# 1st ASTERICS-OBELICS Workshop

12-14 December 2016, Rome, Italy.

# A globally distributed data management solution

WLCG – the LHC's offline computing platform
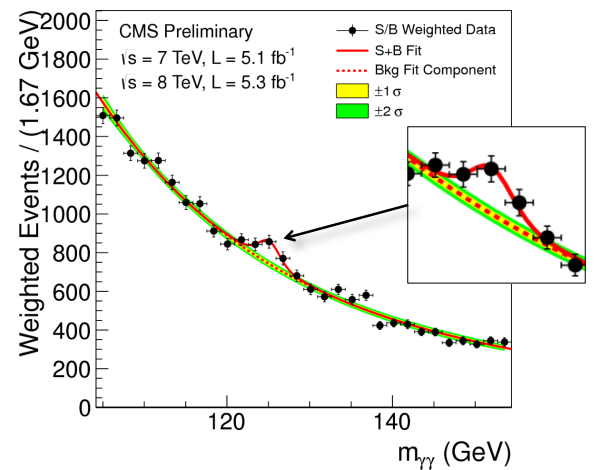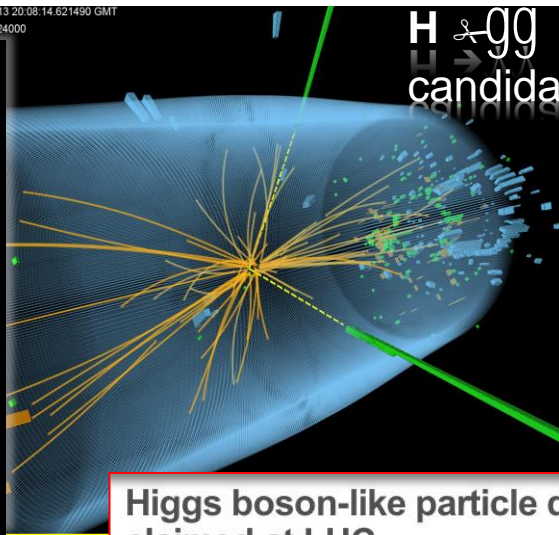
Oliver Keeble

CERN Storage Group

Global Effort → Global Success

Results today only possible due to extraordinary performance of accelerators – experiments – Grid computing

Observation of a new particle consistent with a Higgs Boson (but which one…?)

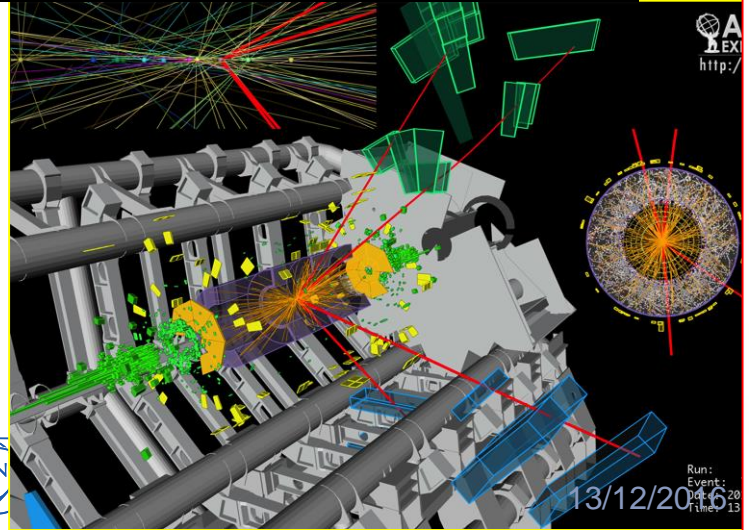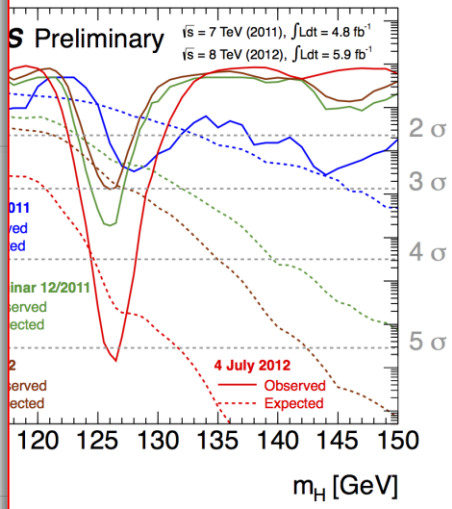Historic Milestone but only the beginning

