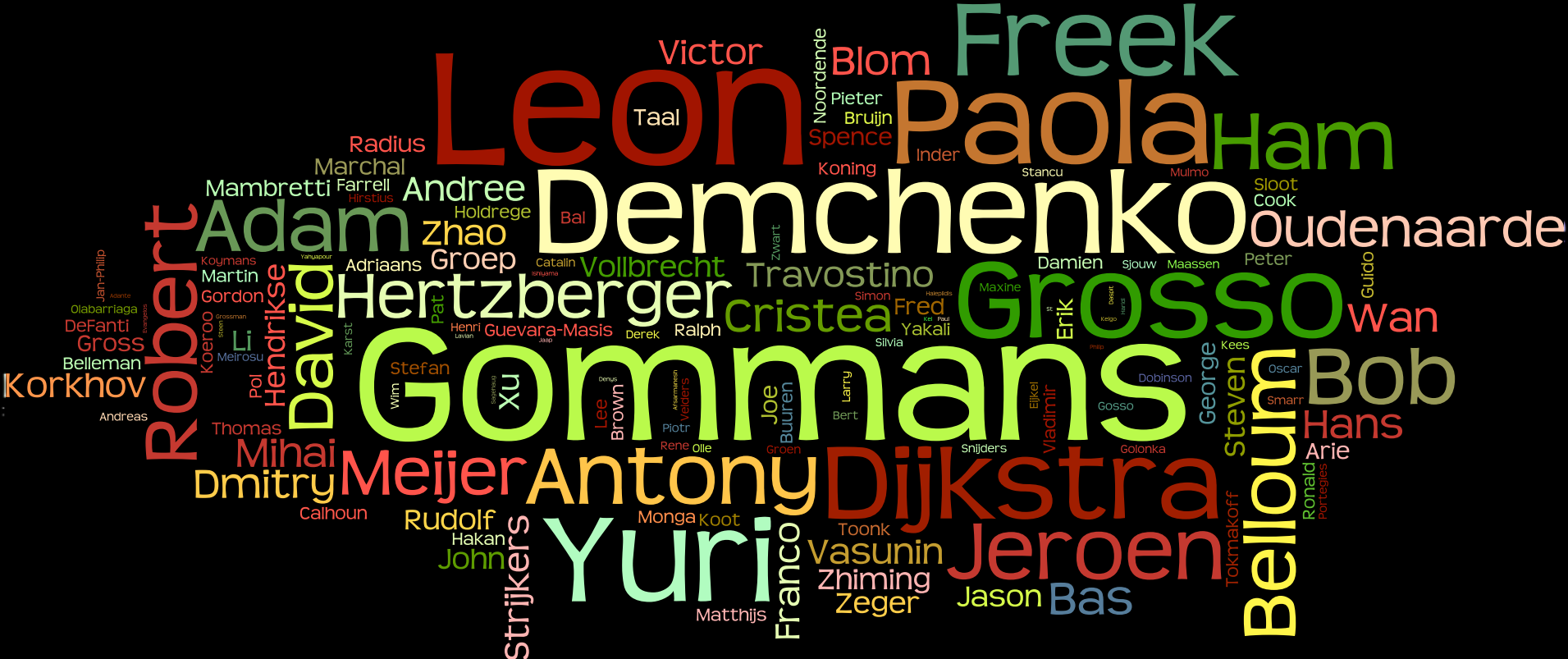


System and Network Engineering Research for Big Data Sciences

Cees de Laat



Mission SNE

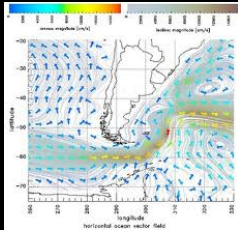
Can we create smart and safe data processing infrastructures that can be tailored to diverse application needs?

- *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*

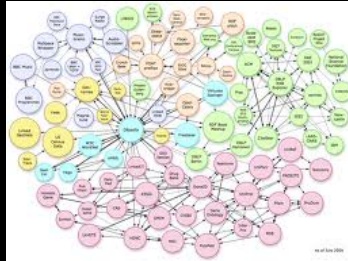
... more data!

Internet developments

Google



DATA



... more realtime!



twitter



myspace
a place for freedom



LinkedIn



SchoolBANK

Hyves

flickr
from YAHOO!



... more users!

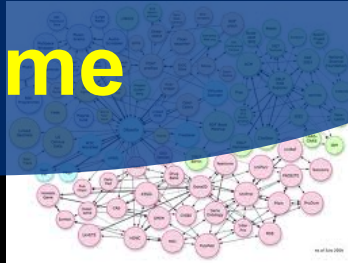
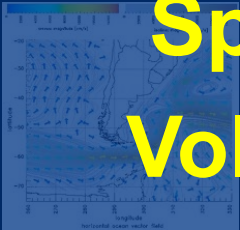
... more data!

Internet developments

Google

Speed
Volume

DATA



Deterministic

Real-time



twitter



Scalable

Secure

Linked in



myspace
SchoolBANK

Hyves

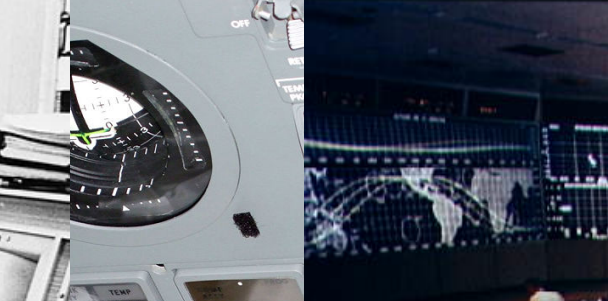
flickr
from YAHOO!



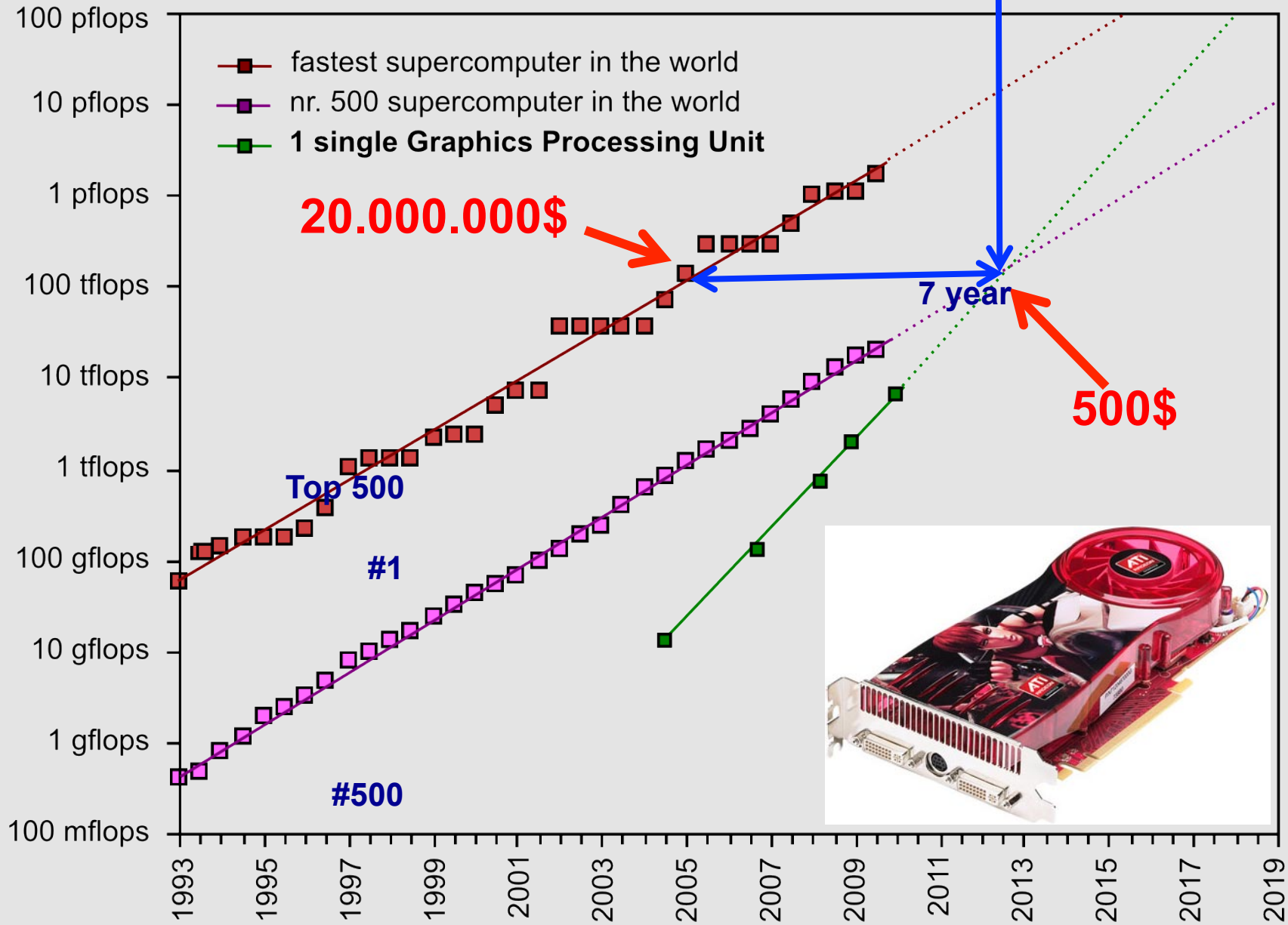
... more users!



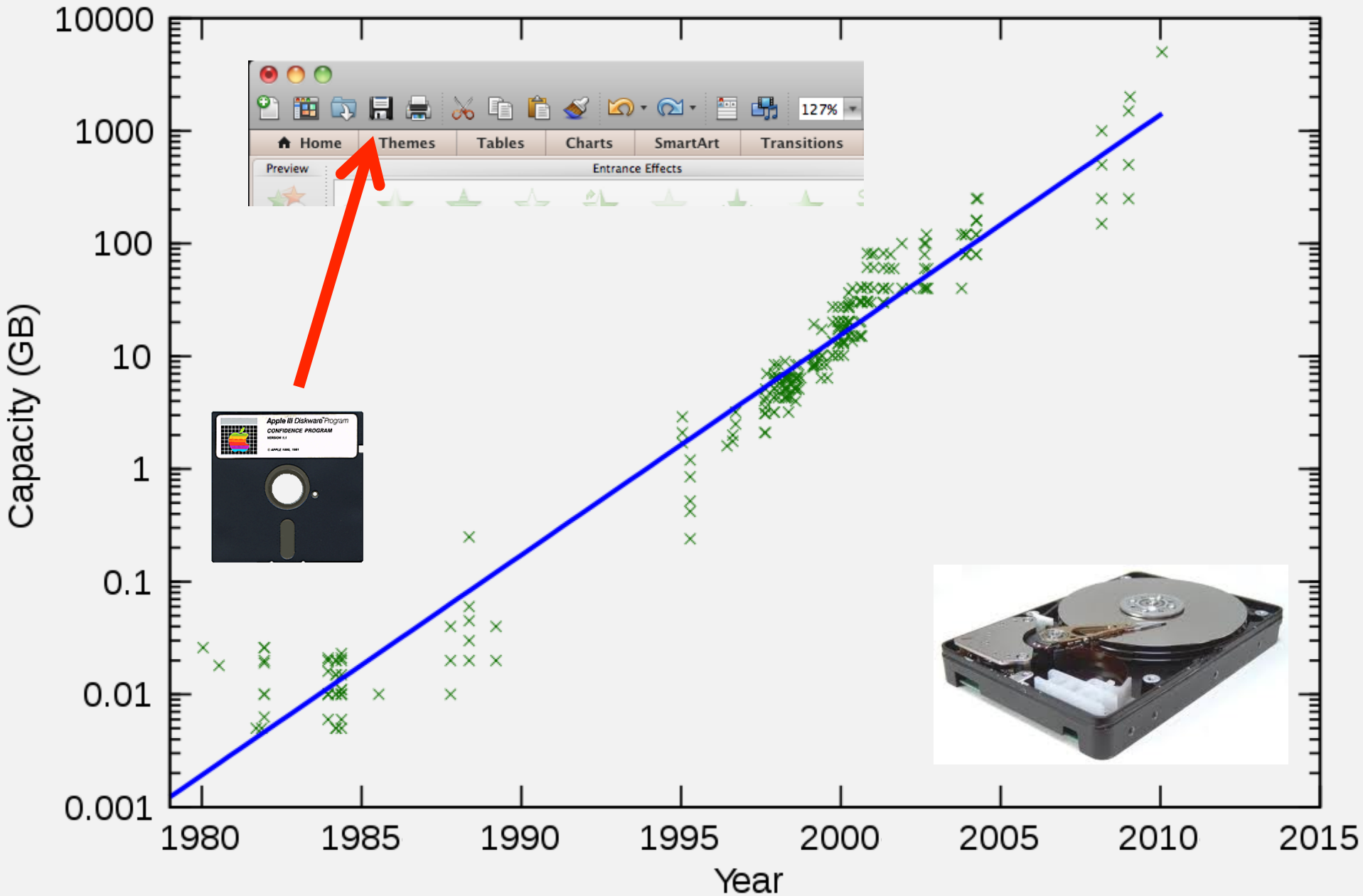




GPU cards are disruptive!

