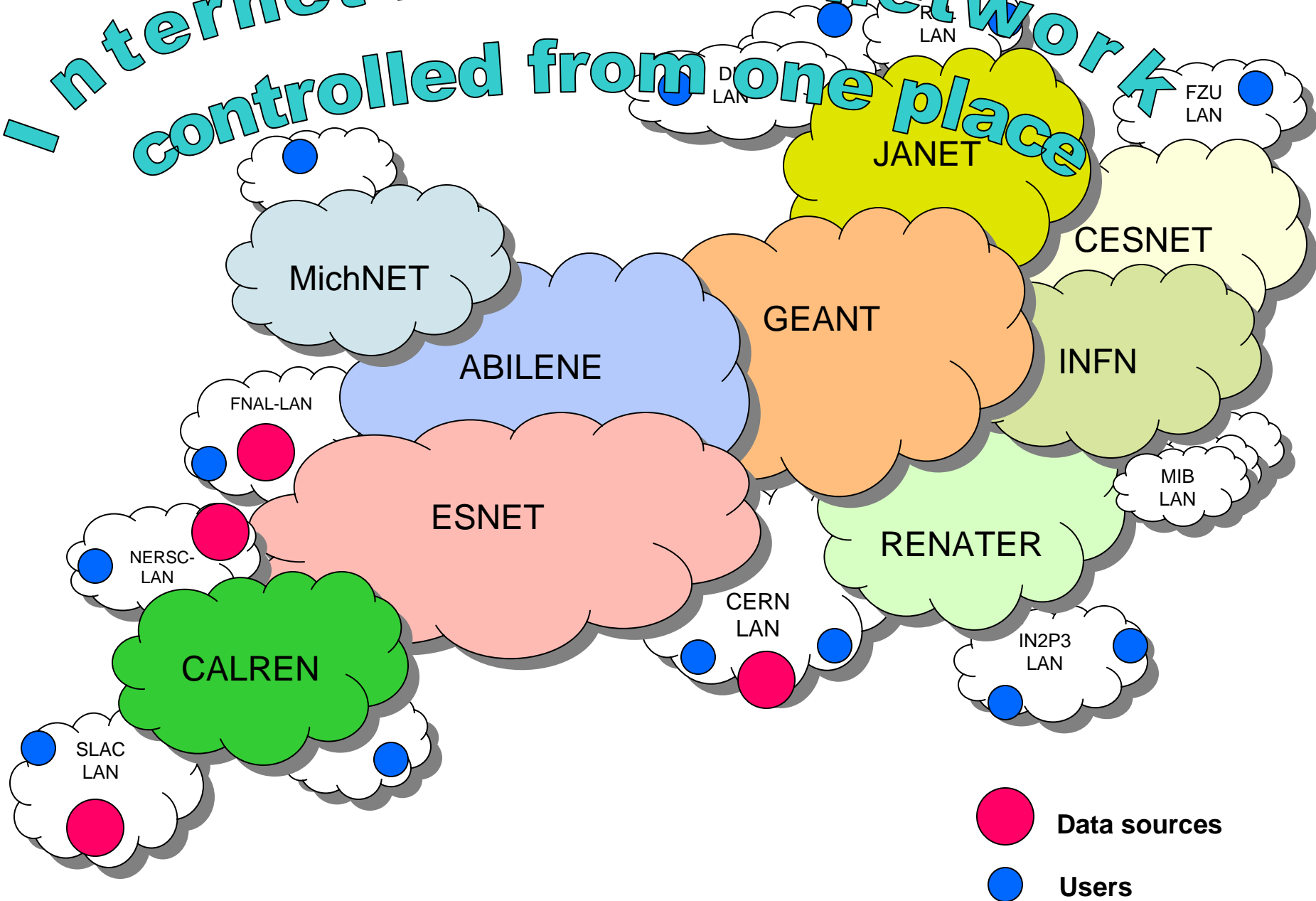
Why E2E tools are needed

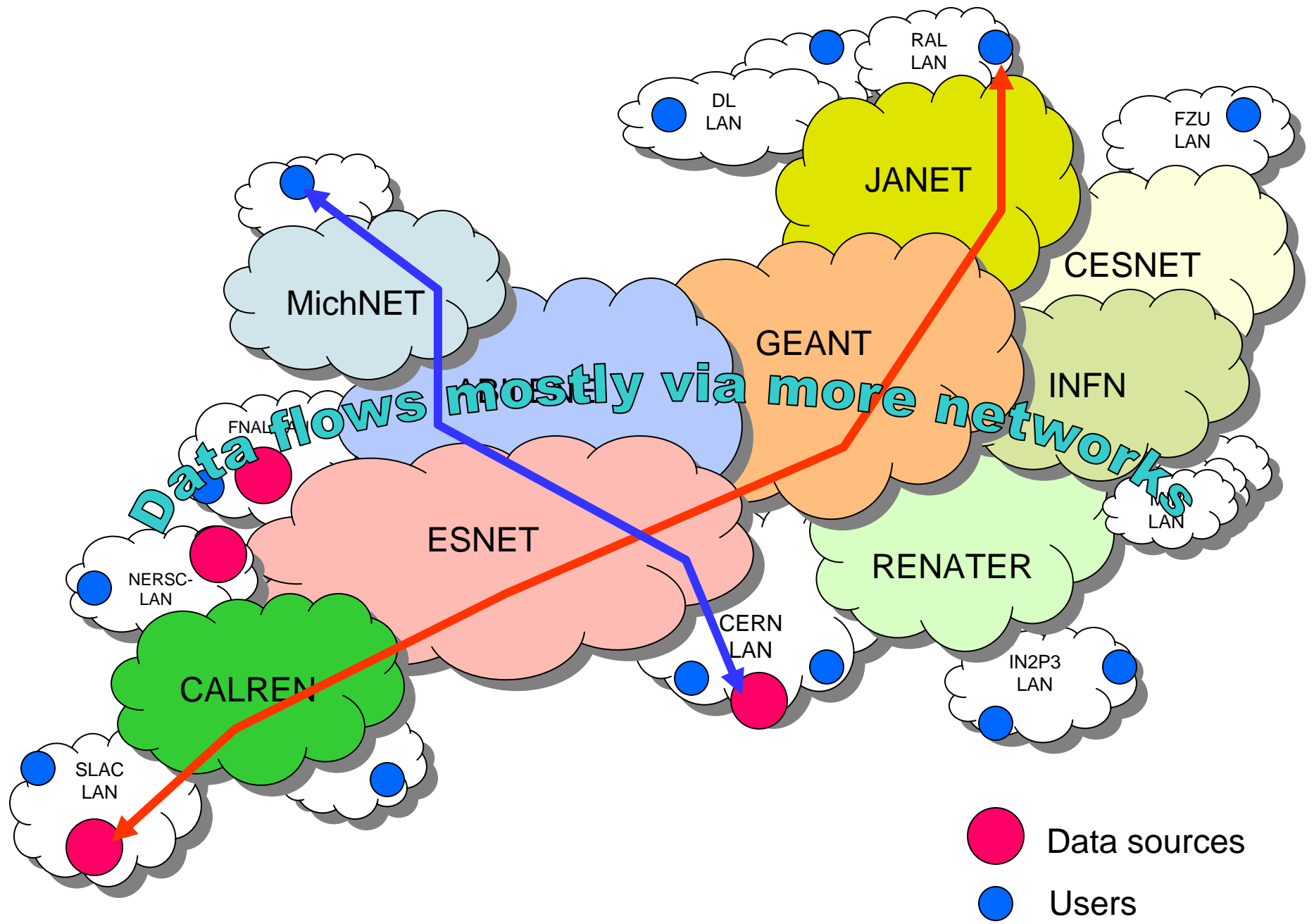
- The **scientific community** is **increasingly dependent** on networking as **international cooperation grows**. **HEP users** (needs transfer huge amount of data between **experimental sites** as **SLAC, FNAL, CERN**, etc. (where data is created) and home institutes spread over the world)

- **What ISPs** (as Abilene, Esnet, Geant..) **can offer to the users for getting information?**

(Not too much because they are only in the middle of the path and they don't cover all parts of connections)

Internet is not one network
controlled from one place





- There must be always **somebody who gives complex information** to the **users of the community**

or

the **users** have to have **a tool which give them such information**

- **How fast I can transfer 20 GB from my experimental site (SLAC,CERN) to my home institute?**
- **Can I run graphical 3D visualization program with data located 1000 miles away?**
- **How stable is line ? (Can I use it in the same conditions for 5 minutes or 2 hours or whole day ?)**

All such questions must be **replied in few seconds** doesn't matter if for **individual user** or for **Grid brokers**

- **Global science has no day and night.**

To reply this we needed the tools that could be used in continuous mode 24 hours a day 7 days a week which can non intrusively detect changes on multiple path or on demand by any user

ABwE: Basic terminology:

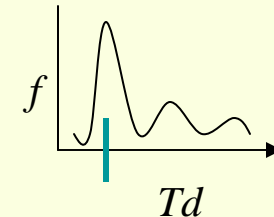
- *Generally:*

$$\text{Available bandwidth} = \text{Capacity} - \text{Load}$$

- **ABwE measure T_d** – Time dispersion P1-P2 (20x PP)

We are trying to distinguish two basic states in our results:

- “**Dominate (free)**” – when $T_d \approx \text{const}$
- “**loaded**” with $T_d = \text{other value}$



T_d results from “**Dominate**” state are used to estimate

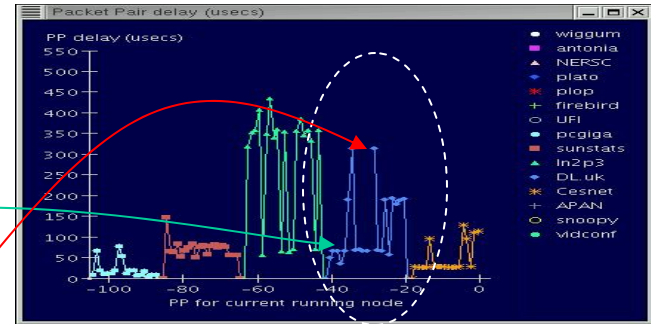
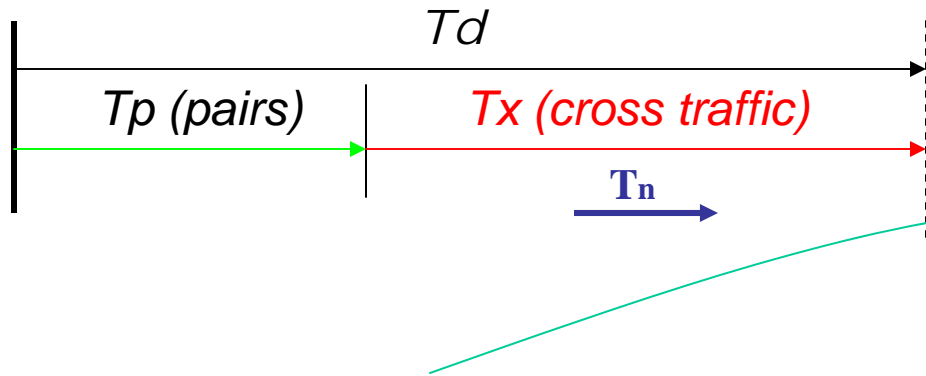
DBC - Dynamic Bottleneck Capacity

T_d measured during the “**loaded**” state is used to estimate the level of

XTR (cross traffic)

$$\text{ABw} = \text{DBC} - \text{XTR}$$

Abing: Estimation principles:



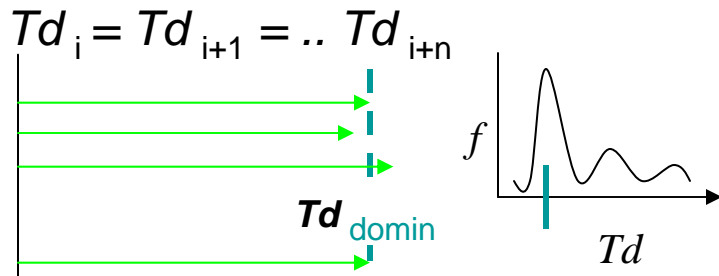
Examples T_d from different paths

"Dominating state"

(when sustained load or no load)

"Load state"

(when load is changing)



$$D_{bc} = L_{pp} / T_{d \text{ domin}}$$

$$q = Tx / T_n \quad (Tx = T_d - T_p)$$

T_x – busy time (transmit time for cross traffic)

T_n – transmit time for average packet

q – **relative queue increment (QDI)**

during decision interval $T_d (h-1)$

$$u = q / (q + 1)$$

$$CT = u * D_{bc}$$

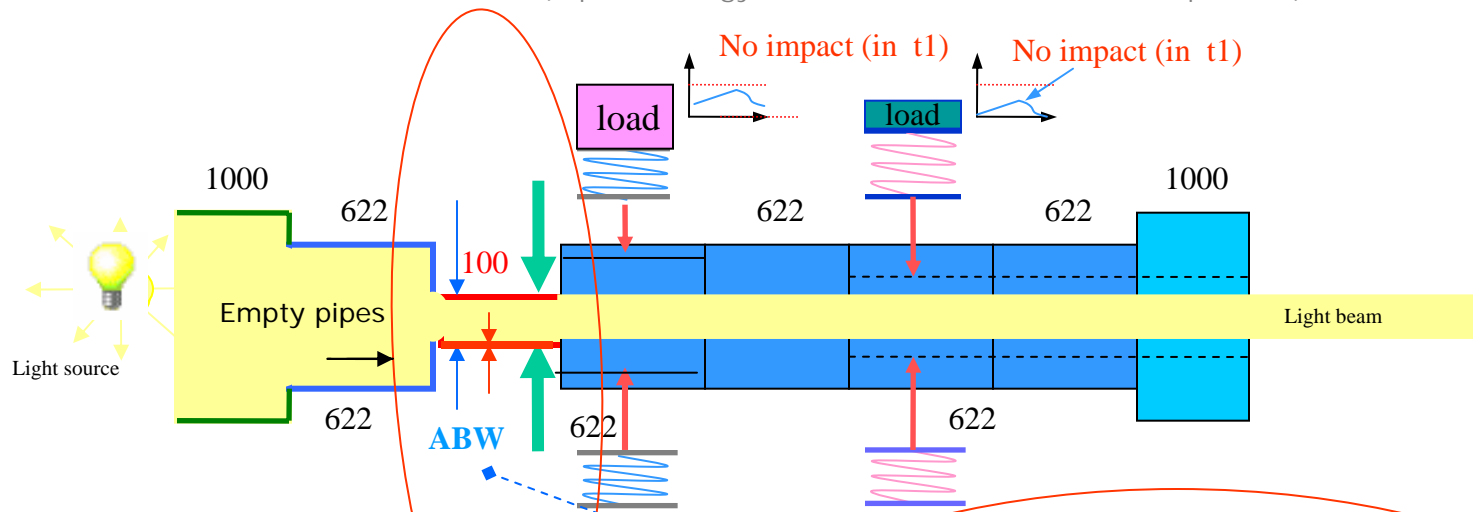
$$Abw = D_{bc} - CT$$

What is DBC

- **DBC** characterize instant high capacity bottleneck that **DOMINATE** on the path
- It covers situations when **routers** in the path **are overloaded** and sending **packets back to back** with its maximal rates
- We discovered that in most cases **only one node dominates** in the instant of our measurements (in our decision interval)

