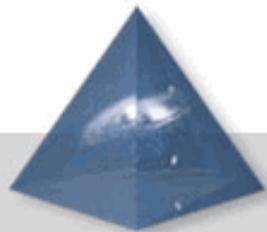# High Performance Networking for Grid Applications

www.science.uva.nl/~delaat

## Cees de Laat

Faculty of Science

# High Performance Networking for Grid Applications

www.science.uva.nl/~delaat

## Cees de Laat

# EU

# SURFnet

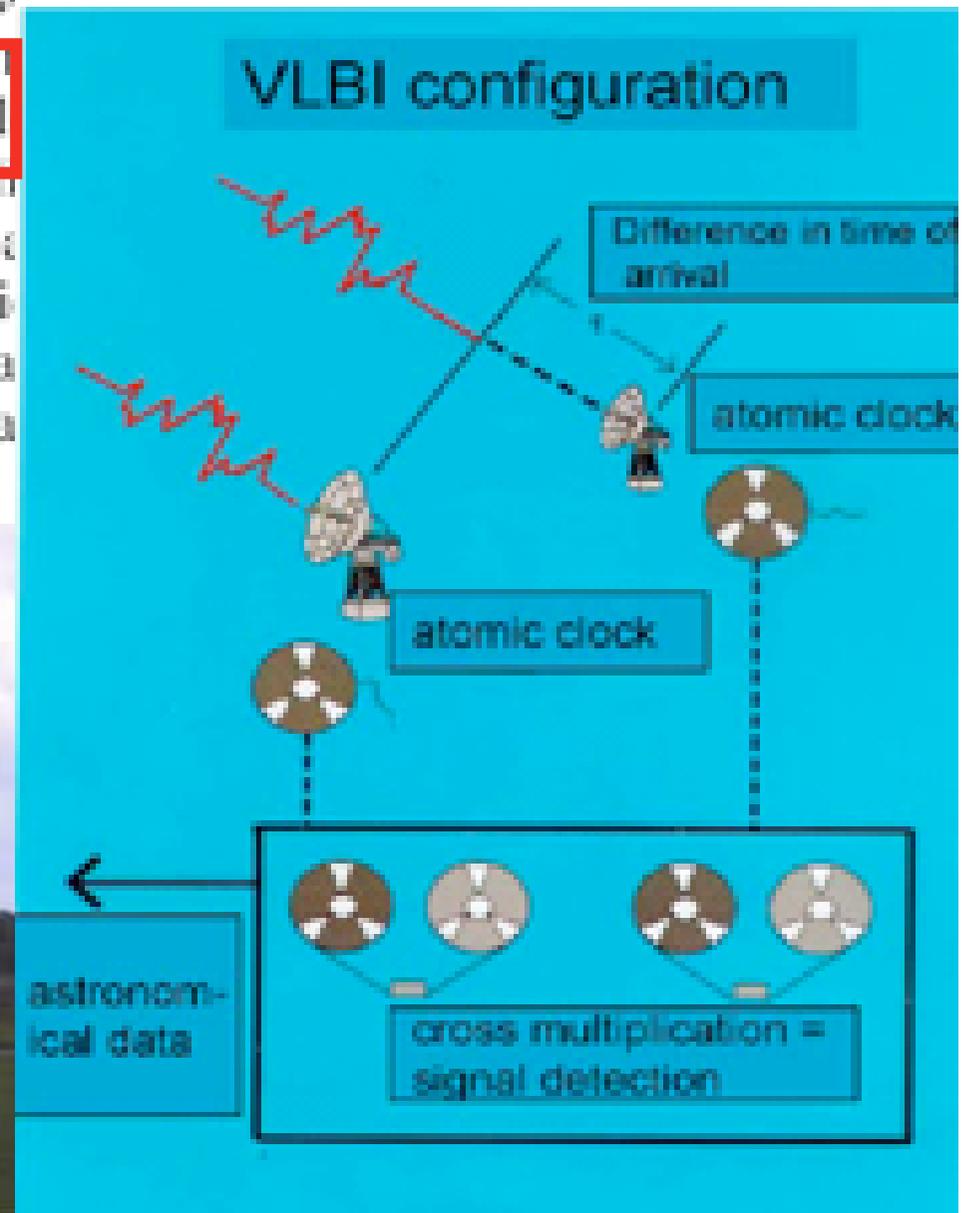## University of Amsterdam

SARA
NIKHEF
NCF

# Contents of this talk

This slide is intentionally left blank

# eVLBI

# VLBI

ger term VLBI is easily capable of generating many Gb of data per
The sensitivity of the VLBI array scales with the square root of the
(=data-rate) and there is a strong push to
. Rates of 8Gb/s or more are entirely feasibl
der development. It is expected that paral
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.

# iGrid 2002

## September 24-26, 2002, Amsterdam, The Netherlands

- 28 demonstrations from 16 countries: Australia, Canada, CERN, France, Finland, Germany, Greece, Italy, Japan, The Netherlands, Singapore, Spain, Sweden, Taiwan, United Kingdom, United States

- Applications demonstrated: art, bioinformatics, chemistry, cosmology, cultural heritage, education, high-definition media streaming, manufacturing, medicine, neuroscience, physics, tele-science

- Grid technologies demonstrated: Major emphasis on grid middleware, data management grids, data replication grids, visualization grids, data/visualization grids, computational grids, access grids, grid portals

- 25Gb transatlantic bandwidth (100Mb/attendee, 250x iGrid2000!)

www.igrid2002.org
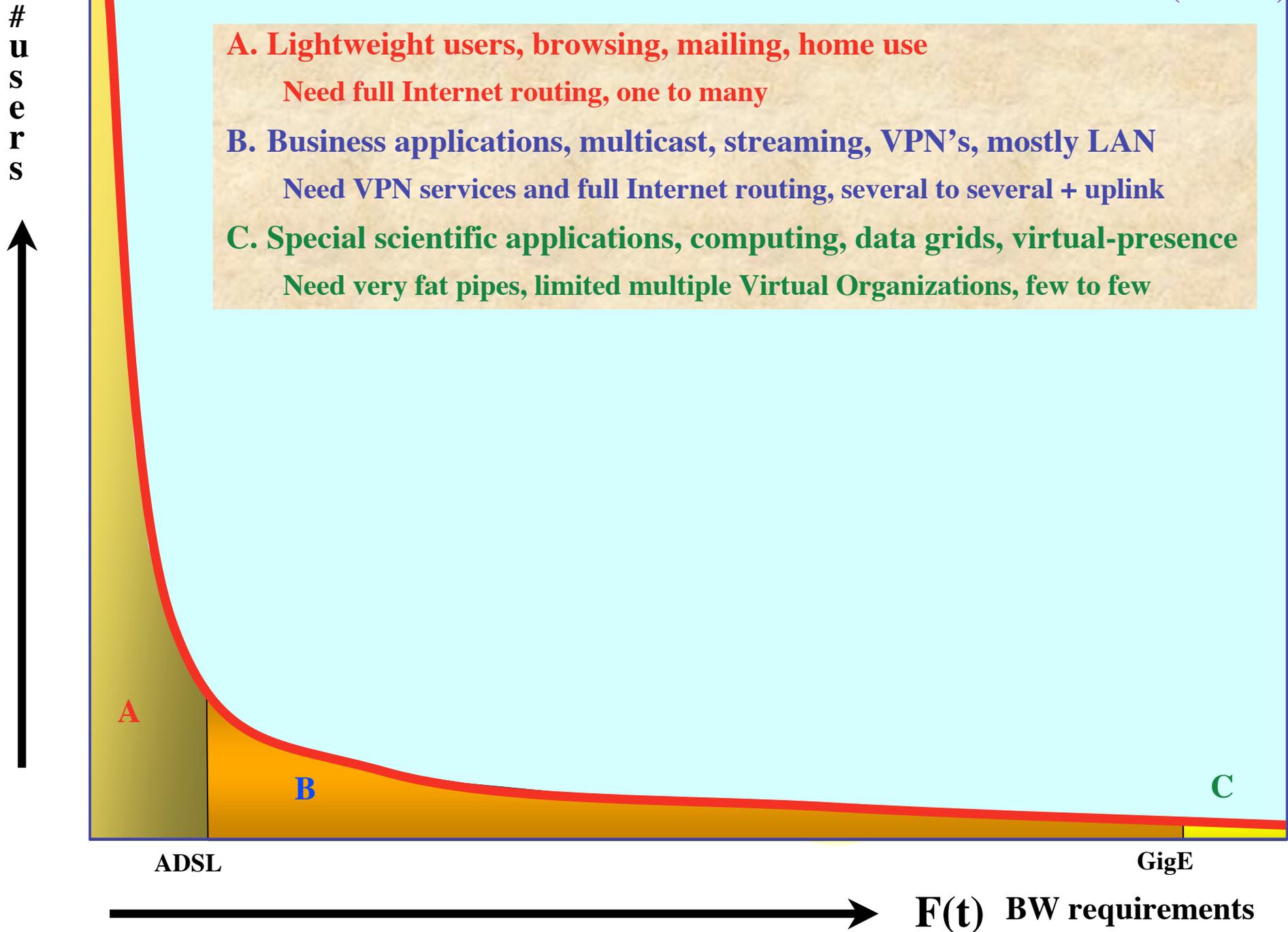
# Experimental Networks

