# The Rationale of the Current Optical Networking Initiatives

www.science.uva.nl/~delaat

## Cees de Laat

Faculty of Science

# The Rationale of the Current Optical Networking Initiatives

www.science.uva.nl/~deλaat

## Cees de Λaat

EU

SURFnet

### University of Amsterdam
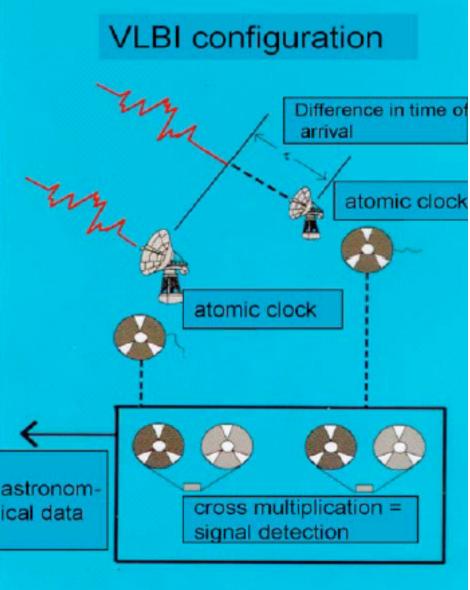
SARA

NIKHEF

optiputer

# VLBI

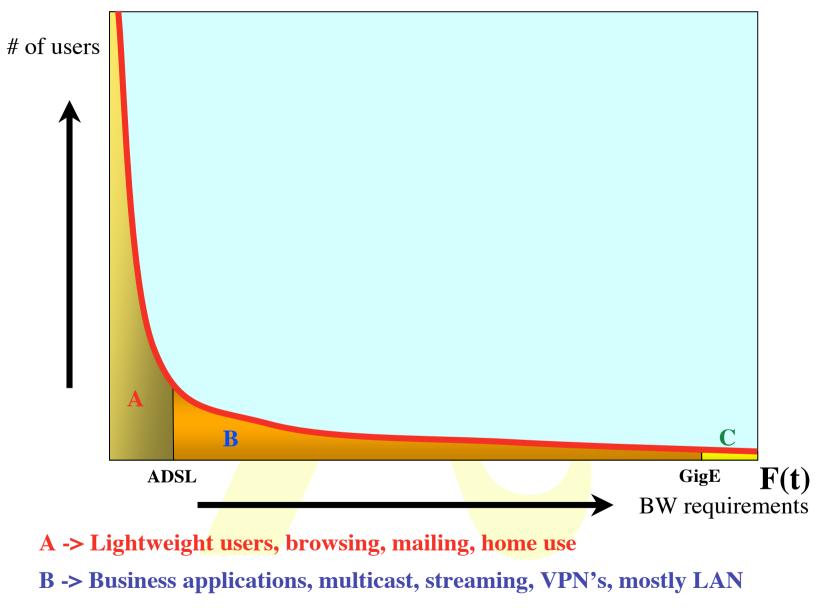ger term VLBI is easily capable of generating many Gb of data per
The sensitivity of the VLBI array scales w
(=data-rate) and there is a strong push to
. Rates of 8Gb/s or more are entirely feasibl
der development. It is expected that parall
orrelator will remain the most efficient approa
s distributed processing may have an appli
lti-gigabit data streams will aggregate into la
or and the capacity of the final link to the da
tor.



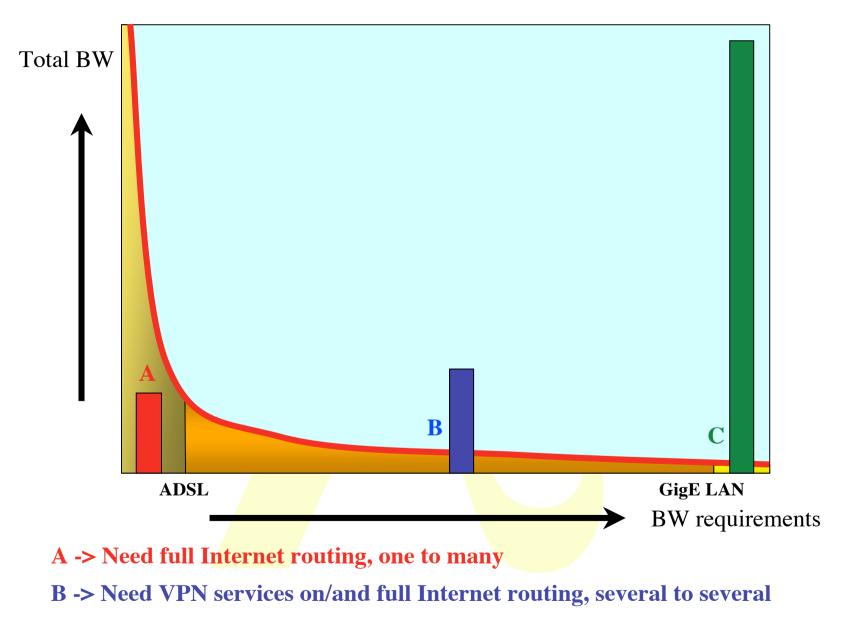*Westerbork Synthesis Radio Telescope -*
*Netherlands*



