# Challenges for enabling eScience over Optical Networks

## Cees de Laat

**SURFnet**

**EU**

**BSIK**

**NWO**

**University of Amsterdam**
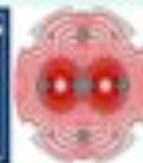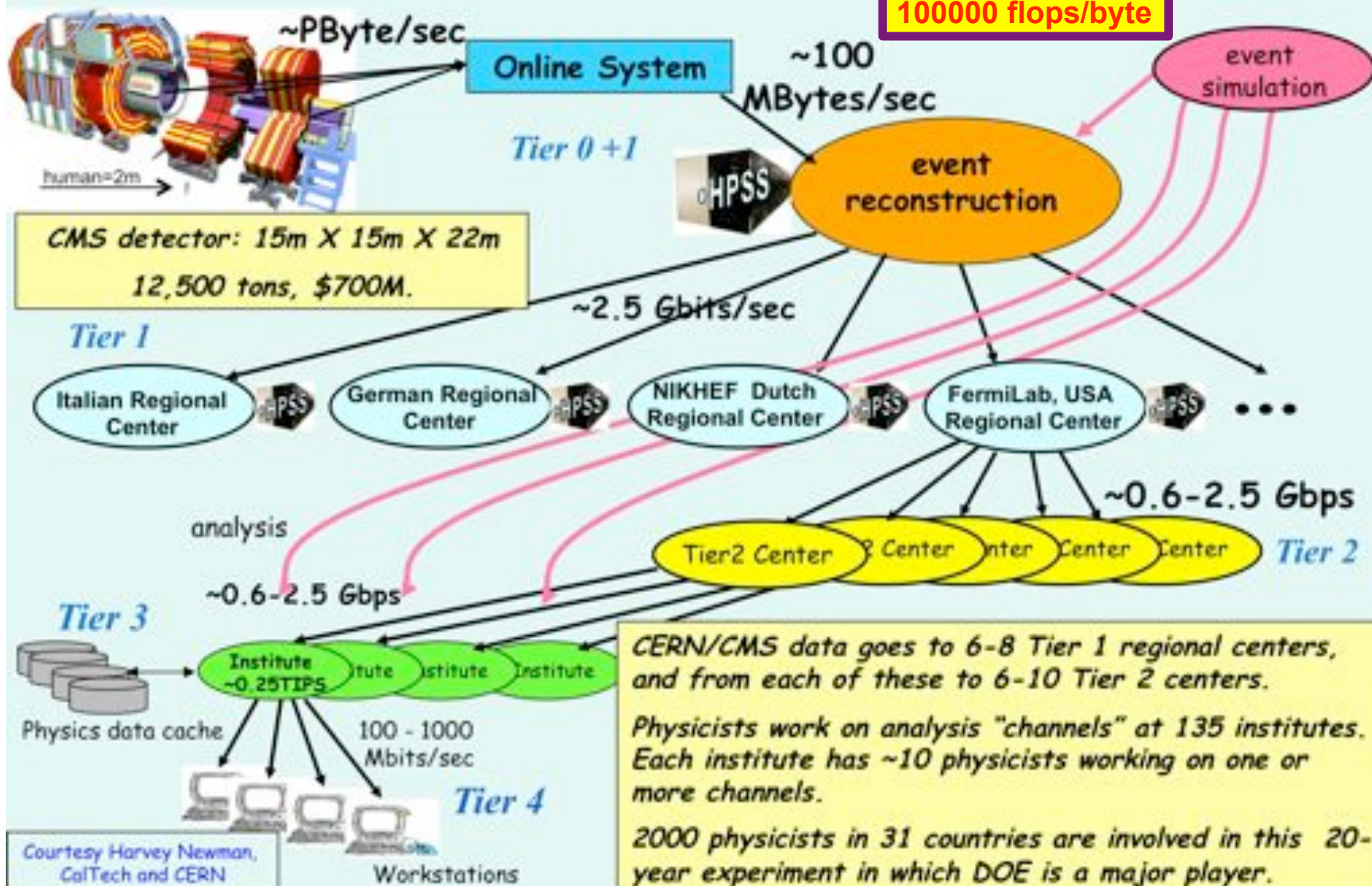
TNO
NCF

# LHC Data Grid Hierarchy
## CMS as example, Atlas is similar
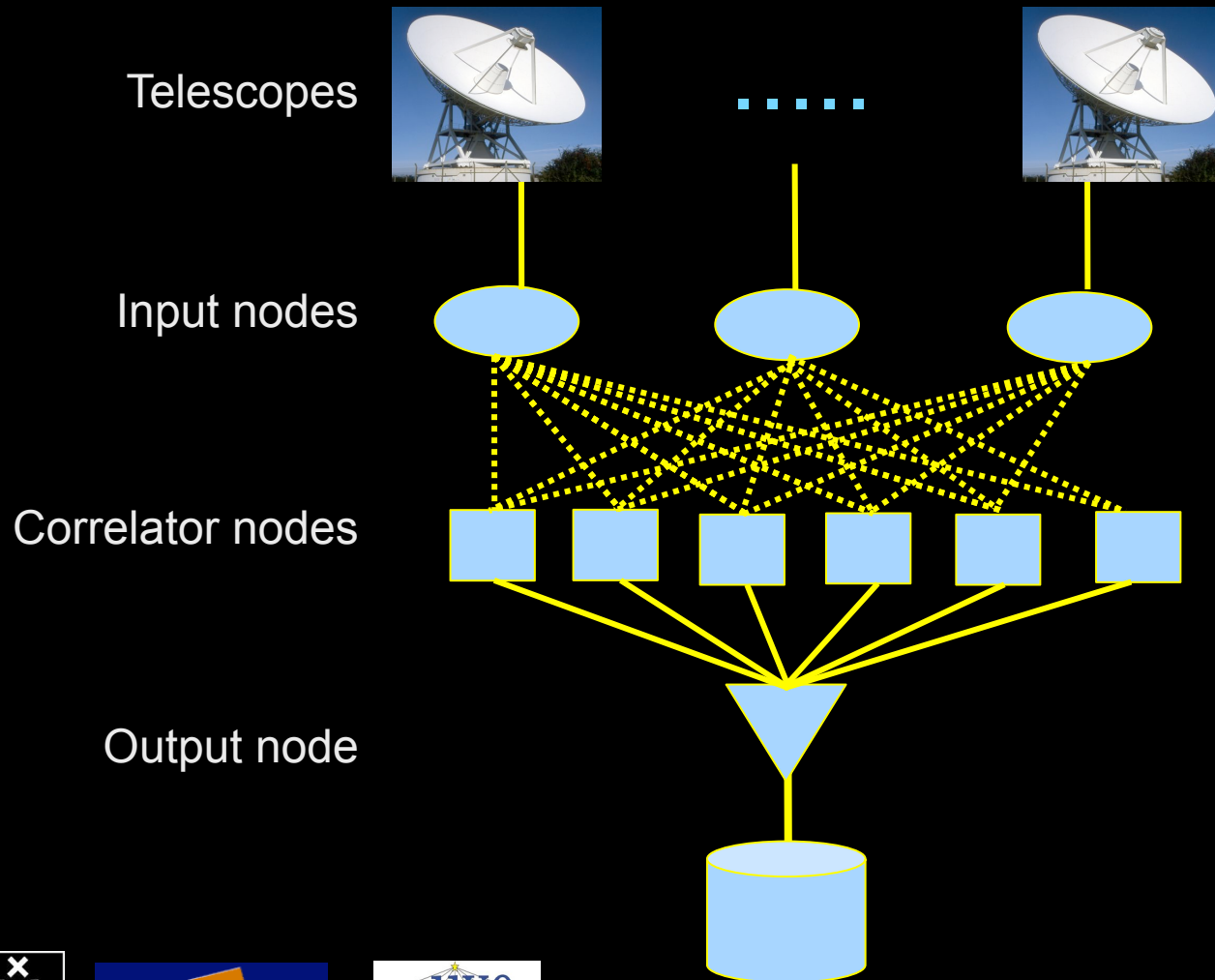
**100000 flops/byte**

