# Lambda-Grid developments
## History - Present - Future

## Cees de Laat

# SURFnet

# EU

## BSIK

### NWO
### University of Amsterdam

TNO
NCF

# Contents

# LHC Data Grid Hierarchy
## CMS as example, Atlas is similar

~PByte/sec

**Online System**

~100 MBytes/sec

*Tier 0 +1*

event simulation

HPSS

**event reconstruction**

human=2m

CMS detector: 15m X 15m X 22m

12,500 tons, $700M.

~2.5 Gbits/sec

*Tier 1*

**Italian Regional Center** — HPSS

**German Regional Center** — HPSS

**NIKHEF Dutch Regional Center** — HPSS

**FermiLab, USA Regional Center** — HPSS

• • •

analysis

~0.6-2.5 Gbps

Tier2 Center   2 Center   nter   Center   Center   *Tier 2*

~0.6-2.5 Gbps

*Tier 3*

**Institute ~0.25TIPS**   tute   stitute   Institute

Physics data cache

100 - 1000 Mbits/sec

*Tier 4*

Courtesy Harvey Newman, CalTech and CERN

Workstations

CERN/CMS data goes to 6-8 Tier 1 regional centers, and from each of these to 6-10 Tier 2 centers.

Physicists work on analysis "channels" at 135 institutes. Each institute has ~10 physicists working on one or more channels.

2000 physicists in 31 countries are involved in this 20-year experiment in which DOE is a major player.

# Data intensive scientific computation through global networks

Nuclear experiments

Belle Experiments

Nobeyama Radio Observatory (VLBI)

X-ray astronomy Satellite ASUKA

Data Reservoir

Very High-speed Network

Digital Sky Survey

Data Reservoir

Distributed Shared files

SUBARU Telescope

Local Accesses

Data Reservoir

Grape6

Data analysis at University of Tokyo

# Sensor Grids

## eVLBI



longer term VLBI is easily capable of generatin
be. The sensitivity of the VLBI array scales with
dth (=data-rate) and there is a strong push to mo
dths. Rates of 8Gb/s or more are entirely feasible.
o under development. It is expected that parallel
ed correlator will remain the most efficient approach
olves dist
multi-gig
relator and
factor.



Westerbork Synthesis Radio Telescope -
Netherlands

~ 40 Tbit/s
www.lofar.org

# Towards Hybrid Networking!

- Costs of photonic equipment 10% of switching 10 % of full routing
  - for same throughput!
  - Photonic vs Optical (optical used for SONET, etc, 10-50 k$/port)
  - DWDM lasers for long reach expensive, 10-50 k$
- Bottom line: look for a hybrid architecture which serves all classes in a cost effective way
  - map A -> L3 , B -> L2 , C -> L1 and L2
- Give each packet in the network the service it needs, but no more !

L1 ≈ 2-3 k$/port
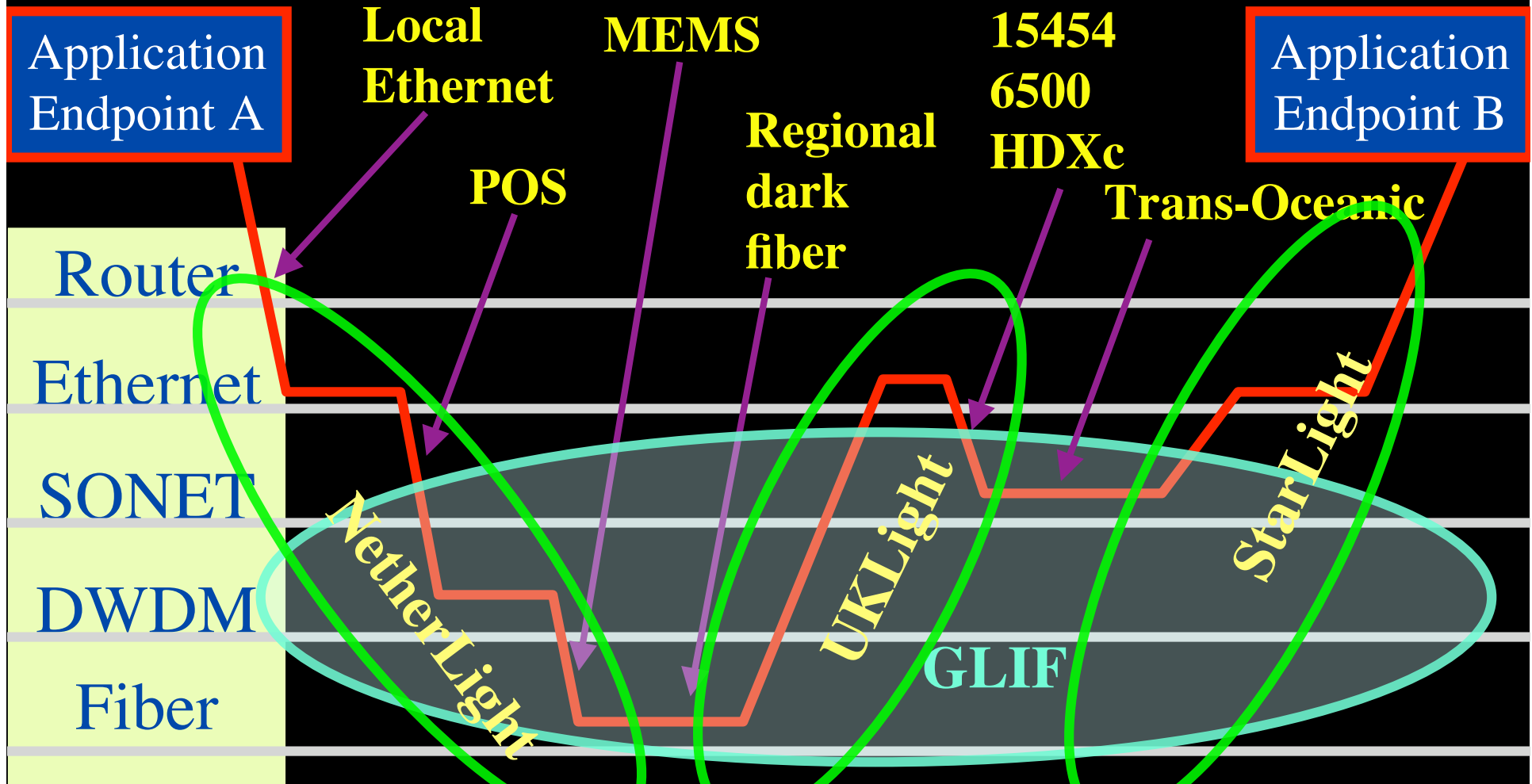
L2 ≈ 5-8 k$/port

L3 ≈ 75+ k$/port

# Trends

- We have made baby-steps on the path to optical networking
  - Still many mails and phone calls
- See several trends:
  - lambda's get fatter and cheaper
  - photonic technology cheap per bandwidth
  - embedded computation capacity increasing
  - latency and high bandwidth congestion avoidance conflict
  - ethernet is getting circuit properties (PBT)
  - applications need more and more predictable behaviour

# How low can you go?

Application Endpoint A

Application Endpoint B

Local Ethernet

POS

MEMS

Regional dark fiber

15454 6500 HDXc

Trans-Oceanic

Router

Ethernet

SONET

DWDM

Fiber

NetherLight

UKLight

StarLight

GLIF

**The playfield => GLIF**

In The Netherlands SURFnet connects between 180:
  - universities;
  - academic hospitals;
  - most polytechnics;
  - research centers.
with an indirect ~750K user base