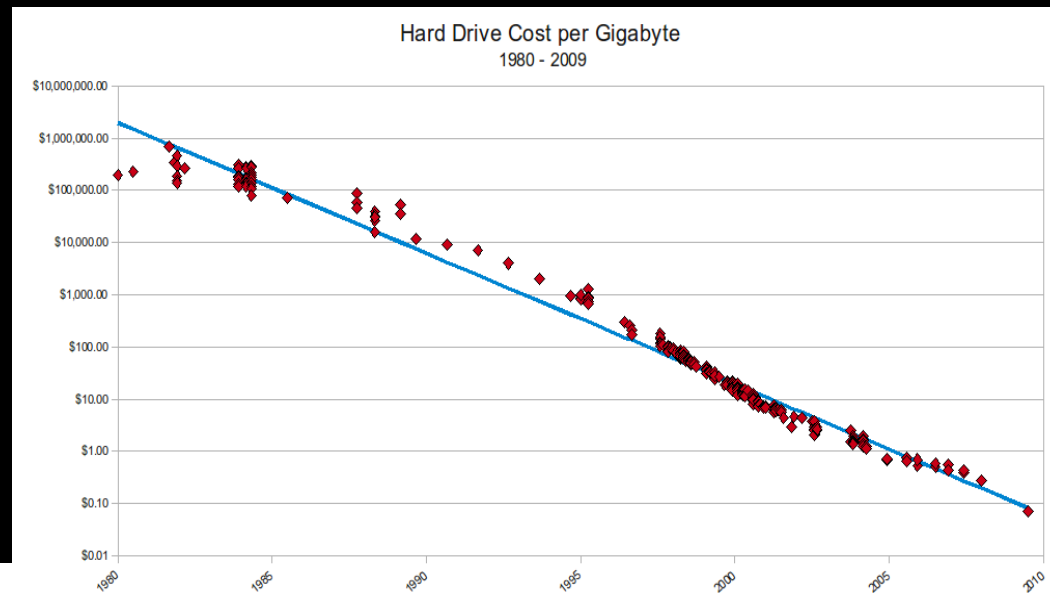


Data storage: doubling every 1.5 year!

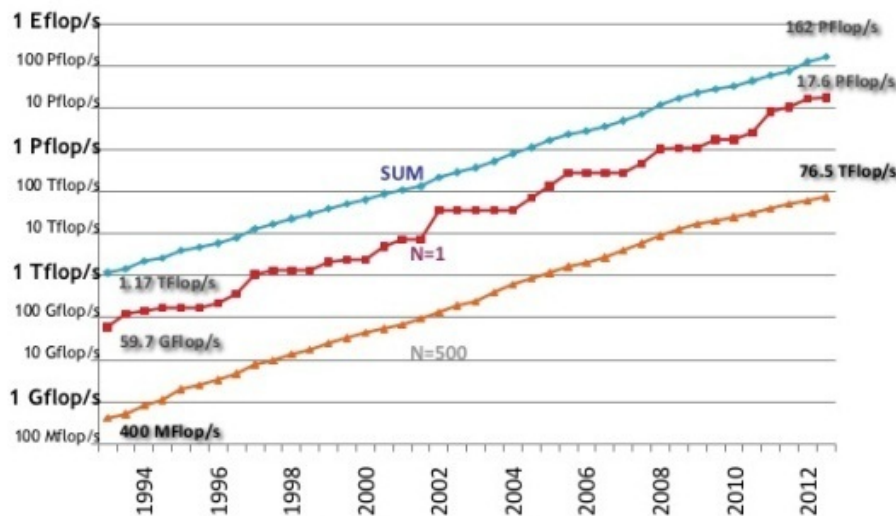


Computing vs Data

Computing per unit cost has doubled roughly every 18 months.



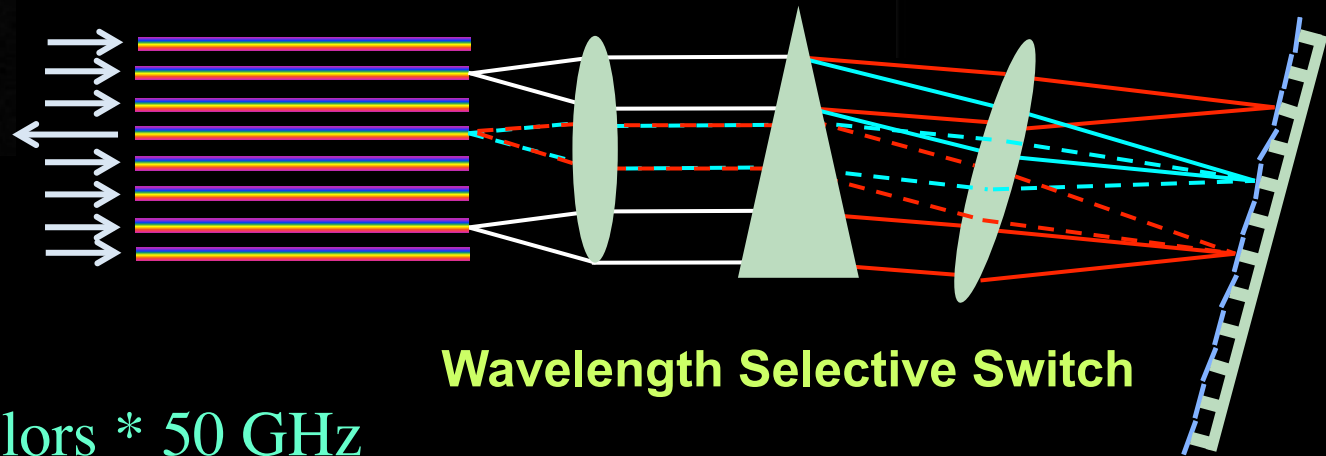
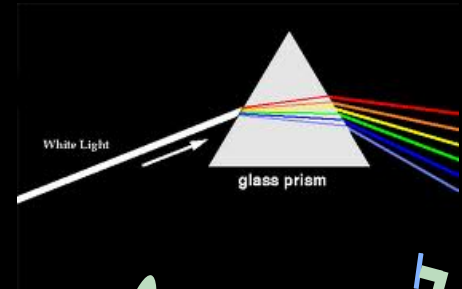
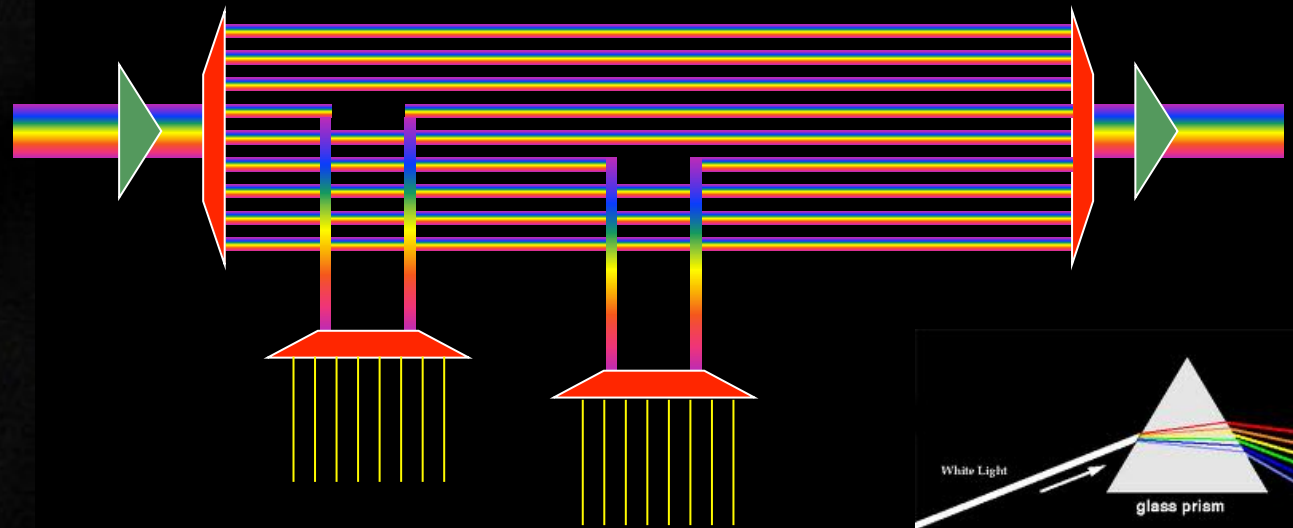
Performance Development



Space per unit cost has doubled roughly every 14 months.

So: data becomes exponentially uncomputable.

Multiple colors / Fiber



Wavelength Selective Switch

Per fiber: $\sim 80-100$ colors * 50 GHz

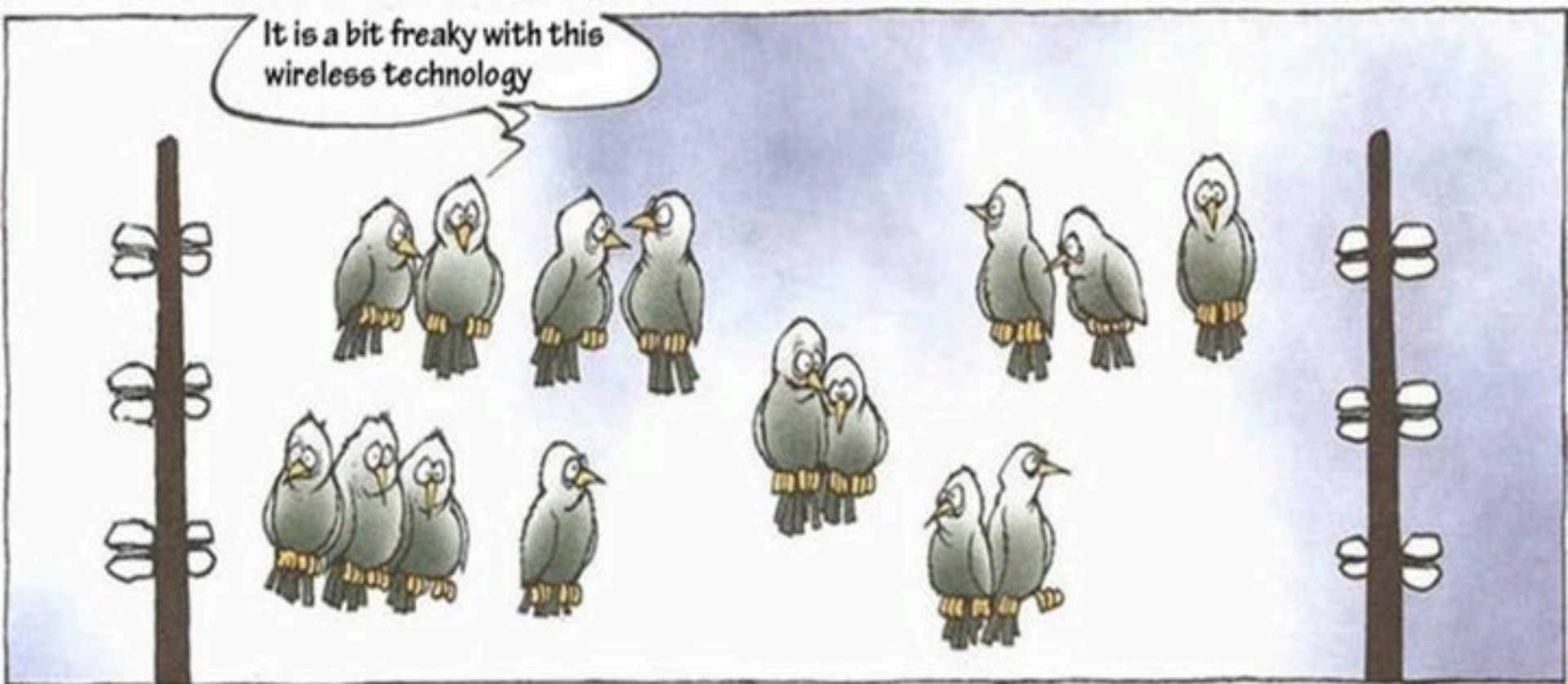
Per color: 10 – 40 – 100 Gbit/s

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

New: Hollow Fiber!

➔ less RTT!

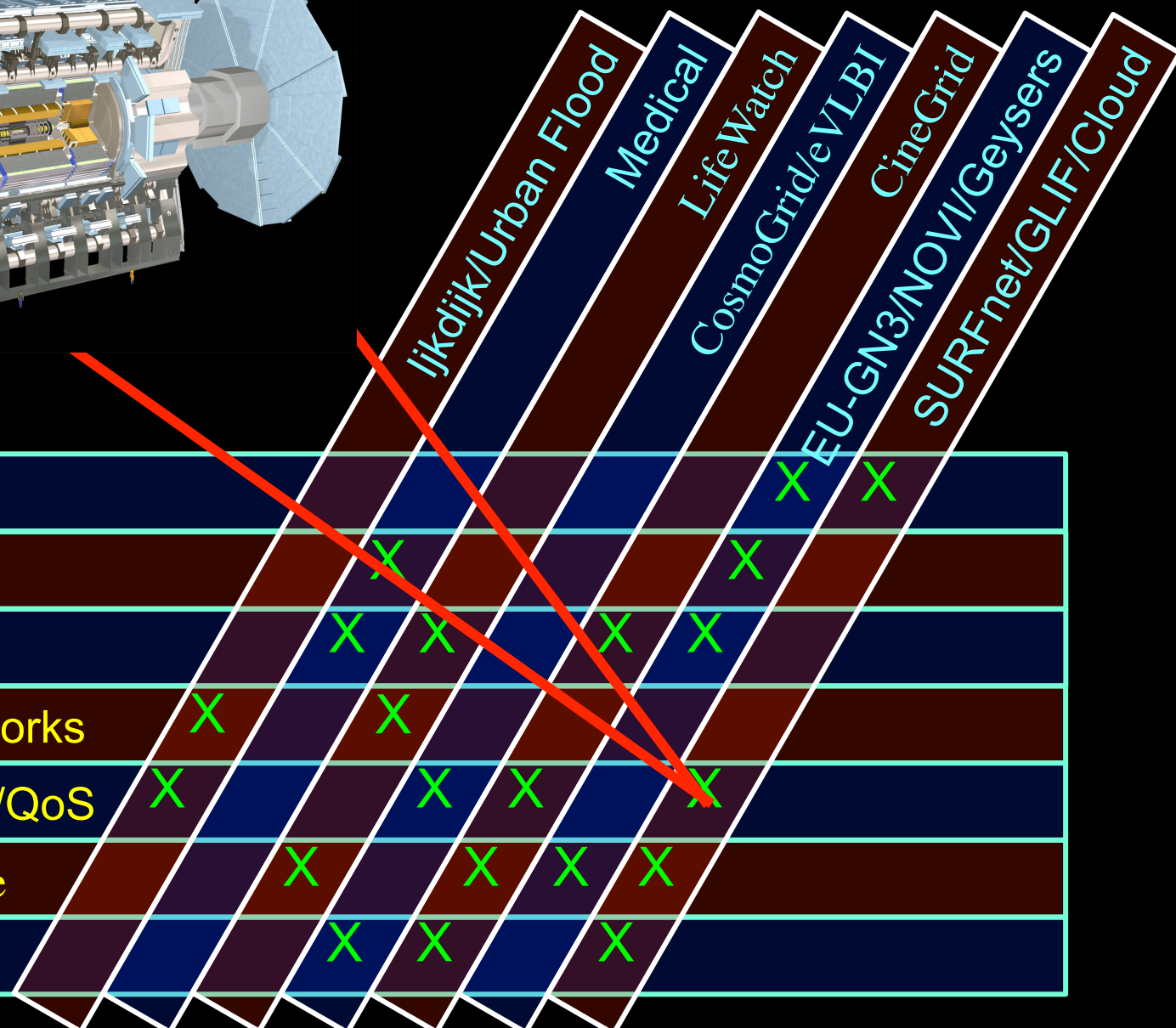
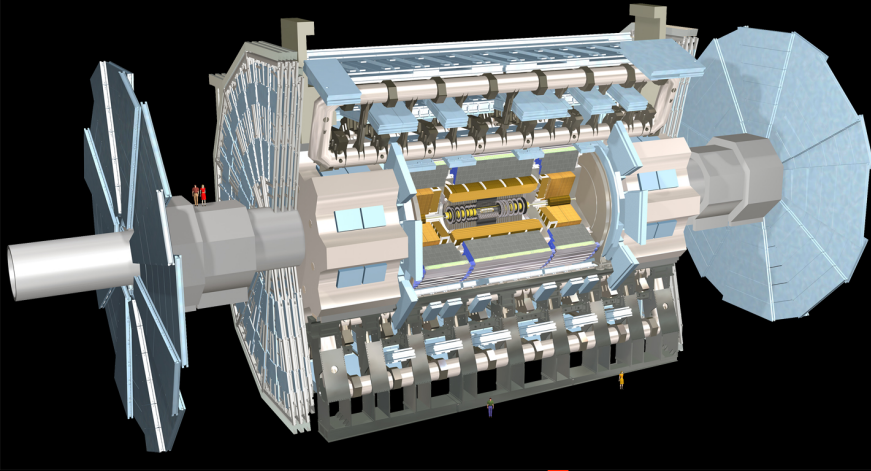
Wireless Networks



