## Cees de Laat

# EU
# COMMIT
# UvA

NWO

PID/EFRO

SURFnet

TNO

NCF

# Science Faculty @ UvA
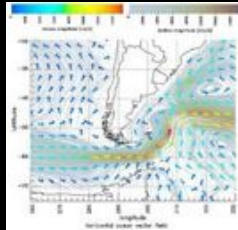


## Informatics Institute

- CSA: Computer Systems Architecture (dr. A.D. Pimentel)

- FCN: Federated Collaborative Networks (Prof. dr. H. Afsarmanesh)

- IAS: Intelligent Autonomous Systems (Prof. dr. ir. F.C.A. Groen)

- ILPS: Information and Language Processing Systems (Prof. dr. M. de Rijke)

- ISIS: Intelligent Sensory Information Systems (Prof. dr. ir. A.W.M. Smeulders)

- SCS: Section Computational Science (Prof. dr. P.M.A. Sloot)

- SNE: System and Network Engineering (Prof. dr. ir. C.T.A.M. de Laat)

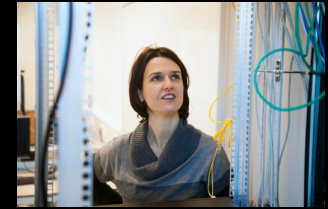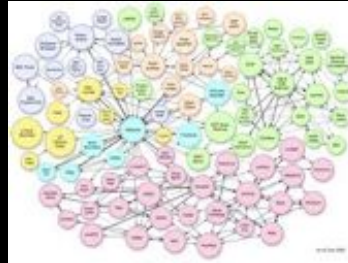- TCS: Theory of Computer Science (Prof. dr. J.A. Bergstra)

... more data!

Internet developments

DATA

... more realtime!

... more users!

# GPU cards are distruptive!



**Legend:**
- fastest supercomputer in the world
- nr. 500 supercomputer in the world
- **1 single Graphics Processing Unit**

**20.000.000$**

**7 year**

**500$**

**Top 500**

**#1**

**#500**

# Data storage: doubling every 1.5 year!

# Multiple colors / Fiber



**Wavelength Selective Switch**

Per fiber: ~ 80-100 colors * 50 GHz

Per color: 10 – 40 – 100 Gbit/s

BW * Distance ~ $2*10^{17}$ bm/s

New: Hollow Fiber!
→ less RTT!

Optical transmission

... more possibilities


CWDM Technology


DWDM Technology


common port — Wavelength Selectable Switch (WSS) — N ports

Virtualization

**# users** (vertical axis label)

A. Lightweight users, browsing, mailing, home use
   Need full Internet routing, one to all

B. Business/grid applications, multicast, streaming, VO's, mostly LAN
   Need VPN services and full Internet routing, several to several + uplink to all

C. E-Science applications, distributed data processing, all sorts of grids
   Need very fat pipes, limited multiple Virtual Organizations, P2P, few to few

For the Netherlands 2011
$\Sigma A = \Sigma B = \Sigma C \approx 1$ Tb/s
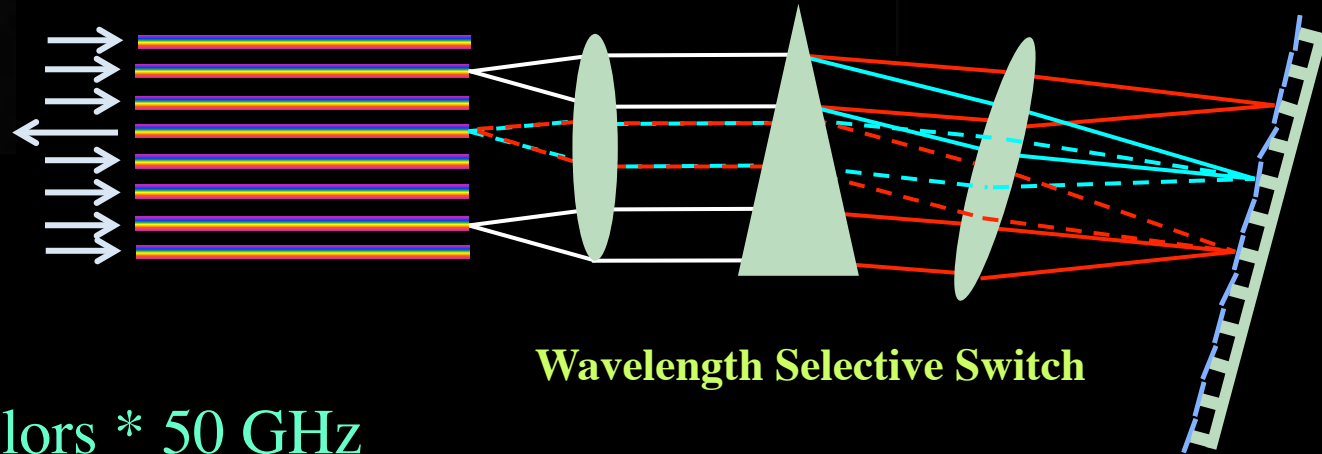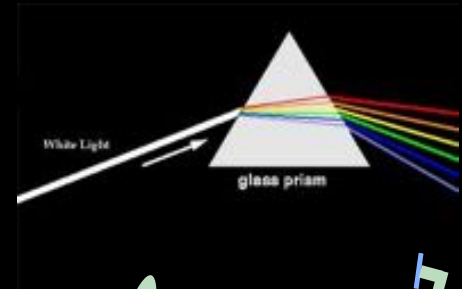However:
 A -> all connects
 B -> on several
 C -> just a few (SP, LHC, LOFAR)