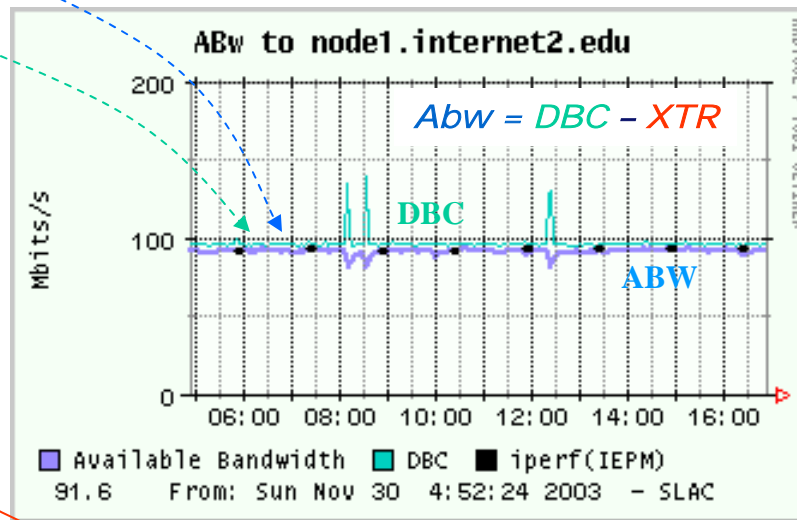
ABwE: Example of narrow link in the path

(Pipes analogy with different diameter and aperture)



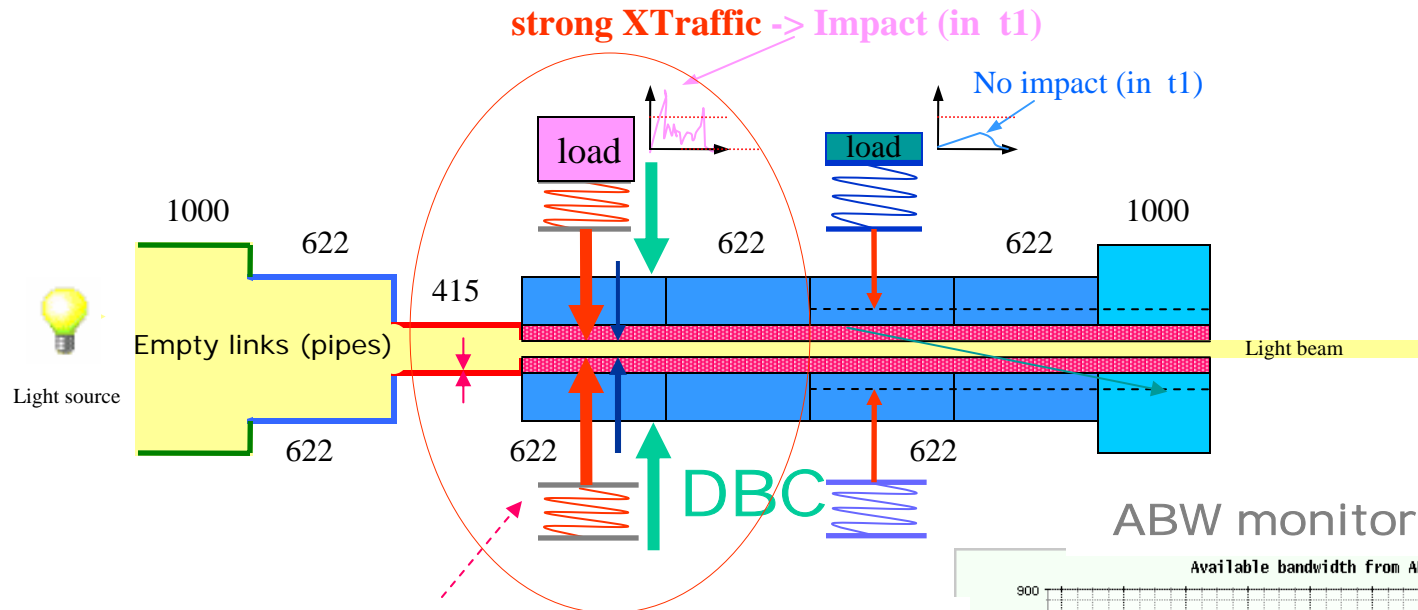
link that has domination effect on bandwidth

ABW monitor SLAC to UFL



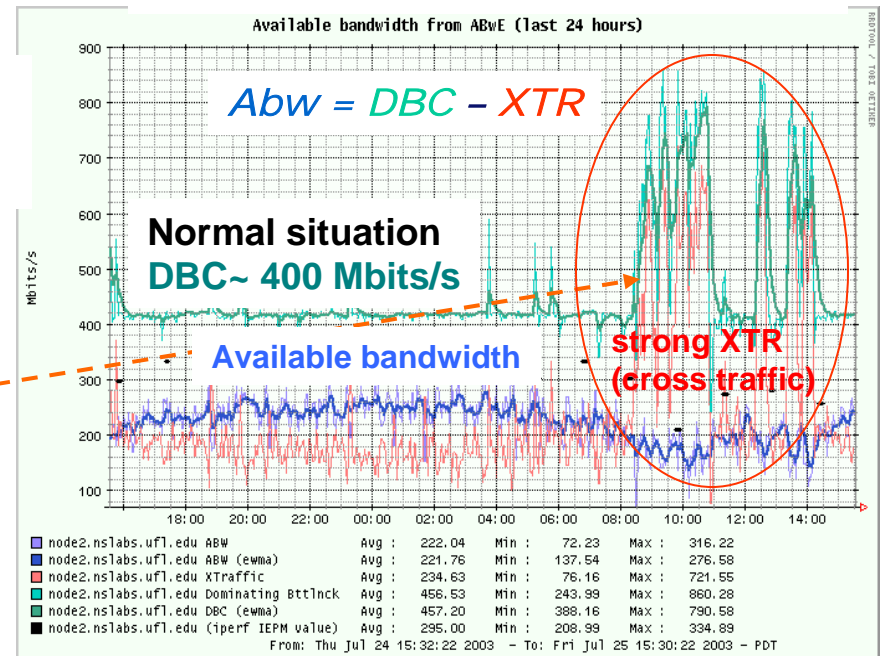
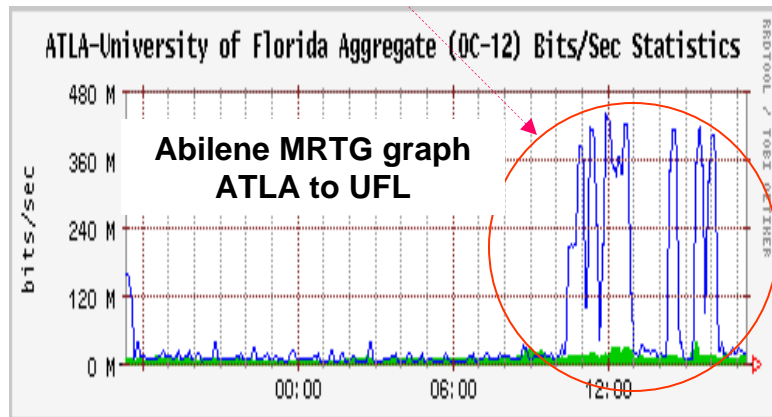
Example of heavy loaded link in the path

(Pipes analogy with different diameter and aperture)

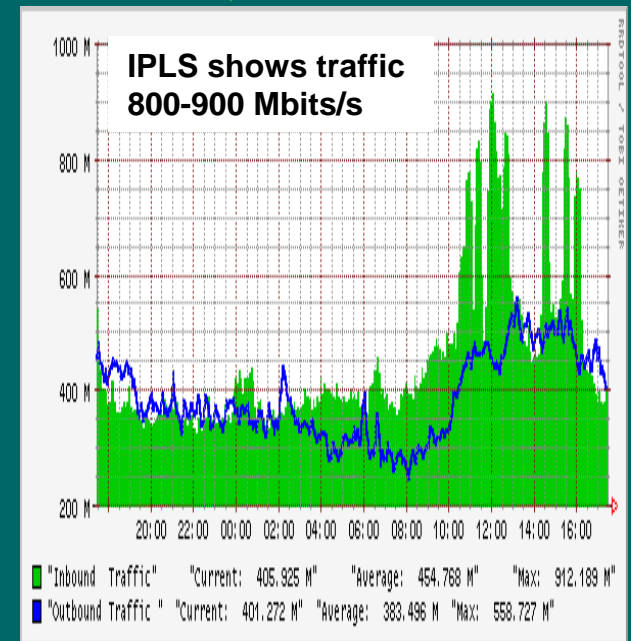
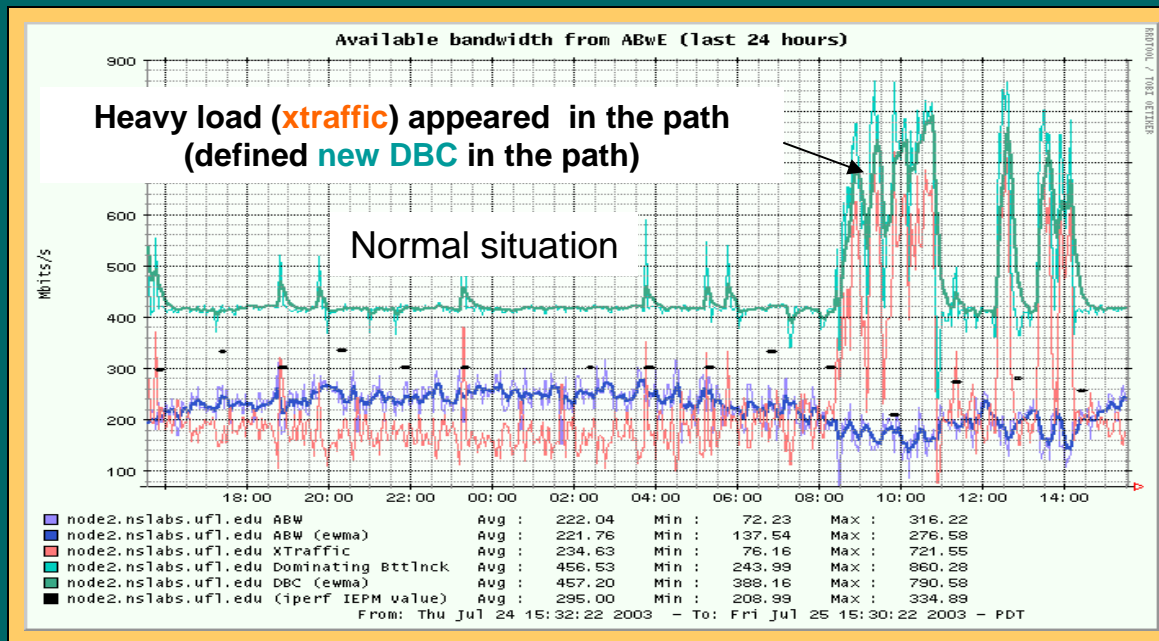
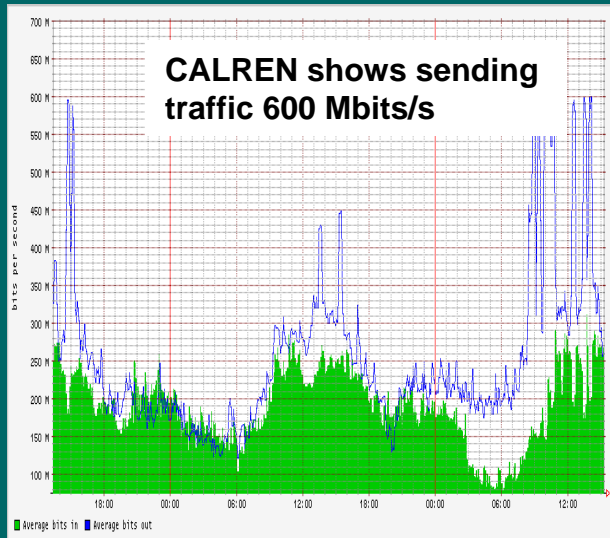


ABW monitor SLAC to UFL

Heavy load (**strong cross traffic**) appeared in the path
It shows new **DBC** in the path because this load dominates in whole path !



