# System and Network Engineering Research for Big Data Sciences

## Cees de Laat

Freek Victor Blom Noordende Paola Ham
Leon Taal Pieter Bruijn Spence Inder Stancu Mulmo Sloot Cook
Radius Koning
Marchal Demchenko Oudenaarde
Mambretti Farrell Andree Bal Zwart Damien Sjouw Maassen Peter Guido
Hirstius
Adam Zhao Groep Vollbrecht Travostino Grosso Wan
Koymans Adriaans Catalin Maxine Erik
Martin Yahyapour Hertzberger Cristea Fred Kelgo Paul Wan
Jan-Philip Gordon Pat Simon Halepidis Kees
Olabarriaga Karst Henri Guevara-Masis Derek Ralph Yakali Silvia George Oscar Bob
DeFanti Koeroo Grossman Lavian Phlip Dobinson Smarr Hans
Gross Li Meirosu Stefan Denus Bert Eljkel Gosso Arie
Belleman Pol Wim Lee Piotr Buuren Joe Larry Vladimir Golonka Ronald Steven
Korkhov Hendrikse David Xu Brown Velders Snijders Portegies
Andreas Rene Olle Groen Belloum
Thomas Mihai Meijer Antony Dijkstra
Dmitry Rudolf Monga Koot Toonk Zhiming Jason Bas Tokmakoff
Calhoun Hakan Vasunin Jeroen
John Strijkers Yuri Franco Zhiming Zeger
Matthijs

From King's Dutch Academy of Sciences
# The Dutch Research Agenda

"Information technology (IT) now permeates all aspects of public, commercial, social, and personal life. bank cards, satnav, and weather radar... IT has become completely indispensable."

"But to guarantee the reliability and quality of constantly bigger and more complicated IT, we will need to find answers to some fundamental questions!"
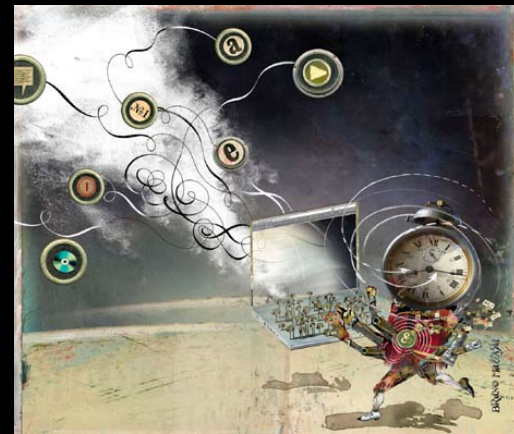
BRuNO MallART

# Reduction of Complexity by Integration

By combining services such as telephony, television, data, and computing capacity within a single network, we can cut down on complexity, energy consumption and maintenance.

- How can we describe and analyze complex information systems effectively?

- How can we specify and measure the quality and reliability of a system?

- How can we combine various different systems?

- How can we design systems in which separate processors can co-operate efficiently via mutual network connections within a much larger whole?

- Can we design information systems that can diagnose their own malfunctions and perhaps even repair them?

- How can we specify, predict, and measure system performance as effectively as possible?

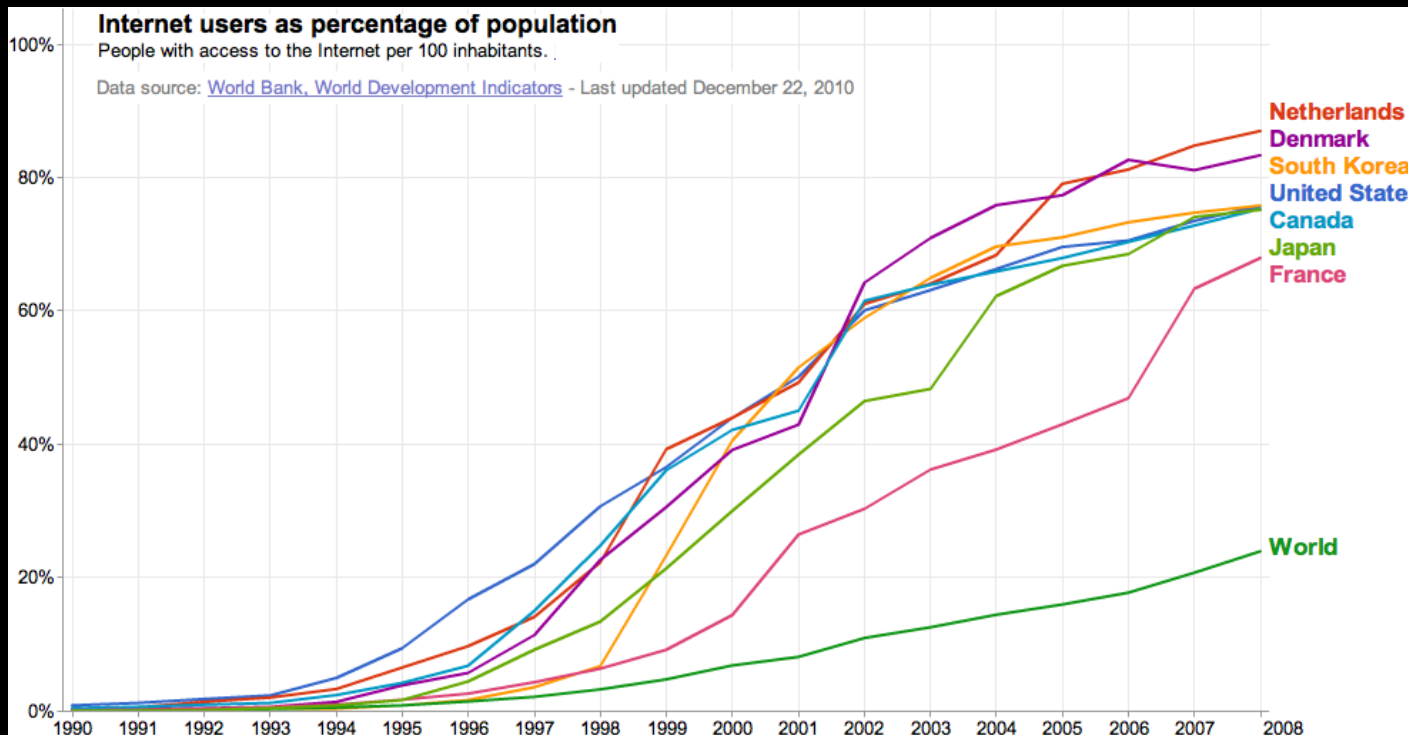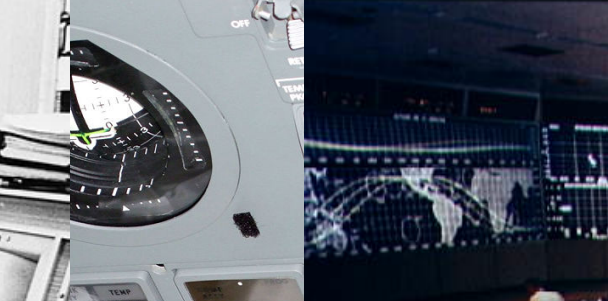SNE addresses a.o. the **highlighted** questions!

# Internet

## From a network experiment that never ended (Vint Cerf)

1974: for the first time the word **internet** (*RFC 675 - Specification of Internet Transmission Control Program) [note -> Open process!]*