A

B

C

ADSL (20 Mbit/s)

GigE

BW

Ref: Cees de Laat, Erik Radius, Steven Wallace, "The Rationale of the Current Optical Networking Initiatives"
iGrid2002 special issue, Future Generation Computer Systems, volume 19 issue 6 (2003)

# 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

# SNE @ UvA



Speed / Volume

Deterministic / Real-time

Scalable / Secure

| | IJkdijk/Urban Flood | Medical | LifeWatch/ENVRI | CosmoGrid/eVLBI | CineGrid | EU-GN3/NOVI/Geysers | SURFnet/GLIF/Cloud |
|---|---|---|---|---|---|---|---|
| Green-IT | | | | | X | X | |
| Privacy/Trust | | X | | | X | | |
| Authorization/policy | | X | X | | X | X | |
| Programmable networks | X | | X | | | | |
| 40-100Gig/TCP/WF/QoS | X | | | X | X | X | |
| Topology/Architecture | | X | | X | X | X | |
| Optical Photonic | | | X | X | | X | |

SNE @ UvA

# SNE @ UvA

Where when will it happen?

| | Ijkdijk/Urban Flood | Medical | LifeWatch/ENVRI | CosmoGrid/eVLBI | CineGrid | EU-GN3/NOVI/Geysers | SURFnet/GLIF/Cloud |
|---|---|---|---|---|---|---|---|
| Green-IT | | | | | | X | X |
| Privacy/Trust | | X | | | X | | |
| Authorization/policy | | X | X | | X | X | |
| Programmable networks | X | | X | | | | |
| 40-100Gig/TCP/WF/QoS | X | | | X | X | | X |
| Topology/Architecture | | X | | X | X | X | |
| Optical Photonic | | X | X | | X | | |

**IJKDIJK**

Sensors: 15000km* 800 bps/m ->12 Gbit/s to cover all Dutch dikes

# Sensor grid: instrument the dikes

First controlled breach occurred on sept 27th '08:

**Many Pflops/s**

**Many small flows -> 12 Gb/s**

# User Programmable Virtualized Networks.

- The network is virtualized as a collection of resources
- UPVNs enable network resources to be programmed as part of the application
- Mathematica interacts with virtualized networks using UPVNs and optimize network + computation

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.

# In the Intercloud virtual servers and networks become software

- Virtual Internets adapt to the environment, grow to demand, iterate to specific designs
- Network support for application specific interconnections are merely opitimizations: Openflow, active networks, cisco distributed switch
- But how to control the control loop?

# Interactive Networks

- SuperComputing 2008
- SuperComputing 2009 (in programmable Grid networks demo)
- SuperComputing 2011
- Next Generation Telecom networks workshop 2011
- LHCONE architecture workshop 2011

# Interactive Networks

Rudolf Strijkers [1,2]

Marc X. Makkes [1,2]

Mihai Christea [1]

Laurence Muller [1]

Robert Belleman [1]

Cees de Laat [1]

Robert Meijer [1,2]

[1] University of Amsterdam, Amsterdam The Netherlands

[2] TNO Information and Communication Technology, Groningen, The Netherlands

# SNE @ UvA



| | Ijkdijk/Urban Flood | | Medical | LifeWatch/ENVRI | CosmoGrid/eVLBI | CineGrid | EU-GN3/NOVI/Geysers | SURFnet/GLIF/Cloud |
|---|---|---|---|---|---|---|---|---|
| Green-IT | | | | | | | X | X |
| Privacy/Trust | | X | | | | X | | |
| Authorization/policy | | X | X | | X | X | | |
| Programmable networks | X | | X | | | | | |
| 40-100Gig/TCP/WF/QoS | X | | | X | X | | X | |
| Topology/Architecture | | X | | X | X | X | | |
| Optical Photonic | | | X | X | | X | | |

# ATLAS detector @ CERN Geneve

# ATLAS detector @ CERN Geneve

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

**100000 flops/byte**

**10 Pflops/s**

~PByte/sec

Online System

~100 MBytes/sec

event simulation

*Tier 0 +1*

cHPSS

event reconstruction

CMS detector: 15m X 15m X 22m

12,500 tons, $700M.