COPYRIGHT : MORTEN INGEMANN

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.

SNE @ UvA



Green-IT

Privacy/Trust

Authorization/policy

Programmable networks

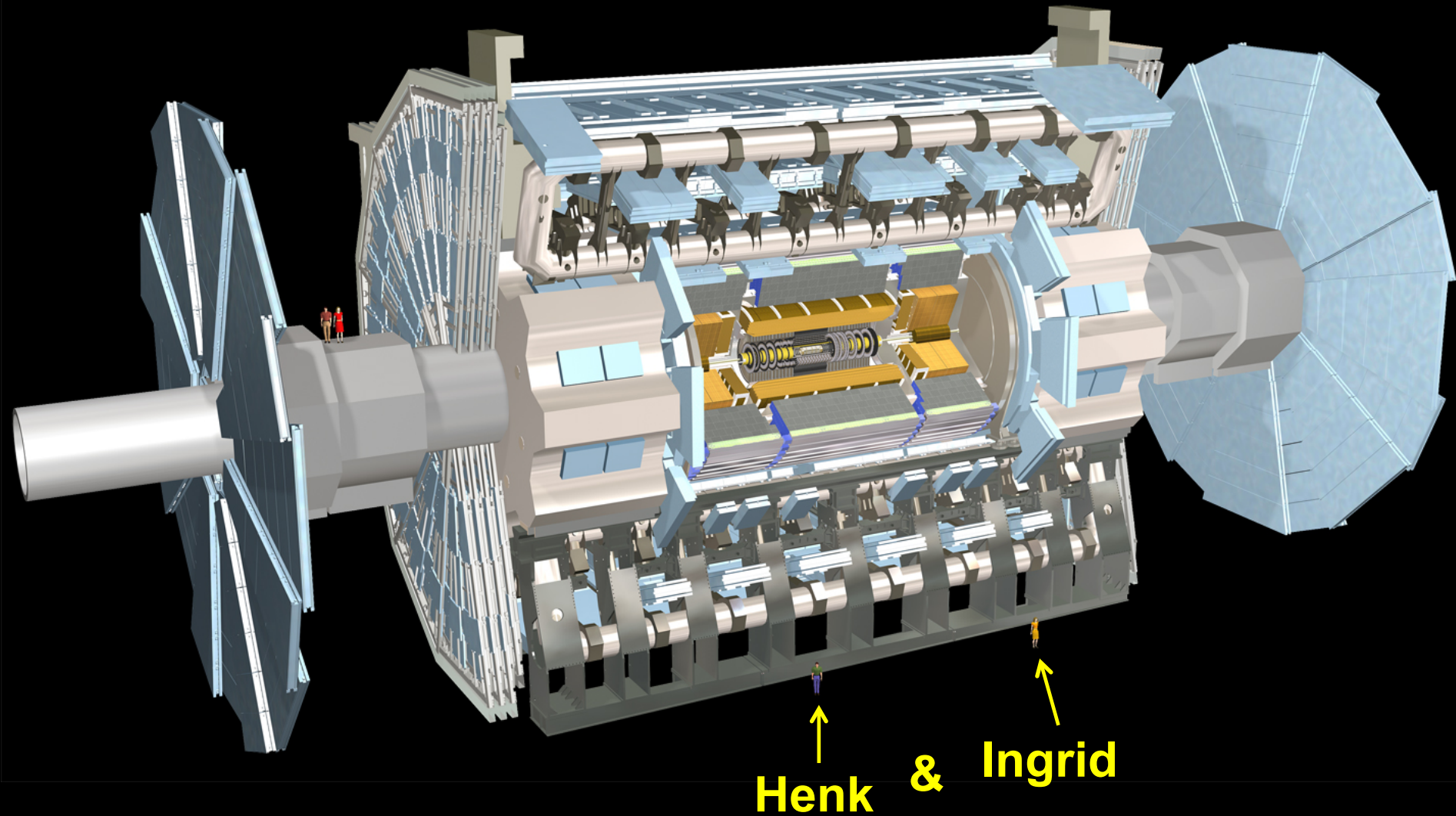
40-100Gig/TCP/WF/QoS

Topology/Architecture

Optical Photonic

					X	X
		X	X		X	
	X		X	X	X	
		X	X	X	X	
	X		X	X	X	
		X	X			

ATLAS detector @ CERN Geneve



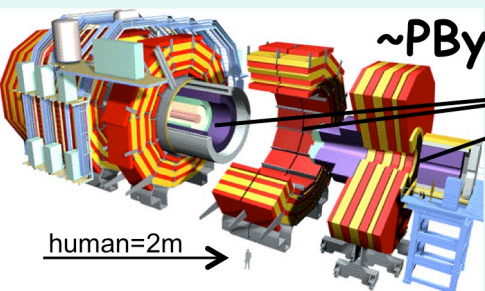
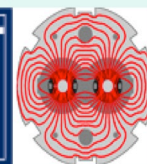
ATLAS detector @ CERN Geneve





LHC Data Grid Hierarchy

CMS as example, Atlas is similar



human=2m →

CMS detector: 15m X 15m X 22m
12,500 tons, \$700M.

Online System

Tier 0 + 1

~100 MBytes/sec

100000 flops/byte

10 Pflops/s

event simulation

event reconstruction

Status 2002!

~2.5 Gbits/sec

Tier 1

Italian Regional Center

German Regional Center

NIKHEF Dutch Regional Center

FermiLab, USA Regional Center

~0.6-2.5 Gbps

Tier 2 Center

Tier 2

analysis

~0.6-2.5 Gbps

Tier 3

Institute ~0.25TIPS

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

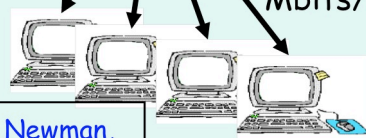
Physicists work on analysis "channels" at 135 institutes. Each institute has ~10 physicists working on one or more channels.

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

Physics data cache

100 - 1000 Mbits/sec

Tier 4



Workstations

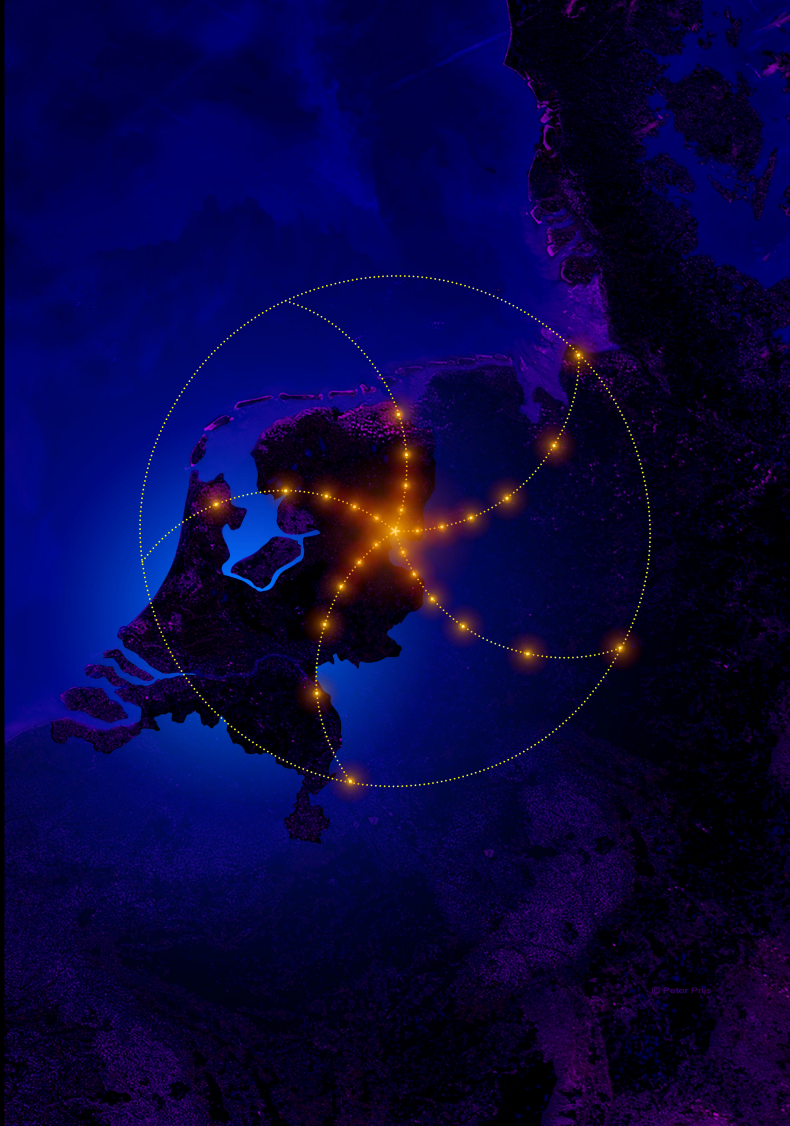
LOFAR as a Sensor Network

20 flops/byte

– LOFAR is a large distributed research infrastructure:

2 Tflops/s

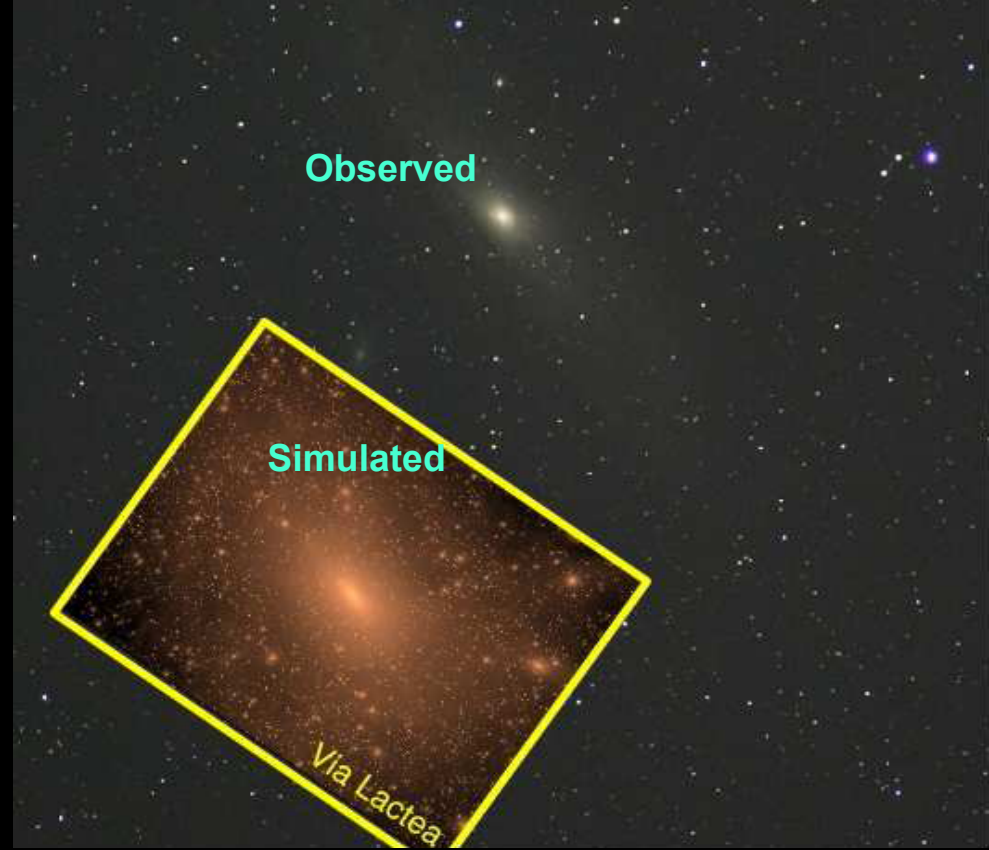
- Astronomy:
 - >100 phased array stations
 - Combined in aperture synthesis array
 - 13,000 small “LF” antennas
 - 13,000 small “HF” tiles
- Geophysics:
 - 18 vibration sensors per station
 - Infrasound detector per station
- >20 Tbit/s generated digitally
- >40 Tflop/s supercomputer
- innovative software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



CosmoGrid

Simon Portegies Zwart et al.