1981: the **TCP/IP** standard was ready to be adopted (*RFC 791,792,793*)
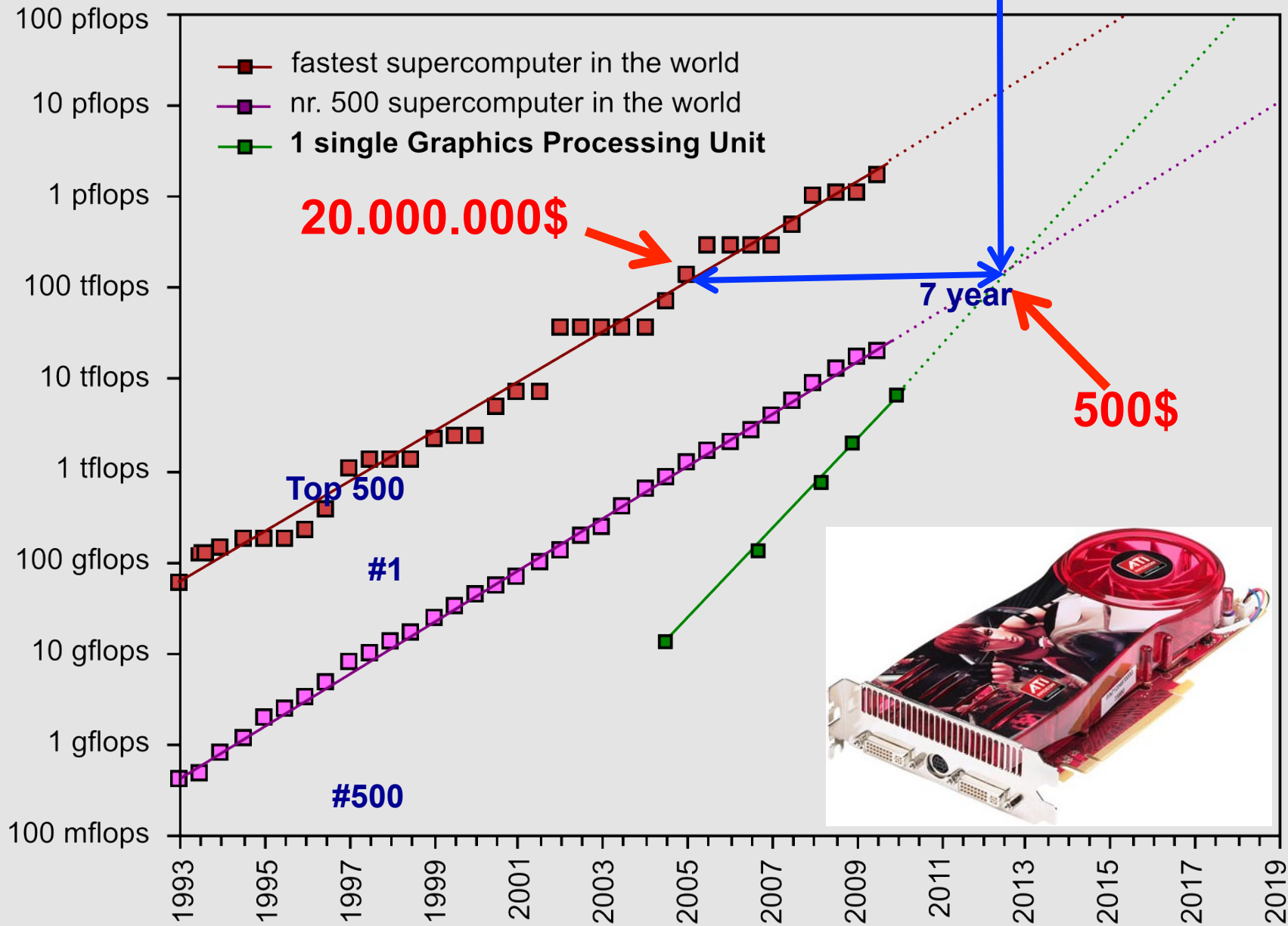
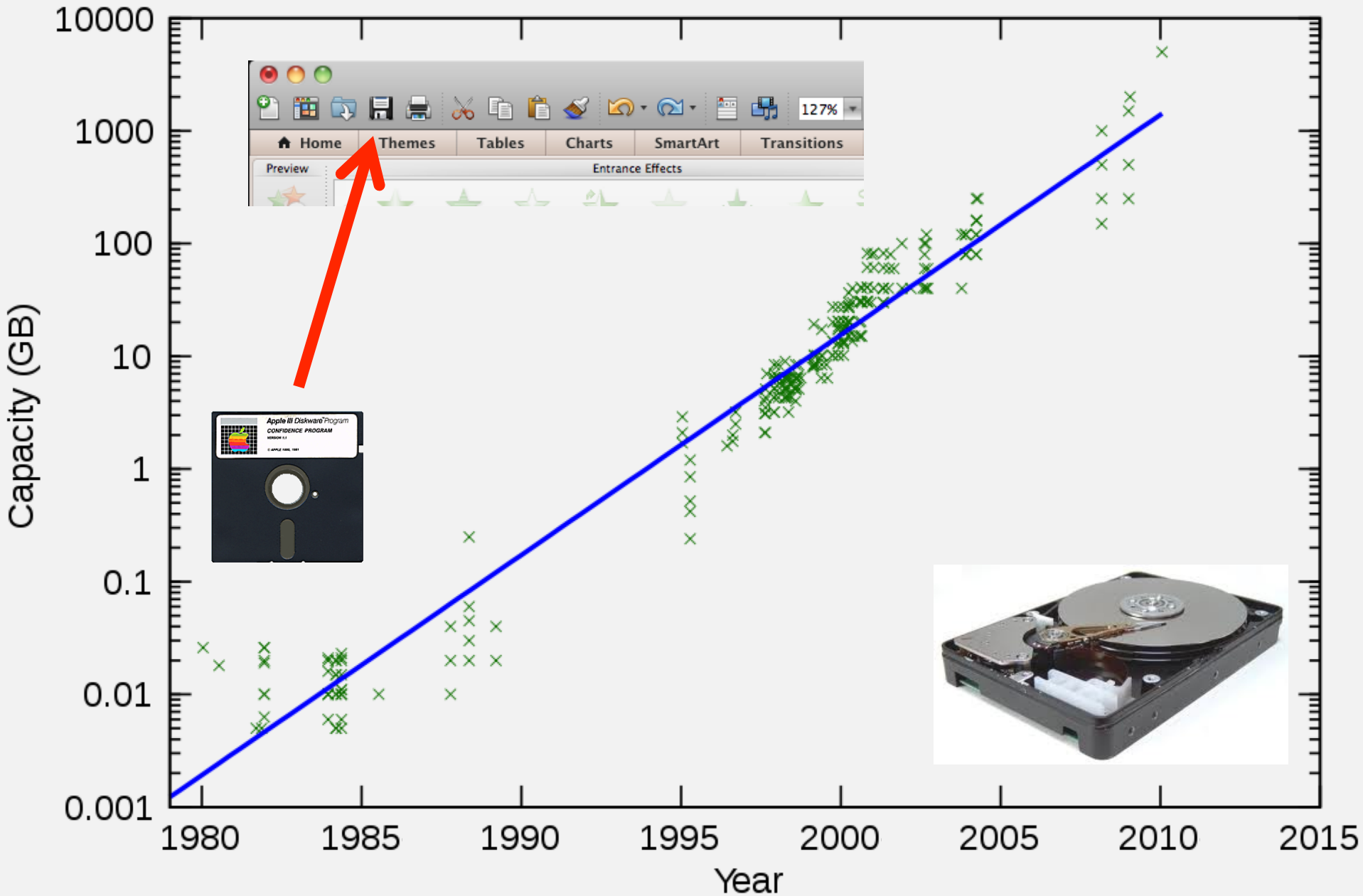## To a network for society

1989: WWW was born

2010



**Internet users as percentage of population**
People with access to the Internet per 100 inhabitants.

Data source: World Bank, World Development Indicators - Last updated December 22, 2010

Netherlands
Denmark
South Korea
United States
Canada
Japan
France

World



IPv4 Exhaustion Counter

▼Present status
Reserved blocks(IANA)

**2%**

7/256 blocks
X-day (estimation)

Jan 20, 2011
Until X-day (estimation)

Today(exhausted?)
Num of IPv4 Address

0(exhausted?)

iNetCore    via IPv4

Ipv6day.nl

# GPU cards are distruptive!
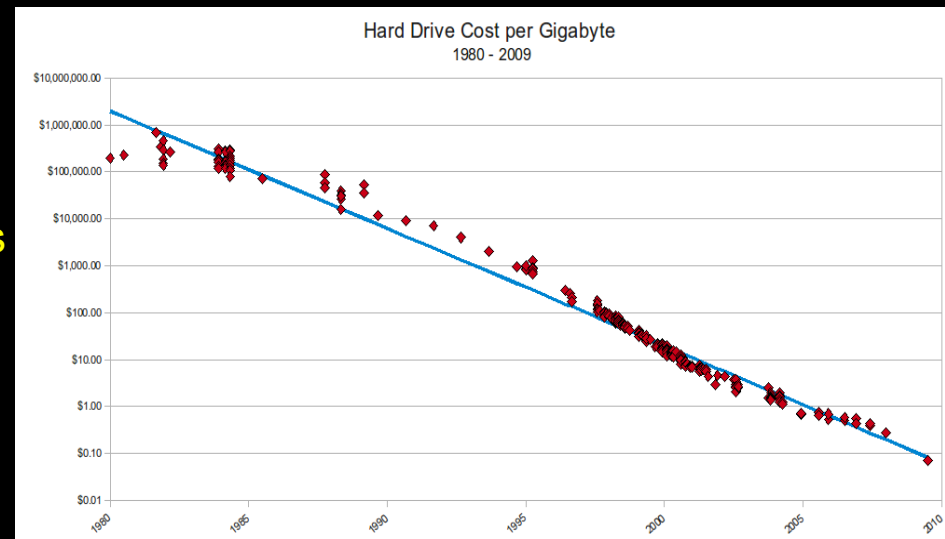
# Data storage: doubling every 1.5 year!

# Reliable and Safe!

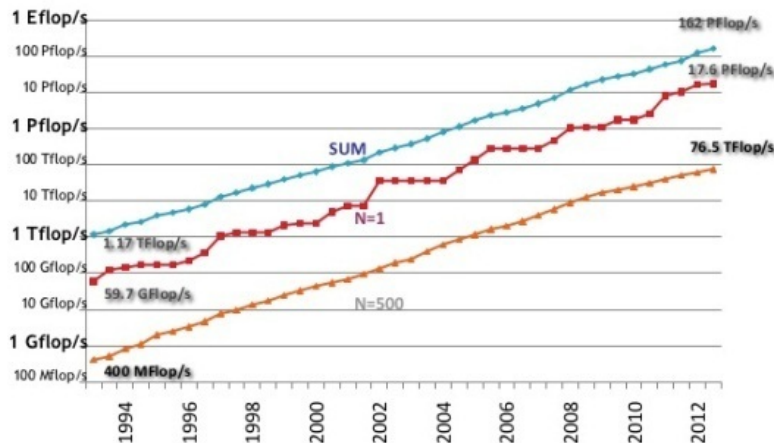This omnipresence of IT makes us not only strong but also vulnerable.