ABwE / MRTG match: TCP test to UFL

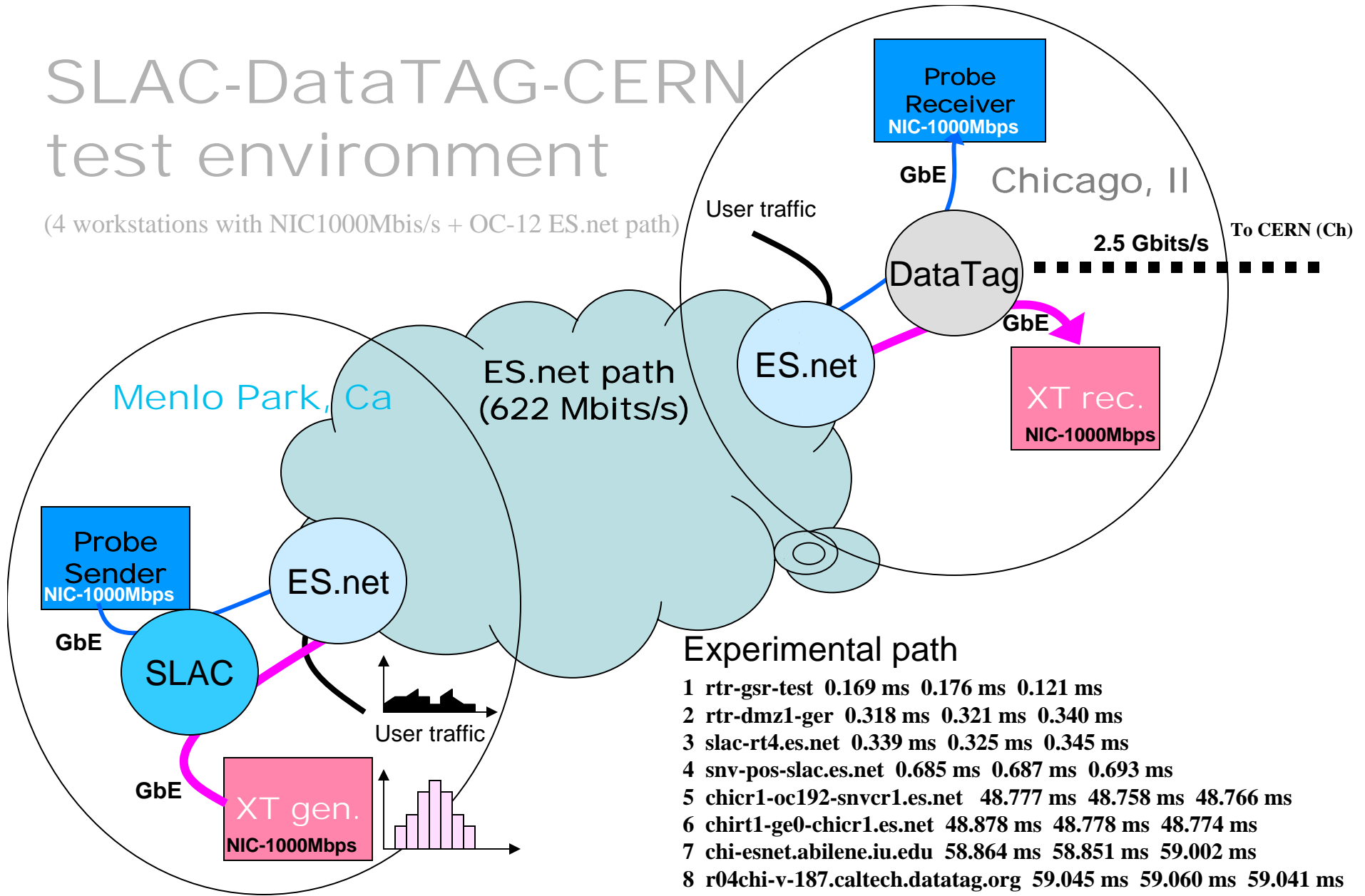


Confront ABwE results with other tools

Iperf, Pathload, Pathchirp

SLAC-DataTAG-CERN test environment

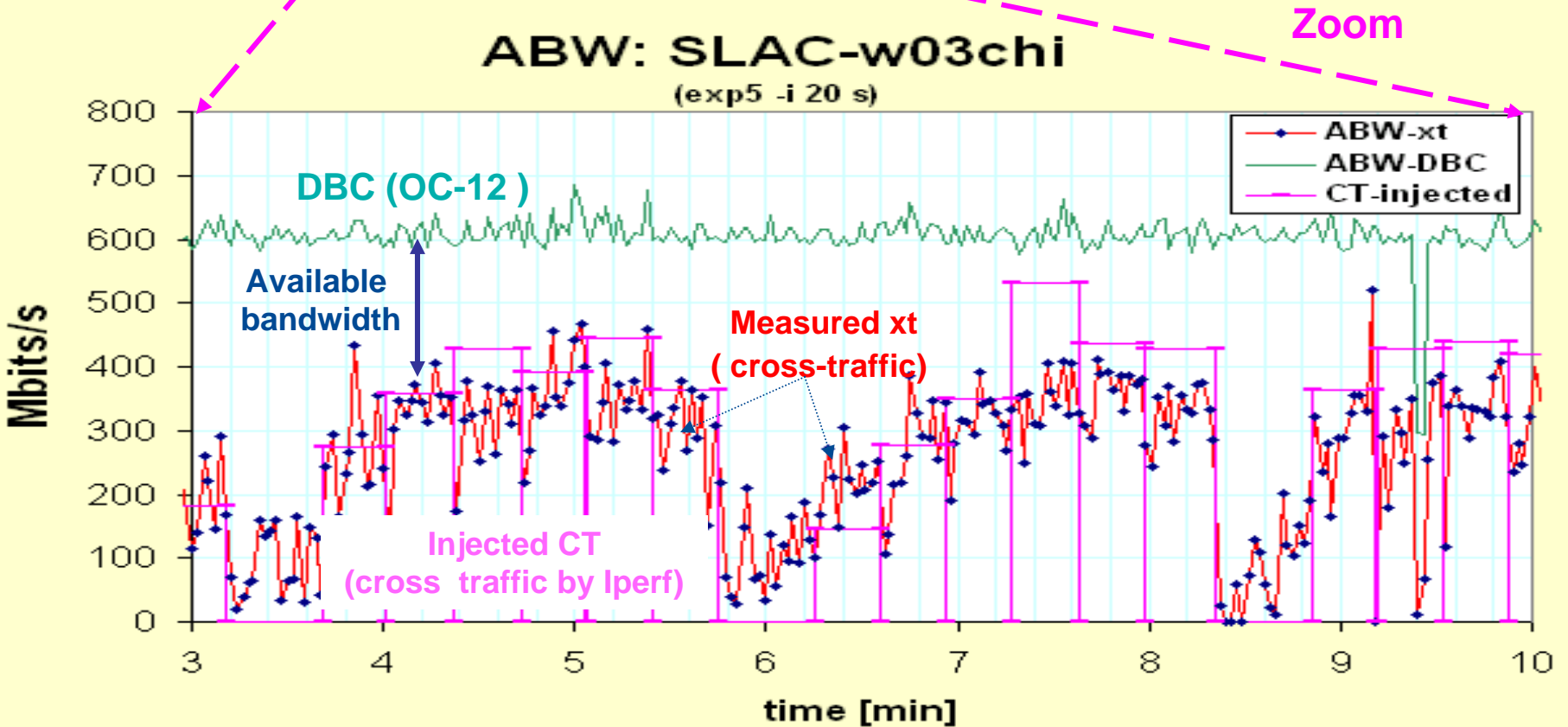
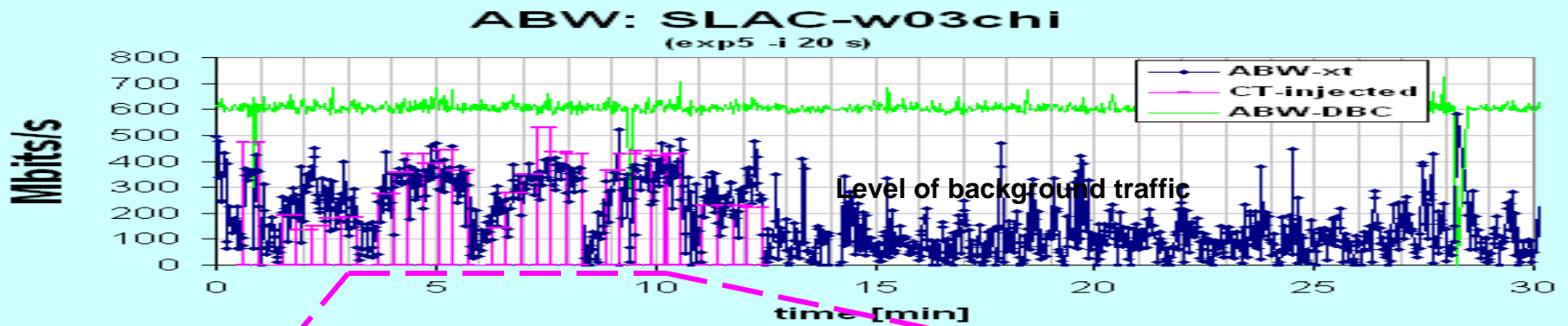
(4 workstations with NIC1000Mbis/s + OC-12 ES.net path)



 **Probing packets**
 **Injected Cross traffic**
 **User traffic (background)**

The match of the cross traffic

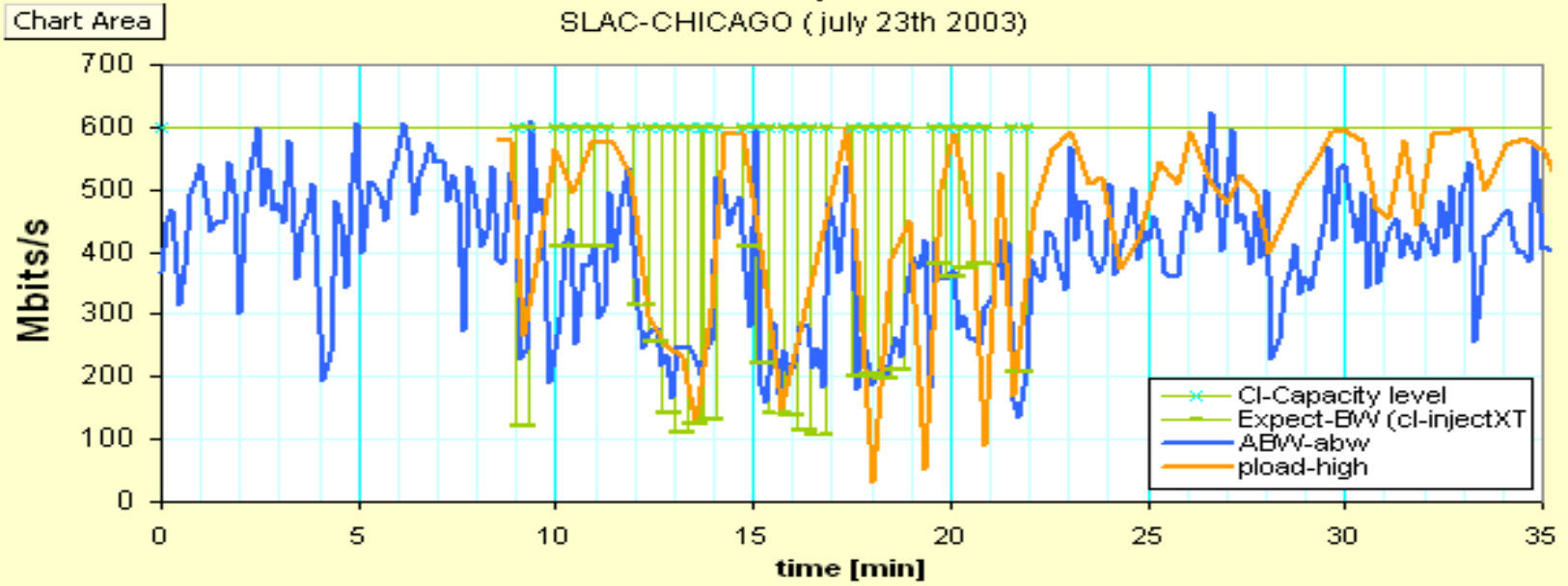
(ABW - XT compare to injection traffic generated by Iperf)



Conclusion: Iperf measure own performance which can approach DBC (in best case)

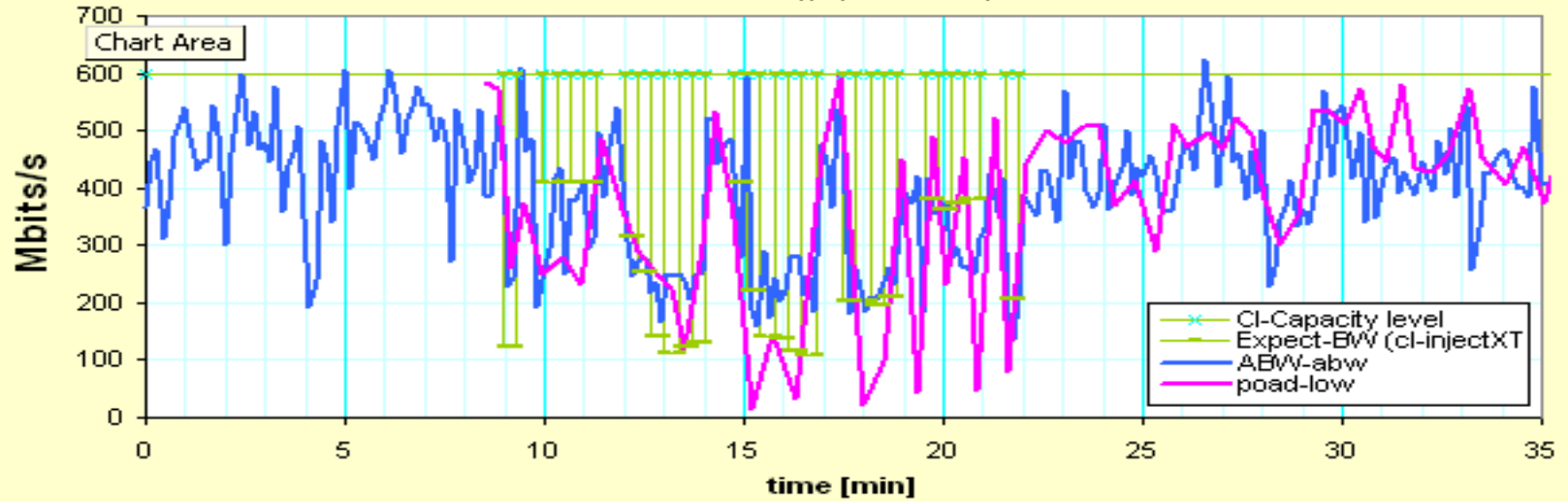
ABW vers. pathload

SLAC-CHICAGO (july 23th 2003)



