Transfer rate: 1.62



PetaByte-scale data management at CERN

Università degli Studi di Trieste Corso di Laurea in Fisica

4th of May 2012

ssimo Lamanna IT Data Storage Services

> © 2011 Europa Tiechnologies US Dept of State Geographer © 2011 MapLink/Tele Atlas © 2011 Google

Computer Centre By Numbers



full name: Computer Centre By Numbers short name: CCBYNUM

group: IT-CF-FPP

site: CERN

email: imre.szebenyi@cern.ch

manager: Imre Szebenyi 😔



Additional service information	1 (more)
Number of processors:	14,972
Number of cores:	64,623
Memory capacity (TiB):	165
Memory modules:	55,729
Raw HDD capacity (TiB):	62,660
Number of HDD's:	62,023
Number of systems:	7,975
Number of RAID controllers:	3,607
Number of enclosures:	1,554
SPEC CPU2006:	503,637
Number of racks:	1,070
Number of virtual machines:	3,908
Number of Fibre channel ports:	742
Number of 1G ports:	16,773
Number of 10G ports:	622
Current power consumption (kW):	2,186
Current power consumption (kVA):	2,305

24x7 operator and system admin support

•

•

Management and Automation framework for large scale Linux clusters

Hardware installation & retirement ~7,000 hardware movements/year; ~1000 disk failures/year

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

The LHC Computing Grid - February 2012

What are all these CERN computers/disks/networks... for?



RC

(Detector) simulation



Data acquisition and processing



Lumi section: 387 Sat Apr 24 2010, 14:00:54 CEST

Electrons p_T = 34.0, 31.9 GeV/c Inv. mass = 91.2 GeV/c²



Data analysis



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

From Physics to Raw Data



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

CERN

From Raw Data to Physics



CERN

Analysis flow (user view)



But how this is done *in practice*? Of course we need CPUs, disks, networks etc.. The problem is orchestrating hardware resources, software and humans :)

"All" data are stored in files (aggregated as "datasets" = collections of files). Only a small fraction of data in real DBs (e.g calibrations). This is one characteristics of HEP computing.

CERN

The role of the CERN Computer CERNIT



RC

CERN IT Department CH-1211 Genève 23

Switzerland www.cern.ch/it







The LHC Computing Grid - April 2011

ATLAS: 7,000 tons, 150 million sensors generating data 40 millions times per second i.e. a petabyte/s (1 million GB/s)

EXPERIMENT

China Slovenia Colombia Spain Czech Republic Sweden Denmark Switzerlan France Ta ATL/ Georgia Tu Fron Greece Us ~100 Israel CE

Japan

ATLAS is around than 3,000 collaborators From 169 universities from 37 countries ~1000 students!!! have at a state of



One of CERN main responsibilities

(LHC scientific programme - LHC data management)

Tier 0 function

- Experimental data to tape (RAW)
- Data distribution (Tier0 \rightarrow Tier1s)
- Data reconstruction and redistribution
 - RAW → ESD/DST,AOD..
- Archival data (e.g. simulation data generated in other centres)

Analysis centre

- Each LHC physicist has access to CERN!
- CERN-based analysis groups



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

CERN Computer Centre: Storage, Distribution CERN and Processing (Reconstruction and Analysis) D



The HEP Data Challenge



LHC will run for 20 years Experiments *are* producing about **15 Million Gigabytes** of data each year (about 20 million CDs!) LHC data analysis requires a computing power equivalent to ~100,000 of today's fastest PC processors

Requires many cooperating computer centres, as CERN can *only* provide ~20% of the capacity

A challenge for physics... ... and a challenge for technology research and industry as well





CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

Big data non HEP specific (... any longer) Department

HEP (LCG):

- •<u>CASTOR</u>
- •<u>EOS</u>
- •Other LCG systems
 - •dCache (most popular among Tier1s) - DESY
 - •DPM (most popular among Tier2s) - CERN
 - •STORM INFN
 - •Xroot SLAC
 - •BestMan Berkeley

Non HEP:

•Google, Facebook, etc...