VLBI configuration

Difference in time of arrival

atomic clock

atomic clock

astronom-ical data

cross multiplication = signal detection

# Know the user

# of users

A

B

C

ADSL

GigE

F(t)

BW requirements

A -> Lightweight users, browsing, mailing, home use

B -> Business applications, multicast, streaming, VPN's, mostly LAN

C -> Special scientific applications, computing, data grids, virtual-presence

# What the user

A -> Need full Internet routing, one to many

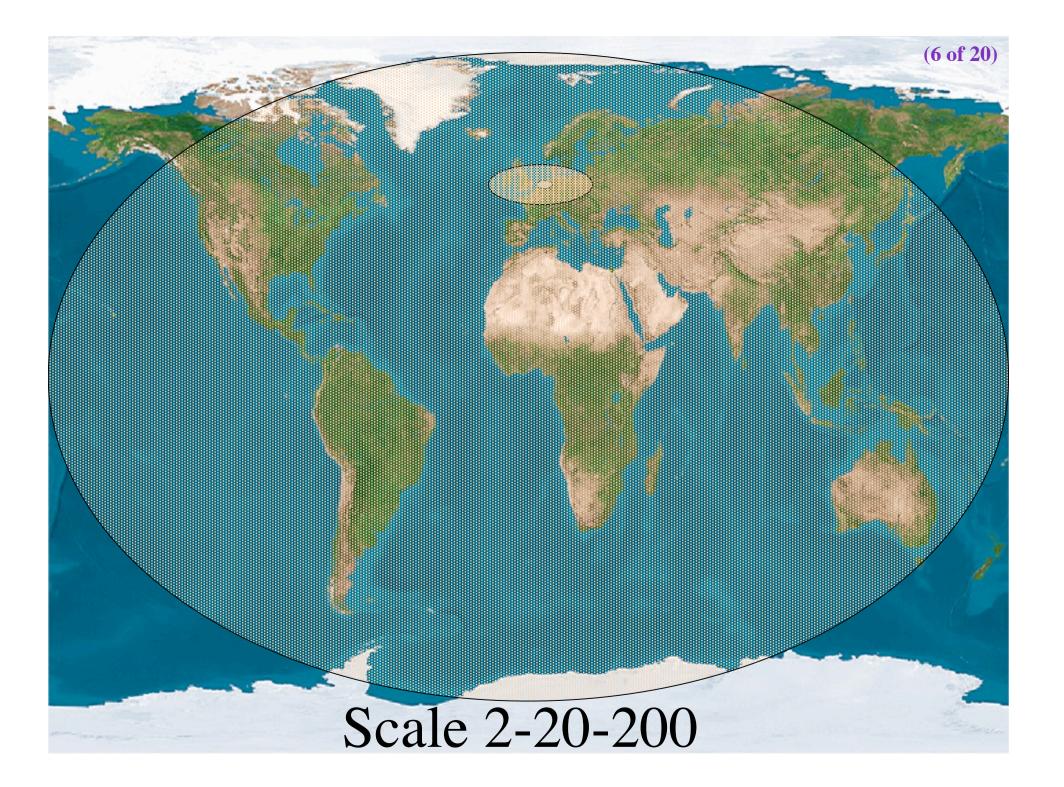B -> Need VPN services on/and full Internet routing, several to several

C -> Need very fat pipes, limited multiple Virtual Organizations, few to few

# So what are the facts

- **Costs of fat pipes (fibers) are one/third of equipment to light them up**
  - Is what Lambda salesmen tell me

- **Costs of optical equipment 10% of switching 10 % of full routing equipment for same throughput**
  - 100 Byte packet @ 10 Gb/s -> 80 ns to look up in 100 Mbyte routing table (light speed from me to you on the back row!)