- A virus, a hacker, or a system failure can instantly send digital shockwaves around the world.

The hardware and software that allow all our systems to operate is becoming bigger and more complex all the time, and the capacity of networks and data storage is increasing by leaps and bounds.
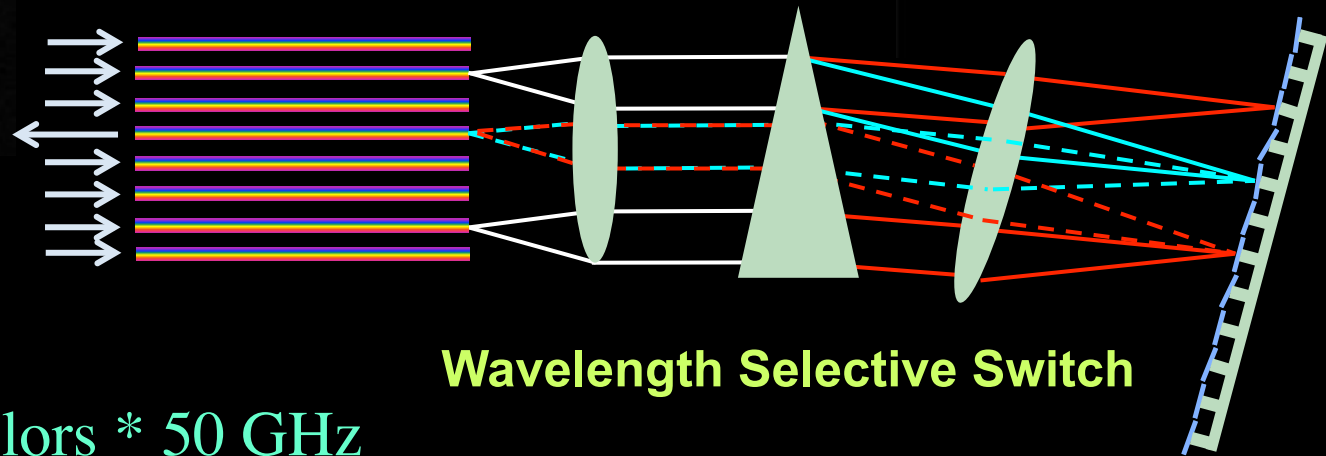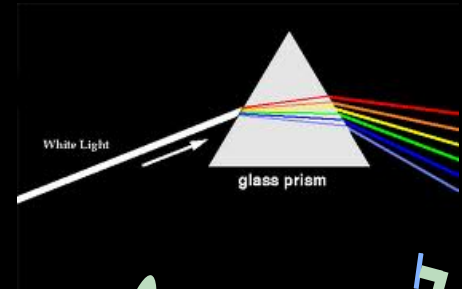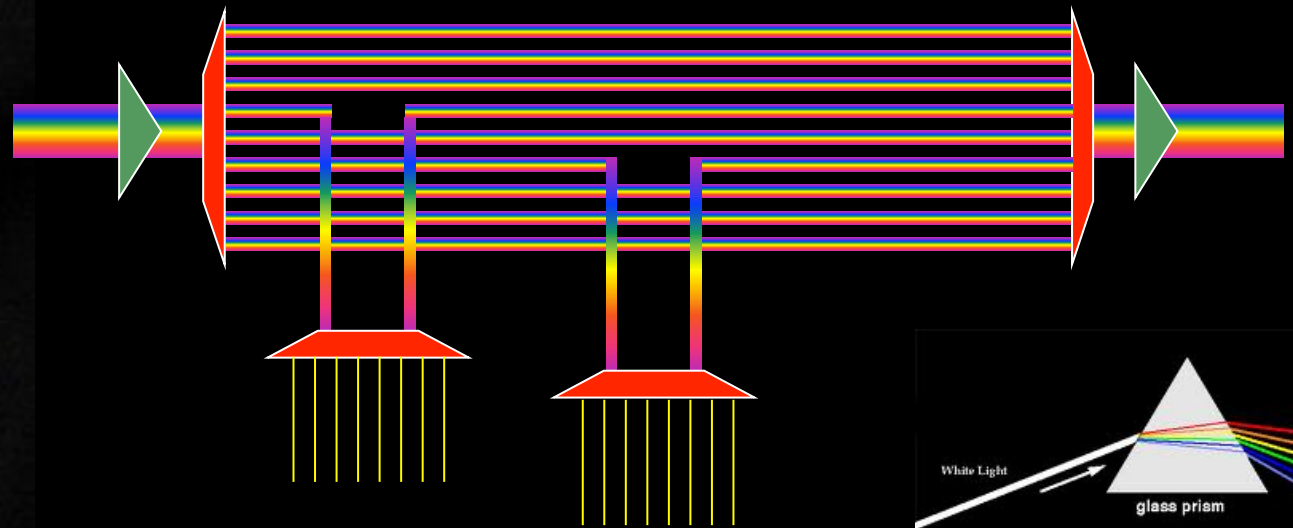


Hard Drive Cost per Gigabyte
1980 - 2009



Performance Development

We will soon reach the limits of what is currently feasible and controllable.

http://www.knaw.nl/Content/Internet_KNAW/publicaties/pdf/20111029.pdf

# Multiple colors / Fiber



**Wavelength Selective Switch**

Per fiber: ~ 80-100 colors * 50 GHz
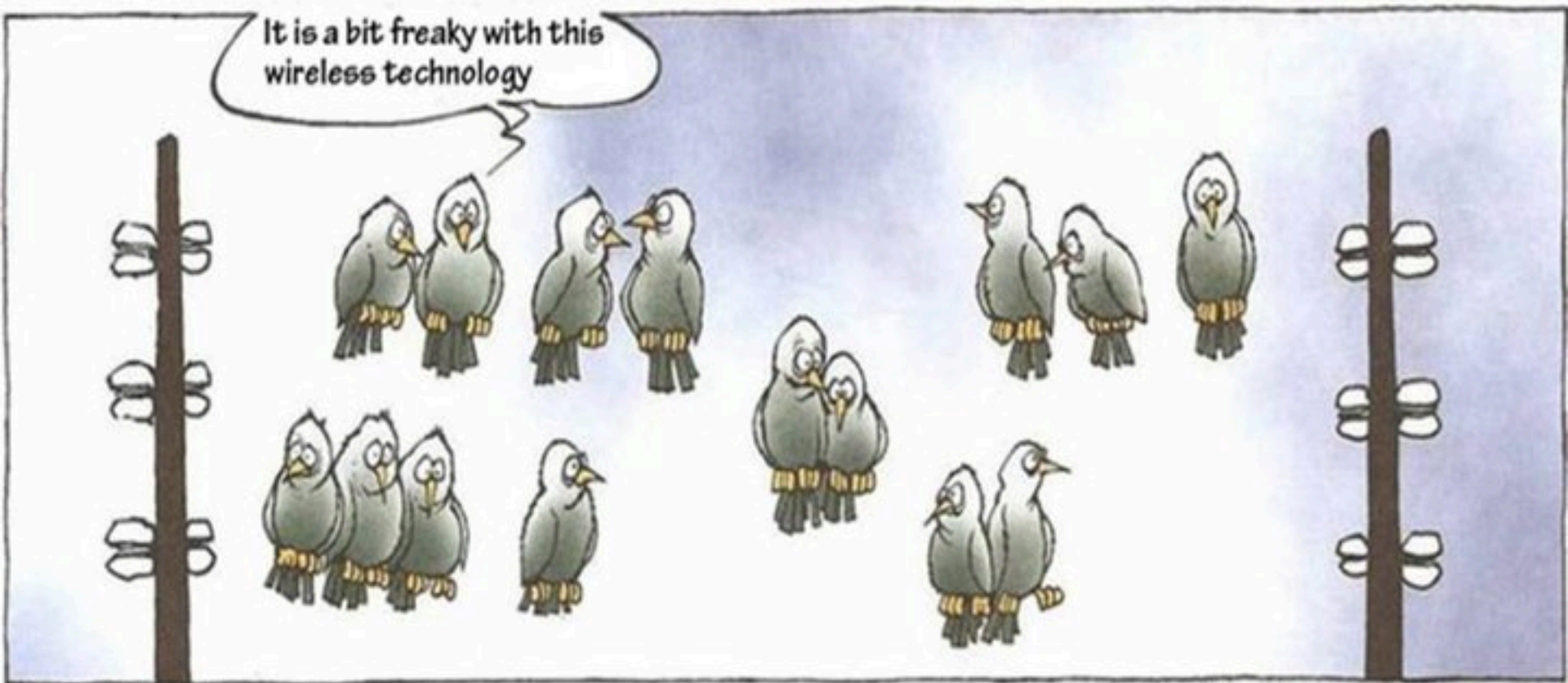
Per color: 10 – 40 – 100 Gbit/s

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

New: Hollow Fiber!