•e.g. HADOOP

•Astronomy:



•LSST: xroot for DB access (large-scale distributed queries)

•Everyone:



Cloud Computing

OpenStack swift

•More file system



•AFS, NFS 4.1, CEPH 13

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

Moore's law

CERN

Department

Microprocessor Transistor Counts 1971-2011 & Moore's Law



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

RC



"Historical" example...

CERN celebrates 20 years of a free, open web

30 Apr 2013

Geneva, 30 April 2013 - Twenty years ago CERN¹ published a statement that made the World Wide Web ("W3", or simply "the web") technology available on a royalty-free basis. By making the software required to run a web server freely available, along with a basic browser and a library of code, the web was allowed to flourish

The technology, invented in 1989 at CERN by Tim Berners-Lee, was originally conceived and developed to meet the demand for information sharing between physicists in universities and institutes around the world.

Other information retrieval systems that used the Internet - such as WAIS and Gopher - were available at the time, but the web's simplicity along with the fact that the tech adoption and development.

"There is no sector of society that has not been transformed by the the web", says Rolf Heuer, CERN Director-General. "From rese has been reshaping the way we communicate, work, innovate an of the way that basic research benefits humankind."

The first website at CERN - and in the world - was dedicated to was hosted on Berners-Lee's NeXT computer. The website desc to access other people's documents and how to set up your own : the original web server - is still at CERN, sadly the world's first v address

To mark the anniversary of the publication of the document that a to use, CERN is starting a project to restore the first website and associated with the birth of the web. To learn more about the pro http://info.cem.ch

http://info.cern.ch/hypertext/WWW/TheProja

The WorldWideWeb (W3) is a wide-area hypermetry and Everything there is online about W3 is linked directly or an

What's out there? Pointers to the world's online information of

World Wide Web

1990s:

The web was invented at CERN! The machine used by Tim Berners-Lee in 1990 to develop and run the first WWW server, multi-media browser and web editor.

CERN

Department

The LHC Computing Grid - April 2011

COMPASS proposal (1996)

- (Data) challenge
 - Semi inclusive → reconstruct "the rest" of the deep inelastic event
 - 35 MB/s
 - 300 TB/year



More than a prototype !!!

- 3 times less data 14 years before LHC (LHC experiment) (7 Moore's cycles hence ~27)
- 40x factor "against" us $(2^{1/3})$

- Use a Linux PC farm
 - Instead of the "usual super computers"

CERN

Department

- C++
 - Instead of good ol' Fortran IV
- All data in a data base (Objectivity/DB)
 - Instead of "plain" files
 - Object store (noSQL)
- First experiment using CASTOR
 - Instead of writing tapes "by ourself"

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

COMP_{ASS}

RC









(m)

RC



nt

CH-1211 Genève 23 Switzerland

www.cern.ch/it

Data management today

- CASTOR for Tier0
 - "RAW" data, tape access
- EOS for analysis
 - Analysis data, disk only
- Optimised for different use cases (T0 and analysis)



CERN

Department

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

RC

Total installed disk capacity : 13.0 PB + 22.4PB Over 100 PB in the CERN store (disk+tape) NB: 1 week of LHC (2012) = 1 PB on tape!



ACRON service

RC

Switzerland www.cern.ch/it

CERN IT Department CH-1211 Genève 23

Instance	Jobs	Efficiency (*)	Unique user/host
@cern	859.6 K	94.8 %	689.0
AES convico			

AFS service

Instance	Capacity	Files	Δ	Size	Δ
@cern	501.7 TB	1.7 G	20.3 M	198.8 TB	3.2 TB

FILER service

Instance	Capacity	Files	Δ	Size	Δ
@cern	192.0 TB	118.4 M	0.7 M	23.7 TB	0.1 TB

CASTOR service

Instance	Files	Δ	Size	Δ	OnTape	Δ
@cern	311.2 M	-0.2 M	86.3 PB	132.0 TB	74.6 PB	121.0 TB