- **Big sciences need fat pipes**

- **Bottom line: create a hybrid architecture which serves all users in one consistent cost effective way**

Scale 2-20-200

# The only formula's

$$\# \lambda \approx \frac{200 * e^{(t-2002)}}{rtt}$$

**Now, as having been a High Energy Physicist we set**

**c = 1**

**e = 1**

**ℏ = 1**

**and the formula reduces to:**

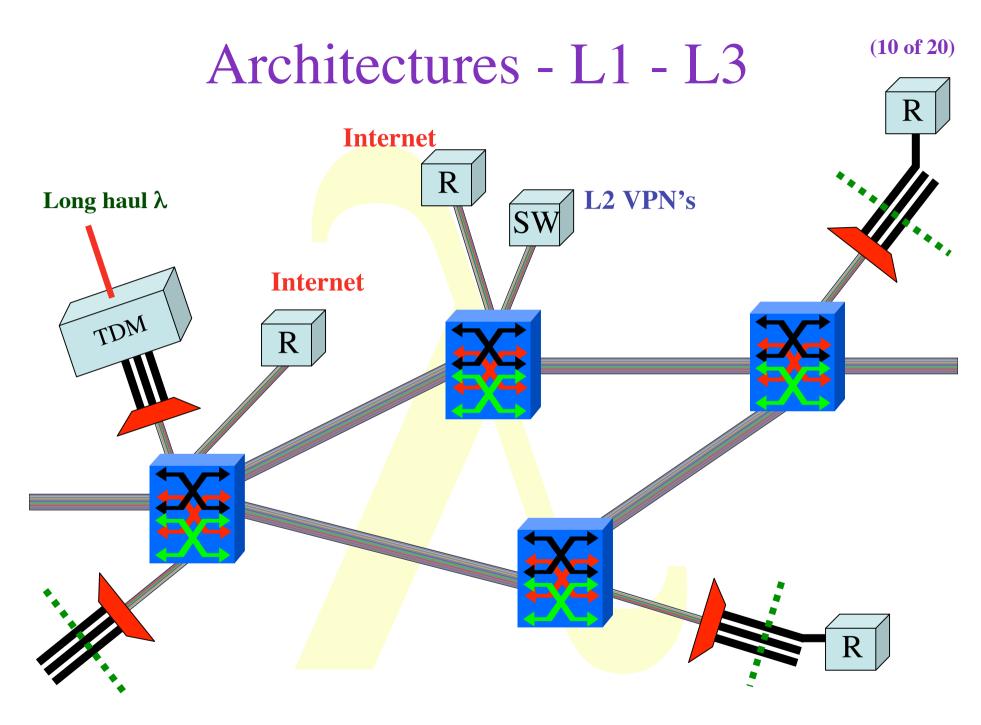$$\# \lambda \approx \frac{200 * e^{(t-2002)}}{rtt}$$

# Services

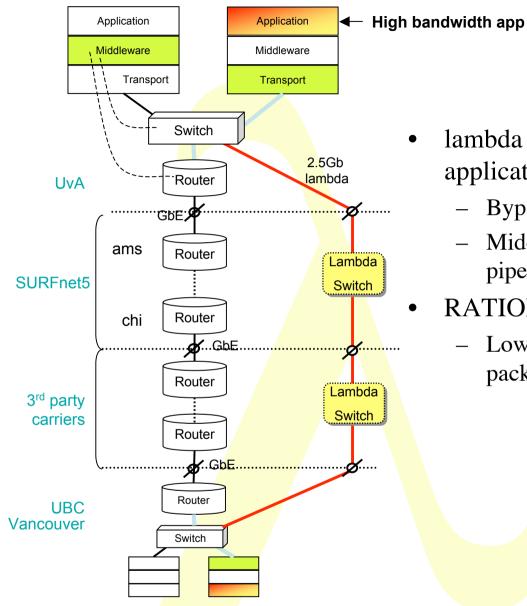|  | 2<br>Metro | 20<br>National/<br>regional | 200<br>World |
|---|---|---|---|
| A | Switching/<br>routing | Routing | ROUTER$ |
| B | VPN's,<br>(G)MPLS | VPN's<br>Routing | Routing |
| C<br>$\# \lambda \approx \dfrac{200 * e^{(t-2002)}}{rtt}$ | dark fiber<br>Optical<br>switching | Lambda<br>switching | Sub-<br>lambdas,<br>ethernet-<br>sdh |

# Current technology + (re)definition

- Current (to me) available technology consists of SONET/SDH switches

- Changing very soon!, optical switch on the way!