- Motivation:
 - previous simulations found >100 times more substructure than is observed!
- Simulate large structure formation in the Universe
- Method: Cosmological N -body code
- Computation: Intercontinental SuperComputer Grid
- Current (2013) problem:
 - 2 PByte data in Oak Ridge!



10 Gb/s dedicated network

270 ms RTT



Moving Cinegrid 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



UHDTV(4K)

3840

2160



Yesterday's Media Transport Method!

8 TByte

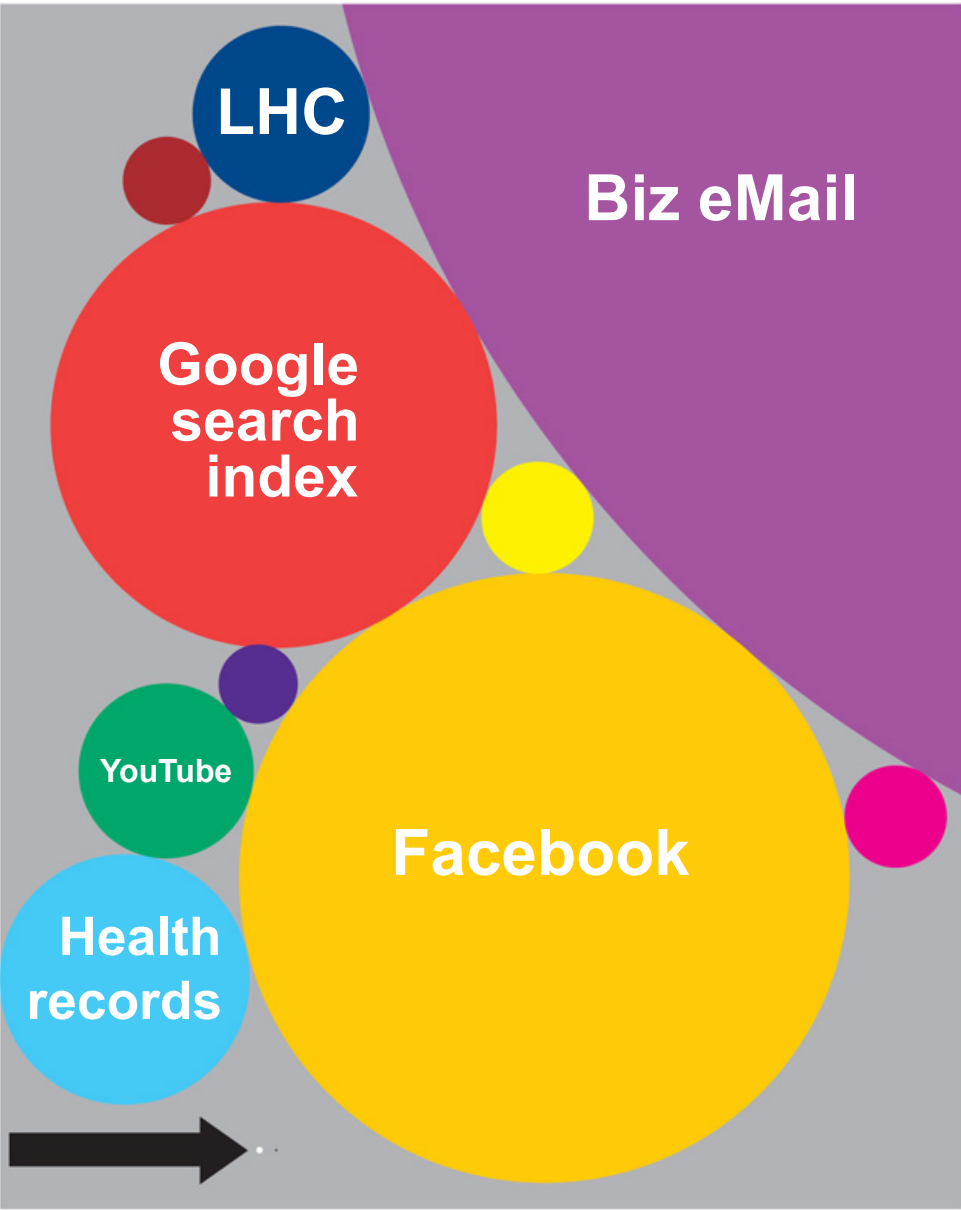


What Happens in an Internet Minute?



And Future Growth is Staggering





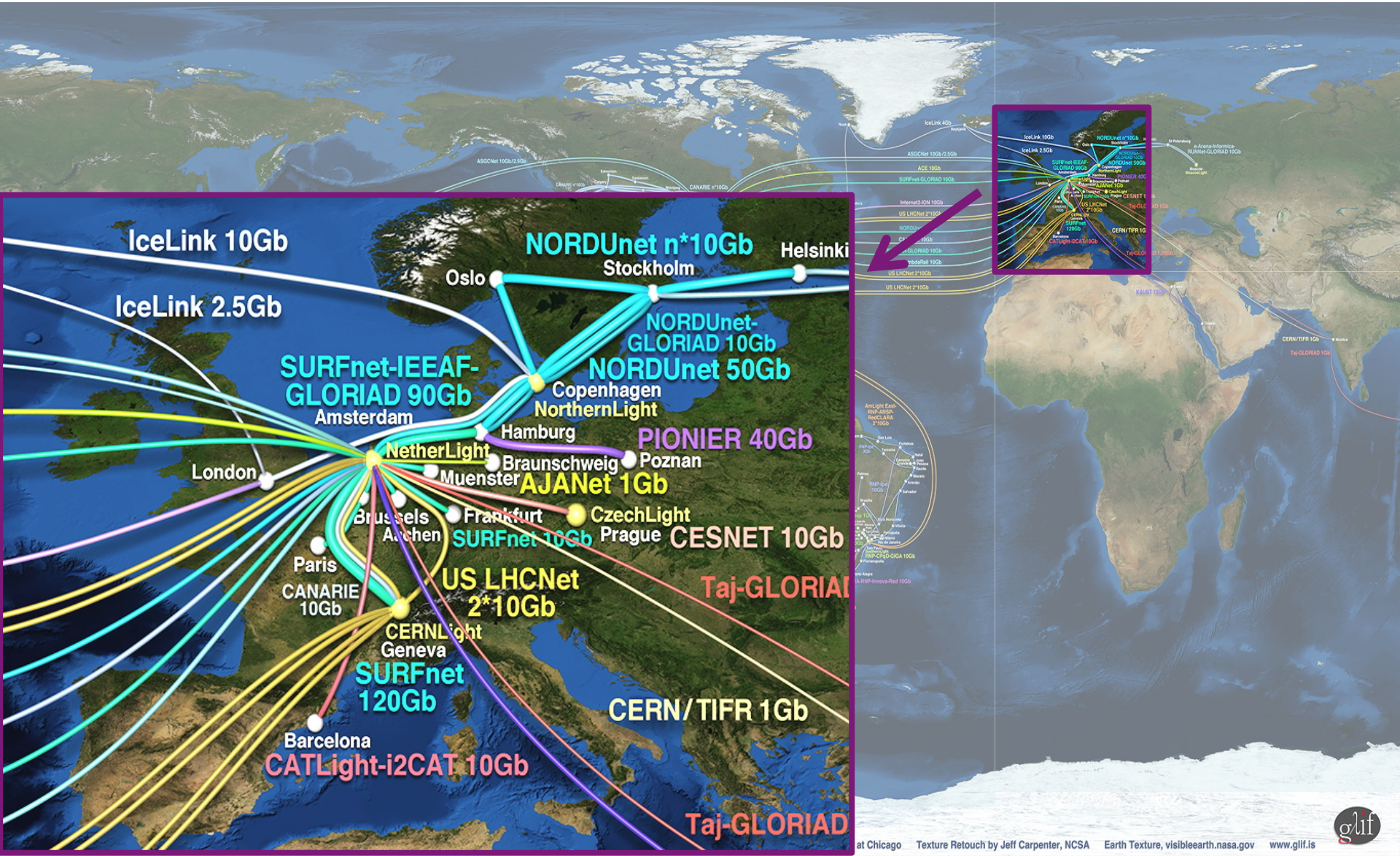
There
is
always
a
bigger
fish

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

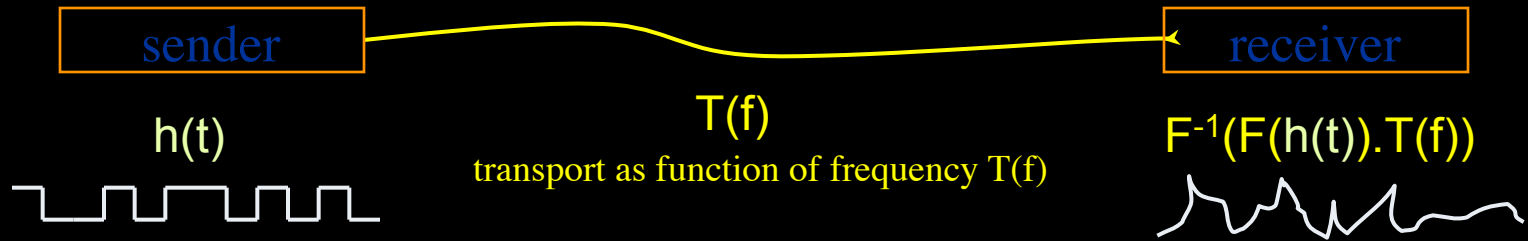
The GLIF – LightPaths around the World

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.



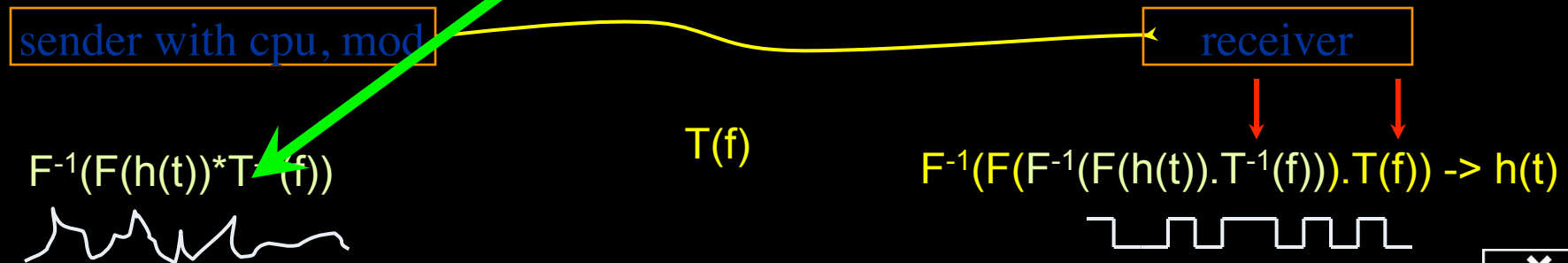
Dispersion compensating modem: eDCO from NORTEL

(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)



ExoGeni @ UvA

Installed and up June 3th 2013



TNC2013 DEMOS JUNE, 2013

DEMO	TITLE	OWNER	AFFILIATION	E-MAIL	A-SIDE	Z-SIDE	PORTS(S) MAN LAN	PORTS(S) TNC2013	DETAILS
1	Big data transfers with multipathing, OpenFlow and MPTCP	Ronald van der Pol	SURFnet	ronald.vanderpol@surfnet.nl	TNC/MECC, Maastricht, NL	Chicago, IL	Existing 100G link between internet2 and ESnet	2x40GE (Juniper)-2x10GE (OME6500)	In this demonstration we show how multipathing, OpenFlow and Multipath TCP (MPTCP) can help in large file transfers between data centres (Maastricht and Chicago). An OpenFlow application provisions multiple paths between the servers and MPTCP will be used on the servers to simultaneously send traffic across all those paths. This demo uses 2x40G on the transatlantic 100G link. ESnet provides 2x40G between MAN LAN and StarLight, ACE and USLHCnet provide additional 10GEs.
2	Visualize 100G traffic	Inder Monga	ESnet	imonga@es.net					Using an SNMP feed from the Juniper switch at TNC2013 and/or Brocade AL35 node in MAN LAN, this demo would visualize the total traffic on the link, or all demos aggregated. The network diagram will show the transatlantic topology and some of the demo topologies.
3	How many modern servers can fill a 100Gbps Transatlantic Circuit?	Inder Monga	ESnet	imonga@es.net	Chicago, Ill	TNC showfloor	1x 100GE	8x 10GE	In this demonstration, we show that with the proper tuning and tools, only 2 hosts on each continent can generate almost 100Gbps of traffic. Each server has 8 10G NICs connected to a 40G virtual circuit, and has perf3 running to generate traffic. ESnet's new "perf3" throughput measurement tool, still in beta, combines the best features from other tools such as iperf, nttop, and netperf. See: https://my.es.net/demos/tnc2013/
4	First European ExoGENI at Work	Jeroen van der Ham	UvA	vdham@uva.nl	RENCI, NC	UvA, Amsterdam, NL	1x 10GE	1x 10GE	The ExoGENI racks at RENC1 and UvA will be interconnected over a 10G pipe and be on continuously, showing GENI connectivity between Amsterdam and the rest of the GENI nodes in the USA.
5	Up and down North Atlantic @ 100G	Michael Enrico	DANTE	michael.enrico@dante.net	TNC showfloor	TNC showfloor	1x 100GE	1x 100GE	The ExoGENI racks at RENC1 and UvA will be interconnected over a 10G pipe and be on continuously, showing GENI connectivity between Amsterdam and the rest of the GENI nodes in the USA.

Connected via the new 100 Gb/s transatlantic

Alien light From idea to realisation!

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



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

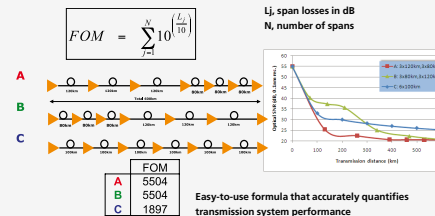
Alien wavelength challenges

- 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.

