# (Possible) HEP Use Case for NDN

Phil DeMar; Wenji Wu
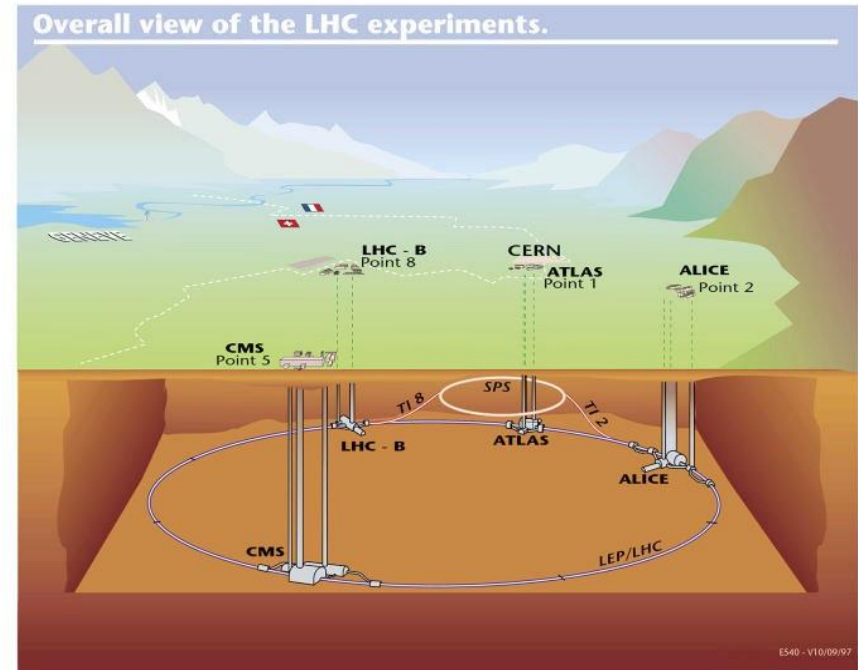NDNComm (UCLA)
Sept. 28, 2015