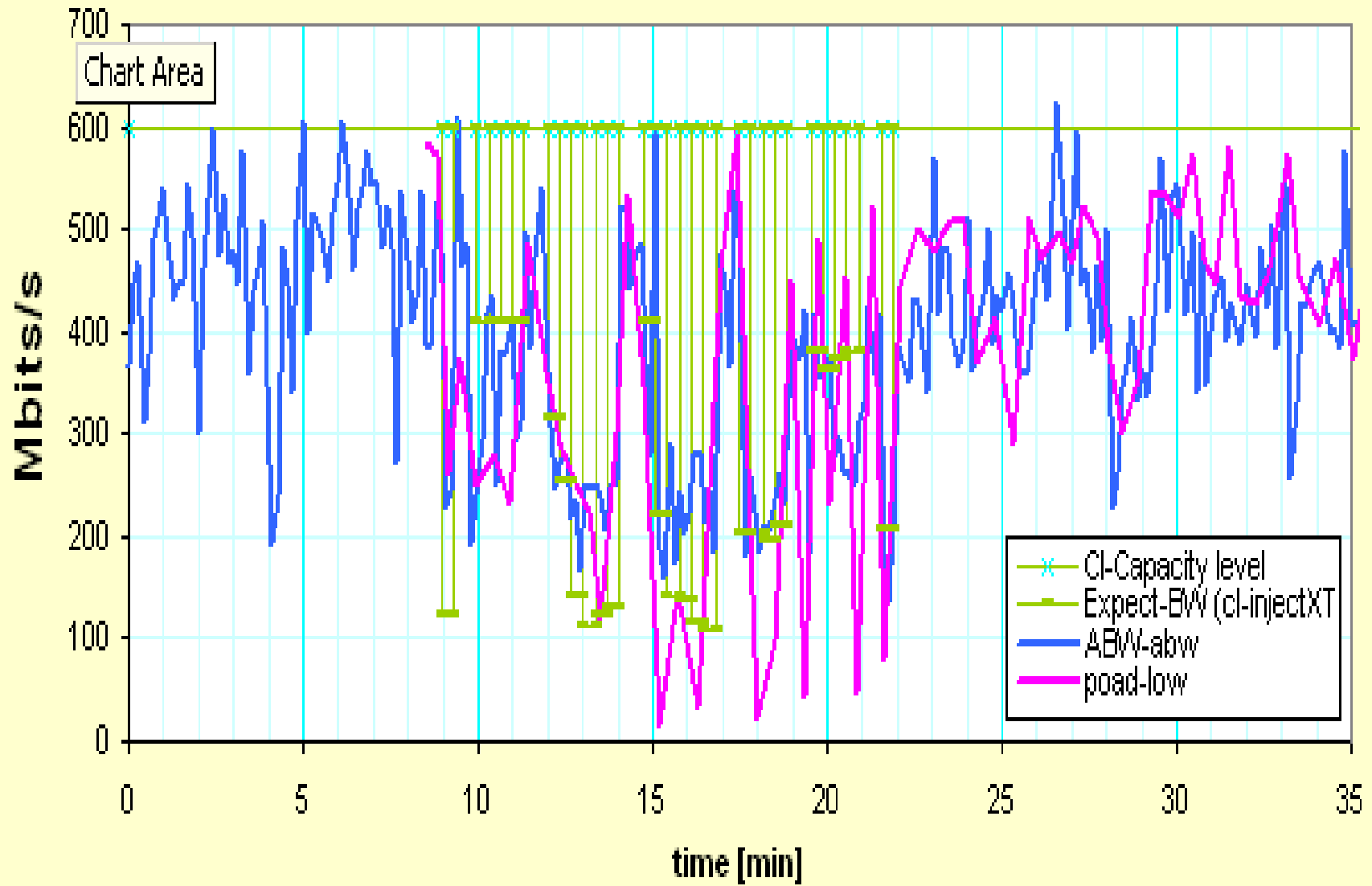
ABW vers. pathload

SLAC-CHICAGO (july 23th 2003)



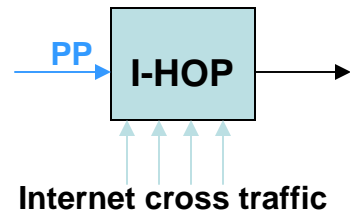
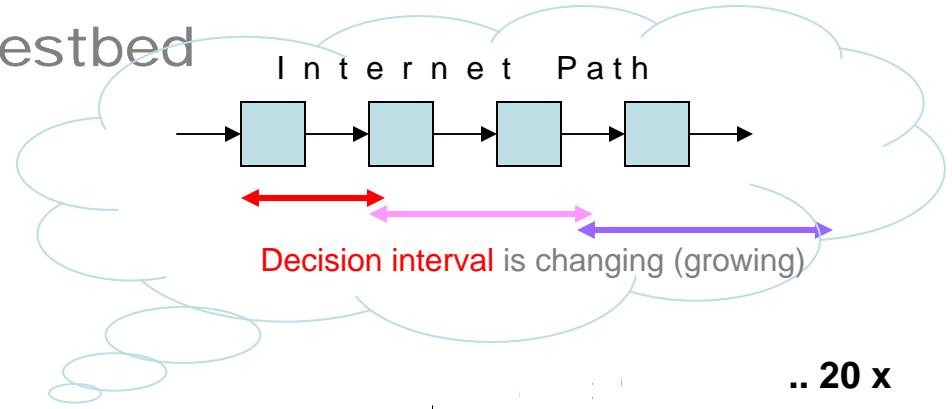
ABW vers. pathload

SLAC-CHICAGO (july 23th 2003)

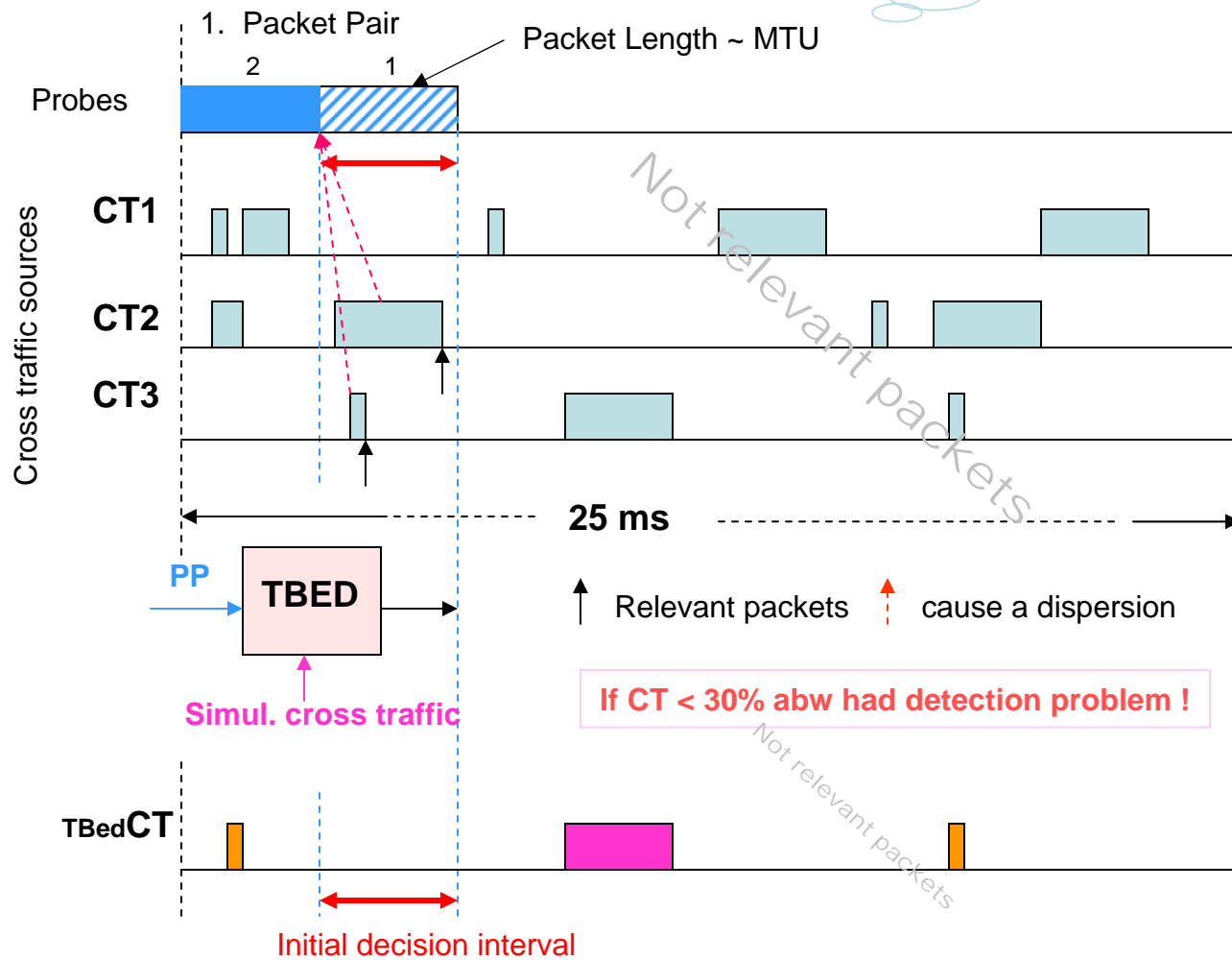


What we learned from CAIDA testbed

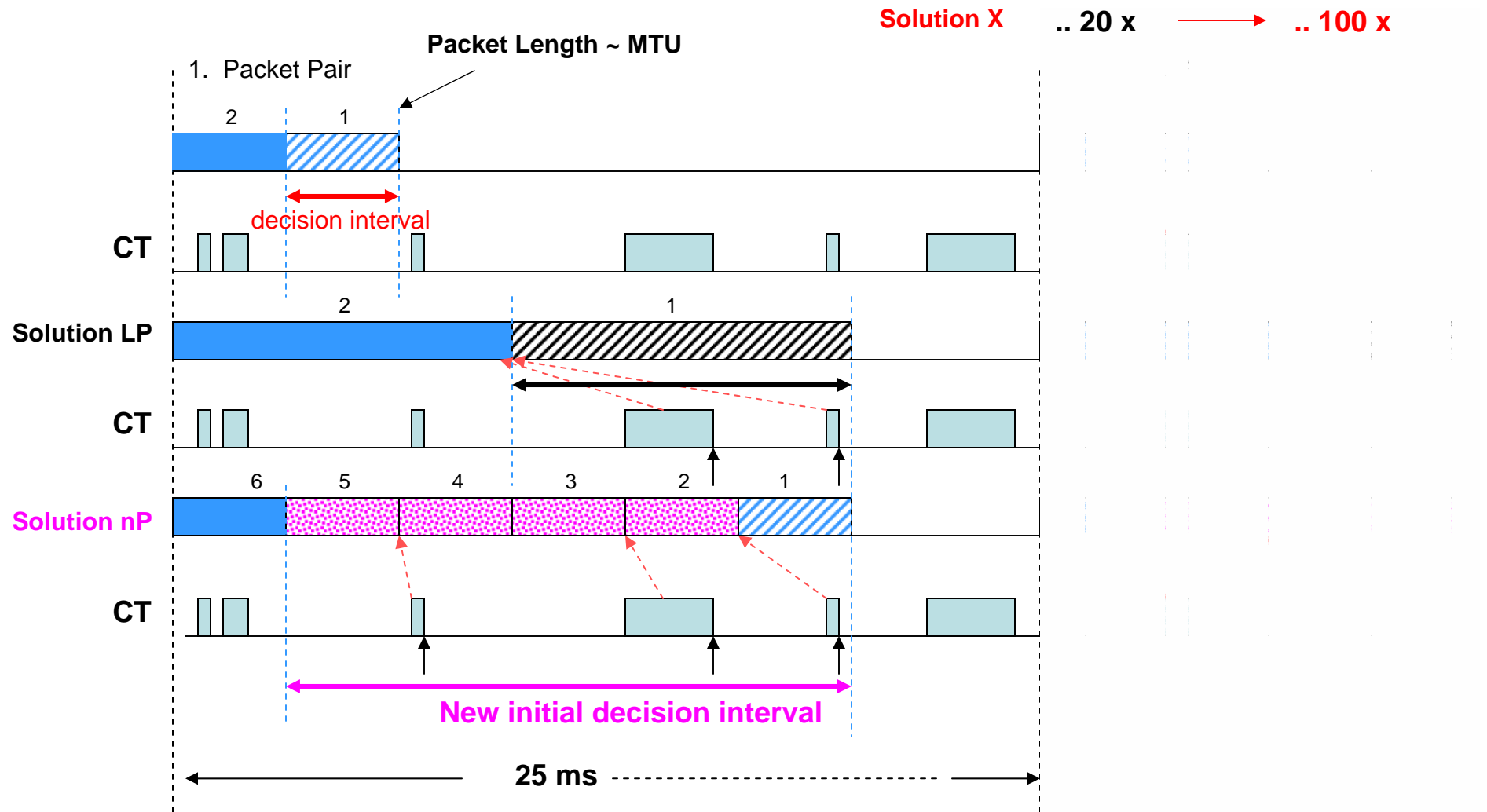
Internet HOP/HOPS vers. Testbed



.. 20 x



How to improve "detection effectiveness"

