# LHC Data Analytics

**Artur Barczyk**
**California Institute of Technology**
**BIG DATA Workshop**
**NIST, Gaithersburg, June 14, 2012**

# DISCLAIMER

**This presentation intends to give an overview of LHC data processing, based on samples and general notions. It is as such intrinsically incomplete, as it's impossible to cover this vast area in a short time without idiosyncratic bias.**

**References to detailed information were intended, and where missing can be obtained from the presenter.**

# OUTLINE

**LHC and its data**

**LHC data processing and analysis chain**

**Data sizes, Data rates**

**Computing Infrastructures**

# THE LHC AND ITS COMPUTING INFRASTRUCTURE

# Large Hadron Collider
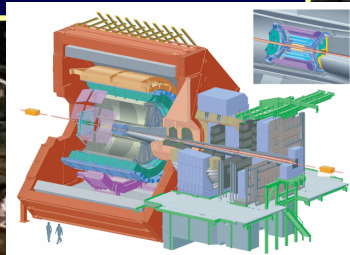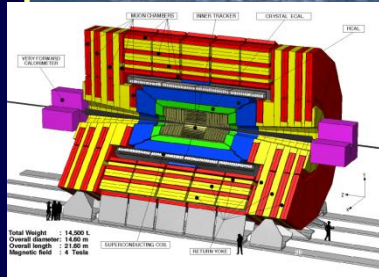# 7 TeV in 2010-11, 8 TeV in 2012

Large Hadron Collider
27 km circumference

Lake Geneva

* *The LHC is a Discovery Machine: High energy and "Luminosity"*
* **The first accelerator to probe deep into the Multi-TeV scale**
* **Many reasons to expect new TeV-scale physics**

**Higgs, SUSY, Substructures, CP Violation, QG Plasma, *… Gravitons, Extra Dimensions, Low Mass Strings, … the Unexpected***
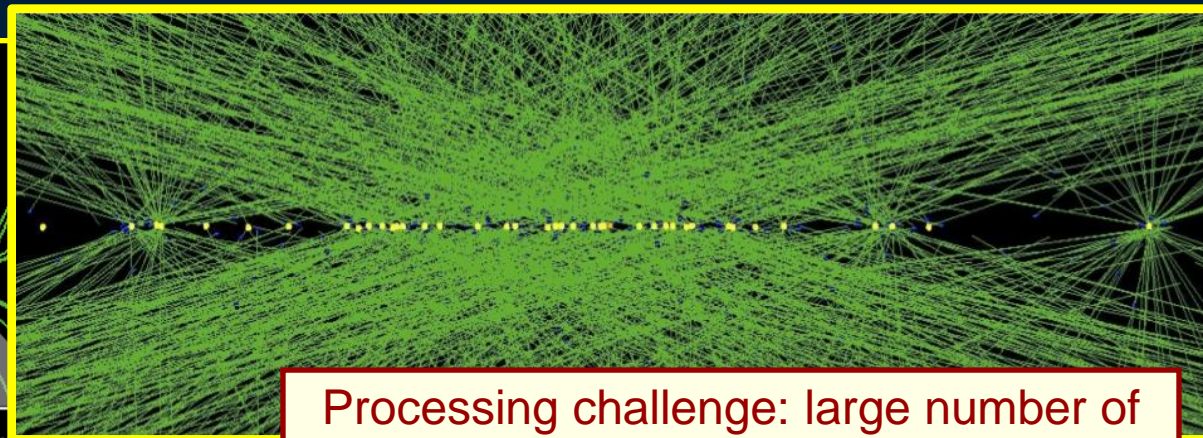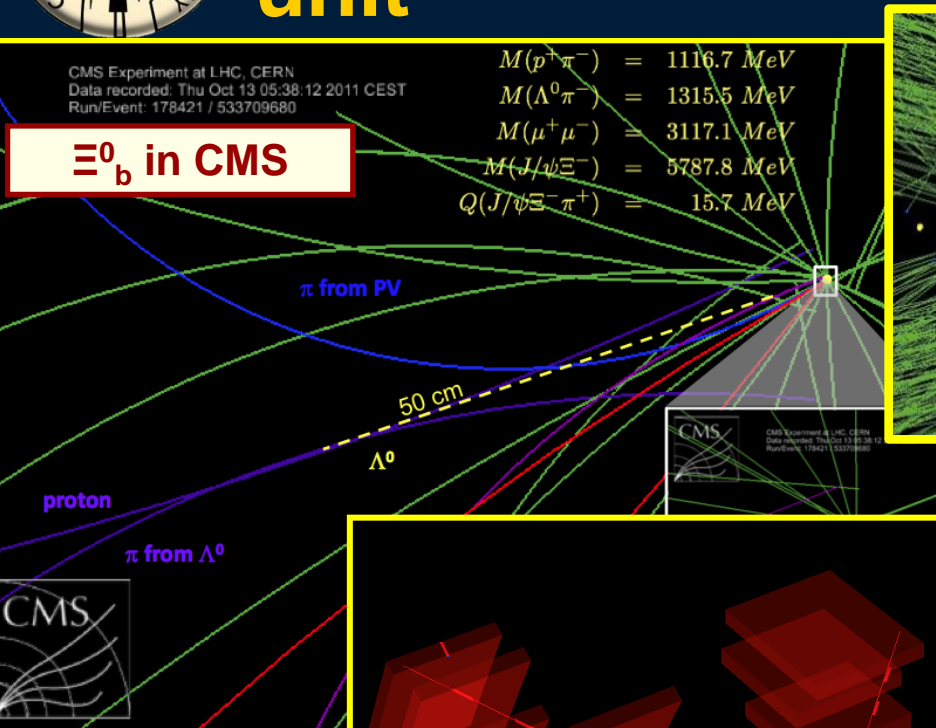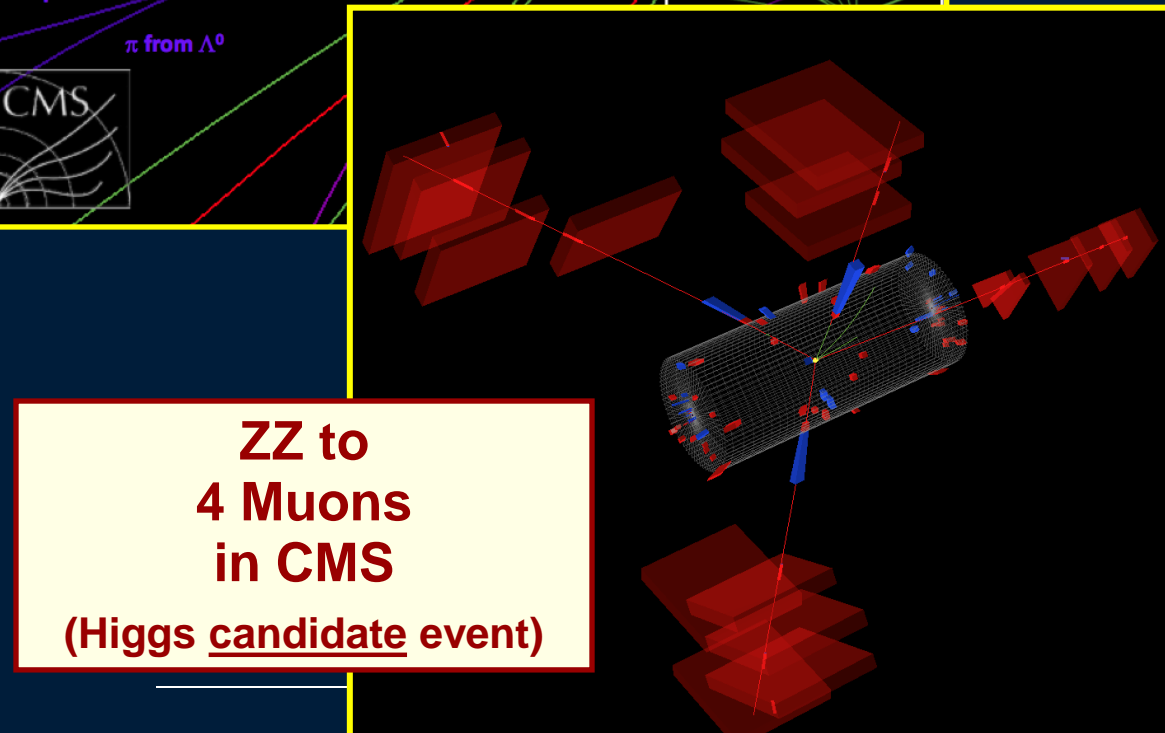
# HEP Events: the basic processing unit