# Outline

➢ LHC  Experiments

➢ LHC Computing Models

➢ CMS Data Federation & AAA

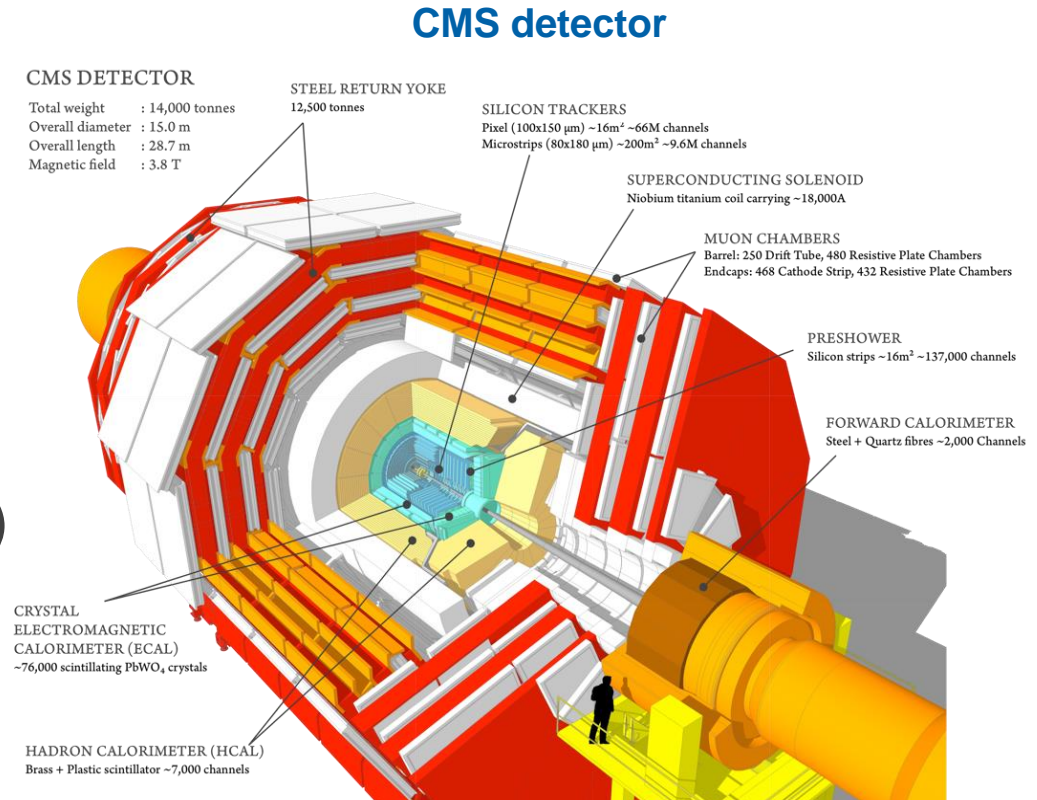➢ Evolving Computing Models & NDN

➢ Summary

**춘춘 Fermilab**

# Large Hadron Collider (LHC) 101

➤ Circumference: ~ 17 Miles

➤ 2 proton beams circulating at 99.9999991% speed of light:

➤ Beams cross and are brought to collision at 4 points:

➤ Experiments built at those points
  – ATLAS
  – CMS
  – ALICE
  – LHCb



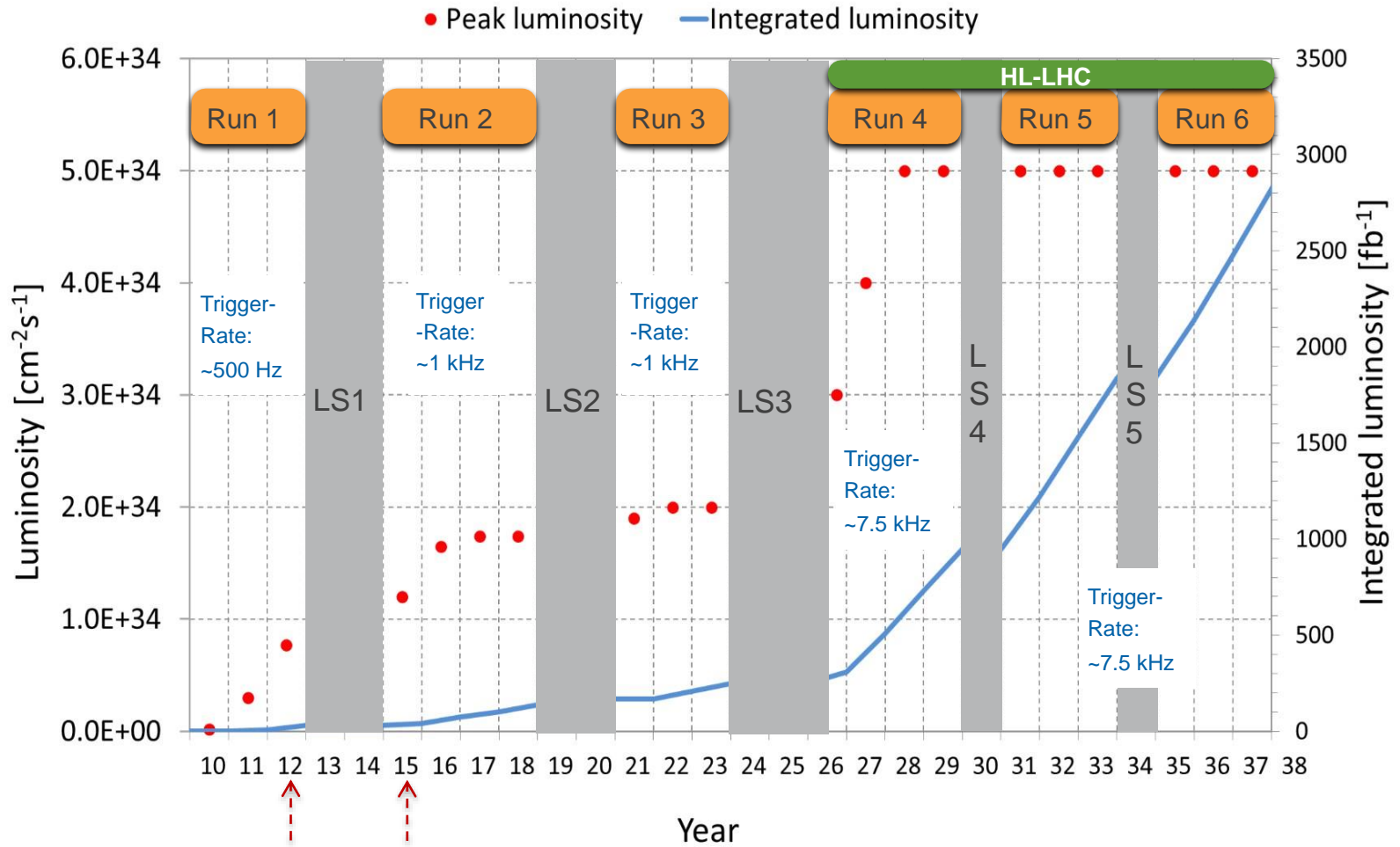Overall view of the LHC experiments.