Red crosses = StarPlane

~ 6000 km

scale comparable to railway system

Common Photonic Layer (CPL) in SURFnet6

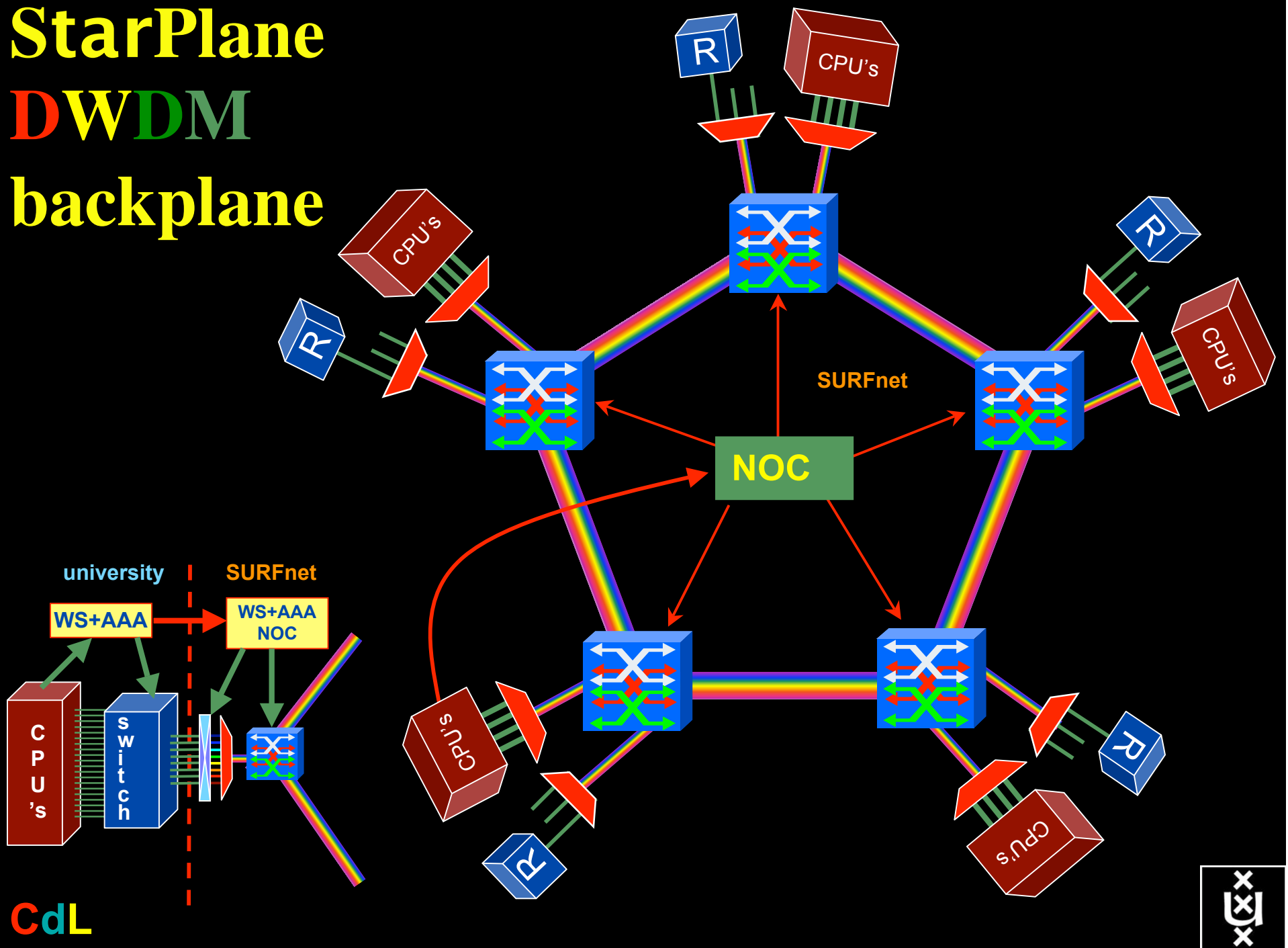supports up to 72 Lambda's of 10 G each 40 G soon.

