# The **UltraLight** Project and the Challenges of Data Storage and Networking for the CMS Experiment

**UltraLight Collaboration**

Caltech (lead inst.)
BNL
Michigan
MIT
Florida
Florida International
FNAL
San Diego
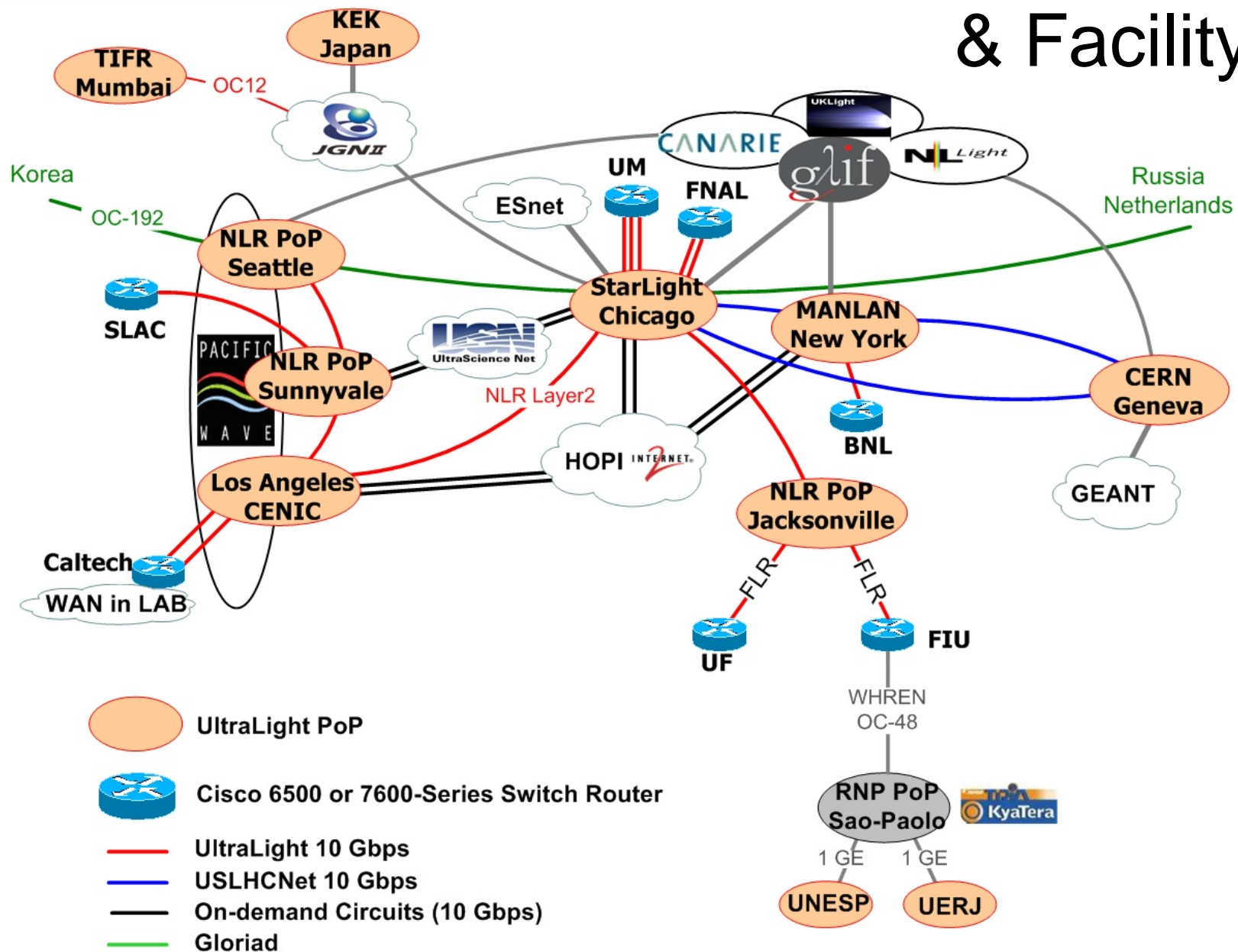SLAC
Vanderbilt