- DWDM+switching coming up

- Starlight uses for the time being VLAN's on Ethernet switches to connect [exactly two] ports (but also routing)

- We want to understand routerless limited environments

- So redefine a $\lambda$ as:

  **"a $\lambda$ is a pipe where you can inspect packets as they enter and when they exit, but principally not when in transit. In transit one only deals with the parameters of the pipe: number, color, bandwidth"**

# Architectures - L1 - L3

Long haul λ

Internet

Internet

L2 VPN's

TDM

R

R

SW

R

R

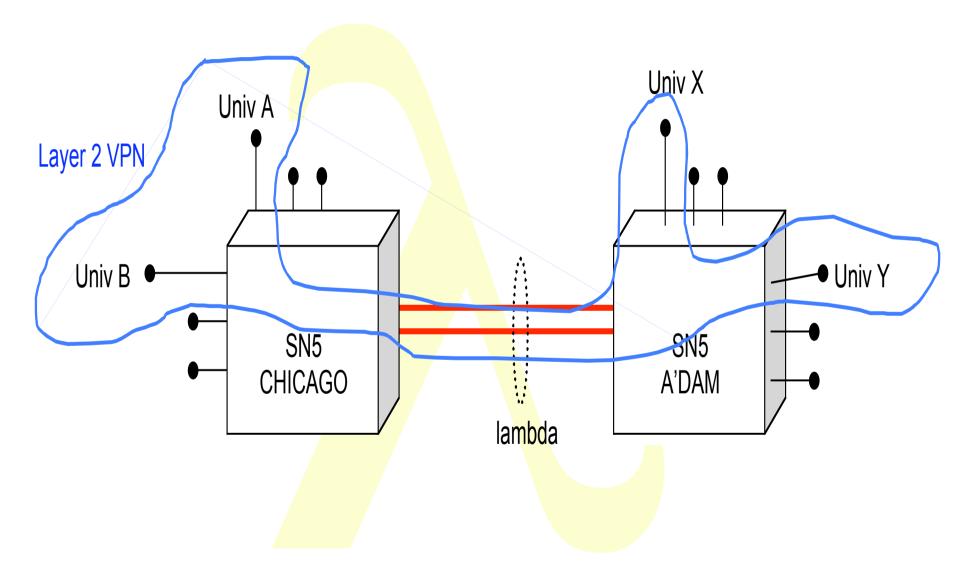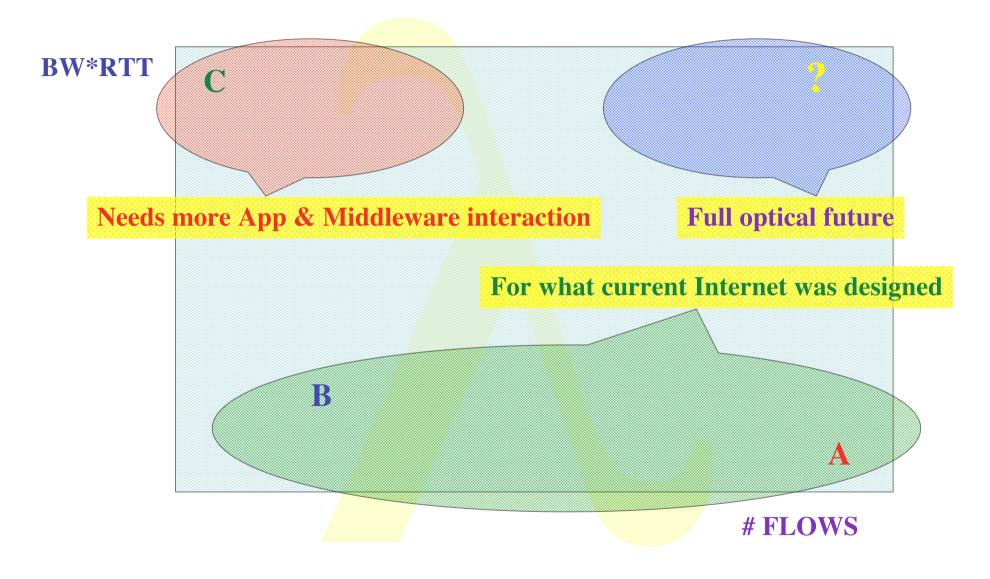**Bring plumbing to the users, not just create sinks in the middle of nowhere**

- lambda for high bandwidth applications
  - Bypass of production network
  - Middleware may request (optical) pipe
- RATIONALE:
  - Lower the cost of transport per packet

