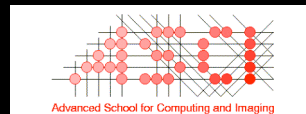


**Lambda-Grid developments**

# **StarPlane**

**Cees de Laat**

**University of Amsterdam**



# users

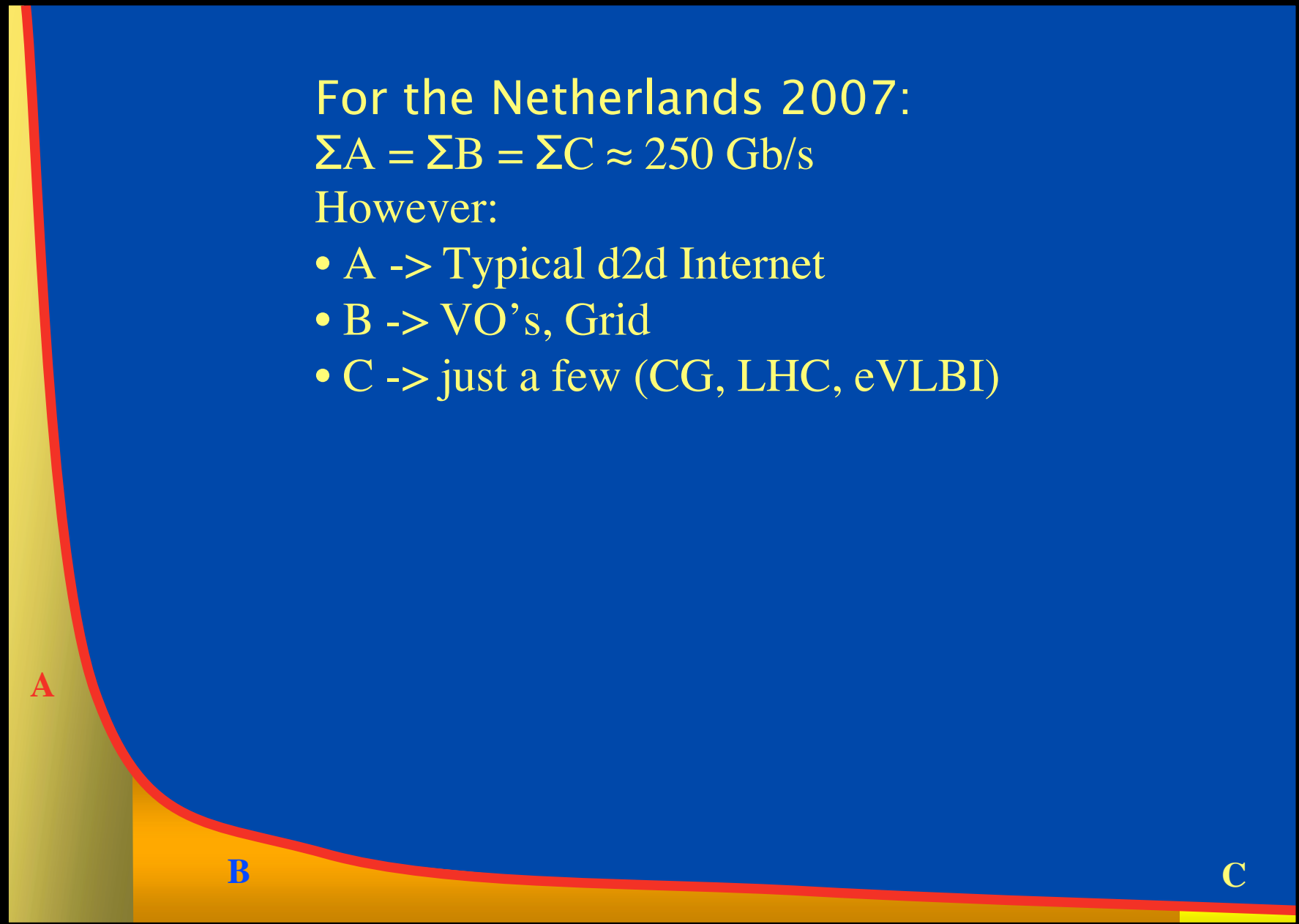


For the Netherlands 2007:

$$\Sigma A = \Sigma B = \Sigma C \approx 250 \text{ Gb/s}$$

However:

- A -> Typical d2d Internet
- B -> VO's, Grid
- C -> just a few (CG, LHC, eVLBI)



ADSL (12 Mbit/s)

BW requirements

GigE

CdL



# Infrastructure

<b>SCALE</b>  <b>CLASS</b>	<b>2</b> <b>Metro</b>	<b>20</b> <b>Regional</b>	<b>200</b> <b>World</b>
<b>A</b>	<b>Switching/ Routing</b>	<b>Routers</b>	<b>ROUTER\$</b>
<b>B</b>	<b>Switches VPN's E-WANPHY</b>	<b>Routing Switches (G)MPLS E-WANPHY</b>	<b>ROUTER\$</b>
<b>C</b>	<b>dark fiber DWDM WSS Photonic switch</b>	<b>DWDM, TDM / SONET Lambda switching</b>	<b>VLAN's TDM SONET Ethernet</b>



In The Netherlands SURFnet connects between 180:

- universities;
- academic hospitals;
- most polytechnics;
- research centers.