➔ less RTT!

# Wireless Networks



protocol LAN due to the easy comparison and convenience in the **digital home**. While consumer PC products has just started to migrate to a much higher bandwidth of 802.11n wireless LAN now working on next-generation standard definition is already in progress.

# Mission SNE

*Can we create smart and safe data processing infrastructures that can be tailored to diverse application needs?*
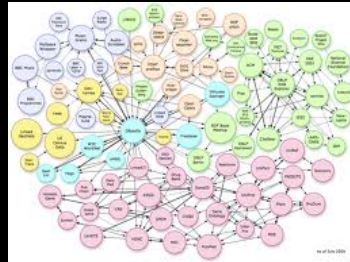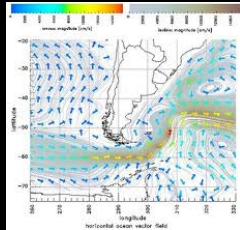
- *Capacity*
  - *Bandwidth on demand, QoS, architectures, photonics, performance*
- *Capability*
  - *Programmability, virtualization, complexity, semantics, workflows*
- *Security*
  - *Authorization, Anonymity, integrity of data in distributed data processing*
- *Sustainability*
  - *Greening infrastructure, awareness*
- *Resilience*
  - *Systems under attack, failures, disasters*

... more data!

Internet developments

DATA

... more realtime!

... more users!

# SNE @ UvA

|  | Ijkdijk/Urban Flood | Medical | LifeWatch | CosmoGrid/eVLBI | EU-GN3/NOVI/Geysers | CineGrid | 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



**Henk** **&** **Ingrid**

# ATLAS detector @ CERN Geneve

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

~PByte/sec

**Online System**

~100 MBytes/sec

**100000 flops/byte**

**10 Pflops/s**

**event simulation**

*Tier 0 +1*

HPSS

**event reconstruction**

**CMS detector: 15m X 15m X 22m**

**12,500 tons, $700M.**

~2.5 Gbits/sec

**Status 2002!**

*Tier 1*

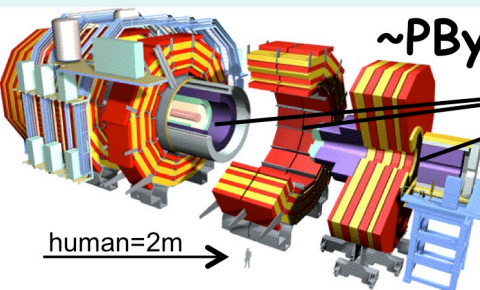| Italian Regional Center | HPSS | German Regional Center | HPSS | NIKHEF Dutch Regional Center | HPSS | FermiLab, USA Regional Center | HPSS | • • • |

~0.6-2.5 Gbps

analysis

Tier2 Center  2 Center  nter  Center  Center   *Tier 2*

~0.6-2.5 Gbps

*Tier 3*

Physics data cache

**Institute ~0.25TIPS**  itute  stitute  Institute

*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.*

100 - 1000 Mbits/sec

*Tier 4*

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

Courtesy Harvey Newman, CalTech and CERN

Workstations

human=2m

Big and small flows don't go well together on the same wire! ☹

# 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**

1 Gbps — MERS

MERS — comp clusters

MERS — Traffic Generators

Content Portal — Streaming Server — 100 TB Storage

*Content Request*

# Experimental Data

**Sage without background traffic**

**Sage with background traffic**



**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

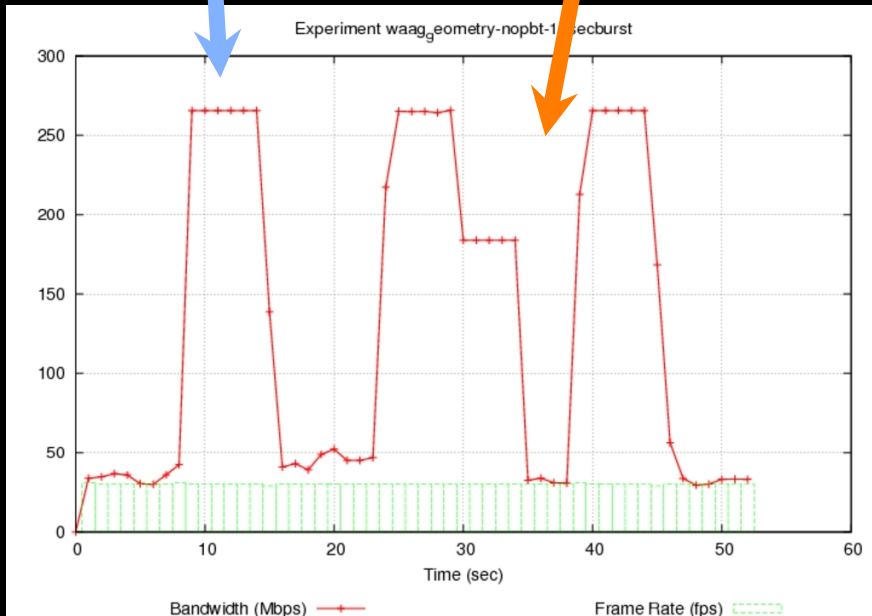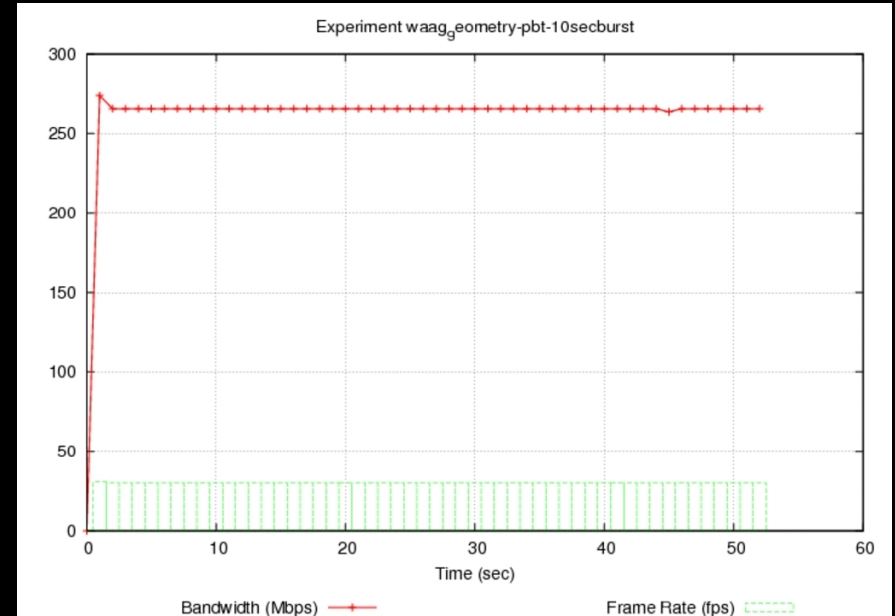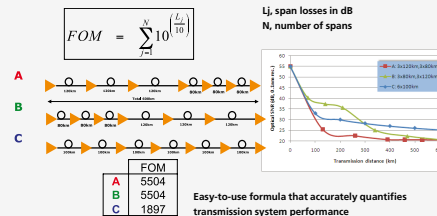- Complex end-to-end optical path engineering in terms of linear (i.e. OSNR, dispersion) and non-linear (FWM, SPM, XPM, Raman) transmission effects for different modulation formats.
- Complex interoperability testing.
- End-to-end monitoring, fault isolation and resolution.
- End-to-end service activation.

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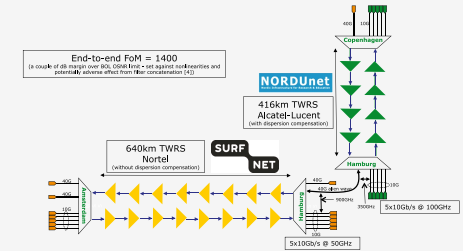
## New method to present fiber link quality, FoM (Figure of Merit)

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.