~PByte/sec

Online System

~100 MBytes/sec

event simulation

Tier 0 +1

HPSS

event reconstruction

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

human=2m

~2.5 Gbits/sec

Tier 1

Italian Regional Center

PSS

German Regional Center

PSS

NIKHEF Dutch Regional Center

PSS

FermiLab, USA Regional Center

PSS

...

~0.6-2.5 Gbps

analysis

~0.6-2.5 Gbps

Tier2 Center  Center  nter  Center  Center    Tier 2

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.
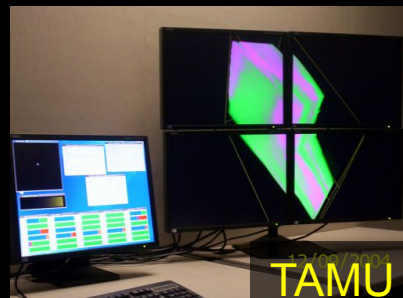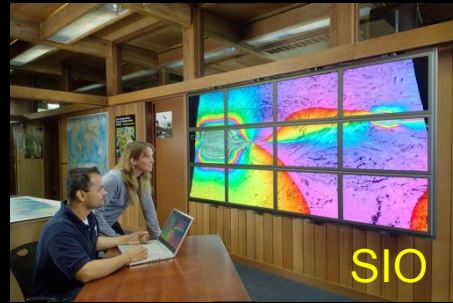
# LOFAR as a Sensor Network

**20 flops/byte**



– LOFAR is a large distributed research infrastructure:

- 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

# US and International OptIPortal Sites



SIO

NCMIR

USGS EDC

NCSA & TRECC

SARA

KISTI

AIST

RINCON & Nortel

TAMU

UCI

UIC

CALIT2

**Real time, multiple 10 Gb/s**

# The "Dead Cat" demo

**1 Mflops/byte**



SC2004,
Pittsburgh,
Nov. 6 to 12, 2004
iGrid2005,
San Diego,
sept. 2005

Many thanks to:
AMC
SARA
GigaPort
UvA/AIR
Silicon Graphics,
Inc.
Zoölogisch Museum

**IJKDIJK**

300000 * 60 kb/s * 2 sensors (microphones) to cover all Dutch dikes

# Sensor grid: instrument the dikes

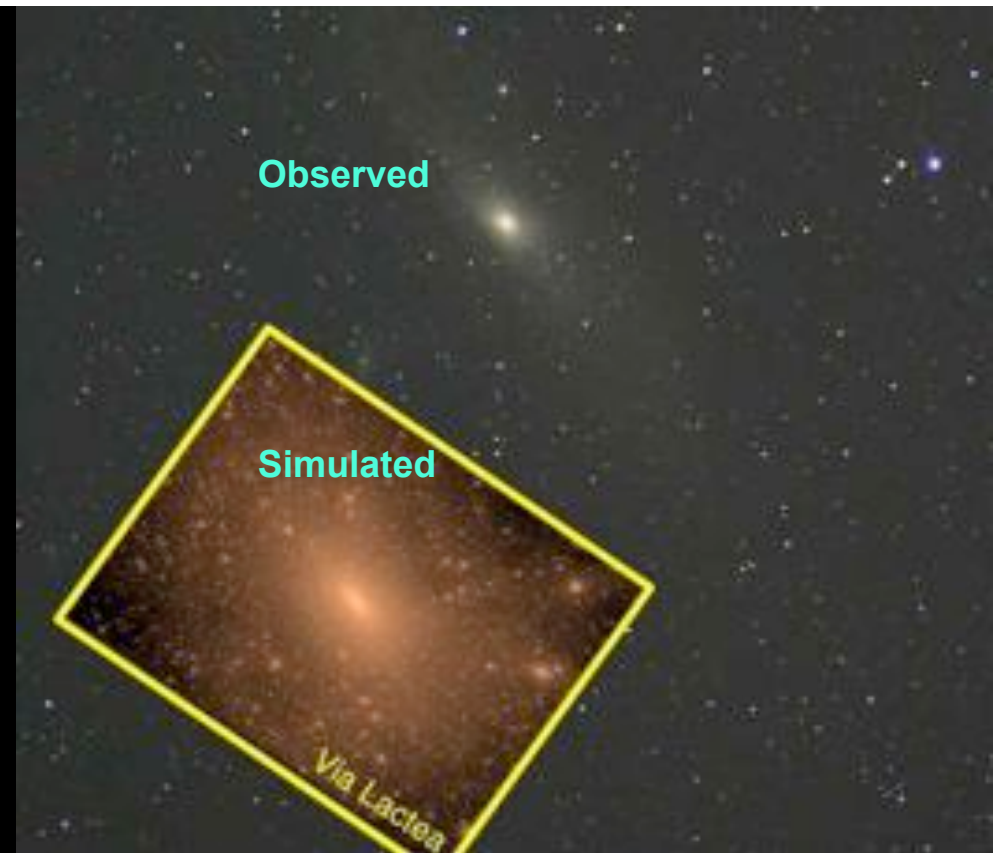## First controlled breach occurred on sept 27th '08:



**Many small flows -> 36 Gb/s**

**Urban_flow !**

# CosmoGrid



Observed

Simulated

Via Lactea

- Motivation:

  **previous simulations found >100 times more substructure than is observed!**

- Simulate large structure formation in the Universe

  – Dark Energy (cosmological constant)

  – Dark Matter (particles)

- Method: Cosmological *N*-body code

- Computation: Intercontinental SuperComputer Grid

# The hardware setup

- 2 supercomputers :

  - 1 in Amsterdam   (60Tflops Power6 @ SARA)

  - 1 in Tokyo (30Tflops Cray XD0-4 @ CFCA)