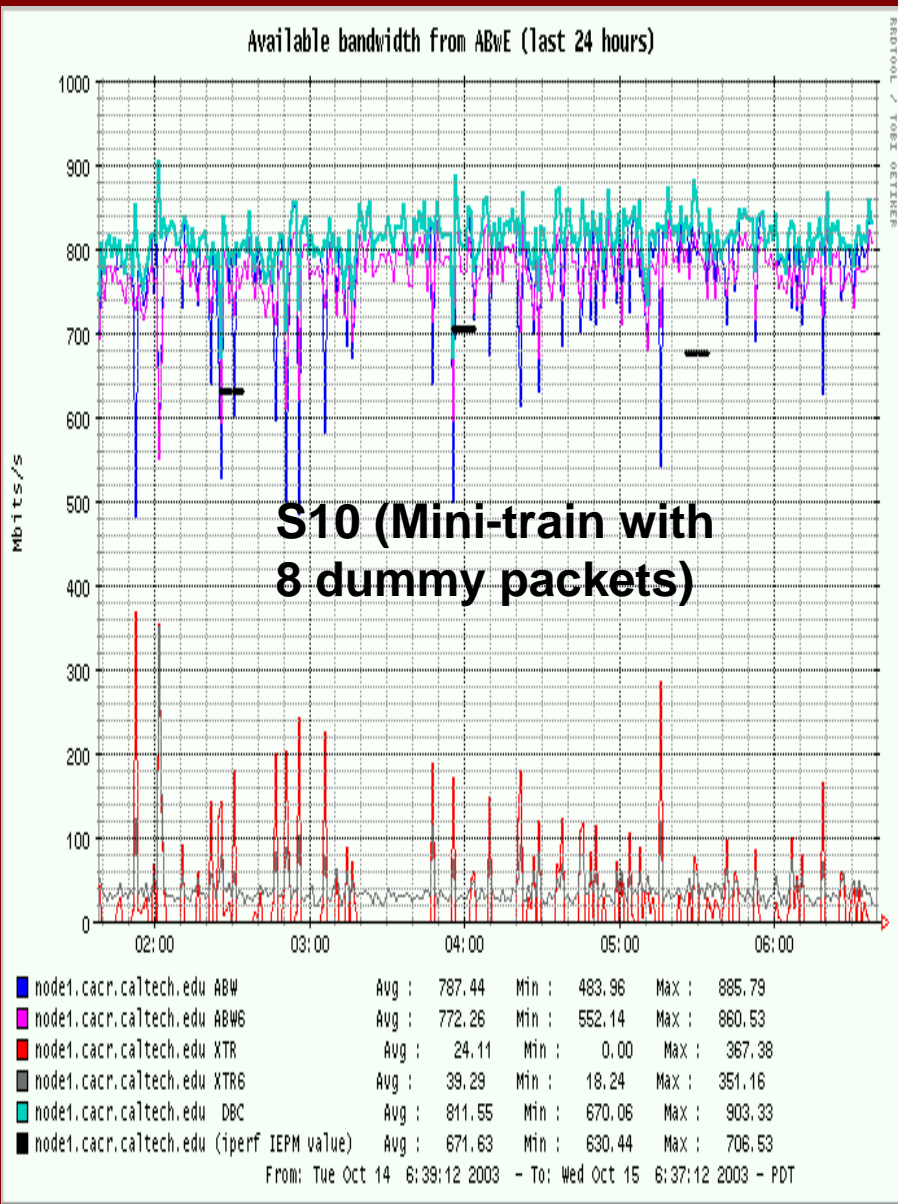
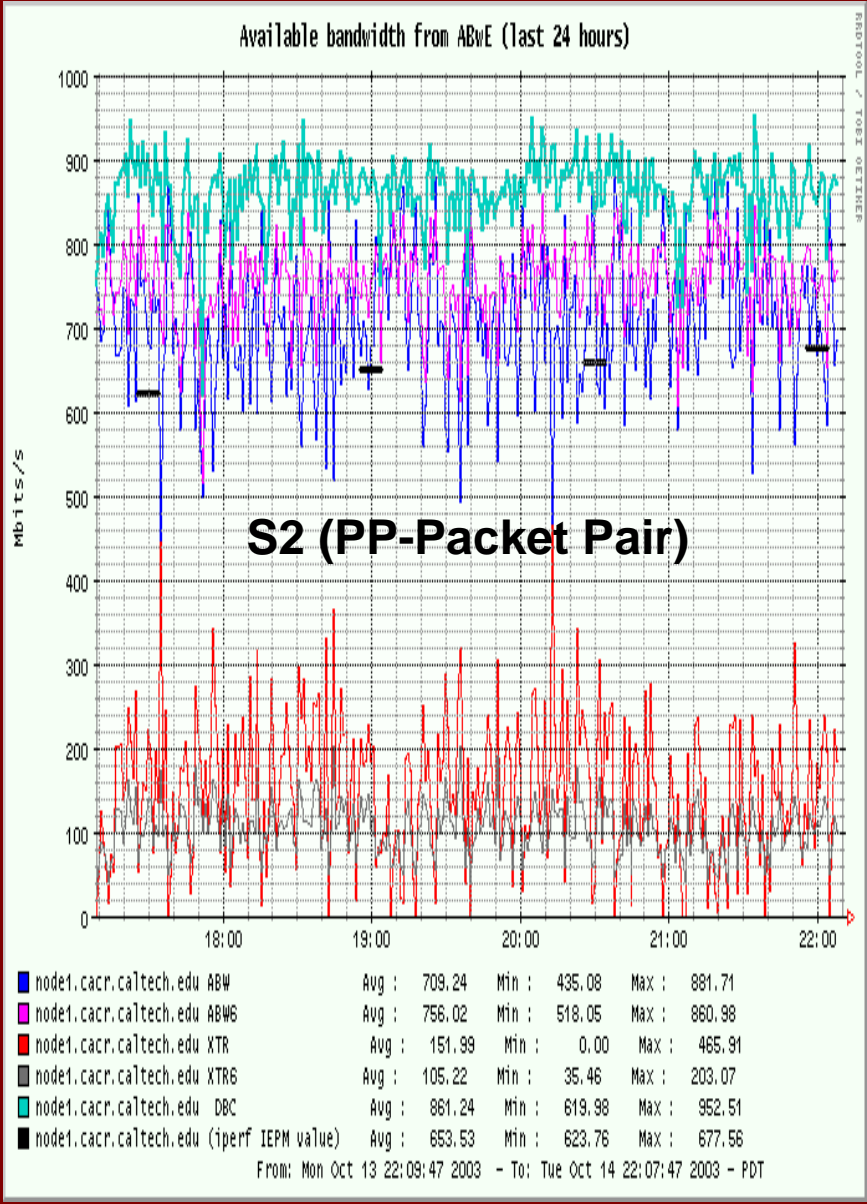


Solution X .. 20 x → .. 100 x

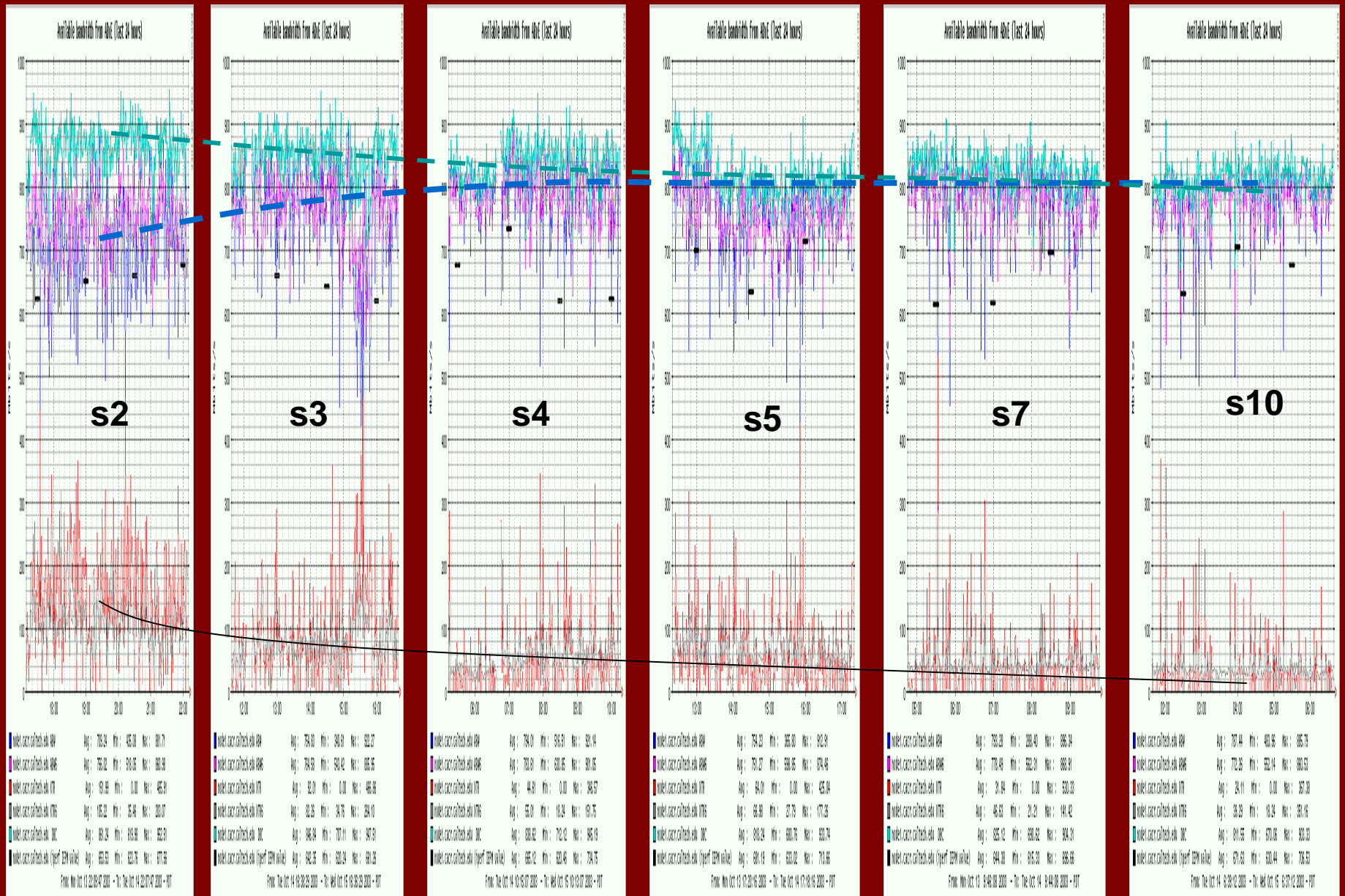
Solution LP – Long packets (9k) ↑ Relevant packets ↑ cause a dispersion
(creates micro-bottlenecks)

Solution nP – n dummy Packets (mini-train)

Measurement time
0.5 s to 2.5 s



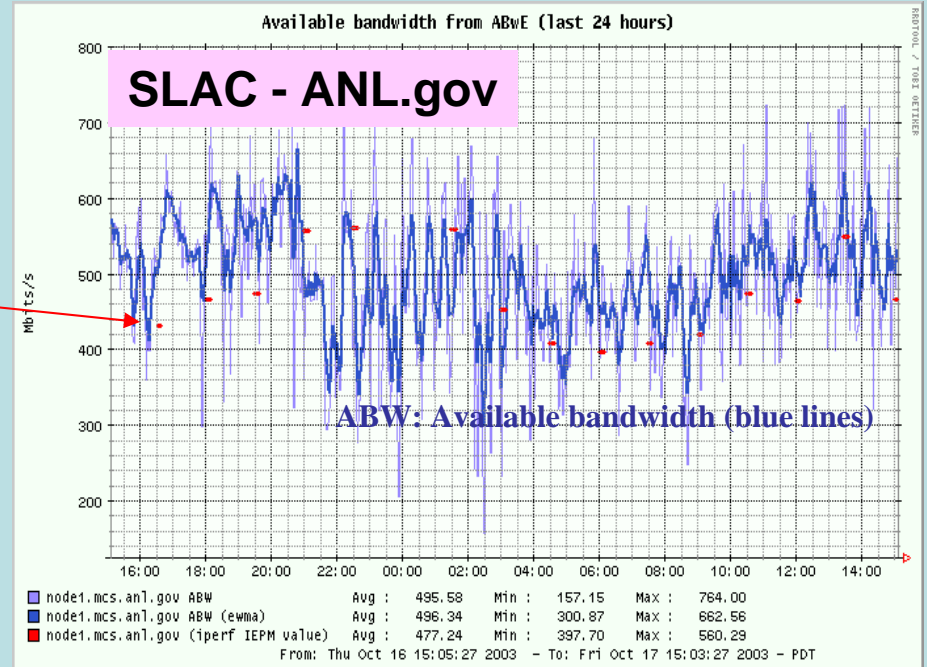
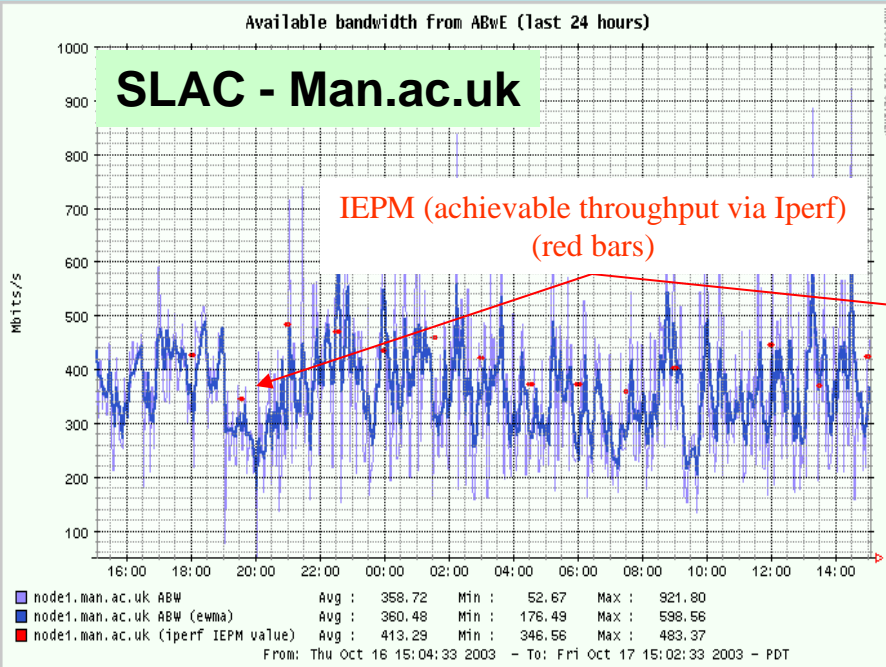
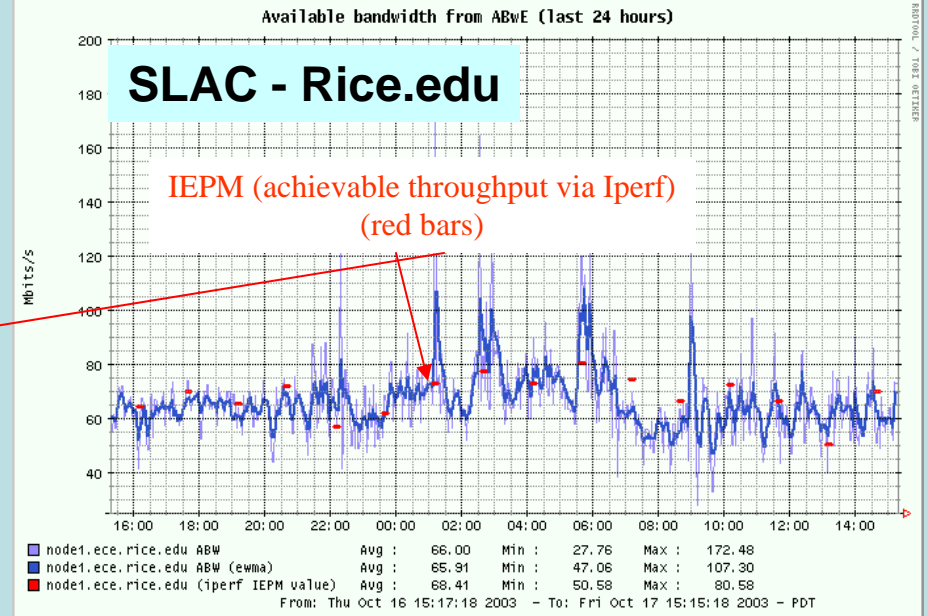
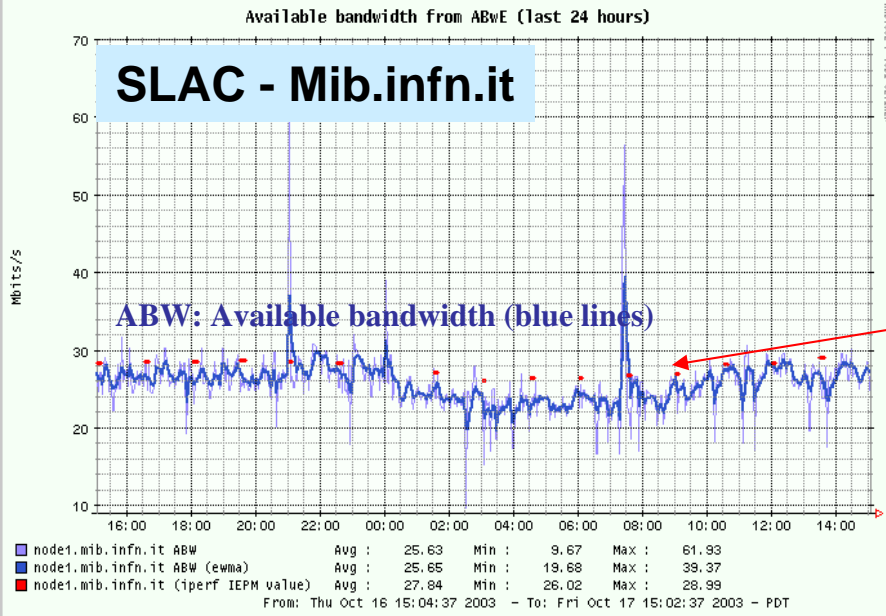
PP versus TRAIN: ABW and DBC merge in TRAIN samples
(SLAC-CALTECH path)