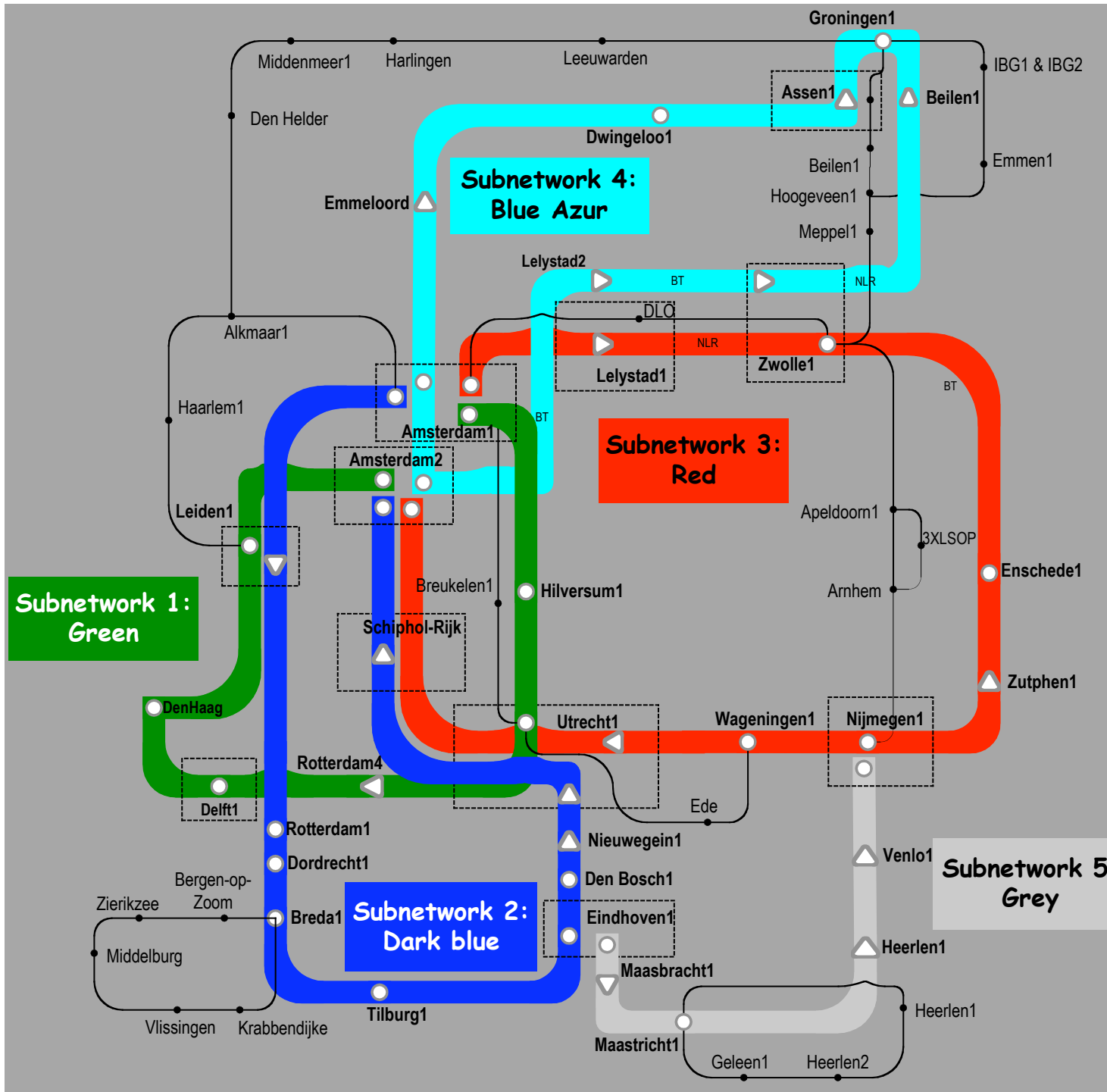
with an indirect ~750K user base

~ 6000 km  
scale  
comparable  
to railway  
system

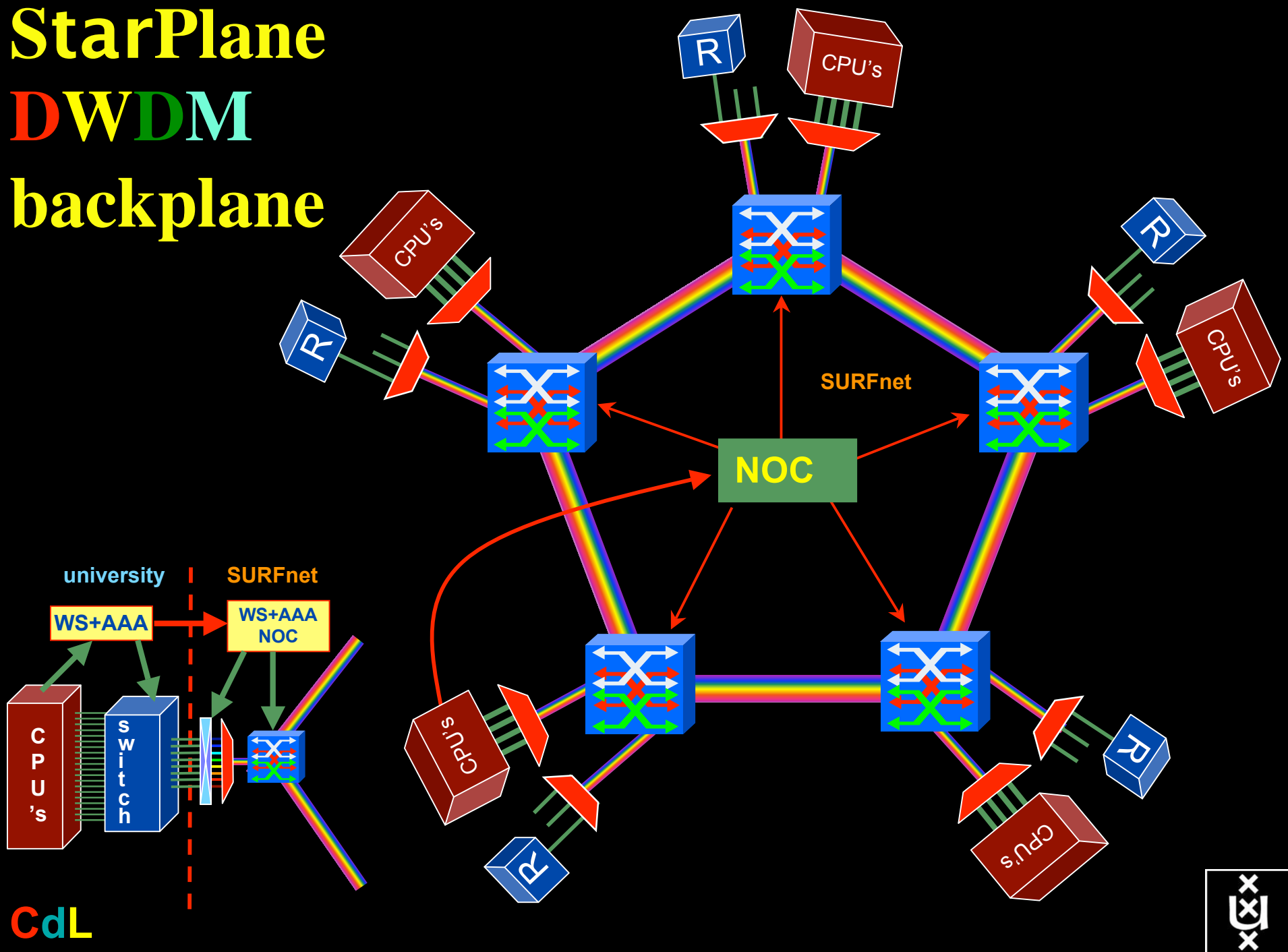
# Common Photonic Layer (CPL) in SURFnet6

supports up to 72 Lambda's of 10 G each  
future:

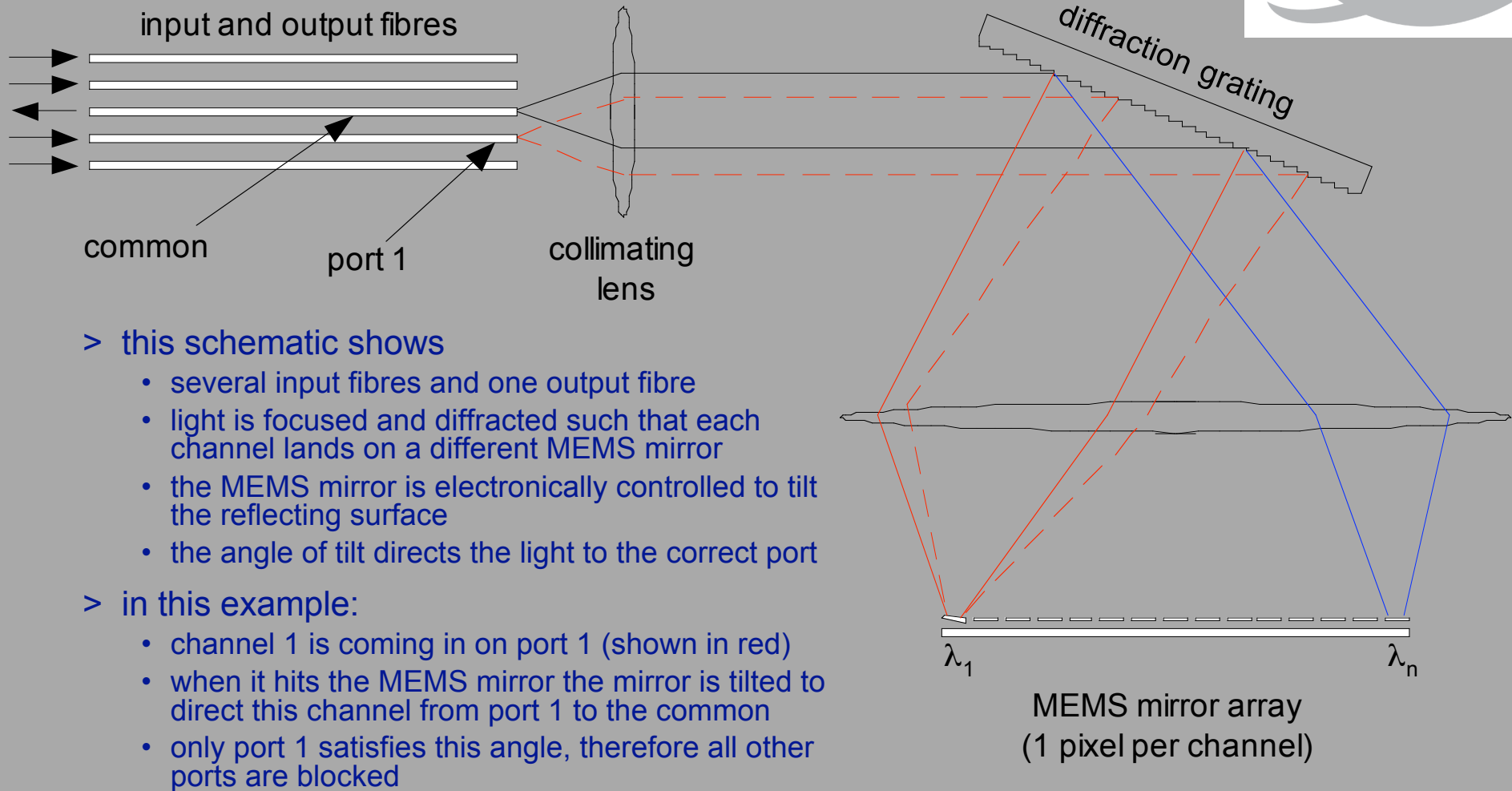
40/100 G.



# StarPlane DWDM backplane



# Module Operation



> this schematic shows

- several input fibres and one output fibre
- light is focused and diffracted such that each channel lands on a different MEMS mirror
- the MEMS mirror is electronically controlled to tilt the reflecting surface
- the angle of tilt directs the light to the correct port

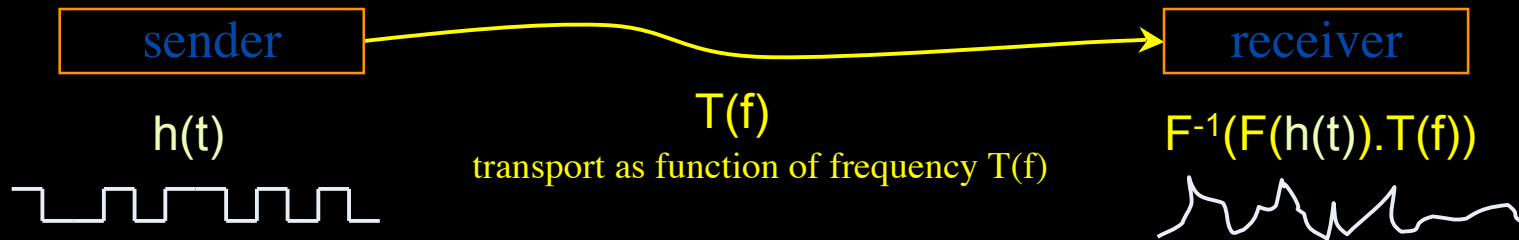
> in this example:

- channel 1 is coming in on port 1 (shown in red)
- when it hits the MEMS mirror the mirror is tilted to direct this channel from port 1 to the common
- only port 1 satisfies this angle, therefore all other ports are blocked

MEMS mirror array  
(1 pixel per channel)

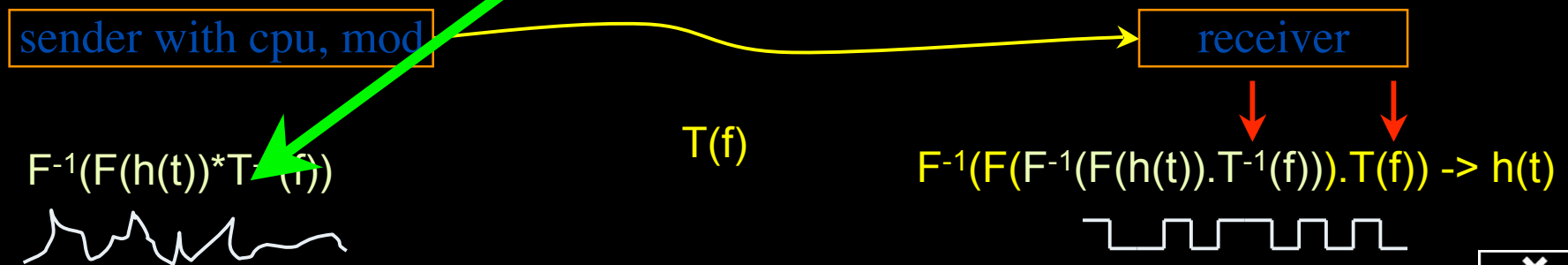
# Dispersion compensating modem: eDCO from NORTEL

(Try to Google eDCO :-)



Solution in 5 easy steps for dummy's :

1. try to figure out  $T(f)$  by trial and error
2. invert  $T(f) \rightarrow T^{-1}(f)$
3. computationally multiply  $T^{-1}(f)$  with Fourier transform of bit pattern to send
4. inverse Fourier transform the result from frequency to time space
5. modulate laser with resulting  $h'(t) = F^{-1}(F(h(t)).T^{-1}(f))$