- **High-performance trials of new technologies that support *application-dictated* development of software toolkits, middleware, computing and networking.**

- **Provide *known and knowable characteristics* with deterministic and repeatable behavior on a persistent basis, while encouraging experimentation with innovative concepts.**

- **Experimental Networks are seen as the *missing link* between Research and Production Networks.**

http://www.evl.uic.edu/activity/NSF/index.html
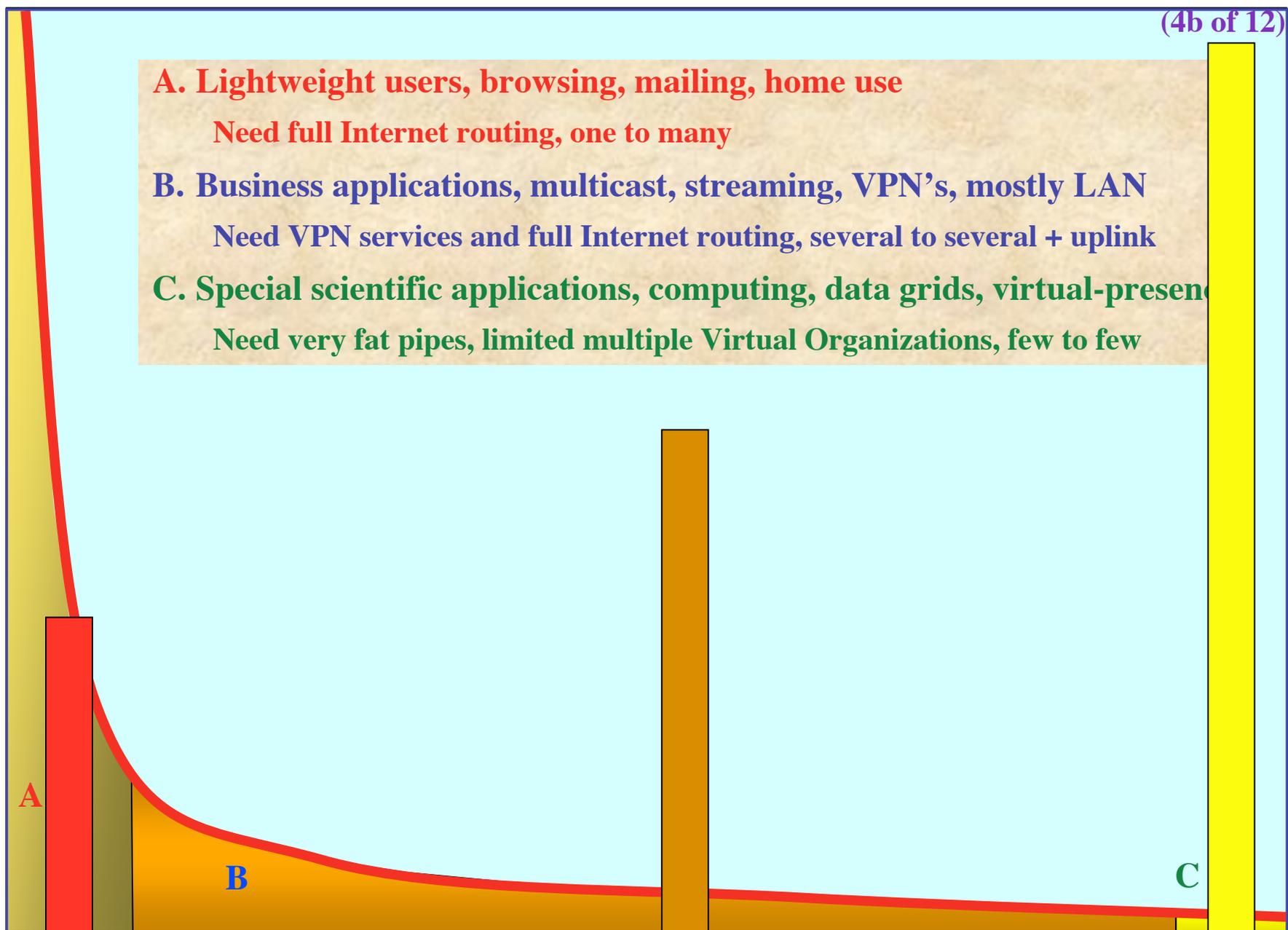http://www.calit2.net/events/2002/nsf/index.html

# What is a LambdaGrid?

- A *grid* is a set of networked, middleware-enabled computing resources.

- A *LambdaGrid* is a grid in which the lambda networks themselves are resources that can be scheduled, like all other computing resources. The ability to schedule and provision lambdas provides *deterministic* end-to-end network performance for real-time or time-critical applications, which cannot be achieved on today's grids.

**#users**

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
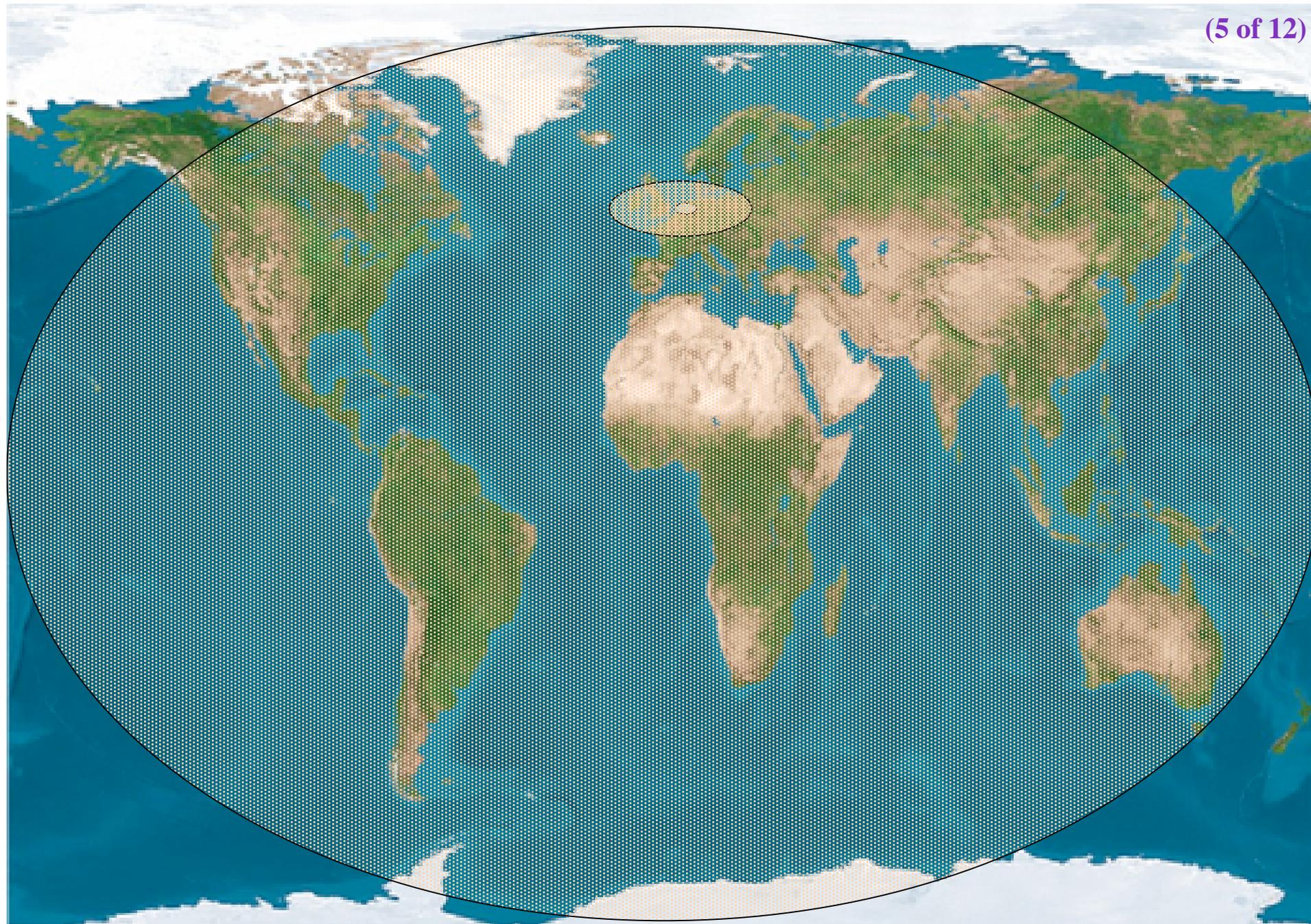
A