R. Cavanaugh

University of Florida

ICFA06

Krakow, Poland

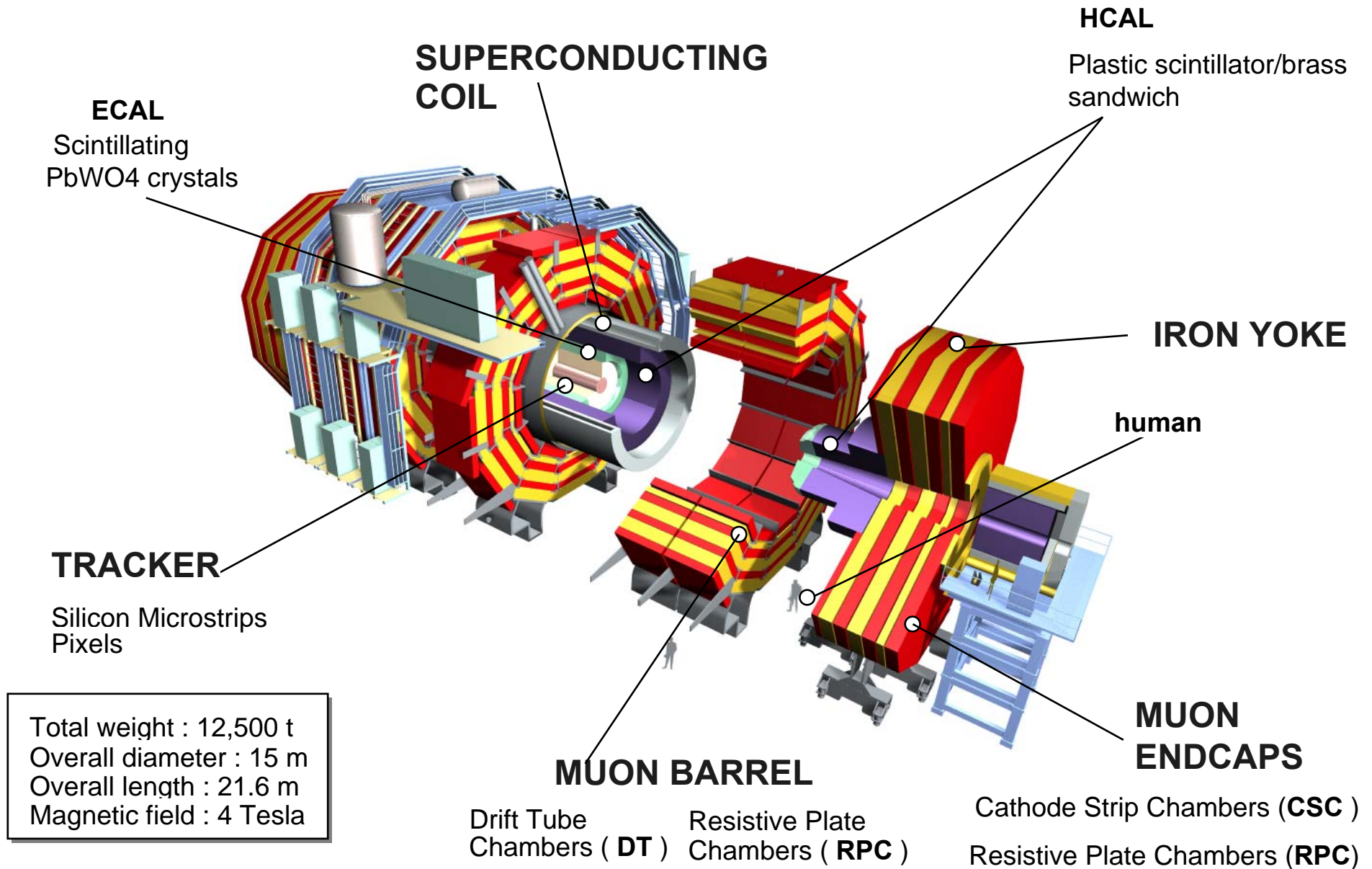# **UltraLight : A New Class of Integrated Information Systems**

- **Expose the Network as an Actively Managed Resource**

- **Based on a "Hybrid" packet- and circuit-switched optical network infrastructure**
  - Ultrascale Protocols (e.g. FAST) and Dynamic Optical Paths
- **Monitor, Manage and Optimize resources in real-time**
  - Using a set of Agent-Based Intelligent Global Services
- **Leverages already-existing, developing software infrastructure in round-the-clock operation:**
  - MonALISA, GAE/Clarens, OSG

- **Exceptional Support from**
  - Industry: Cisco & Calient
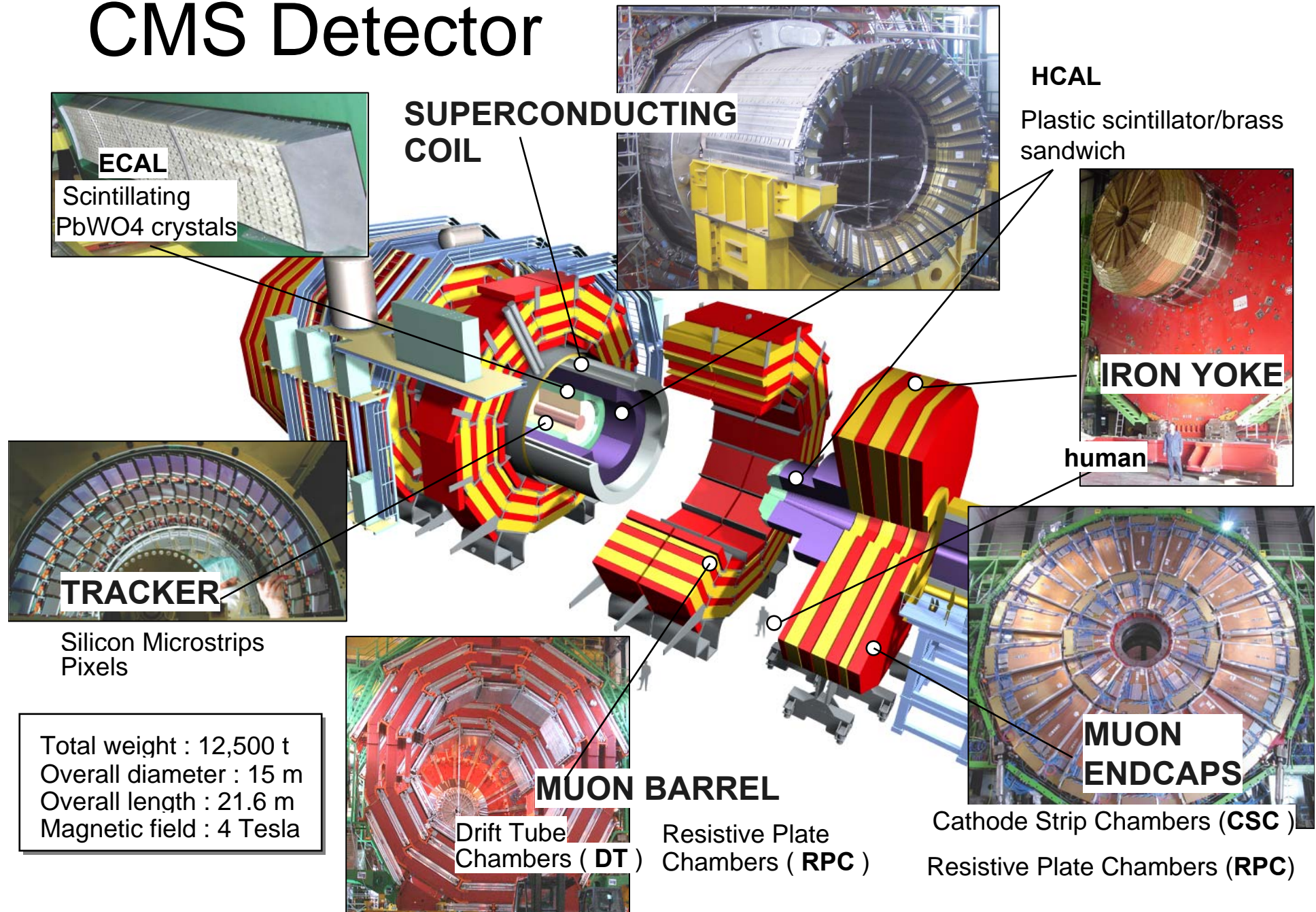  - Research community: NLR, CENIC, Internet2/Abilene, ESnet