EOS service

Instance	Files	Δ	Size	Δ
@cern	134.5 M	-1.7 M	15.7 PB	462.1 TB
alice	96.8 M	-2.4 M	4.1 PB	-86.4 TB
atlas	29.9 M	0.5 M	4.2 PB	132.8 TB
cms	6.8 M	0.2 M	4.8 PB	400.0 TB
Ihcb	1.0 M	0.1 M	2.5 PB	15.8 TB

Data durability on a disk farm

% df -H

CERNIT

11 RAID1 pairs



CERN IT Departmen
CH-1211 Genève 23
Switzerland
www.cern.ch/it

% df -H				
Filesystem	Size	Used	Avail	Use% Mounted on
/dev/sdb	1.9T	2.2G	1.9T	1% /srv/castor/01
/dev/sdc	2.0T	2.2G	2.0T	1% /srv/castor/02
/dev/sdd	2.0T	2.2G	2.0T	1% /srv/castor/03
/dev/sde	2.0T	2.2G	2.0T	1% /srv/castor/04
/dev/sdf	2.0T	2.2G	2.0T	1% /srv/castor/05
/dev/sdg	2.0T	2.2G	2.0T	1% /srv/castor/06
/dev/sdh	2.0T	2.2G	2.0T	1% /srv/castor/07
/dev/sdi	2.0T	2.2G	2.0T	1% /srv/castor/08
/dev/sdj	2.0T	2.2G	2.0T	1% /srv/castor/09
/dev/sdk	2.0T	2.2G	2.0T	1% /srv/castor/10
/dev/sdl	2.0T	2.2G	2.0T	1% /srv/castor/11

2 TB * 11 filesystem (* 2 copies)

~ 44 TB raw storage in one box

1000 boxes \rightarrow 2000 disks (JBOD)

- + compact (~1 PB in 4 boxes)
- potentially "too little spindles" (analysis)
- potentially "too small network I/F" (analysis)
- data placement problems?

Different setups (mirror)



- Durability (data on HD1)
 - $\approx 1 p(HD1) * p(HD1')$
 - First approximation dominated by the number of replicas
 - If anything, in set up 1 there is a correlation ("positive") which accounts for correlated failures in HD1 and HD1'
- Availability (data on HD1)
 - Setup 1
 - 1– p(host1)
 - Setup 2
 - $\sim \approx 1 p(host1) * p(host2)$

As today:

CASTOR: mainly Hardware RAID (RAID1) EOS: tunable number of full copies (n>1)

HD2



HD1'

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

RC

Reliability... Good



- File loss is not nice but unavoidable with a certain probability
 - RAID-1 does not protect against controller or machine problem, filesystem corruption and finger trouble
 - typically important files can be recovered from offsite
- In case of backup (CASTOR) the tape reliability is helping the disk one





www.cern.ch/it

EOS replica mechanics



Network IO for file creations with 3 replicas:



CER

Department

500 MB/s injection result in

- 1 GB/s output on eth0 of all disk servers (write out of FS1 0.5 GB/s twice)

- 1.5 GB/s input on eth0 of all disk servers

(0.5 GB/s x 3 copies)

Plain (no replica) Replica (here 3 replicas) More sophisticated redundant storage (RAID5, RAID6, LDPC) 25



(asynchronous operations)

www.cern.ch/it

Bringing all together...

The World Wide Web provides seamless access to information that is stored in many millions of different geographical locations

The **Grid** is an infrastructure that provides seamless access to computing power and data storage capacity distributed over the globe

CERN IT Department

CH-1211 Genève 23

Switzerland

www.cern.ch/it



CERN



WLCG Tiers Organization



Tier-0 (CERN):

Data recording Initial data

CERN

Department

- reconstruction
- Data distribution

Tier-1 (11 centres): Permanent storage Re-processing Analysis

Tier-2 (~130 centres):

- Simulation
- End-user analysis

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

How does it work?

