

Lambda-Grid developments

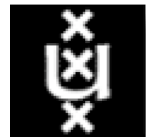
www.science.uva.nl/~deLaat

Cees de Laat

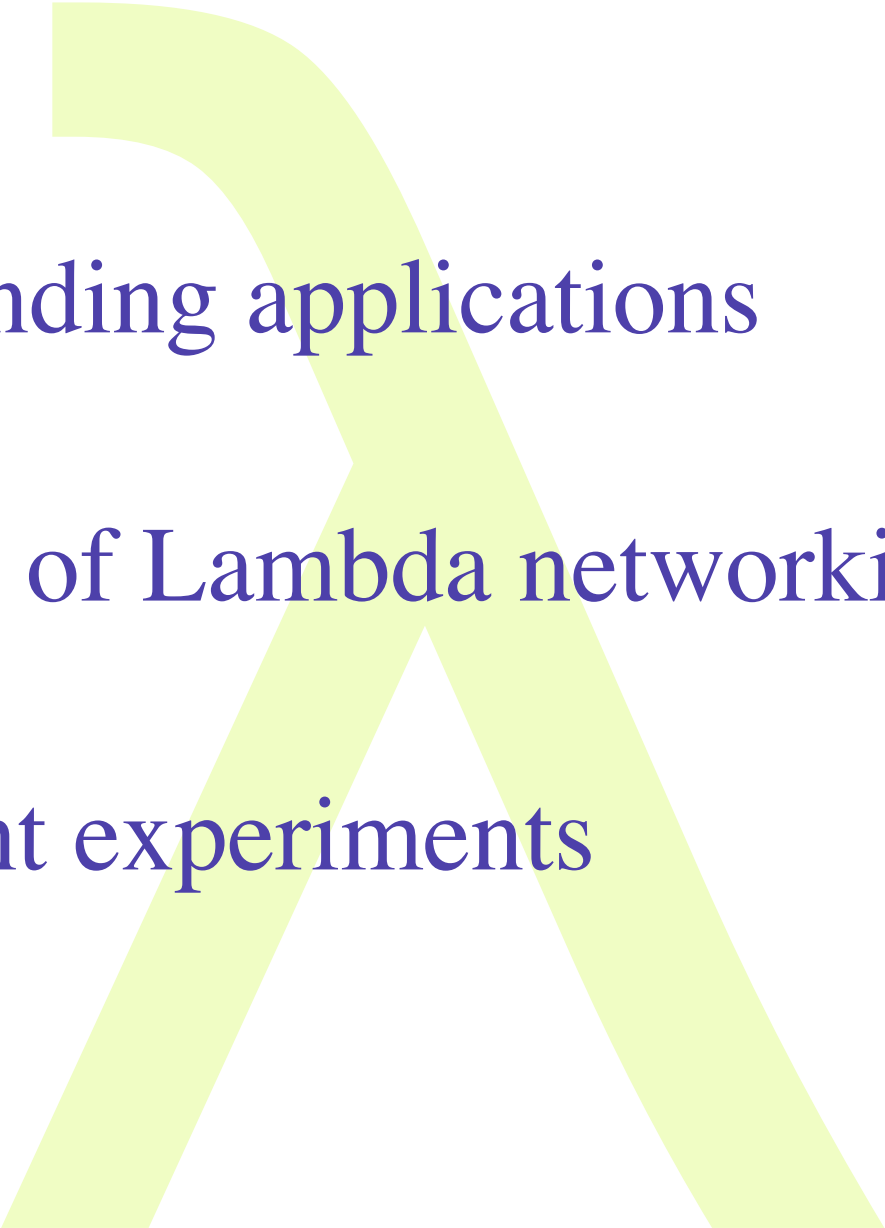
GigaPort
EU

University of Amsterdam

SARA
NCF

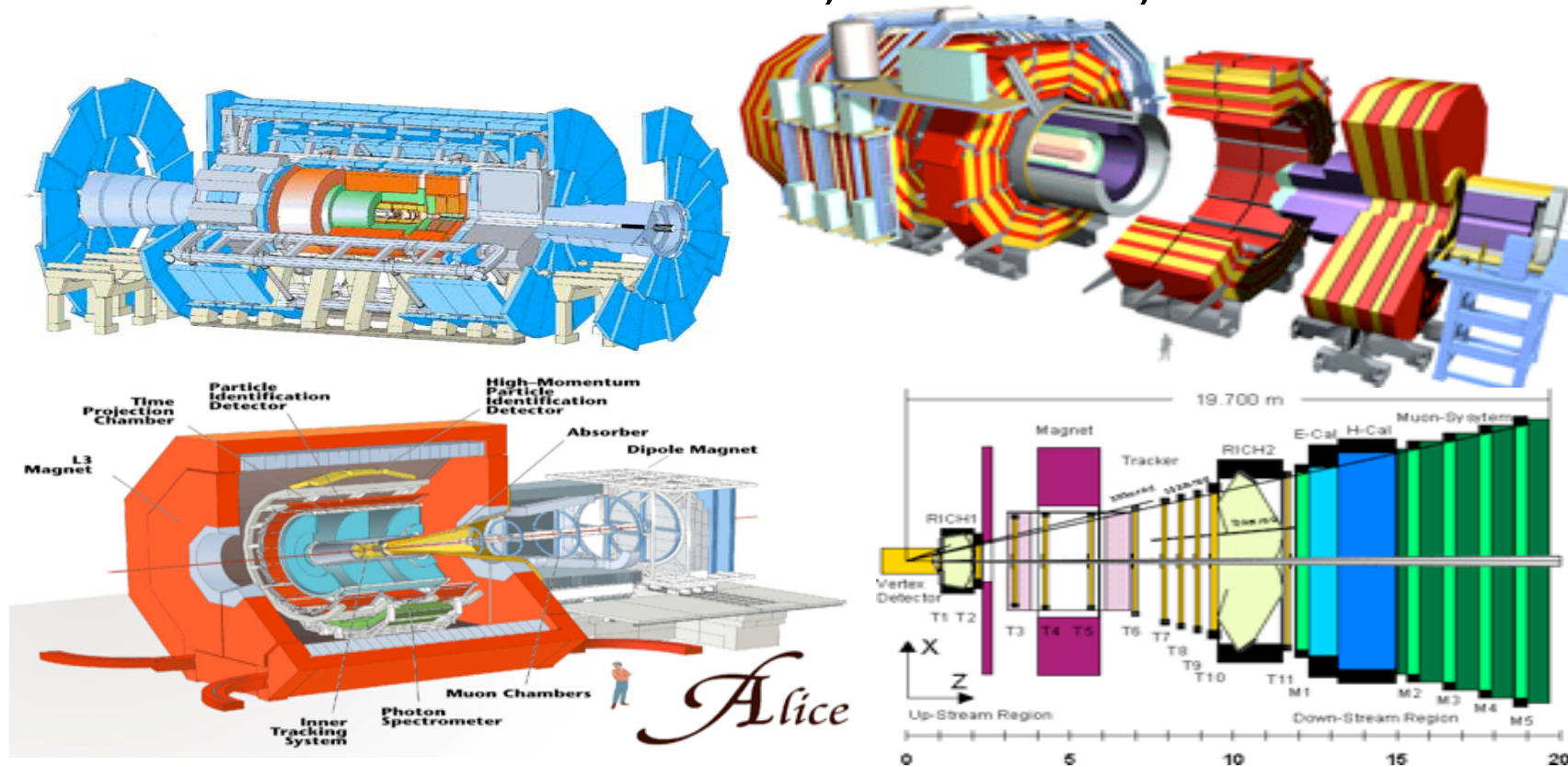


Contents of this talk

- 
- Demanding applications
 - Model of Lambda networking
 - Current experiments

Four LHC Experiments: The Petabyte to Exabyte Challenge

- **ATLAS, CMS, ALICE, LHCb**



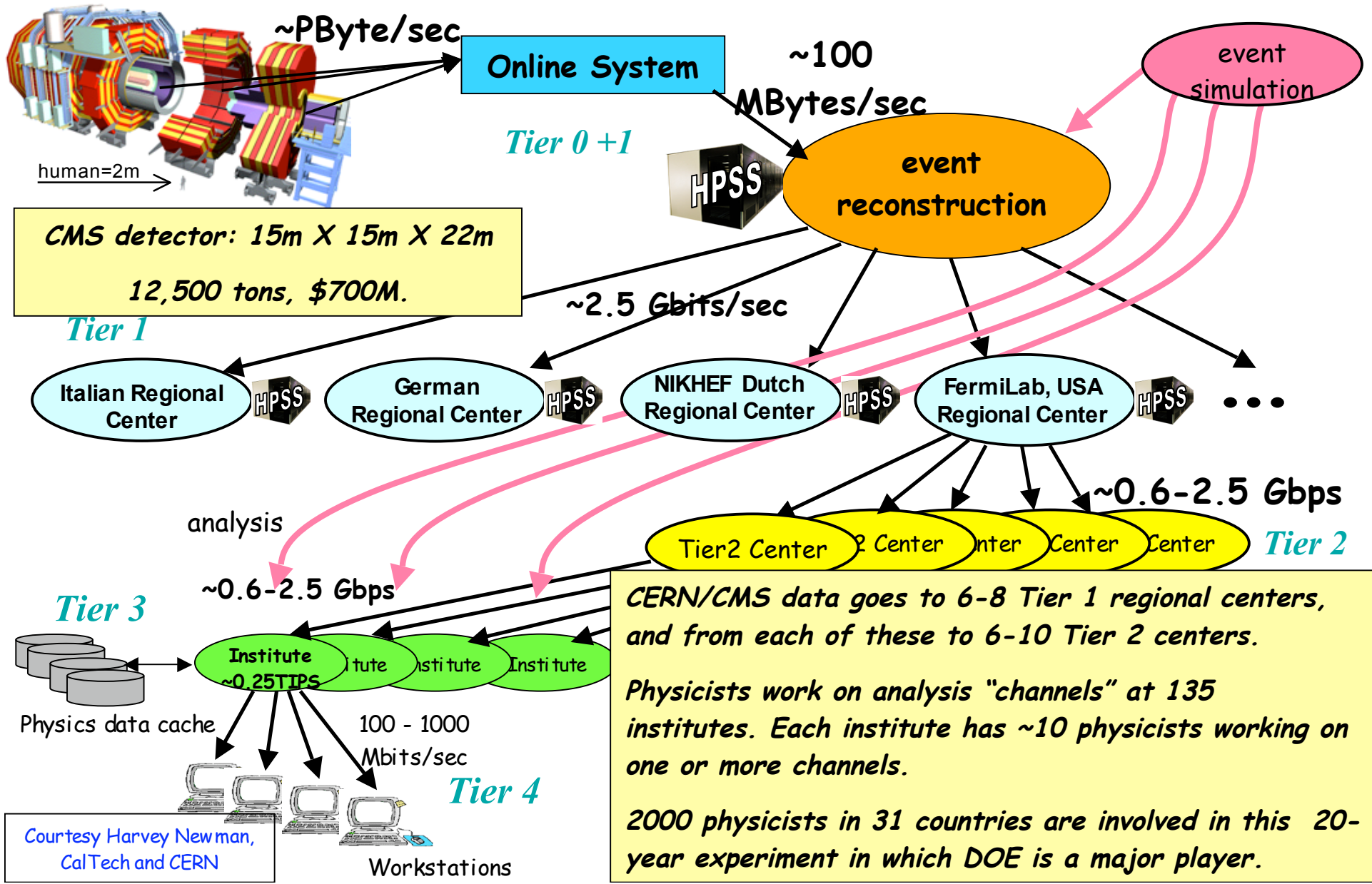
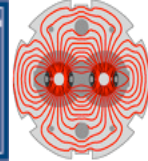
6000+ Physicists & Engineers; 60+ Countries; 250 Institutions