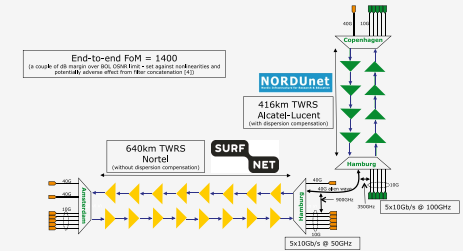
New method to present fiber link quality, FoM (Figure of Merit)

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.

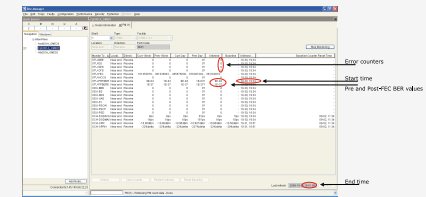


Transmission system setup

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.



REFERENCES
ACKNOWLEDGEMENTS

[1] "OPERATIONAL SOLUTIONS FOR AN OPEN DWDM LAYER", O. GERSTEL ET AL. OFC2009 | [2] "AT&T OPTICAL TRANSPORT SERVICES", BARBARA E. SMITH, OFC'09
 [3] "OPEX SAVINGS OF ALL-OPTICAL CORE NETWORKS", ANDREW LORD AND CARL ENGINEER, ECCO2009 | [4] NORTEL/SURFNET INTERNAL COMMUNICATION
 WE ARE GRATEFUL TO NORDUNET FOR PROVIDING US WITH BANDWIDTH ON THEIR DWDM LINK FOR THIS EXPERIMENT AND ALSO FOR THEIR SUPPORT AND ASSISTANCE DURING THE EXPERIMENTS. WE ALSO ACKNOWLEDGE TELINDUS AND NORTEL FOR THEIR INTEGRATION WORK AND SIMULATION SUPPORT

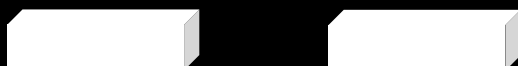
ClearStream @ TNC2011

Setup
codename:
FlightCees



UvA

iPerf 17 3.2 GHz Q-core Amd Ph II 3.6 GHz HexC



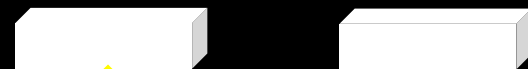
Mellanox

40G E



Copenhagen

iPerf 2* dual 2.8 GHz Q-core



Mellanox



CERN

CIENA DWDM

17 ms RTT

Hamburg

Alcatel DWDM

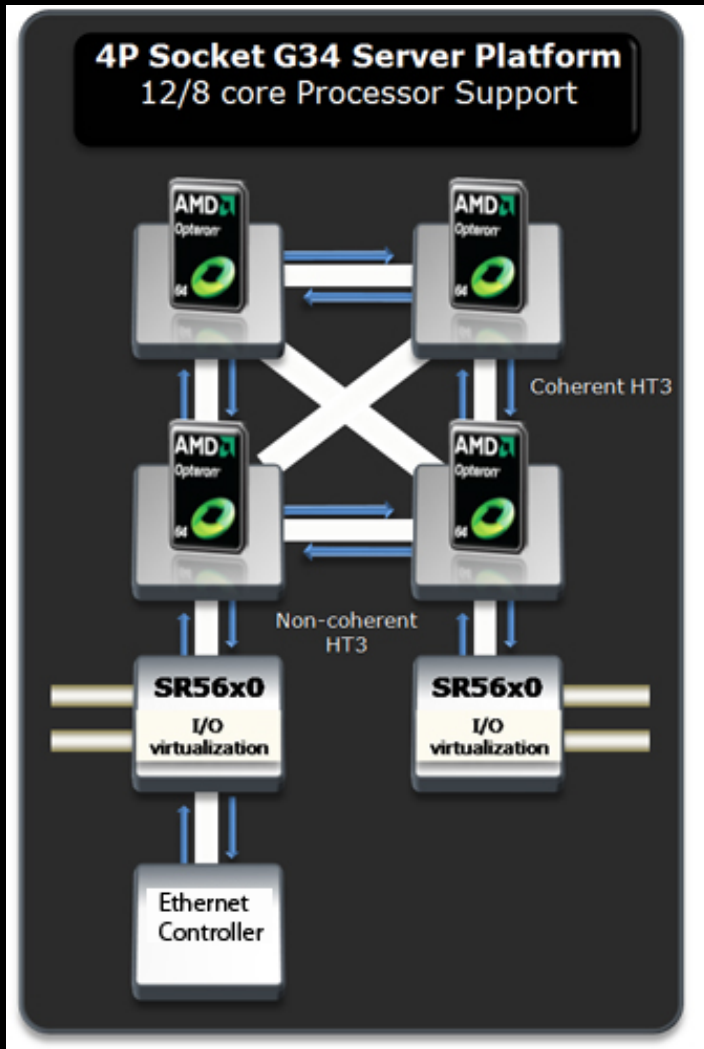
27 ms RTT

Amsterdam – Geneva (CERN) – Copenhagen – 4400 km (2700 km alien light)

Results (rtt = 17 ms)

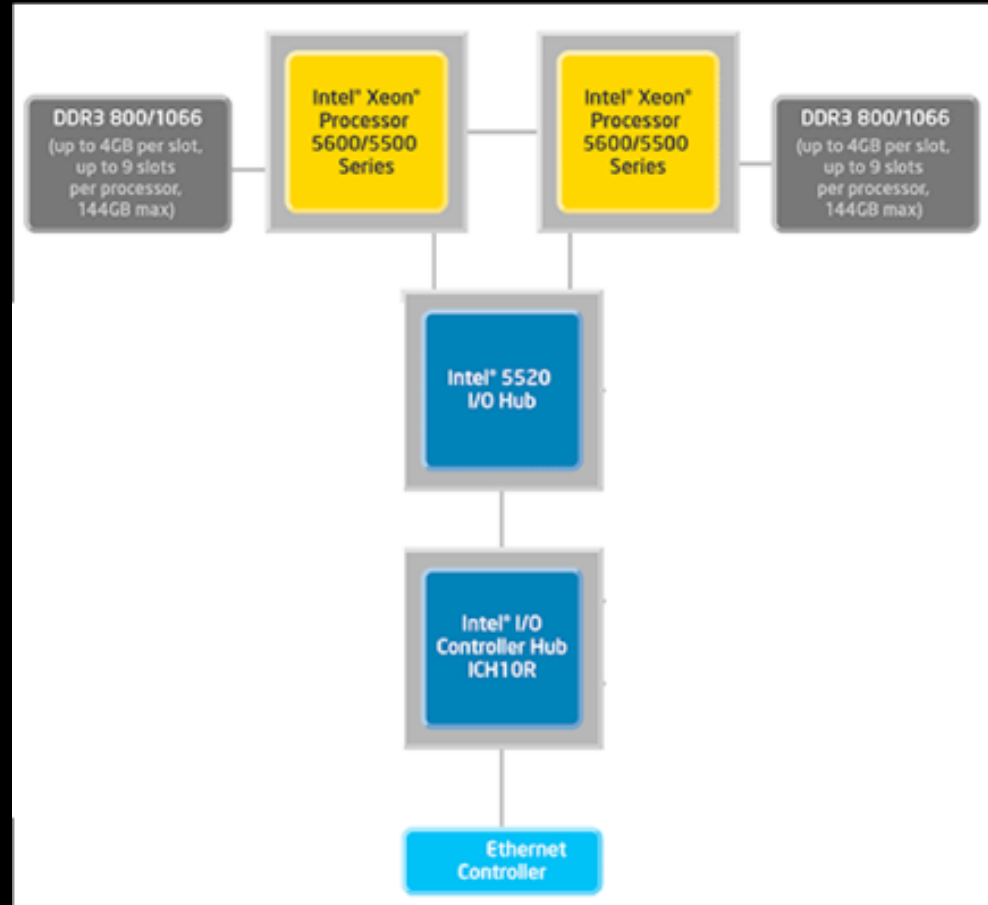
- ❑ Single flow iPerf 1 core -> 21 Gbps
- ❑ Single flow iPerf 1 core <> -> 15+15 Gbps
- ❑ Multi flow iPerf 2 cores -> 25 Gbps
- ❑ Multi flow iPerf 2 cores <> -> 23+23 Gbps
- ❑ DiViNe <> -> 11 Gbps
- ❑ Multi flow iPerf + DiVine -> 35 Gbps
- ❑ Multi flow iPerf + DiVine <> -> 35 + 35 Gbps