🔁 **Fermilab**

# Compact Muon Solenoid (CMS) Experiment

➢ Detector built around collision point

➢ Records flight path and energy of all particles produced in a collision

➢ 100 Million individual measurements (channels)

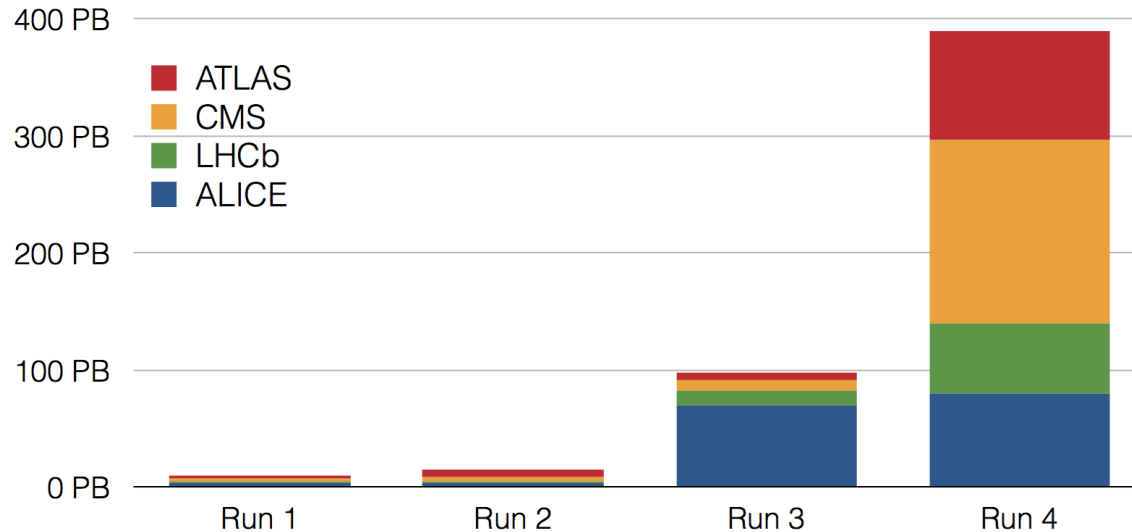➢ All measurements of a collision together are called: **event**

**CMS detector**

CMS DETECTOR
Total weight      : 14,000 tonnes
Overall diameter : 15.0 m
Overall length    : 28.7 m
Magnetic field    : 3.8 T

STEEL RETURN YOKE
12,500 tonnes

SILICON TRACKERS
Pixel (100x150 μm) ~16m² ~66M channels
Microstrips (80x180 μm) ~200m² ~9.6M channels

SUPERCONDUCTING SOLENOID
Niobium titanium coil carrying ~18,000A

MUON CHAMBERS
Barrel: 250 Drift Tube, 480 Resistive Plate Chambers
Endcaps: 468 Cathode Strip, 432 Resistive Plate Chambers

PRESHOWER
Silicon strips ~16m² ~137,000 channels

FORWARD CALORIMETER
Steel + Quartz fibres ~2,000 Channels

CRYSTAL ELECTROMAGNETIC CALORIMETER (ECAL)
~76,000 scintillating PbWO₄ crystals

HADRON CALORIMETER (HCAL)
Brass + Plastic scintillator ~7,000 channels

🟰 **Fermilab**

# LHC schedule

M. Girone (CERN)

🔷 Fermilab

# Projected LHC data volumes

RAW

M. Girone (CERN)



Legend:
- ATLAS
- CMS
- LHCb
- ALICE

Y-axis: 0 PB, 100 PB, 200 PB, 300 PB, 400 PB
X-axis: Run 1, Run 2, Run 3, Run 4

**Exabyte era...**

➤ Raw data = generated by detector(s)

➤ Derived data = reconstructed data, simulation data, summary data sets, etc…)