# Network Laboratory Testbed & Facility

# CMS Detector

**ECAL**
Scintillating
PbWO4 crystals

**SUPERCONDUCTING COIL**

**HCAL**
Plastic scintillator/brass sandwich

**IRON YOKE**

**human**

**TRACKER**

Silicon Microstrips
Pixels

Total weight : 12,500 t
Overall diameter : 15 m
Overall length : 21.6 m
Magnetic field : 4 Tesla

**MUON BARREL**

Drift Tube
Chambers ( **DT** )

Resistive Plate
Chambers ( **RPC** )

**MUON ENDCAPS**

Cathode Strip Chambers (**CSC** )

Resistive Plate Chambers (**RPC**)

# CMS Detector

**ECAL**

Scintillating
PbWO4 crystals

**SUPERCONDUCTING COIL**

**HCAL**

Plastic scintillator/brass sandwich

**IRON YOKE**

**human**

**TRACKER**

Silicon Microstrips
Pixels

Total weight : 12,500 t
Overall diameter : 15 m
Overall length : 21.6 m
Magnetic field : 4 Tesla

**MUON BARREL**

Drift Tube
Chambers ( **DT** )

Resistive Plate
Chambers ( **RPC** )

**MUON ENDCAPS**

Cathode Strip Chambers (**CSC** )

Resistive Plate Chambers (**RPC**)

# CMS Detector

HCAL

Plastic scintillator/brass sandwich

SUPERCONDUCTING COIL

ECAL

Scintillating PbW

```
Detector        No. Channels   Sensors
-----------     ------------   -------------------------------

Vertex           80 000 000    Pixels
Tracker          16 000 000    Silicon Microstrips
Preshower           512 000    Silicon
Calorimeters        125 000    Scintillating Crystals
                               Scintillator / Brass sandwich
Muons             1 000 000    Drift/Strip/Plate Chambers
-----------     ------------   -------------------------------

Total          ~100 000 000    Channels (most not read out!)



Event Size                     ~1-2  MB (with selective readout)
DAQ Bandwidth                  ~200  MB / sec
```

TRA

Silico Pixel

Total w
Overal
Overall length : 21.6 m
Magnetic field : 4 Tesla

MUON BARREL

Drift Tube Chambers ( **DT** )

Resistive Plate Chambers ( **RPC** )

Cathode Strip Chambers (**CSC** )

Resistive Plate Chambers (**RPC**)

# Main Science Problem



**"Haystack"**

**"Needle"**

σ    LHC    √s=14TeV    L=10³⁴cm⁻²s⁻¹    rate    ev/year

barn

mb

μb

nb

pb

fb

jet E_T or particle mass (GeV)

**ON-line** ⟶

**OFF-line**

LEVEL-1 Trigger
Hardwired processors (ASIC, FPGA)
Pipelined massive parallel

**Tightly coupled**      **Loosely coupled**

HIGH LEVEL Triggers
Farms of
processors

Reconstruction&ANALYSIS
TIER0/1/2
Centers

**Very Powerful
Very Efficient**

**"Sift"**

| 25ns | 3μs | ms | sec | hour | year |
|------|-----|-----|-----|------|------|
| $10^{-9}$ | $10^{-6}$ | $10^{-3}$ | $10^{0}$ | $10^{3}$ | $10^{6}$ sec |
| | | | **Giga** | **Tera** | **Petabit** |

# Main Science Problem

σ  LHC  √s=14TeV  L=$10^{34}$cm$^{-2}$s$^{-1}$  rate  ev/year

barn

**LEVEL-1 Trigger**
**Hardwired processors  (ASIC, FPGA)**
**Pipelined massive parallel**

σ inelastic  L1 input  GHz  $10^{16}$

$10^{15}$

mb  b$\bar{\text{b}}$

MHz  $10^{13}$  $10^{14}$

max Detector output
max L1 output
max HLT input

**HIGH LEVEL Triggers**
**Farms of processors**

**Individual TB transactions should finish in *minutes to hours*, rather than *hours to days***

μb  jets

kHz  $10^{10}$  $10^{11}$  $10^{12}$

W
Z
W→lν

HLT output  $10^{9}$

Z→lν
t$\bar{\text{t}}$

nb

Hz  $10^{7}$  $10^{8}$

**Reconstruction&ANALYSIS**
**TIER0/1/2 Centers**

gg→H$_{SM}$

SUSY q̃q̃+q̃g̃+g̃g̃
tanβ=2, μ=m$_{\tilde{g}}$=m$_{\tilde{q}}$/2
tanβ=2, μ=m$_{\tilde{g}}$=m$_{\tilde{q}}$