Tens of PB 2008; To 1 EB by ~2015
Hundreds of TFlops To PetaFlops



LHC Data Grid Hierarchy

CMS as example, Atlas is similar



Courtesy Harvey Newman, CalTech and CERN

VLBI

VLBI is easily capable of generating many Gb of data per

The sensitivity of the VLBI array scales with

(data-rate) and there is a strong push to

Rates of 8Gb/s or more are entirely feasible

development. It is expected that parallel

correlator will remain the most efficient approach

s distributed processing may have an application

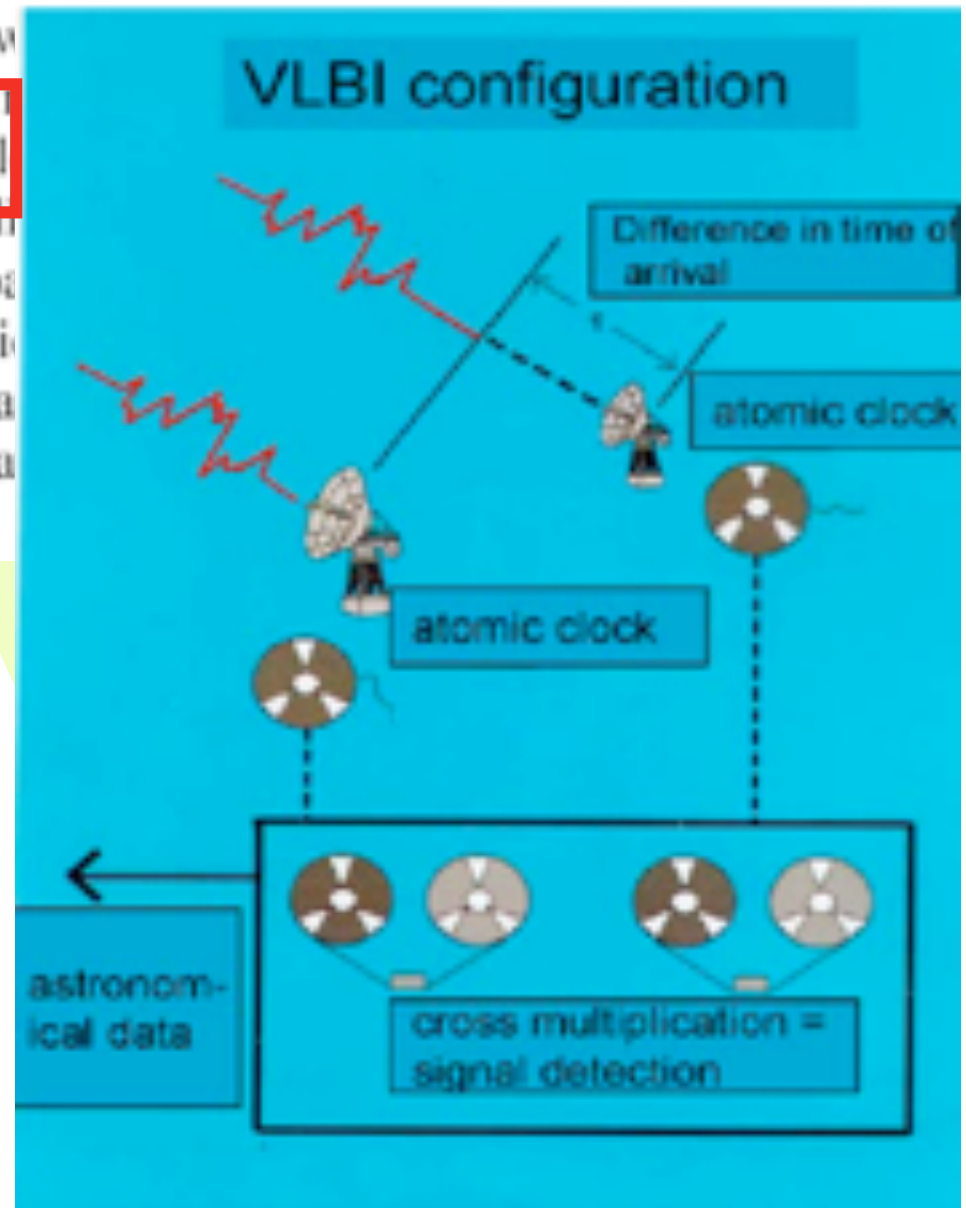
ulti-gigabit data streams will aggregate into larger

or and the capacity of the final link to the data

center.



Westerbork Synthesis Radio Telescope - Netherlands



Lambdas as part of instruments

GigaPort



www.lofar.org

SURF/net
/

OptIPuter Project Goal: Scaling to 100 Million Pixels

- **JuxtaView (UIC EVL) for PerspecTile LCD Wall**
 - Digital Montage Viewer
 - 8000x3600 Pixel Resolution~30M Pixels
- **Display Is Powered By**
 - 16 PCs with Graphics Cards
 - 2 Gigabit Networking per PC



Grids

Showned you:

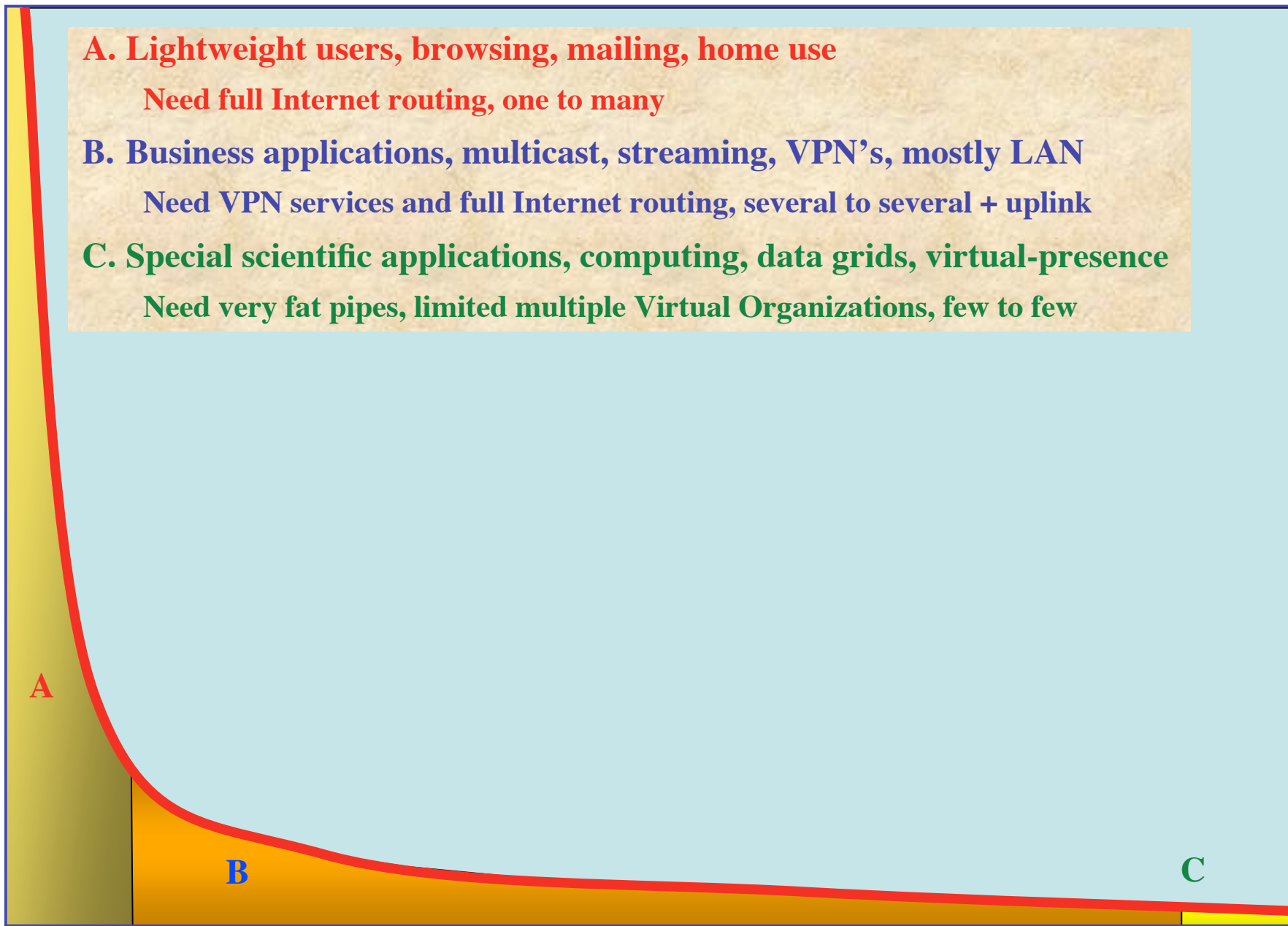
- **Computational Grids**
 - HEP and LOFAR analysis requires massive CPU capacity
- **Data Grid**
 - Storing and moving HEP, Bio and Health data sets is major challenge
- **Instrumentation Grids**
 - Several massive data sources are coming online
- **Visualization Grids**
 - Data object (TByte sized) inspection, anywhere, anytime

Contents of this talk

- 
- Demanding applications
 - Model of Lambda networking
 - Current experiments

U
S
E
R
S

- A. Lightweight users, browsing, mailing, home use**
Need full Internet routing, one to many
- B. Business applications, multicast, streaming, VPN's, mostly LAN**
Need VPN services and full Internet routing, several to several + uplink
- C. Special scientific applications, computing, data grids, virtual-presence**
Need very fat pipes, limited multiple Virtual Organizations, few to few



ADSL

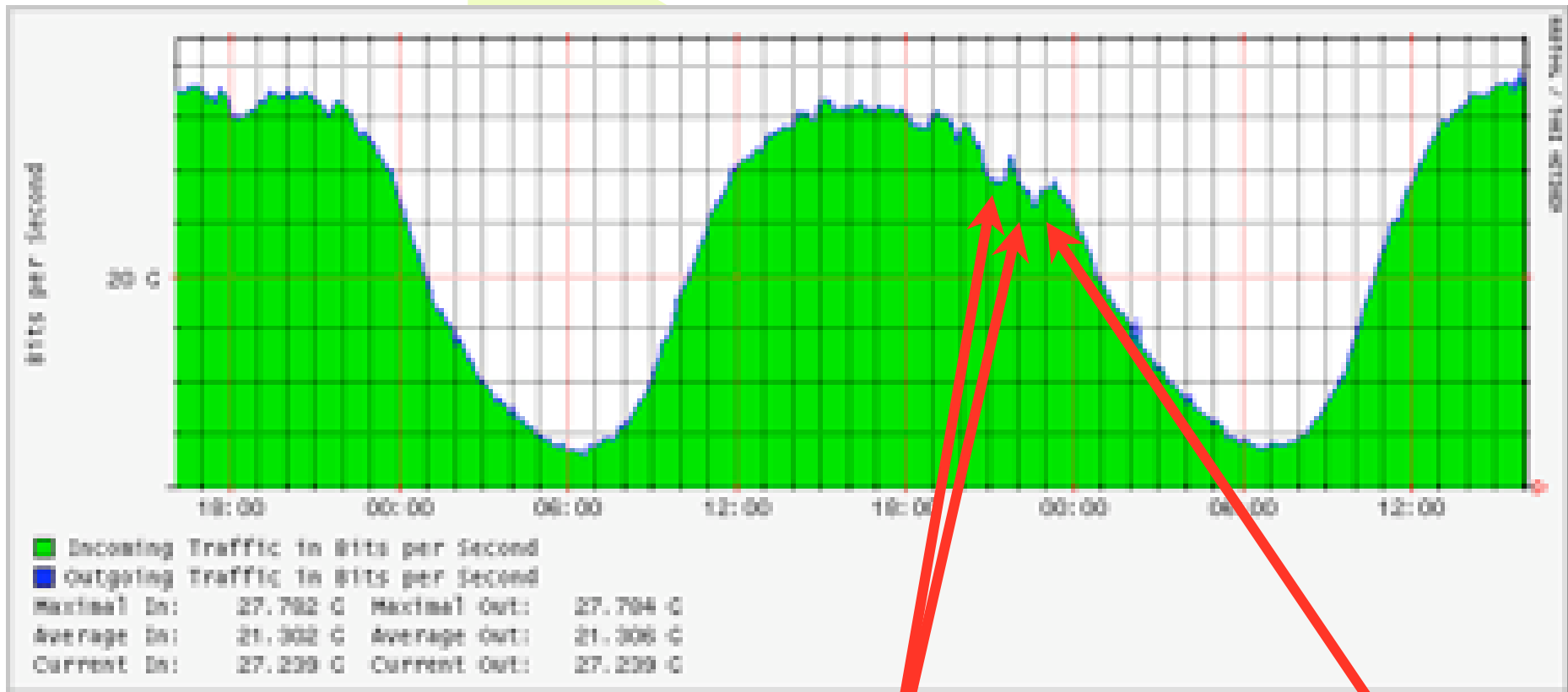
GigE

BW requirements

The Dutch Situation

- **Estimate A**
 - 17 M people, 6.4 M households, 25 % penetration of 0.5-2.0 Mb/s ADSL, 40 times under-provisioning ==> 20 Gb/s