$\Xi^0_b$ in CMS

$$M(p^+\pi^-) = 1116.7\ MeV$$
$$M(\Lambda^0\pi^-) = 1315.5\ MeV$$
$$M(\mu^+\mu^-) = 3117.1\ MeV$$
$$M(J/\psi\Xi^-) = 5787.8\ MeV$$
$$Q(J/\psi\Xi^-\pi^+) = 15.7\ MeV$$

CMS Experiment at LHC, CERN
Data recorded: Thu Oct 13 05:38:12 2011 CEST
Run/Event: 178421 / 533709680

$\pi$ from PV

50 cm

$\Lambda^0$

proton

$\pi$ from $\Lambda^0$

Processing challenge: large number of primary vertices in single bunch crossing

**ZZ to 4 Muons in CMS**

(Higgs candidate event)

**Pb-Pb event in Alice**

ALICE

6

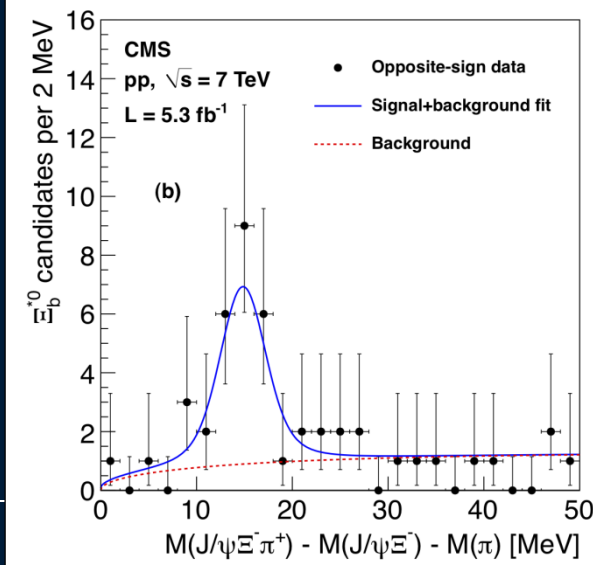# How do we search for New Physics? (Simplified)