  – (derived data volumes) ~= (raw data volumes) x 8

🔷 Fermilab

# CMS Collaboration

➢ 186 institutions (globally distributed)
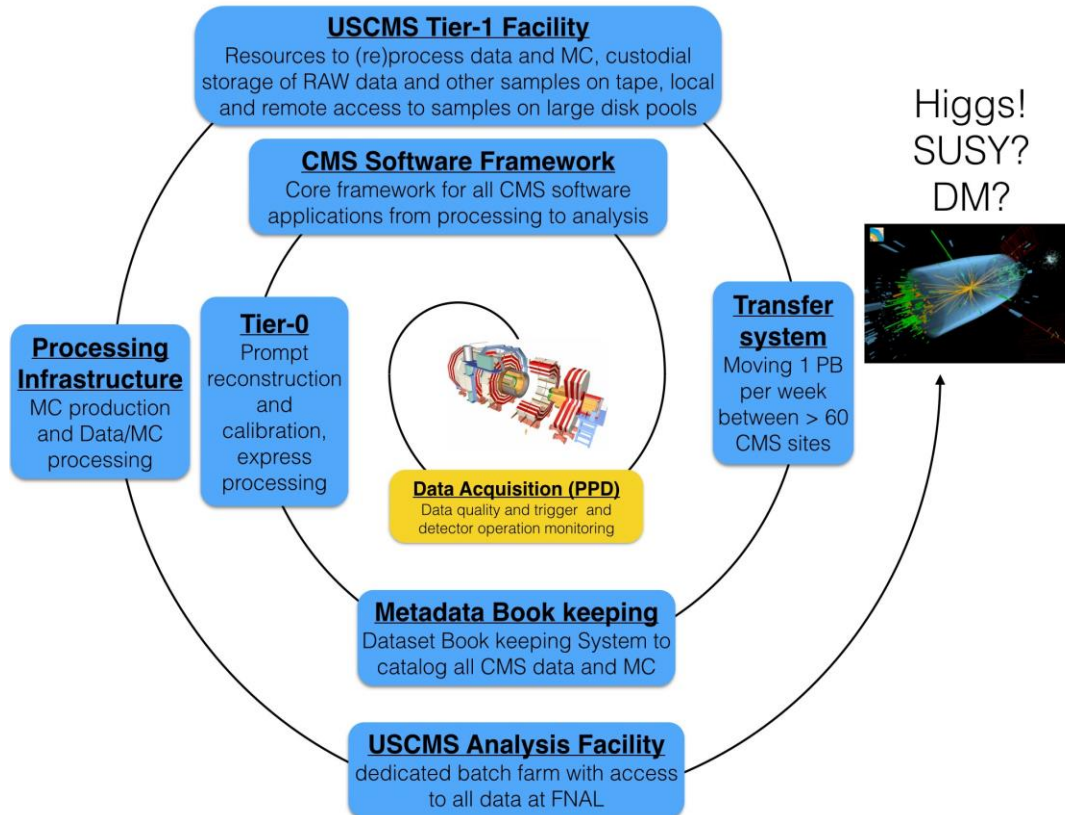  – High b/w R&E networks support experiment data movement

🔰 Fermilab

# LHC Computing Models

**춘춘 Fermilab**

# Computing Lifecycle: CMS



**USCMS Tier-1 Facility**
Resources to (re)process data and MC, custodial storage of RAW data and other samples on tape, local and remote access to samples on large disk pools

**CMS Software Framework**
Core framework for all CMS software applications from processing to analysis

**Processing Infrastructure**
MC production and Data/MC processing

**Tier-0**
Prompt reconstruction and calibration, express processing

**Data Acquisition (PPD)**
Data quality and trigger and detector operation monitoring

**Transfer system**
Moving 1 PB per week between > 60 CMS sites

Higgs! SUSY? DM?