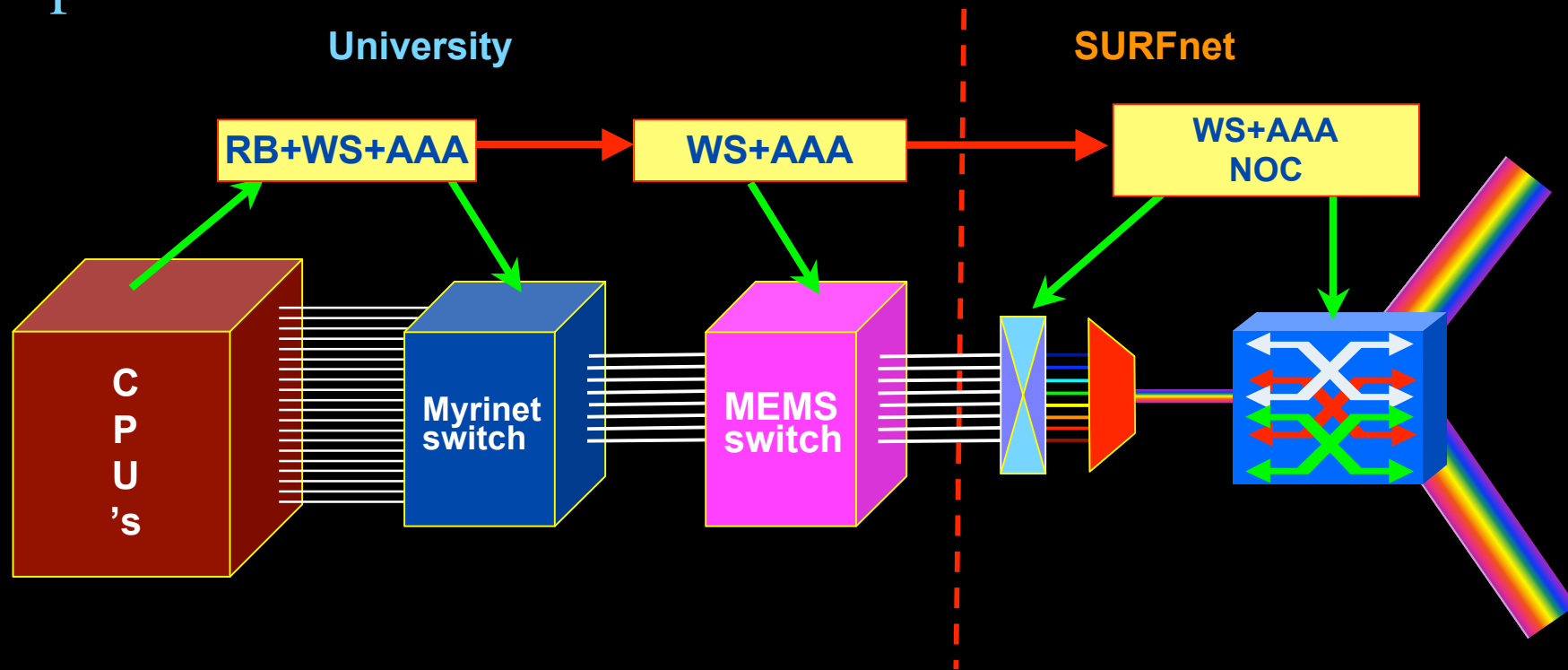
(ps. due to power  $\sim$  square E the signal to send **looks** like uncompensated received but is not)



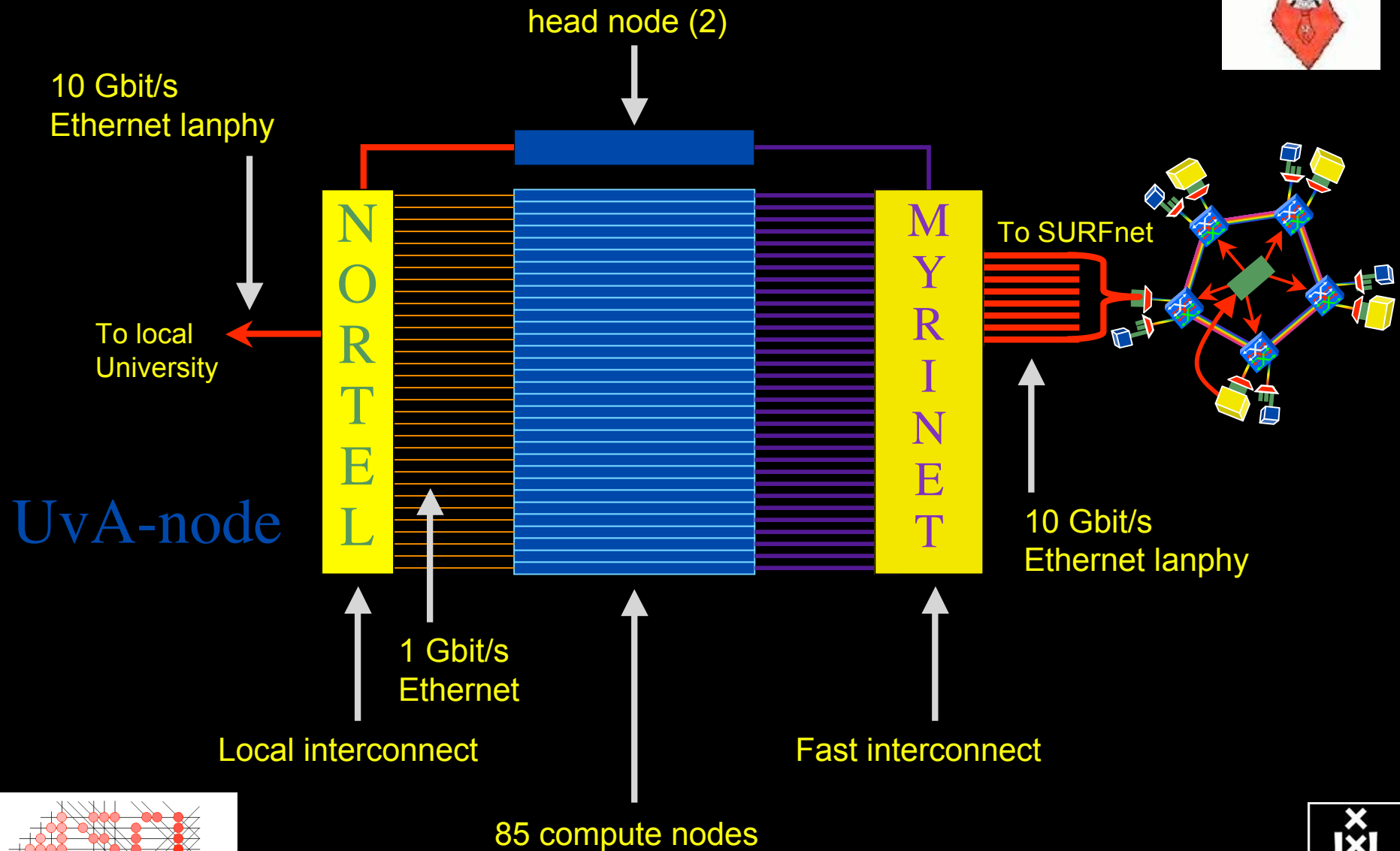


# The challenge for sub-second switching

- bringing up/down a  $\lambda$  takes minutes
  - this was fast in the era of old time signaling (phone/fax)
  - $\lambda \rightarrow \lambda$  influence (Amplifiers, non linear effects)
  - however minutes is historically grown, 5 nines, up for years
  - working with Nortel to get setup time significantly down
- plan B:



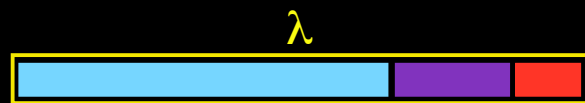
# DAS-3 Cluster Architecture





# QOS in a non destructive way!

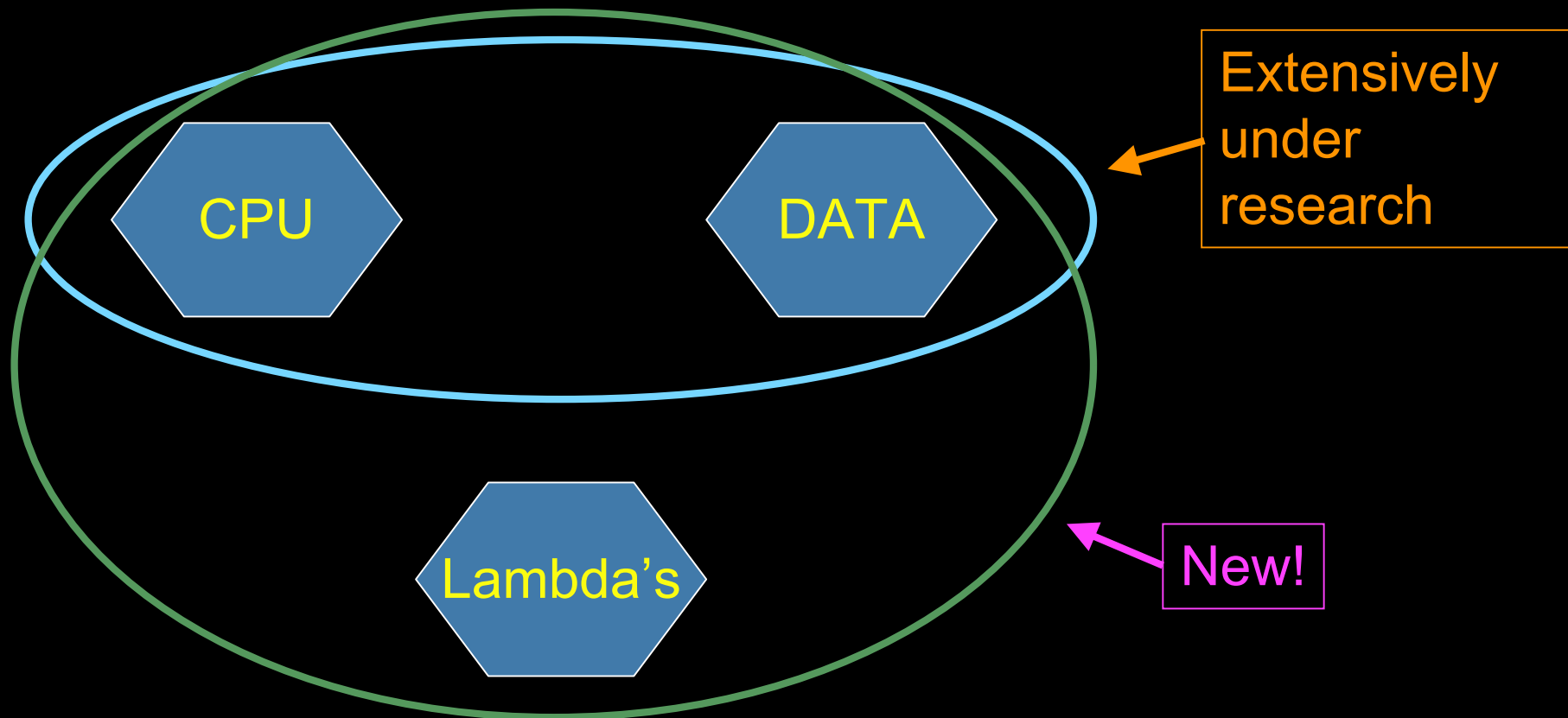
- Destructive QOS:
  - have a link or  $\lambda$
  - set part of it aside for a lucky few under higher priority
  - rest gets less service



- Constructive QOS:
  - have a  $\lambda$
  - add other  $\lambda$ 's as needed on separate colors
  - move the lucky ones over there
  - rest gets also a bit happier!



# GRID Co-scheduling problem space



The StarPlane vision is to give flexibility directly to the applications by allowing them to choose the logical topology in real time, ultimately with sub-second lambda switching times on part of the SURFnet6 infrastructure.



# What makes StarPlane fly?

- Wavelength Selective Switches
  - for the “low cost” photonics
- Sandbox by confining StarPlane to one band
  - for experimenting on a production network
- Optimization of the controls to turn on/off a Lambda
  - direct access to part of the controls at the NOC
- electronic Dynamically Compensating Optics (eDCO)
  - to compensate for changing lengths of the path
- traffic engineering
  - to create the OPN topologies needed by the applications
- Open Source GMPLS
  - to facilitate policy enabled cross domain signalling





Status: [Overview](#) [Throughput](#) [Scroll time](#) [Last 7 days](#)  
 Repeat: [Load](#) [Ping](#) [UDP](#) [Plot](#) [<<](#) [<<](#) [>>](#) [23:30:01](#) [30 min.](#)

## Overview Net Tests between DAS-3 Hosts

- [Authorise here](#) to store the current table settings in your cookies file.
- See the [getting started](#) introduction or the [user guide](#) for a description of the table below.
- Some [observations](#) about the package and the required bandwidth.

Select ping value: [min](#), [avg](#), [max](#), [all](#), [host](#).

Select UDP value: [rate](#), [host](#).

### DAS-3 Net Test Results

Date: 15/05/2007

Time: 23:30:01

#### Load

VU-083	VU-085	LIACS-125	LIACS-127	UvA-236	UvA-239
0	0	0.087	0	0.013	0.05

#### Ping Min (ms)

(see in column)

	VU-083	VU-085	LIACS-125	LIACS-127	UvA-236	UvA-239
VU-083	---	---			0.695	
VU-085	---	---	1.380			
LIACS-125		1.380	---	---		
LIACS-127			---	---		1.220
UvA-236	0.695				---	---
UvA-239				1.230	---	---

#### Throughput [Mbit/s]

(see in column)

	VU-083	VU-085	LIACS-125	LIACS-127	UvA-236	UvA-239
--	--------	--------	-----------	-----------	---------	---------

Status: **Overview** Throughput Scroll time: Last 7 days  
 Report: Load Ping UDP Plot <<< << >> 23.10.01 30 min.

## Throughput [Mbit/s]

(row vs column)

	YU-083	YU-085	LIACS-125	LIACS-127	UvA-236	UvA-239
YU-083	---	---			4267.46	
YU-085	---	---	4674.64			
LIACS-125		5143.93	---	---		
LIACS-127			---	---		4284.89
UvA-236	3829.06				---	---
UvA-239				4445.64	---	---

## UDP Data Rate [Mbit/s]

(row vs column)

	YU-083	YU-085	LIACS-125	LIACS-127	UvA-236	UvA-239
YU-083	---	---			6440.39	
YU-085	---	---	6549.51			
LIACS-125		6548.28	---	---		
LIACS-127			---	---		6528.95
UvA-236	6554.22				---	---
UvA-239				6551.25	---	---

The load, roundtrip, throughput and UDP data series are each scaled with their private color distributions as is displayed below:

load	0	0.25	0.5	0.75	1	1.25	1.5	1.75	2
ping min [ms]	0.695	0.781	0.866	0.952	1.037	1.123	1.209	1.294	1.38
throughput [Mbit/s]	3829.06	3993.419	4157.778	4322.136	4486.495	4650.854	4815.213	4979.571	5143.93

# Heterogeneous clusters

(# of unused ports)

	LU	TUD	UvA-VLE	UvA-MN	VU	TOTALS
<b>Head</b>						
* storage	10TB	5TB	2TB	2TB	10TB	29TB
* CPU	2x2.4GHz DC	2x2.4GHz DC	2x2.2GHz DC	2x2.2GHz DC	2x2.4GHz DC	46.4 GHz
* memory	16GB	16GB	8GB	16GB	8GB	64GB
* Myri 10G	1		1	1	1	40 Gb/s
* 10GE	1	1	1	1	1	50 Gb/s
<b>Compute</b>	32	68	40 (+1)	46	85	271
* storage	400GB	250GB	250GB	2x250GB	250GB	84 TB
* CPU	2x2.6GHz	2x2.4GHz	2x2.2GHz DC	2x2.4GHz	2x2.4GHz DC	1.9 THz
* memory	4GB	4GB	4GB	4GB	4GB	1048 GB
* Myri 10G	1		1	1	1	2030 Gb/s
<b>Myrinet</b>						
* 10G ports	33 (7)		41	47	86 (2)	2070 Gb/s
* 10GE ports	8		8	8	8	320 Gb/s
<b>Nortel</b>						
* 1GE ports	32 (16)	136 (8)	40 (8)	46 (2)	85 (11)	339 Gb/s
* 10GE ports	1 (1)	9 (3)	2	2	1 (1)	