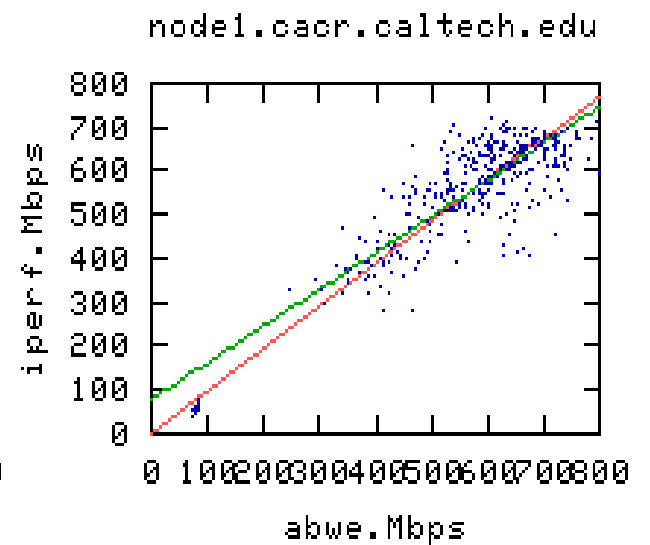
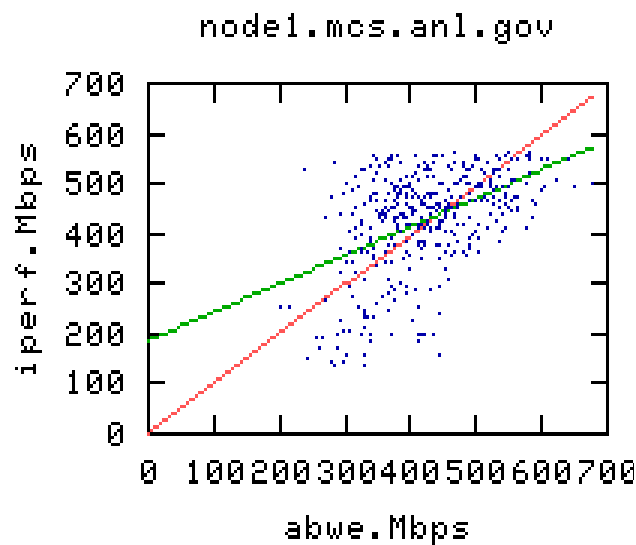
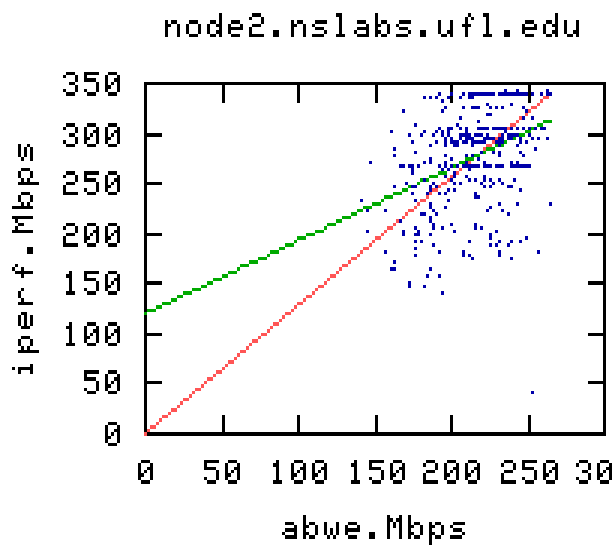
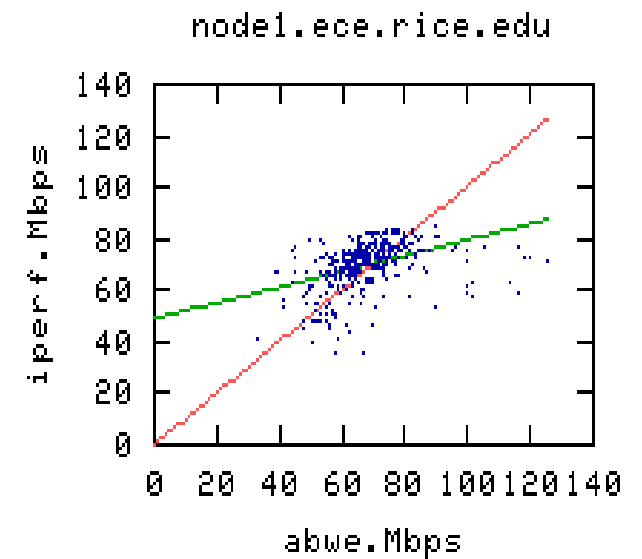
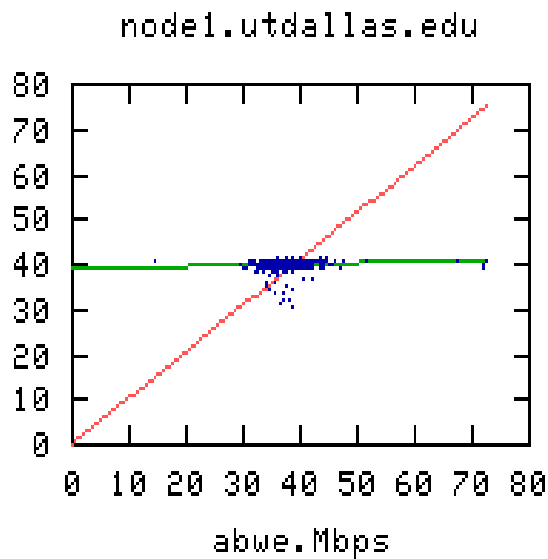
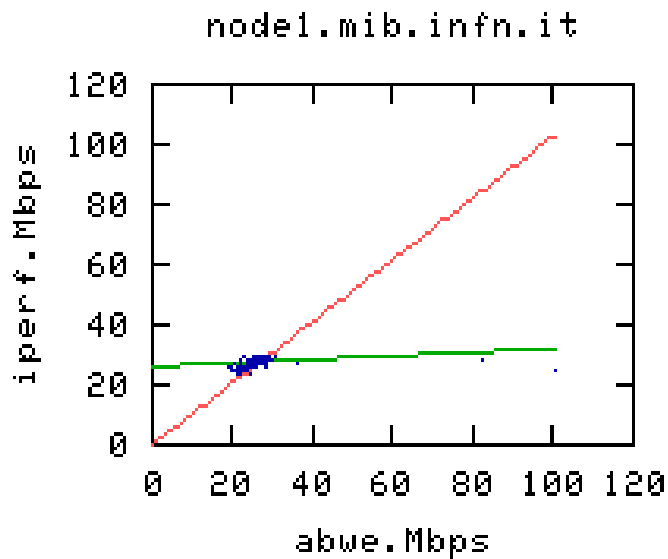
PP versus TRAIN: **ABW** and **DBC merge** in TRAIN samples
(SLAC-CALTECH path)

Compare long term
Bandwidth statistics
on real paths

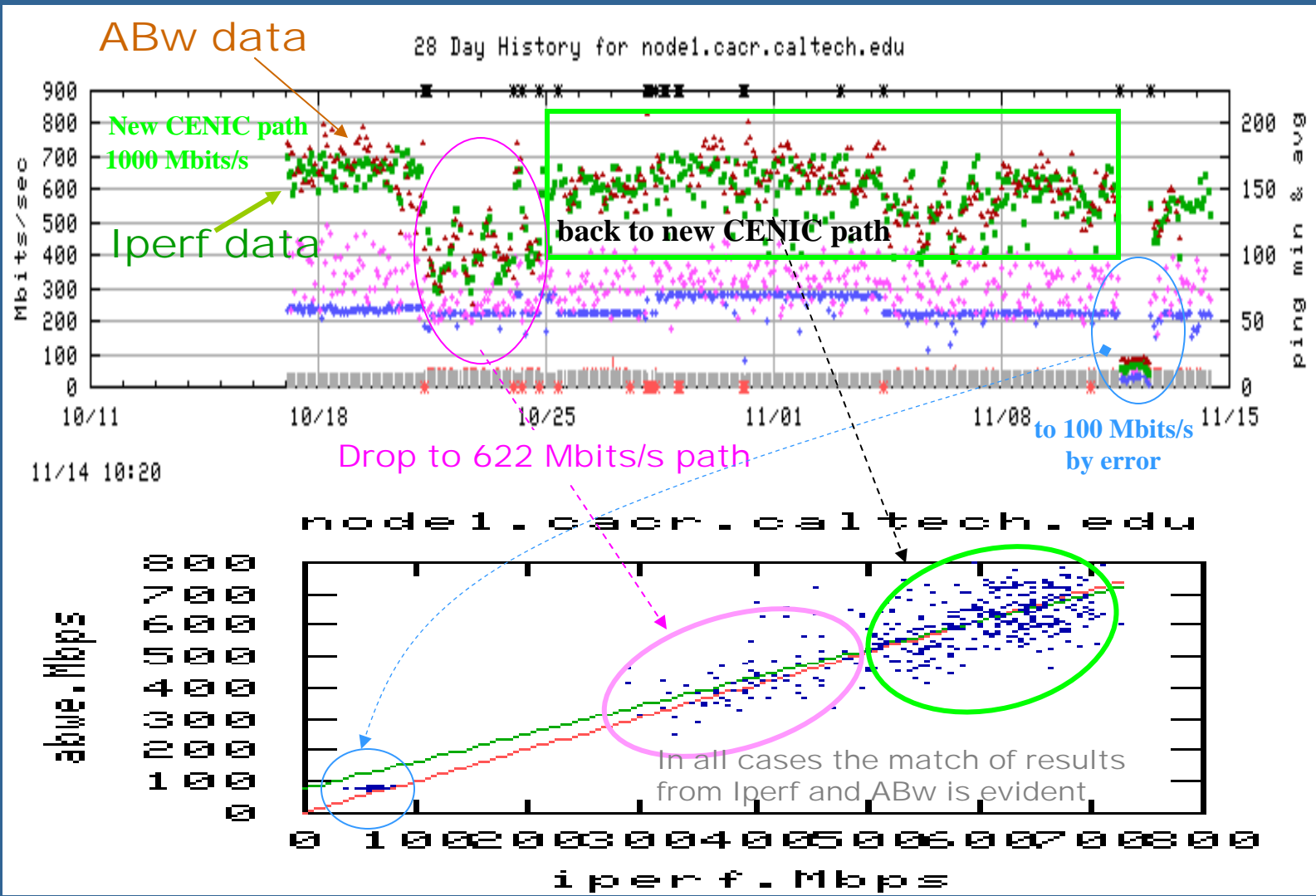
ESNET, Abilene, Europe



IEPM-Iperf vers. ABW (24 hours match)



Scatter plot graphs
 Achievable throughput via Iperf versus ABw
 on different paths (range 20–800 Mbits/s)
 (28 days history)