$10^{6}$

q$\bar{\text{q}}$→qqH$_{SM}$

$10^{5}$

*Requests from Multiple users for Multiple types of… …Multiple times!*

H$_{SM}$→γγ
h→γγ

$10^{4}$

**New Physics Searches ⇒ *multi-Terabyte* scale Datasets!**

fb

H→2Z→4l

$10^{2}$

$10$

Z$_{SM}$→3γ  scalar LQ  Z$_{\eta}$→2l  μHz

$1$

50  100  200  500  1000  2000  5000
jet E$_T$ or particle mass (GeV)

| 25ns | 3μs | ms | sec | hour | year |
|---|---|---|---|---|---|
| $10^{-9}$ | $10^{-6}$ | $10^{-3}$ | $10^{0}$ | $10^{3}$ | $10^{6}$ sec |
| | | | Giga | Tera | Petabit |

# LHC Computing Grid Hierarchy

**CMS Experiment**

**A lot of effort has been devoted to understand T0-T1 & T1-T1 interactions (manageable)**

**Online System**

100-1500 MBytes/s

**Most IT resources outside of CERN**

*Tier 0*

**CERN Computer Center > 20 TIPS**

10-40 Gbps

*Tier 1*

**Germany**    **France**    **Italy**    **USA**    **• • •**

**"Production" Data Processing**

2.5-40 Gbps

Data Analysis "Chaotic"

*Tier 2*    CalTech    Florida    San Diego    • • •

1-2.5 Gbps

*Tier 3*    Institute    Institute    Institute    Institute

Physics cache    1-10 Gbps

*Tier 4*    PCs

- ~10s of Petabytes/yr by 2007-8
- ~1000 Petabytes in < 10 yrs

**2008 resources**

Chart (percentages 0%–100% for CPU, Disk, Tape):
- CERN
- Tier1
- Tier2

# LHC Computing Grid Hierarchy



CMS Experiment

Online System

100-1500 MBytes/s

*Tier 0* — CERN Computer Center > 20 TIPS

10-40 Gbps

*Tier 1* — Germany, France, Italy, USA, ...

2.5-40 Gbps

*Tier 2* — CalTech, Florida, San Diego, ...

1-2.5 Gbps

*Tier 3* — Institute, Institute, Institute, Institute

Physics cache

*Tier 4* PCs

1-10 Gbps

Most IT resources outside of CERN

"Production" Data Processing

Data Analysis "Chaotic"

2008 resources

- ~10s of Petabytes/yr by 2007-8
- ~1000 Petabytes in < 10 yrs

**Much less effort has been devoted to understand T1-T2 & T2-T2 interactions (less manageable & more complex)**

# LHC Computing Grid Hierarchy



CMS Experiment

Online System

100-1500 MBytes/s

*Tier 0*

CERN Computer Center > 20 TIPS

Most IT resources outside of CERN

10-40 Gbps

*Tier 1*

Germany   France   Italy   USA   •••

"Production" Data Processing

Future Goal (partially achieved)

UltraLight

2.5-40 Gbps

Data Analysis "Chaotic"

*Tier 2*

CalTech   Florida   San Diego   Michigan

1-2.5 Gbps

*Tier 3*

Institute   Institute   Institute   Institute

Physics cache   *Tier 4*   PCs

1-10 Gbps

- ~10s of Petabytes/yr by 2007-8
- ~1000 Petabytes in < 10 yrs
- UltraLight Network Provides a Data "Bus"
  - T1-T2 & T2-T2 Non-hierarchical data flows
  - Natural from Data Analysis point of view!

2008 resources

Legend: CERN, Tier1, Tier2

CPU   Disk   Tape

# Recent Experience with CMS Transfers across the LHC Grid Hierarchy

# CMS Transfer Sinks



**Data distributed somewhat democratically across the different Tier-1/2s…**

**… aggregate average ~ 1 Pb / Month** (meets LHC baseline)

LCG Service Challenge 4

Data Transferred (TB)

4000
3500
3000
2500
2000
1500
1000
500
0

2006-05-22    2006-06-18    2006-07-16    2006-08-13    2006-09-10

Day

**"not bad" :  many (most?) Tier2 analysis centres not starved…**
**…Critical for delivering Physics to the Collaboration!!**

# CMS Transfer Sources



LCG Service Challenge 4

**However, not all Data Providers are equal!** **Importance of :**
**(1) non-CERN resources**
**(2) any T2 – any T1 routes**
**(3) some T2s**

FNAL

CERN

Data Transferred (TB)

4000
3500
3000
2500
2000
1500
1000
500
0

2006-05-22    2006-06-18    2006-07-16    2006-08-13    2006-09-10

Day

**Fault tolerance suggests usefulness of dynamically (not statically) routed flows**

# CMS Transfer Quality of Service

## LCG Service Challenge 4



Fraction of successful Transfers

# CMS Transfer Quality of Service

## LCG Service Challenge 4



*Quality has improved* significantly (amount of green), but much work remains (white gaps and amount of red/yellow)

Day

**End-host readiness is the biggest problem! Very labour intensive…**

Service outages: 2 days out of 120

| 0-10% | 10-20% | 20-30% | 30-40% | 40-50% | 50-60% | 60-70% | 70-80% | 80-90% | 90-100% | 100+% |

Fraction of successful Transfers

# What is Achievable in the *Ultra Light* Lab?

- **Supercomputing 2005**
- **151 Gbps peak rate**
- **100+ Gbps sustained throughput for hours**
- **475 Terabytes of physics data transported in less than 24 hours**
- **Sustained rate of 100+ Gb/s translates to > 1 Petabyte per day**