**Status 2002!**

*Tier 1*

~2.5 Gbits/sec

Italian Regional Center

German Regional Center

NIKHEF Dutch Regional Center

FermiLab, USA Regional Center

...

analysis

~0.6-2.5 Gbps

Tier2 Center   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

human=2m

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.

# Diagram for SAGE video streaming to ATS



**Lab 10, Nortel**

SAGE Display  
SAGE Servers  
MERS  
1 Gbps

User  
Regular Browser

Netherlight Canarie

Internet  
*Content Choice*

**UvA, Amsterdam**

MERS  
MERS  
MERS  
comp clusters  
Traffic Generators

1 Gbps

Content Portal  
Streaming Server  
100 TB Storage

*Content Request*

# Experimental Data

**Sage without background traffic**

**Sage with background traffic**



Experiment waag geometry-nopbt-10secburst



Experiment waag geometry-pbt-10secburst

**10 Second Traffic bursts with No PBT**

**10 Second Traffic bursts with PBT**

PBT is *SIMPLE* and *EFFECTIVE* technology to build a shared Media-Ready Network

# Alien light From idea to realisation!

# 40Gb/s alien wavelength transmission via a multi-vendor 10Gb/s DWDM infrastructure

**N C F**

## Alien wavelength advantages

- Direct connection of customer equipment[1]
  → cost savings
- Avoid OEO regeneration → power savings
- Faster time to service[2] → time savings
- Support of different modulation formats[3]
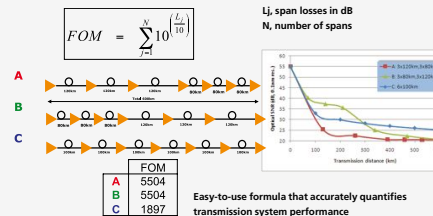  → extend network lifetime

## Alien wavelength challenges

- Complex end-to-end optical path engineering in terms of linear (i.e. OSNR, dispersion) and non-linear (FWM, SPM, XPM, Raman) transmission effects for different modulation formats.
- Complex interoperability testing.
- End-to-end monitoring, fault isolation and resolution.
- End-to-end service activation.

**In this demonstration we will investigate the performance of a 40Gb/s PM-QPSK alien wavelength installed on a 10Gb/s DWDM infrastructure.**
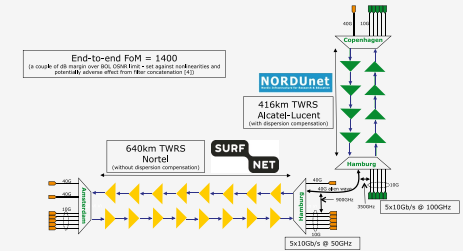
## New method to present fiber link quality, FoM (Figure of Merit)

In order to quantify optical link grade, we propose a new method of representing system quality: the FOM (Figure of Merit) for concatenated fiber spans.
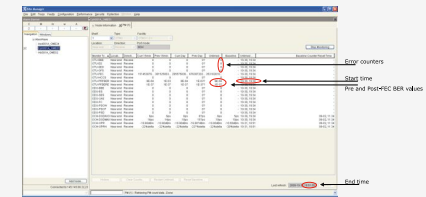


$$FOM = \sum_{j=1}^{N} 10^{\left(\frac{L_j}{10}\right)}$$

Lj, span losses in dB
N, number of spans

| | FOM |
|---|---|
| A | 5504 |
| B | 5504 |
| C | 1897 |

**Easy-to-use formula that accurately quantifies transmission system performance**

## Transmission system setup

JOINT SURFnet/NORDUnet 40Gb/s PM-QPSK alien wavelength DEMONSTRATION.