End-to-end FoM = 1400
(a couple of dB margin over BOL OSNR limit + set against nonlinearities and potentially adverse effect from fiber concatenation [4])

NORDUnet
416km TWRS Alcatel-Lucent (with dispersion compensation)
640km TWRS Nortel (without dispersion compensation)
SURF NET
Copenhagen
Hamburg
5x10Gb/s @ 100GHz
5x10Gb/s @ 50GHz

## 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.

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

From GLIF October 2010 @ CERN

# 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

cpu3

cpu0

Bandwidth (bps)

Bandwidth (bps) for single iperf thread - testcees

CoreID

We used numactl to bind iperf to cores

SNE @ UvA

Columns: Ijkdijk/Urban Flood, Medical, LifeWatch, CosmoGrid/eVLBI, CineGrid, EU-GN3/NOVI/Geysers, SURFnet/GLIF/Cloud

| | Ijkdijk/Urban Flood | Medical | LifeWatch | 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 | |

# e -Very Large Base Interferometer

# eEVN: European VLBI Network



Dec 4 — Dec 5 — Dec 6

Deadline for submitting eVLBI observing proposals

Program committee decides if eVLBI science can be justified

eVLBI Observing Run

Correlation at JIVE

Scientist downloads data from www.jive.nl

12:00 | 18:00 | 24:00 | 06:00 | 12:00 | 18:00 | 24:00 | 06:00 | 12:00

# The SCARIe project

**SCARIe:** a research project to create a Software Correlator for e-VLBI.
**VLBI Correlation:** signal processing technique to get high precision image from spatially distributed radio-telescope.

Telescopes

Input nodes

Correlator nodes

Output node

16 Gbit/s - 2 Tflop ➜

THIS IS A DATA FLOW PROBLEM !!!

Research:



Figure 2. Grid architecture that includes programmable network services.

# LOFAR as a Sensor Network

**20 flops/byte**

– LOFAR is a large distributed research infrastructure:

**2 Tflops/s**

- Astronomy:
  - >100 phased array stations
  - Combined in aperture synthesis array
  - 13,000 small "LF" antennas
  - 13,000 small "HF" tiles
- Geophysics:
  - 18 vibration sensors per station
  - Infrasound detector per station
- >20 Tbit/s generated digitally
- >40 Tflop/s supercomputer
- innovative software systems
  - new calibration approaches
  - full distributed control
  - VO and Grid integration
  - datamining and visualisation

# 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)

7680
4320

UHDTV(4K)

3840
2160

**60°**

1.5 × Picture Height

0.75 × Picture Height   **100°**

CineGrid

Yutaka TANAKA
SHARP CORPORATION
Advanced Image Research Laboratories

# Moving Big Data Objects Globally

- ☐ **Digital Motion Picture for Audio Post-Production**
  - ■ 1 TV Episode Dubbing Reference ~ 1 GB
  - ■ 1 Theatrical 5.1 Final Mix ~ 8 GB
  - ■ 1 Theatrical Feature Dubbing reference ~ 30 GB

- ☐ **Digital Motion Picture Acquisition**
  - ■ 4K RGB x 24 FPS x 10bit/color: ~ 48MB/Frame uncompressed *(ideal)*
  - ■ 6:1 ~ 20:1 shooting ratios => 48TB ~ 160TB digital camera originals

- ☐ **Digital Dailies**
  - ■ HD compressed MPEG-2 @ 25 ~ 50 Mb/s

- ☐ **Digital Post-production and Visual Effects**
  - ■ Gigabytes - Terabytes to Select Sites Depending on Project

- ☐ **Digital Motion Picture Distribution**
  - ■ Film Printing in Regions
    - ☐ Features ~ 8TB
    - ☐ Trailers ~ 200GB
  - ■ Digital Cinema Package to Theatres
    - ☐ Features ~ 100 - 300GB per DCP
    - ☐ Trailers ~ 2 - 4GB per DCP

Yesterday's Media Transport Method!

8 TByte

# What Happens in an **Internet Minute?**

639,800 GB of global IP data transferred

**20** New victims of identity theft

**47,000** App downloads

**61,141** Hours of music

**204 million** Emails sent

**$83,000** In sales

**20 million** Photo views

**3,000** Photo uploads

**320+** New Twitter accounts

**100,000** New tweets

**135** Botnet infections

**6** New Wikipedia articles published

**1,300** New mobile users

**100+** New Linkedin accounts

**277,000** Logins

**6 million** Facebook views

**2+ million** Search queries

**30** Hours of video uploaded

**1.3 million** Video views

## And **Future Growth** is **Staggering**

**Today**, the number of **networked devices** = the global population

By **2015**, the number of **networked devices** = **2x** the global population

In **2015**, it would take you **5 years** to view all video crossing IP networks each **second**

IP

(intel)

There

is

always

a

bigger

fish

Size of data sets in terabytes

| | |
|---|---|
| Business email sent per year ...................................2,986,100 | National Climactic Data Center database.........................6,144 |
| Content uploaded to Facebook each year.................182,500 | Library of Congress' digital collection............................5,120 |
| Google's search index ..............................................97,656 | US Census Bureau data ..................................................3,789 |
| Kaiser Permanente's digital health records ...............30,720 | Nasdaq stock market database .......................................3,072 |
| Large Hadron Collider's annual data output ................15,360 | Tweets sent in 2012.............................................................19 |
| Videos uploaded to YouTube per year ........................15,000 | Contents of every print issue of WIRED ...........................1.26 |

# The GLIF – LightPaths around the World

F Dijkstra, J van der Ham, P Grosso, C de Laat, "A path finding implementation for multi-layer networks", Future Generation Computer Systems 25 (2), 142-146.



GLIF Map 2011: Global Lambda Integrated Facility    Visualization by Robert Patterson, NCSA, University of Illinois at Urbana–Champaign    Data Compilation by Maxine D. Brown, University of Illinois at Chicago    Texture Retouch by Jeff Carpenter, NCSA    Earth Texture, visibleearth.nasa.gov    www.glif.is

# The GLIF – LightPaths around the World

F Dijkstra, J van der Ham, P Grosso, C de Laat, "A path finding implementation for multi-layer networks", Future Generation Computer Systems 25 (2), 142-146.

# The GLIF – LightPaths around the World



We investigate: TomTom for complex networks!

# 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):

Predicate

**Subject** → **Object**

Subject → Object Subject → Object Subject → Object Subject

Object Subject → Object Subject

| Location | Device | Interface | Link |
|----------|--------|-----------|------|
| name → | description → | locatedAt → | hasInterface → |
| connectedTo → | capacity → | encodingType → | 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
    <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



IP-layer

Ethernet layer

STS — layer

OC-192 — layer

UTP — layer

fiber — layer

**End host** — **SONET switch with Ethernet intf.** — **Ethernet & SONET switch** — **SONET switch** — **SONET switch with Ethernet intf.** — **End host**

Université du Quebec — CA★Net Canada — StarLight Chicago — MAN LAN New York — NetherLight Amsterdam — Universiteit van Amsterdam

# Multi-layer Network PathFinding



Path between interfaces A1 and E1:
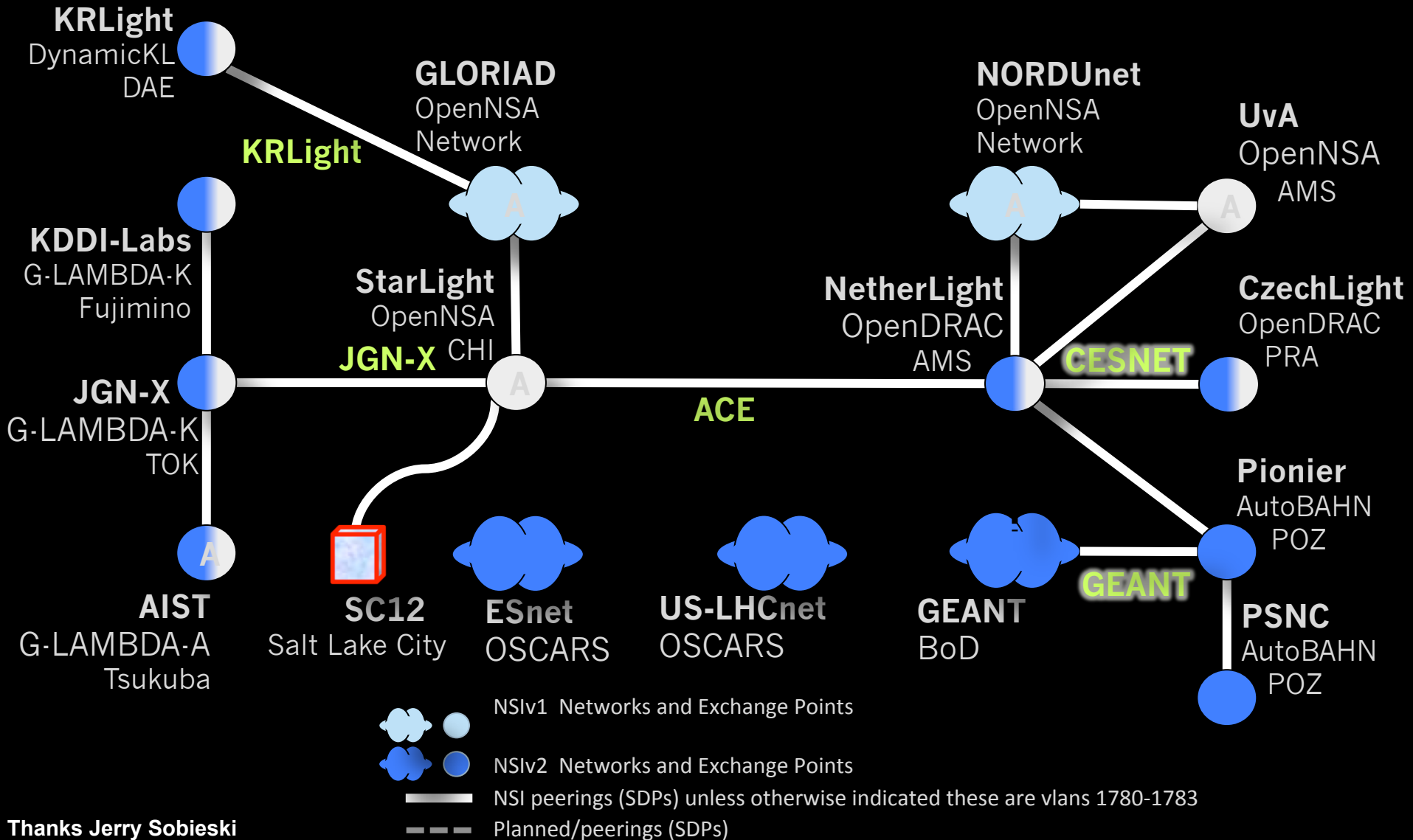A1-A2-B1-B4-D4-D2-C3-C4-C1-C2-B2-B3-D3-D1-E2-E1
Scaling: Combinatorial problem