# Contents

1. The need for hybrid networking

2. StarPlane; a grid controlled photonic network

3. RDF/Network Description Language

4. Tera-networking

5. Programmable networks

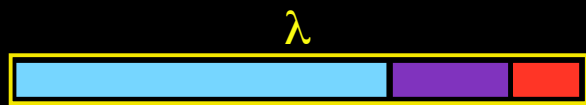# StarPlane
# DWDM
# backplane



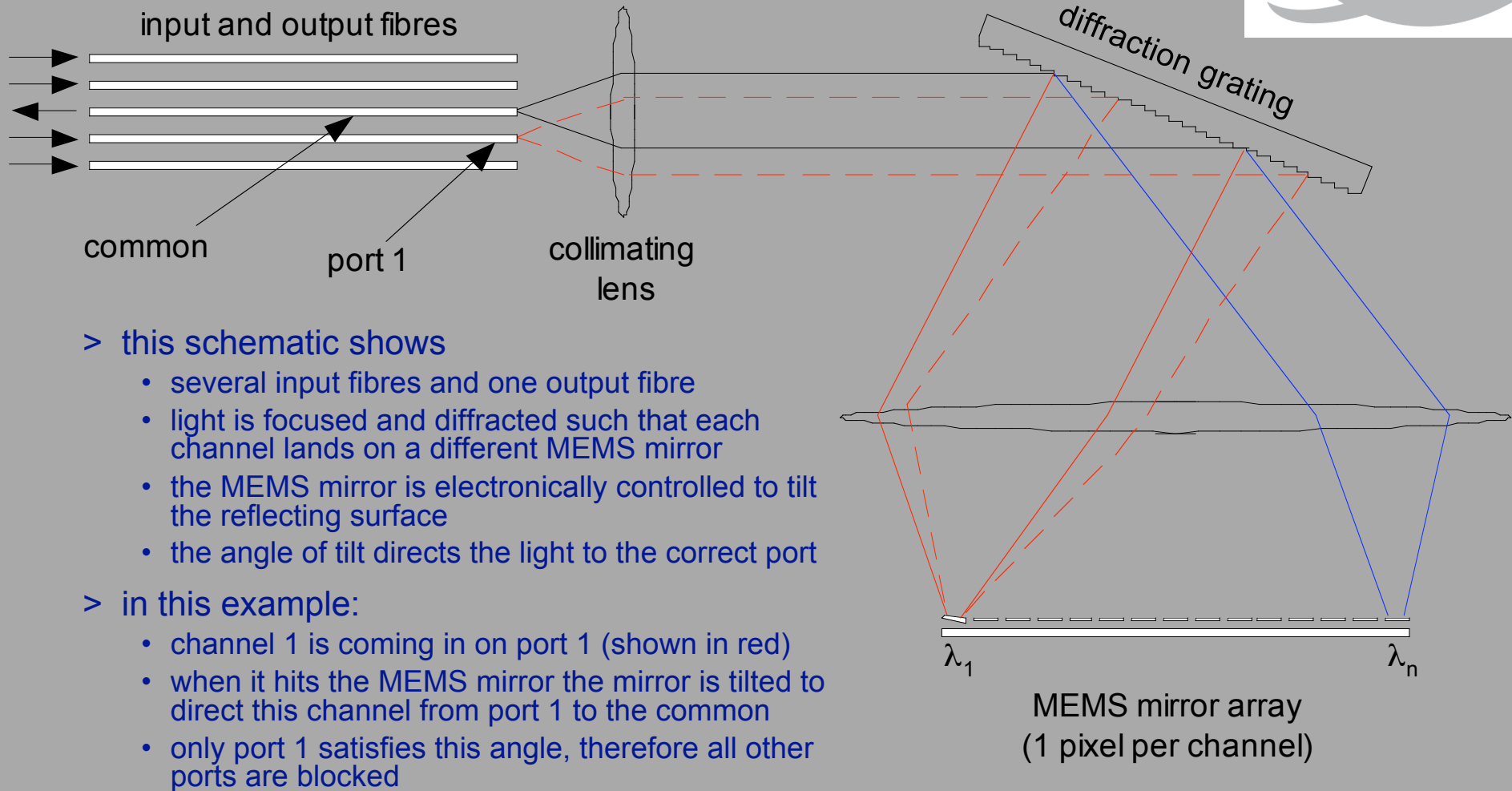**CdL**

# 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

$\lambda$

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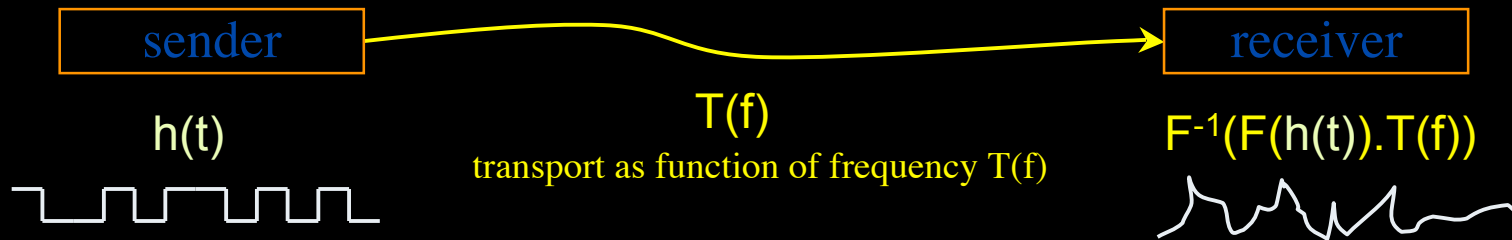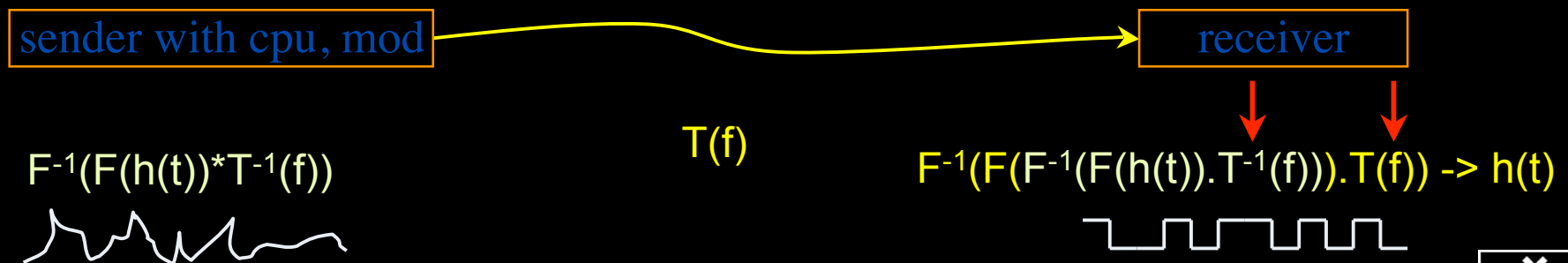  - move the lucky ones over there
  - rest gets also a bit happier!

$\lambda$          $\lambda$          $\lambda$

# Module Operation

input and output fibres

common

port 1

collimating lens

diffraction grating

> 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

$\lambda_1$

$\lambda_n$

MEMS mirror array
(1 pixel per channel)

# Dispersion compensating modem: eDCO from NORTEL
## (Try to Google eDCO :-)

sender $\longrightarrow$ receiver

h(t)

T(f)
transport as function of frequency T(f)

$F^{-1}(F(h(t)).T(f))$

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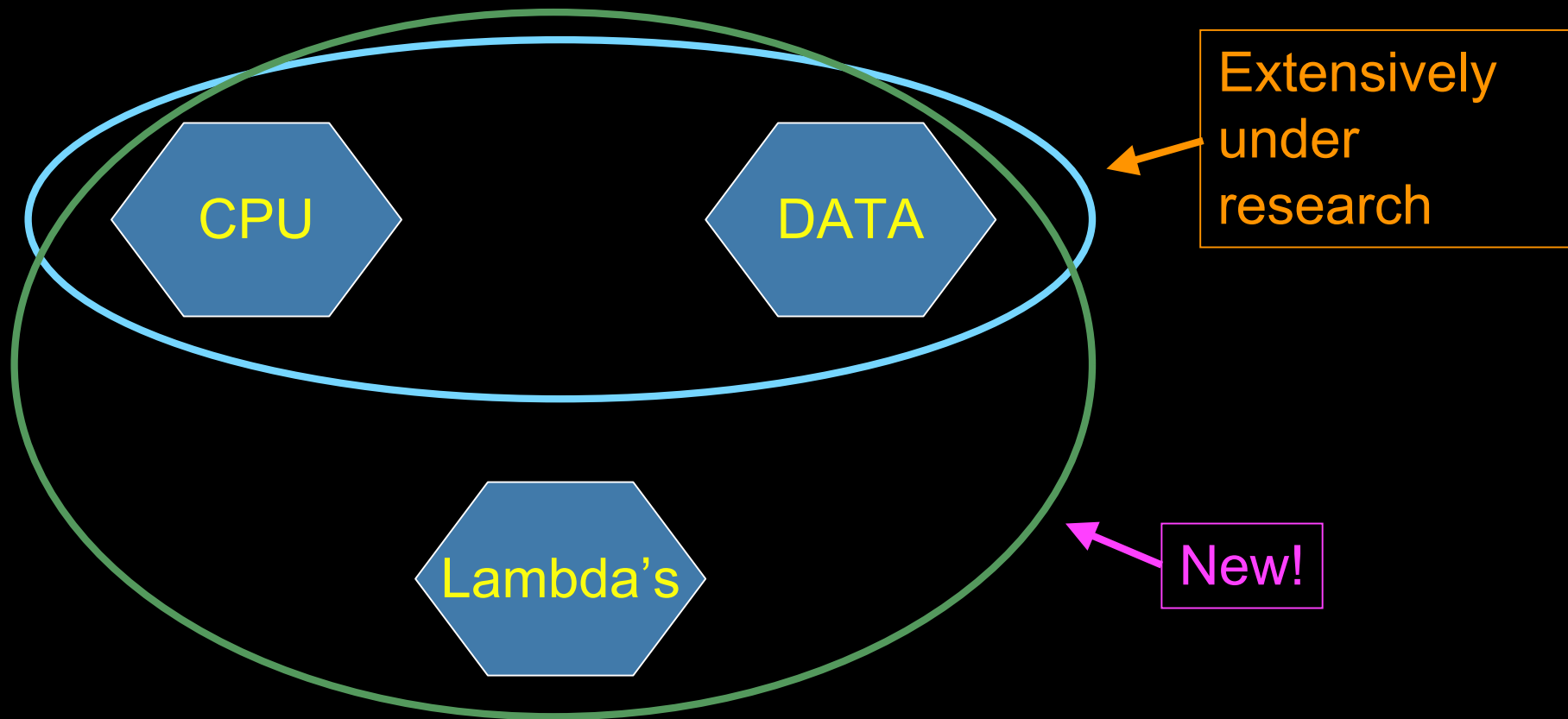
1. try to figure out T(f) by trial and error
2. invert T(f) -> $T^{-1}(f)$
3. computationally multiply $T^{-1}(f)$ with Fourier transform of bit pattern to send
4. inverse Fourier transform the result from frequency to time space
5. modulate laser with resulting h'(t) = $F^{-1}(F(h(t)).T^{-1}(f))$

sender with cpu, mod $\longrightarrow$ receiver

$F^{-1}(F(h(t))*T^{-1}(f))$

T(f)

$F^{-1}(F(F^{-1}(F(h(t)).T^{-1}(f))).T(f))$ -> h(t)

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

# GRID Co-scheduling problem space



CPU

DATA

Lambda's

Extensively under research

New!

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.