Global Implications for the future

R-D Heuer

**Higgs boson-like particle discovery claimed at LHC**

💬 COMMENTS (1665)

By Paul Rincon
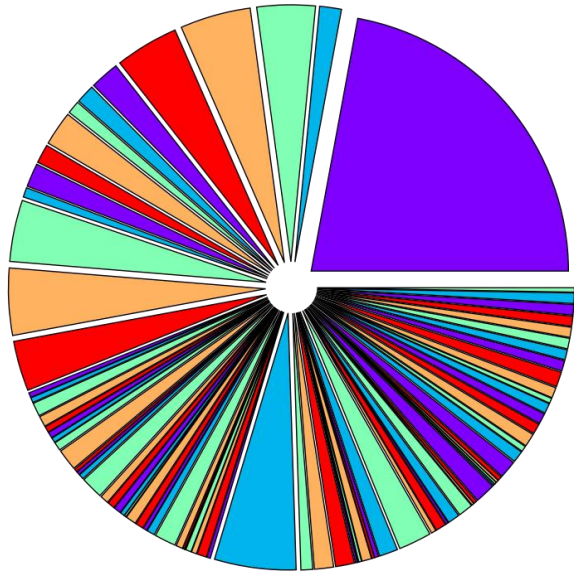Science editor, BBC News website, Geneva

The moment when Cern director Rolf Heuer confirmed the Higgs results

**Cern scientists reporting from the Large Hadron Collider (LHC) have claimed the discovery of a new particle consistent with the Higgs boson.**
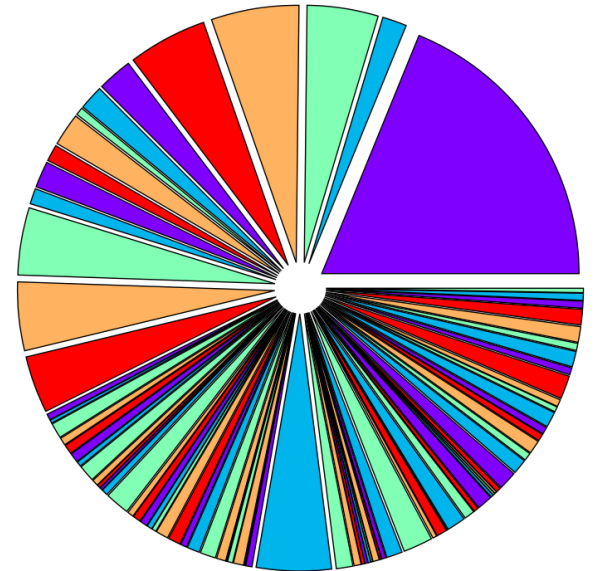
CMS Preliminary
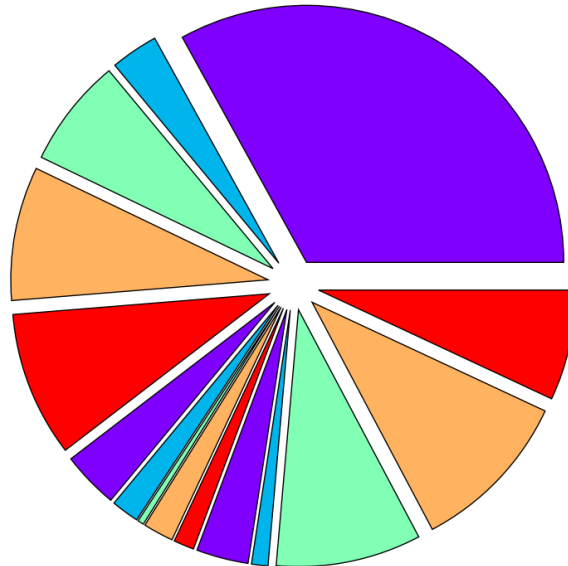√s = 7 TeV, L = 5.1 fb⁻¹
√s = 8 TeV, L = 5.3 fb⁻¹