AMS-IX



June 19th 2004

Lost :-)

European championship football **Holland -- Czech Republic**

The Dutch Situation

- **Estimate A**

- 17 M people, 6.4 M households, 25 % penetration of 0.5-2.0 Mb/s ADSL, 40 times under-provisioning ==> 20 Gb/s

- **Estimate B**

- SURFnet has 10 Gb/s to about 12 institutes and 0.1 to 1 Gb/s to 180 customers, estimate same for industry (overestimation) ==> 20-40 Gb/s

The Dutch Situation

- **Estimate A**

- 17 M people, 6.4 M households, 25 % penetration of 0.5-2.0 Mb/s ADSL, 40 times under-provisioning ==> 20 Gb/s

- **Estimate B**

- SURFnet has 10 Gb/s to about 12 institutes and 0.1 to 1 Gb/s to 180 customers, estimate same for industry (overestimation) ==> 20-40 Gb/s

- **Estimate C**

- Leading HEF and ASTRO + rest ==> 80-120 Gb/s
- LOFAR ==> \approx 26 Tbit/s

u
s
e
r
s

A. Lightweight users, browsing, mailing, home use

Need full Internet routing, one to many

B. Business applications, multicast, streaming, VPN's, mostly LAN

Need VPN services and full Internet routing, several to several + uplink

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

Need very fat pipes, limited multiple Virtual Organizations, few to few

$\Sigma C \gg 100 \text{ Gb/s}$

$\Sigma B \approx 40 \text{ Gb/s}$

$\Sigma A \approx 20 \text{ Gb/s}$

A

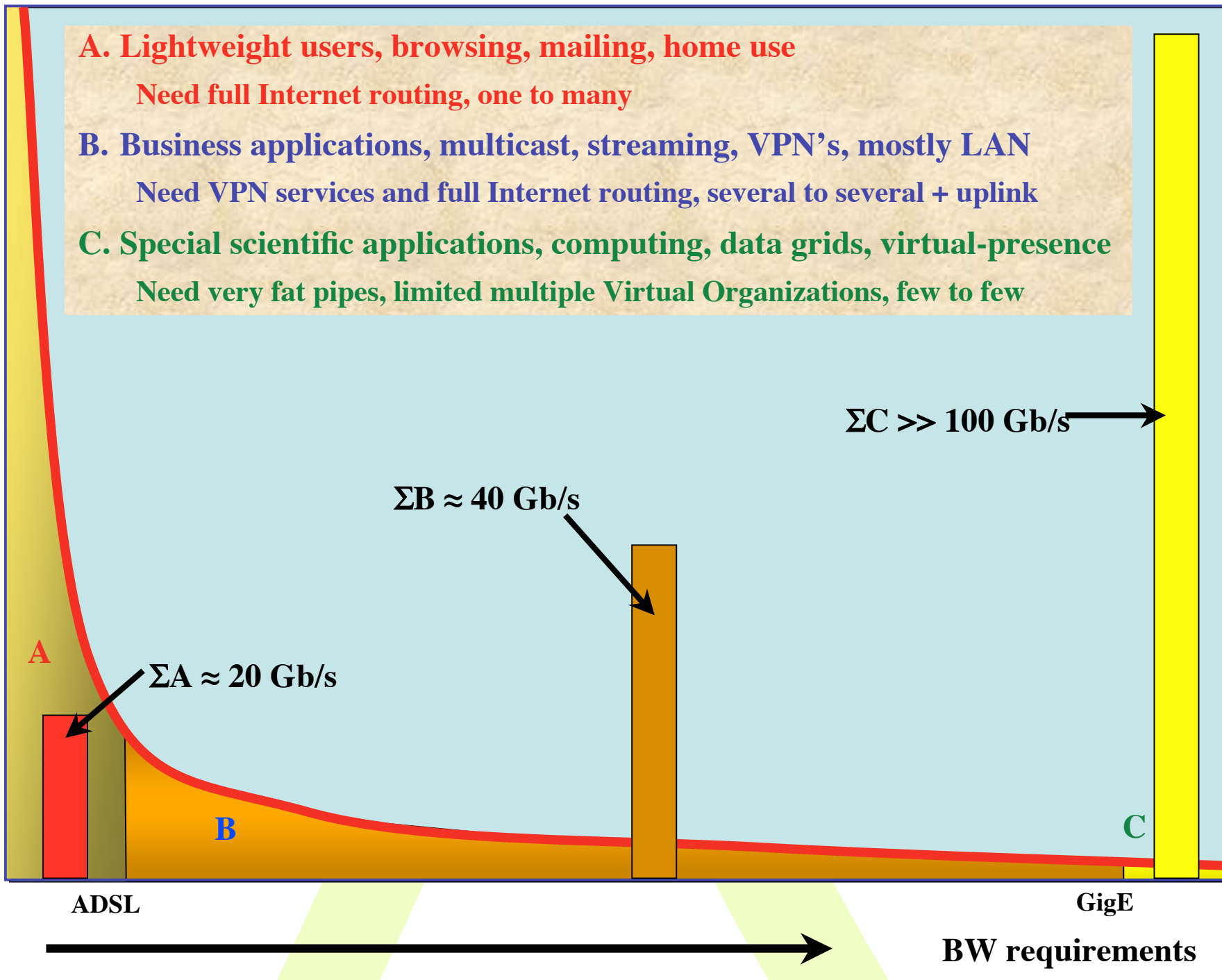
B

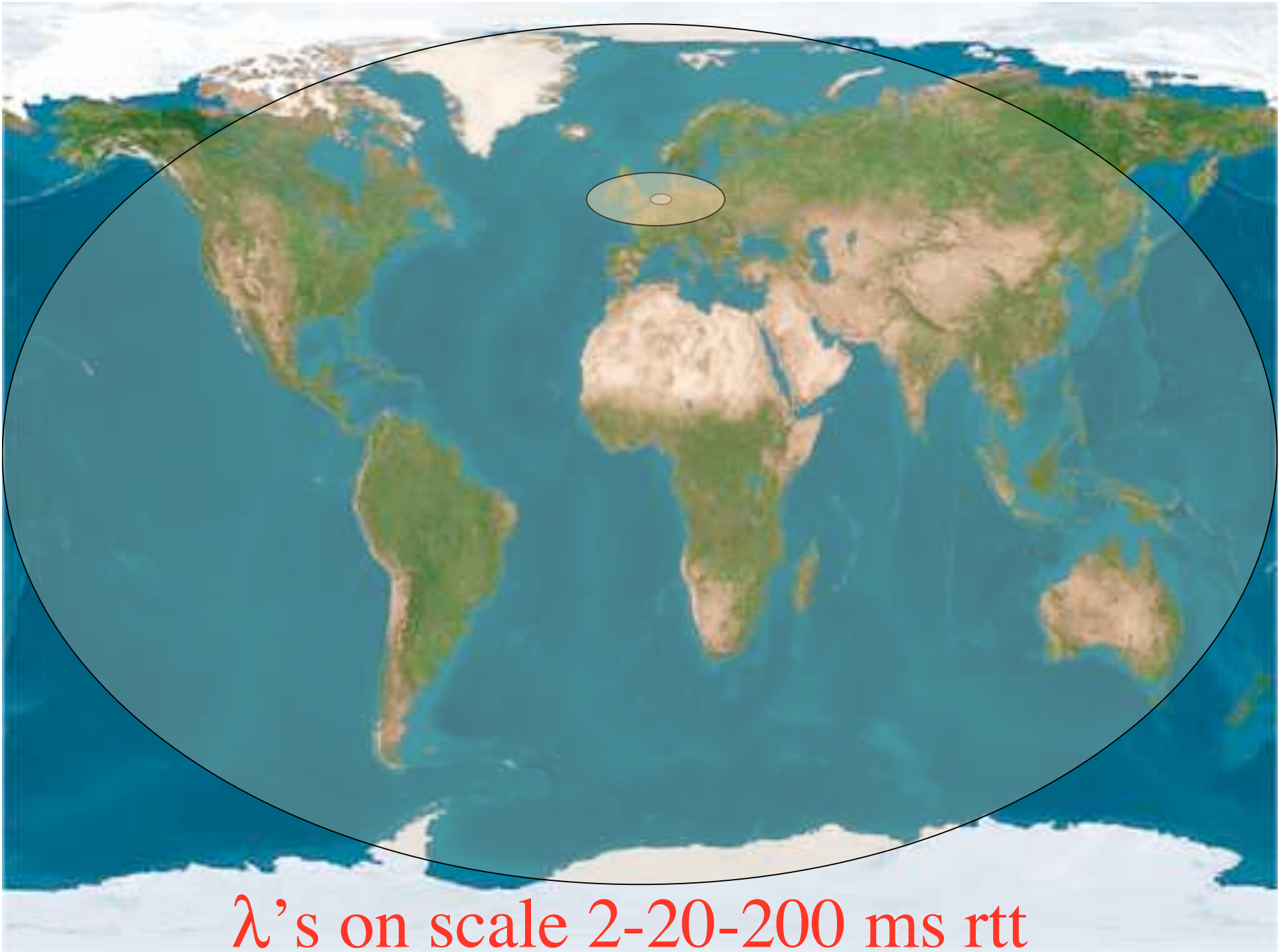
C

ADSL

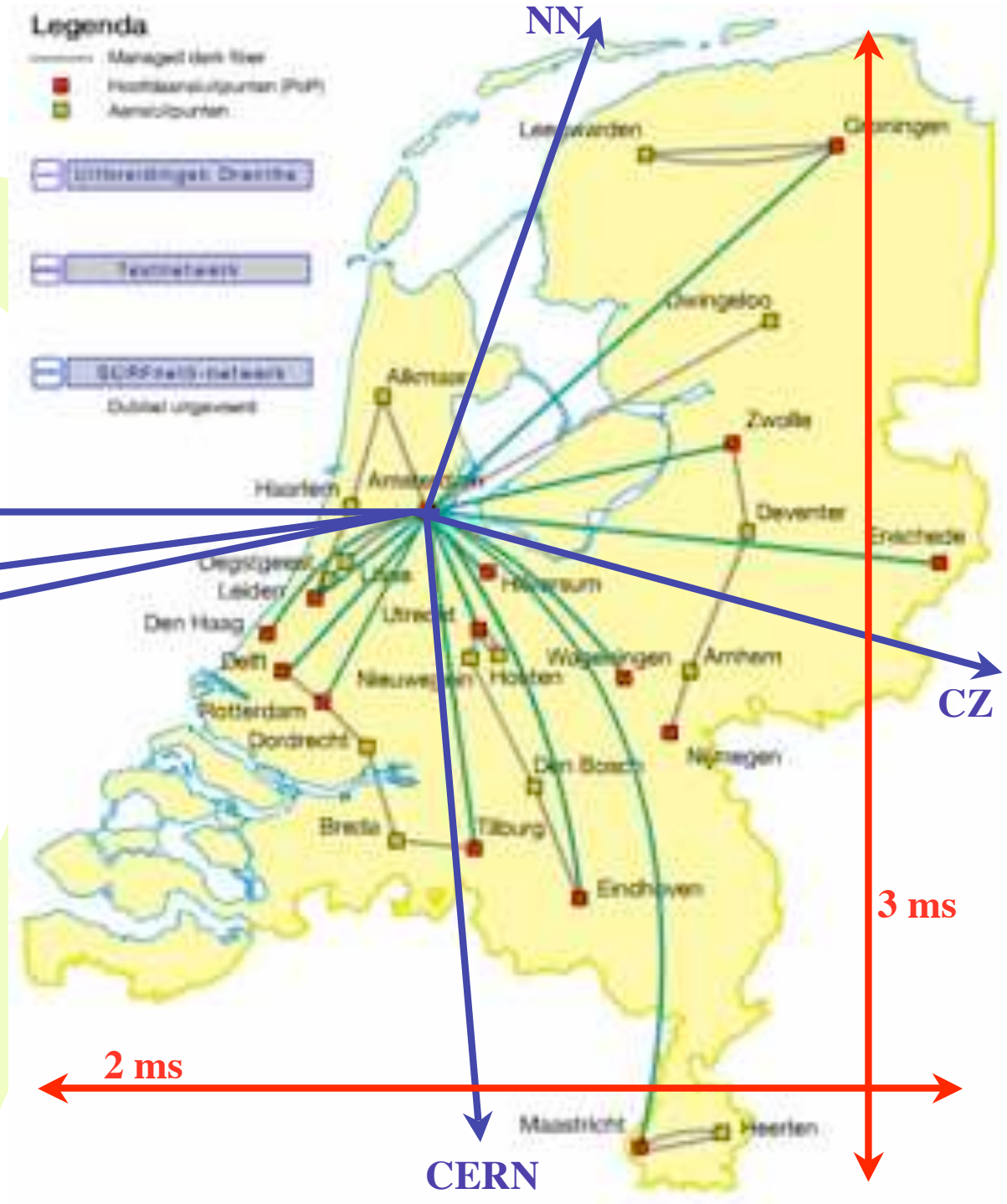
GigE

BW requirements





λ 's on scale 2-20-200 ms rtt



StarLight

NY

UK

CZ

CERN

3 ms

2 ms

SURFnet
fibers
(old pict by now)