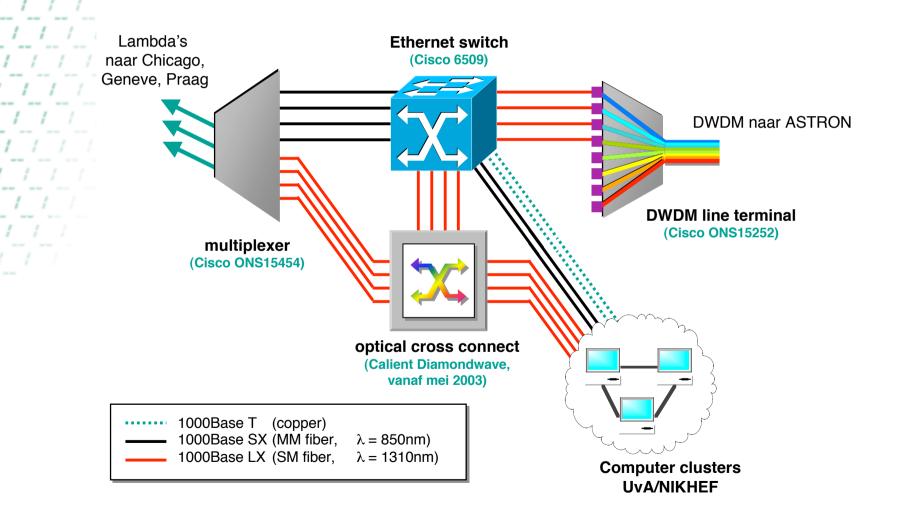# Distributed L2

# Transport in the corners

# NetherLight Amsterdam setup

Lambda's
naar Chicago,
Geneve, Praag

**Ethernet switch**
**(Cisco 6509)**

DWDM naar ASTRON

**DWDM line terminal**
**(Cisco ONS15252)**

**multiplexer**
**(Cisco ONS15454)**

**optical cross connect**
**(Calient Diamondwave,**
**vanaf mei 2003)**

······· 1000Base T    (copper)
——— 1000Base SX (MM fiber,    $\lambda$ = 850nm)
——— 1000Base LX  (SM fiber,    $\lambda$ = 1310nm)

**Computer clusters**
**UvA/NIKHEF**

15

SURFnet

6509

DAS: 32*2 CPU's
Myrinet

1 Gbs    100Mbs

FORCE10

switch

10 Gbs

15454

calient

server

Lambda's
to Chicago,
Geneve, Praag

Fat pc

4 HP servers

Dark fiber
To
Dwingeloo

# GigaPort
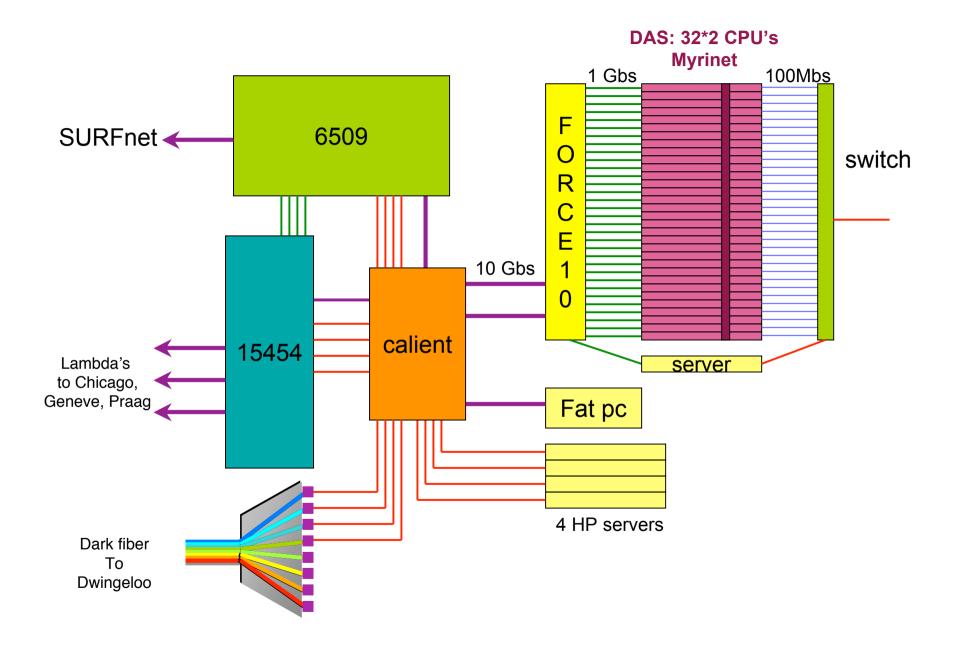
**SURFnet and GigaPort Next Generation**
*Creating the innovation engine*