S/B Weighted Data
S+B Fit
Bkg Fit Component
±1 σ
±2 σ

$m_{\gamma\gamma}$ (GeV)

CMS Preliminary
√s = 7 TeV (2011), ∫Ldt = 4.8 fb⁻¹
√s = 8 TeV (2012), ∫Ldt = 5.9 fb⁻¹

4 July 2012

Observed
Expected

$m_H$ [GeV]

13/12/2016

Astenics Obelics Roma 2016

3

# Pledged Resources

CPU by computer centre
Total ~400k cores

Tape by computer centre
Total ~400 PB

Disk by computer centre
Total ~300 PB

Running jobs: 305236
Active CPU cores: 435925
Transfer rate: 14.40 GiB/sec

# Workflows



**Simulation**

**Reconstruction**

**Analysis**

**Detector**

**Distribution & Storage**

**Archival**

More than half the CPU goes on simulation.

Most of the rest is reconstruction.

The remainder is analysis.

# From yesterday…

## Different probes/methods/specifications

We are here!

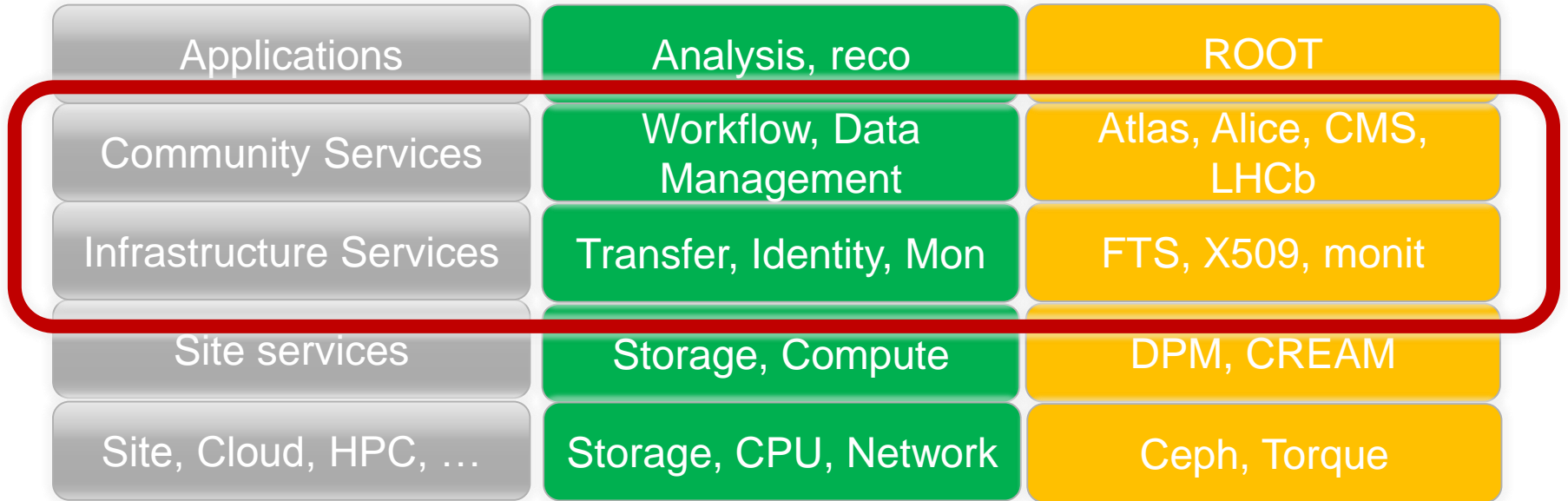| Projects | Data Processing | Main requirements/challenges |
|---|---|---|
| EVENT-BASED (γ-rays, CR, ν) <br><br> CTA. KM3Net … | Evt-builder, calib. and reconstruction; reduction, real-time science. | Raw big-data. Data formats. Algorithms. On-site operation and reduction. Cooperative science tools. Observatory (A&A). Multi-λ. […] |
| IMAGE-BASED (far-IR, VIS) <br><br> EUCLID. LSST … | Surveys/deep observation; combining photometer and spectrograph info.; Catalogue of objects. | Big-data products: data base challenges. Graphical processing, Algorithms. Images format. Catalogue preservation and query. A&A. […] |
| SIGNAL-BASED (Radio, GW) <br><br> SKA, LIGO-Virgo … | Noise cleaning; time-series, mathematical processing (FT) converting signal in images. | Algorithms. New computing architectures and data centres. Combination of HPC and HTC. Fast soft reduction. Data mining and preservation. A&A […] |