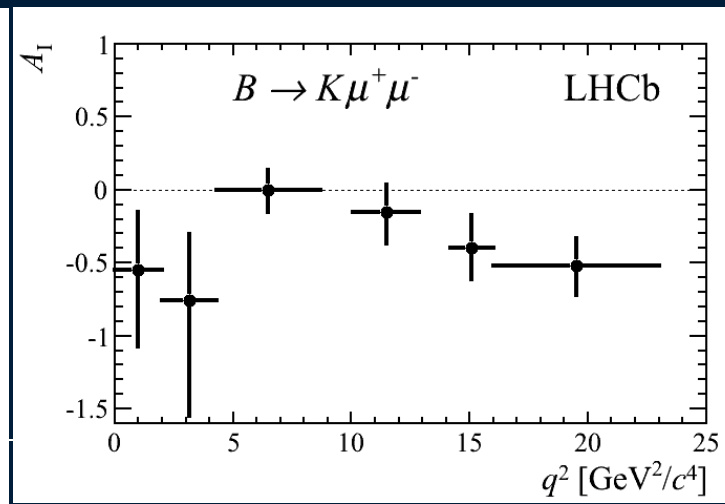- **Statistical Analysis**
  - Very low S/B ratio ($10^{-9} - 10^{-13}$)
- **Looking for the predicted/expected**
  - Select events with expected signatures; Evaluate signal over background estimate (Monte Carlo simulation)
- **Looking for the unexpected**
  - Use known processes, do precision measurement, look out for deviations



CMS Preliminary, $\sqrt{s}$ = 7 TeV
Combined, $L_{int}$ = 1.1-1.7 fb$^{-1}$

Observed
Expected ± 1σ
Expected ± 2σ

CMS excluded: 145-216, 226-288, 310-400

95% CL limit on $\sigma/\sigma_{SM}$
Higgs boson mass (GeV/c$^2$)



Theory  Data
$B \rightarrow K^* \mu^+ \mu^-$    LHCb
$A_I$
$q^2$ [GeV$^2$/c$^4$]



$B \rightarrow K \mu^+ \mu^-$    LHCb
$A_I$
$q^2$ [GeV$^2$/c$^4$]



CMS
pp, $\sqrt{s}$ = 7 TeV
L = 5.3 fb$^{-1}$

Opposite-sign data
Signal+background fit
Background

(b)

$\Xi_b^{*0}$ candidates per 2 MeV
M(J/$\psi \Xi^- \pi^+$) - M(J/$\psi \Xi^-$) - M($\pi$) [MeV]

# The LHC Computing Challenge

Signal/Noise: $10^{-13}$ ($10^{-9}$ offline)

Data volume

- High rate * large number of channels * 4 experiments

→ **15 PetaBytes of new data each year**  → 23 PB in 2011

Compute power

- Event complexity * Nb. events * thousands users

→**200 k CPUs**  → 250 k CPU

→**45 PB of disk storage**  → 150 PB

Worldwide analysis & funding

- Computing funding locally in major regions & countries
- Efficient analysis

→ **GRID technology**   Ian Bird, CERN



H → γγ

Higgs signal



High Level-1 Trigger (1 MHz)

High No. Channels High Bandwidth ( 1000 Gbit/s)

High Data Archives (PetaBytes)

LHCb

ATLAS CMS

KTev

HERA-B

KLOE

CDF II D0 II

BaBar

CDF, D0

H1 ZEUS

UA1

NA49

ALICE

LEP

Level-1 Rate (Hz)

Event Size (byte)

WLCG
Worldwide LHC Computing Grid

8

# LHC Computing Infrastructure



LHC Experiment

Online System — 100-200 MBytes/s

**Tier 0** — CERN Computer Center > 20 TIPS

2.5 - 10 Gbits/s

**Tier 1** — Japan, UK, France, USA

2.5 Gbits/s

**Tier 2** — Tier2 Center, 2 Center, 2 Center, 2 Center

~0.6 Gbits/s

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

Physics cache

1 Gbits/s

**Tier 4** — PCs, other portals

**WLCG in brief:**
- 1 Tier-0 (CERN)
- 11 Tiers-1s; more under discussion
- 68 Tier-2 Federations; > 140 sites

**Plus O(300) Tier-3s worldwide**

2000 km
1000 mi

Provided by GStat 2.0

Scale = 1 : 111M

Permalink

148.29688, 48.86719

9

# COMPUTING MODELS

**The transition**

# Computing Site Roles (so far)



**Prompt calibration and alignment**
**Reconstruction**
**Store complete set of RAW data**

**Data Reprocessing**
**Archive RAW and Reconstructed data**

**Monte Carlo Production**
**Physics Analysis**
**Store Analysis Objects**

**Physics Analysis, Interactive Studies**

Tier 0 (CERN)

Tier 1

Tier 1

Tier 2

Tier 2

Tier 3

US LHCNet

## PanDA System Overview

Production managers

PanDA server

Data Management System (DQ2)

production job

https

define

submitter (bamboo)

https

ATLAS Dashboard

task/job repository (Production DB)

analysis job

https

submit

End-user

EGEE/EGI

Wor...