B

C

ADSL

GigE

**F(t)** BW requirements

**#users**

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

A

B

C

ADSL

GigE

F(t) BW requirements

Scale 2-20-200

# The only formula's

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

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

**c = 1**

**e = 1**

**h̄ = 1**

**and the formula reduces to:** $\# \lambda \approx \frac{200 * e^{(t-2002)}}{rtt}$

SURFnet
Lambda's
fibers
(old already)

# Services

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

# Current technology + (re)definition

- Current (to me) available technology consists of SONET/SDH switches, 10 gig ethernet and dark fiber environments

- Optical switch installed (this week)!

- DWDM+switching included

- Starlight/NetherLight deploy 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"**
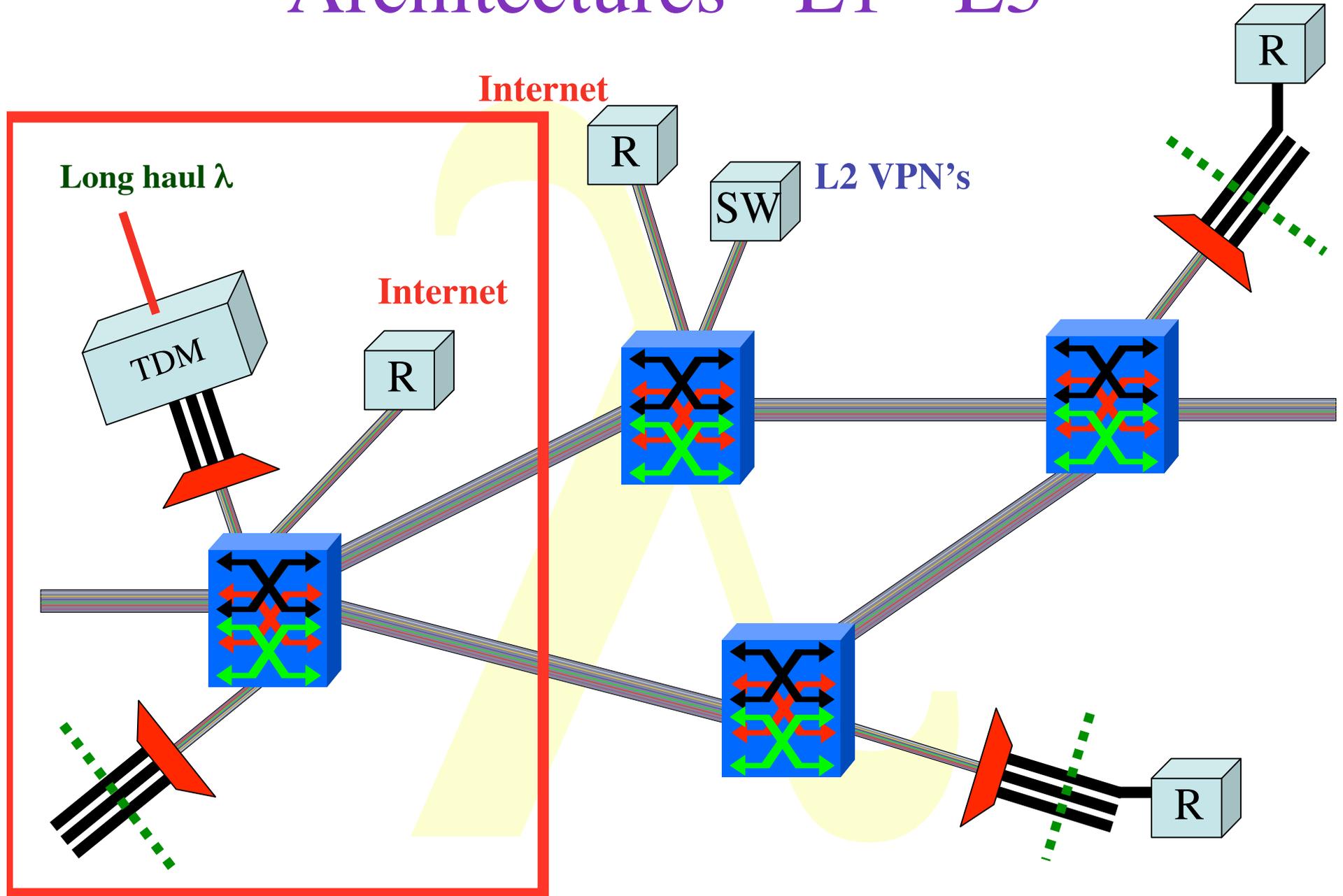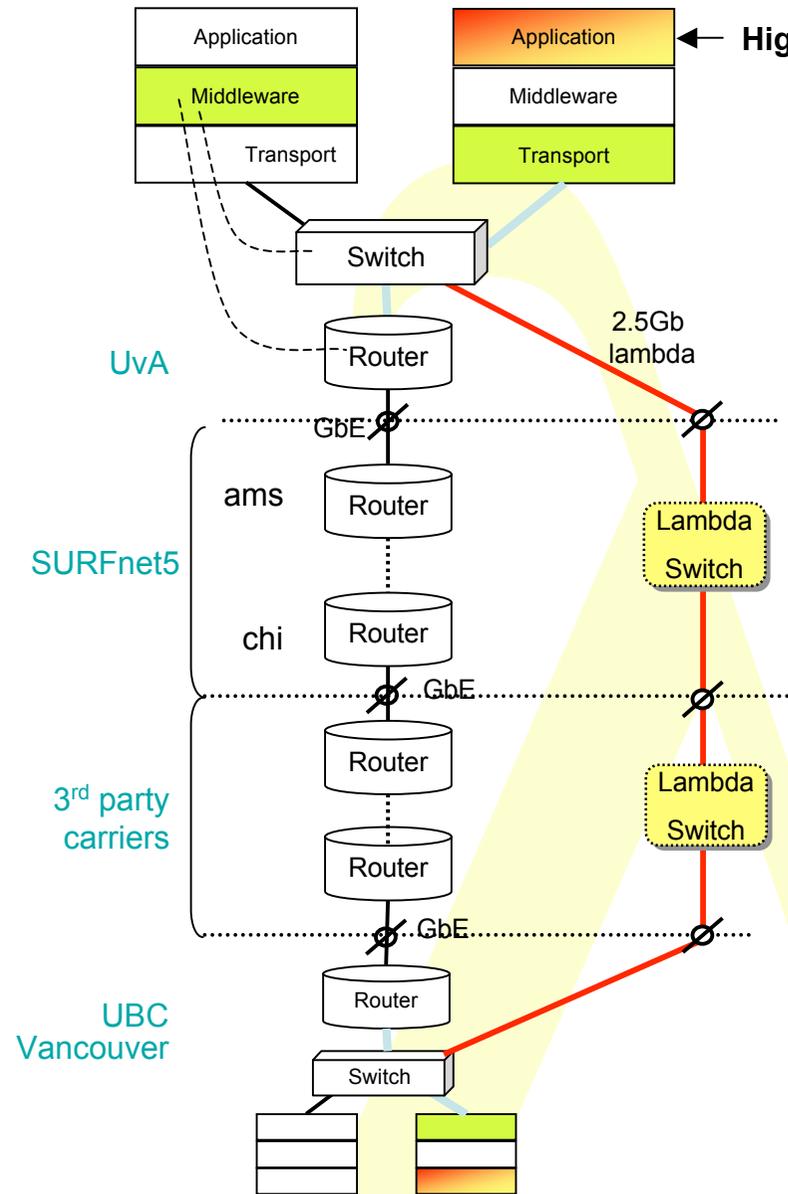
# MEMS optical switch (CALIENT)