Kees Neggers

Managing Director SURFnet

Praha, 20 February 2003

SURFnet

**GigaPort**

- **Provides the Dutch National Research Network**

- **Not for profit company**

- **200 connected organisations, 500.000 users**

- **Turnover (2002): 35M€**

- **Infrastructure services:**
  - innovation paid for by government
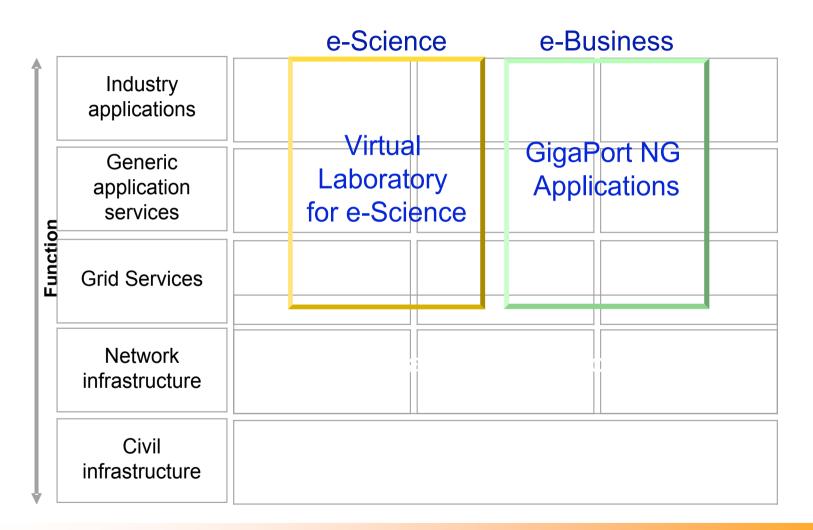  - cost effective exploitation for higher education and research

**SURFnet**

# Effects of GigaPort

- **World leading research infrastructure in NL**
  - **15 PoPs connected by thirty 10Gbps lambdas**
  - **Dual stack IPv4 and IPv6**

- **Helps transition in Telecoms market**
  - **GigaMAN**
  - **Fiber to the dormitories**
  - **Access pilots/ mobility**

- **Advanced Optical Exchange: NetherLight**
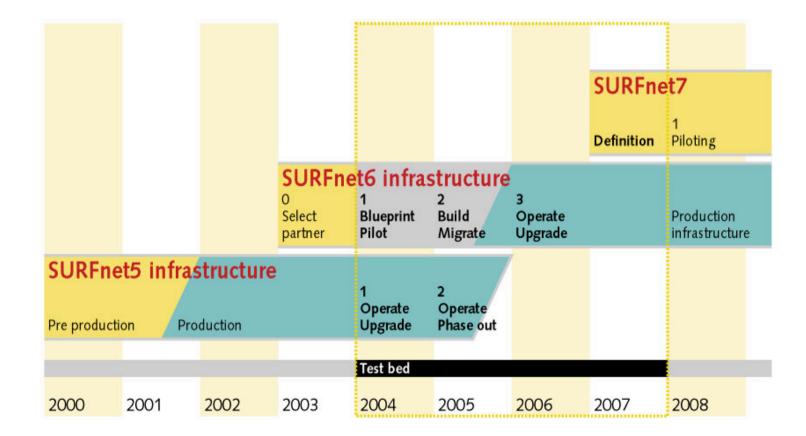
- **Playground for new applications**

# Scope GigaPort Next Generation

**GigaPort**

e-Science          e-Business

| | | | |
|---|---|---|---|
| Industry applications | | | |
| Generic application services | Virtual Laboratory for e-Science | GigaPort NG Applications | |
| Grid Services | | | |
| Network infrastructure | | | |
| Civil infrastructure | | | |

**Function** (vertical axis label, left side)

**SURFnet**

- **optical transmission and switching**

- **integrating light paths in network**

- **routing: new internet features and protocols**

- **monitoring & network management**

- **Testing methodology**

- **network access management (roaming, security, usability, personalized service provisioning)**
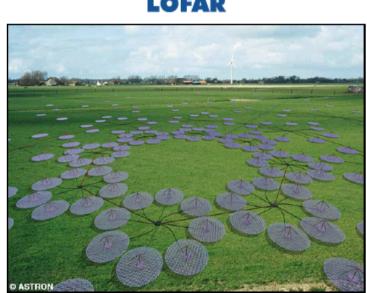
- **service grids**