**~35'000 Atlas jobs running concurrently at any moment**

### Completed jobs
137 Days from Week 01 of 2012 to Week 20 of 2012

2012 p-p run



Legend: US, UK, DE, FR, IT, CA, NL, ES, ND, TW, other

Maximum: 935,926 , Minimum: 68,559 , Average: 348,999 , Current: 278,368

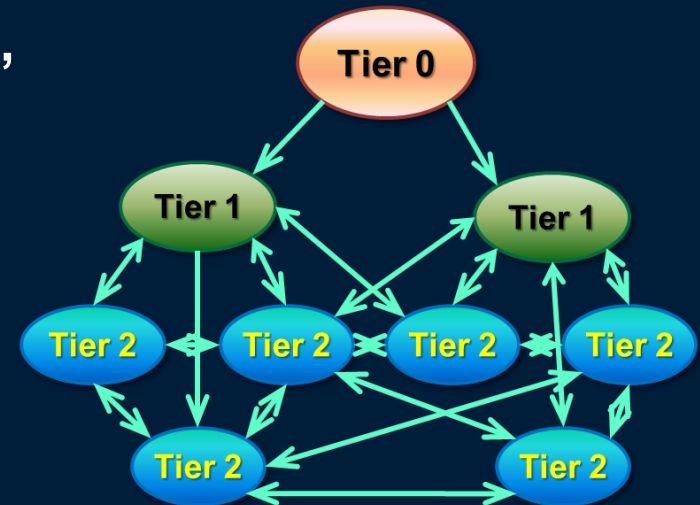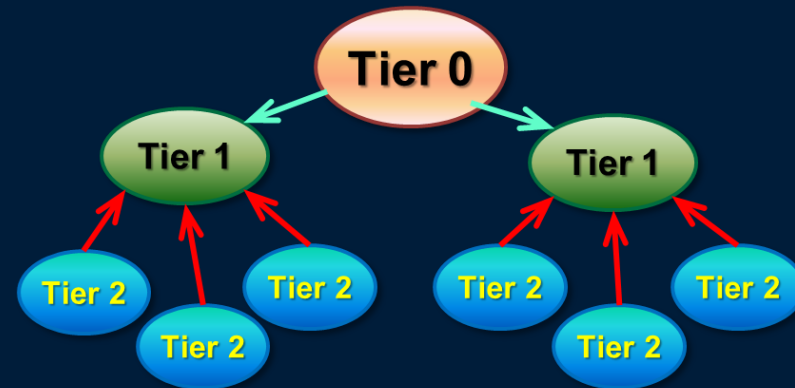# Parallel processing

- **At event granularity:**
  - LHC experiments typically use 1 core per job with sequential event processing
  - The Grid is perfectly matching this process workflow
    - Jobs dispatched to the Grid, running over event files (present at the executing site)
  - No inter-node or inter-site process synchronization (distributed, but independent computation)
- **New approaches being investigated:**
  - Make more efficient use of multi-core, multi-cpu architectures
  - Possibly make use of massively parallel hardware (GPUs)
    - E.g. in event reconstruction
  - Granularity remains at single event level
    - No clear advantage of processing single event on multiple nodes

# Computing Models Evolution

- **Moving away from the strict MONARC model**
- **Introduced gradually since 2010**
- **3 recurring themes:**
  - **Flat(ter) hierarchy: Any site can use any other site as source of data**
  - **Dynamic data caching: Analysis sites will pull datasets from other sites "on demand", including from Tier2s in other regions**
    - **Possibly in combination with strategic pre-placement of data sets**
  - **Remote data access: jobs executing locally, using data cached at a remote site in quasi-real time**
    - **Possibly in combination with local caching**
- **Variations by experiment**
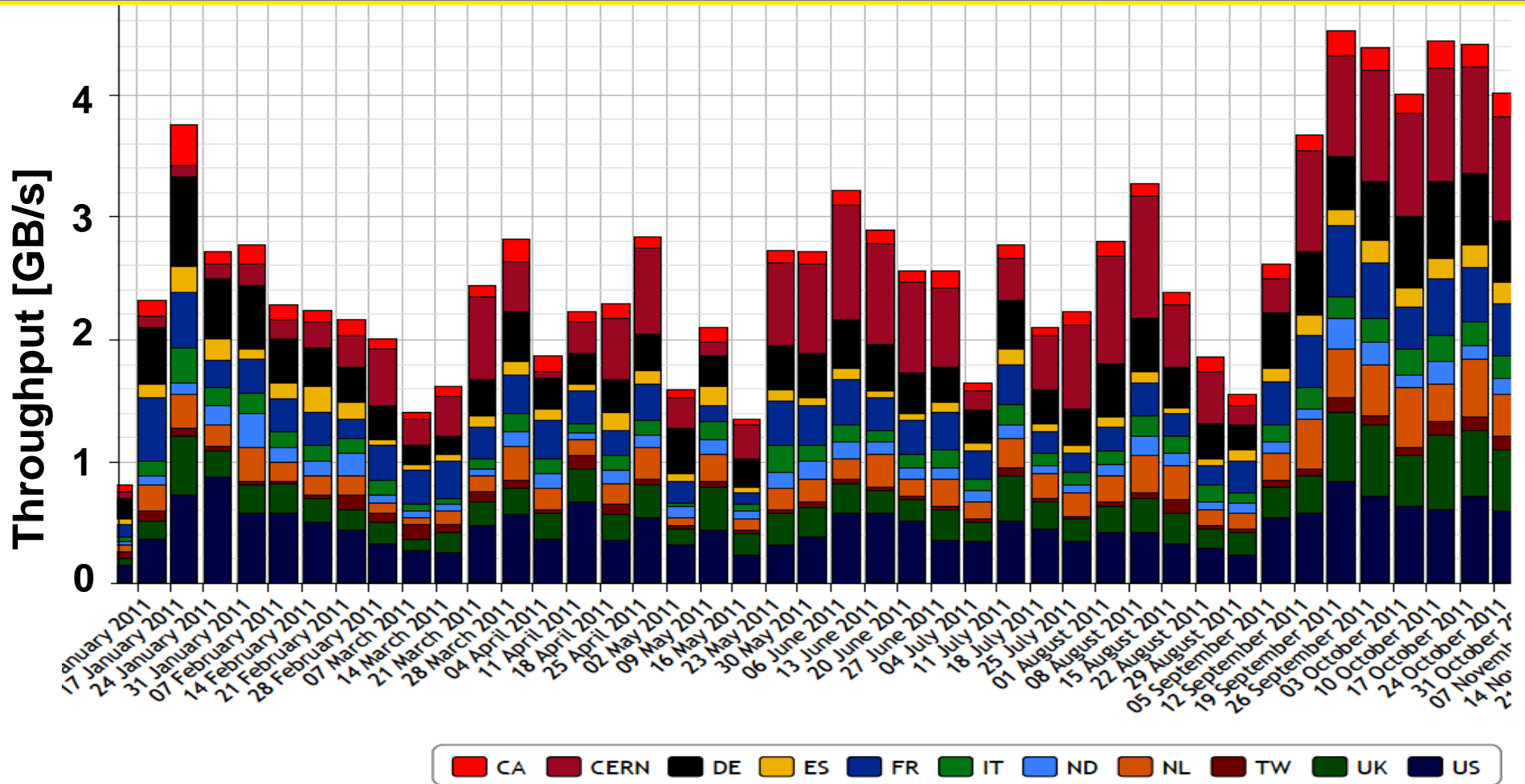- **Increased reliance on network performance!**