**WAN Total Traffic**

**WAN Total Traffic**

**Not a production-level exercise…**
**…Required tremendous amount of continuous manual attention! (~1 ♱ /λ)**

# *Ultra Light* & CMS Experiences Expose <u>The End-host Challenge</u>

- Intense level of manual interventions (current reality)
  - The plague of data transfer performance!!
  - Requires distributed knowledge (expensive & continually out-of-date)
  - Introduces long failure and upgrade response-times
- Capable End-Host hardware often not deployed!!
- Storage technologies improving – still relatively young
  - SRM (WLCG, EGEE, OSG)
    - standard interface (only)
  - CASTOR (CERN)
    - Large scale (complex) storage (disk+tape) solution
    - Fully managed (queues, priorities, etc) data management
  - DPM (EGEE), dCache (DESY / FNAL) & LSTORE (Vanderbilt)
    - Smaller scale (simpler, but still complex), aggregate disk management
    - No queues or prioritized data management
- Clear need for E2E management, automation and self-healing!!
  - Has proven to be difficult in a loosely coupled "social" environment

# Network R&D
# Global Planning Services

- ## VINCI :
  - Virtual Intelligent Networks for Computing Infrastructures
  - Based on existing MonALISA framework
  - What does vinci do?

- ## LISA :
  - Localhost Information Service Agent
  - Monitors end-systems
    - User
    - servers



- ## UltraLight contributing to a VINCI prototype
  - production level system will come from elsewhere

http://ultralight.caltech.edu/web-site/gae/movies/ml_optical_path/ml_os.htm

BC NEWS | News Fro...  The top news headline...  CNN.com - Breaking N...  Google News  Telecourses/e-Learnin...

umble!  All  I like it!  Search or Tag here  Share  Pages  Friends  Menu

3D Map

LSU_Calient

MCNC_OS

■ 10.12a....
● 10.14a....
▲ 10.10a....
♦ 10.14a....

Groups

TabPan

GMap

AS371

Topology

Load

WAN

VO JOBS

OS GMap

Multi-view

la-x3

calient3

calient2  calient1

la-x2

B - VMware Workstation

Edit  View  VM  Team  Windows  Help

FC3

cil@la-opt2:~/SCRIPTS

Edit  View  Terminal  Tabs  Help

```
and [ bbcp -F -f -S "ssh -x -a -oFallBackToRsh=
I -l %U %H bin/bbcp" -T "ssh -x -a -oFallBackTo
no %I -l %U %H bin/bbcp" -v -P 1 /dev/zero la-x
ev/null ] FINISHED    releasing Optical Path
ed by

cal P

ration took 56 ms

@la-opt2 SCRIPTS]$ ./MLBBCOPY la-x3 file1.data
```

Copy data from host la-x2 to la-x3

plications  Actions

Once path is set, transfer data

3D Map

LSU_Calient
MCNC_OS

■ 10.12a....
● 10.14a....
▲ 10.10a....
◆ 10.14a....

Groups

TabPan

GMap

AS371

Topology

Load

WAN

VO JOBS

OS GMap

Multi-view

B - VMware Workstation

Edit  View  VM  Team  Windows  Help

FC3

cil@la-opt2:~/SCRIPTS

Edit  View  Terminal  Tabs  Help

```
... g "ssh -x -a -oFallBackToRsh=
: Creating
DEST_DIR/
94.1 KB/s
le copied

and [ bbcp -F -f -S "ssh -x -a -oFallBack...Rsh=
I -l %U %H bin/bbcp" -T "ssh -x -a -oFallBackTo
no %I -l %U %H bin/bbcp" -v -P 1 file1.data la-
EST_DIR/file1.data ] FINISHED ... releasing Opt
Path
```

Finished with transfer

applications  Actions

la-x3

calient3

calient2

calient1

la-x2

http://ultralight.caltech.edu/web-site/gae/movies/ml_optical_path/ml_os.htm

BC NEWS | News Fro...  The top news headline...  CNN.com - Breaking N...  Google News  Telecourses/e-Learnin...

umble!  All  I like it!  Search or Tag here  Share  Pages  Friends  Menu

3D Map

Groups

■ 10.12a....
● 10.14a....
▲ 10.10a....
◆ 10.14a....

TabPan

GMap

AS371

Topology

Load

WAN

VO JOBS

OS GMap

Multi-view

LSU_Calient
MCNC_OS

la-x3

Releasing path

calient3

calient2          calient1

la-x2

B - VMware Workstation

Edit  View  VM  Team  Windows  Help

FC3

cil@la-opt2:~/SCRIPTS

Edit  View  Terminal  Tabs  Help

```
and [ bbcp -F -f -S "ssh -x -a -oFallBackToRsh=
I -l %U %H bin/bbcp" -T "ssh -x -a -oFallBackTo
no %I -l %U %H bin/bbcp" -v -P 1 /dev/zero la-x
ev/null ] FINISHED ... releasing Optical Path
ed by signal 2.

cal Path Released ...

ration took 44 ms

@la-opt2 SCRIPTS]$
```

Applications  Actions

http://ultralight.caltech.edu/web-site/gae/movies/ml_optical_path/ml_os.htm

BC NEWS | News Fro...  The top news headline...  CNN.com - Breaking N...  Google News  Telecourses/e-Learnin...

umble!  All  |  I like it!  |  Search or Tag here  |  Share  Pages  Friends  Menu

3D Map

LSU_Calient
MCNC_OS

Groups

■ 10.12a....
● 10.14a....
▲ 10.10a....
◆ 10.14a....