The only formula

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

Compares very well with SURFnet's resources and Lambda's @ NetherLight

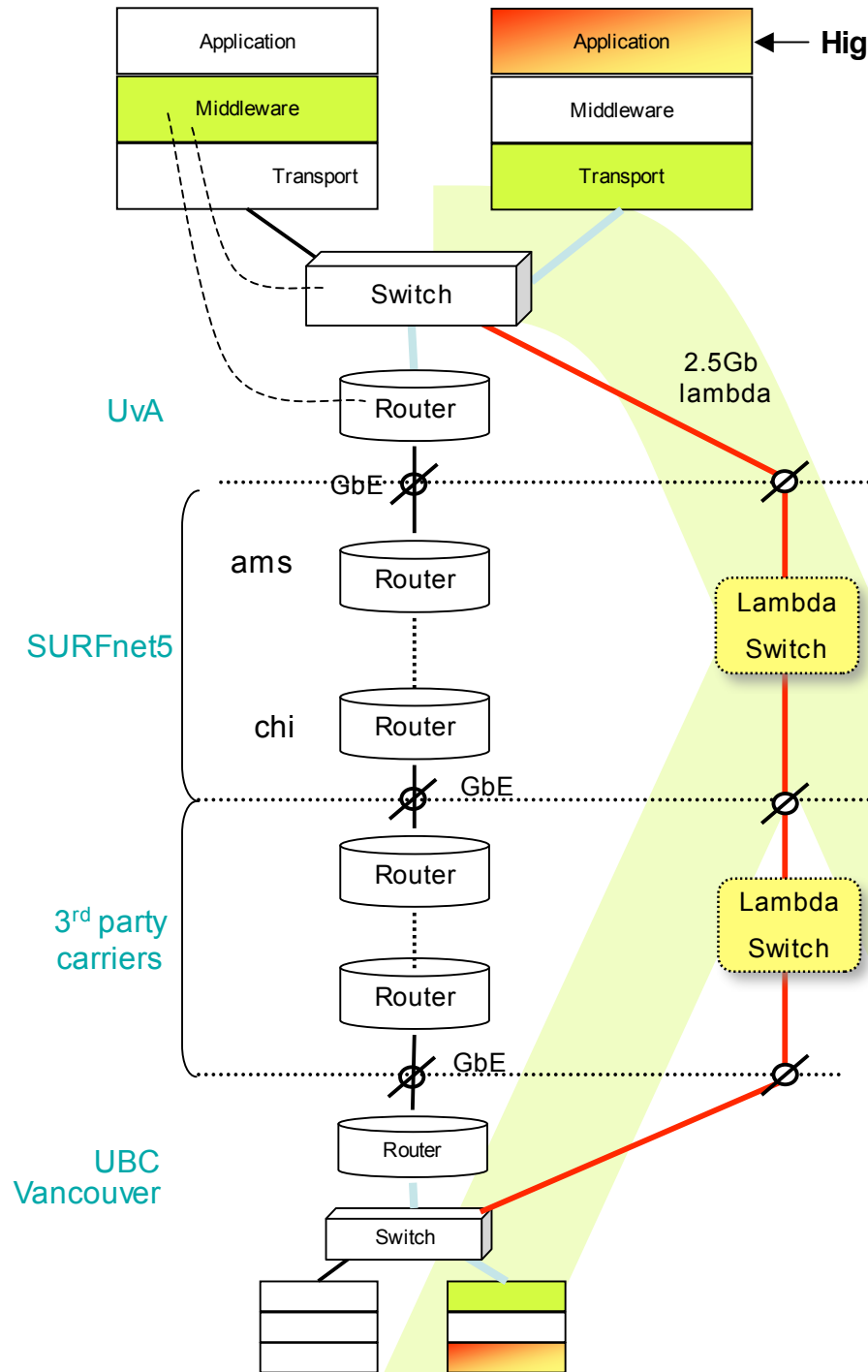
- 1 Transatlantic Lambda in 2002, now ~10 from EU+US
- 5300 km dark fiber in Holland \approx railway net

So what?

- **Costs of optical equipment 10% of switching 10 % of full routing equipment for same throughput**
 - 10G routerblade -> 100-300 k\$, 10G switch port -> 10-20 k\$, MEMS port -> 0.7 k\$
 - DWDM lasers for long reach expensive, 10-50k\$ (???)
 - 64 Byte packet @ 10 Gbit/s -> 52 ns -> time to look up destination in 140 kEntries routing table (light speed from me to you (15 meter)!)
- **Bottom line: look for a hybrid architecture which serves all classes in a cost effective way (A -> L3 , B -> L2 , C -> L1)**
- **Give each packet in the network the service it needs, but no more**
- **Look at worldwide ethernet infrastructure:**
 - Tested 10 Gbit/s Ethernet WANPHY Amsterdam-CERN
 - <http://www.surfnet.nl/en/publications/pressreleases/021003.html>
- **Look at worldwide lambda structure:**
 - <http://www.glif.is/>

UVA/EVL's
64*64
Optical Switch
@ NetherLight
in SURFnet POP @
SARA
Costs 1/100th of a
similar throughput
router
or 1/10th of an
Ethernet switch but
with specific services!

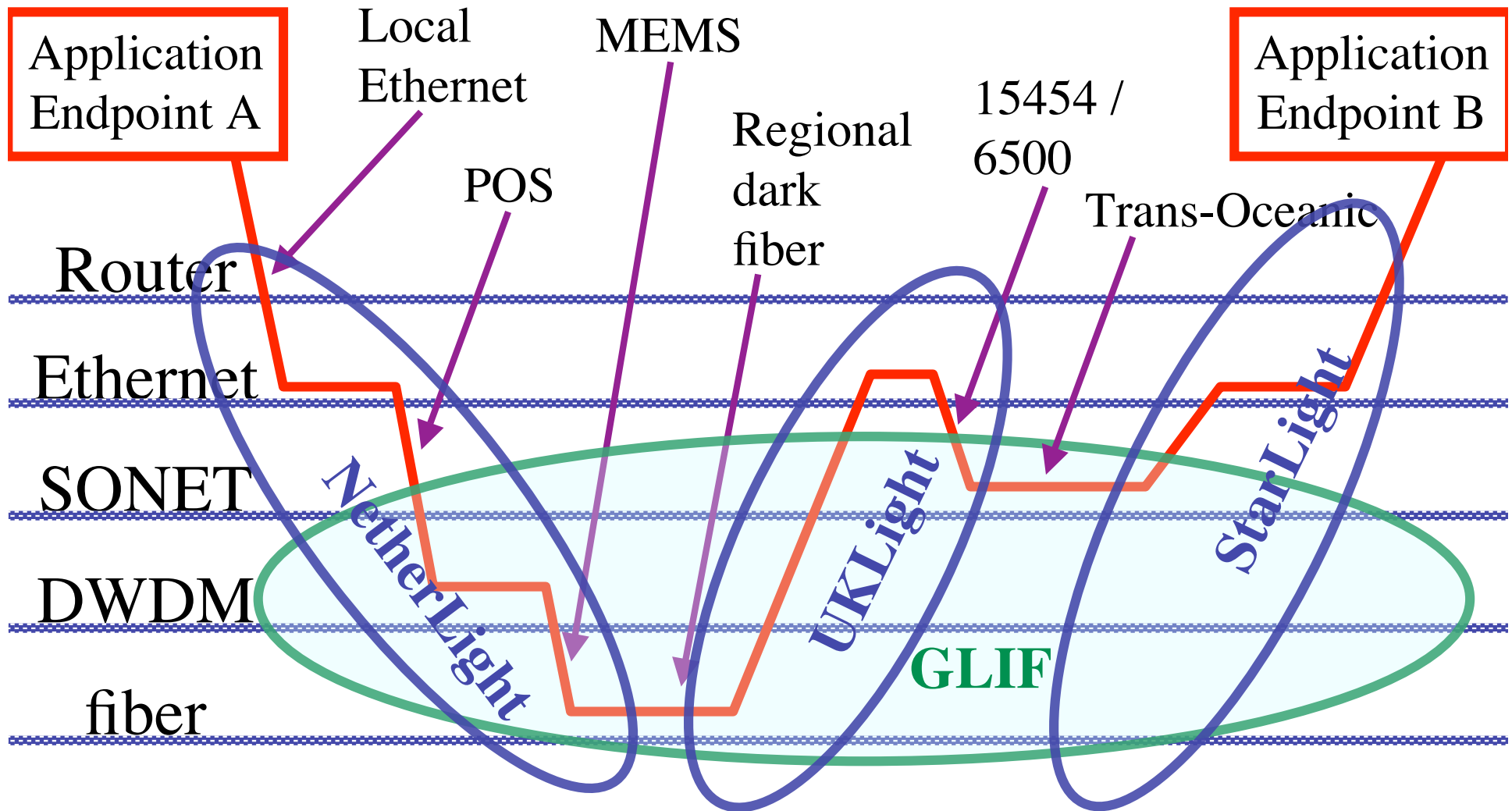




- lambda for high bandwidth applications
 - Bypass of production network
 - Middleware may request (optical) pipe
- RATIONALE:
 - Lower the cost of transport per packet
 - Use Internet as controlplane!

