# LHC Computing in der Hochenergiephysik

Hartmut Stadie
Universität Hamburg

LSDMA 2013 Spring Meeting
11-13 March 2013

Universität Hamburg
DER FORSCHUNG | DER LEHRE | DER BILDUNG

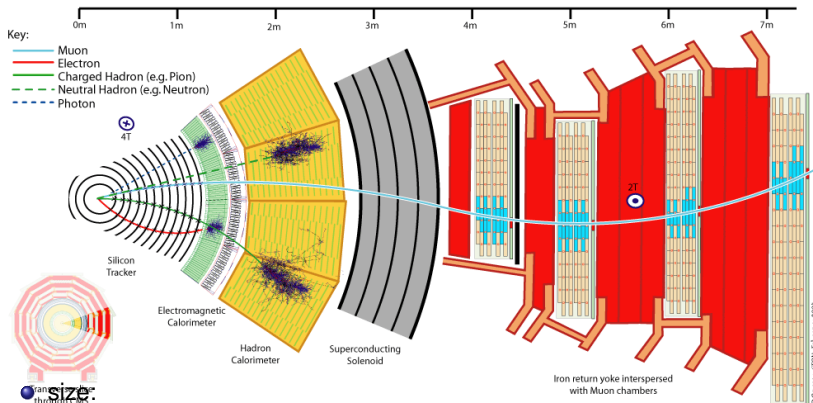SPONSORED BY THE

Federal Ministry
of Education
and Research

**Introduction**
0000

**Computing model**
000000

**Current status**
000000000

**Data management at a site**
000

**Conclusion**
0

# Outline

# The Large Hadron Collider (LHC)



- proton-proton collider
- circumference: 26.66 km
- $\sqrt{s} = 7 - 8$ TeV
- crossing rate: 50 ns rate: 20 MHz
- data per crossing (event): $\sim 1$ MB
- collider experiments: ALICE, ATLAS, CMS, LHCb,...
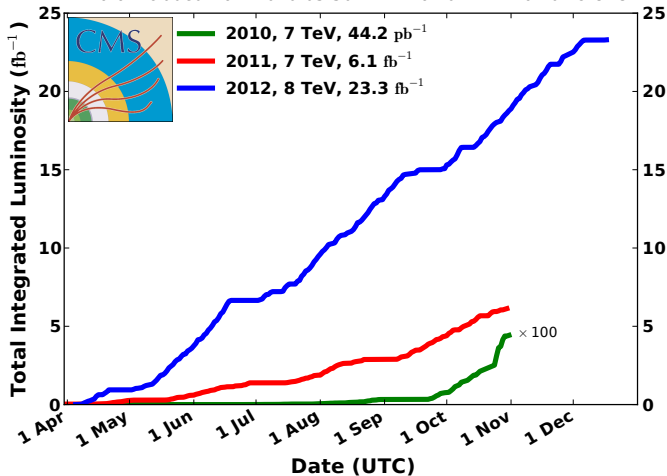- online filter: output rate $\sim 500$ MB/s

# The CMS Experiment



Key:
- Muon
- Electron
- Charged Hadron (e.g. Pion)
- Neutral Hadron (e.g. Neutron)
- Photon

4T

2T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with Muon chambers

D. Barney, CERN, February 2004

- Transverse slice through CMS
- size:
  length:     21.6 m
  diameter:   14.6 m

- mass: 12,500 t

- magnetic field:
  solenoid: 3.8 T

**Introduction**
○○●○

**Computing model**
○○○○○○

**Current status**
○○○○○○○○○

**Data management at a site**
○○○

**Conclusion**
○

# Current Status



**CMS Integrated Luminosity, pp**

# Data Organization

Physicists analyze the data requiring specific signatures (final states/channels).