# Automated GOLE + NSI
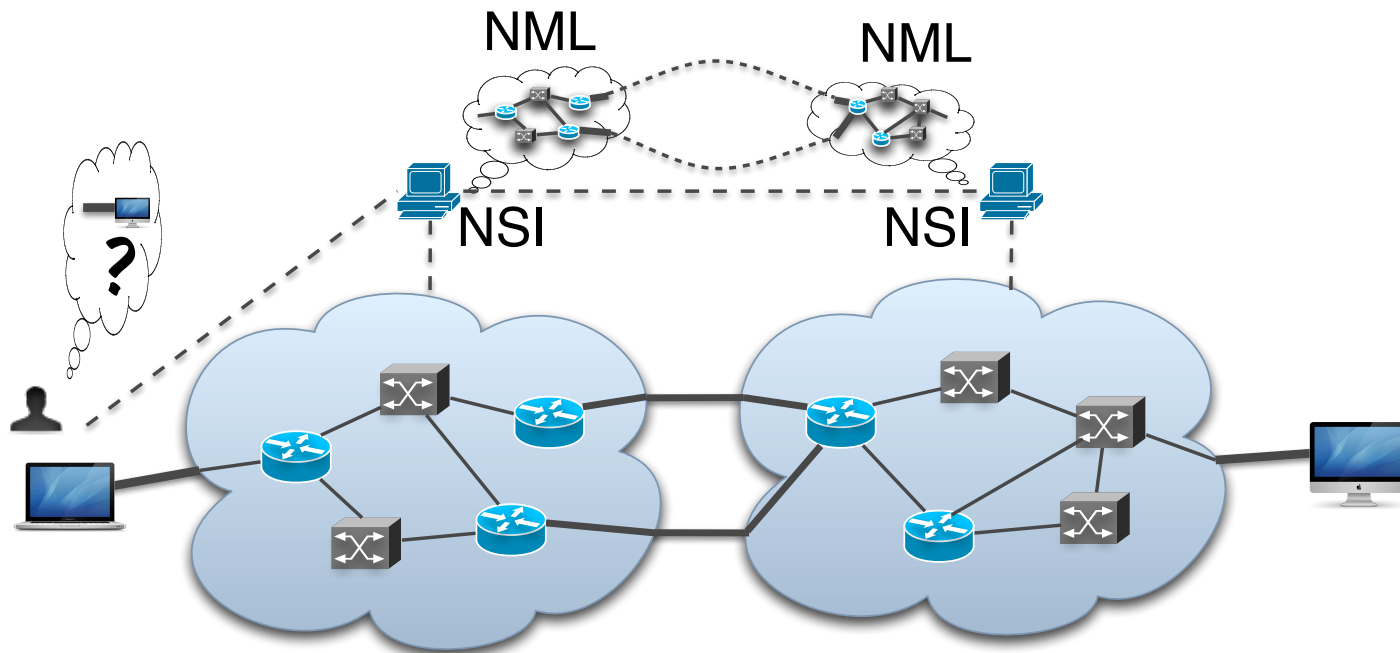
## Joint NSI v1+v2 Beta Test Fabric     Nov 2012
## Ethernet Transport Service

**KRLight**
DynamicKL
DAE

**KRLight**

**GLORIAD**
OpenNSA
Network

**NORDUnet**
OpenNSA
Network

**UvA**
OpenNSA
AMS

**KDDI-Labs**
G-LAMBDA-K
Fujimino

**StarLight**
OpenNSA
CHI

**JGN-X**

**NetherLight**
OpenDRAC
AMS

**CESNET**

**CzechLight**
OpenDRAC
PRA

**JGN-X**
G-LAMBDA-K
TOK

**ACE**

**Pionier**
AutoBAHN
POZ

**AIST**
G-LAMBDA-A
Tsukuba

**SC12**
Salt Lake City

**ESnet**
OSCARS

**US-LHCnet**
OSCARS

**GEANT**
BoD

**GEANT**

**PSNC**
AutoBAHN
POZ

NSIv1  Networks and Exchange Points

NSIv2  Networks and Exchange Points

NSI peerings (SDPs) unless otherwise indicated these are vlans 1780-1783

Planned/peerings (SDPs)

**Thanks Jerry Sobieski**
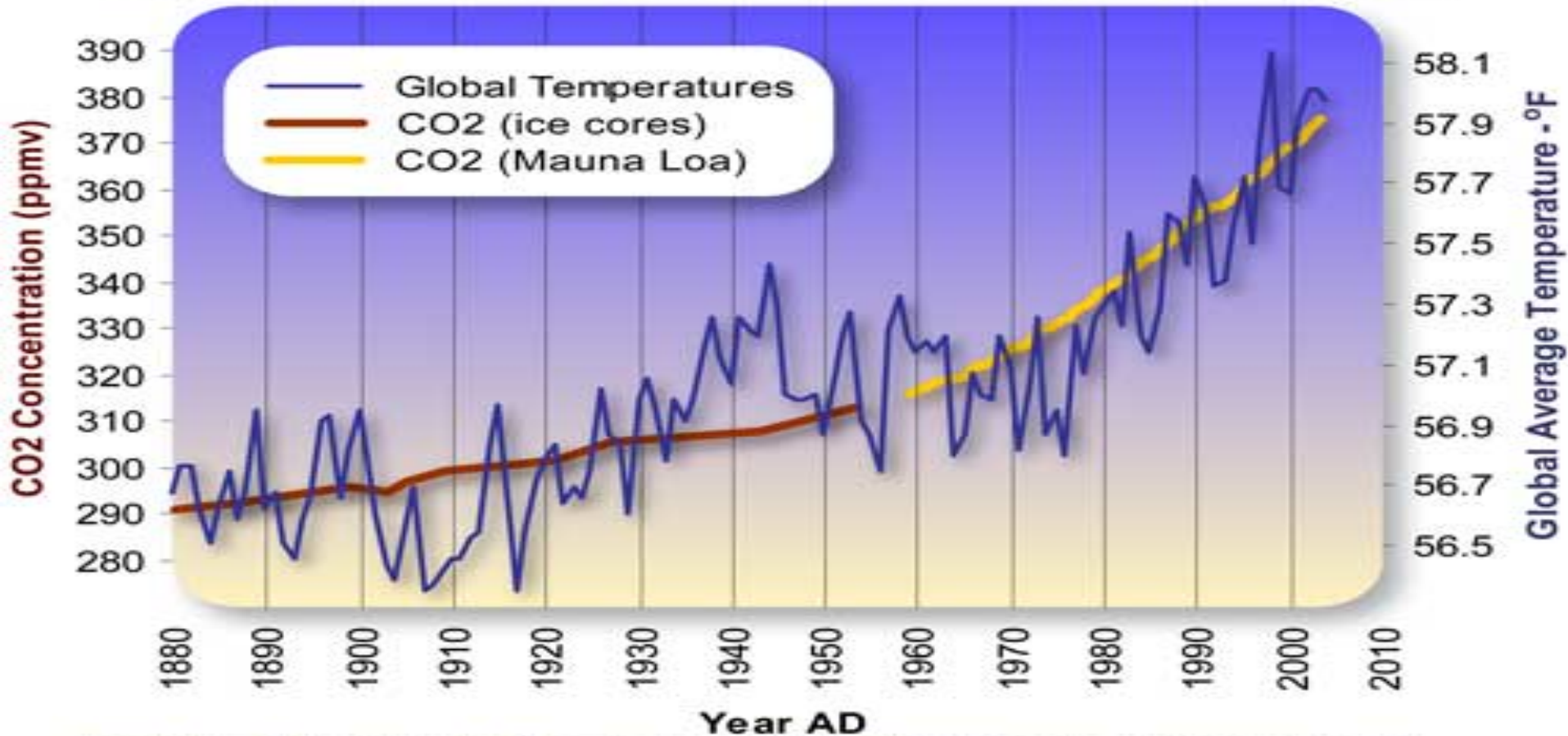
# Network Topology Description

Network topology research supporting automatic network provisioning
- Inter-domain networks
- Multiple technologies
- Based on incomplete information
- Possibly linked to other resources