Server Architecture



DELL R815

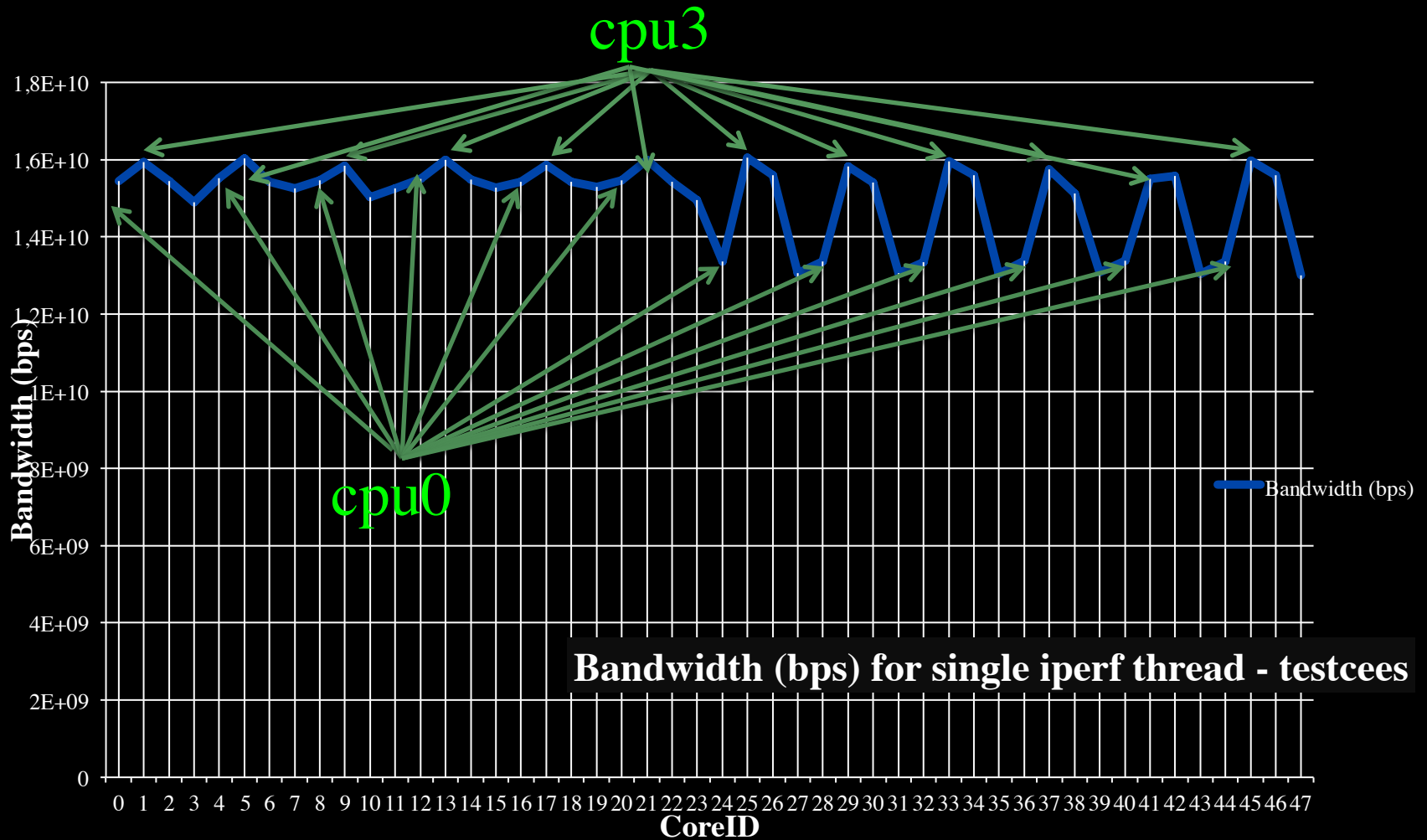
4 x AMD Opteron 6100



Supermicro X8DTT-HIBQF

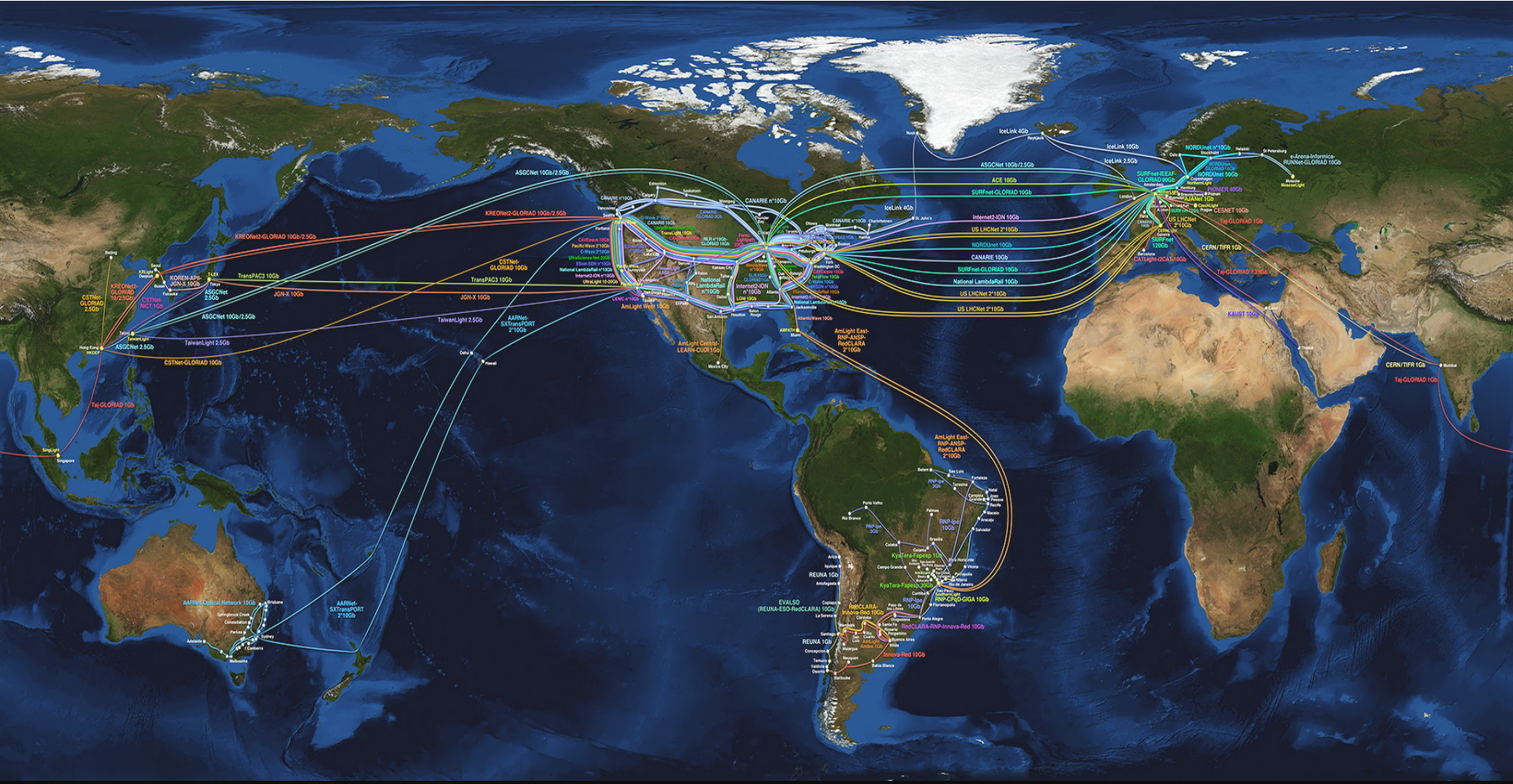
2 x Intel Xeon

CPU Topology benchmark



We used numactl to bind iperf to cores

The GLIF – LightPaths around the World



We investigate:
complex networks!



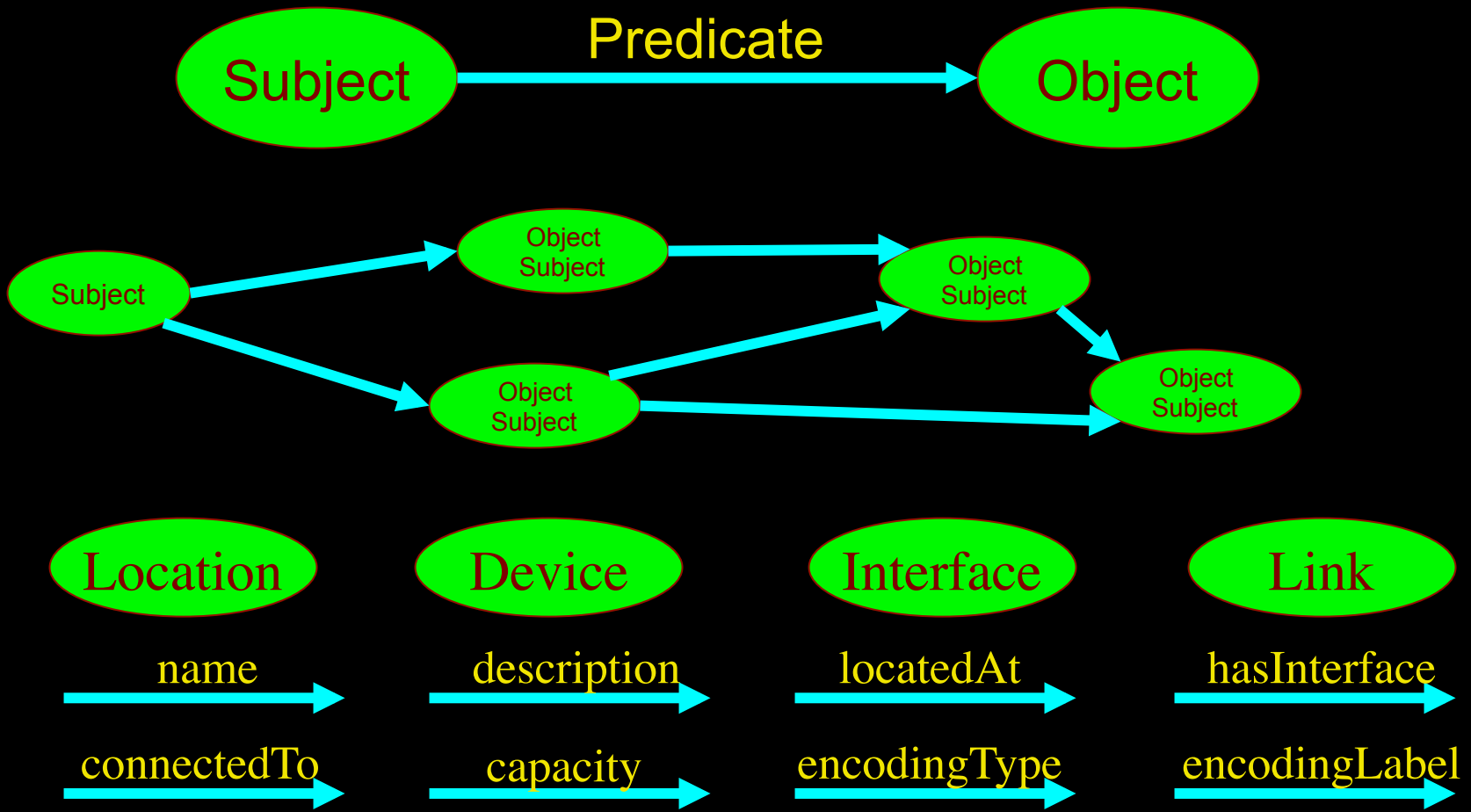
for



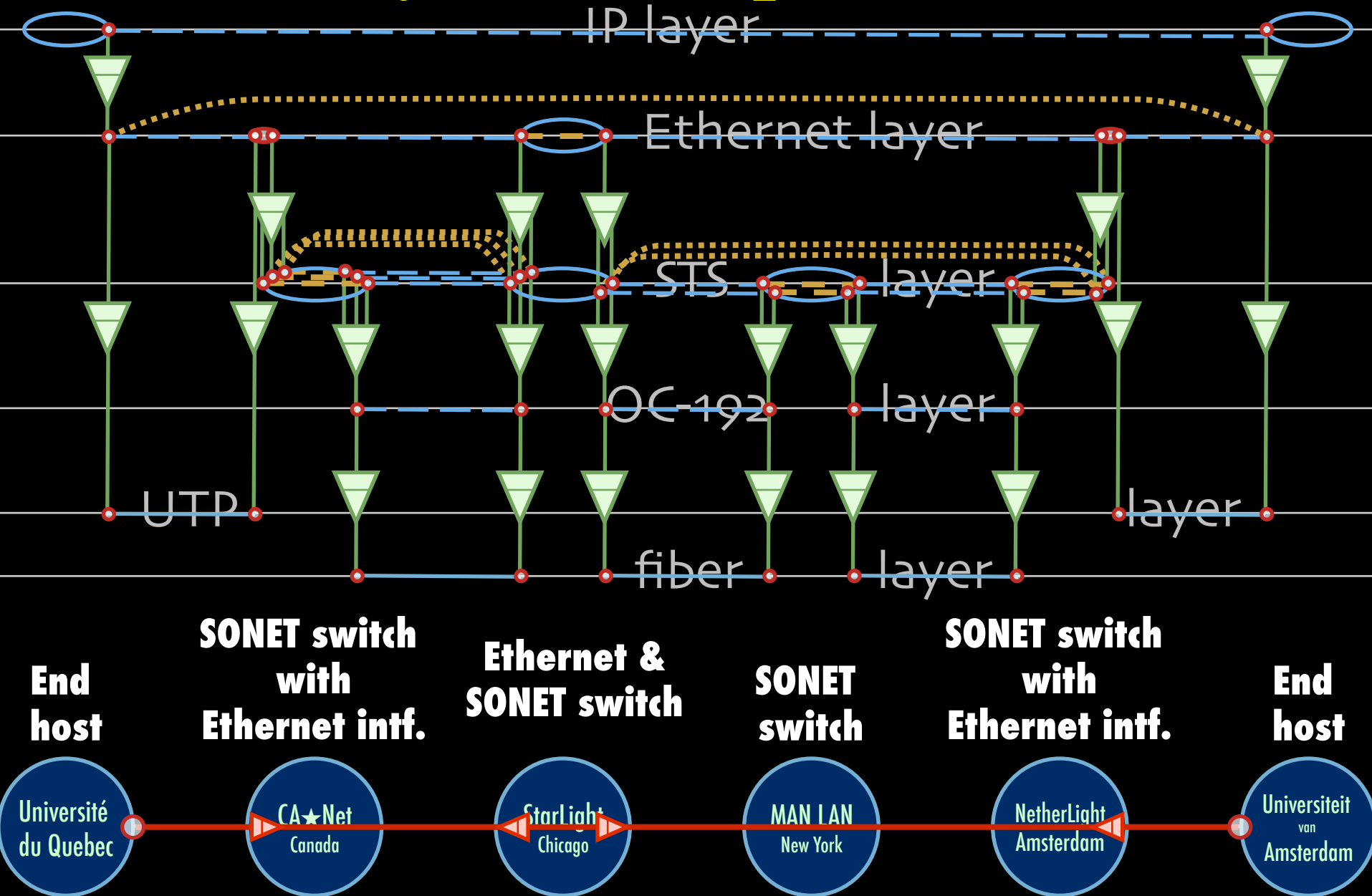
LinkedIn for Infrastructure



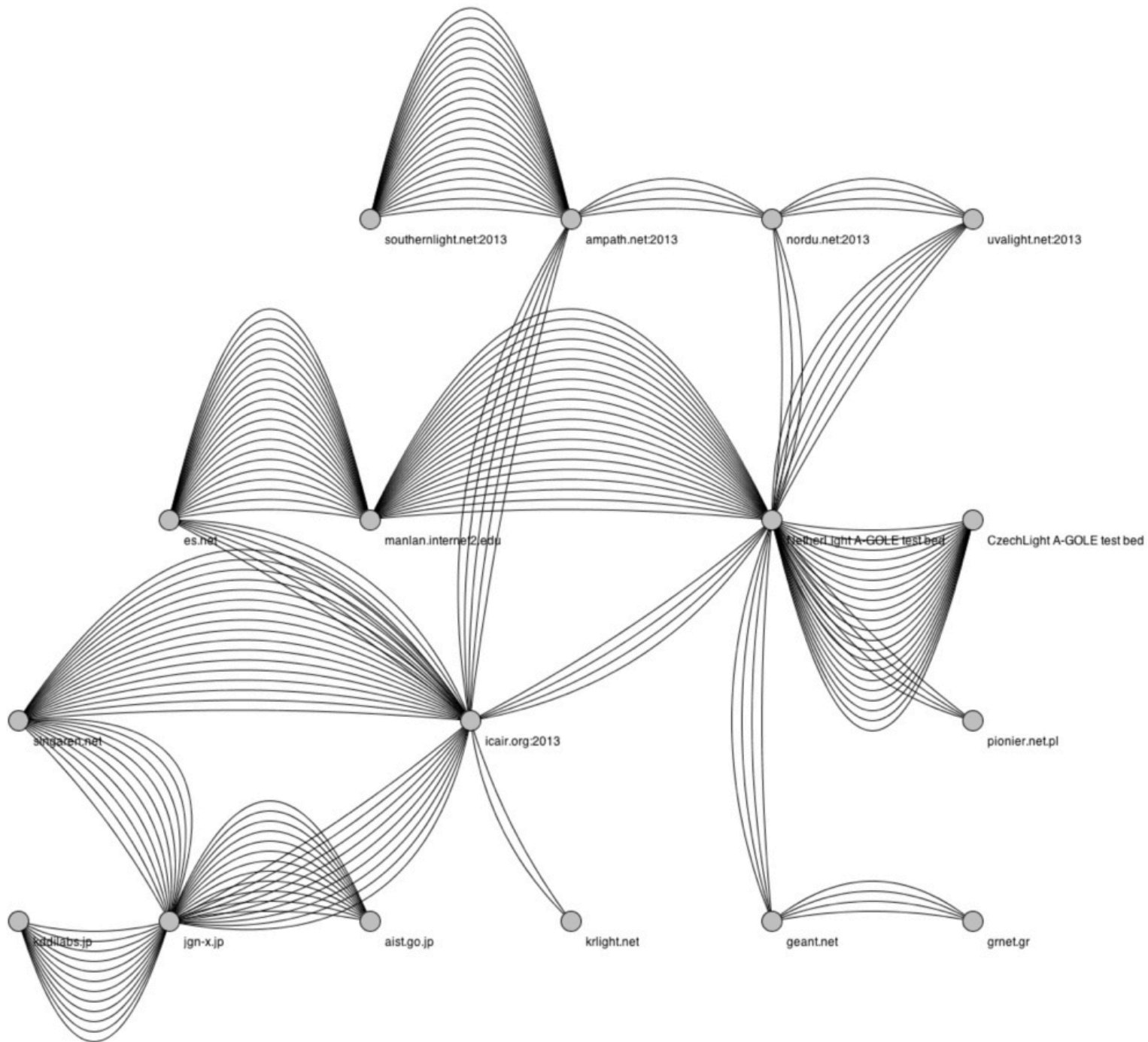
- From semantic Web / Resource Description Framework.
- The RDF uses XML as an interchange syntax.
- Data is described by triplets (Friend of a Friend):