**Metadata Book keeping**
Dataset Book keeping System to catalog all CMS data and MC

**USCMS Analysis Facility**
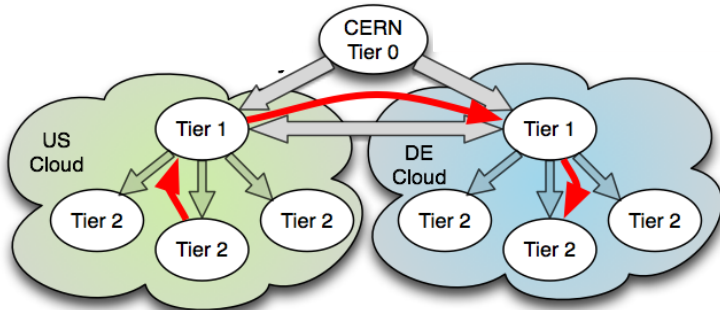dedicated batch farm with access to all data at FNAL

O. Gutsche (FNAL)

➢ Tier structure for computing (MONARC):

  ➢ Tier 0 = CERN

  ➢ Tier 1 = National data centers for event reconstruction & archiving

  ➢ Tier 2 = Computing facilities for Monte Carlo production & event analysis

  ➢ Tier 3 = Collaboration sites
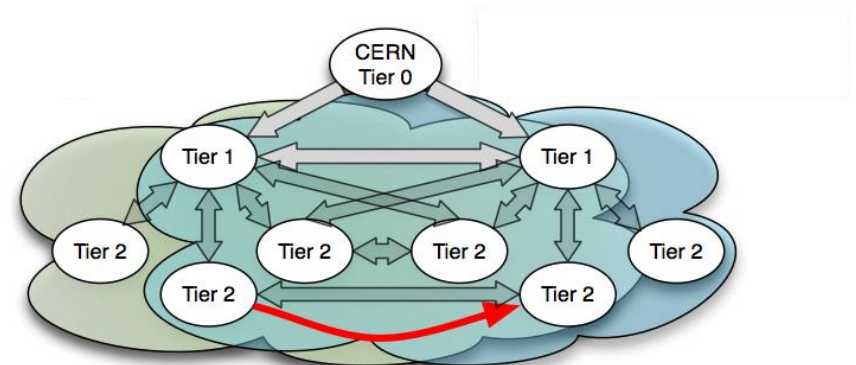
  ➢ Tier 4 = Physicist desktops

Fermilab

# CMS Computing GRID infrastructure

- CERN (T0) at the center

- 7 Tier-1 centers:
  - Connected to T0 by a "dedicated" network

- 54 Tier-2 facilities
  - Connected to T1s by R&E networks

- ~120,000 cores

- ~75PB disk

- ~100PB tape



**54 T2 sites**

France
UK
Germany
Italy
Spain
Russia
T1
USA (FNAL)
T0 @ CERN

General Purpose Scientific Networks between all T1 and T2 sites
**GPN**

Dedicated Optical Private Network between T0 and all T1 sites
**LHCOPN**

O. Gutsche (FNAL)

**Fermilab**

# Tier Model for Data Movement Abandoned



- ➢ MONARC hierarchical model
- ➢ Based on expectation of low b/w & modest storage at T2s
- ➢ CMS abandoned MONARC before the LHC even started…
  - ➢ ATLAS followed suit during Run I

- ➢ Any CMS T1/T2 site could be used as a data source
- ➢ Encouraged more flexible data placement & replication
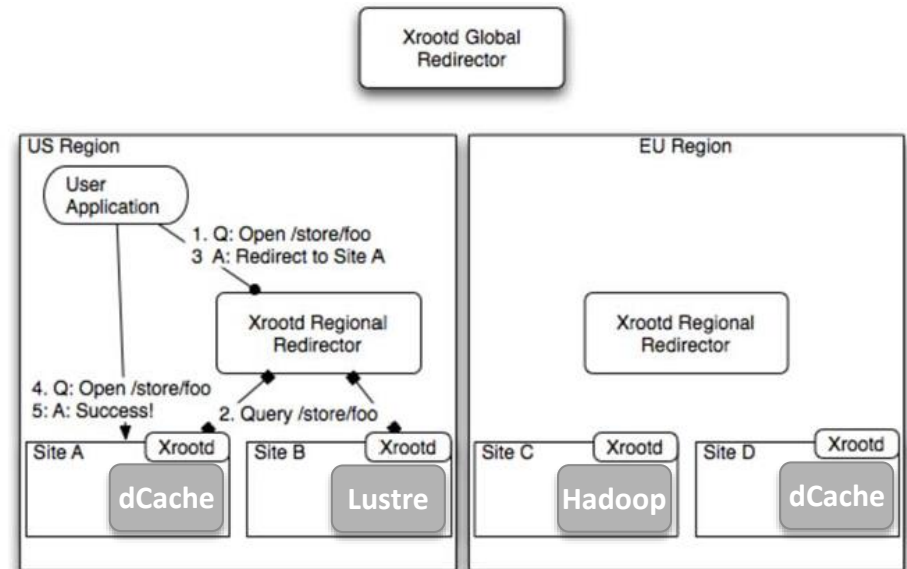- ➢ Enabled more efficient utilization of available resources