http://redmine.ogf.org/projects/nml-wg
http://redmine.ogf.org/projects/nsi-wg

http://sne.science.uva.nl/ndl

# Need for GreenIT



Global Average Temperature and Carbon Dioxide Concentrations, 1880 – 2004

Data Source Temperature: ftp://ftp.ncdc.noaa.gov/pub/data/anomalies/annual_land.and.ocean.ts
Data Source CO2 (Siple Ice Cores): http://cdiac.esd.ornl.gov/ftp/trends/co2/siple2.013
Data Source CO2 (Mauna Loa): http://cdiac.esd.ornl.gov/ftp/trends/co2/maunaloa.co2

Graphic Design: Michael Ernst, The Woods Hole Research Center

# Greening the Processing System

# Turn Green Tech into Greenbacks
## IT Certifications for Jobs That Make a Difference

**ENERGY STAR**

green tech

## Uptime Institute Accredited Tier Designer

The Uptime Institute has long been a proponent for green data center design and implementation. Its certification course on data center design embeds green principles into the curriculum.

**CERTIFIED GREEN COMPUTING USER SPECIALIST**

**GCI** TM
Green Computing Initiative

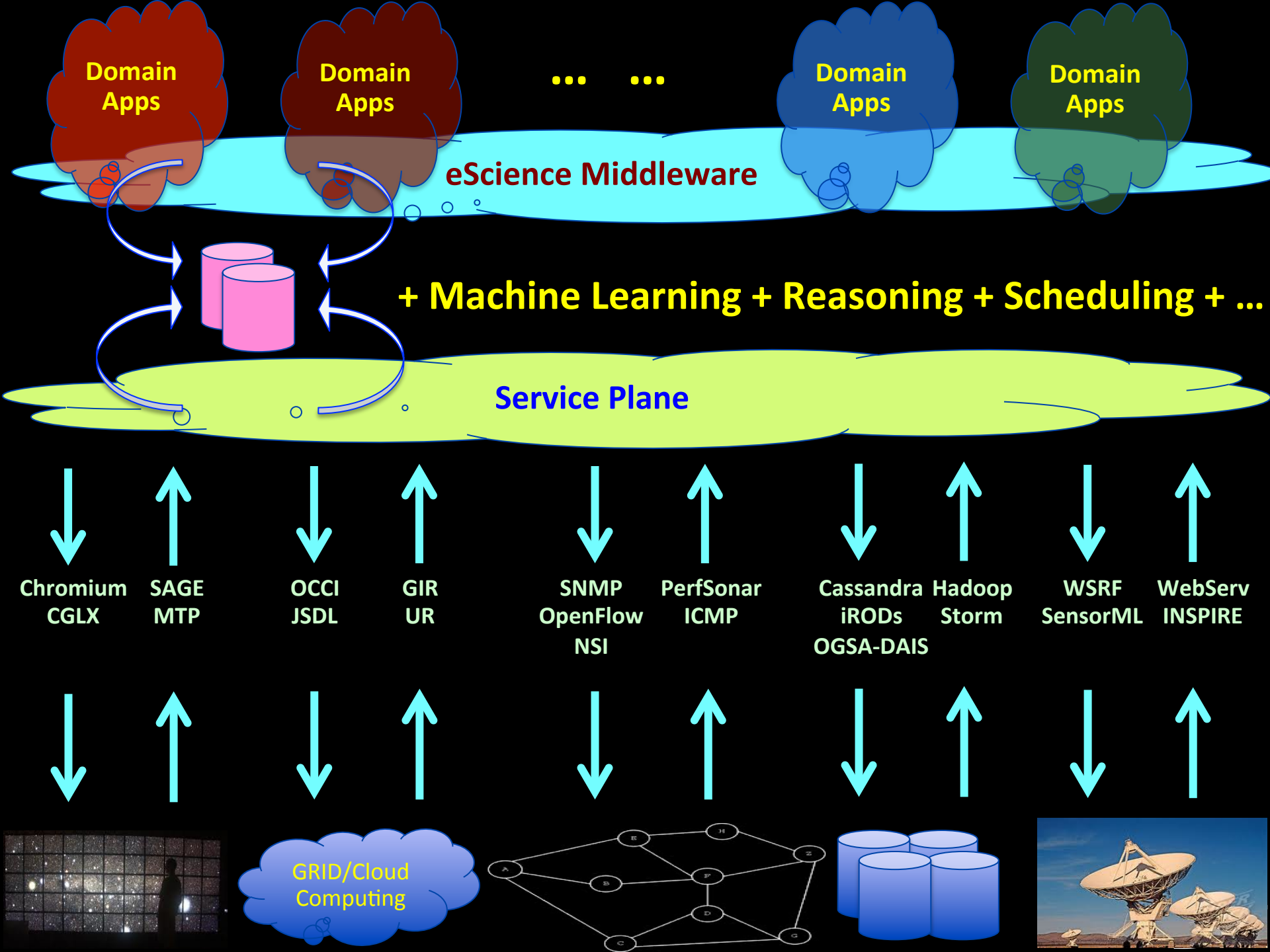**SUSTAINABILITY**
**Your Career**

# ECO-Scheduling

I want to

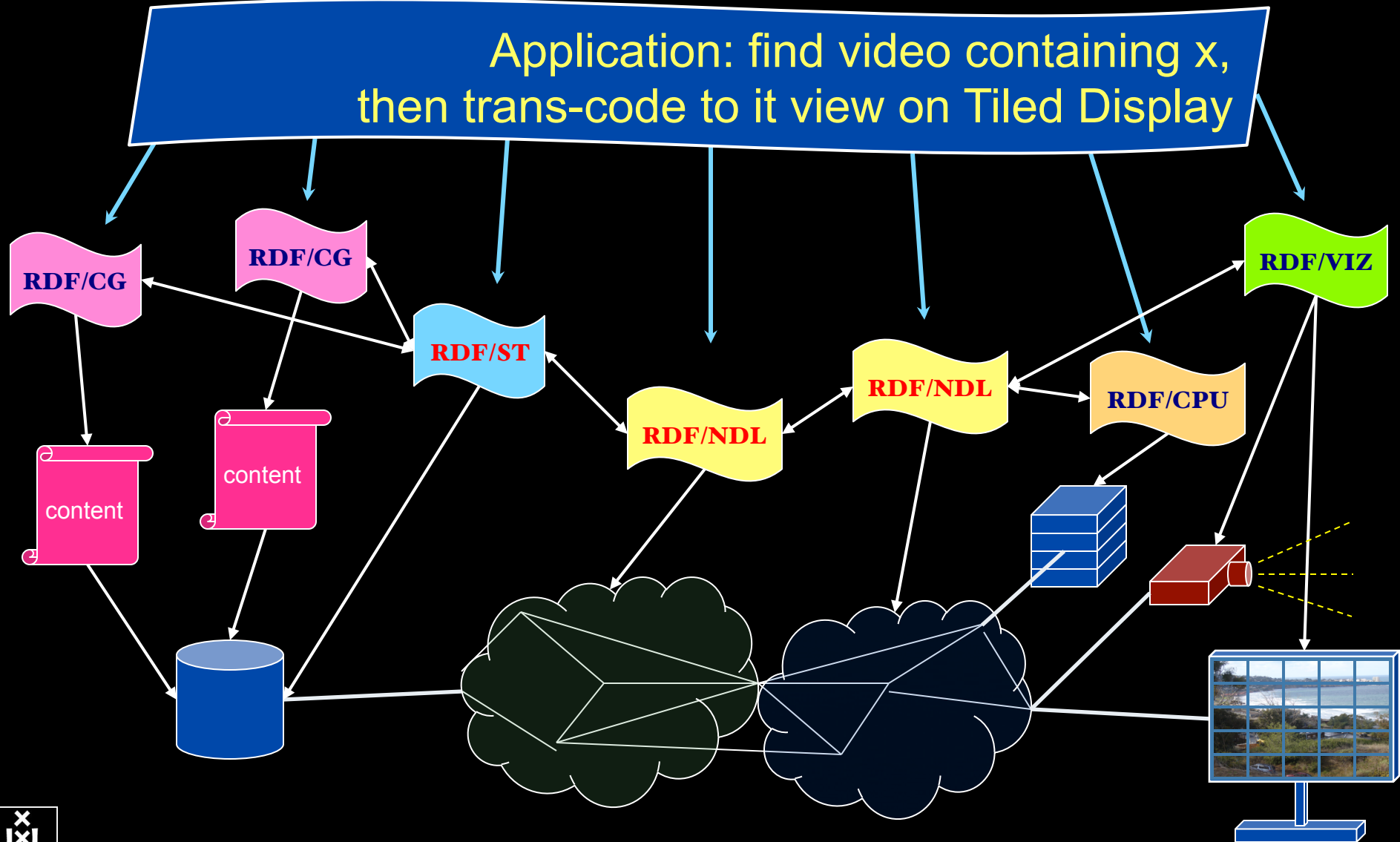**"Show Big Bug Bunny in 4K on my Tiled Display using green Infrastructure"**