# So what are the facts

- **Costs of fat pipes (fibers) are one/third of cost 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 @ 40 Gb/s -> 20 ns to look up in 140 kEntries routing table (light speed from me to you!)

- **Big sciences need fat pipes**

- **Bottom line: look for a hybrid architecture which serves all users in a cost effective way**

# Architectures - L1 - L3

**Internet**

**Long haul λ**

**Internet**
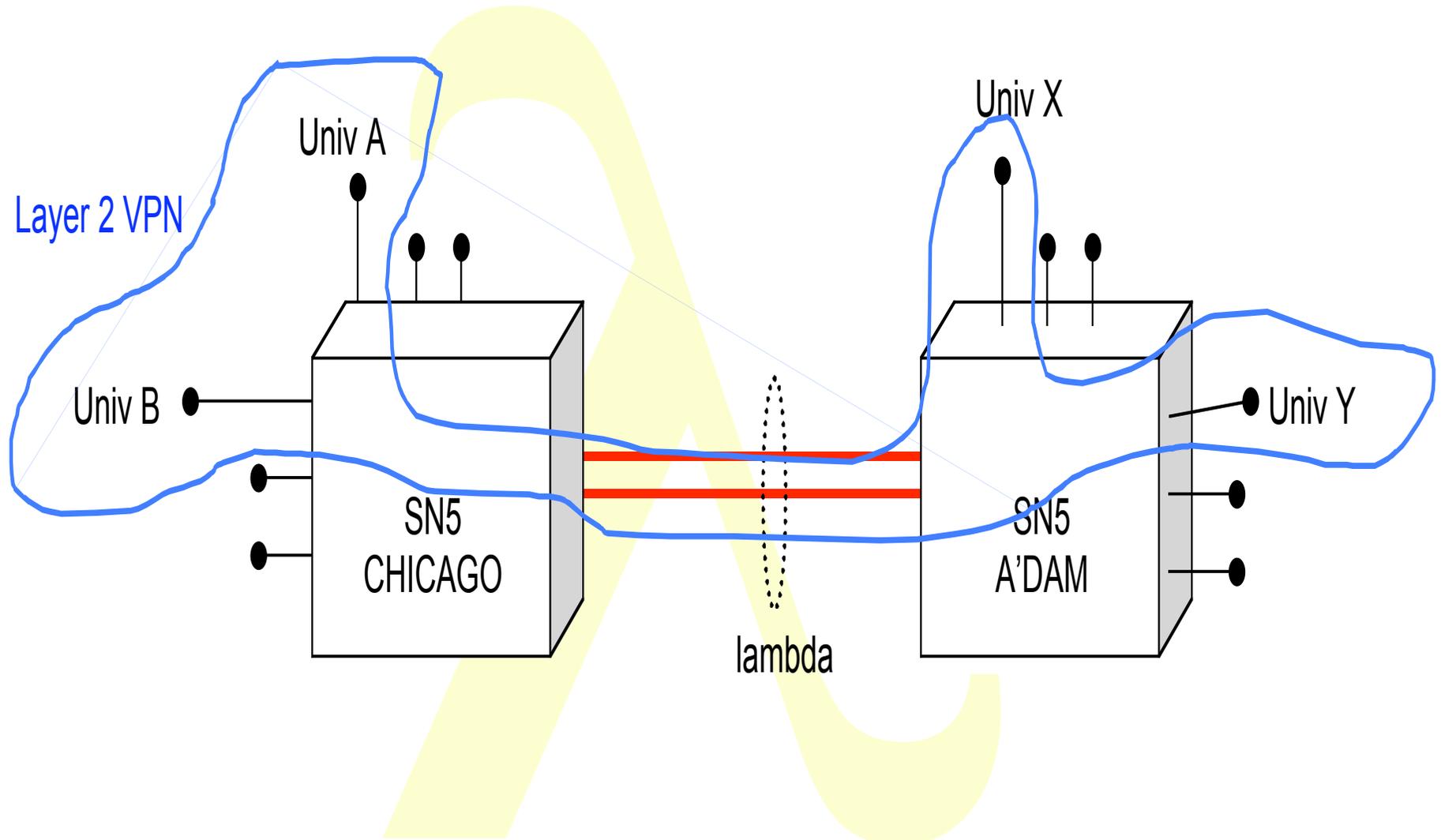
**L2 VPN's**

TDM

R

SW

R

R

R

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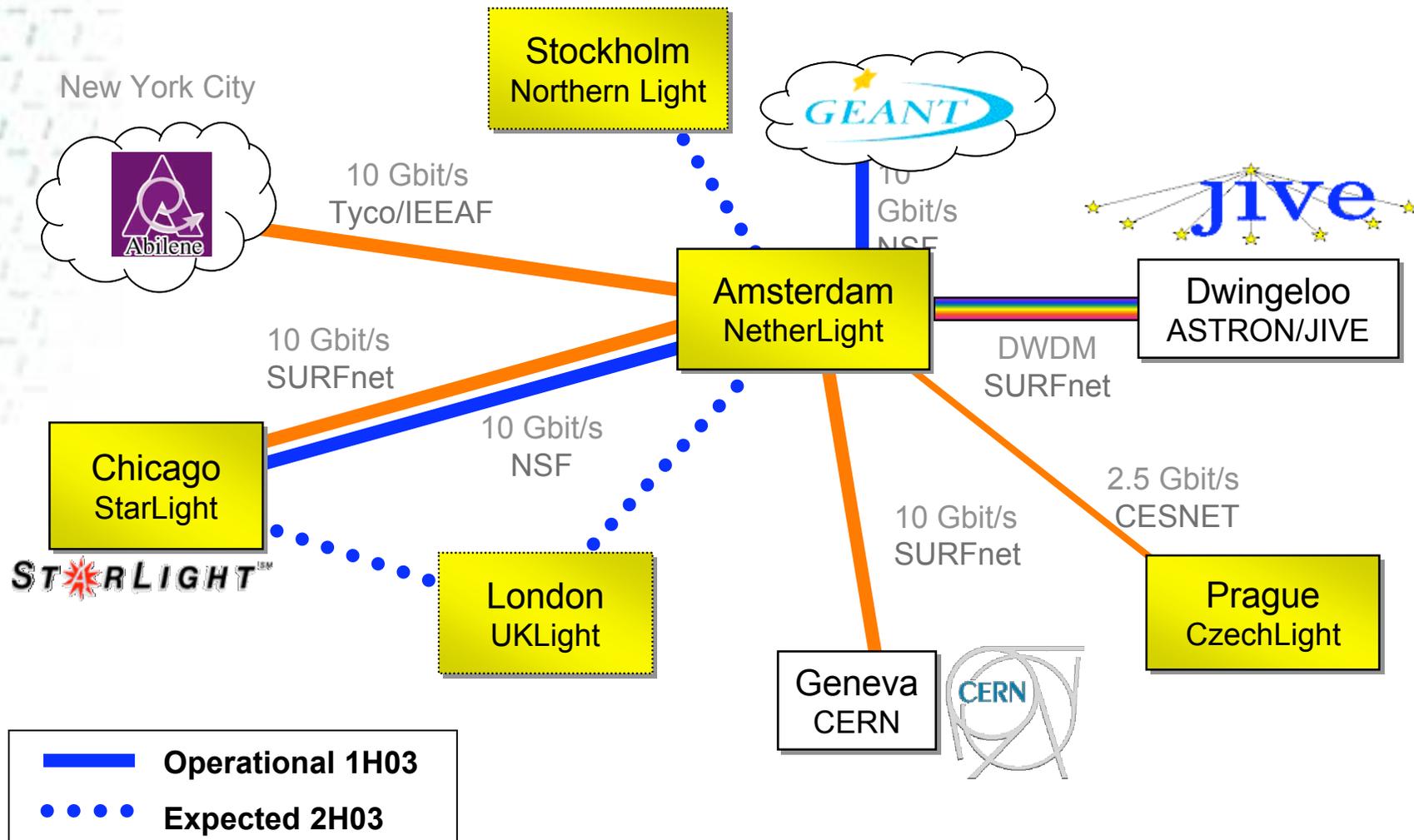  - Lower the cost of transport per packet