# ATLAS Data Flow by Region: Jan. – Nov. 2011

**~2.8 Gbytes/sec Average, 4.5 GBytes/sec Peak**

***~100 Petabytes Transferred During 2011***

# CMS data movement overview



CMS PhEDEx - Transfer Rate
52 Weeks from Week 24 of 2011 to Week 24 of 2012

**CMS: all sites
Weekly data rates
0.5 - 2.5 GB/s**

CMS PhEDEx - Transfer Rate
52 Weeks from Week 24 of 2011 to Week 24 of 2012

**CMS: T1-T2
Weekly data rates 0.2 – 0.9 GB/s**

CMS PhEDEx - Transfer Rate
52 Weeks from Week 24 of 2011 to Week 24 of 2012

**CMS: T2-T2
Weekly rates: 0.1-0.6 GB/s**

# NETWORK INFRASTRUCTURES FOR LHC DATA PROCESSING

**CHEP 2012: Network was 2nd most discussed topic in the Computer Facilities, Production Grids and Networking track**



Andreas Heiss, CHEP 2012,
Track Summary: Computer Facilities, Production Grids and Networking

# R&E Networking Landscape

- **R&E networks as substrate for LHC data movement**
  - National (e.g. ESnet, SURFNet, RNP, …)
  - Regional (e.g. MiLR, CENIC, NORDUNet, …)
  - International (e.g. NORDUnet, GEANT, ACE, …)
  - Open Exchanges (e.g. Starlight, MANLAN, NetherLight, CERNLight, …)
  - Dedicated, mission–oriented network (US LHCNet)

For more complete listing see e.g. ICFA-SCIC report at http://icfa-scic.web.cern.ch/



http://es.net/



http://glif.is

# R&E Networks are vital for modern science



ESnet Accepted Traffic: Jan 1990 - Apr 2012 (Log Scale)

**ESnet Traffic increases 10X Every 47 Months on Avg. (since 1991)**

**ESnet's New 100G Backbone, ETTC: Nov 2012**

Petabytes / month

Month

# US LHCNet: Dedicated Network Infrastructure

- **Mission-oriented network**
- **Dedicated to LHC data movement between Europe and US**
- **High availability goal (99.95+), despite challenges of submarine environment**

# 100 & 40 Gbps System Integration

US LHCNet

**SuperComputing 2012 Caltech demo:**
**100Gbps single wave; 40 GE servers**
**Seattle - University of Victoria (212km)**

Network Traffic - Disk to Disk

60 Gbps

40 Gbps

Server (in black) moves from
2 x 10G to 3 x 10G receive

Second Gen 3, 40 G server
added.