## Test results



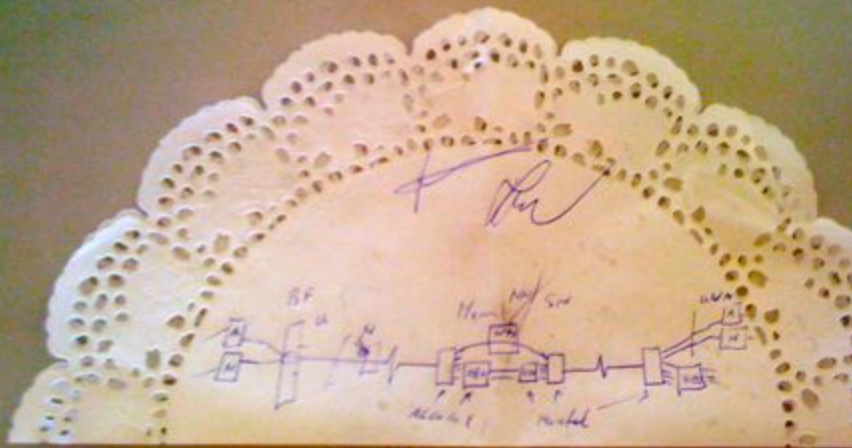Error-free transmission for 23 hours, 17 minutes → BER < 3.0 10^-16

## Conclusions

- We have investigated experimentally the all-optical transmission of a 40Gb/s PM-QPSK alien wavelength via a concatenated native and third party DWDM system that both were carrying live 10Gb/s wavelengths.
- The end-to-end transmission system consisted of 1056 km of TWRS (TrueWave Reduced Slope) transmission fiber.
- We demonstrated error-free transmission (i.e. BER below 10-15) during a 23 hour period.
- More detailed system performance analysis will be presented in an upcoming paper.

**NORTEL** · **NORDUnet** Nordic Infrastructure for Research and Education · **telindus** together with **IBIT** · **SURF NET**

# Alien light From idea to realisation!



# 40Gb/s alien wavelength transmission via a multi-vendor 10Gb/s DWDM infrastructure

NCF

### Alien wavelength advantages

- Direct connection of customer equipment[1]
  → cost savings
- Avoid OEO regeneration → power savings
- Faster time to service[2] → time savings
- Support of different modulation formats[3]
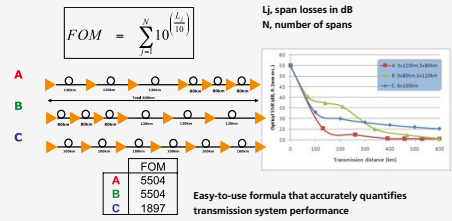  → extend network lifetime

### Alien wavelength challenges

- Complex end-to-end optical path engineering in terms of linear (i.e. OSNR, dispersion) and non-linear (FWM, SPM, XPM, Raman) transmission effects for different modulation formats.
- Complex interoperability testing.
- End-to-end monitoring, fault isolation and resolution.
- End-to-end service activation.

**In this demonstration we will investigate the performance of a 40Gb/s PM-QPSK alien wavelength installed on a 10Gb/s DWDM infrastructure.**
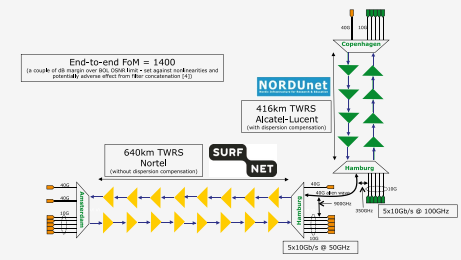
### New method to present fiber link quality, FoM (Figure of Merit)

In order to quantify optical link grade, we propose a new method of representing system quality: the FOM (Figure of Merit) for concatenated fiber spans.
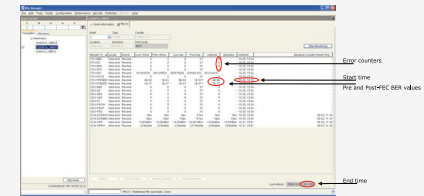


$$FOM = \sum_{j=1}^{N} 10^{\left(\frac{L_j}{10}\right)}$$

Lj, span losses in dB
N, number of spans

| | FOM |
|---|---|
| A | 5504 |
| B | 5504 |
| C | 1897 |

**Easy-to-use formula that accurately quantifies transmission system performance**

### Transmission system setup

JOINT SURFnet/NORDUnet 40Gb/s PM-QPSK alien wavelength DEMONSTRATION.



### Test results



Error-free transmission for 23 hours, 17 minutes → BER < 3.0 10⁻¹⁶

### Conclusions

- We have investigated experimentally the all-optical transmission of a 40Gb/s PM-QPSK alien wavelength via a concatenated native and third party DWDM system that both were carrying live 10Gb/s wavelengths.
- The end-to-end transmission system consisted of 1056 km of TWRS (TrueWave Reduced Slope) transmission fiber.
- We demonstrated error-free transmission (i.e. BER below 10-15) during a 23 hour period.
- More detailed system performance analysis will be presented in an upcoming paper.