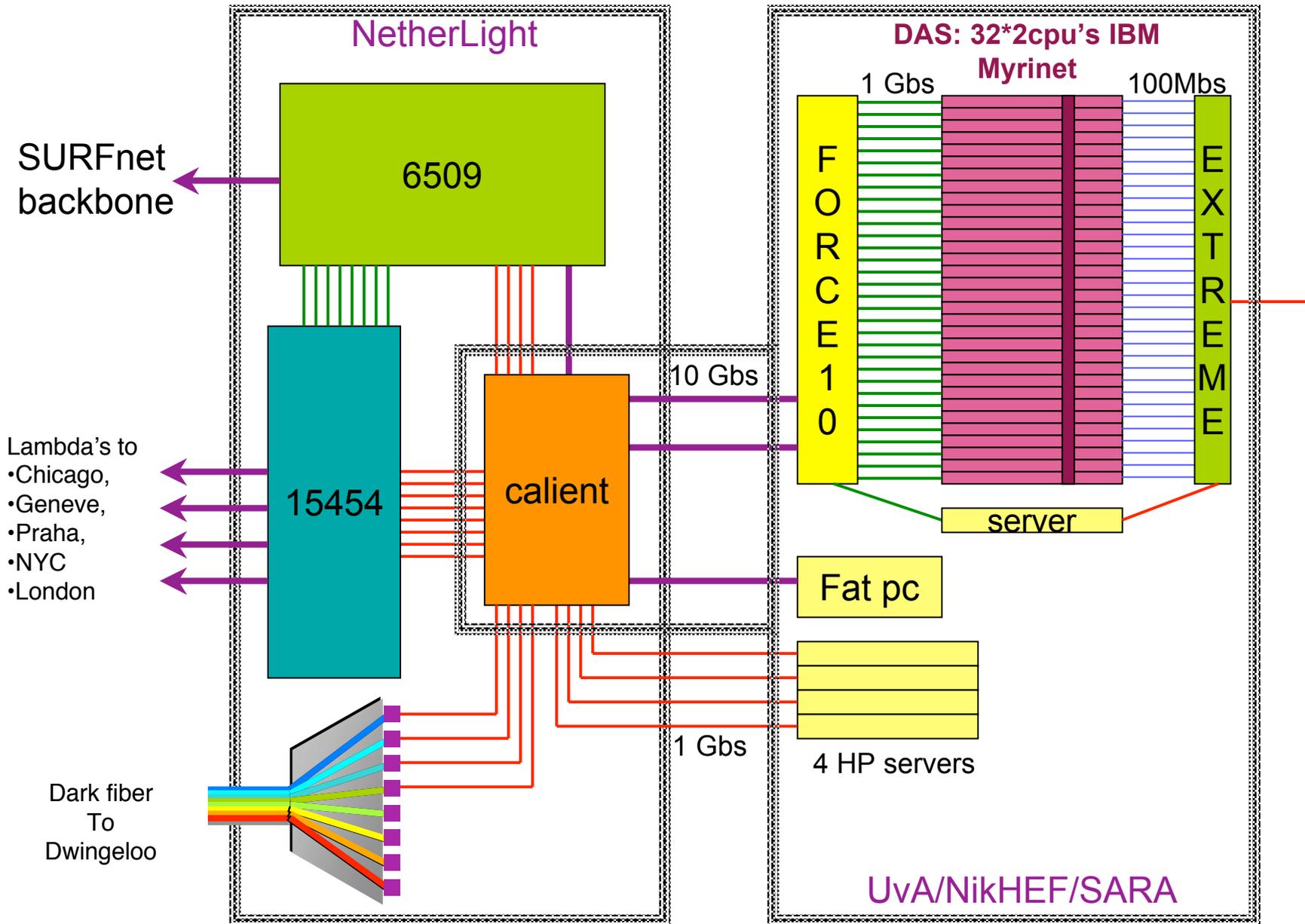# How low can you go?