- Both computers are connected via an intercontinental optical 10 Gbit/s network
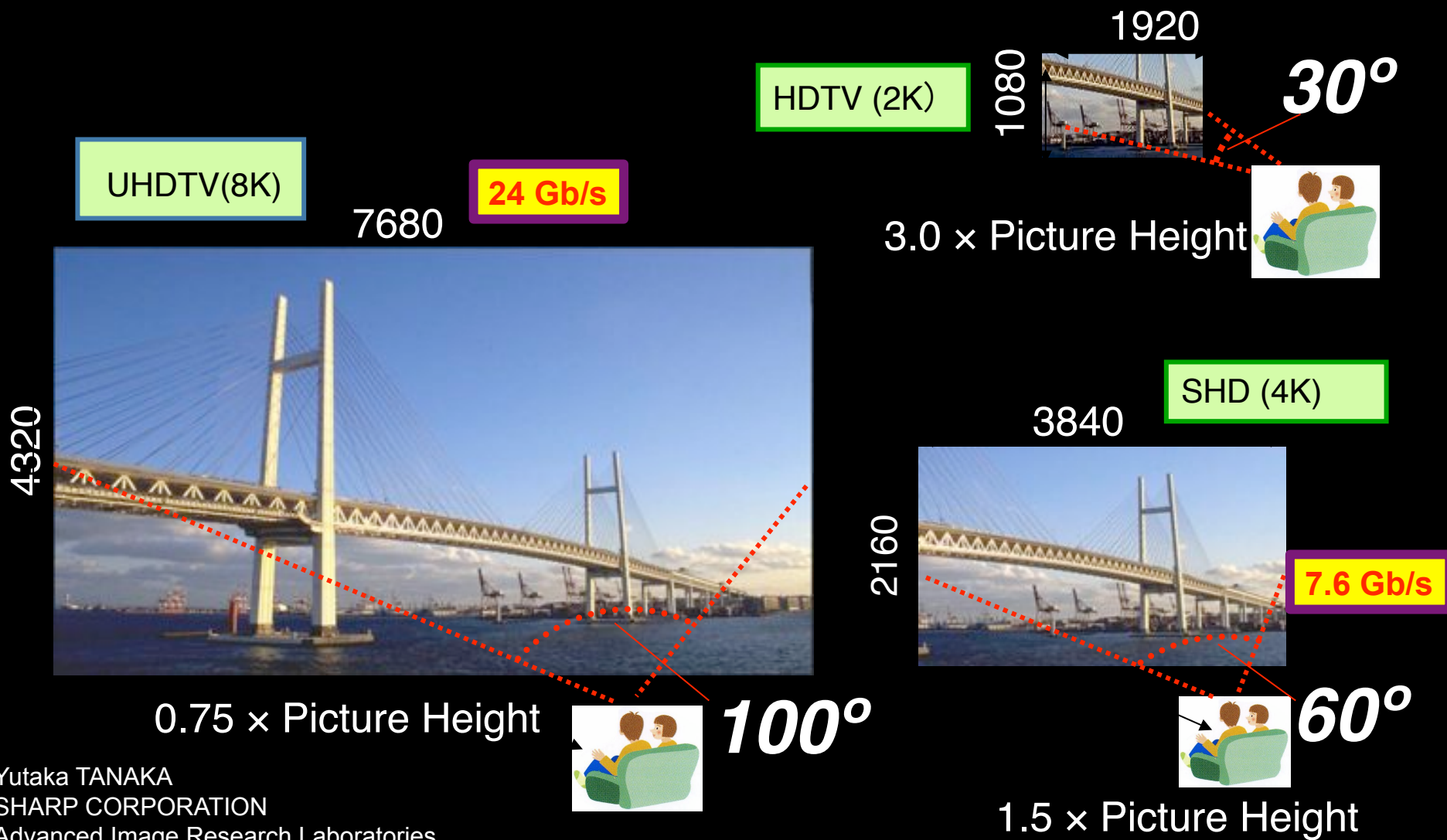


10 Gb/s dedicated network

270 ms RTT

7.6 Gb/s

CineGrid @ Holland Festival 2007

# CineGrid: 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

1920

HDTV (2K)

1080

**30°**

3.0 × Picture Height

UHDTV(8K)

**24 Gb/s**

7680

SHD (4K)

3840

4320

2160

**7.6 Gb/s**

0.75 × Picture Height

**100°**

**60°**

1.5 × Picture Height

Yutaka TANAKA
SHARP CORPORATION
Advanced Image Research Laboratories

# CineGrid portal

CineGrid **distribution center Amsterdam**

Home | About | Browse Content | cinegrid.org | cinegrid.nl

## Amsterdam Node Status:

node41:
Disk space used: 8 GiB
Disk space available: 10 GiB

## Search node:

[ Search ]

## Browse by tag:

amsterdam animation
antonacci blender boat
bridge bunny cgi delsa holland
hollandfestival
leidschestraat
muziekgebouw
nieuwmarkt opera prague ship
train tram trams waag

UvA [logo] Universiteit van Amsterdam

# CineGrid Amsterdam

Welcome to the Amsterdam CineGrid distribution node. Below are the latest additions of super-high-quality video to our node.

For more information about CineGrid and our efforts look at the about section.

# Latest Additions

## Wypke

Wypke

**Available formats:**
4k dst (4.8 KB)
**Duration:** 1 hour and 8 minutes
**Created:** 1 week, 2 days ago
**Author:** Wypke
**Categories:**

## Prague Train

Steam locomotive in Prague.

**Available formats:**
4k dst (3.9 KB)
**Duration:** 27 hours and 46 minutes
**Created:** 1 week, 2 days ago
**Author:** CineGrid
**Categories:** delsa prague train