NORTEL    NORDUnet    telindus    together with IBIT    SURFNET

# Visit CIENA Booth
# surf to http://tnc11.delaat.net



**ClearStream**
End-to-End Ultra Fast Transmission Over a Wide Area 40 Gbit/s Lambda

Incoming Amsterdam 25.5 Gbps

Incoming Copenhagen 20.97 Gbps

**Total Throughput 46.47 Gbps RTT 44.032 ms**

# Results (rtt = 17 ms)

- Single flow iPerf  1 core            ->      21 Gbps

- Single flow iPerf  1 core    <>    ->      15+15 Gbps

- Multi flow iPerf 2 cores            ->      25 Gbps

- Multi flow iPerf 2 cores    <>    ->      23+23 Gbps

- DiViNe                            <>    ->      11 Gbps

- Multi flow iPerf + DiVine          ->      35 Gbps

- Multi flow iPerf + DiVine <>    ->      35 + 35 Gbps

# Performance Explained

- Mellanox 40GE card is PCI-E 2.0 8x (5GT/s)
- 40Gbit/s raw throughput but ….
- PCI-E is a network-like protocol
  - 8/10 bit encoding -> 25% overhead -> 32Gbit/s maximum data throughput
  - Routing information
- Extra overhead from IP/Ethernet framing
- Server architecture matters!
  - 4P system performed worse in multithreaded iperf

# Server Architecture



DELL R815
4 x AMD Opteron 6100



Supermicro X8DTT-HIBQF
2 x Intel Xeon

# CPU Topology benchmark



We used numactl to bind iperf to cores

GLIF 2011

**Visualization courtesy of Bob Patterson, NCSA
Data collection by Maxine Brown.**

We investigate:  for complex networks!

# The GLIF – lightpaths around the world

# LinkedIN for Infrastructure

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

Subject → **Predicate** → Object

Subject → Object Subject → Object Subject → Object Subject
Subject → Object Subject → Object Subject

Location — name → / connectedTo →

Device — description → / capacity →

Interface — locatedAt → / encodingType →

Link — hasInterface → / encodingLabel →

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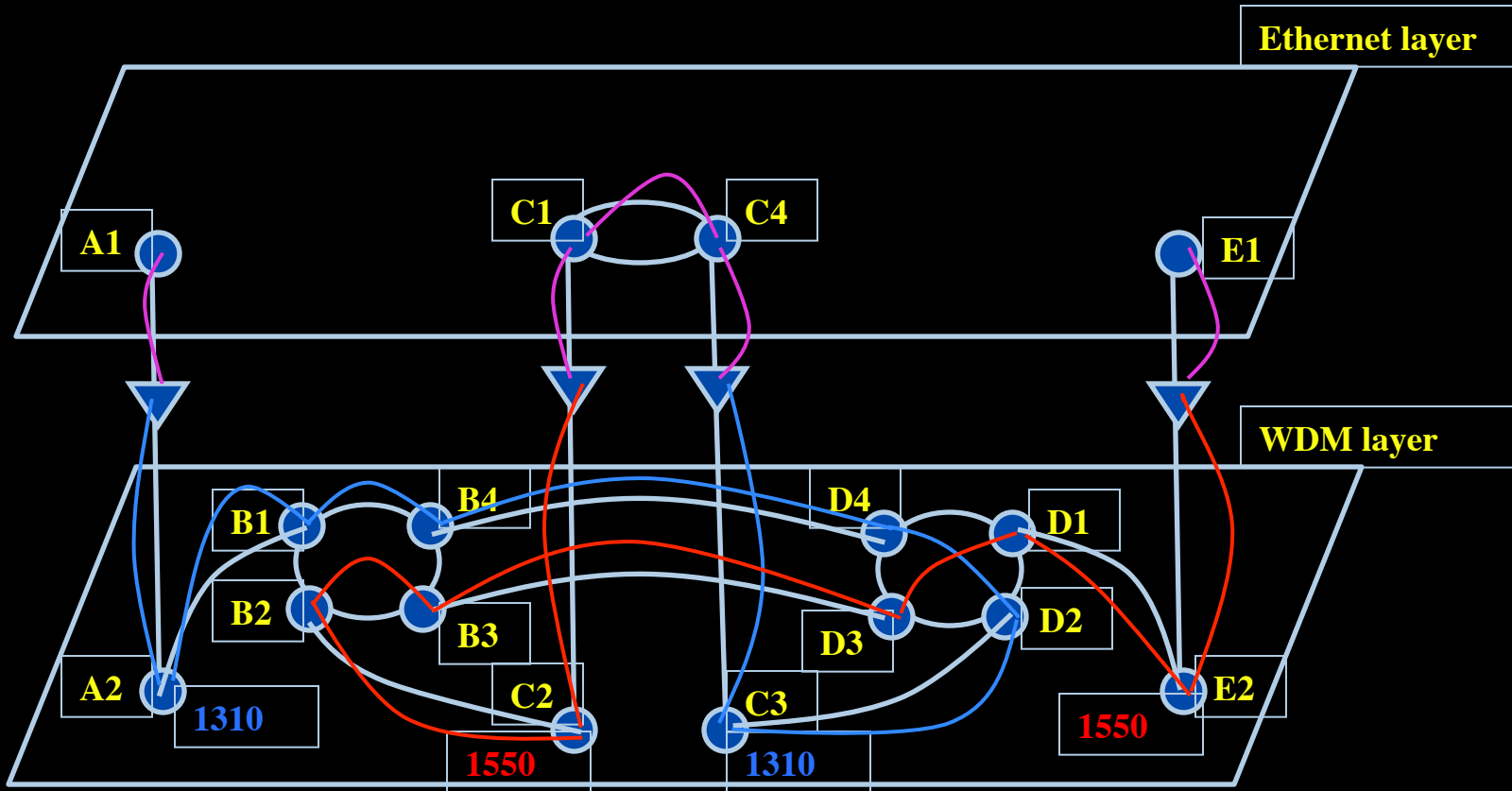
```xml
<!-- all the interfaces of TDM3.amsterdam1.netherlight.net -->

<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>
```