Net Tests between DAS-3 Hosts

Ping AB [ms] from / to node125.das3.liacs.nl (LIACS-125)

Skipped tests: UvA-236-M, UvA-239-M

| Date | Time | >> VU-083 | << VU-083 | >> VU-085 | << VU-085 | >> LIACS-127 | << LIACS-127 | >> UvA-236 | << UvA-236 | >> UvA-239 | << UvA-239 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 31/05/2007 | 12:30:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.383 / 1.420 | | | | | | |
| 31/05/2007 | 12:00:01 | | | 1.380 / 1.383 / 1.410 | 1.380 / 1.384 / 1.450 | | | | | | |
| 31/05/2007 | 11:30:01 | | | 1.380 / 1.383 / 1.410 | 1.380 / 1.382 / 1.390 | | | | | | |
| 31/05/2007 | 11:00:02 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.382 / 1.400 | | | | | | |
| 31/05/2007 | 10:30:01 | | | 1.380 / 1.383 / 1.390 | 1.380 / 1.382 / 1.390 | | | | | | |
| 31/05/2007 | 10:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.383 / 1.410 | | | | | | |
| 31/05/2007 | 09:30:01 | | | 1.380 / 1.384 / 1.410 | 1.380 / 1.382 / 1.400 | | | | | | |
| 31/05/2007 | 09:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.383 / 1.400 | | | | | | |
| 31/05/2007 | 08:30:02 | | | 1.380 / 1.383 / 1.410 | 1.380 / 1.382 / 1.400 | | | | | | |
| 31/05/2007 | 08:00:01 | | | 1.380 / 1.383 / 1.410 | 1.380 / 1.383 / 1.410 | | | | | | |
| 31/05/2007 | 07:30:02 | | | 1.380 / 1.382 / 1.390 | 1.380 / 1.381 / 1.390 | | | | | | |
| 31/05/2007 | 07:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.383 / 1.400 | | | | | | |
| 31/05/2007 | 06:30:01 | | | 1.380 / 1.383 / 1.410 | 1.380 / 1.382 / 1.390 | | | | | | |
| 31/05/2007 | 06:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.382 / 1.420 | | | | | | |
| 31/05/2007 | 05:30:01 | | | 1.380 / 1.382 / 1.400 | 1.380 / 1.382 / 1.410 | | | | | | |
| 31/05/2007 | 05:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.382 / 1.390 | | | | | | |
| 31/05/2007 | 04:30:01 | | | 1.380 / 1.381 / 1.390 | 1.380 / 1.381 / 1.390 | | | | | | |
| 31/05/2007 | 04:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.384 / 1.410 | | | | | | |
| 31/05/2007 | 03:30:02 | | | 1.380 / 1.384 / 1.410 | 1.380 / 1.382 / 1.400 | | | | | | |
| 31/05/2007 | 03:00:02 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.382 / 1.400 | | | | | | |
| 31/05/2007 | 02:30:01 | | | 1.380 / 1.382 / 1.400 | 1.380 / 1.382 / 1.400 | | | | | | |
| 31/05/2007 | 02:00:01 | | | 1.380 / 1.383 / 1.410 | 1.380 / 1.384 / 1.410 | | | | | | |
| 31/05/2007 | 01:30:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.382 / 1.390 | | | | | | |
| 31/05/2007 | 01:00:01 | | | 1.380 / 1.382 / 1.410 | 1.380 / 1.383 / 1.400 | | | | | | |

**Very constant and predictable!**

# Contents

# Architecture SC06

# Network Description Language

- From semantic Web / Resource Description Framework.
- The RDF uses XML as an interchange syntax.
- Data is described by triplets:

```
  Subject  ──Predicate──▶  Object
```

| Location | Device | Interface | Link |
|----------|--------|-----------|------|

name ──▶        description ──▶        locatedAt ──▶        hasInterface ──▶

connectedTo ──▶        capacity ──▶        encodingType ──▶        encodingLabel ──▶

# The Modelling Process

Network Elements → Functional Elements → Syntax



```
<ndl:Device rdf:about="#Force10">
  <ndl:hasInterface rdf:resource=
    "#Force10:te6/0"/>
</ndl:Device>
<ndl:Interface rdf:about="#Force10:te6/0">
  <rdfs:label>te6/0</rdfs:label>
  <ndl:capacity>1.25E6</ndl:capacity>
  <ndlconf:multiplex>
    <ndlcap:adaptation rdf:resource=
      "#Tagged-Ethernet-in-Ethernet"/>
    <ndlconf:serverPropertyValue
      rdf:resource="#MTU-1500byte"/>
  </ndlconf:multiplex>
  <ndlconf:hasChannel>
    <ndlconf:Channel rdf:about=
      "#Force10:te6/0:vlan4">
      <ndleth:hasVlan>4</ndleth:hasVlan>
      <ndlconf:switchedTo rdf:resource=
        "#Force10:gi5/1:vlan7"/>
    </ndlconf:Channel>
  </ndlconf:hasChannel>
</ndl:Interface>
```
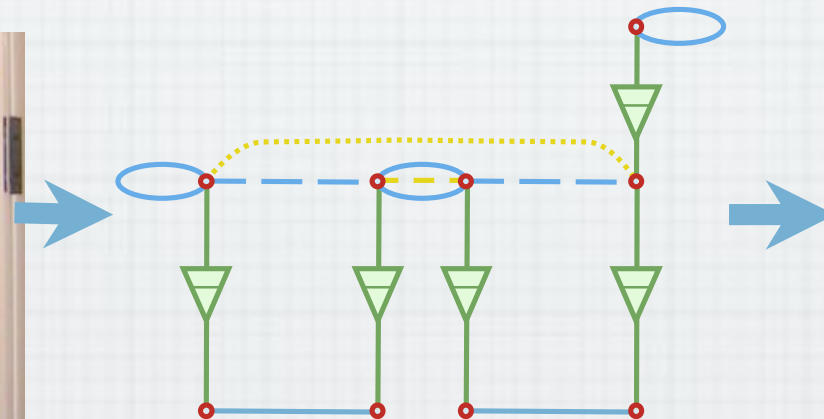
# NetherLight in RDF

```xml
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
    xmlns:ndl="http://www.science.uva.nl/research/air/ndl#">
<!-- Description of Netherlight -->
<ndl:Location rdf:about="#Netherlight">
    <ndl:name>Netherlight Optical Exchange</ndl:name>
</ndl:Location>
<!-- TDM3.amsterdam1.netherlight.net -->
<ndl:Device rdf:about="#tdm3.amsterdam1.netherlight.net">
    <ndl:name>tdm3.amsterdam1.netherlight.net</ndl:name>
    <ndl:locatedAt rdf:resource="#amsterdam1.netherlight.net"/>
    <ndl:hasInterface rdf:resource="#tdm3.amsterdam1.netherlight.net:501/1"/>
    <ndl:hasInterface rdf:resource="#tdm3.amsterdam1.netherlight.net:501/3"/>
    <ndl:hasInterface rdf:resource="#tdm3.amsterdam1.netherlight.net:501/4"/>
    <ndl:hasInterface rdf:resource="#tdm3.amsterdam1.netherlight.net:503/1"/>
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
    <ndl:hasInterface rdf:resourc
```

```xml
<!-- all the interfaces of TDM3.amsterdam1.netherlight.net -->

<ndl:Interface rdf:about="#tdm3.amsterdam1.netherlight.net:501/1">
            <ndl:name>tdm3.amsterdam1.netherlight.net:POS501/1</ndl:name>
            <ndl:connectedTo rdf:resource="#tdm4.amsterdam1.netherlight.net:5/1"/>
</ndl:Interface>
<ndl:Interface rdf:about="#tdm3.amsterdam1.netherlight.net:501/2">
            <ndl:name>tdm3.amsterdam1.netherlight.net:POS501/2</ndl:name>
            <ndl:connectedTo rdf:resource="#tdm1.amsterdam1.netherlight.net:12/1"/>
</ndl:Interface>
```

# NDL Generator and Validator



**Step 1 - Location**

Indicate the name and a short description of the network that is going to be described in NDL.

Name: Lighthouse    Description: SNE Lab

Provide also the latitude and the longitude of this location: this will aid the visualization programs.
Both latitude and longitude should use **floating point** notation.

Latitude: 52.3651    Longitude: 4.9527

**Step 2 - Devices**

Indicate the name of all the devices present in the network. If you need to describe more than 3 devices just "Add a Device"

Device: Rembrandt3
Device: Speculaas
Device:

Add a Device

## NDL for the GLIF - NDL Validator

NDL - Network Description Language - is an ontology for description of (hybrid) networks, aim provisioning. The GLIF collaboration makes use of NDL to describe each individual domain, maps.