## VLC: Big Buck Bunny

(c) copyright Blender Foundation | http://www.bigbuckbunny.org

**Available formats:**
1080p HPEG4 (1.1 GB)
**Duration:** 1 hour and 0 minutes
**Created:** 1 month, 1 week ago
**Author:** Blender Foundation
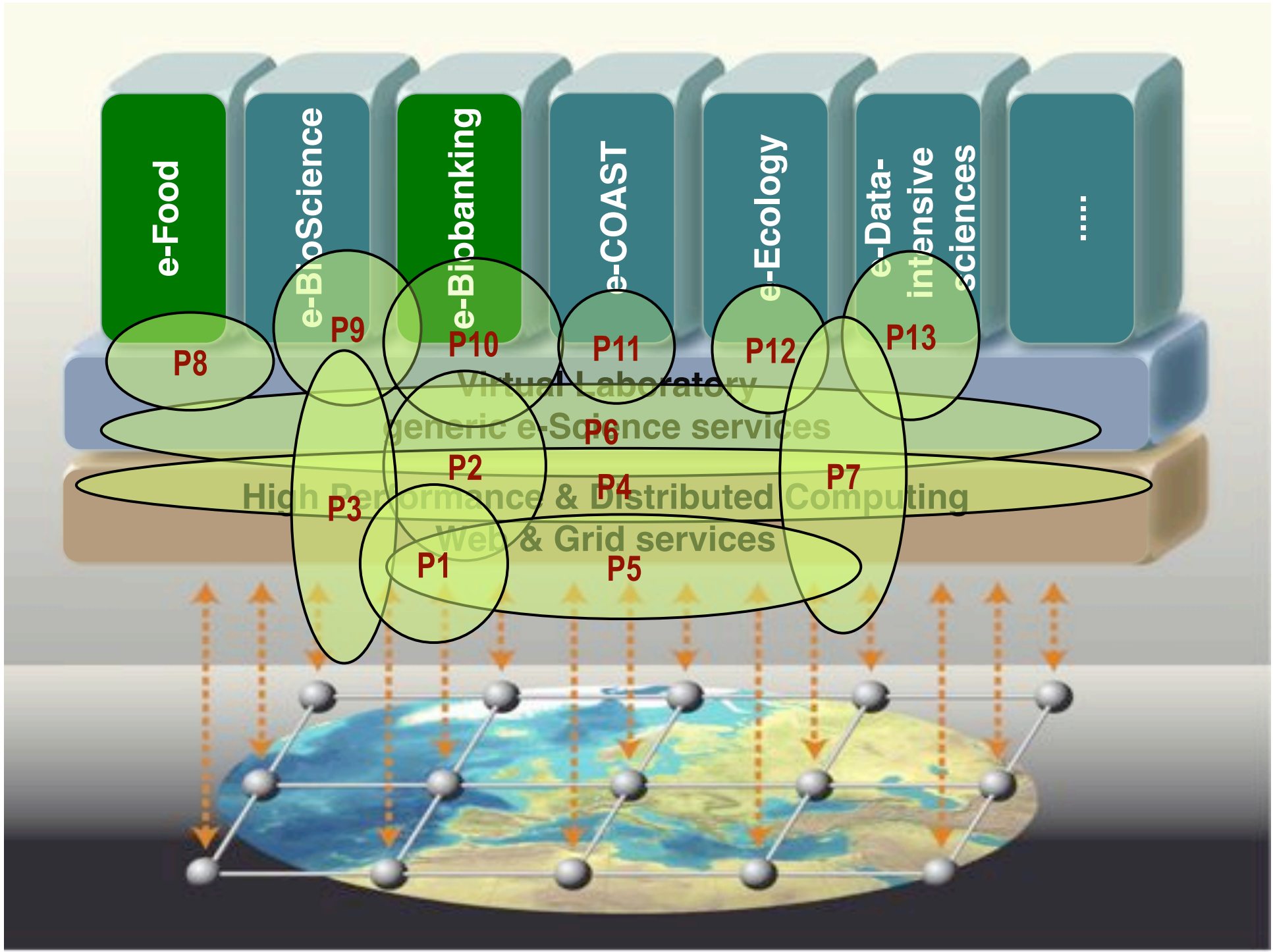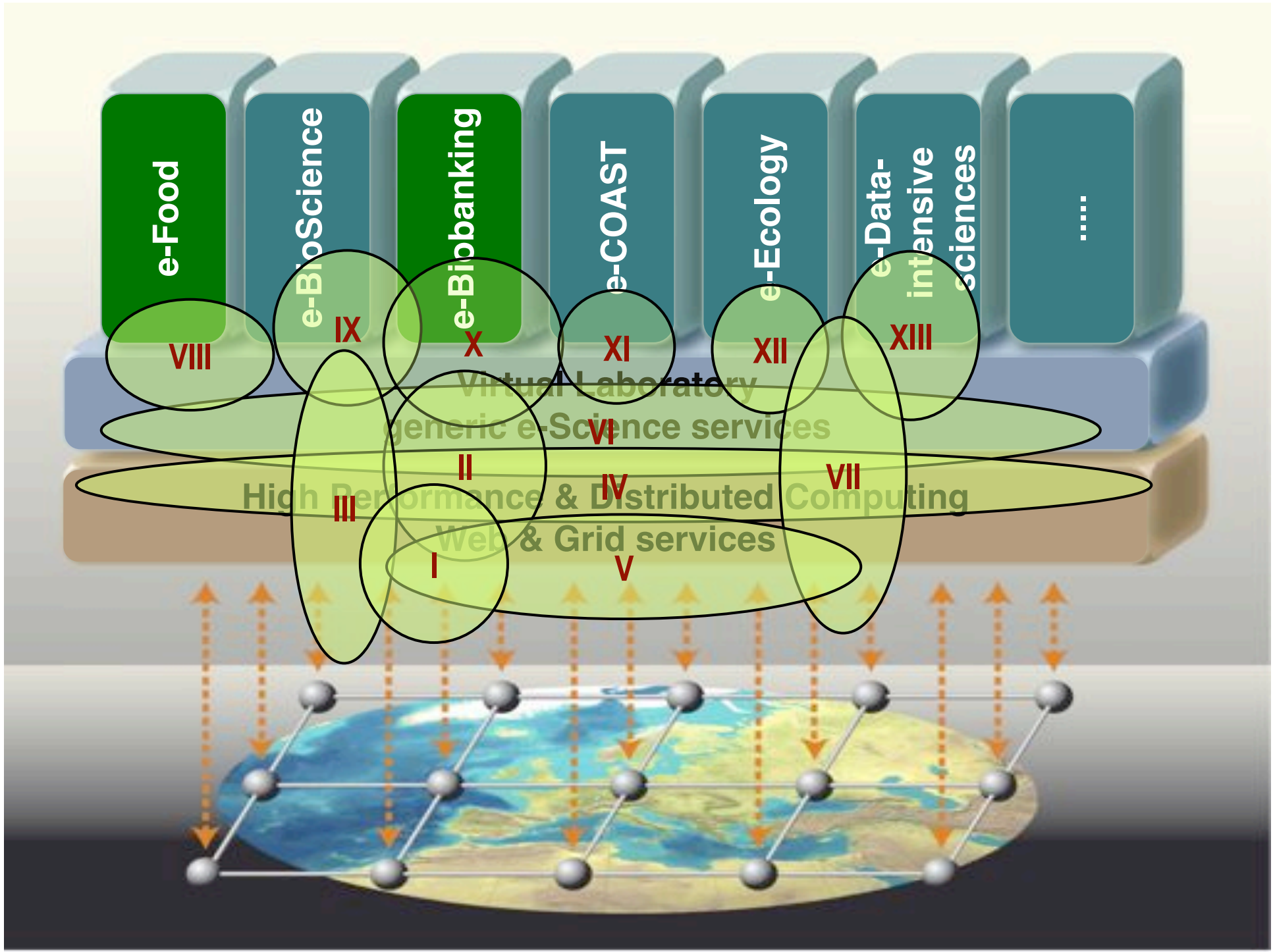**Categories:** animation blender bunny
cgi

e-Food

e-BioScience

e-Biobanking

e-COAST

e-Ecology

e-Data-intensive sciences

.....