TabPan

GMap

AS371

Topology

3 - VMware Workstation

Edit  View  VM  Team  Windows  Help

FC3

cil@la-opt2:~/SCRIPTS

Edit  View  Terminal  Tabs  Help

o la-x3:/dev/null ]

and [ bbcp -F -f -S "ssh -x -a -oFallBackToRsh=
I -l %U %H bin/bbcp" -T "ssh -x -a -oFallBackTo
no %I -l %U %H bin/bbcp" -v -P 1 /dev/zero la-x
ev/null ] FINISHED ... releasing Optical Path

cal Path Released ...

ration took 44 ms

@la-opt2 SCRIPTS]$

Load

WAN

VO JOBS

OS GMap

Multi-view

Applications  Actions

la-x3

calient3

calient2          calient1
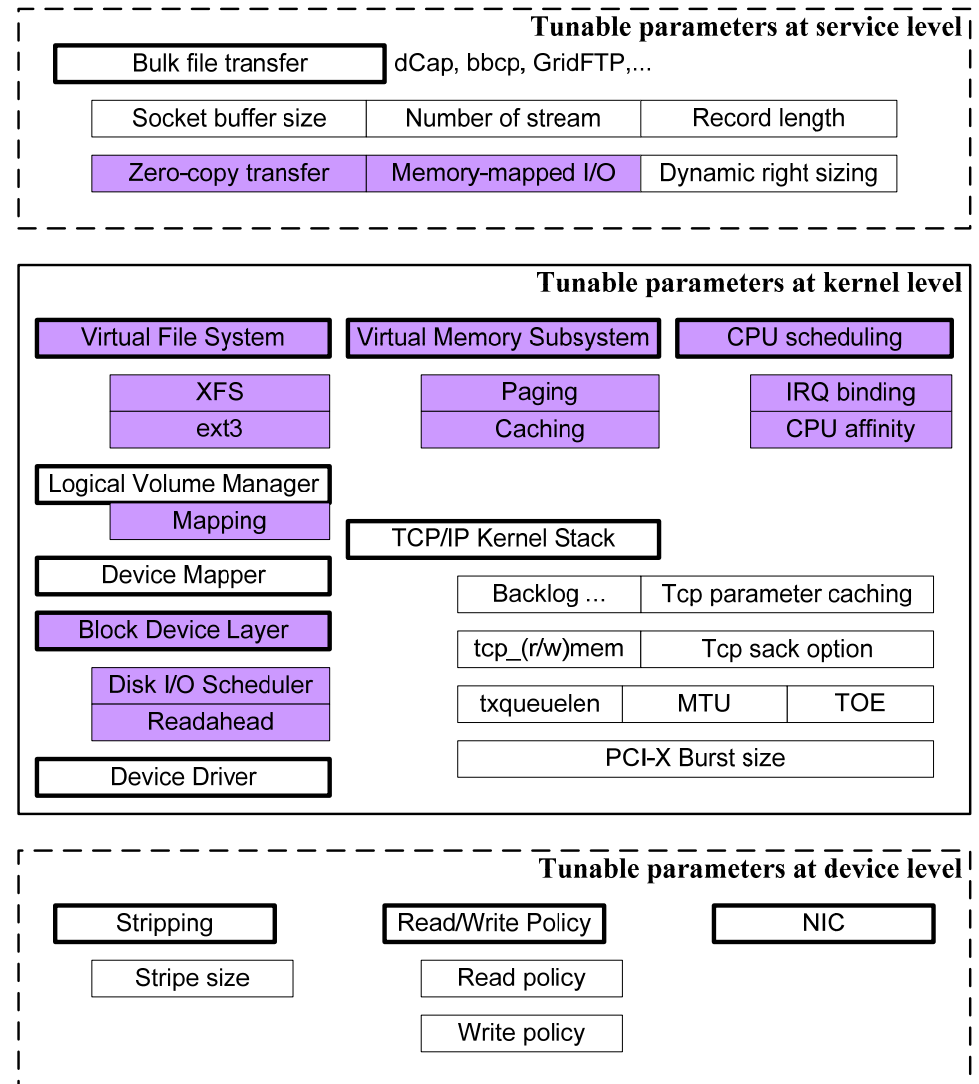
la-x2

# Storage R&D
# Global Planning Services

- UltraLight (optical networks in general) moving towards a managed "control plane"
  - Light-paths will be allocated/scheduled to data-flow requests via policy based priorities, queues, and advanced reservations
  - Clear need to match "Network Resource Management" with "Storage Resource Management"
    - Well known co-scheduling problem!
    - In order to develop an effective NRM, must understand and interface with SRM!
- End systems remain the current bottle-neck for large scale data transport over the WAN
  - Key to effective filling/draining of the pipe
  - Need highly capable hardware (servers, etc)
  - Need carefully tuned software (kernel, etc)

# UltraLight using a Wholistic, Multi-level Approach

- **End-host Device Technologies**
  - Choosing right H/W platform for the price ($20K)

- **End-host Software Stacks**
  - Tuning storage server for stable and high throughput

- **End-Systems Management**
  - Specifying quality of service (QoS) model for Ultralight storage server
  - SRM/dCache
  - LSTORE (& SRM/LSTORE)