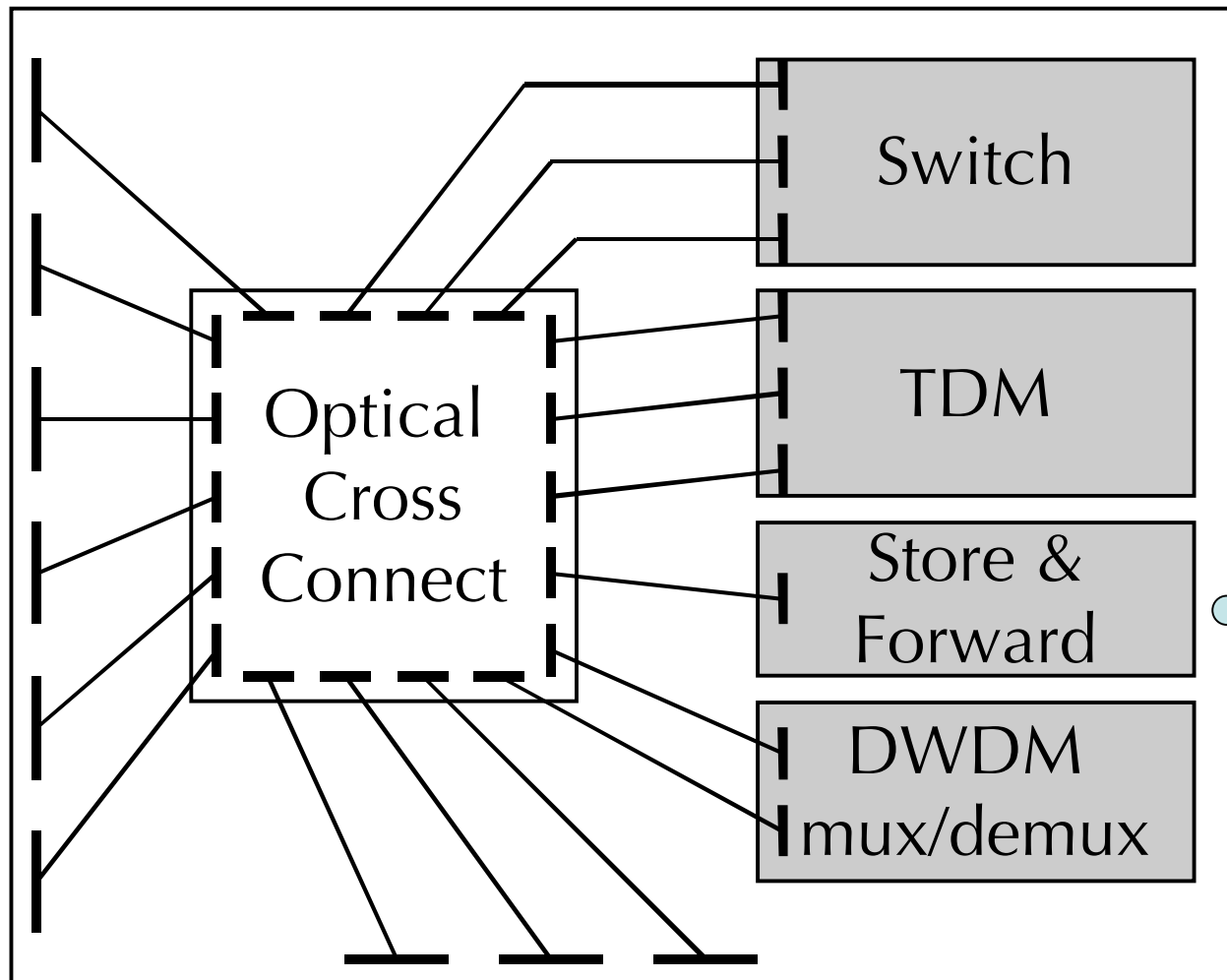


How low can you go?



Optical Exchange as Black Box

Optical Exchange



TeraByte
Email
Service

Contents of this talk

- 
- Demanding applications
 - Model of Lambda networking
 - Current experiments

GLIF: Global Lambda Integrated Facility

- Established at the 3rd Lambda Grid Workshop, August 2003 in Reykjavik, Iceland
- Collaborative initiative among worldwide NRENs, institutions and their users
- A world-scale Lambda-based Laboratory for application and middleware development

GLIF vision:

To build a new grid-computing paradigm, in which the central architectural element is optical networks, not computers, to support this decade's most demanding e-science applications.

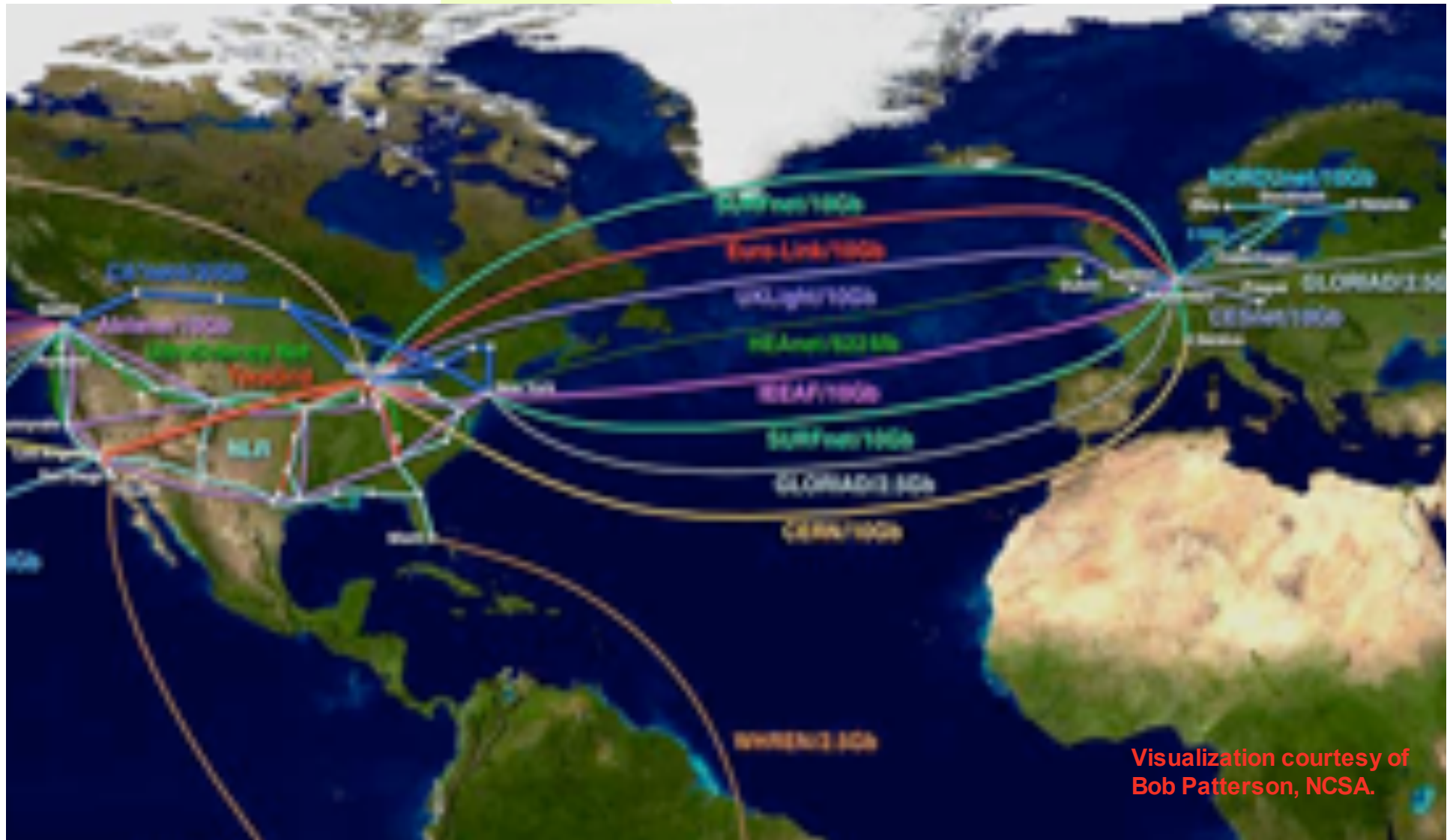


GLIF Q3 2004



Visualization courtesy of
Bob Patterson, NCSA.

GLIF Q3 2004



Little GLORIAD

<http://www.nsf.gov/od/lpa/news/03/pr03151.htm>



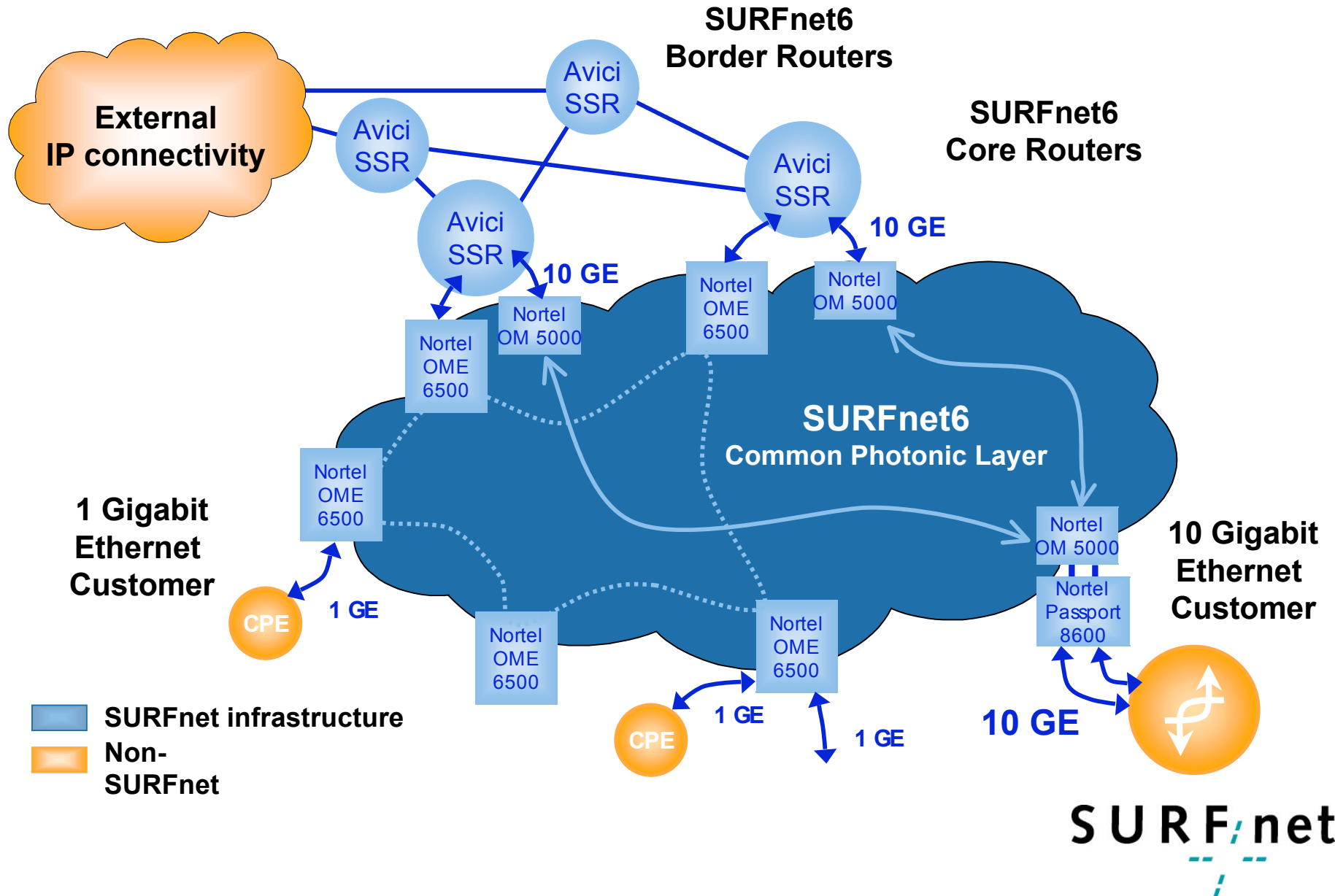
T. Schindler / National Science Foundation