VIII

IX

X

XI

XII

XIII

Virtual Laboratory
generic e-Science services

VI

II

IV

VII

III

High Performance & Distributed Computing
Web & Grid services

I

V

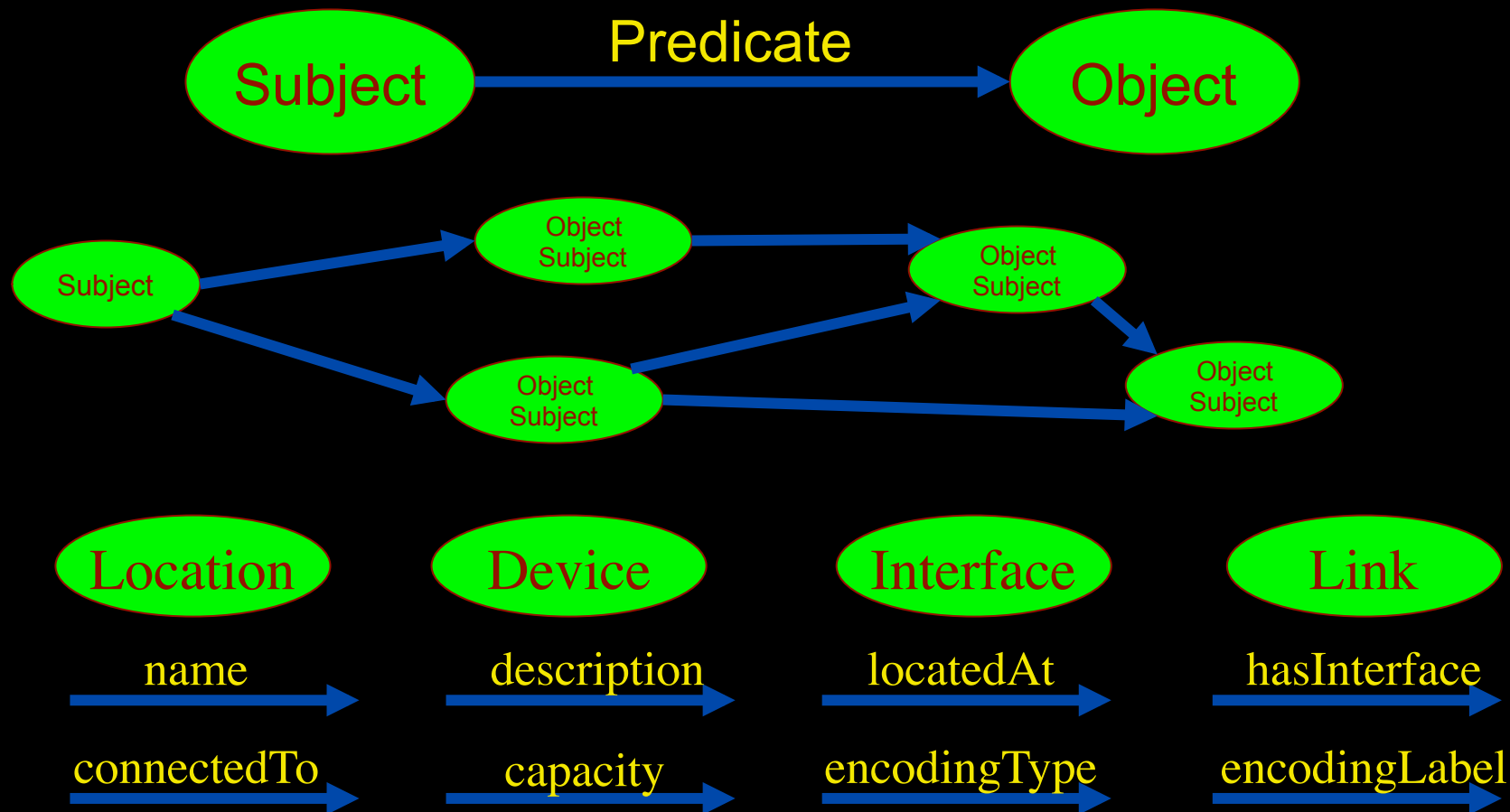# Network Description Language

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

Subject —— Predicate ——▶ Object

Subject
Object Subject
Object Subject
Object Subject
Object Subject

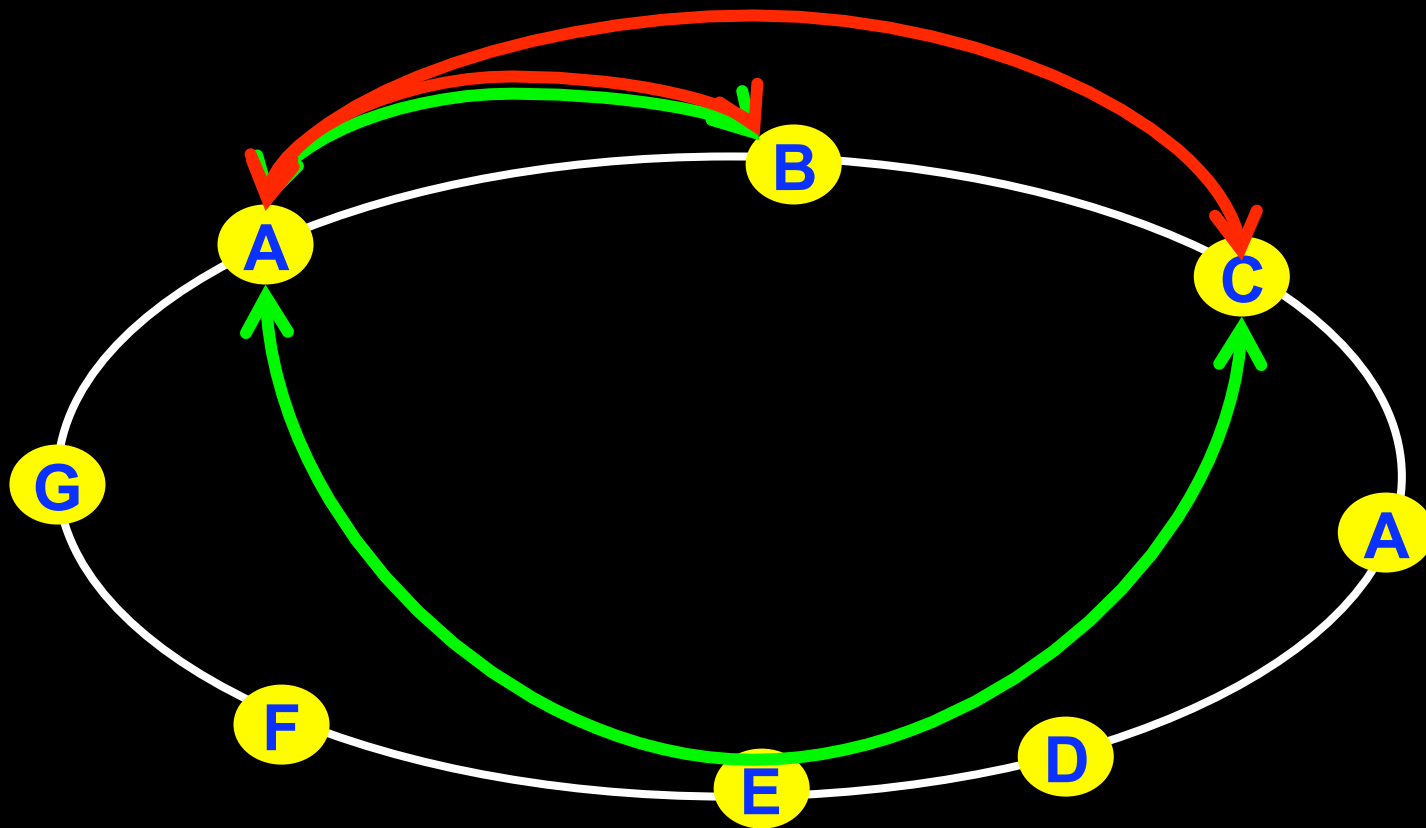| Location | Device | Interface | Link |
|---|---|---|---|
| name | description | locatedAt | hasInterface |
| connectedTo | capacity | encodingType | encodingLabel |