# What advantages do we have?

- Event independence
- "Read-only" data
- File-based data
- Scientific Linux
- Coarse grained (VO) authorisation
- X509 acceptance
- Large proportion of CPU intensive work

# WLCG Stack

| | | |
|---|---|---|
| Applications | Analysis, reco | ROOT |
| Community Services | Workflow, Data Management | Atlas, Alice, CMS, LHCb |
| Infrastructure Services | Transfer, Identity, Mon | FTS, X509, monit |
| Site services | Storage, Compute | DPM, CREAM |
| Site, Cloud, HPC, … | Storage, CPU, Network | Ceph, Torque |

# WLCG Stack

| Applications | Analysis, reco | ROOT |
| --- | --- | --- |
| Community Services | Workflow, Data Management | Atlas, Alice, CMS, LHCb |
| Infrastructure Services | Transfer, Identity, Mon | FTS, X509, monit |
| Site services | Storage, Compute | DPM, CREAM |
| Site, Cloud, HPC, … | Storage, CPU, Network | Ceph, Torque |

# Not (only) a grid

- The WLCG infrastructure comprises
  - Grid - pledged
  - Cloud - rented
  - HPC - allocated
  - Volunteer - donated
  - Concepts
    - Opportunistic resources
    - Pre-emptibility

# Volunteer



**Slots of Running Jobs**
662 Hours from 2015-02-03 to 2015-03-03 UTC

- 2nd largest simulation site
- Running 4-5k parallel jobs
- 20M events simulated
- 5M CPU hours

http://cern.ch/go/v6jG

Legend:
- BNL-ATLAS
- UKI-LT2-BRUNEL
- MWT2
- UKI-NORTHGRID-MAN-HEP
- BU_ATLAS_TIER2
- WUPPERTALPROD
- IFIC-LCG2
- UNI-FREIBURG
- CERN-P1
- BOINC
- CERN-PROD
- INFN-NAPOLI-ATLAS
- UKI-SOUTHGRID-RALPP
- SIGNET
- UKI-NORTHGRID-LANCS-HEP
- UTA_SWT2
- GRIF-IRFU
- GRIF-LPNHE
- RAL-LCG2
- TRIUMF-LCG2
- UKI-SCOTGRID-GLASGOW
- NDGF-T1
- UKI-LT2-QMUL
- SWT2_CPB
- TAIWAN-LCG2
- HU_ATLAS_TIER2
- IAAS
- IN2P3-CC
- DESY-HH
- FZK-LCG2
- PIC
- UKI-SOUTHGRID-OX-HEP
- INFN-ROMA1
- PRAGUELCG2
- ARNES
- NIKHEF-ELPROD
- INFN-T1
- LRZ-LMU
- AGLT2
- WT2
- UKI-LT2-RHUL
- INFN-MILANO-ATLASC
- CYFRONET-LCG2
- UKI-NORTHGRID-LIV-HEP
- ... plus 72 more

Maximum: 112,630 , Minimum: 0.00 , Average: 63,358 , Current: 62,935

# HPC – backfill on Titan



**ATLAS production running on Titan in 2016**

ATLAS Titan Usage Per Month

Pure opportunistic backfill mode, no project allocation, ATLAS Geant4 simulations

Sergey Panitkin                                                          11

# 2016 data volumes

**Transfered Data Amount per Virtual Organization for WRITE Requests**



**Transfered Data Amount per Virtual Organization for WRITE Requests**



~160 PB on tape at CERN
500 M files

LHC data – Continue to break records:
10.7 PB recorded in July
CERN archive ~160 PB

June-Aug 2016
>500 TB / day
(Run 1 peak for HI was 220 TB)