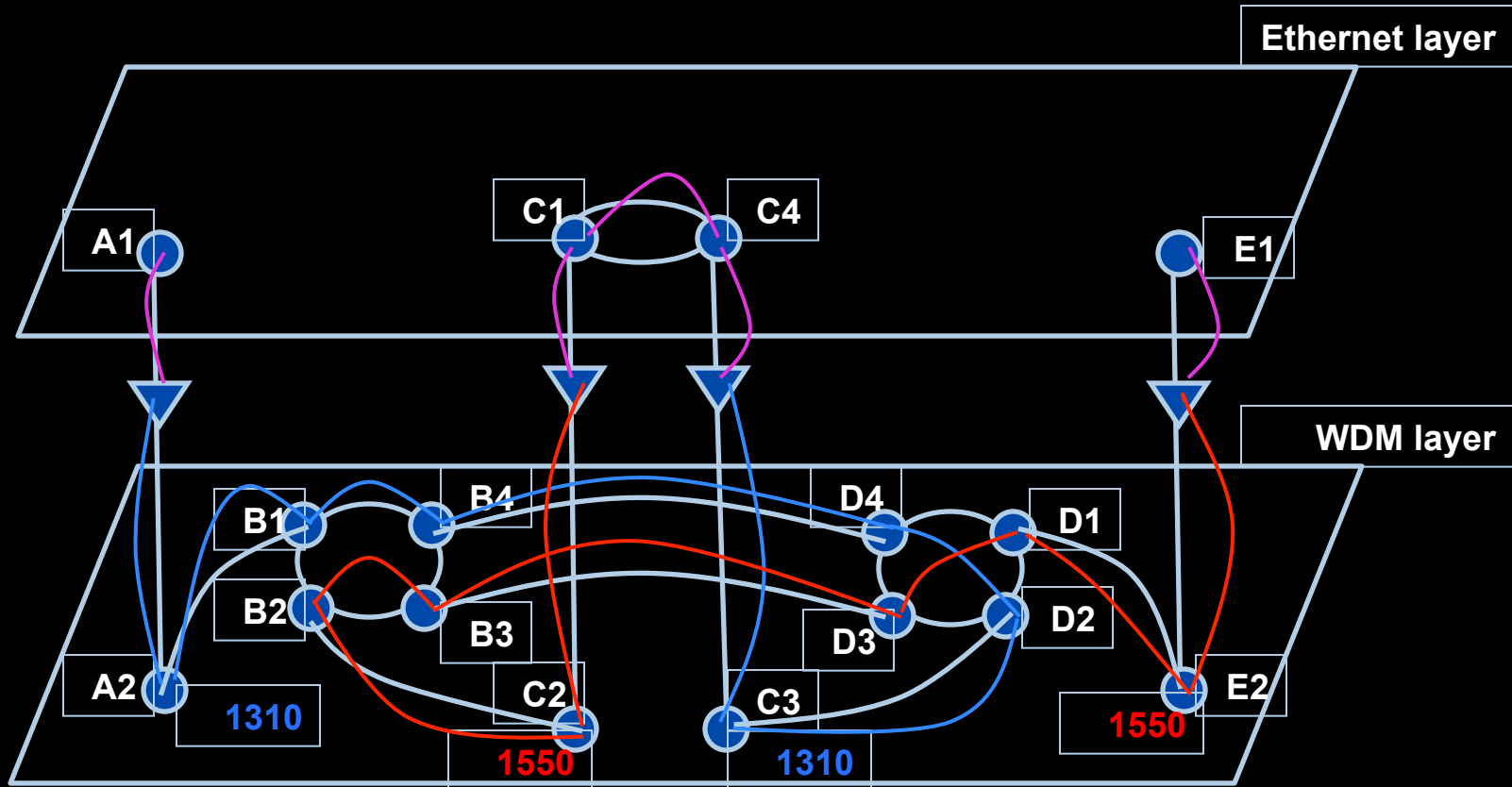
Multi-layer descriptions in NDL



GLIF 2013



Multi-layer Network PathFinding



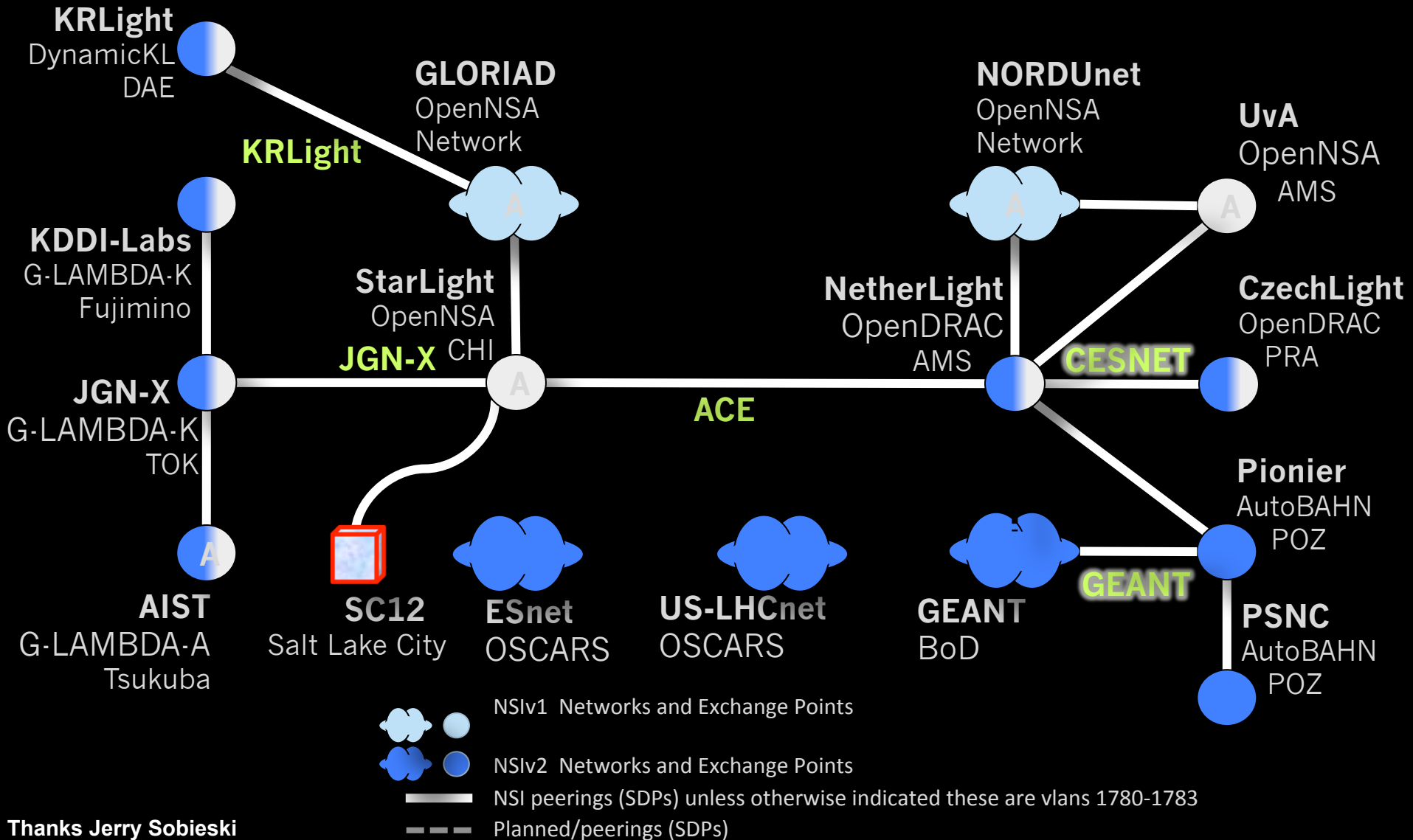
Path between interfaces A1 and E1:
A1-A2-B1-B4-D4-D2-C3-C4-C1-C2-B2-B3-D3-D1-E2-
E1

Scaling: Combinatorial problem

Automated GOLE + NSI

Joint NSI v1+v2 Beta Test Fabric Nov 2012

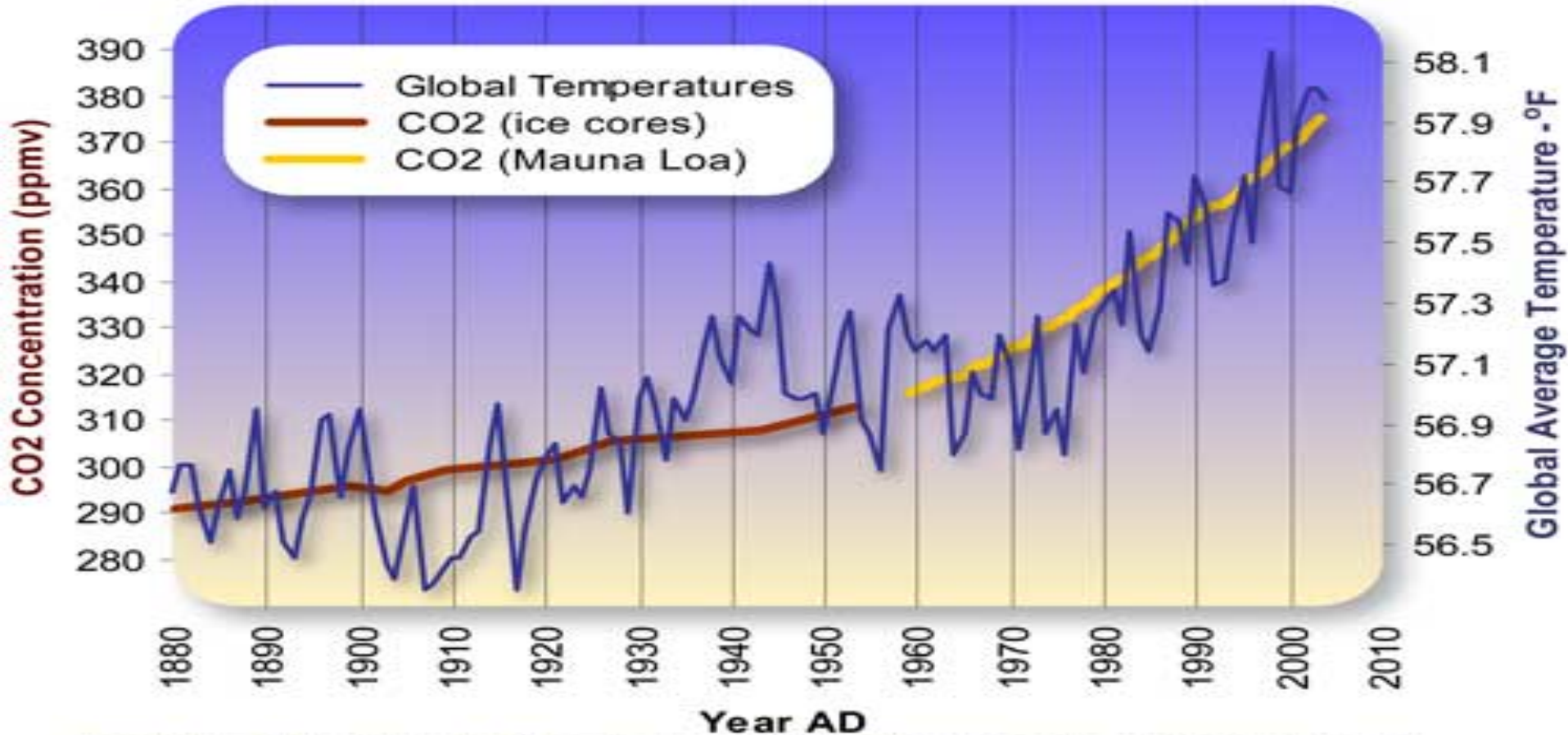
Ethernet Transport Service



Thanks Jerry Sobieski

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



Need for GreenIT

Global Average Temperature and
Positive proof of global warming.



Data Source CO₂ (Mauna Loa): <http://cdiac.esd.ornl.gov/ftp/trends/co2/maunaloa.co2>

Graphic Design: Michael Ernst, The Woods Hole Research Center



ECO-Scheduling



Bits-Nets-Energy

Bits to Energy or Energy to Bits

a calculator for a road to cleaner computing

Choose a service scenario

PUE of source and destination data center
 Src: Dest:

Transport network between source and destination data center

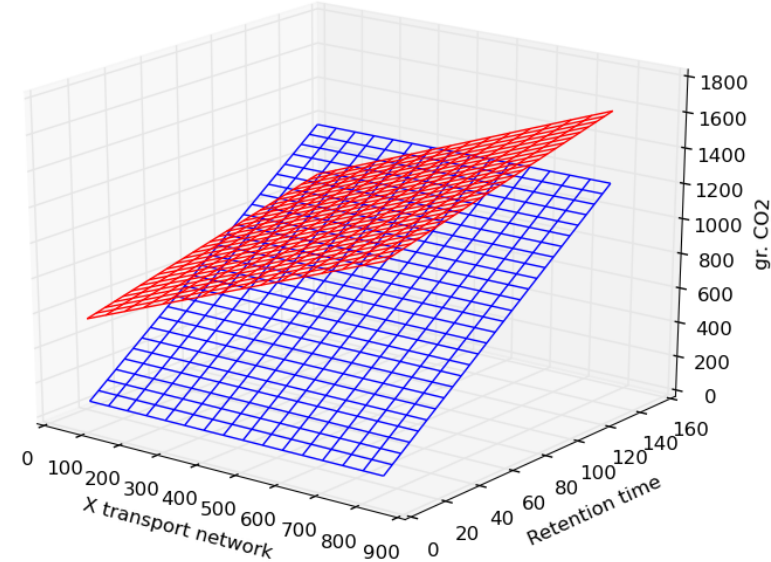
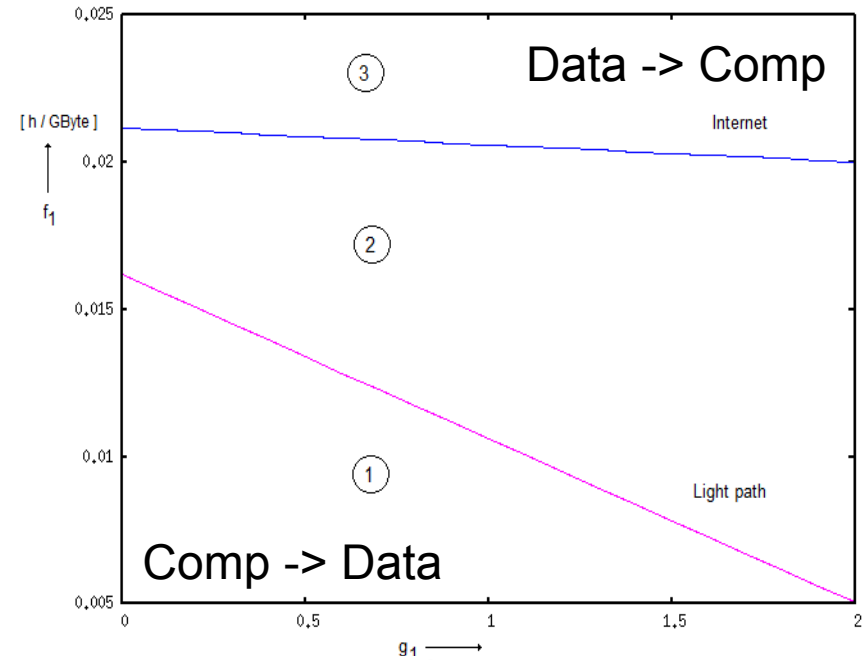
Energy production X [gr CO₂/kWh]

source datacenter X: location energy production:

dest. datacenter X: location energy production:

transport network X:

Calculate cost in gr CO₂



Storage to energy:

- When should you move hot or cold data to a green remote data center for storage?
- Given different network paths what are the decision boundaries as function of the task complexity.