# The Problem

I want AC and AB

Success depends on the order of requests
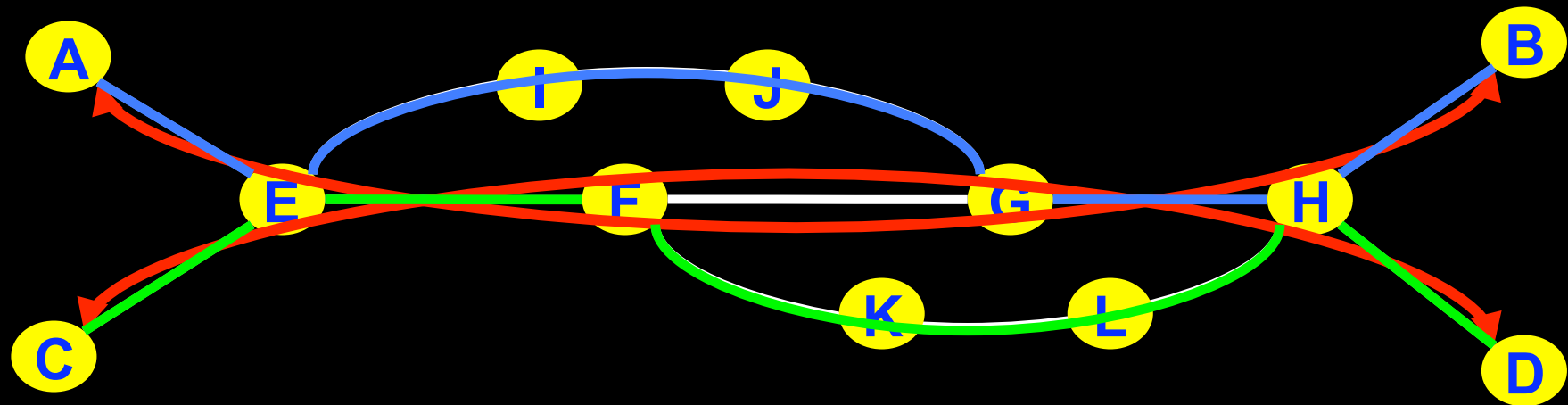
Wouldn't it be nice if I could request [AB, AC, ...]

# Another one ☺

I want AB and CD

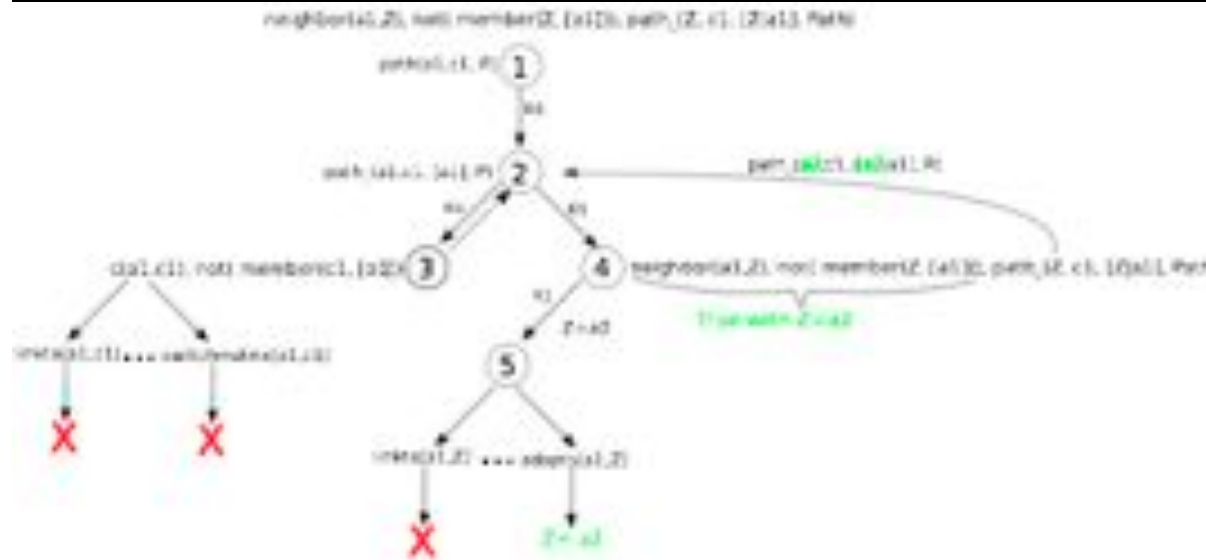Success does not even depend on the order!!!