ERN**IT** Department

ATLAS

Not substantially different for the other HEP experiments Heavily simplified...

What do we want to achieve

The user wants to specify a subset of the data and run applications on it (chain of programs reading intermediate outputs)

Only at the end of the chain data sizes and computational complexity this can be (possibly) done on a laptop

1000+ of physicists worldwide after the <u>same</u> data

Behind the scenes...

CERN







5 PB Tape storage

10 application domains

More info: INFN (Istituto Nazionale Fisica Nucleare): http://www.infn.it IGI (Italian Grid Initiative): http:// www.italiangrid.org/



Outlook

CERN data for LHC:

- Solid foundation (CASTOR)
- Successful new product (EOS)
- CASTOR and EOS coexist (complementary)

Goals

- Stability with low operations costs
- Open to adapt to new ideas (also from non-HEP areas)
- Open directions ahead:
 - Data management at the heart of our activity
 - Critical for the success of physics research
 - An exciting field of study by itself
 - New stuff coming
 - » Is it any good?



CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it



Questions?





CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

RC

	CERN	e neerannene nome rage			-
	https://ert.cern.ch/browse_www	//wd_pds?p_web_site_id=1	☆ □ ▼)	• (G* twiki computer cent	٩)
I GCal Docs News ▼	Bank - RW - Dict - Mix - Old - Ub	untu and Free S 🔻 Latest Headlines	৯ Apple ▼ Google Ma	aps YouTube Popular - P	ytho
e-Rec	ruitment	• Vacancies O All CERN C	HR Department	Silte map Contact us CERN Hor Search	ne
HR General Information	Recruitment Training	Staff Career	Services	CERN Official Documen	ts
Welcome Page					
Intranet	CERN is the European Organization fo Franco-Swiss border near Geneva (mo	r Nuclear Research, based on the pre).	News		
Posictor in o.PT			News		
Login to o-PT	We are at the forefront of technologie opportunities for both working and lea	s in many fields, and there are irning at CERN, including student and	Welcome to our e-Recri	uitment website! We rely on	
Search Vacancies	graduate programmes, as well as vaca electricity mechanics electronics and	ancies in many domains such as	your feedback to contin Problems can be report	ue improving this site. ed by mail.	
Sedicit vacancies		comparing, etc.			
Full Search	Please use the left-hand menu to sear opportunities or look for further inform	ch and apply for our current nation on our recruitment conditions	CERN staff please click menu to access interna	on "Internal posts" in the I vacancy notices (AIS login	
	and programmes, recruitment events	or to contact us.	required)! (more info)		
Recently Published	Cliquez ici pour une explication en	français	Important technical info	ormation.	
By Reference		i u i guis			
Employment Conditions			0		
Information for			Deadlines		
Staff					
Fellows			Technical &	07-AUG-09	
Associates			Doctoral Students	07 650 00	
Students		Contraction of the second s	Fellows	07-SEP-09	
Marie Curie Actions			Scientific Associates	19-MAK-09	
Special Programs			Do not wait until the la application as addition	st day to send your al information will be	
Apprentices			requested by CERN on	ce your application form has	
Your Feedback		RECTIVE 25	been received! More In	10	
Contact Us					7
FAQ		642 HAT	Focus on		
			Are you an undergradu. State nationality in a te practical training period project? CERN has a Te that could interest you. least 18 months of you studies, and your cours	ate student of a CERN Member schnical field looking for a I or a place to do your final chnical Student programme If you have completed at r technical undergraduate e requires a practical trajning	r



CERN options for students

University level (BS/Master)

- Summer student
- OpenLab summer students

Master thesis

Technical student (non physicist)

PhD students

Doctoral students

Young scientists/engineers

- Fellowship
- Other programmes

CERN IT Department CH-1211 Genève 23 Switzerland www.cern.ch/it

Done

34

¥

period of 6 to 12 months, which you wish to spend at CERN, apply to the Technical Student Programme.