T. Wenaus (BNL)

🟦 **Fermilab**

# CMS Data Federation & AAA

# Data Federation - XrootD

➤ LHC experiments have implemented federated data storage, made possible by:
  – High bandwidth WAN connectivity across all tiers
  – Global data namespace(s)

➤ Based on XrootD:
  – "Hides" local file storage systems
  – Hierarchical, w/ regional, global, & local redirectors
  – Maintains catalog of known file locations
    • Negative cache as well
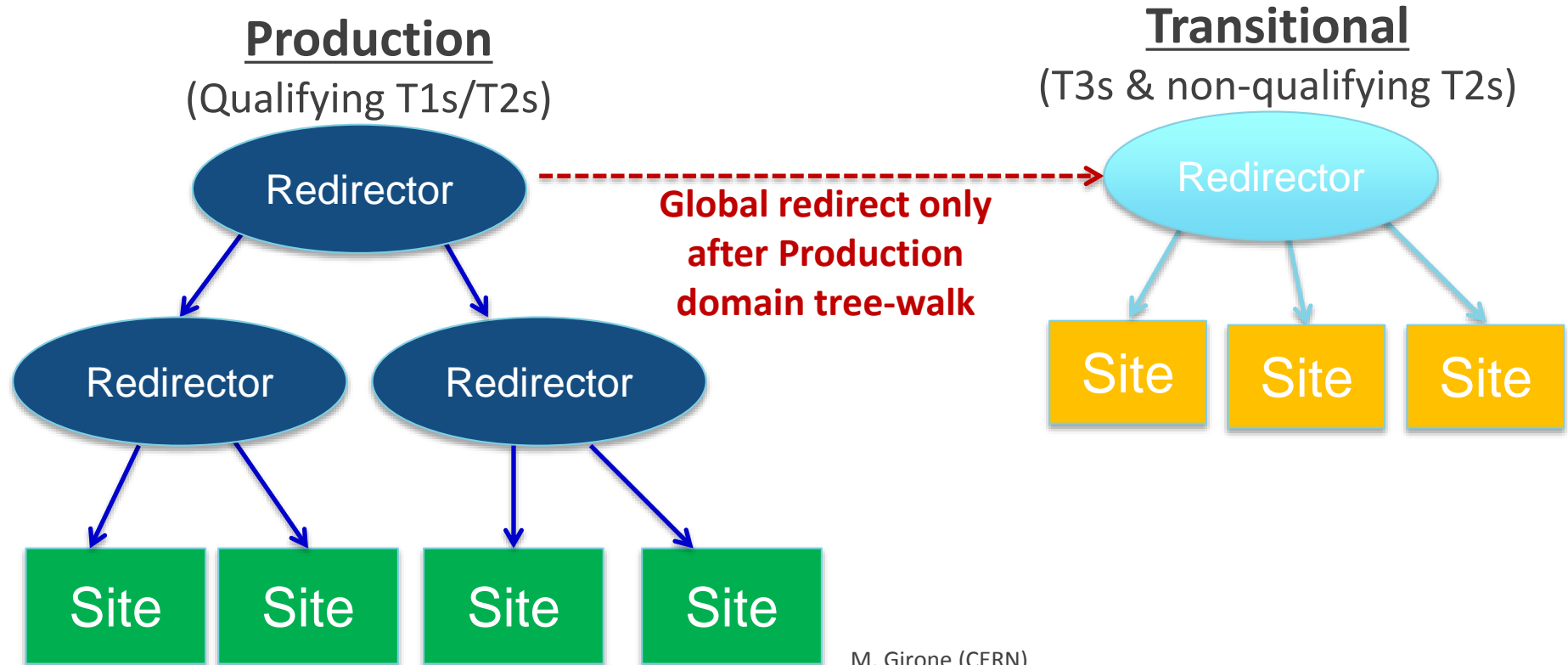  – Tree-walk redirects to locate file

# Any Data, Any Time, Anywhere (AAA)

➢ AAA is CMS's implementation of federated storage:
  – Based on XrootD
  – Finds data based on logical file name
  – Transfers data to application