# NDL + PROLOG

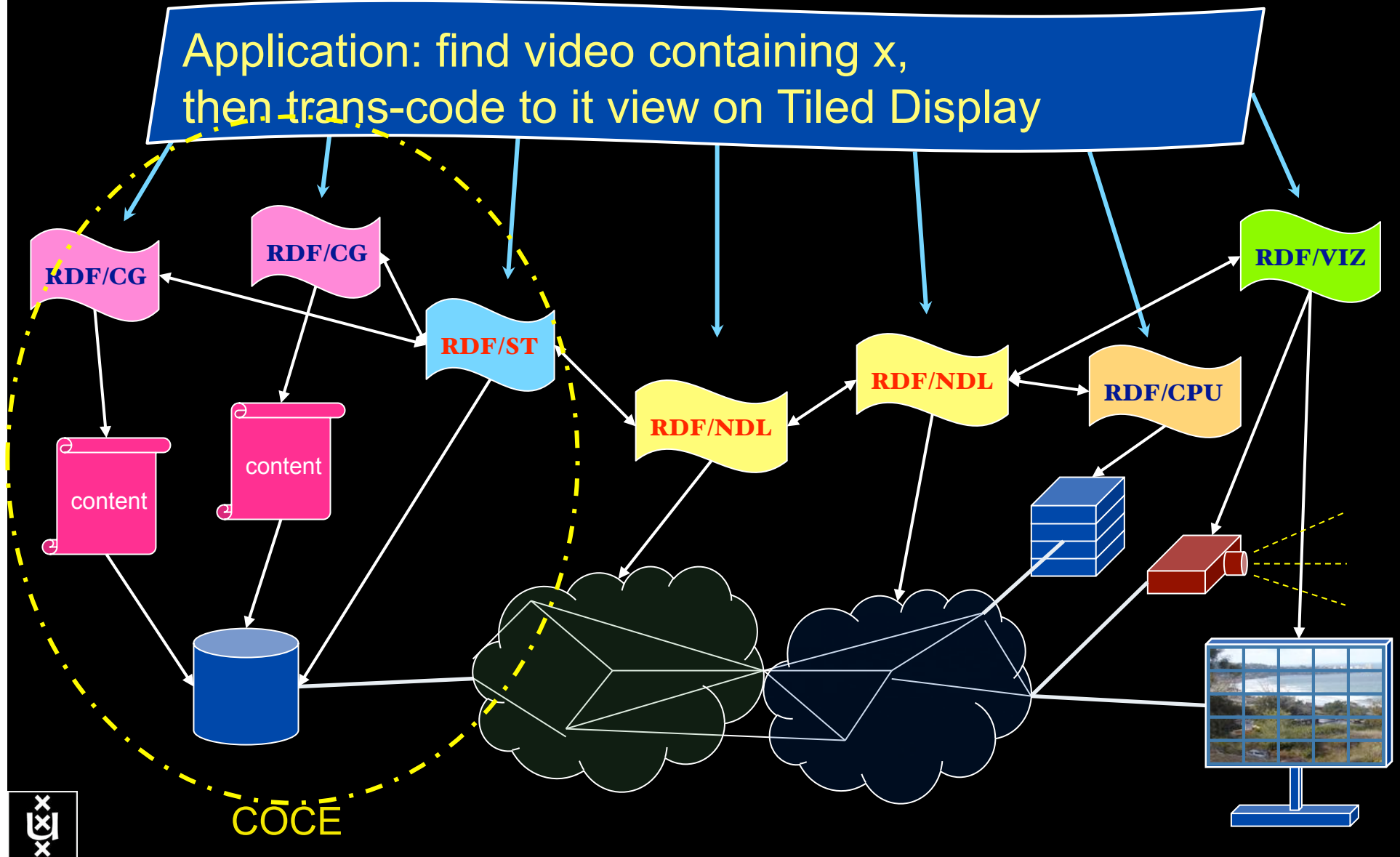Research Questions:
- order of requests
- complex requests
- usable leftovers



- Reason about graphs

- Find sub-graphs that comply with rules

- It finds solutions to previous slides!

# RDF describing Infrastructure "I want"

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

RDF/CG

RDF/CG

RDF/CG

RDF/ST

RDF/NDL

RDF/NDL

RDF/CPU

RDF/VIZ

content

content

COCE

# 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 turns a room full of pc's in 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      ->      ?

# 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}
```
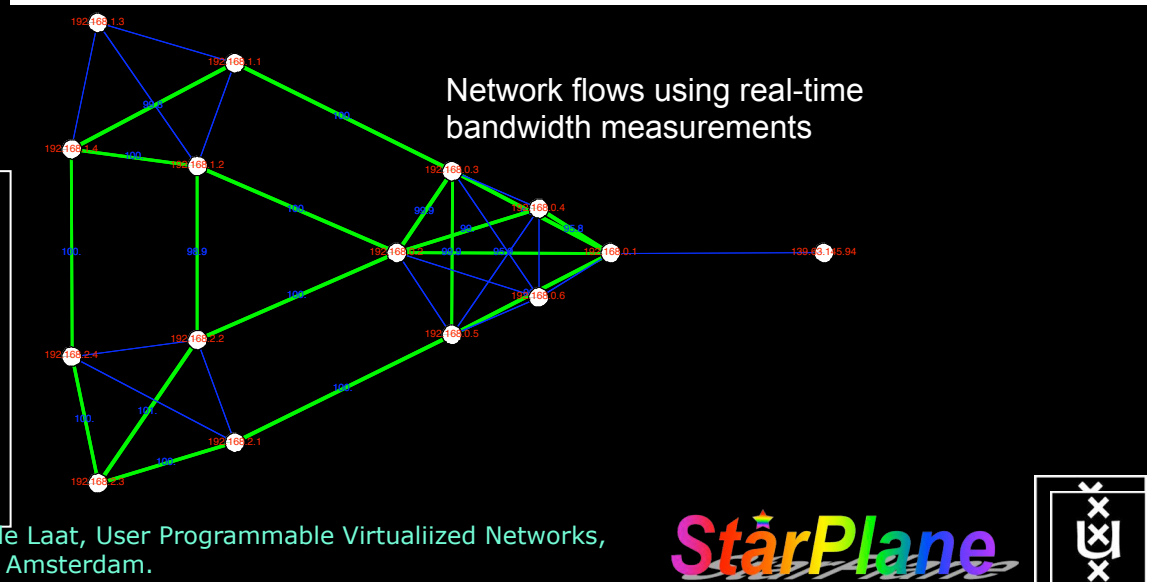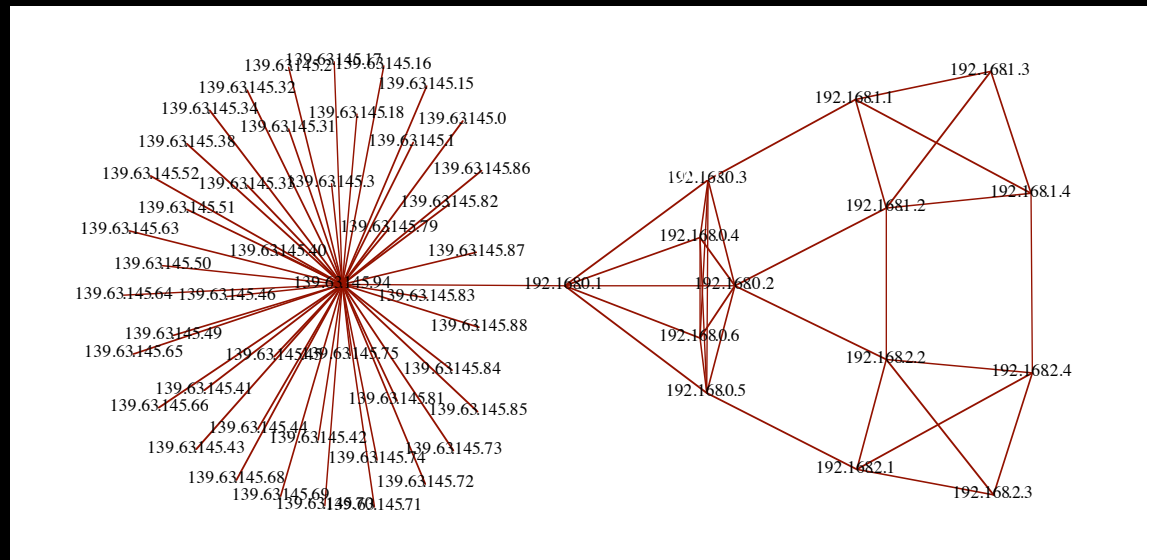
### Transaction on shortest path with tokens