### Collider data:

- split directly into O(10) primary data sets (PD) based on event signatures
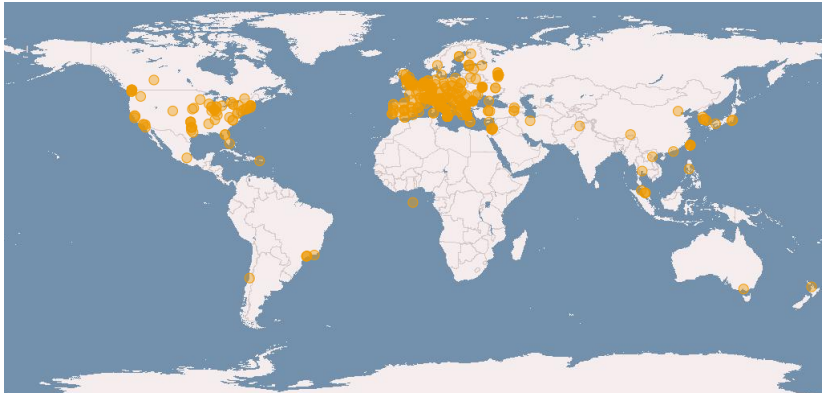- every PD in 2012 consists of roughly 100 TB

### Simulated events:

- comparisons with simulation to find deviations from the standard theory (e.g., new particles like the Higgs boson), measure properties
- overall size (in transfer DB): at least 10 PB

A typical analysis of the 2012 data might run over 200 TB a couple of times.
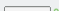
# Worldwide LHC Computing Grid (WLCG)

- provides resources and services to LHC experiments
- highly distributed system
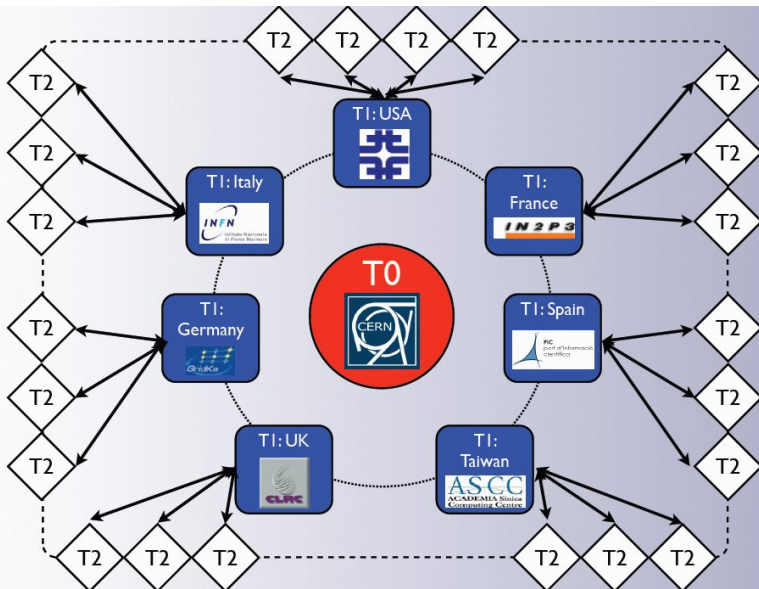
# Worldwide LHC Computing Grid (WLCG)

current overall resources sorted by the different classes of computing sites (Tiers):

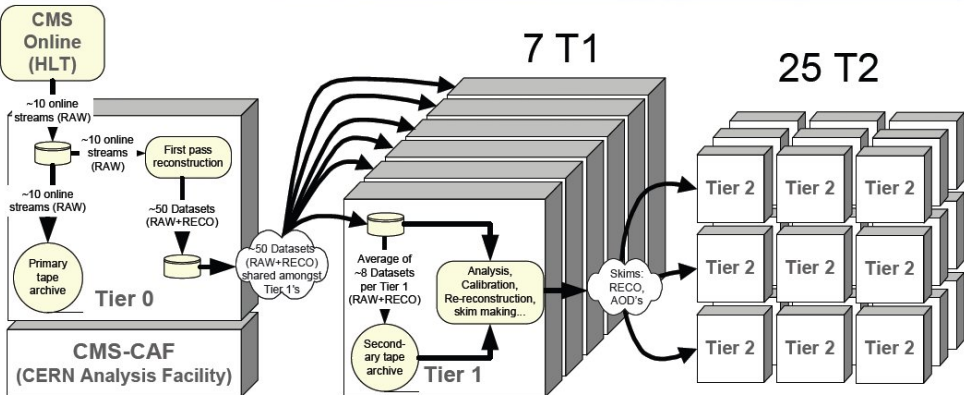| Name | Sites | CPUs | | | Online Storage Space (GB) | | Nearline Storage Space (GB) | |
|---|---|---|---|---|---|---|---|---|
| | | Physical | Logical | SI2000 | TotalSize | UsedSize | TotalSize | UsedSize |
| 0 | 1 | 5,228 | 27,564 | 171,558,336 | 33,423,802 | 73% | 80,583,197 | 91% |
| 1 | 12 | 30,705 | 92,091 | 233,409,658 | 88,468,065 | 70% | 102,050,988 | 45% |
| 2 | 107 | 31,856 | 171,365 | 418,632,474 | 137,697,045 | 44% | 1,235,440 | 68% |
| 3 | 29 | 1,562 | 8,143 | 16,341,756 | 2,913,321 | 59% | 0 | 0% |
| Total | 149 | 69,351 | 299,163 | 839,942,224 | 262,502,233 | 150,186,250 | 183,869,625 | 121,311,029 |