This page will provide you with tools to validate an NDL file. We provide here two types of va

- Syntax validation
- Content validation

### Syntax validation

We can validate that the NDL file you generated is written following the latest NDL schema. Y will get back feedback on its validity.

Please paste your NDL file below:

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
xmlns:ndl="http://www.science.uva.nl/research/sne/ndl#"
xmlns:geo="http://www.w3.org/2003/01/geo/wgs84_pos#">

<!-- Description of foo-->
<ndl:Location rdf:about="#foo">
<ndl:name>bar</ndl:name>
<geo:lat>0</geo:lat>
<geo:long>0</geo:long>
</ndl:Location>

<!--Rem2-->
<ndl:Device rdf:about="#Rem2">
<ndl:name>Rem2</ndl:name>
    <ndl:locatedAt rdf:resource="#foo"/>
    <ndl:hasInterface rdf:resource="#Rem2:eth0"/>
</ndl:Device>

<!--GLIF-->
```

Submit

### Content validation

Often NDL files reference information contained in other files managed by others. Such as for example when an interface on a local device connects to an interface to a remote device. The content validator performs a few basic checks to see that the information contained in cross-referencing NDL files is consistent.
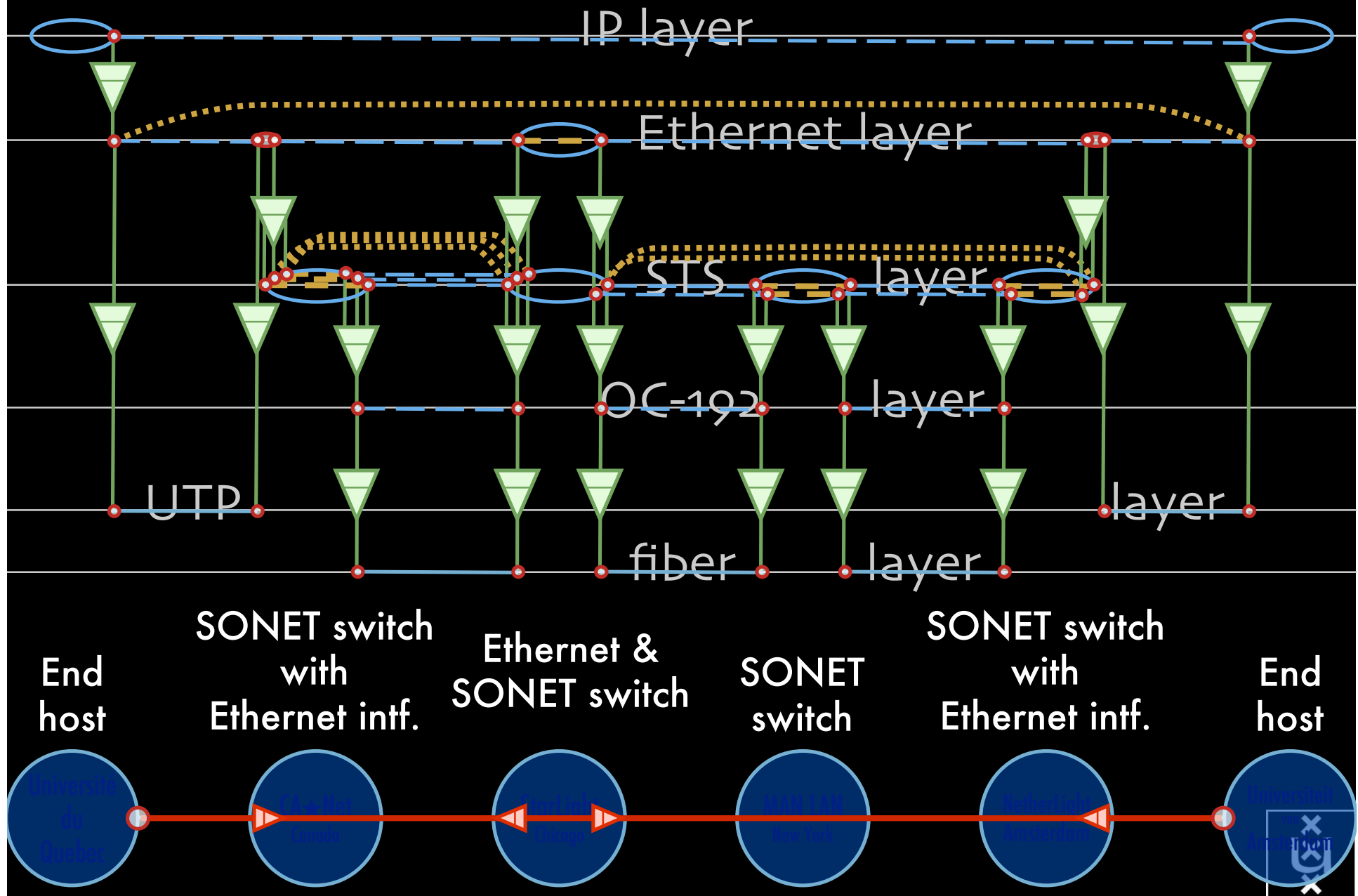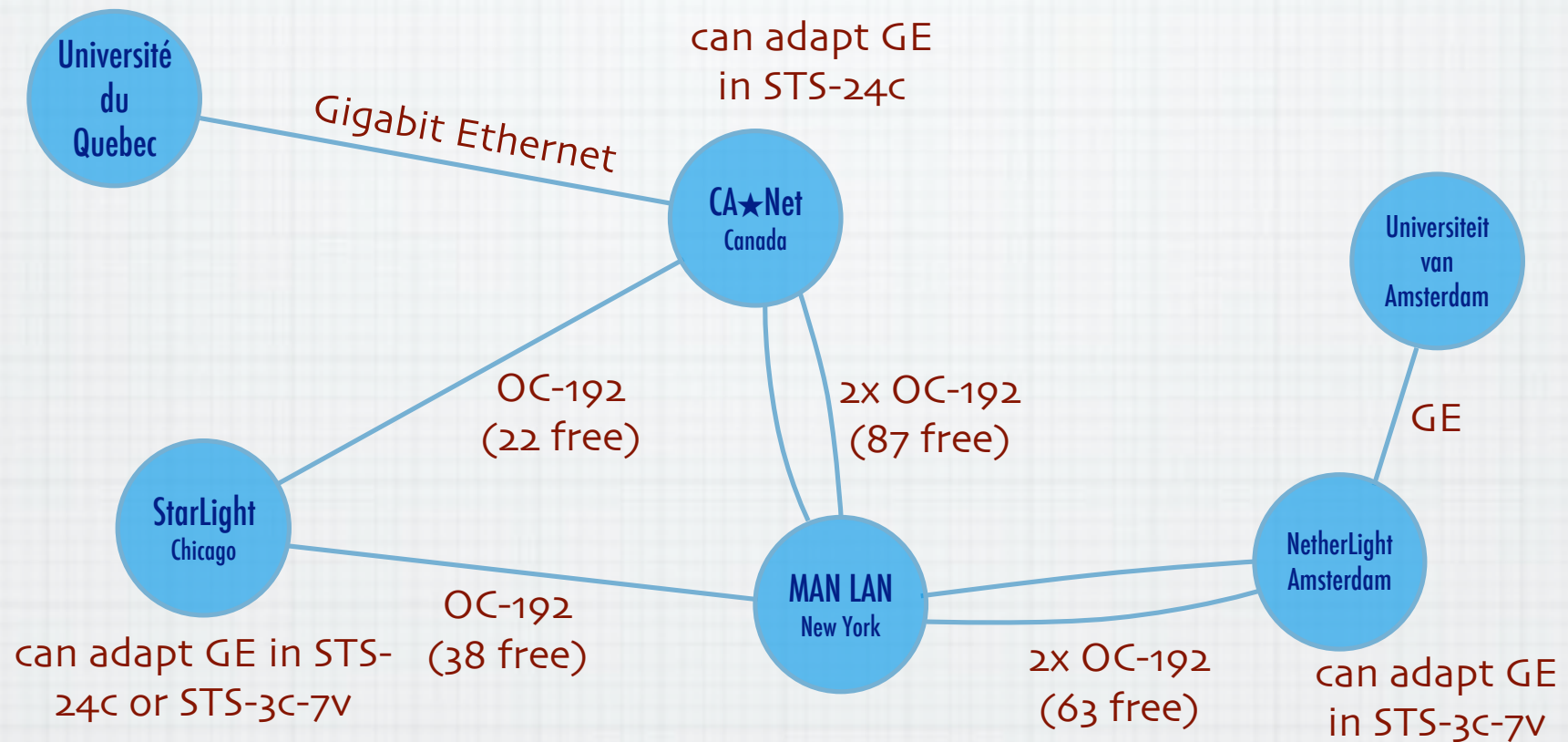
Please enter the URL of the NDL file to be validated

Submit

see http://trafficlight.uva.netherlight.nl/NDL-demo/