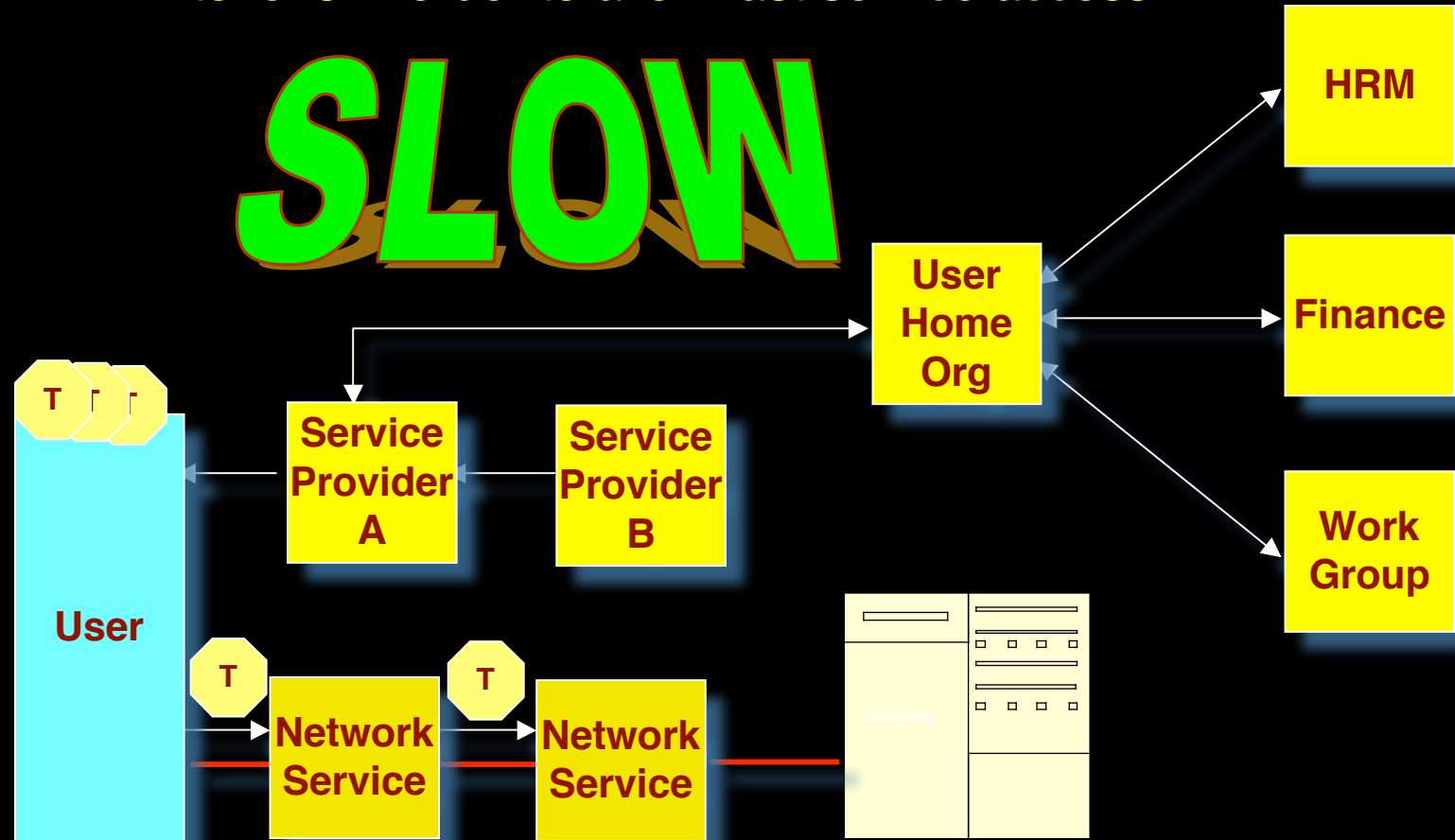


# Power is a big issue

- UvA cluster uses (max) 30 kWh
- 1 kWh ~ 0.1 €
- per year -> 26 k€/y
- add cooling 50% -> 39 k€/y
- Emergency power system -> 50 k€/y
- per rack 10 kWh is now normal
- **YOU BURN ABOUT HALF THE CLUSTER OVER ITS LIFETIME!**
- Terminating a 10 Gb/s wave costs about 200 W
- Entire loaded fiber -> 16 kW
- Wavelength Selective Switch : few W!