on top of bare resources:

- fabric: batch systems, storage system
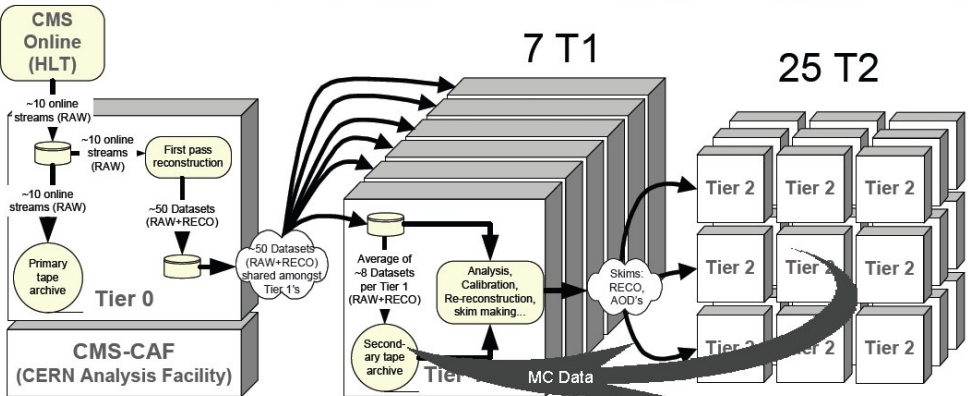- tools for resource sharing
- collective tools

# CMS Computing Model

# Data Flow

Introduction
○○○○

**Computing model**
○○○○○●○

Current status
○○○○○○○○○

Data management at a site
○○○

Conclusion
○

# Data Flow (MC)

**Introduction**
0000

**Computing model**
000000●

**Current status**
000000000

**Data management at a site**
000

**Conclusion**
0

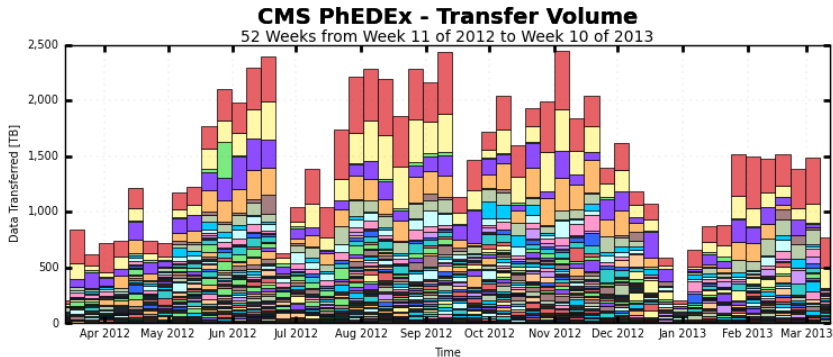# Experiment-Specific Services

### Experiment-Specific Services:

needed for distributed computing:

- production agents based on grid tools (WLCG job submission)
- dataset database (DBS) and trivial file catalog at sites
- dataset transfer service (PhEDEx)
  - uses grid tools (FTS, SRM) (as of 5 years ago)
  - special interfaces to Castor, dCache, etc) for file validation (checksums) and integrity tests
  - DB for transfers, dataset locations, commissioned links
  - agents for scheduling transfers, consistency checks, deletion, etc
  - transfers requests need to be approved by data manager of destination site
- (distributed) calibration database (FroNTier) squid web cache
- analysis job submission tool (CRAB) grid tools (WLCG job submission), bridges to local batch systems