Physics Data in CASTOR

# Data - transfer





**Transfer Throughput**
2016-12-09 04:40 to 2016-12-09 08:40 UTC

Most LHC transfers are managed by the File Transfer Service (FTS)

Try it at https://webfts.cern.ch

# Data – storage systems

**Instances**

Pie chart (Instances) segments: DPM, DCACHE, STORM, CLASSICSE, BESTMAN, CASTOR, XROOTD, EOS, HDFS, ARC

**Storage**

Pie chart (Storage) segments: DPM, DCACHE, STORM, BESTMAN, CASTOR, XROOTD, EOS, HDFS

## Basic Architecture

```
Client
  → Head
  → Disk    Disk
              → Tape
```

Systems are evolving towards standards. All now offer HTTPS access.

# Global Data Federation



```
/dir1
/dir1/file1
/dir1/file2
/dir1/file3
```

With 2 replicas

*Site A*

```
.../dir1/file1
.../dir1/file2
```

*Site B*

```
.../dir1/file2
.../dir1/file3
```

In use by
- Atlas (FAX)
- CMS (AAA)
- Main uses
  - Failover
  - Overflow
  - Diskless sites

# Data – what is it?

- ROOT files

- Typically a few GB each

- Column-like structured storage

- Lots of I/O optimisation

  - WAN access

```
Open https://server/data.root
While (next event) {
  do stuff;

}
```

# Software distribution

- CVMFS

- r/o cached fs

- >350M files

```
[lxplus109] ls /cvmfs
alice.cern.ch           clicdp.cern.ch          ilc.desy.de
alice-ocdb.cern.ch      cms.cern.ch             lhcb.cern.ch
ams.cern.ch             cms-ib.cern.ch          na61.cern.ch
atlas.cern.ch           cvmfs-config.cern.ch    sft.cern.ch
atlas-condb.cern.ch     geant4.cern.ch
atlas-nightlies.cern.ch grid.cern.ch
[lxplus109]
```

# The road ahead

# High-Lumi LHC resource estimates

## Data estimates for 1st year of HL-LHC (PB)

Legend: ALICE, ATLAS, CMS, LHCb

(Stacked bar chart with categories "Raw" and "Derived", y-axis from 0 to 1000)

## CPU Needs for 1st Year of HL-LHC (kHS06)

Legend: ALICE, ATLAS, CMS, LHCb

(Stacked bar chart with category "CPU (HS06)", y-axis from 0 to 250000)

Data:
- X10 from 2016
  - Raw 2016: 50 PB → 2027: 600 PB
  - Derived (1 copy): 2016: 80 PB → 2027: 900 PB

CPU:
- x60 from 2016

Technology at ~20%/year will bring x6-10 in 10-11 years

# HL-LHC Solutions

- Raw data
- Triggers
- Detector design
- …

- Reconstruction and simulation algorithms

Parameters

Core Algorithms

Infrastructure

Software Performance

- Performance/architectures/memory
- Tools
- Concurrency
- Vectorisation
- Collaboration with externals – via HSF
- …

- New grid/cloud models
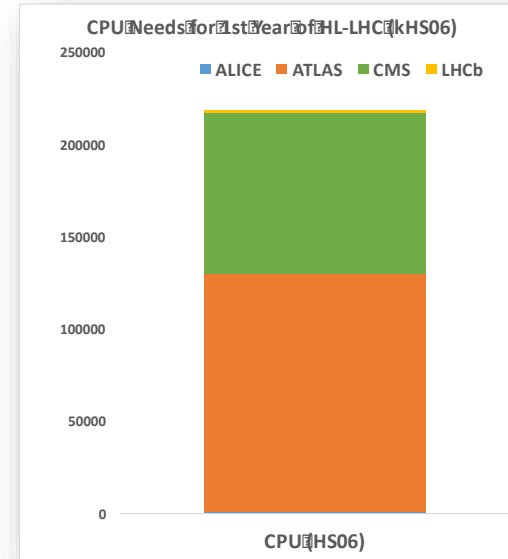- Optimise CPU/disk/network
- Economies of scale via clouds, joint procurements etc.
- Opportunistic resources
- Pre-emptible jobs
- Storage consolidation
  - WAN access
- Data strategies
- Caching solutions
- …