Mission

Can we create smart and safe data processing infrastructures that can be tailored to diverse application needs?

- *Capacity*
 - *Bandwidth on demand, QoS, architectures, analytics, performance*
- *Capability*
 - *Integration, virtualization, complexity, semantics, workflows*
- *Security*
 - *Anonymity, integrity of data in distributed data processing*
- *Sustainability*
 - *Greening infrastructure, awareness*
- *Resilience*
 - *Systems under attack, failures, disasters*

SMART

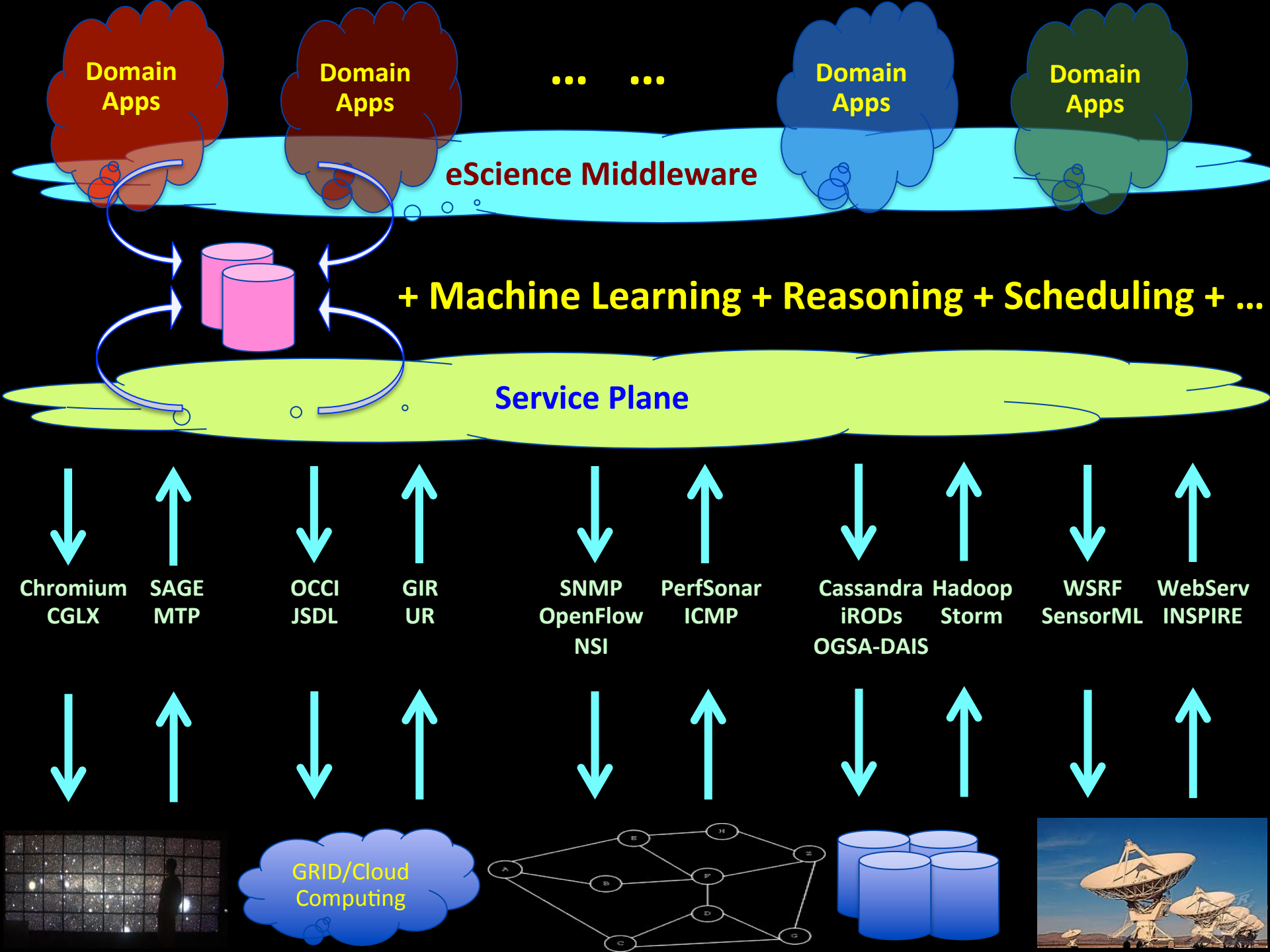


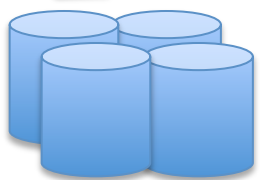
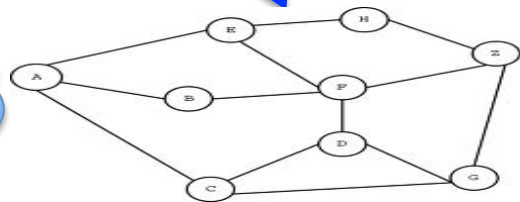
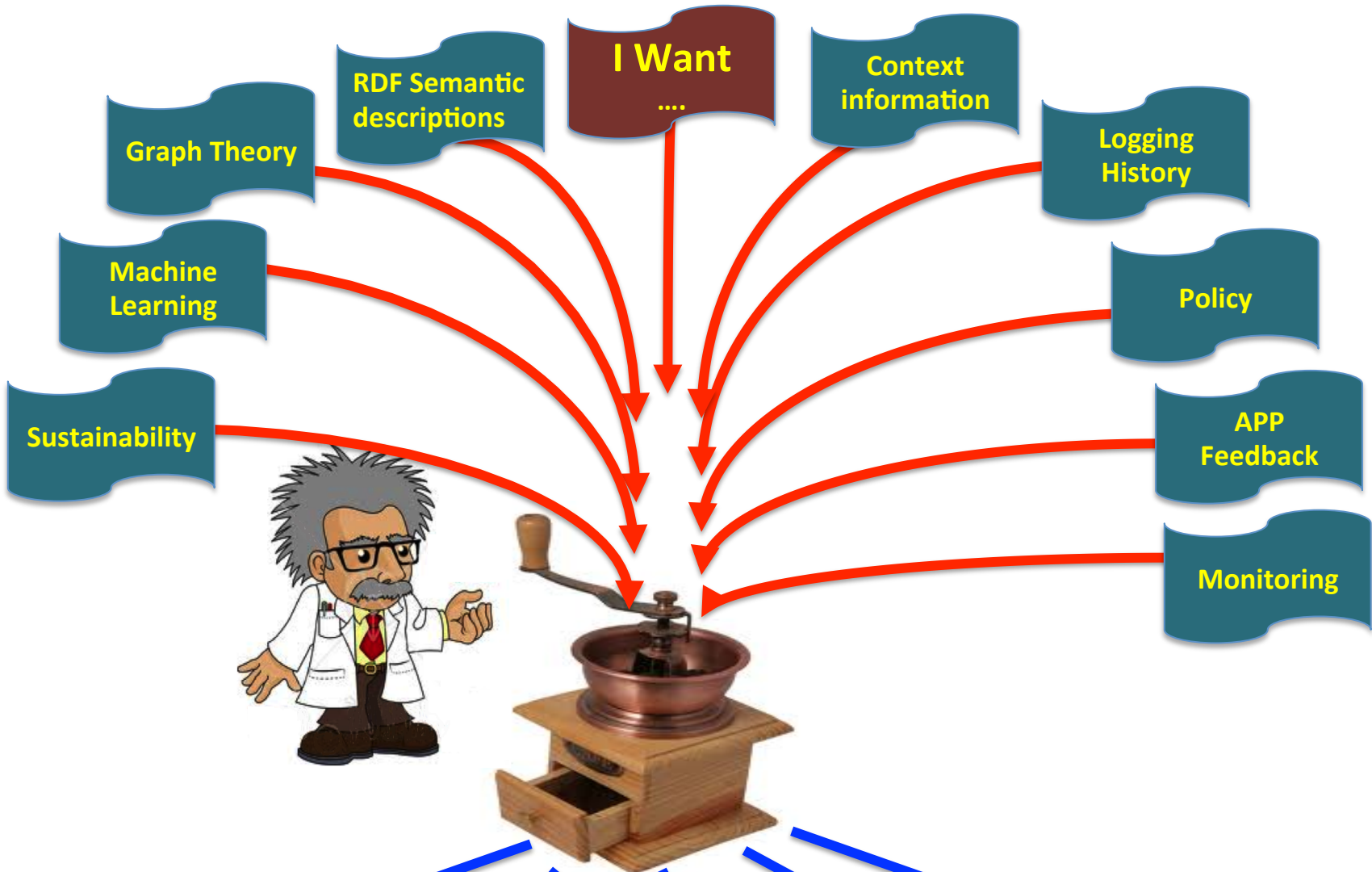
I want to



“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!

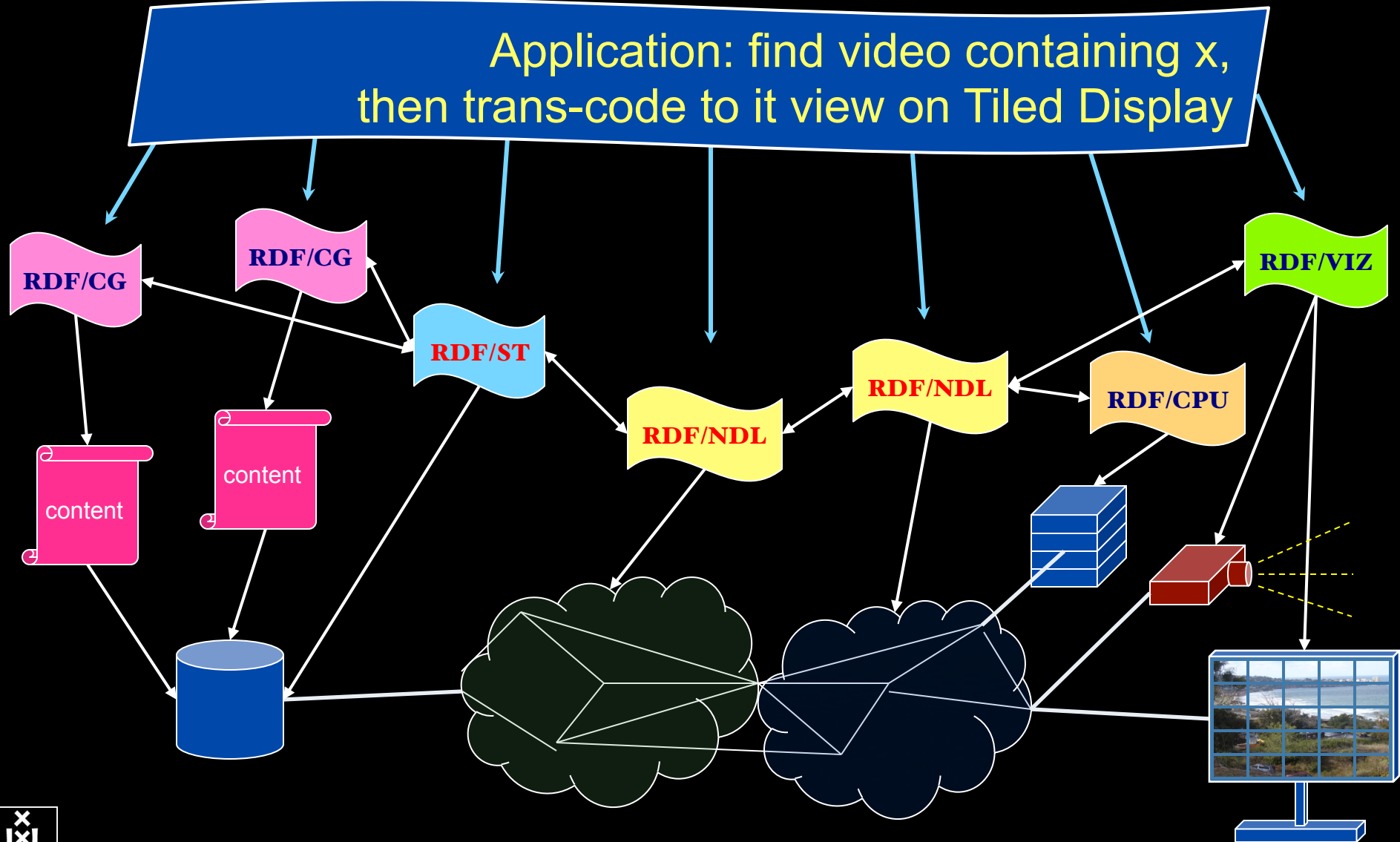




RDF describing Infrastructure

“I want”

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



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, Data³, App

to:

DDOS attacks destroying Banks and Bitcoins.

Conclusion:

Need for Safe, Smart, Resilient Sustainable Infrastructure.