# Multi-layer descriptions in NDL

# Multi-layer Network PathFinding



Path between interfaces A1 and E1:

A1-A2-B1-B4-D4-D2-C3-C4-C1-C2-B2-B3-D3-D1-E2-E1

Scaling: Combinatorial problem

# Virtualisatie van infrastructuur & QoS

# Information Modeling

Define a common information model for **infrastructures** and **services**. Base it on Semantic Web.

J. van der Ham, F. Dijkstra, P. Grosso, R. van der Pol, A. Toonk, C. de Laat
*A distributed topology information system for optical networks based on the semantic web*,
In: Elsevier Journal on Optical Switching and Networking, Volume 5, Issues 2-3, June 2008, Pages 85-93

R.Koning, P.Grosso and C.de Laat
*Using ontologies for resource description in the CineGrid Exchange*
In: Future Generation Computer Systems (2010)

# SNE @ UvA

| | IJkdijk/Urban Flood | Medical | LifeWatch/ENVRI | CosmoGrid/eVLBI | CineGrid | EU-GN3/NOVI/Geysers | SURFnet/GLIF/Cloud |
|---|---|---|---|---|---|---|---|
| Green-IT | | | | | X | X | |
| Privacy/Trust | | | X | | X | | |
| Authorization/policy | | | X | X | X | X | |
| Programmable networks | X | | X | | | | |
| 40-100Gig/TCP/WF/QoS | X | | | X | X | X | |
| Topology/Architecture | | | X | X | X | X | |
| Optical Photonic | | | X | X | | X | |

# Why is more resolution is better?
1. More Resolution Allows Closer Viewing of Larger Image
2. Closer Viewing of Larger Image Increases Viewing Angle
3. Increased Viewing Angle Produces Stronger Emotional Response

HDTV (2K)

1920

1080

***30°***

3.0 × Picture Height

UHDTV(8K)

**24 Gb/s**

7680

4320

SHD (4K)

3840

2160

**7.6 Gb/s**

0.75 × Picture Height

***100°***

***60°***

1.5 × Picture Height

Yutaka TANAKA
SHARP CORPORATION
Advanced Image Research Laboratories

# RDF describing Infrastructure



Application: find video containing x,
then trans-code to it view on Tiled Display

RDF/CG

RDF/CG

RDF/VIZ

RDF/ST

RDF/NDL

RDF/NDL

RDF/CPU

content

content

PG&CdL

# Why?

## I want to:

"Show Big Bug Bunny in 4K on my Tiled Display using green Infrastructure"

- Big Bugs Bunny can be on multiple servers on the Internet.
- Movie may need processing / recoding to get to 4K for Tiled Display.
- Needs deterministic Green infrastructure for Quality of Experience.
- Consumer / Scientist does not want to know the underlying details.
  ➔ His refrigerator also just works.