**GigaPort**

LOFAR

- **Grid computing**

- **Data mining**

- **Data visualization**

- **Virtual reality**

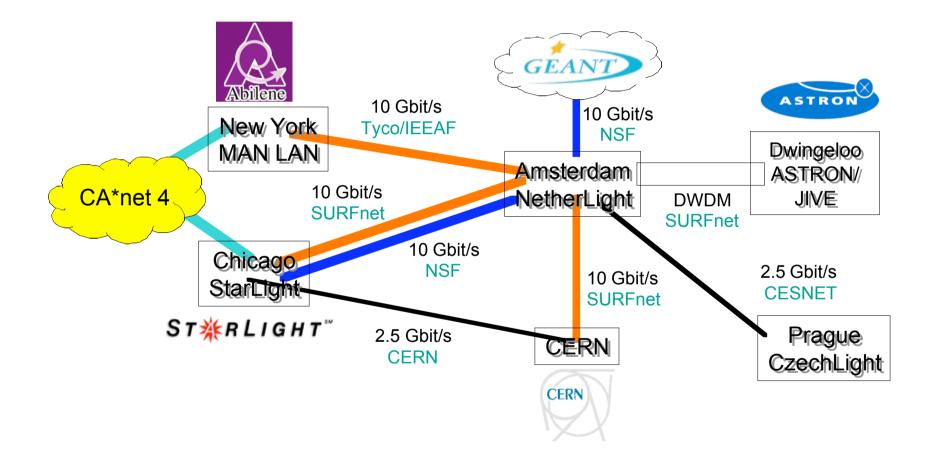- **Remote cooperation**

- **Radioastronomy**



© ASTRON

**Telecommunication infrastructures become part of scientific instruments**

SURF net

# Emerging international lambda grid

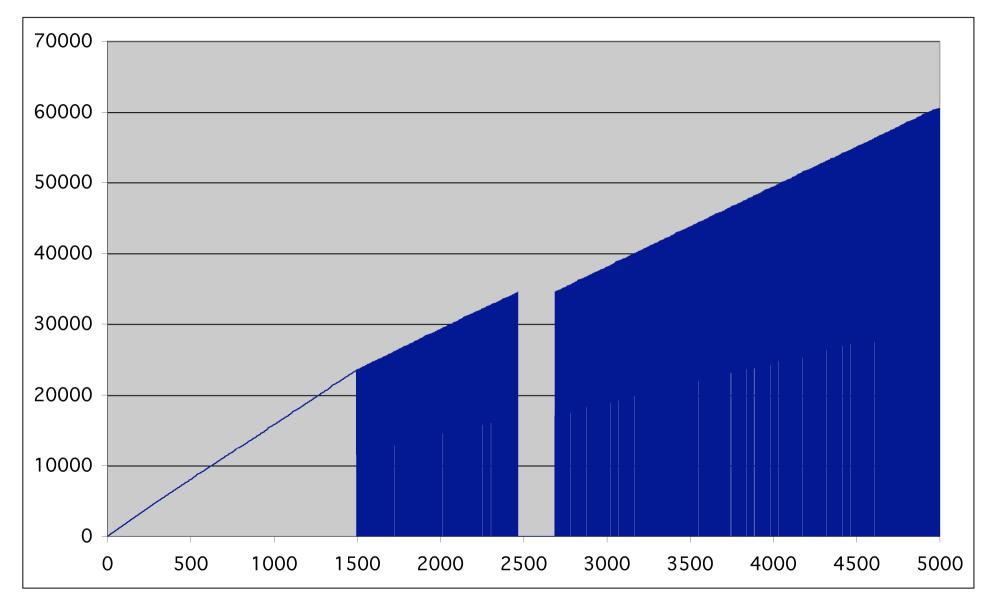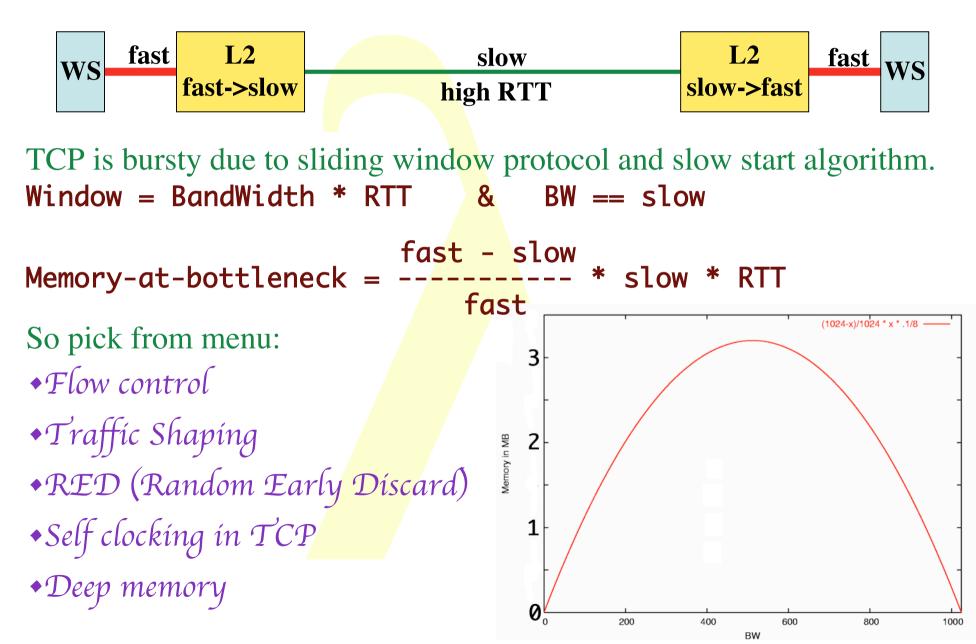# Early Lambda/LightPath usage experiences