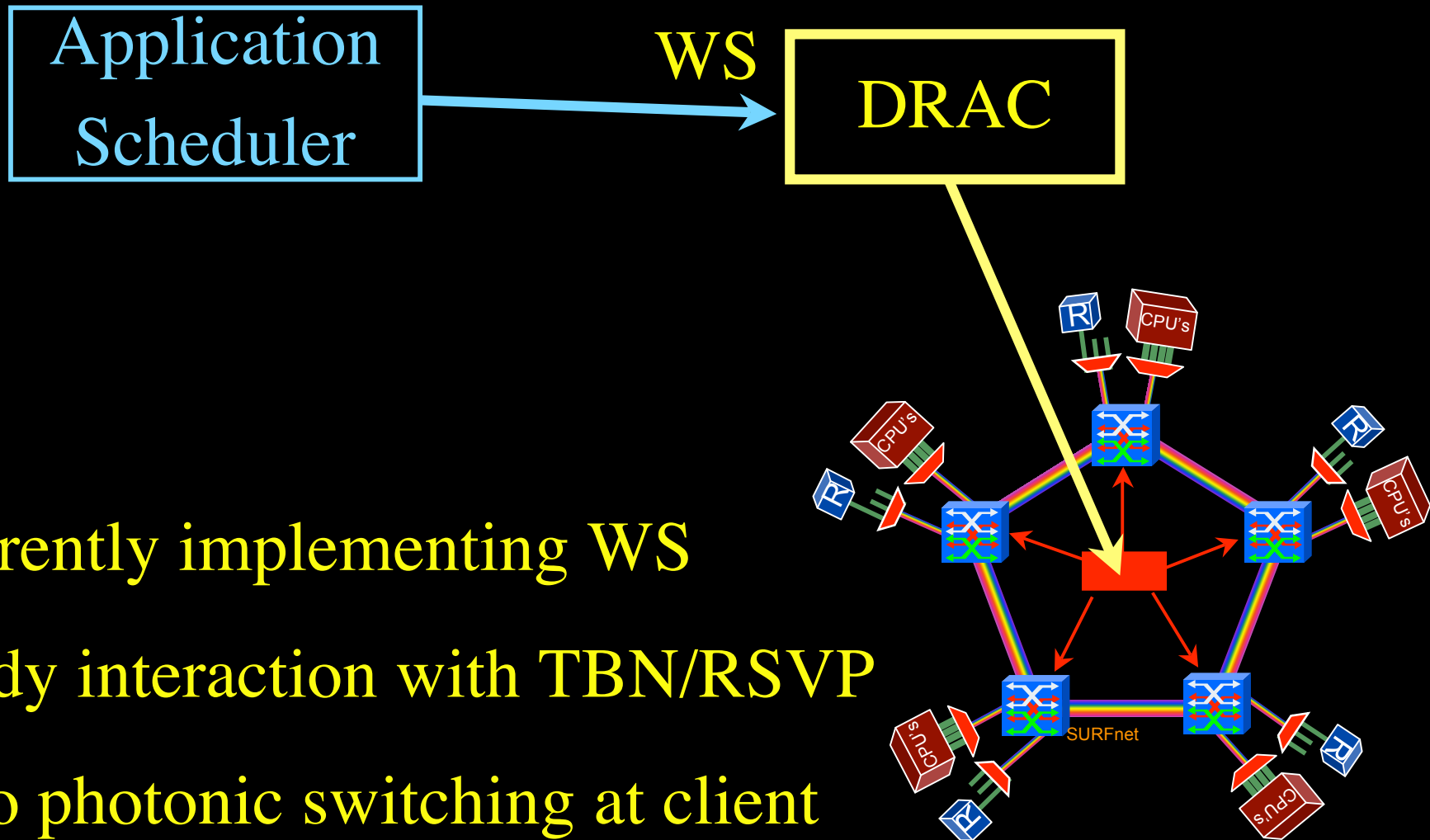


Use AAA concept to split (time consuming) service authorization process from service access using secure tokens in order to allow fast service access.



**FAST**

# Control Plane

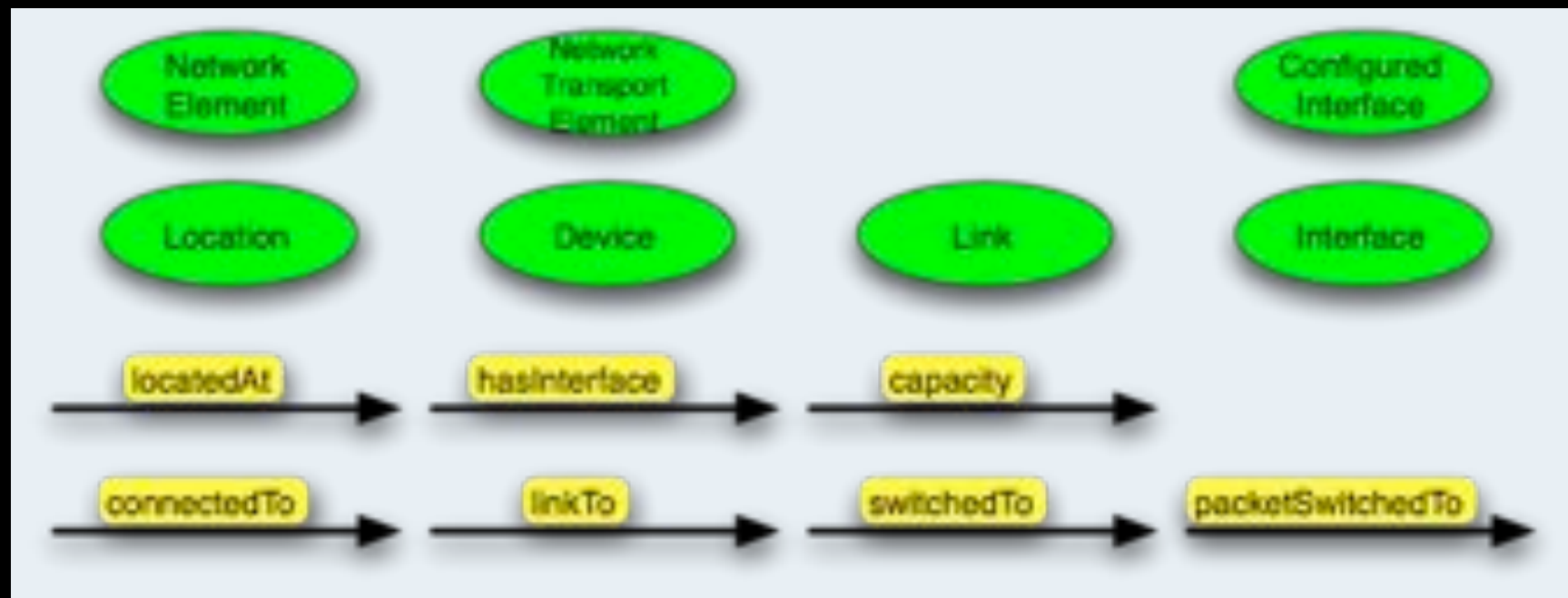


- currently implementing WS
- study interaction with TBN/RSVP
- also photonic switching at client

# StarPlane and NDL

While on topologies. SNE group is working on NDL - Network Description Language.

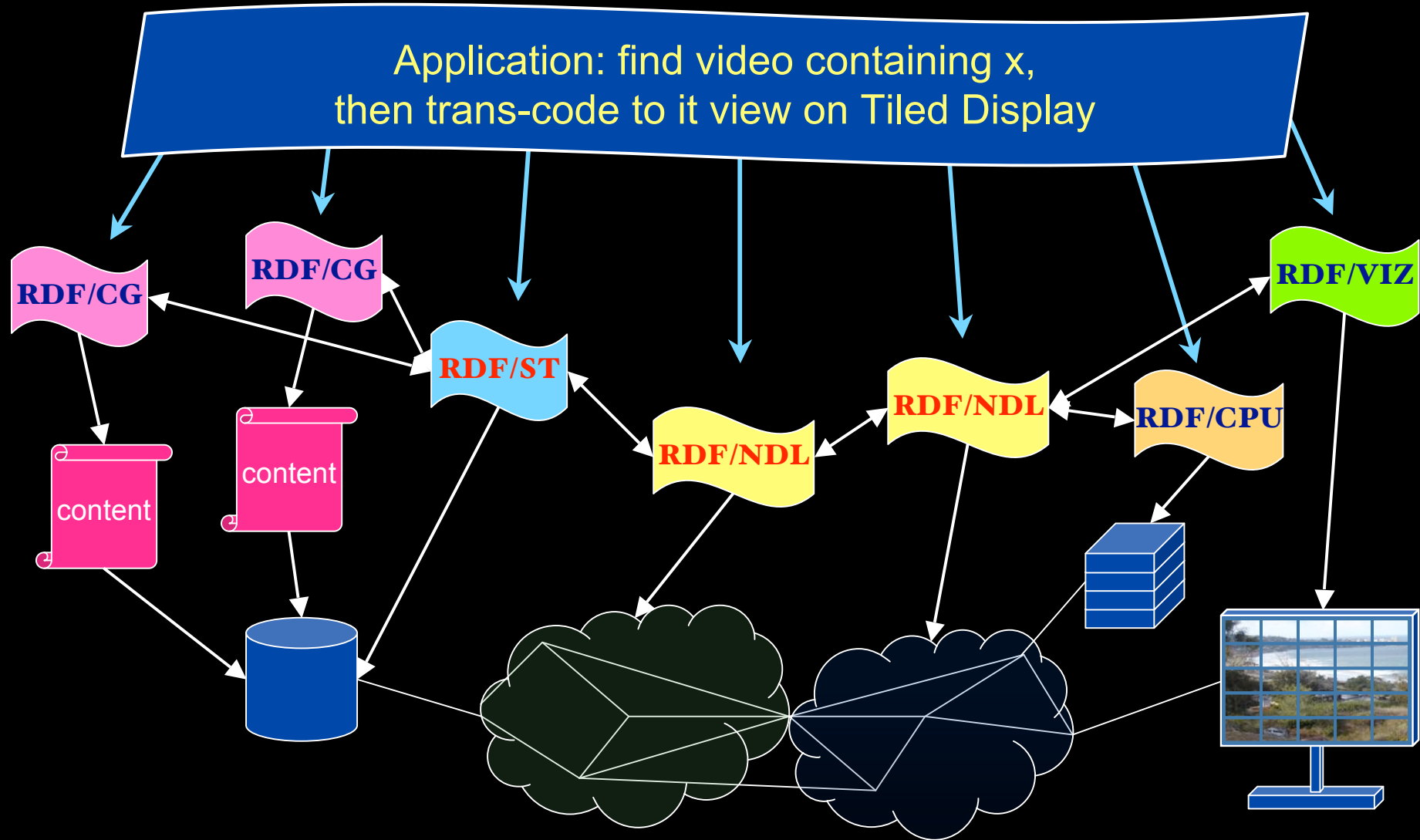
NDL is an RDF data model, based on idea of Semantic Web, for network topology descriptions.