# Applications and Networks become aware of each other!

# The Ten Problems with the Internet
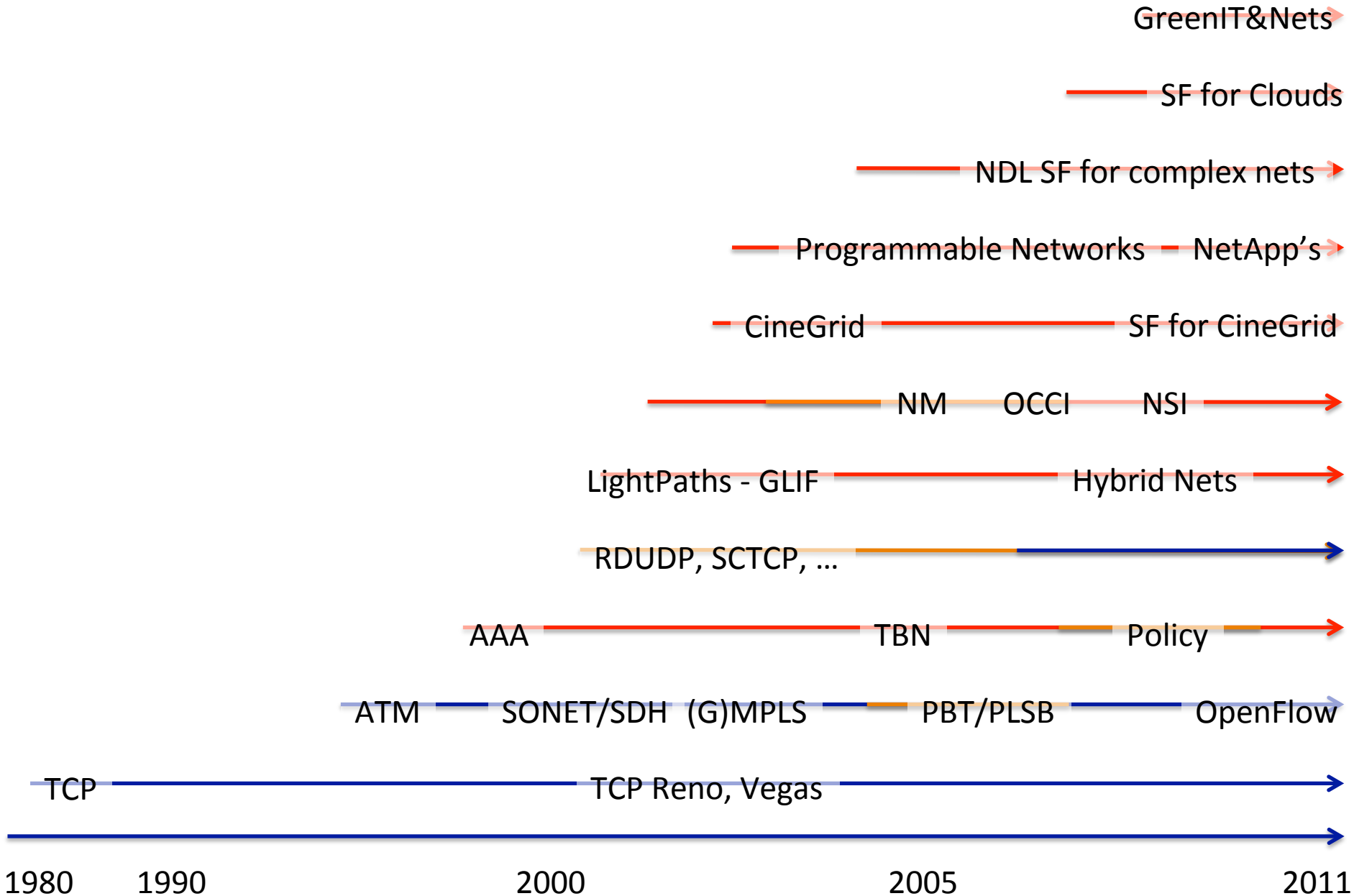
1. **Energy Efficient Communication**
2. Separation of Identity and Address
3. Location Awareness
4. **Explicit Support for Client-Server Traffic and Distributed Services**
5. Person-to-Person Communication
6. Security
7. **Control, Management, and Data Plane separation**
8. **Isolation**
9. Symmetric/Asymmetric Protocols
10. **Quality of Service**

*Nice to have:*
- Global Routing with Local Control of Naming and Addressing
- **Real Time Services**
- **Cross-Layer Communication**
- Manycast
- Receiver Control
- Support for Data Aggregation and Transformation
- **Support for Streaming Data**
- **Virtualization**

ref: Raj Jain, "Internet 3.0: Ten Problems with Current Internet Architecture and Solutions for the Next Generation", Military Communications Conference, 2006. MILCOM 2006. IEEE