- 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.
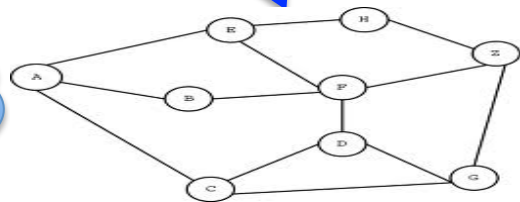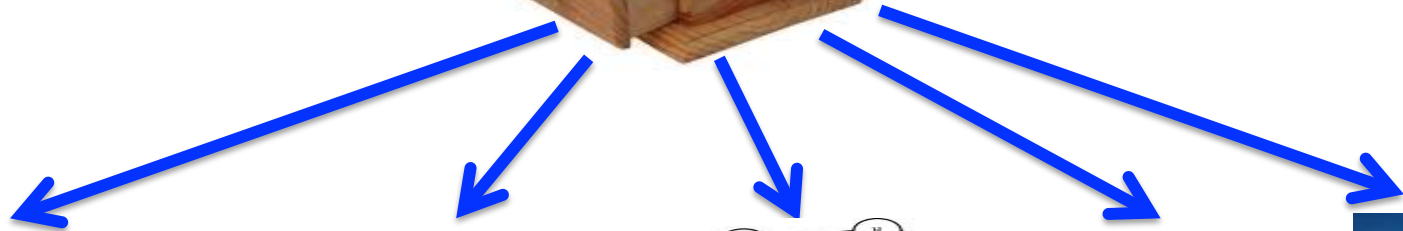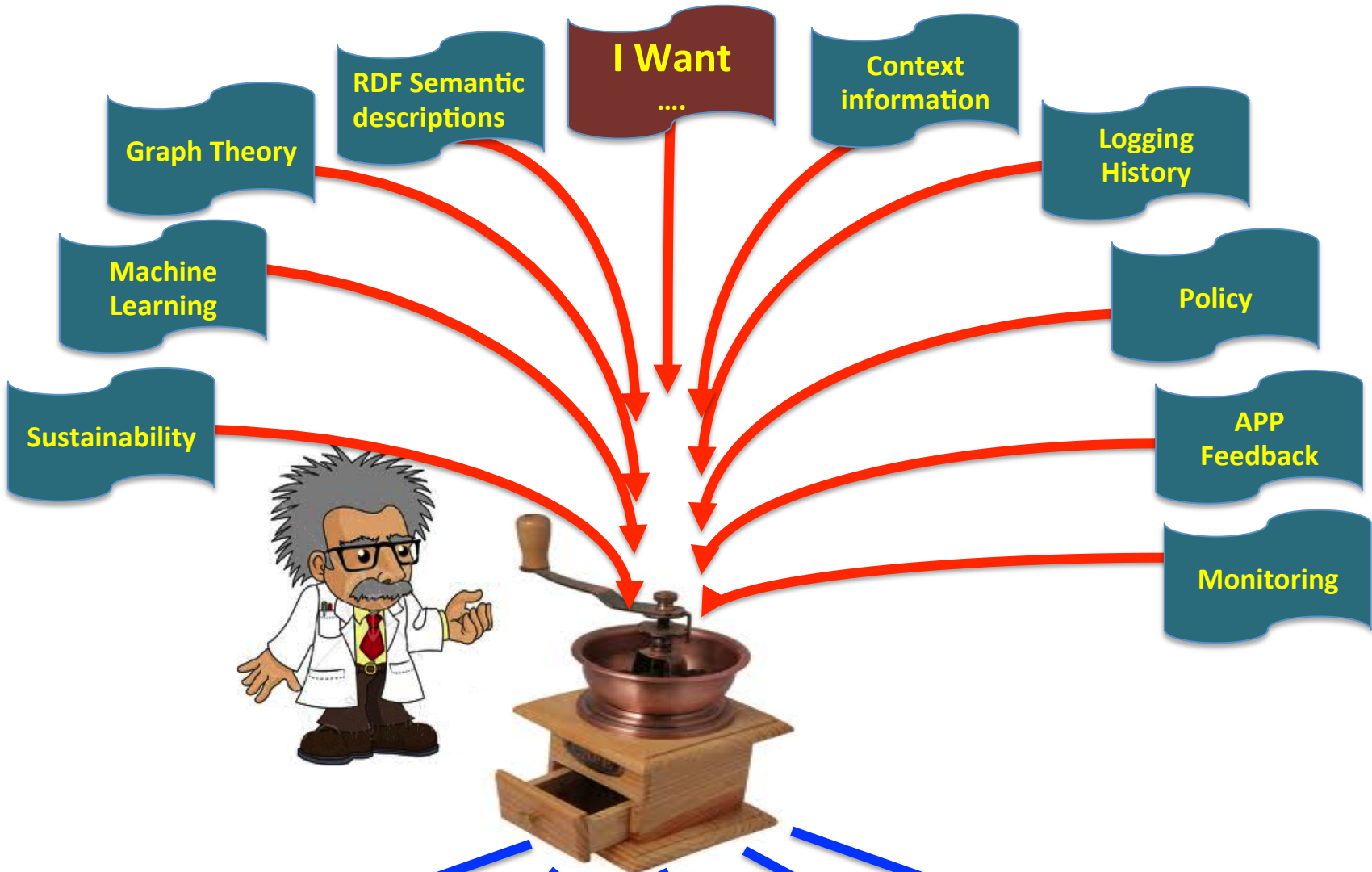  - ➔ His refrigerator also just works.

# RDF describing Infrastructure "I want"



**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/CPU

RDF/NDL

content

content

Graph Theory

RDF Semantic descriptions

I Want ….

Context information

Logging History

Machine Learning

Policy

Sustainability

APP Feedback

Monitoring

Cloud Computing

# Paper #1 + Q's

*This global experiment wants to see if high-end applications needing transport capacities of multiple Gbps for up to hours at a time can be handled through an optical bypass network.*

Tom DeFanti, Cees de Laat, Joe Mambretti, Kees Neggers, Bill St. Arnaud.

# Paper #1 + Q's

- Q1: This article is now 10 years old. Back then Twitter did not exist.  What do you think will be the drivers for network capacity demand in Science and Society 10 years from now?

- Q2:  List arguments why one would use photonic networks directly in science applications and arguments why not tu use photonics directly but use current Internet.

- Q3: This question is not directly from this paper but fun to figure out via search on the web: Fiber cable systems under the ocean are very expensive and cost 100's of millions to put in place. How many fibers do they put in one cable and why that amount?

# Paper #2 + Q's

A distributed topology information system for optical networks based on the semantic web.

Jeroen van der Ham, Freek Dijkstra, Paola Grosso, Ronald van der Pol, Andree Toonk, Cees de Laat

http://delaat.net/pubs/2008-j-4.pdf

# Paper #2 + Q's

- Q1: Suppose this method of describing networks is a total worldwide success and allows to find superfast networking paths through the CI (CyberInfrastructure). The question becomes: Does it scale? Can you find reasons why and/or why not it could scale up to the size of the internet?

- Q2: Are the described methods and framework fault tolerant? If not, then list the issues in your view. What do you see best ways to do something about it.

- Q3: List advantages of NDL, or more generically, using semantic web methods for describing cyber infrastructure?

# The constant factor in our field is Change!

The 50 years it took Physicists to find one particle, the Higgs, we came from:

"Fortran goto", Unix, c, SmallTalk, DECnet, TCP/IP, c++, Internet, WWW, Semantic Web, Photonic networks, Google, grid, cloud, Data^3, App

to:

DDOS attacks destroying Banks and Bitcoins.

Conclusion:

Need for Safe, Smart, Resilient Sustainable Infrastructure.