Introduction
0000

Computing model
000000

Current status
●000000000

Data management at a site
000

Conclusion
0

# Data transfers by destination (volume per week)

Introduction
oooo

Computing model
oooooo

Current status
ooooooooo

Data management at a site
ooo

Conclusion
o

# Data transfers by destination (integrated)



**CMS PhEDEx - Cumulative Transfer Volume**
52 Weeks from Week 11 of 2012 to Week 10 of 2013

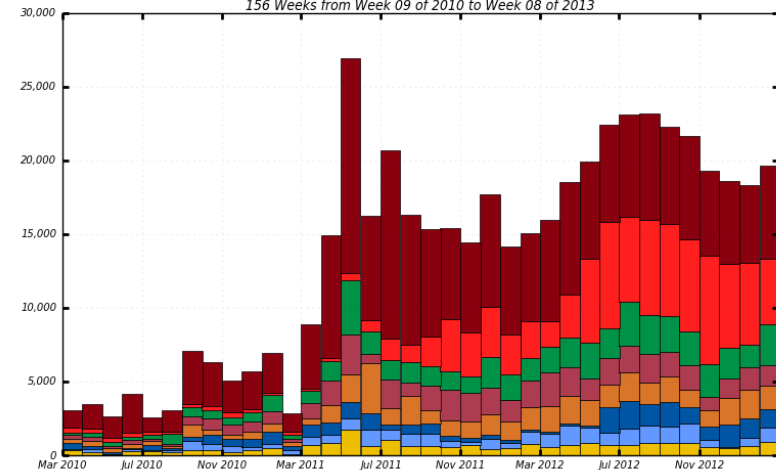Total: 76,426 TB, Average Rate: 0.00 TB/s

Introduction
○○○○
Computing model
○○○○○○
Current status
○○●○○○○○○○
Data management at a site
○○○
Conclusion
○

# Tier-0/1 data processing obs



**Running jobs**
*156 Weeks from Week 09 of 2010 to Week 08 of 2013*

Legend:
- T1_US_FNAL
- T1_FR_CCIN2P3
- T0_CH_CERN
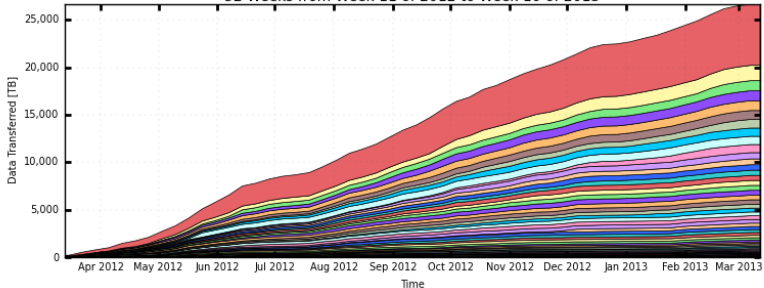- T1_TW_ASGC
- T1_IT_CNAF
- T1_ES_PIC
- T1_DE_KIT
- T1_UK_RAL