Application Endpoint A

Local Ethernet

MEMS

Regional dark fiber

ONS 15454

Application Endpoint B

POS

Trans-Atlantic

Router

Ethernet

SONET

DWDM

fiber

NetherLight

NetherLight

TransLight

# Virtual Organization on L2

**SURF net** High-quality Internet for higher education and research

# NetherLight Network: 2003
## Emerging international lambda grid

NetherLight

DAS: 32*2cpu's IBM
Myrinet

6509

SURFnet
backbone

1 Gbs

100Mbs

F
O
R
C
E
1
0

E
X
T
R
E
M
E

10 Gbs

15454

calient

Lambda's to
•Chicago,
•Geneve,
•Praha,
•NYC
•London

server

Fat pc

1 Gbs

4 HP servers

Dark fiber
To
Dwingeloo

UvA/NikHEF/SARA

NetherLight

# Early Lambda/LightPath TDM experiences

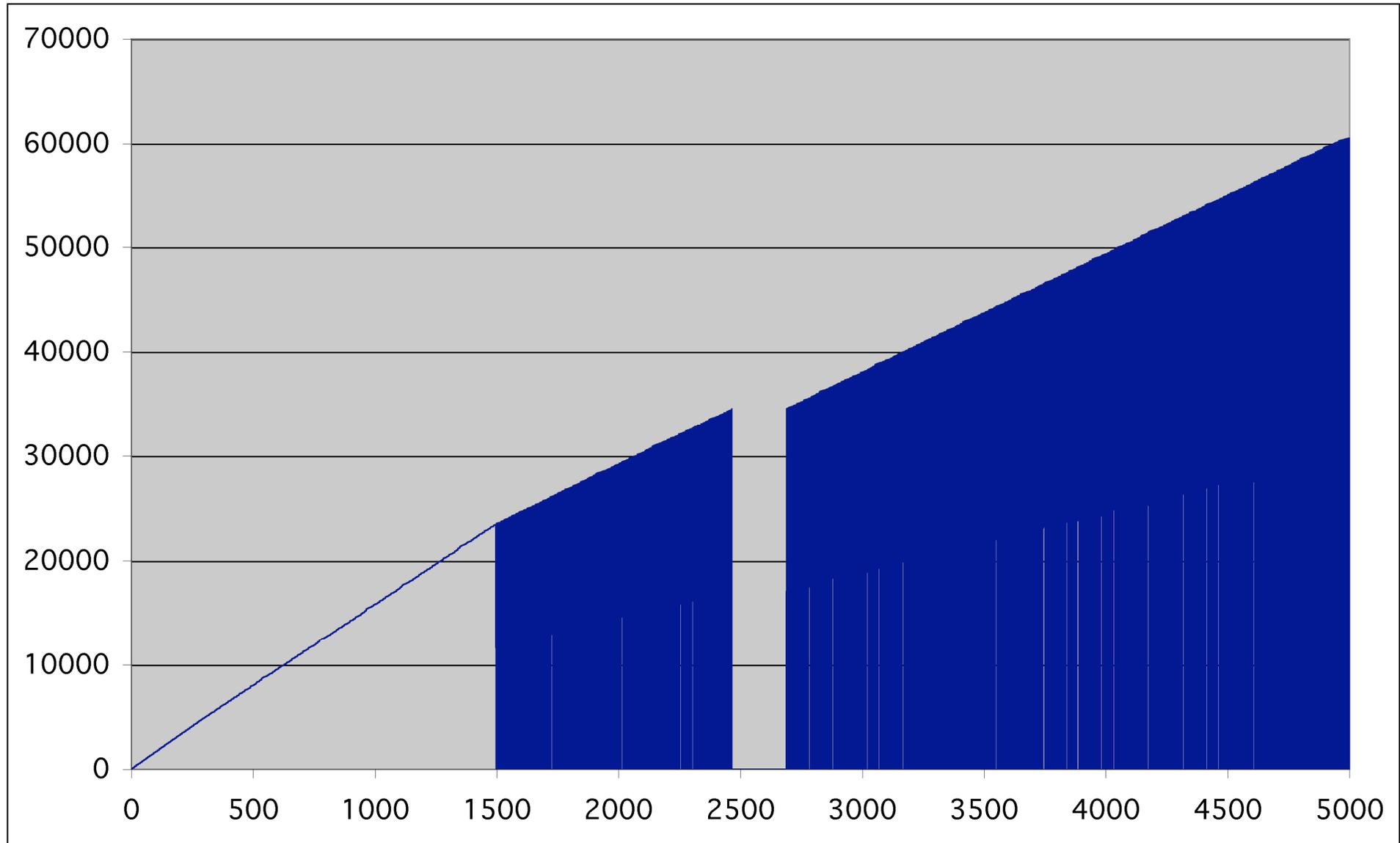| WS | **fast** | **L2** **fast->slow** | **slow** **high RTT** | **L2** **slow->fast** | **fast** | WS |