```
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

# DAS-3 Cluster Architecture

head node (2)

Fast interconnect

10 Gb/s
Ethernet lanphy

**NORTEL**

**MYRINET**

To local
University

To SURFnet

UvA-node

10 Gb/s
Ethernet lanphy

1 Gb/s
Ethernet

8 * 10 Gb/s from
bridgenodes

Local interconnect

85 (40+45) compute nodes
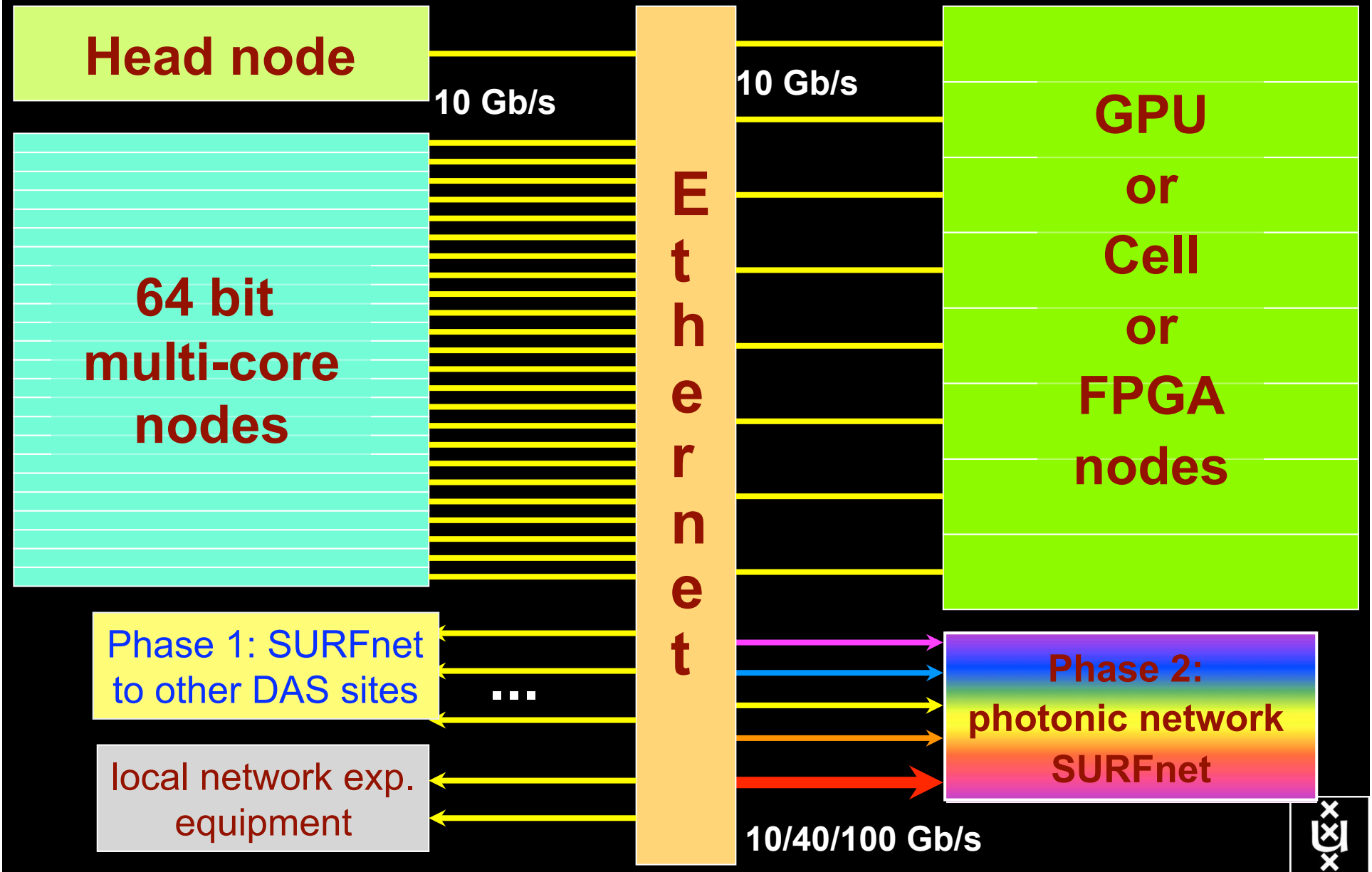
DAS-3

DAS-4 Proposed Architecture

# Themes for next years

- Network modeling and simulation
- Cross domain Alien Light switching
- eScience infrastructure virtualization (NSI)
- Photonic networking -> Tb/s
- Capacity & Capability
- Data handling, integrity, security, privacy
- Reasoning about services
- Fault tolerance, Fault isolation, monitoring
- eScience Data and Media specific services
- Cloud paradigm, green compute&store&net&viz
- ENERGY dependency! (2009: 1Wy=1€)

# Quotes from OnVector 2008

prof. Ken-Ichi Sato:

- It is very difficult to predict future services, however, video is expected to be the king media used for bit rate demanding services. High-quality video technologies are rapidly advancing.

- TCP/IP bottleneck is becoming more and more tangible. It will limit the future envisaged network expansion -the energy bottleneck and throughput bottleneck need to be resolved.

- Fast optical circuit/path switching will play the key role to create cost effective and bandwidth abundant future networks.

- Hierarchical optical path network and the node technologies are very important, and hence they need to be fully developed soon.

# Quotes from OnVector 2008

- dr. Kazuo Hagimoto:
- NTT is developing a system that automatically generates metadata such as title, summary, and key words that are extracted from voice or subtitles.

dr. Shimizu:

- Applications for Tbit networks:
  - High Resolution Simulation
  - Weather Forecast
  - Earthquake Forecast
  - City Planning
  - Digital Engineering
  - Nano Device Engineering
  - Protein Structural Analysis

# Quotes from OnVector 2008

prof. Larry Smarr:

- Interconnecting Regional Optical Networks
  Is Driving Campus Optical Infrastructure Deployment

prof. Ed Seidel:

- Petascale computing will not only provide huge data, but will demand new computing modalities
- Will place new demands on networking, data management, visualization, resource co- allocation
- Applications need to be configurable for the new type of infrastructure, need to be aware of environment
- If we don't solve these problems, people will use machines anyway, but science will suffer!

Bill s'Arnaud:

- "Optical networks (as opposed to electronic routed networks) have much smaller carbon footprint"

# Interactive programmable networks

# Questions ?

A Declarative Approach to Multi-Layer Path Finding Based on Semantic Network Descriptions.

http://delaat.net:/~delaat/papers/declarative_path_finding.pdf

Thanks: Paola Grosso & Jeroen vd Ham & Freek
Dijkstra & team for several of the slides.

SURF NET