In StarPlane we are researching use of NDL for topology exchange and topology requests from clients.

ref: Talk from Paola Grosso on NDL/RDF

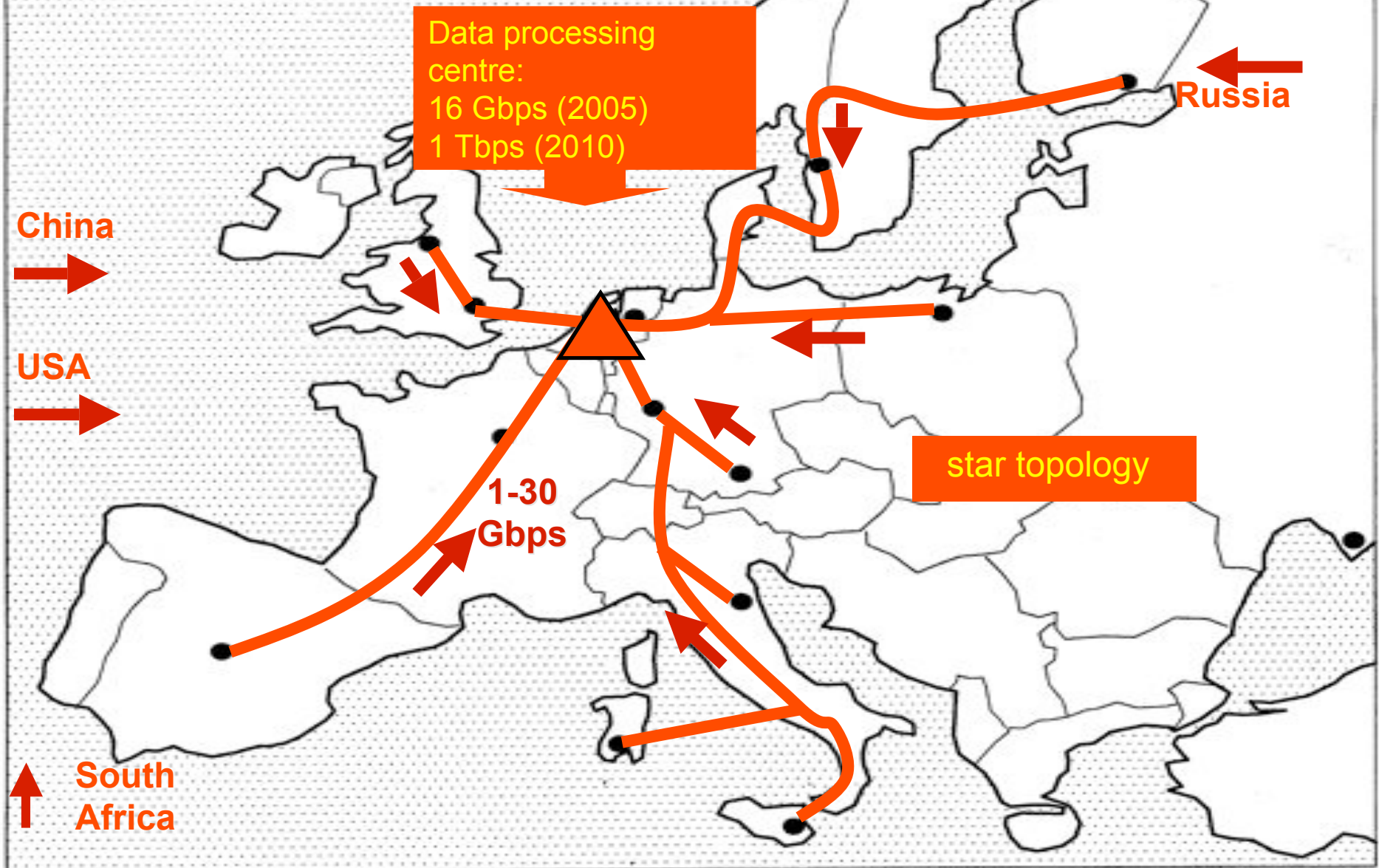
# RDF describing Infrastructure



# StarPlane Applications

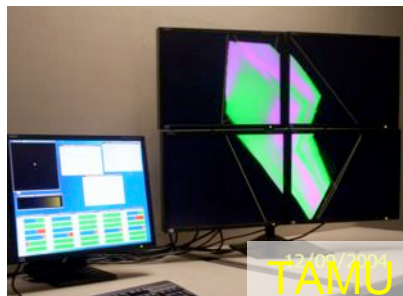
- Large 'stand-alone' file transfers
  - User-driven file transfers
  - Nightly backups
  - Transfer of medical data files (MRI)
- Large file (speedier) Stage-in/Stage-out
  - MEG modeling (Magneto encephalography)
  - Analysis of video data
- Application with static bandwidth requirements
  - Distributed game-tree search
  - Remote data access for analysis of video data
  - Remote visualization
- Applications with dynamic bandwidth requirements
  - Remote data access for MEG modeling
  - SCARI

# eEVN: European VLBI Network



This slide courtesy of Richard Schilizzi <[schilizzi@jive.nl](mailto:schilizzi@jive.nl)>

# US and International OptIPortal Sites





# CineGrid@SARA



# Tera-Thinking

- What constitutes a Tb/s network?
- CALIT2 has 8000 Gigabit drops ?->? Terabit Lan?
- look at 80 core Intel processor
  - cut it in two, left and right communicate 8 TB/s
- think back to teraflop computing!
  - MPI makes it a teraflop machine
- massive parallel channels in hosts, NIC's
- TeraApps programming model supported by
  - TFlops -> MPI / Globus
  - TBytes -> OGSA/DAIS
  - TPixels -> SAGE
  - TSensors -> LOFAR, LHC, LOOKING, CineGrid, ...
  - Tbit/s -> ?



# Questions ?

Thanks to:

SURFnet, NWO (grant 643.000.504), NORTEL

Team: Li Xu, Jason Maasen, JP Velders, Leon Gommans, Paola Grosso, Herbert Bos, Henri Bal

Special thanks to Kees Neggers.