➢ High-level philosophy:  remote storage ~= local storage:
  – In practice:  CPU efficiency slightly lower w/ remote data

➢ Principally driven by (macro) economics:
  – Maximizes efficiency of collaboration computing resources
  – Fallback data access & overflow job redistribution capabilities

➢ A few numbers:
  – Nearly all (95%+) CMS data available via AAA
  – Projection is 20%+ of CMS Run II data access through AAA
    • Local storage access is not through AAA…

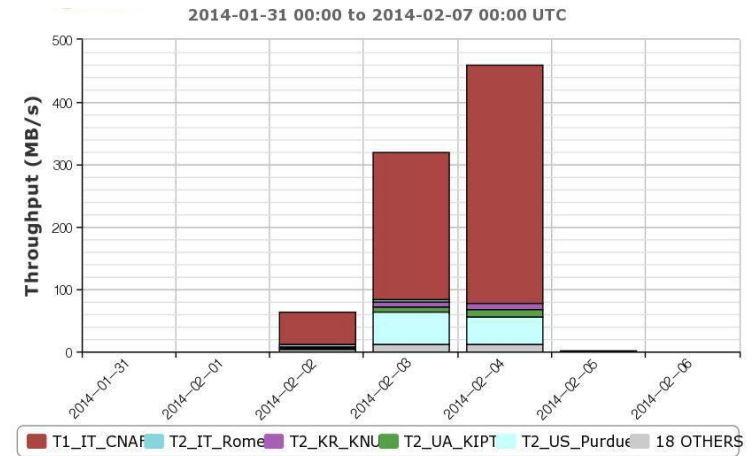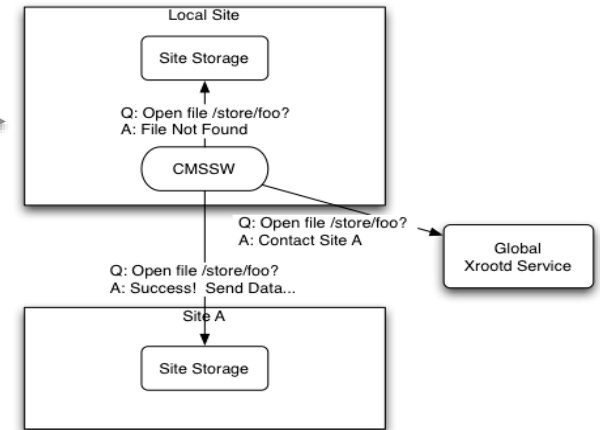**춘 Fermilab**

# AAA's Two-domain Federation

- ➢ Production domain for AAA performance-certified sites
  - – Transition domain for sites not meeting performance standards
  - – All CMS T1s and most T2s are now Production-certified

**Production**
(Qualifying T1s/T2s)

**Transitional**
(T3s & non-qualifying T2s)

Redirector - - - - - - - - - - - - - - - - - - - - → Redirector

**Global redirect only
after Production
domain tree-walk**

Redirector      Redirector

Site   Site   Site   Site

Site   Site   Site

M. Girone (CERN)

**Fermilab**

# AAA Fallback Mode

➤ Job unable to access local data:
  – AAA fallback capability locates remote copy of data
  – Job is able to complete…

➤ Useful in redirecting jobs to other sites in overflow situations

➤ Real life example:
  – DB error results in "missing" local data at FNAL
  – AAA failover locates replica at CNAF (Italy)
  – Jobs run for 2 days using CNAF data, without anyone noticing…

# Evolving Computing Models & NDN

**莽 Fermilab**

# Additional Trends in CMS Computing Model…

➢ Dynamic data placement (ALICE/ATLAS):
  – Distributing/redistributing (abbreviated) data sets by popularity
  – Subset of larger trend for dynamic data management in general

➢ Cloud & High Performance Computing (HPC) cycles:
  – Amazon Web Service spot CPU cycles already highly economic
  – Next gen. super computers will have massive computing power