# 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
```
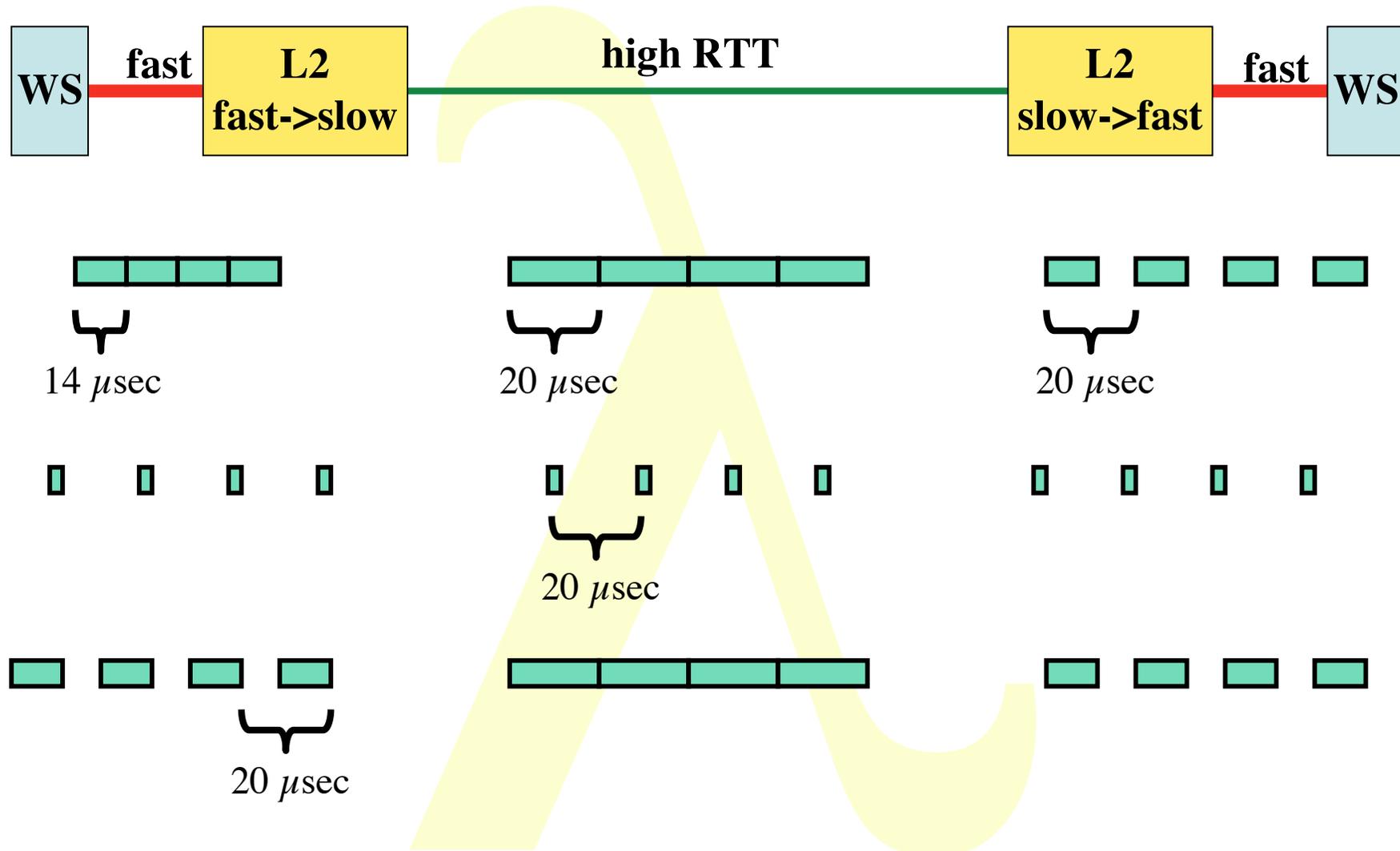
So pick from menu:

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

# Self-clocking of TCP

# Layer - 2 requirements from 3/4

WS — **fast** — L2 **fast->slow** — **high RTT** — L2 **slow->fast** — **fast** — WS

Window = BandWidth * RTT    &    BW == slow

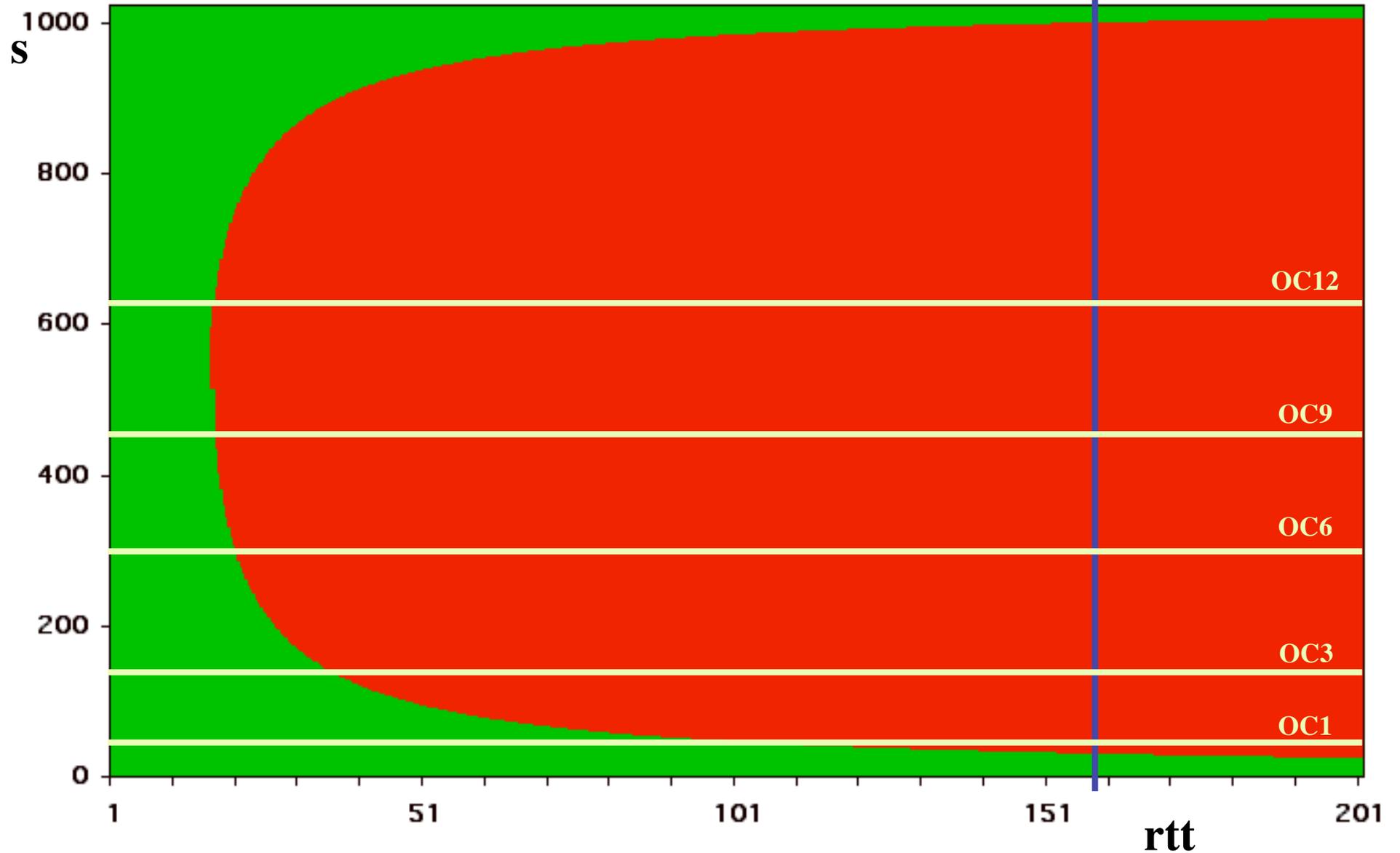$$\text{Memory-at-bottleneck} = \frac{\text{fast - slow}}{\text{fast}} * \text{slow} * \text{RTT}$$

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

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

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

Forbidden area, solutions for s when f = 1 Gb/s, M = 0.5 Mbyte (19c of 20)
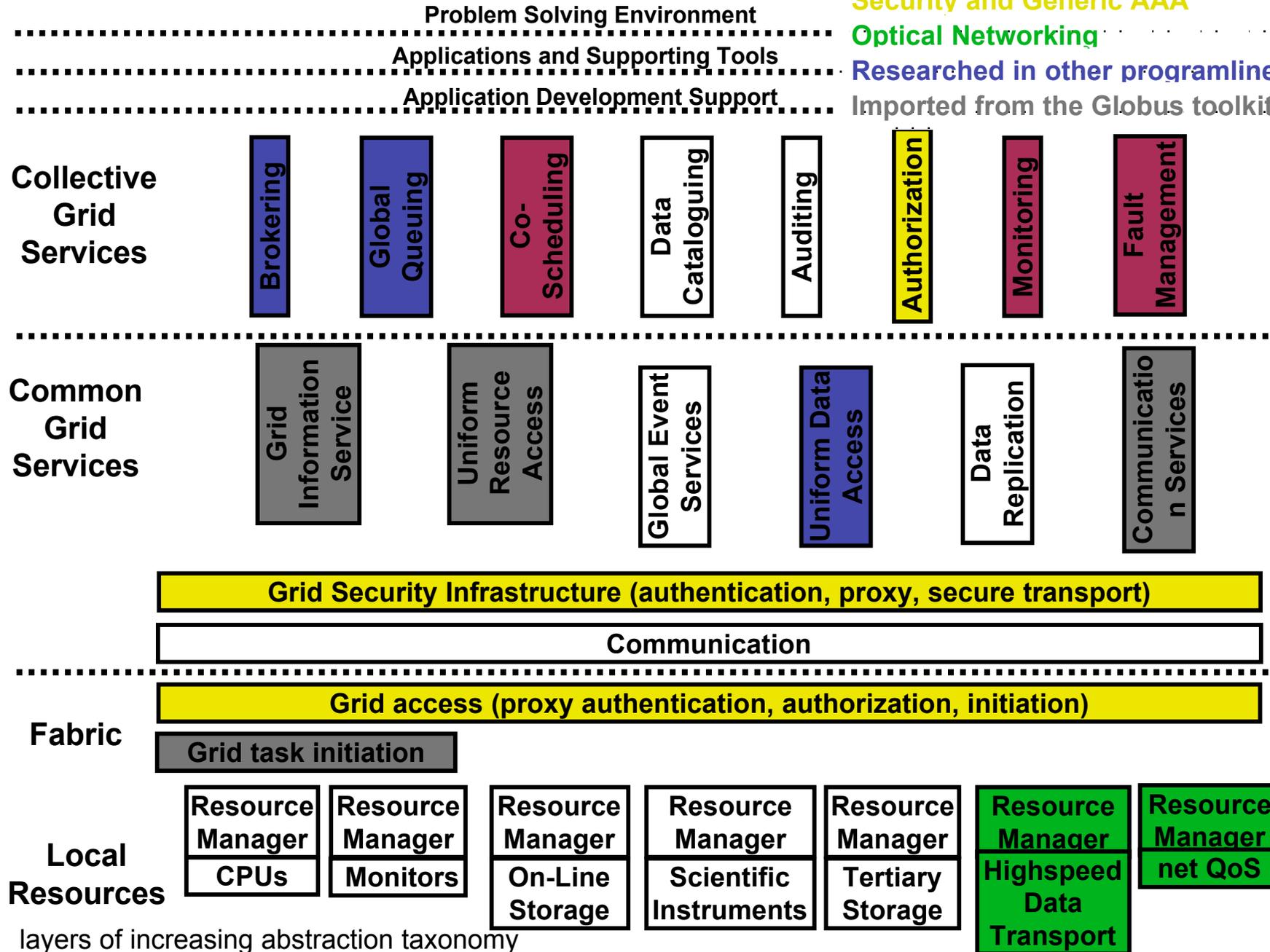AND NOT USING FLOWCONTROL

High performance computing and
Processor memory co-allocation

Security and Generic AAA

Optical Networking

Researched in other programlines
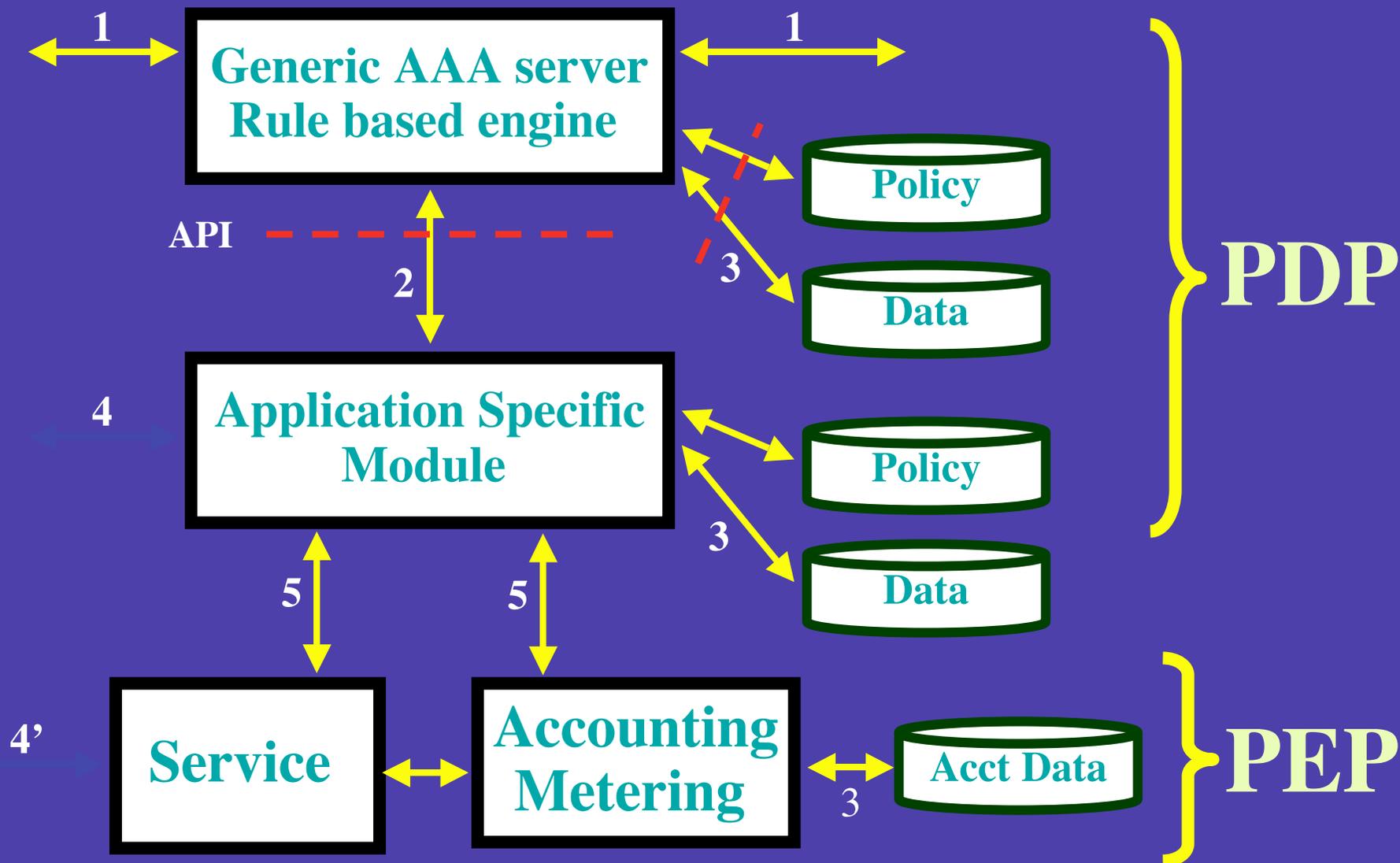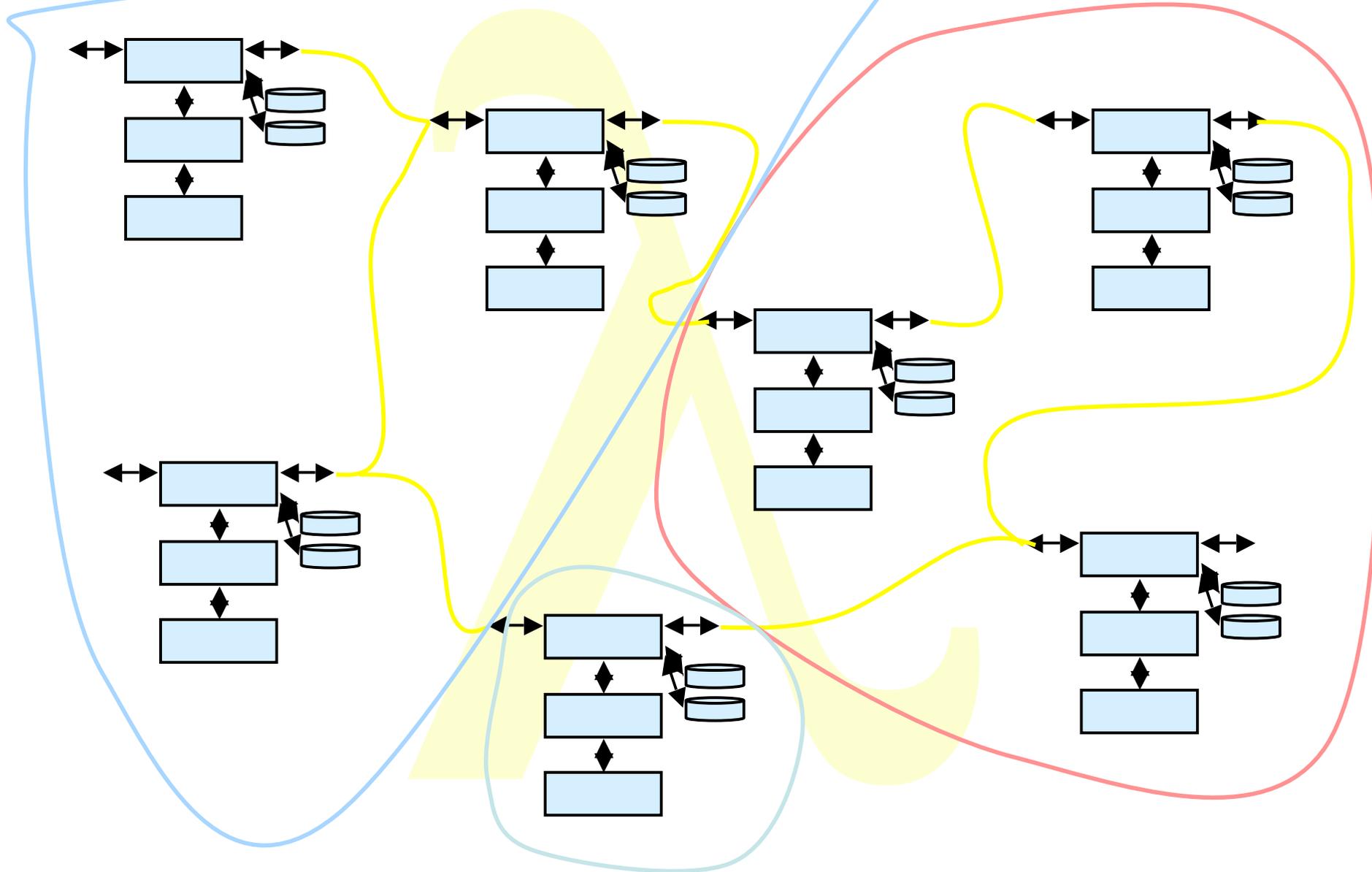
Imported from the Globus toolkit

Problem Solving Environment

Applications and Supporting Tools

Application Development Support

**Collective Grid Services**

- Brokering
- Global Queuing
- Co-Scheduling
- Data Cataloguing
- Auditing
- Authorization
- Monitoring
- Fault Management

**Common Grid Services**

- Grid Information Service
- Uniform Resource Access
- Global Event Services
- Uniform Data Access
- Data Replication
- Communication Services

**Grid Security Infrastructure (authentication, proxy, secure transport)**

**Communication**

**Grid access (proxy authentication, authorization, initiation)**

**Fabric**

**Grid task initiation**

**Local Resources**

| Resource Manager | Resource Manager | Resource Manager | Resource Manager | Resource Manager | Resource Manager | Resource Manager |
|---|---|---|---|---|---|---|
| CPUs | Monitors | On-Line Storage | Scientific Instruments | Tertiary Storage | Highspeed Data Transport | net QoS |

layers of increasing abstraction taxonomy

# Starting point



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

# Multi domain case

# (Future) Projects

- **National:**
  - **NCF Grid project**
  - **VLE**
  - **GigaPort-NG**
  - **LOFAR**
- **European**
  - **DataGrid**
  - **DataTAG**
- **International**
  - **NetherLight**
  - **StarLight**
  - **AnyLight, LowLight, BackLight**
  - **Optiputer**

**Research:**

**Models of Lambda networking**

**Transport**

**AAA**

# The END

RESERVED

Case
Delaat