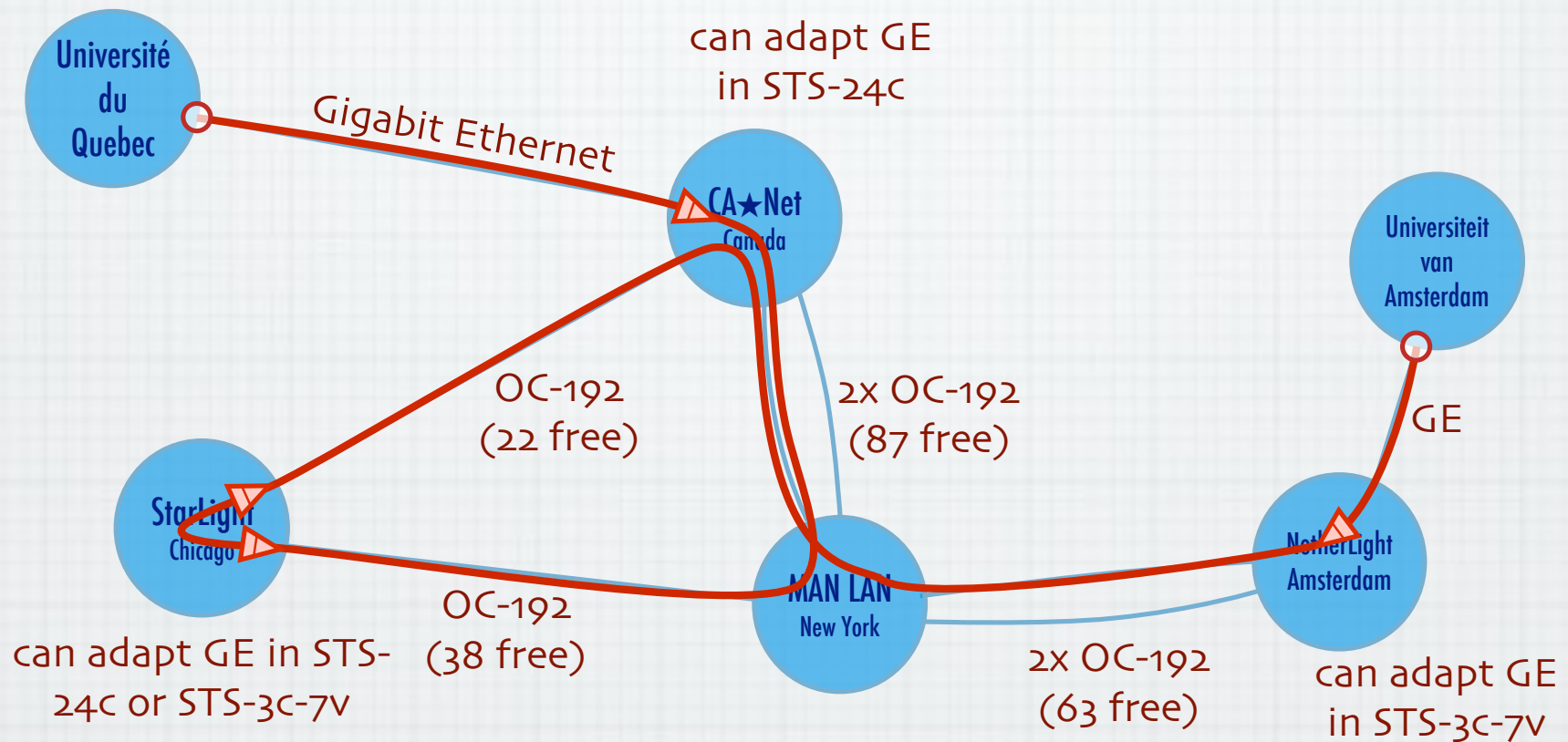
# NDL SN6 Visualisation

# Multi-layer extensions to NDL

# A weird example



Université du Quebec

CA★Net
Canada

can adapt GE
in STS-24c

Gigabit Ethernet

Universiteit van Amsterdam

OC-192
(22 free)

2x OC-192
(87 free)

GE

StarLight
Chicago

MAN LAN
New York

NetherLight
Amsterdam

can adapt GE in STS-
24c or STS-3c-7v

OC-192
(38 free)

2x OC-192
(63 free)

can adapt GE
in STS-3c-7v

# The result :-)



Université du Quebec

can adapt GE in STS-24c

CA★Net Canada

Gigabit Ethernet

Universiteit van Amsterdam

OC-192 (22 free)

2x OC-192 (87 free)

GE

StarLight Chicago

NetherLight Amsterdam

can adapt GE in STS-24c or STS-3c-7v

OC-192 (38 free)

MAN LAN New York

2x OC-192 (63 free)

can adapt GE in STS-3c-7v

Thanks to Freek Dijkstra & team

# MultiDomain MultiLayer pathfinding in action

# MultiDomain MultiLayer pathfinding in action

# OGF NML-WG
## *Open Grid Forum - Network Markup Language workgroup*

Chairs:

Paola Grosso – Universiteit van Amsterdam
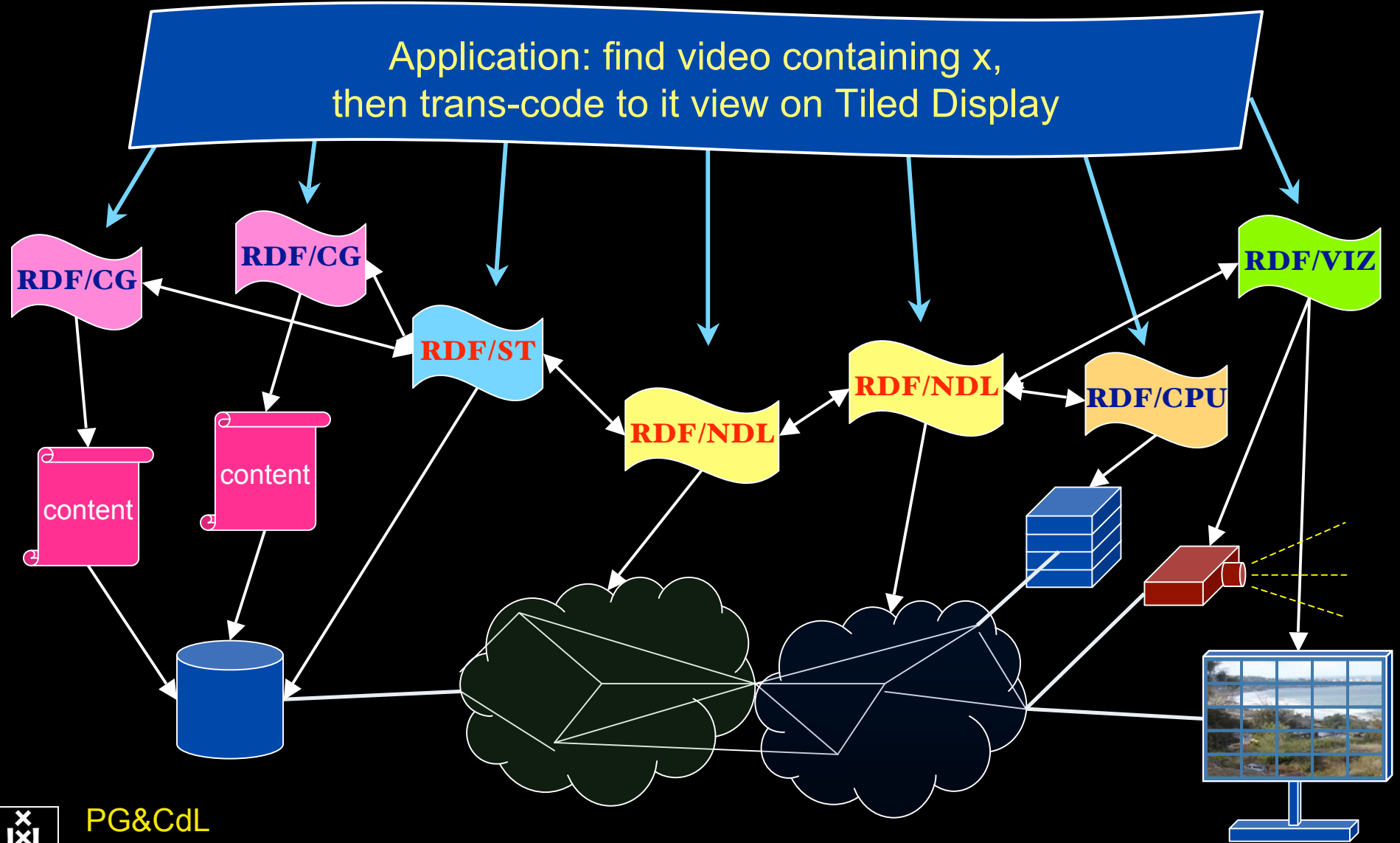
Martin Swany – University of Delaware

Purpose:

*To describe network topologies, so that the outcome is a standardized network description ontology and schema, facilitating interoperability between different projects.*

https://forge.gridforum.org/sf/projects/nml-wg