SURFnet6 on dark fiber

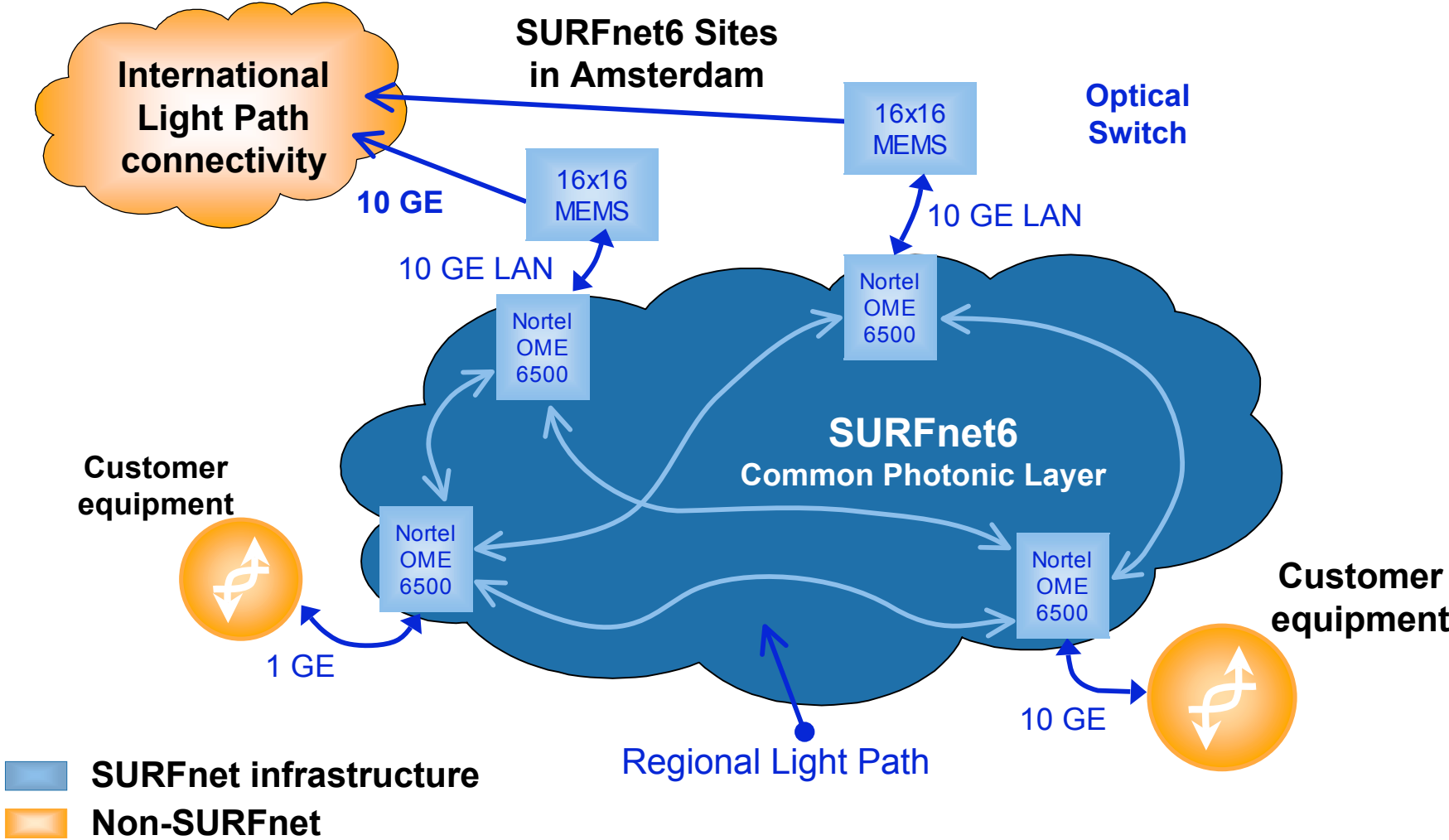


- SURFnet6 will be entirely based on own dark fiber
- Over 5300 km fiber pairs available today; average price paid for 15 year IRUs: < 6 EUR/meter per pair
- Managed dark fiber infrastructure will be extended with new routes, to be ready for SURFnet6