- **Wide Area Testbeds (REDDnet)**
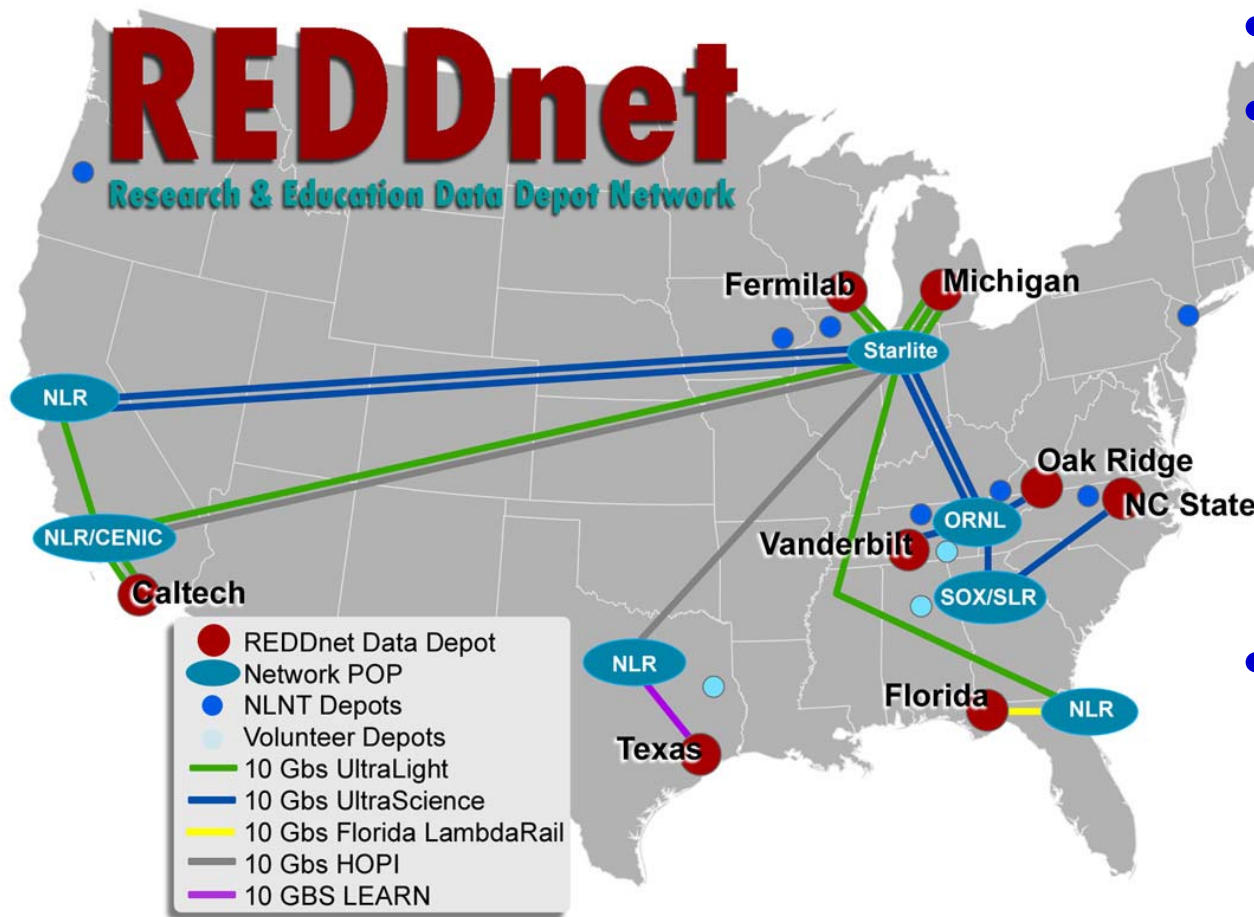
# Tunable Parameter Space

- **Multiple layers**
  - Service/AP level
  - Kernel level
  - Device level
- **Complexity of tuning**
  - Fine tuning is very complex task
  - Now investigating the possibility of *auto-tuning* daemon for storage server

**Tunable parameters at service level**

| Bulk file transfer | dCap, bbcp, GridFTP,... | |
|---|---|---|
| Socket buffer size | Number of stream | Record length |
| Zero-copy transfer | Memory-mapped I/O | Dynamic right sizing |

**Tunable parameters at kernel level**

| Virtual File System | Virtual Memory Subsystem | CPU scheduling |
|---|---|---|
| XFS / ext3 | Paging / Caching | IRQ binding / CPU affinity |