Network Traffic

100

70 Gbps
60 Gbps
50 Gbps
40 Gbps
30 Gbps
20 Gbps

Traffic IN

2 x 30 G
**Gen 3** PCI 40 G
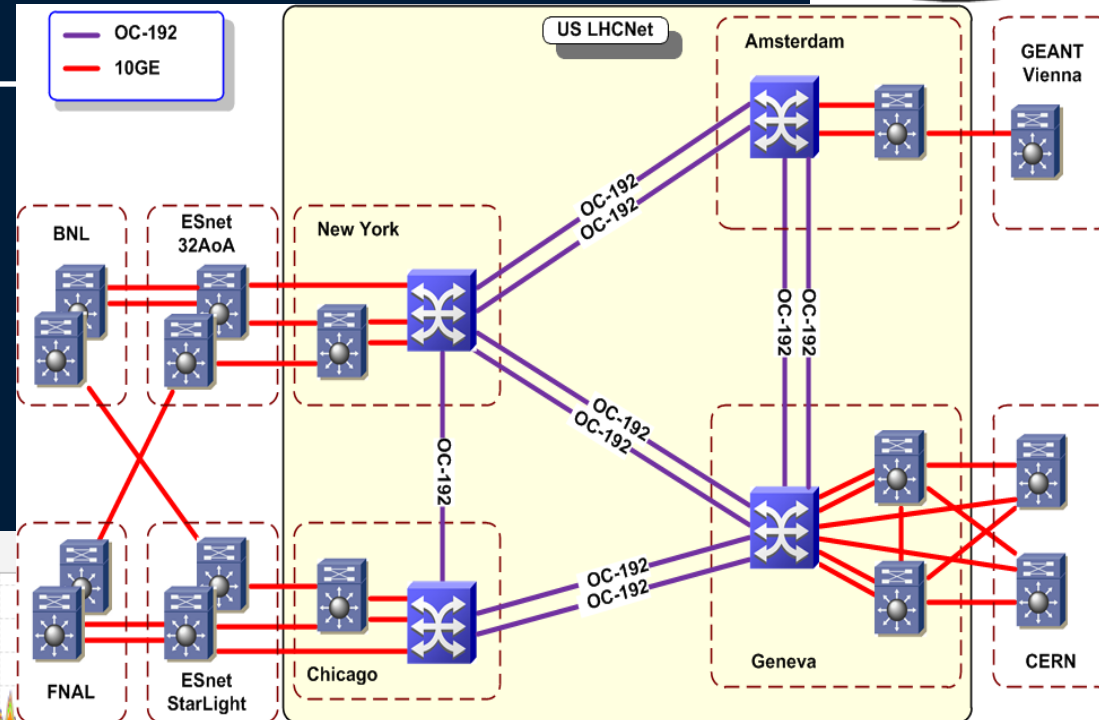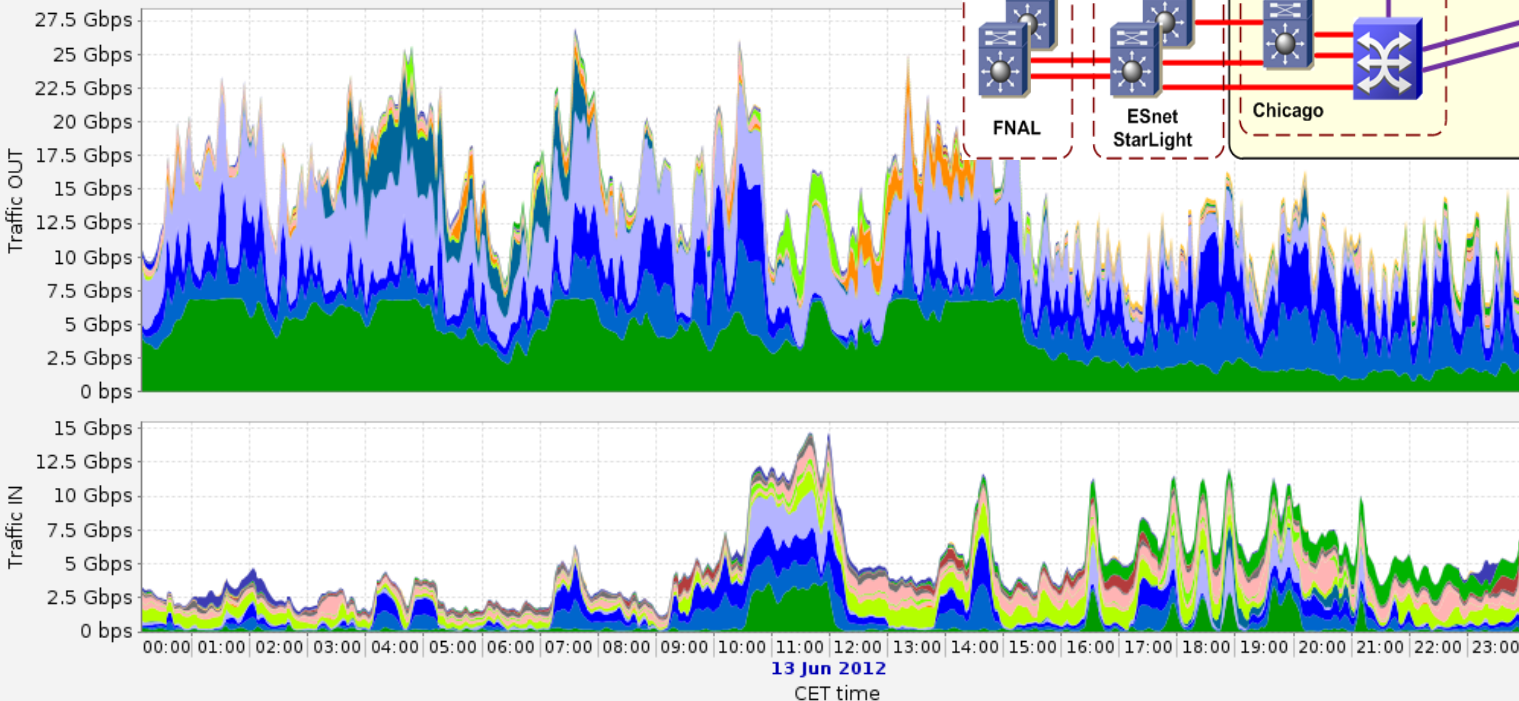Network cards

**SSD-to-SSD: up to ~60Gbps**

0

20 Gbps
30 Gbps
40 Gbps
50 Gbps
60 Gbps
70 Gbps

Traffic OUT

3:00  00:00  01:00  02:00  03:00  04:00  05:00  06:00  07:00  08:00
17 Nov 2011
PST time

Network Traffic

**40 Gbps**

30 Gbps
25 Gbps
20 Gbps
15 Gbps

IN

90

14:30  14:33  14:36  14:39  14:42  14:4

■ sc-fdt-dynes ■ sc1-g2-1 ■ sc10 ■ sc11 ■ sc12
■ sc29-g2-5 ■ sc3-g2-3 ■ sc30-g2-6 ■ sc31-g2-7

**Geneva-Amsterdam**
**(Caltech, CERN, SURFnet)**
**Single Server-Pair, 40GE NIC, 1650 km**
**Mem-mem: 40Gbps**
**SSD-SSD: ~18Gbps**