Introduction
0000

Computing model
000000

Current status
000●00000

Data management at a site
000

Conclusion
0

# Tier-2 data transfers by destination

Introduction
○○○○

Computing model
○○○○○○

Current status
○○○○●○○○○

Data management at a site
○○○

Conclusion
○

# CMS Computing Model

Introduction
0000

Computing model
000000

Current status
000000●0000

Data management at a site
000

Conclusion
0

19 / 26

# Tier-2 data transfers between Tier-2s by source



T2-T2 fraction: 25%

Introduction
0000

Computing model
000000

Current status
0000000●00

Data management at a site
000

Conclusion
0

# Tier-2 data transfers (volume per week)



**CMS PhEDEx - Transfer Volume**
52 Weeks from Week 11 of 2012 to Week 10 of 2013

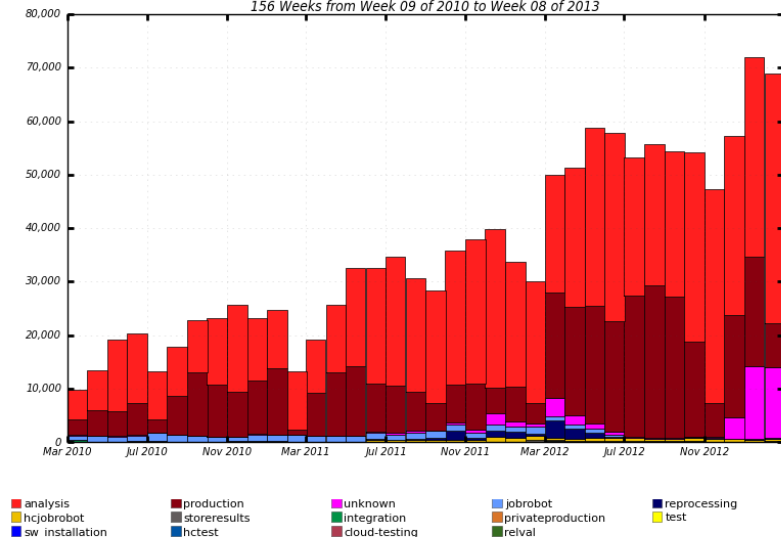Maximum: 956.37 TB, Minimum: 55.73 TB, Average: 502.88 TB, Current: 185.70 TB

Introduction
0000

Computing model
000000

Current status
00000000●0

Data management at a site
000

Conclusion
0

# Tier-2 activities



*Running jobs*
156 Weeks from Week 09 of 2010 to Week 08 of 2013

# Analysis jobs per site



*Running jobs*
*156 Weeks from Week 09 of 2010 to Week 08 of 2013*

**Introduction**
0000

**Computing model**
000000

**Current status**
000000000

**Data management at a site**
●00

**Conclusion**
○

# Data Management at a Tier-2 (e.g. DESY)



DESY houses National Analysis Facility: local space > 1 PB
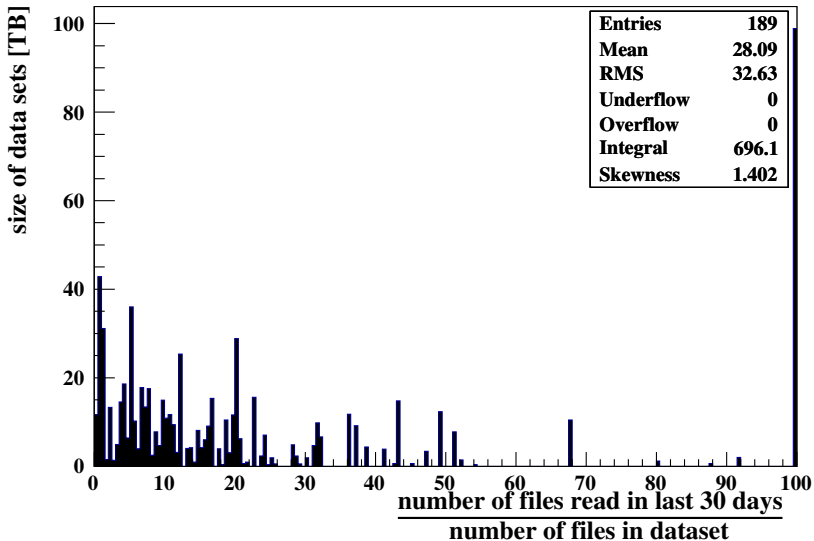
# Management of local data

## Workflow

- users request datasets
- data manager approves request if justified
- data manager also identifies and requests samples of common interest
- every month popularity is evaluated based on dCache access logs
- unused samples are scheduled for deletion after a grace period

## Problems

- PhD thesis lasts three years (requests for outdated data)
- users hardly request the deletion of old data
- no trust in file transfers, do not see Tier-2 storage as a cache
- want to have all data at one place

# Example for January

## Conclusions

### Lay person's conclusions

- LHC experiments deal with large data
- use "divide et impera" to break problems down
- resources and middleware from WLCG
- experiment specific tools needed
- good separation and interfaces between sites for scaling needed
- data management very difficult (centralized systems to rigid, by-demand/request system needs resources to scale with data)