# Contents

# TeraThinking

- 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, ...
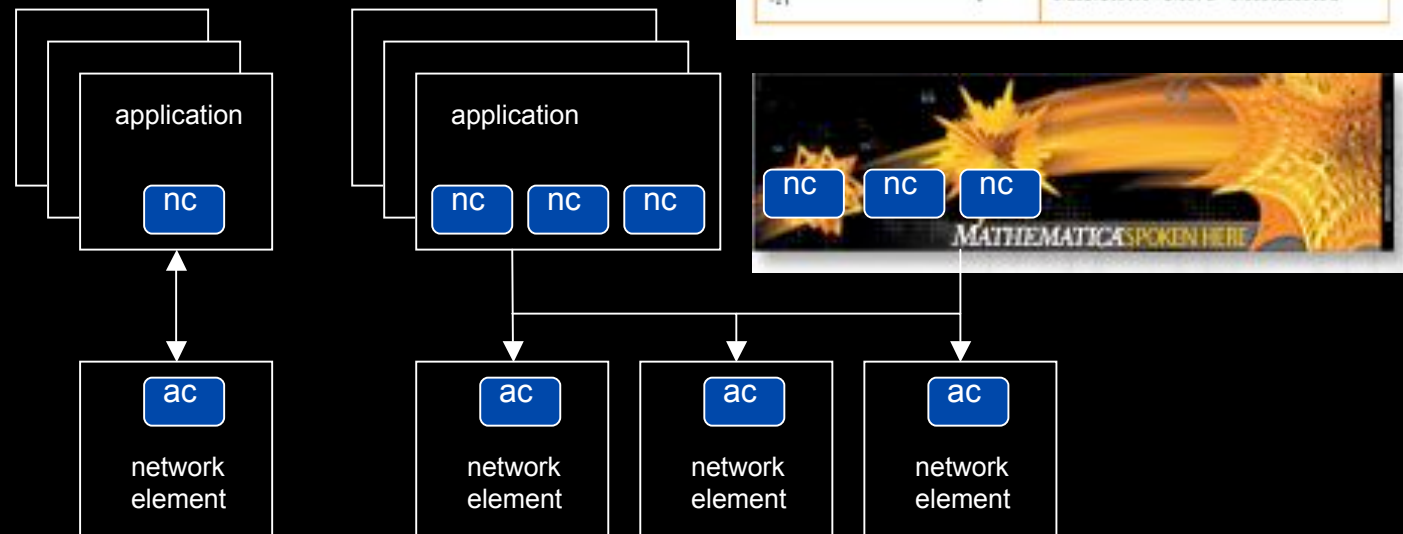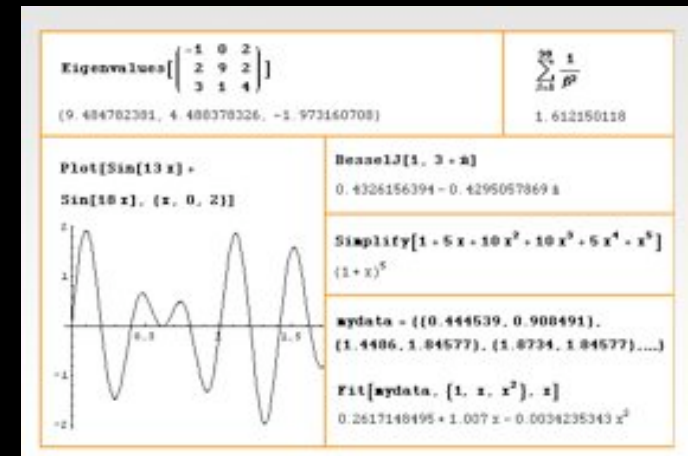  - Tbit/s          ->          ?

# Need for discrete parallelism

- it takes a core to receive 1 or 10 Gbit/s in a computer

- it takes one or two cores to deal with 10 Gbit/s storage

- same for Gigapixels

- same for 100's of Gflops

- Capacity of every part in a system seems of same scale
- look at 80 core Intel processor
  - cut it in two, left and right communicate 8 TB/s
- massive parallel channels in hosts, NIC's

- Therefore we need to go massively parallel allocating complete parts for the problem at hand!

# User Programmable Virtualized Networks allows the results of decades of computer science to handle the complexities of application specific networking.

- The network is virtualized as a collection of resources
- UPVNs enable network resources to be programmed as part of the application
- Mathematica, a powerful mathematical software system, can interact with real networks using UPVNs

# Mathematica enables advanced graph queries, visualizations and real-time network manipulations on UPVNs

## Topology matters can be dealt with algorithmically
## Results can be persisted using a transaction service built in UPVN

### Initialization and BFS discovery of NEs