FDT: http://monalisa.cern.ch/FDT/

http://supercomputing.caltech.edu

# The LHC Optical Private Network
## Serving LHC Tier0 and Tier1 sites

- **Dedicated network resources for Tier0 and Tier1 data movement**
- **Layer 2 overlay on R&E infrastructure**
- **130 Gbps total Tier0-Tier1 capacity (today)**
- **Simple architecture**
  - Point-to-point Layer 2 circuits
  - Flexible and scalable topology
- **Grew organically**
  - From star to partial mesh
- **Open to technology choices**
    - have to satisfy requirements
    - OC-192/SDH-64, EoMPLS, …
- **Federated governance model**
  - Coordination between stakeholders



https://twiki.cern.ch/twiki/bin/view/LHCOPN/WebHome

# LHC Open Network Environment - The Background

- **So far, T1-T2, T2-T2, and T3 data movements have been using General Purpose R&E Network infrastructure**
  - **Shared resources (with other science fields)**
  - **Mostly best effort service**
- **Increased reliance on network performance → need more than best effort**
  - **Separate large LHC data flows from routed R&E GPN**
- **Collaboration on global scale, diverse environment, many parties**
  - **Solution has to be Open, Neutral and Diverse**
  - **Agility and Expandability**
    - Scalable in bandwidth, extent and scope
- **Organic activity, growing over time according to needs**
- **LHCONE Services being constructed:**
  - **Multipoint, virtual network (logical traffic separation and TE possibility)**
  - **Static/dynamic point-to-point Layer 2 circuits (guaranteed bandwidth, for high-throughput data movement)**
  - **Monitoring/diagnostic**

**http://lhcone.net**

# Following Important New Initiatives (Networking sample)

- **LHCONE Goal: Manage LHC's large data flows**
  - Workflow efficiency
  - As site capabilities progress (Nx10G, soon 100G)
  - Avoid triggering DOS alarms (and counter actions)
- **Lightpath technologies**
  - ESnet OSCARS, Internet2 ION, SURFnet/Ciena DRAC…
  - DYNES (reaching end-sites)
  - OGF NSI standards
- **Network virtualization**
  - Data center and WAN
  - Multipoint, multipath topologies
- **Software Defined Networking**
  - OpenFlow, …



NORDUnet — Current AutomatedGOLE + NSI

Demo NetworkSupercomputing 2011

Jerry Sobieski, NORDUnet

# US: DYNES Project, supporting LHC data movement

- **NSF funded; Internet2, Caltech, UoMichigan, Vanderbilt**
- **Nation-wide Cyber-instrument extending hybrid & dynamic capabilities to campus & regional networks**
- **Provides 2 basic capabilities at campuses and regional networks:**
1. **Network resource allocation such as bandwidth to ensure transfer performance**
2. **Monitoring of the network and data transfer performance**
- **Extending capability existing in backbone networks like ESnet and Internet2**
- **Tier2 and Tier3 sites need in addition**
3. **Hardware at the end sites capable of making optimal use of the available network resources**

http://internet2.edu/dynes



*Two typical transfers that DYNES supports: one Tier2 - Tier3 and another Tier1-Tier2.*

*The clouds represent the network domains involved in such a transfer.*

**US LHCNet**

DYNES is currently scaling up to full size (Phase 3, until August 31, 2012), and will transition to routine O&M in 2012-2013



U of Michigan
U of Wisconson-Madison
U of Chicago
UIUC
Indiana University
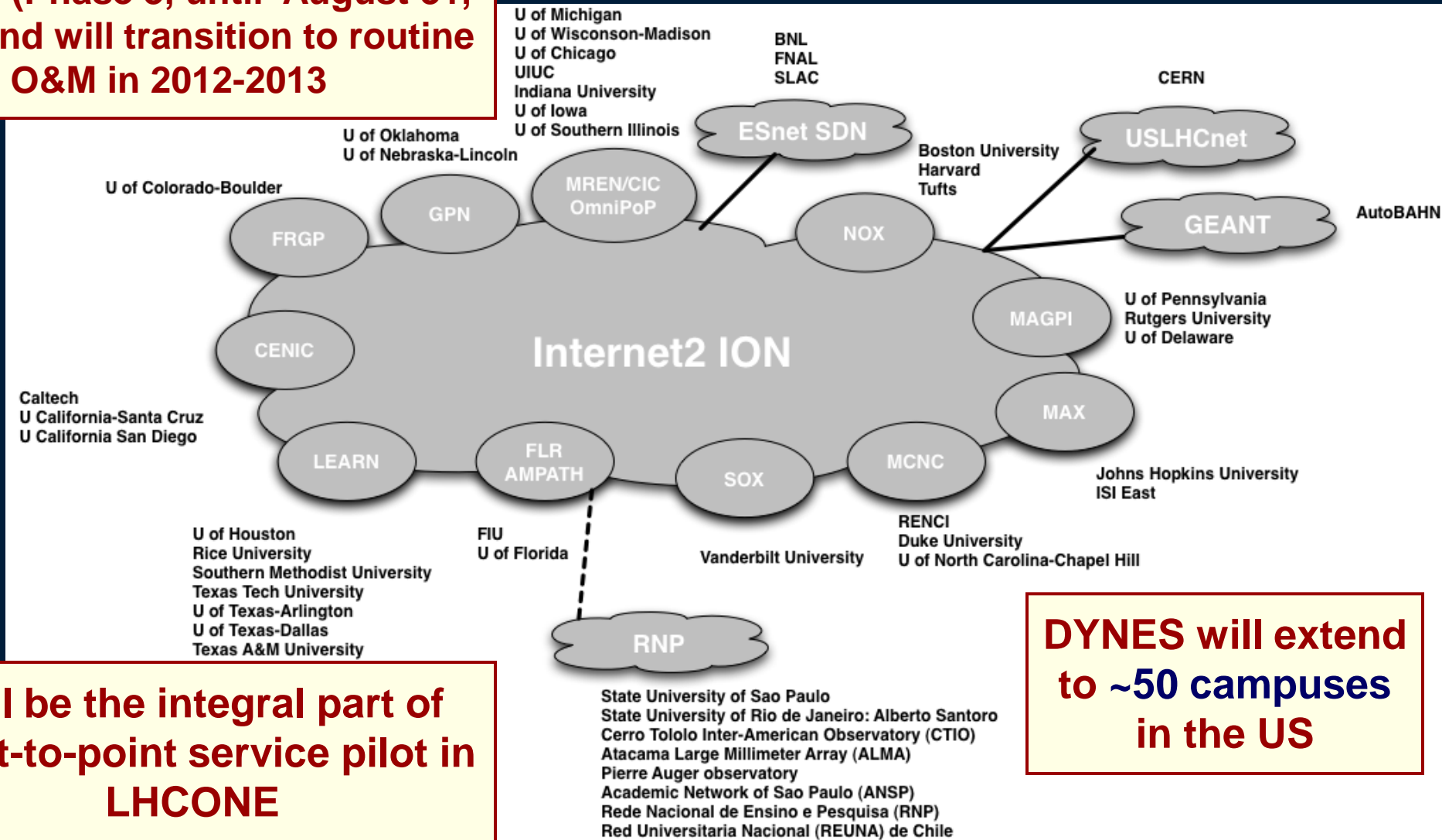U of Iowa
U of Southern Illinois