# TimeLine

GreenIT&Nets

SF for Clouds

NDL SF for complex nets

Programmable Networks — NetApp's

CineGrid — SF for CineGrid

NM    OCCI    NSI

LightPaths - GLIF    Hybrid Nets

RDUDP, SCTCP, …

AAA    TBN    Policy

ATM    SONET/SDH  (G)MPLS    PBT/PLSB    OpenFlow

TCP    TCP Reno, Vegas

1980    1990    2000    2005    2011

# TimeLine

GreenIT&Nets — **Sustainable Internet**

SF for Clouds

NDL SF for complex nets

Programmable Networks — NetApp's

**Cognitive Nets and clouds**

CineGrid — SF for CineGrid

**Machine Learning** **+** → **"I Want" Internet 3.0**

NM   OCCI   NSI

aths - GLIF — Hybrid Nets

OP, SCTCP, …

**Virtualized Internet**

TBN — Policy

(G)MPLS — PBT/PLSB — OpenFlow

TCP — **Good Old Trucking**

2005 — 2011 — 2020

# TimeLine

**Sustainable Internet**

**Cognitive Nets and clouds**

**Machine Learning** **+** → **"I Want" Internet 3.0**

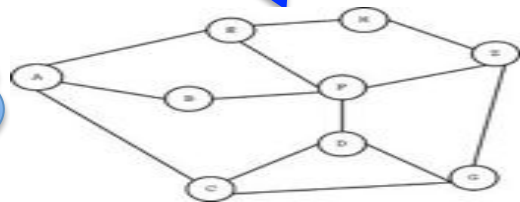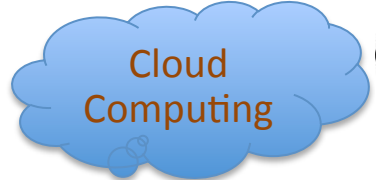**Virtualized Internet**

**Good Old Trucking**

I retire

2020

2040

# Complex e-Infrastructure!

# Complex e-Infrastructure!

California Institute for Telecommunications and Information Technology (Calit2)
Wide Area Network

**Domain Apps**   **Domain Apps**   ...   ...   **Domain Apps**   **Domain Apps**

**eScience Middleware**

**+ ML + reasoning (ProLog?) + Scheduling + …**

**Service Plane**

**SAGE CGLX Cromium** | **SAGE** | **OCCI JSDL SAGA** | **GIR UR** | **NSI NetConf SNMP OpenFlow** | **PerfSonar** | **DIAS ByteIO iRODs** | **DIAS ByteIO** | **OGSA** | **WebServ**

Cloud Computing
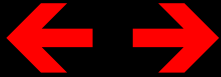
# ECO-Scheduling

# Hybrid Networking <->  Computing

Routers       ← →       Supercomputers

Ethernet switches   ← →   Grid  & Cloud

Photonic transport   ← →   GPU's

What matters:

    Energy consumption/multiplication

    Energy consumption/bit transported

# Challenges

- **Data – Data – Data**
  - Archiving, publication, searchable, transport, self-describing, DB innovations needed, multi disciplinary use

- **Virtualisation**
  - Another layer of indeterminism

- **Greening the Infrastructure**
  - e.g. Department Of Less Energy: http://www.ecrinitiative.org/pdfs/ECR_3_0_1.pdf

- **Disruptive developments**
  - BufferBloath, Revisiting TCP, influence of SSD's & GPU's
  - Multi layer Glif Open Exchange model
  - Invariants in LightPaths (been there done that ☺)
    - X25, ATM, SONET/SDH, Lambda's, MPLS-TE, VLAN's, PBT, OpenFlow, ….
  - Authorization & Trust & Security and Privacy

# The Way Forward!

- Nowadays scientific computing and data is dwarfed by commercial & cloud, there is also no scientific water, scientific power.
  - Understand how to work with elastic clouds
  - Trust & Policy & Firewalling on VM/Cloud level
- Technology cycles are 3 – 5 year
  - Do not try to unify but prepare for diversity
  - Hybrid computing & networking
  - Compete on implementation & agree on interfaces and protocols
- Limitation on natural resources and disruptive events
  - Energy becomes big issue
  - Follow the sun
  - Avoid single points of failure (aka Amazon, Blackberry, …)
  - Better very loosely coupled than totally unified integrated…

# Q & A

Slides thanks to:

- Paola Grosso
- Sponsors see slide 1. ☺
- SNE Team & friends, see below