```
Needs["WebServices`"]
<<DiscreteMath`Combinatorica`
<<DiscreteMath`GraphPlot`
InitNetworkTopologyService["edge.ict.tno.nl"]

Available methods:
{DiscoverNetworkElements,GetLinkBandwidth,GetAllIpLinks,Remote,
NetworkTokenTransaction}

Global`upvnverbose = True;
AbsoluteTiming[nes = BFSDiscover["139.63.145.94"];][[1]]
AbsoluteTiming[result = BFSDiscoverLinks["139.63.145.94", nes];][[1]]

Getting neigbours of: 139.63.145.94
Internal links: {192.168.0.1, 139.63.145.94}
(...)
Getting neigbours of:192.168.2.3
Internal links: {192.168.2.3}
```
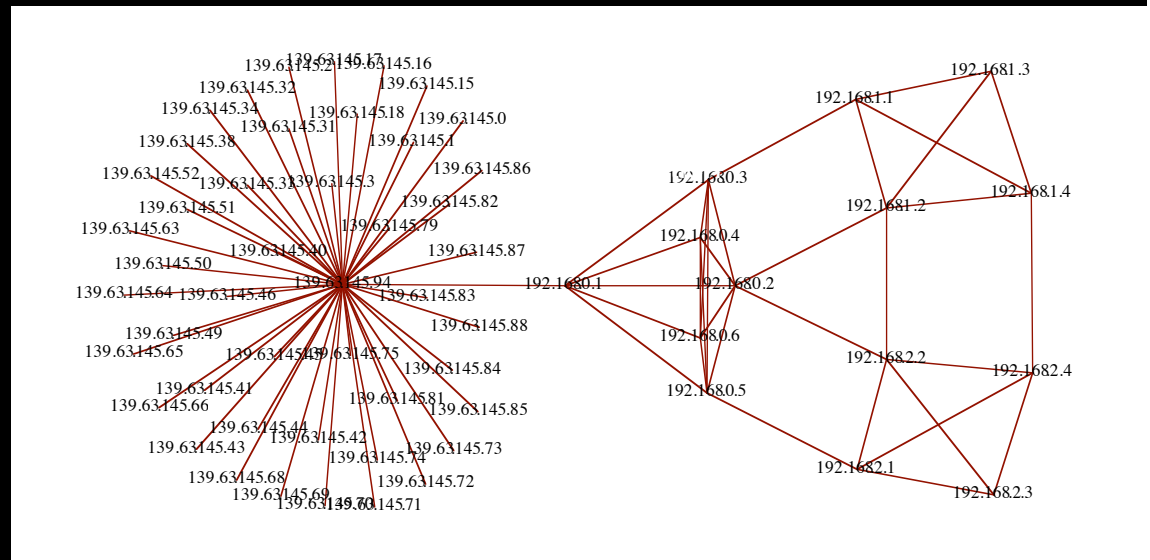
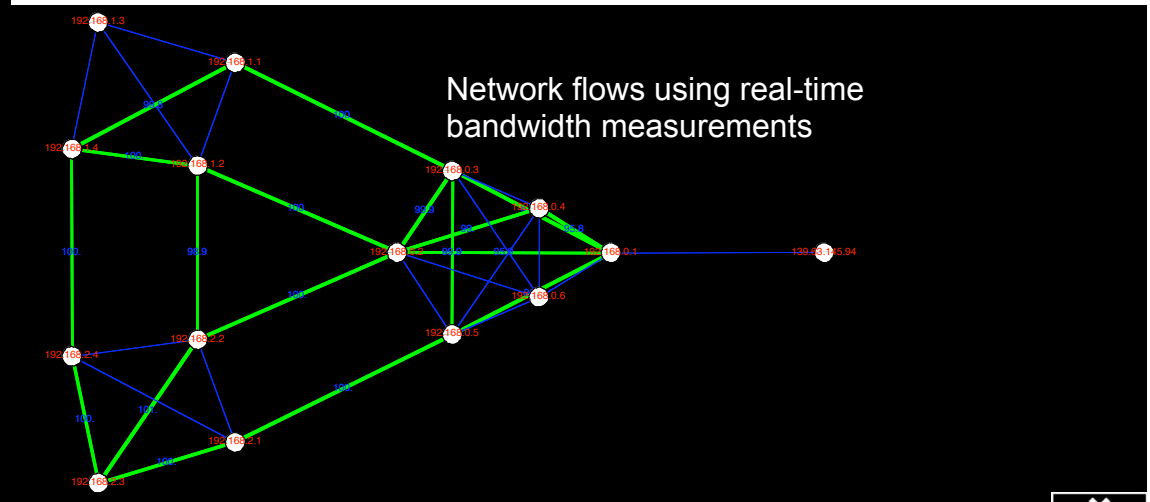### Transaction on shortest path with tokens

```
nodePath = ConvertIndicesToNodes[
            ShortestPath[ g,
                Node2Index[nids,"192.168.3.4"],
                Node2Index[nids,"139.63.77.49"]],
                nids];
Print["Path: ", nodePath];
If[NetworkTokenTransaction[nodePath, "green"]==True,
    Print["Committed"], Print["Transaction failed"]];

Path:
{192.168.3.4,192.168.3.1,139.63.77.30,139.63.77.49}

Committed
```



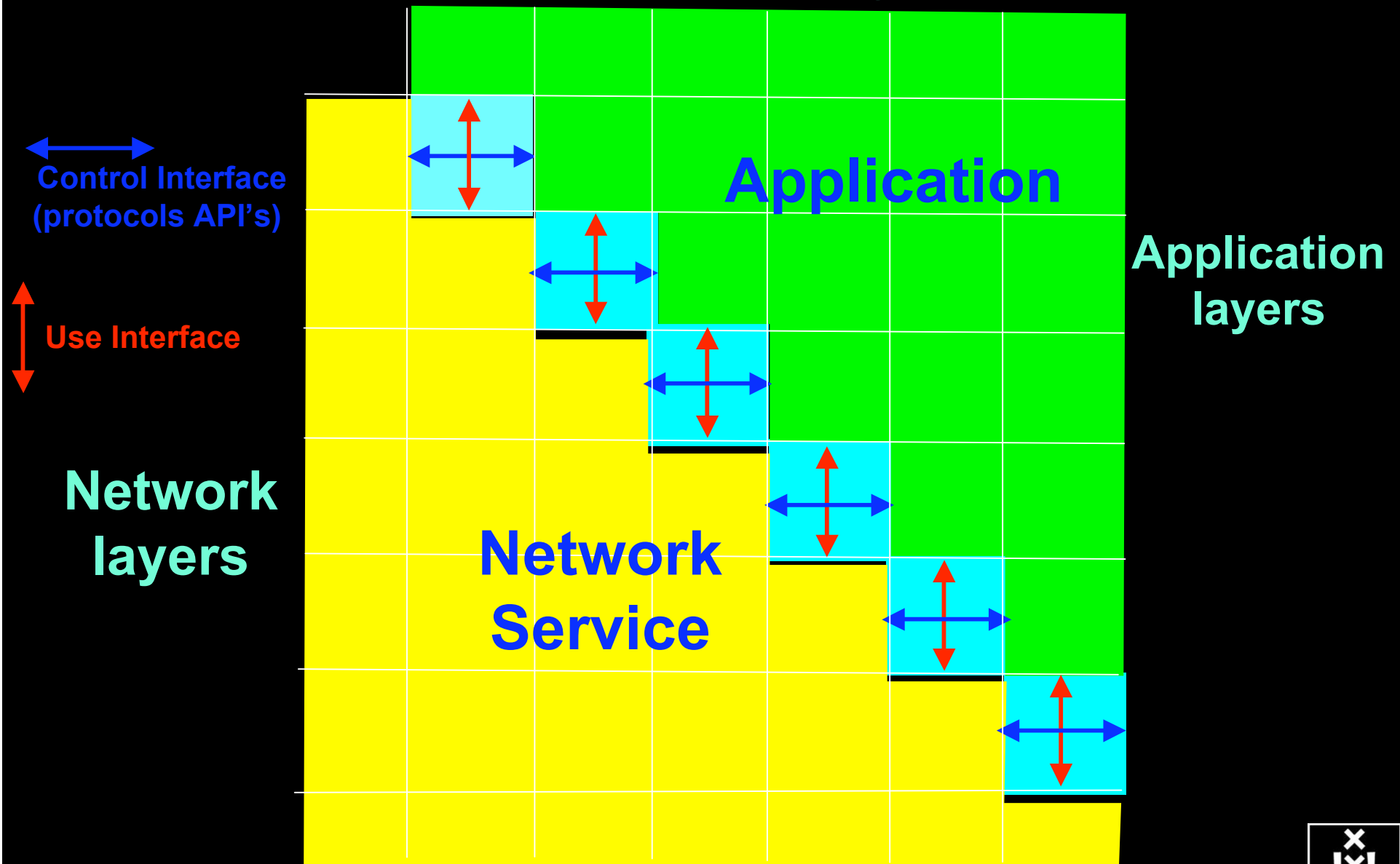Network flows using real-time bandwidth measurements

ref: Robert J. Meijer, Rudolf J. Strijkers, Leon Gommans, Cees de Laat, User Programmable Virtualiized Networks, accepted for publication to the IEEE e-Science 2006 conference Amsterdam.

**StarPlane**

# Functional building blocks



Control Interface (protocols API's)

Use Interface

Network layers

Application layers

Application

Network Service

# 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!

# *Questions ?*

I did not talk about **StarPlane**

...