# Summary

- WLCG is the production offline computing platform for the 4 LHC experiments
- Can process multiple PB / month
- In 2025 we will have a new accelerator with new experiments
  - Order of magnitude more load at fixed cost
  - Technology sharing = sustainability & reduced costs

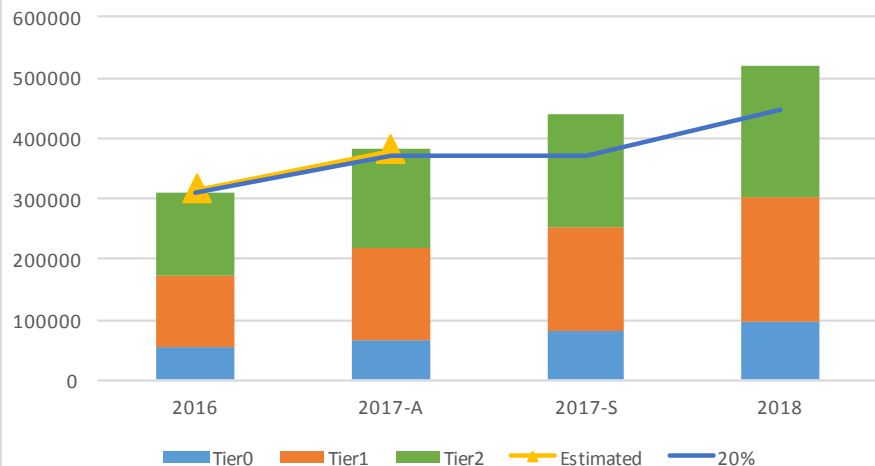# Acknowledgement

# Supplementary Slides
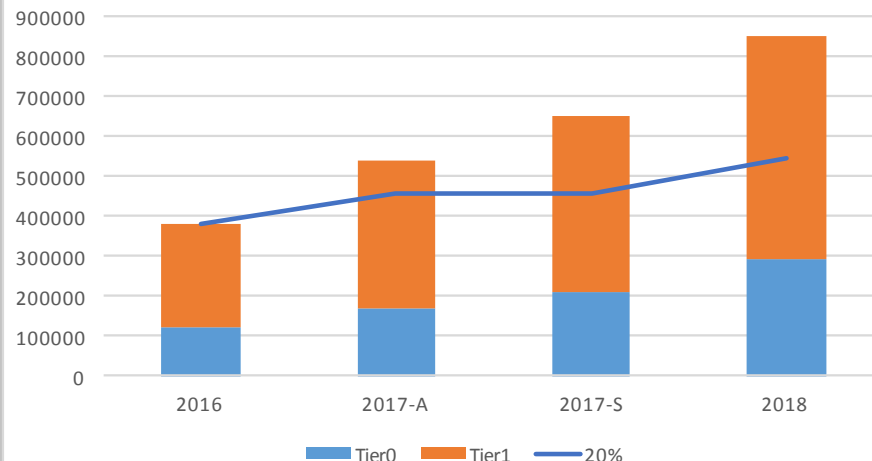
# Data Management Directions

- Reduce cost/volume
  - cost of storage management
    - integrating standard (non HEP) solutions e.g. ceph
    - protocol zoo, SRM-less operation
    - T2 storage as cache
    - multi-site storage
    - regional federations
    - cloud storage
    - system manageability
  - storage overheads
    - redundancy
      - replication, erasure, RAID levels etc
    - reduce system reliability requirements?
      - reduce cost/impact of data loss
    - component technology
      - shingled disks
      - consumer/enterprise disks

- Reduce volume used
  - reduced number of global replicas
    - remote access
    - latency hiding
      - applications, overcommitting
    - global federations
    - CPU-only resources (inc cloud)
  - data formats and lifecycle, intermediate products
  - resource reporting
    - monitoring usage
    - eliminating dark data
  - data "enrichment"
    - popularity
    - caching, avoiding unused data
      - promoting locality in workflows
  - trading disk for…
    - tape
      - data parking
    - CPU
      - maintain metadata enabling regeneration of data on demand

# LHC: "outstanding performance"



Estimated: Estimates made in 2014 for Run 2 up to 2017

20%: Growth of 20%/yr starting in 2016 ("flat budget")