IP network implementation



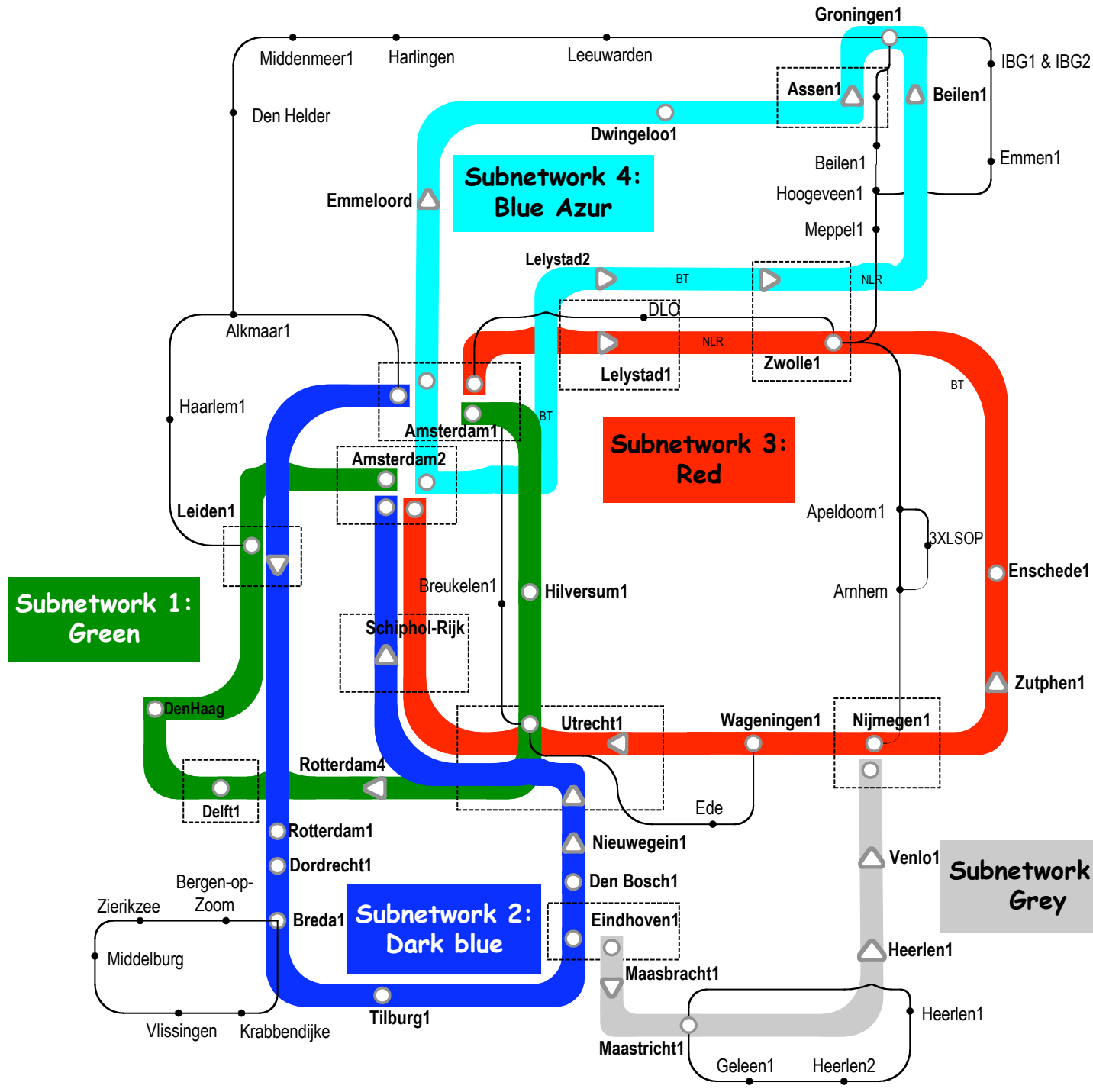
Light Paths provisioning implementation



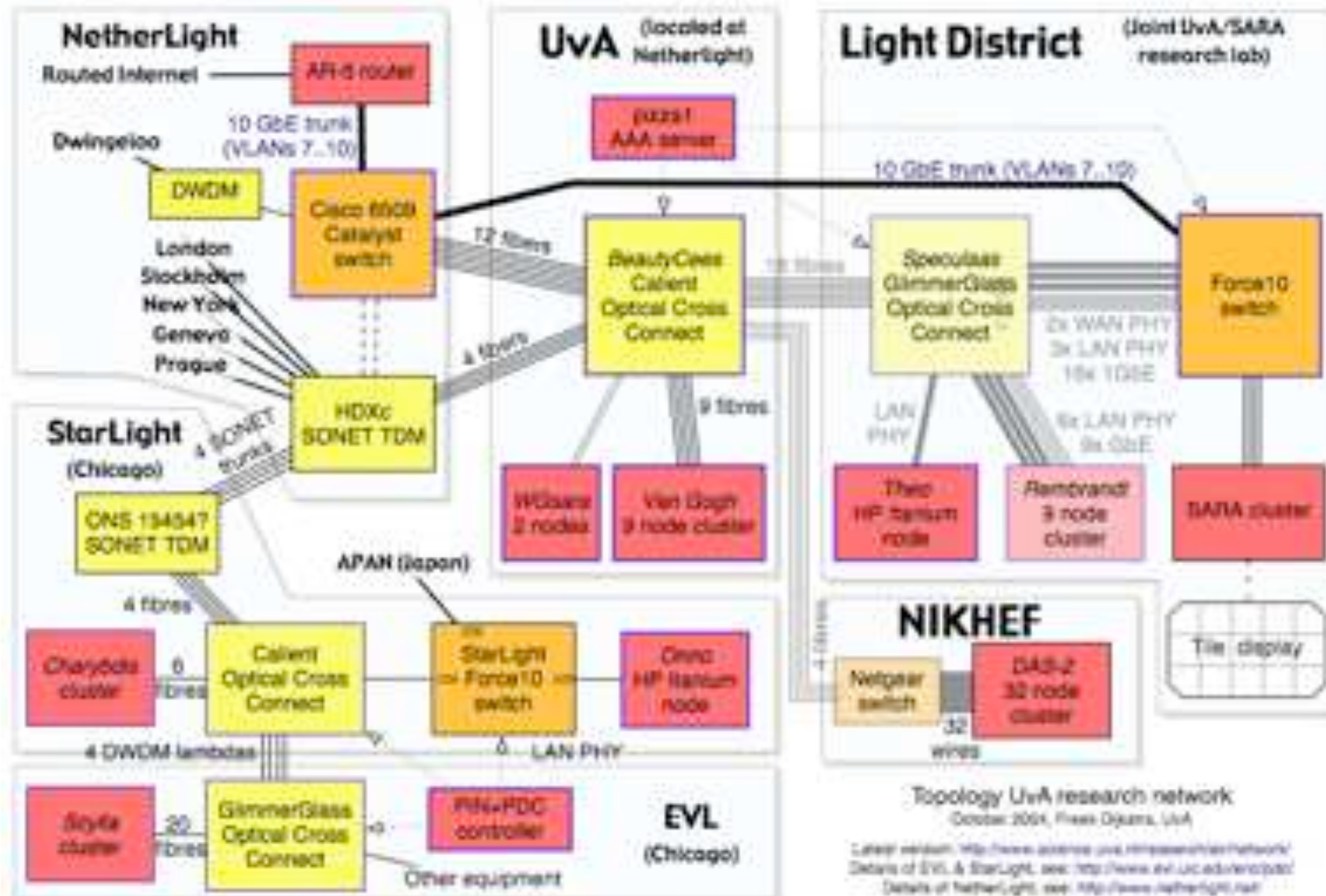
GigaPort

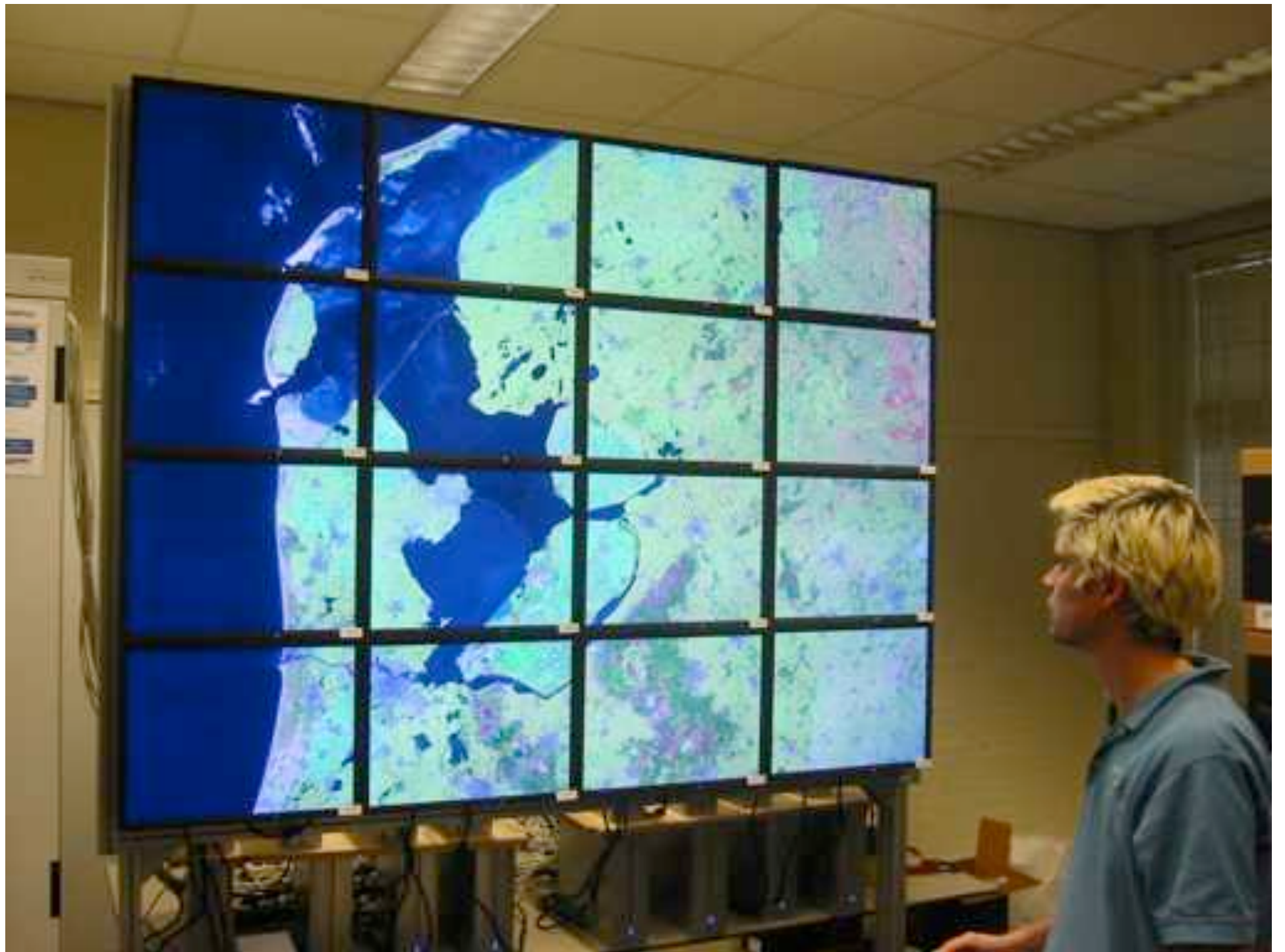
Common Photonic Layer (CPL) in SURFnet6

SURFnet



LightHouse







SURFnet is looking after your Lambda's in Science Park Amsterdam



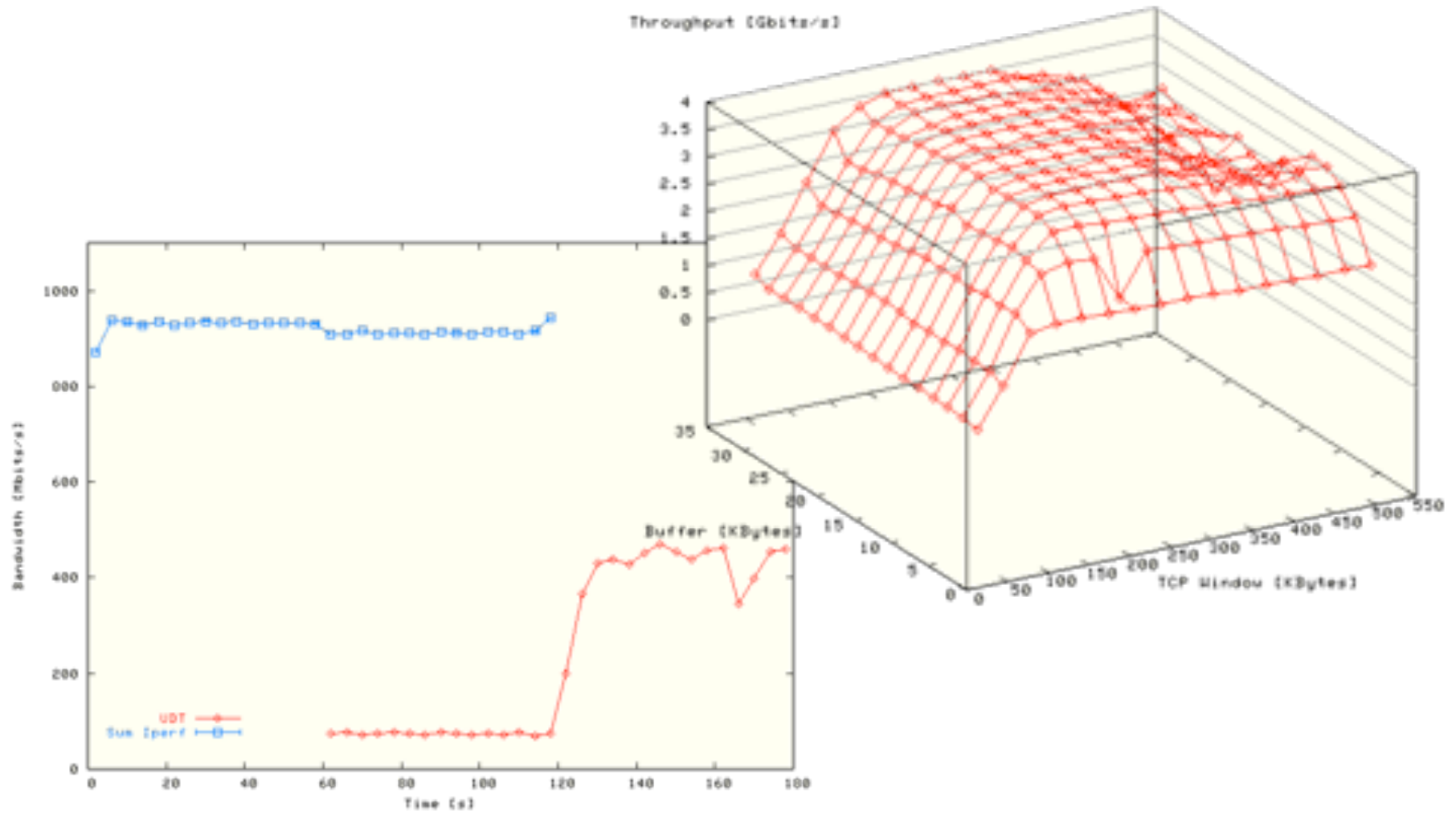
- **Optical Networking:** What innovation in architectural models, components, control and light path provisioning are needed to integrate dynamically configurable optical transport networks and traditional IP networks to a generic data transport platform that provides end-to-end IP connectivity as well as light path (lambda and sub-lambda) services?
- **High performance routing and switching:** What developments need to be made in the Internet Protocol Suite to support data intensive applications, and scale the routing and addressing capabilities to meet the demands of the research and higher education communities in the forthcoming 5 years?
- **Management and monitoring:** What management and monitoring models on the dynamic hybrid network infrastructure are suited to provide the necessary high level information to support network planning, network security and network management?
- **Grids and access; reaching out to the user:** What new models, interfaces and protocols are capable of empowering the (grid) user to access, and the provider to offer, the network and grid resources in a uniform manner as tools for scientific research?
- **Testing methodology:** What are efficient and effective methods and setups to test the capabilities and performance of the new building blocks and their interworking, needed for a correct functioning of a next generation network?



Research topics

- Optical networking architectures and models for usage
- Transport protocols for massive amounts of data
- Authorization of complex resources in multiple domains
- Embedding in Grid environments

Example Measurements



Layer - 2 requirements from 3/4



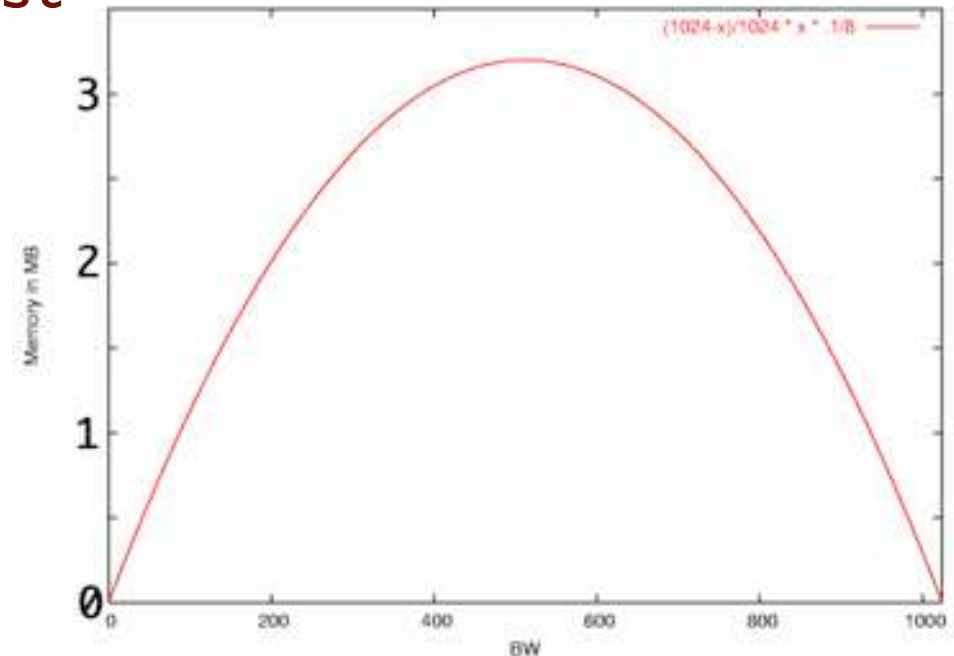
TCP is bursty due to sliding window protocol and slow start algorithm.

Window = BandWidth * RTT & BW == slow

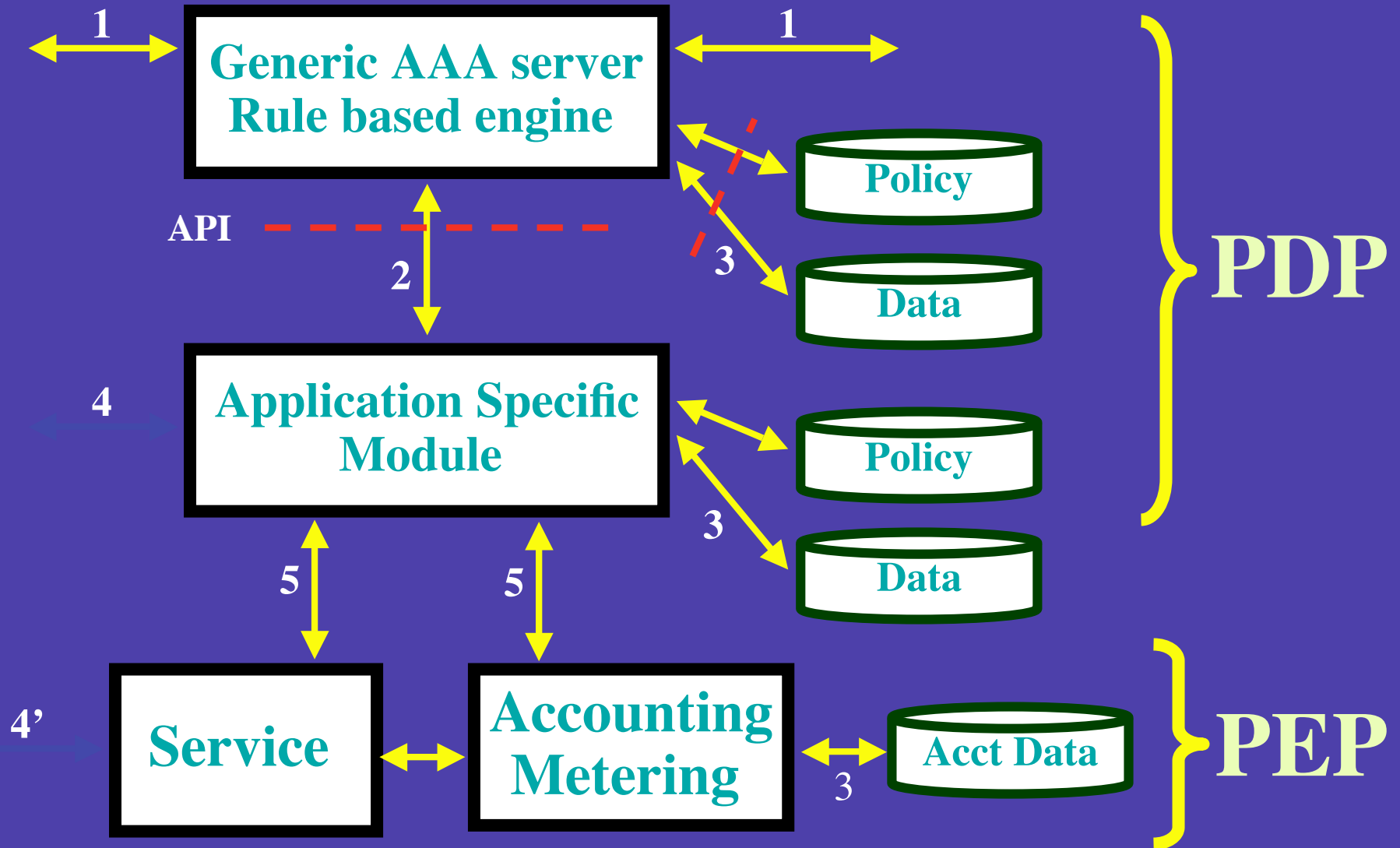
Memory-at-bottleneck = $\frac{\text{fast} - \text{slow}}{\text{fast}} * \text{slow} * \text{RTT}$