28 days bandwidth history

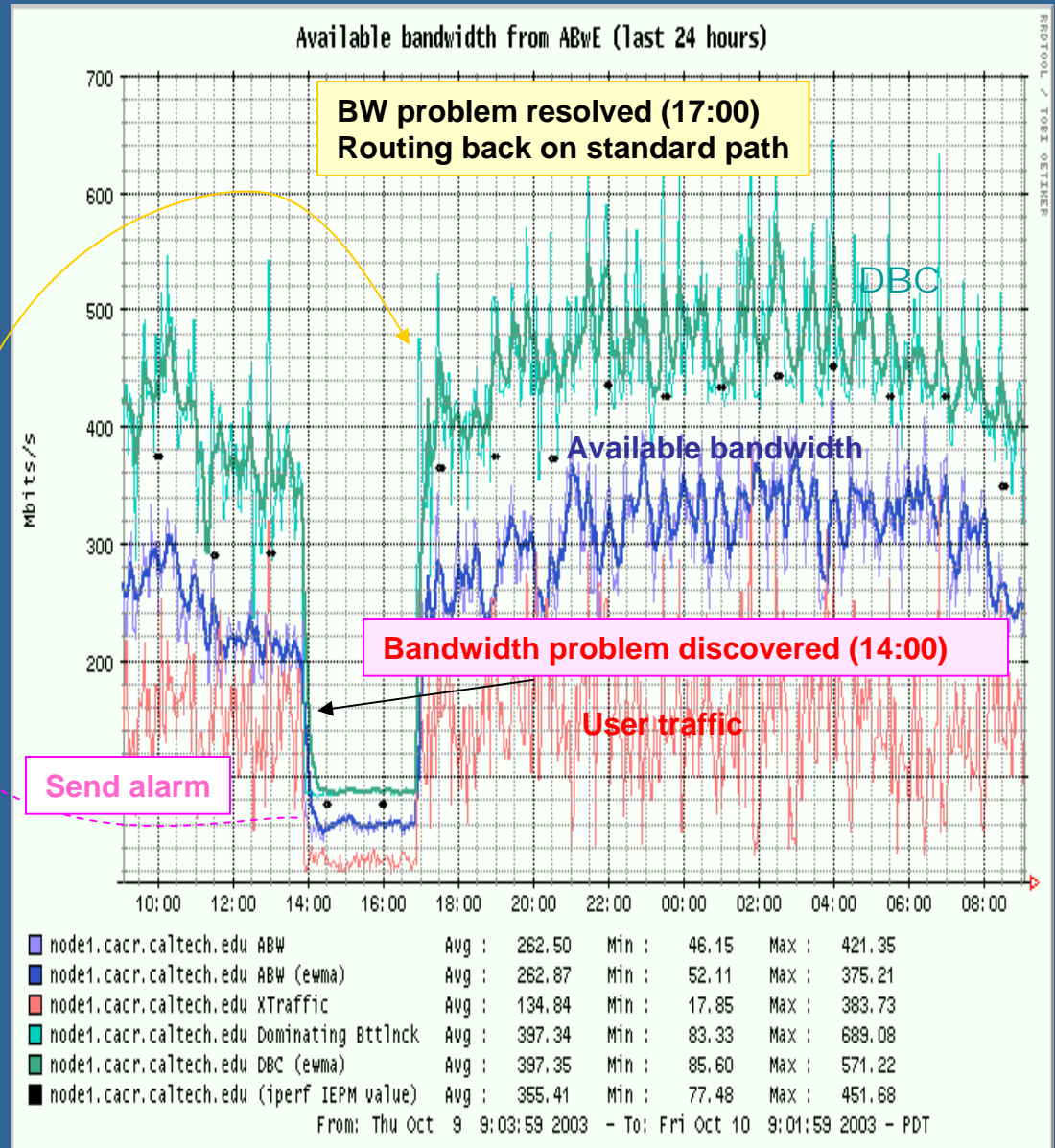
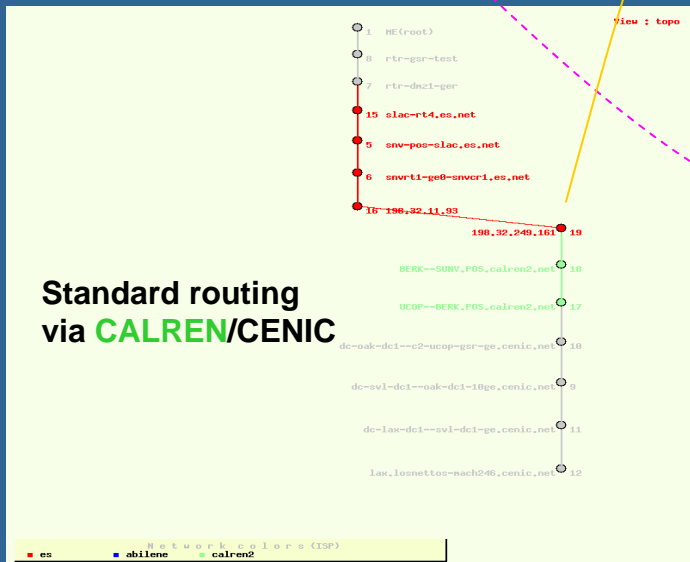
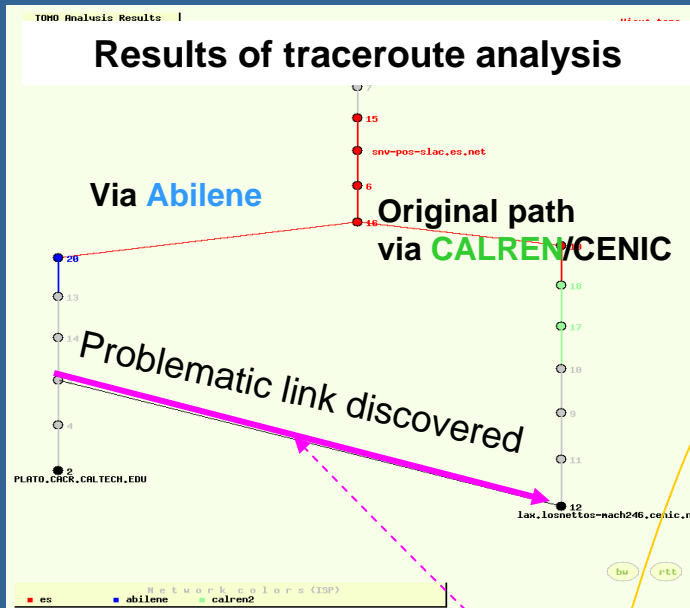
During this time we can see several different situations caused by different routing from SLAC to CALTECH

What we can detect with continues bandwidth monitoring

- Immediate bandwidth on the path
- Automatic routing changes when line is broken (move to backup lines)
- Unexpected Network changes (Routing changes between networks, etc.)
- Line updates (155 -> 1Giga, etc.)
- Extreme heavy load

ABw as Troubleshooting tool

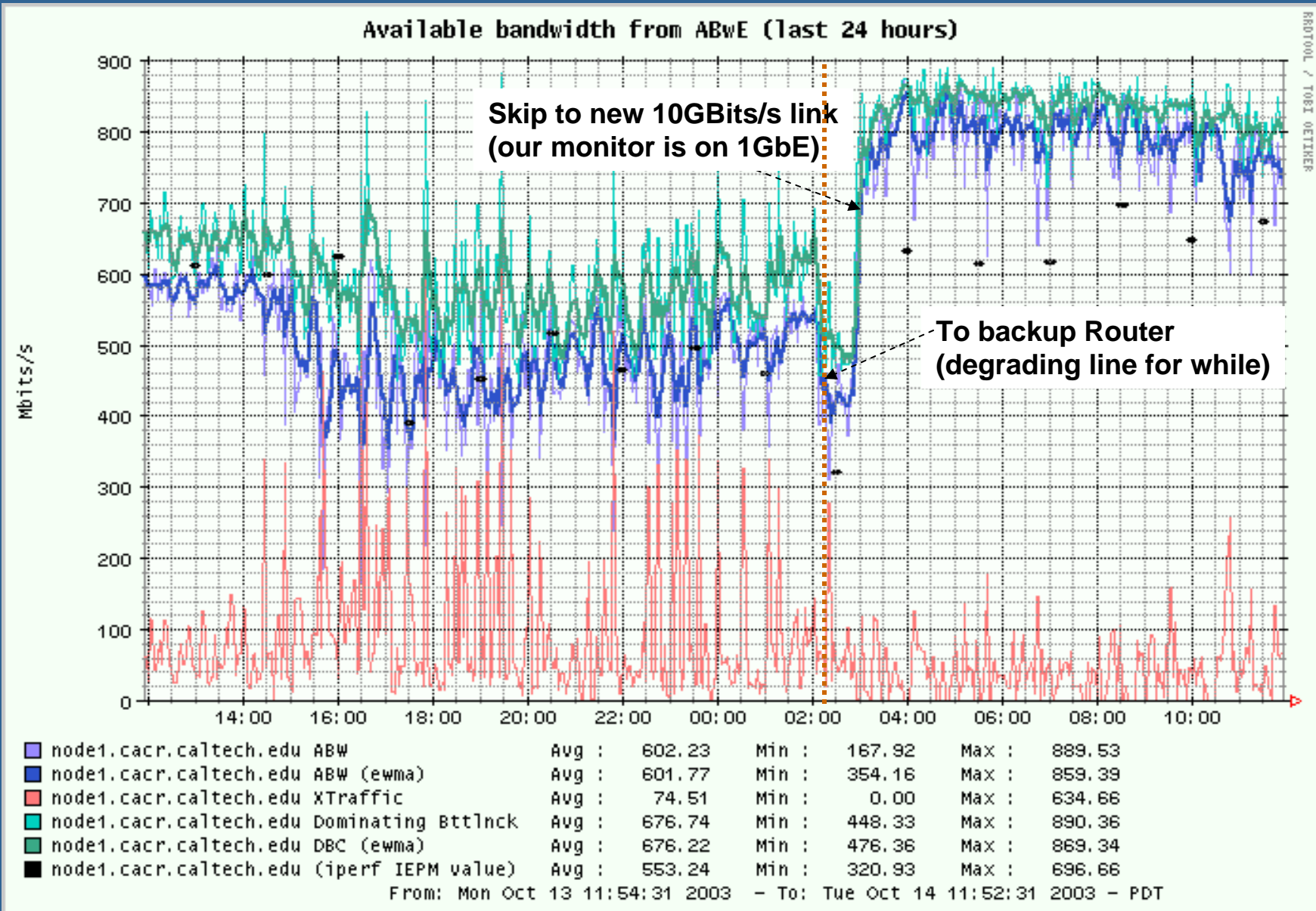
(Discovering Routing problems and initiate alarming)



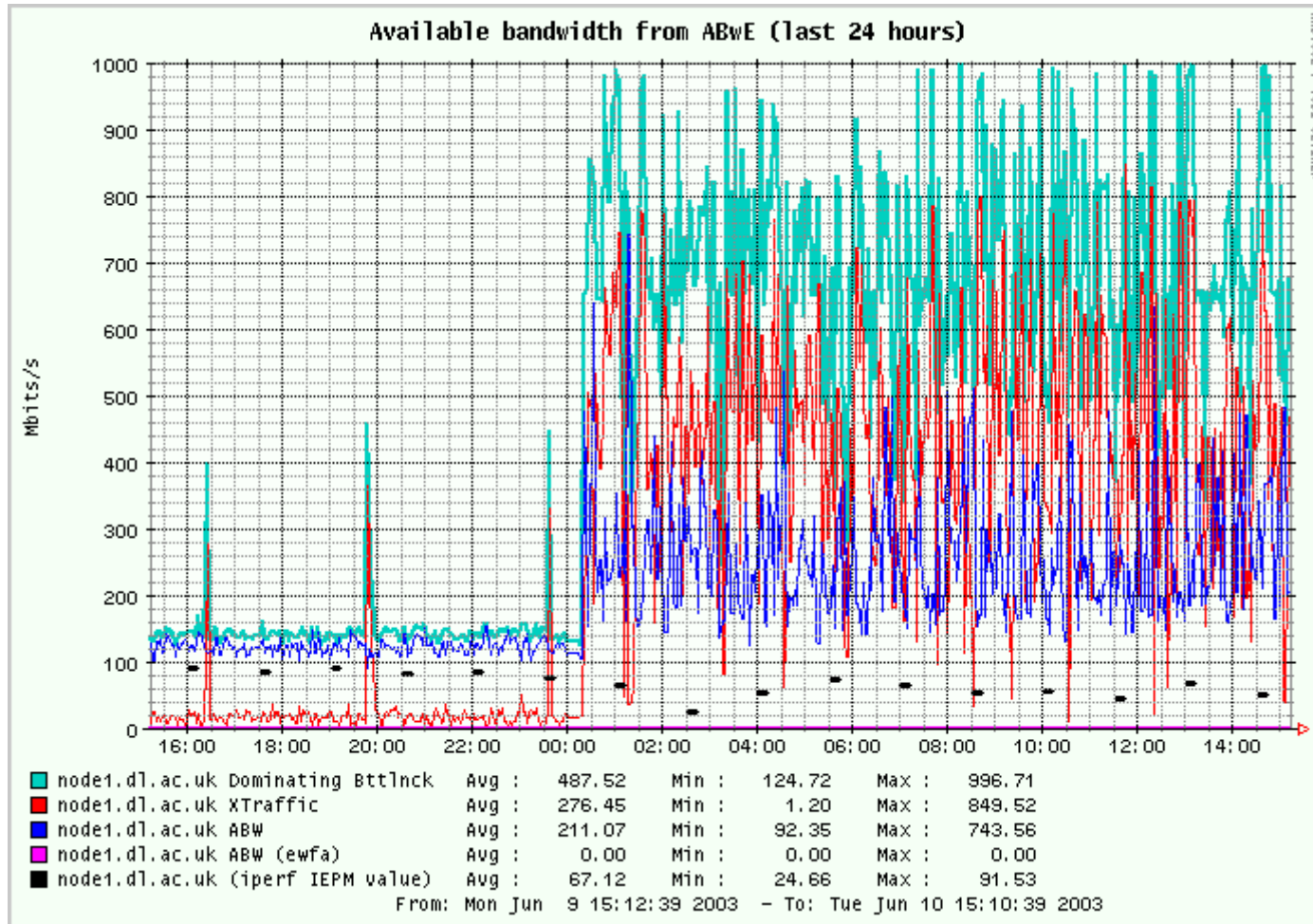
(Example from SLAC – CENIC path)