Logical Volume Manager
Mapping

TCP/IP Kernel Stack

Device Mapper

| Backlog ... | Tcp parameter caching |
|---|---|

Block Device Layer

| tcp_(r/w)mem | Tcp sack option |
|---|---|

Disk I/O Scheduler
Readahead

| txqueuelen | MTU | TOE |
|---|---|---|

| PCI-X Burst size | |
|---|---|

Device Driver

**Tunable parameters at device level**

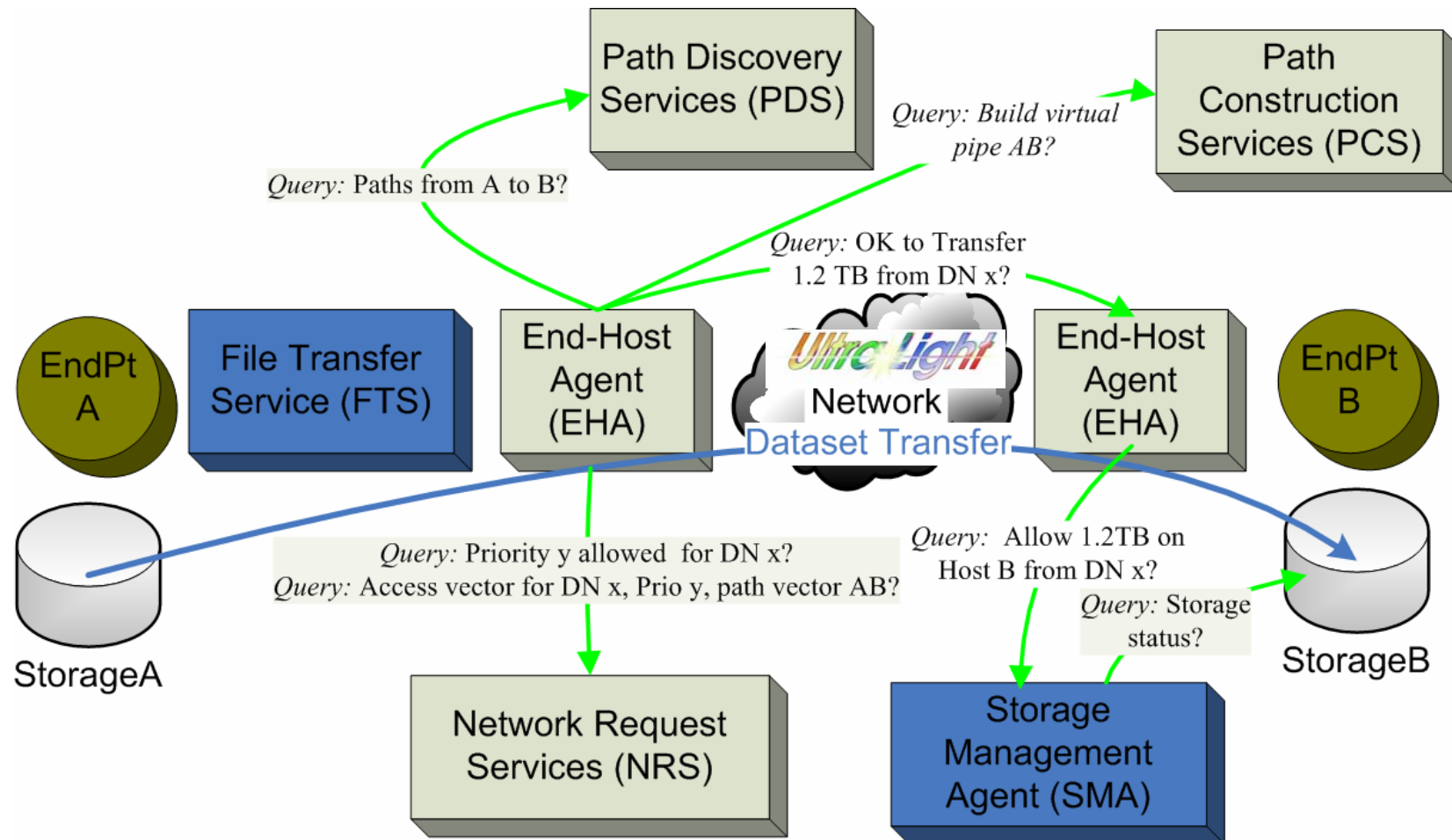| Stripping | Read/Write Policy | NIC |
|---|---|---|
| Stripe size | Read policy | |
| | Write policy | |

# REDDnet
## Research and Education Data Depot Network

**Powered by UltraLight**



- Led by Vanderbilt
- 8 initial sites
- Multiple disciplines
  - Sat imagery (AmericaView)
  - HEP
  - Terascale Supernova Initative
  - Structural Biology
  - Bioinformatics
- Storage
  - 500TB disk
  - 200TB tape

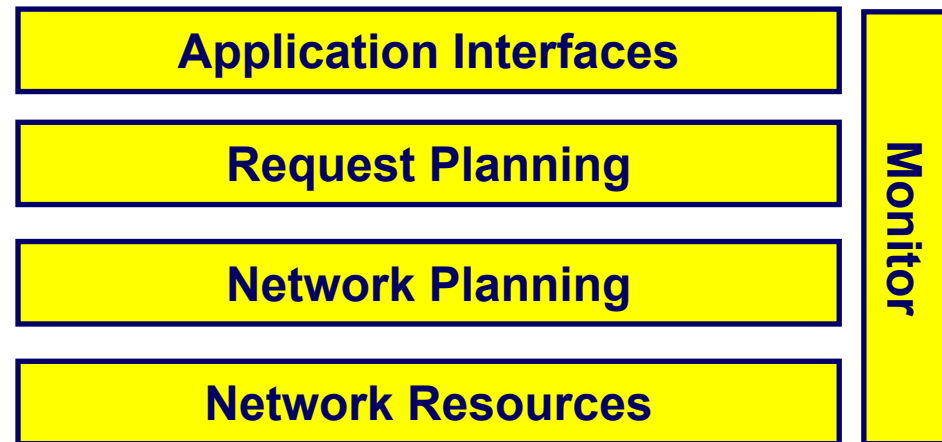# *Sets the Stage for Testing Fully Integrated System*

# Summary

- UltraLight is a global Laboratory, uniquely positioned
  - Spans Tier-0, some Tier-1s, several Tier-2s, and some Tier-3s
  - Includes participation from ATLAS (not discussed in this talk)
- End-hosts remain serious bottleneck in delivering CMS data to the higher Tiers for physics data analysis
  - Human in the loop problem
  - Incapable hardware (sometimes, perhaps even often)
  - Fine tuning of services
- UltraLight working to address these (generic) problems by
  - Researching and developing Global Planning Services
  - Using a wholistic approach (devices, parameters, services, WAN)

- Final thought:
  - Not only critical for LHC, also important for preparing for SLHC! (HEP always asks: how much time to 4x our data sample?)

# Make UltraLight available to Physics applications and their environments

- **Unpredictable multi user analysis**
- **Overall demand typically fills the capacity of the resources**
- **Real time monitor systems for networks, storage, computing resources,… : E2E monitoring**

| Application Interfaces |
| Request Planning |
| Network Planning |
| Network Resources |

**Monitor**

Support data transfers ranging from predictable movement of large scale (simulated and real) data, to highly dynamic analysis tasks initiated by rapidly changing teams of scientists