AWS Planned 100G to PNWG

PACIFIC GIGAPOP

amazon web services

Direct Connect ESnet Pilot 2x10G

Direct Connect ESnet Pilot 1x10G

100G R&E Exchange

ESnet

M. Ernst (BNL)

| System attributes | NERSC Now | OLCF Now | ALCF Now | NERSC Upgrade | OLCF Upgrade | ALCF Upgrades | |
|---|---|---|---|---|---|---|---|
| Name Planned Installation | Edison | TITAN | MIRA | Cori 2016 | Summit 2017-2018 | Theta 2016 | Aurora 2018-2019 |
| System peak (PF) | 2.6 | 27 | 10 | > 30 | 150 | >8.5 | 180 |
| Peak Power (MW) | 2 | 9 | 4.8 | < 3.7 | 10 | 1.7 | 13 |
| Total system memory | 357 TB | 710TB | 768TB | ~1 PB DDR4 + High Bandwidth Memory (HBM) +1.5PB persistent memory | > 1.74 PB DDR4 + HBM + 2.8 PB persistent memory | >480 TB DDR4 + High Bandwidth Memory (HBM) | > 7 PB High Bandwidth On-Package Memory Local Memory and Persistent Memory |
| Node performance (TF) | 0.460 | 1.452 | 0.204 | > 3 | > 40 | > 3 | > 17 times Mira |
| Node processors | Intel Ivy Bridge | AMD Opteron Nvidia Kepler | 64-bit PowerPC A2 | Intel Knights Landing many core CPUs Intel Haswell CPU in data partition | Multiple IBM Power9 CPUs & multiple Nvidia Voltas GPUS | Intel Knights Landing Xeon Phi many core CPUs | Knights Hill Xeon Phi many core CPUs |
| System size (nodes) | 5,600 nodes | 18,688 nodes | 49,152 | 9,300 nodes 1,900 nodes in data partition | ~3,500 nodes | >2,500 nodes | >50,000 nodes |
| System Interconnect | Aries | Gemini | 5D Torus | Aries | Dual Rail EDR-IB | Aries | 2nd Generation Intel Omni-Path Architecture |
| File System | 7.6 PB 168 GB/s, Lustre | 32 PB 1 TB/s, Lustre | 26 PB 300 GB/s GPFS™ | 28 PB 744 GB/s Lustre | 120 PB 1 TB/s GPFS™ | 10PB, 210 GB/s Lustre initial | 150 PB 1 TB/s Lustre |

**Fermilab**

# CMS Computing (today…) vs NDN

*Warning!!! My interpretation only!*
*Subject to large error bars on both ends…*

| | CMS (today) | NDN |
|---|---|---|
| Namespace | Global logical file names | Hierarchical data name space |
| Content-based data retrieval | Middleware service | Basic network service |
| Routing optimization | Some architectural & middleware optimizations | Basic network service |
| Caching optimizations | Middleware optimizations | Basic network service (?) (not clear how this would work with LHC scale data volumes) |
| Scalable Repository | Open Science Grid Stashcache (middleware) [?] | Repo-Se (?) |

**Fermilab**

# But Don't Confuse Us with NetFlix…

➢ NetFlix delivers streaming video content to ~20M users
  – Regarded as largest content provider for internet traffic

➢ CMS much smaller user base & generates only a fraction of NetFlix's traffic
  – But CMS aggregate amount of data is 1000X NetFlix
  – NetFlix deals with much lower amount of data, which is much easier to efficiently replicate or cache

|  | NetFlix | CMS |
|---|---|---|
| **Users** | 20M | 100K |
| **Total Data** | 20TB | 20PB |

O. Gutsche (FNAL)

**Fermilab**

# NDN Activities in High Energy Physics (HEP)...

➢ Climate Data Sciences NDN test bed (C. Papadopoulos, etc.) has ties with HEP community

  – Caltech Network Research group (H. Newman) is involved

➢ Imperial College London (D. Rand, etc.) evaluating NDN in a local test bed:

  – Application-level (ROOT)
  – Repository-level

➢ Caltech & FNAL funded to create small NDN test bed for CMS app evaluations

**⚛ Fermilab**

# Summary…

➢ **LHC experiments heading toward exascale data volumes:**
  – Terabit networks will be needed to handle that data

➢ **LHC computing models are becoming increasingly distributed in nature:**
  – Both data storage & CPU
  – This creates greater demands on network services beyond b/w

➢ **LHC computing is already implementing content-based data services at the middleware level**

➢ **There seems to be a natural fit for NDN with LHC computing:**
  – Performance optimizations within the exascale data / terabit network environment will be key

**🎗 Fermilab**

# Questions?

**춘 Fermilab**