So pick from menu:

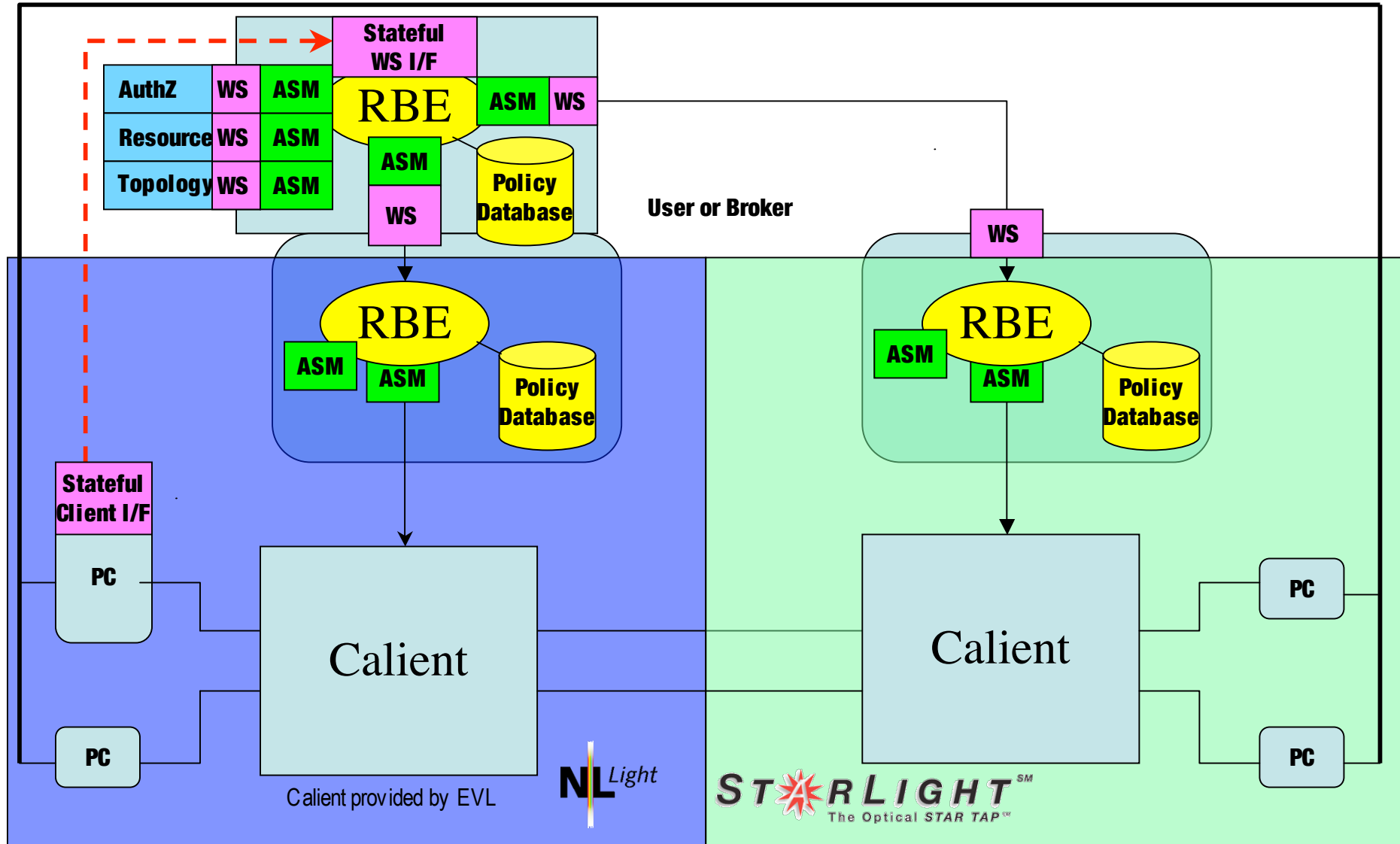
- ◆ Flow control
- ◆ Traffic Shaping
- ◆ RED (Random Early Discard)
- ◆ Self clocking in TCP
- ◆ Deep memory



Starting point



RFC 2903 - 2906 , 3334 , policy draft

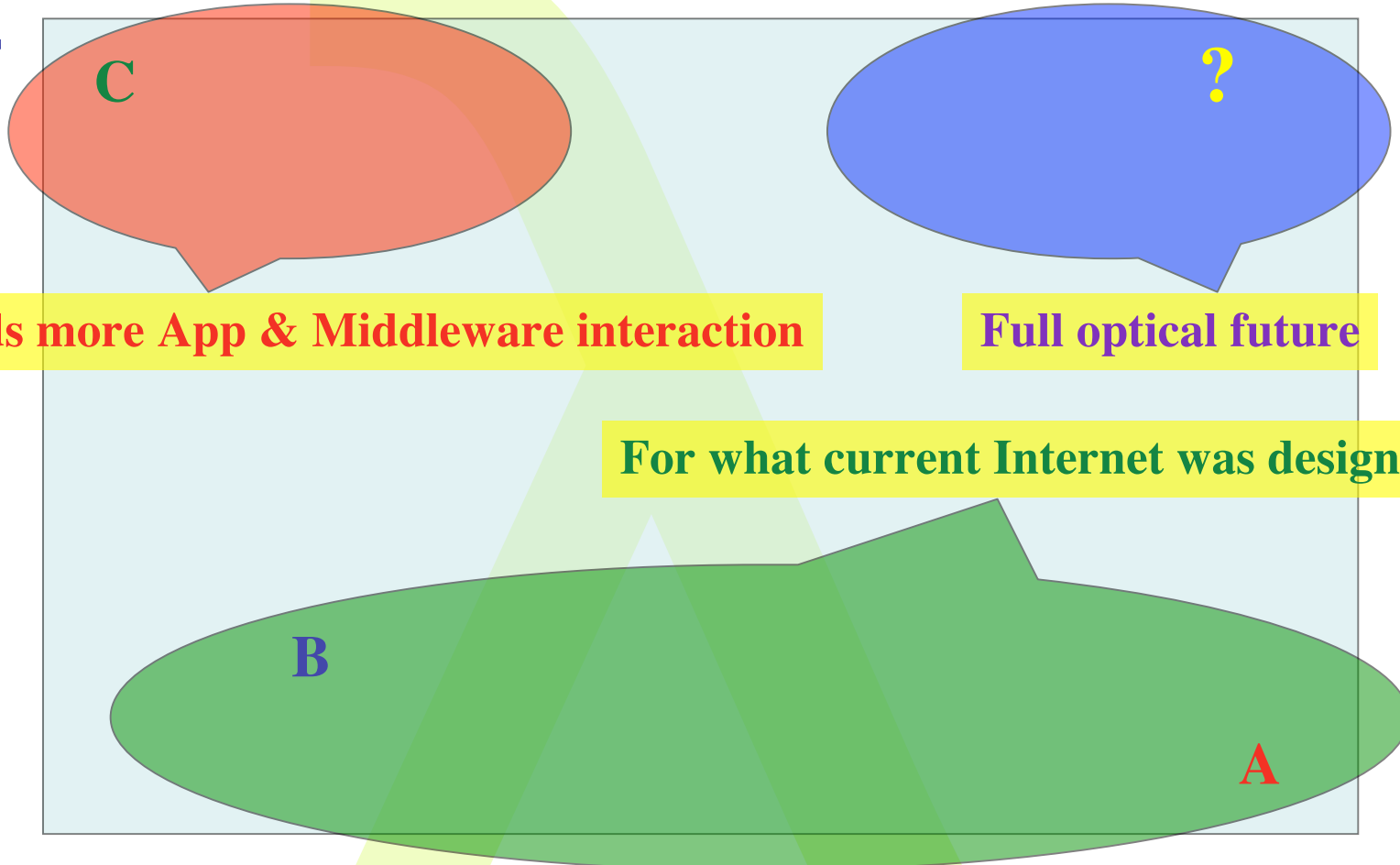


Conclusions

- Demanding applications
 - (Science) data repositories mirroring
 - Instrumentation grids
 - Visualisation and collaboration support
- Model of Lambda networking
 - Identify traffic types
 - Scales of infrastructure
 - Map efficiently to lower the cost/packet
- Current experiments
 - NetherLight
 - VLE/eScience Amsterdam
 - Networking research
(control plane, transport protocols, optical net models)

Transport in the corners

$BW * RTT$



Needs more App & Middleware interaction

Full optical future

For what current Internet was designed

FLOWS

Not quite The END

Thanks to

SURFnet: Kees Neggens, UIC&iCAIR: Tom DeFanti, Joel Mambretti, CANARIE: Bill St. Arnaud

Freek Dijkstra, Hans Blom, Leon Gommans, Bas van oudenaarde, Arie Taal, Pieter de Boer, Bert Andree, Martijn de Munnik, Antony Antony, Rob Meijer, VL-team.



Partially complete list:

- Caas
- Chase
- Cess
- Kess
- Case

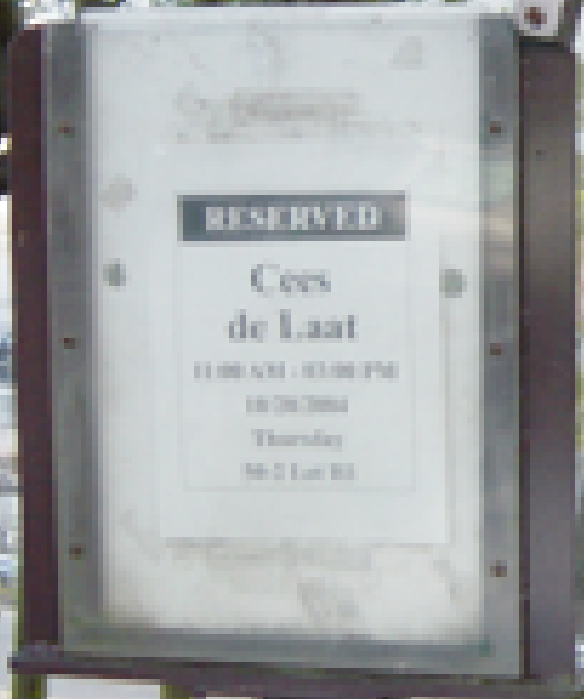


The END

Thanks to

SURFnet: Kees Neggers, UIC&iCAIR: Tom DeFanti, Joel Mambretti, CANARIE: Bill St. Arnaud

Frek Dijkstra, Hans Blom, Leon Gommans, Bas van Oulenaarde, Arie Taal, Pieter de Boer, Bert Andree, Martijn de Munnik, Antony Antony, Rob Meijer, VL-team



Partially complete list:

- Caas
- Chase
- Cess
- Kess
- Case