SLAC - CENIC path upgrade from 1 to 10 Gigabit

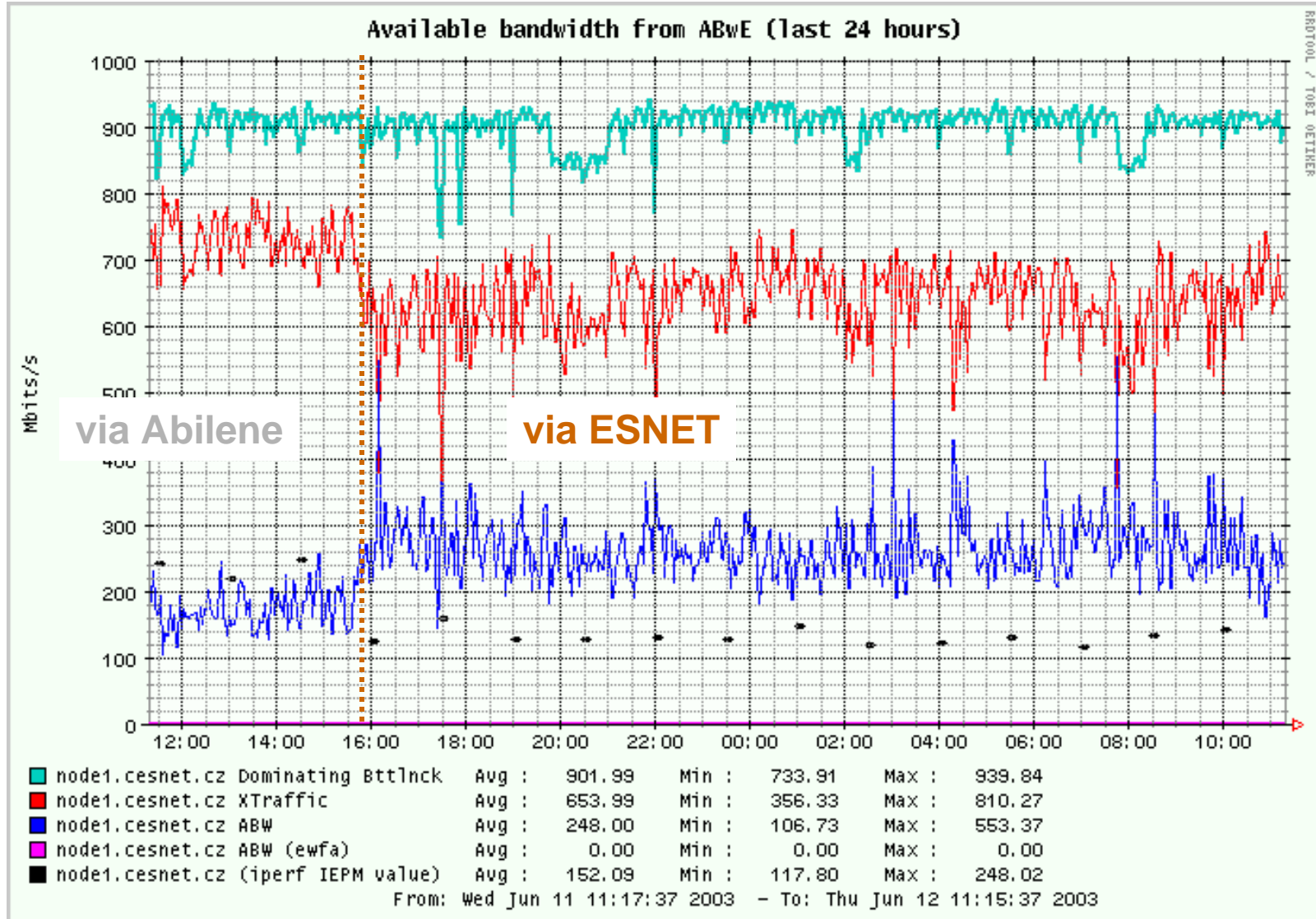
(Current monitoring machines allow monitor traffic in range $1 < 1000$ Mbits only)

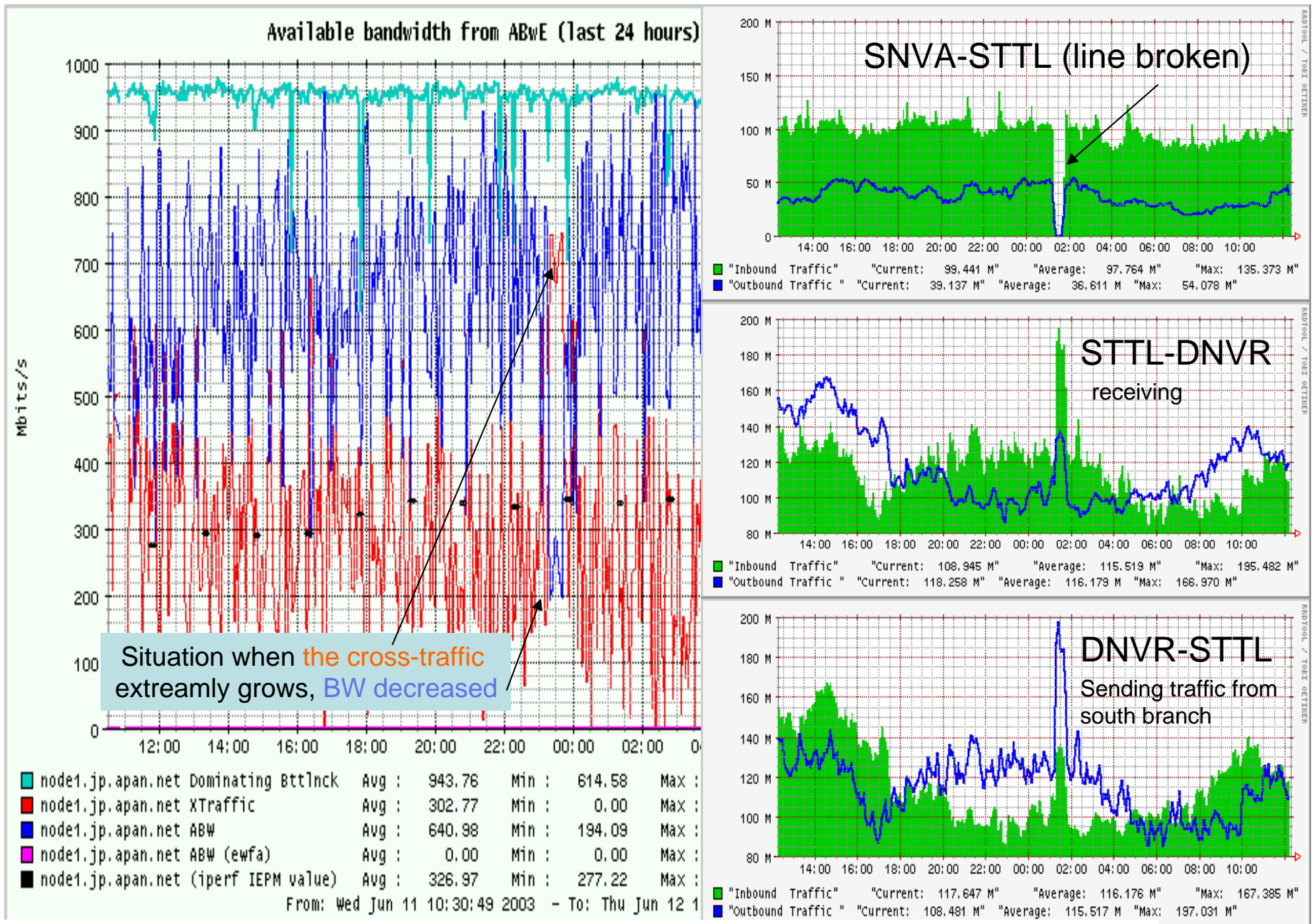


Upgrade 155Mbits/s line to 1000Mbits/s at dl.uk



SLAC changed routing to CESNET





Abilene - automatic rerouting - June 11, 2003

Typical SLAC traffic (long data transfer when physical experiment ends)

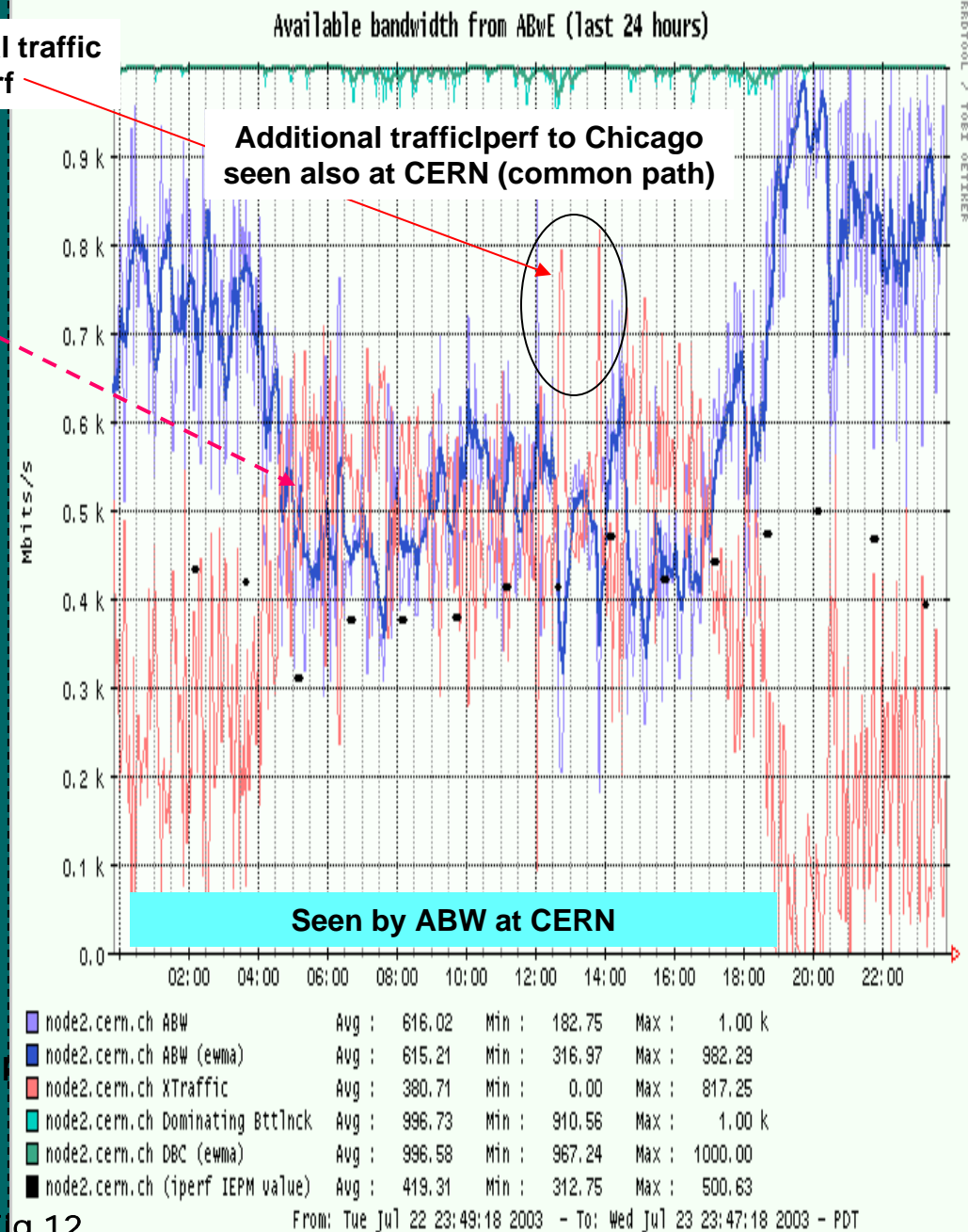
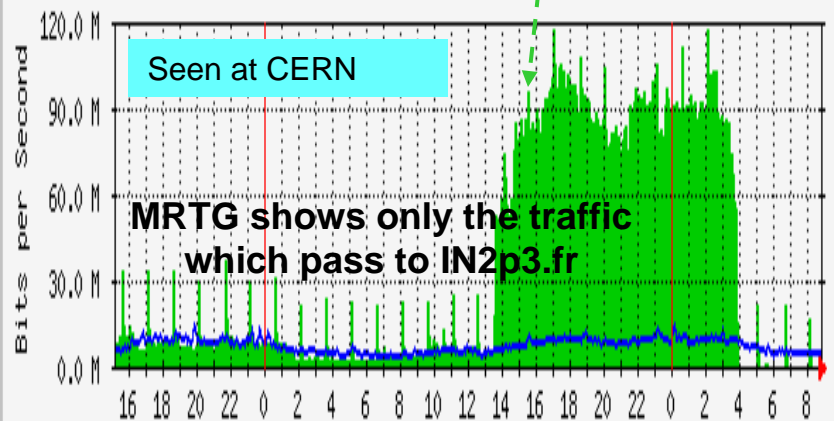
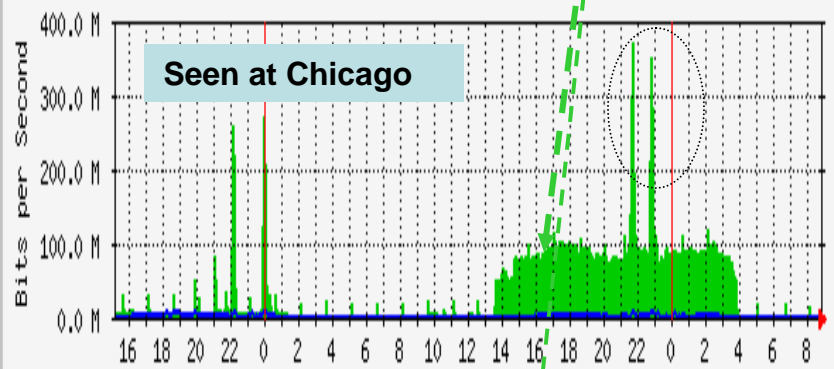
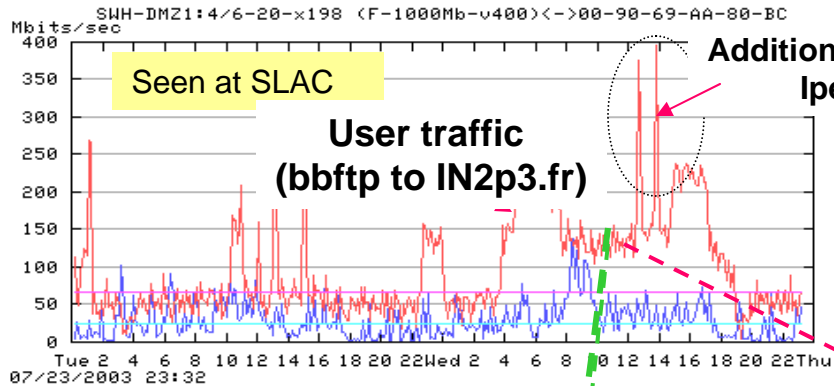


Fig.12

Abing new ABwE tool

- Interactive (reply < 1 second)
- Very low impact on the network traffic (40 packets to get value for destination)
- Simple and robust (responder can be installed on any machine on the network)
- Keyword function for protecting the client-server communication
- Measurements in both directions
- Same resolution as other similar methods

<http://www-iepm.slac.stanford.edu/tools/abing>

Thank you

References:

<http://moat.nlanr.net/PAM2003/PAM2003papers/3781.pdf>

<http://www-iepm.slac.stanford.edu/tools/abing>