U of Oklahoma
U of Nebraska-Lincoln

U of Colorado-Boulder

BNL
FNAL
SLAC

CERN

Boston University
Harvard
Tufts

AutoBAHN

ESnet SDN

USLHCnet

GEANT

MREN/CIC OmniPoP

GPN

NOX

FRGP

Internet2 ION

MAGPI

U of Pennsylvania
Rutgers University
U of Delaware

CENIC

MAX

Caltech
U California-Santa Cruz
U California San Diego

LEARN

FLR AMPATH

SOX

MCNC

Johns Hopkins University
ISI East

U of Houston
Rice University
Southern Methodist University
Texas Tech University
U of Texas-Arlington
U of Texas-Dallas
Texas A&M University

FIU
U of Florida

Vanderbilt University

RENCI
Duke University
U of North Carolina-Chapel Hill

RNP

State University of Sao Paulo
State University of Rio de Janeiro: Alberto Santoro
Cerro Tololo Inter-American Observatory (CTIO)
Atacama Large Millimeter Array (ALMA)
Pierre Auger observatory
Academic Network of Sao Paulo (ANSP)
Rede Nacional de Ensino e Pesquisa (RNP)
Red Universitaria Nacional (REUNA) de Chile

Will be the integral part of point-to-point service pilot in LHCONE

DYNES will extend to ~50 campuses in the US
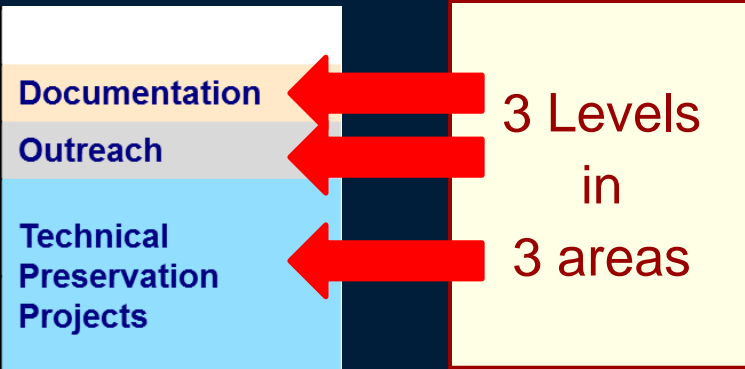
# A WORD ON ARCHIVING

# Long-term Data Preservation

- **ICFA\* Study Group formed in 2009**
  - DPHEP: Study Group for Data Preservation and Long Term Analysis in High Energy Physics
- **Recent end of several experiments**
  - @ LEP, HERA, PEP-II, KEKB, Tevatron
- **There is need to preserve information**

| | Preservation Model | Use Case | | |
|---|---|---|---|---|
| 1 | Provide additional documentation | Publication related info search | **Documentation** | 3 Levels |
| 2 | Preserve the data in a simplified format | Outreach, simple training analyses | **Outreach** | in |
| 3 | Preserve the analysis level software and data format | Full scientific analysis, based on the existing reconstruction | **Technical Preservation Projects** | 3 areas |
| 4 | Preserve the reconstruction and simulation software as well as the basic level data | Retain the full potential of the experimental data | | |

- **Recently published report: http://arxiv.org/pdf/1205.4667**

\* International Committee on Future Accelerators, http://www.fnal.gov/directorate/icfa/

# Summary

- **Excellent performance of the LHC and its Experiments**
  - Producing 10s of Petabytes of new data each year
  - Large statistics are necessary for discovery of new physics
  - Datasets distributed between >140 sites world-wide
  - WLCG is the underlying global, distributed computing infrastructure
- **Data and Computing Models are evolving**
  - More dynamic, more optimized
  - More reliant on network performance
- **Requires new approaches to networking**
  - Intelligent, holistic, systems approach is emerging
    - End-systems, dynamic optical network architectures, monitoring
  - DYNES will extend dynamic bandwidth allocation capability to 50 US campuses, and connect to partner networks abroad
  - LHCONE: virtual multi-domain network for traffic engineering LHC flows

# QUESTIONS…

**Artur.Barczyk@cern.ch**