EVL => WCW ⎯⎯⋄⎯⎯

Sum Throughput [Mbit/s]



# Streams

Sun TCP Window [Mbytes]

# 5000  1 kByte UDP packets

# Layer - 2 requirements from 3/4

TCP is bursty due to sliding window protocol and slow start algorithm.

```
Window = BandWidth * RTT    &    BW == slow


                           fast - slow
Memory-at-bottleneck = ------------- * slow * RTT
                              fast
```

So pick from menu:

- *Flow control*
- *Traffic Shaping*
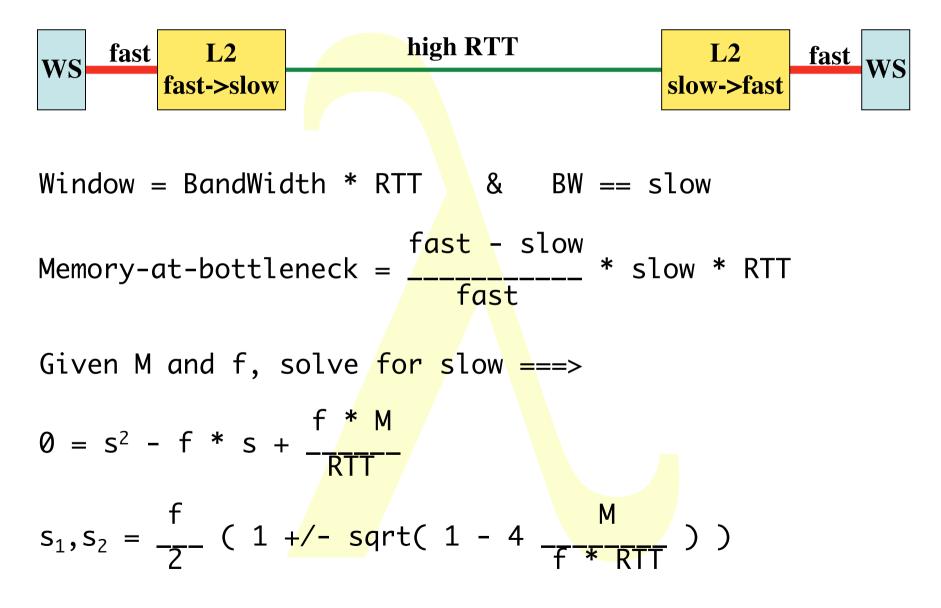- *RED (Random Early Discard)*
- *Self clocking in TCP*
- *Deep memory*

# Self-clocking of TCP

# Layer - 2 requirements from 3/4



```
Window = BandWidth * RTT    &    BW == slow
```

$$\text{Memory-at-bottleneck} = \frac{\text{fast} - \text{slow}}{\text{fast}} * \text{slow} * \text{RTT}$$

```
Given M and f, solve for slow ===>
```

$$0 = s^2 - f * s + \frac{f * M}{RTT}$$

$$s_1, s_2 = \frac{f}{2} \left( 1 +/- \sqrt{1 - 4 \frac{M}{f * RTT}} \right)$$

Forbidden area, solutions for s when f = 1 Gb/s, M = 0.5 Mbyte (20 of 25) AND NOT USING FLOWCONTROL

layers of increasing abstraction taxonomy

# Starting point



**Generic AAA server Rule based engine**

**Application Specific Module**

**Service**

**Accounting Metering**

**Policy**

**Data**

**Policy**

**Data**

**Acct Data**

API

PDP

PEP

1     1     2     3     4     5     5     3     4'     3

**RFC 2903 - 2906 , 3334 , policy draft**

# Multi domain case

# iGrid2002

- www.igrid2002.org
- 25 demonstrations
- 16 countries (at least)
- Level3, Tyco, IEEAF Lambda's
- CISCO, Hp equipment sponsoring
- Shipping nightmare, debugging literally
- ~30 Gbit/s International connectivity
- Huge networking collaboration
- Smelly NOC in the iGrid preparation weekend

# Lessons learned

- **Most applications could not cope with the network!!!**

- **No bottleneck whatsoever in the network**

- **Many got about 50 - 100 mbit/s singlestream tcp**

- **On Sunday evening my laptop had the highest single stream to Chicago (~ 340 Mbit/s)**

- **NIC's, Linux implementation and timing problem**

- **Gridftp severely hit**

- **~ 22 papers to be published**

# The END

Thanks to

SURFnet: Kees Neggers

UIC&iCAIR: Tom DeFanti, Joel Mambretti

CANARIE: Bill St. Arnaud

This work is supported by:

SURFnet

EU-IST project DATATAG