





International Collaboration for **Data Preservation** and
Long Term Analysis in High Energy Physics

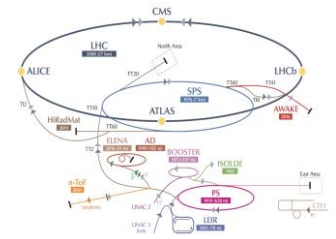
“Big Data” at the LHC (and LEP, and the FCC...)

Big Data & IoT Summit

Jamie.Shiers@cern.ch



Outline



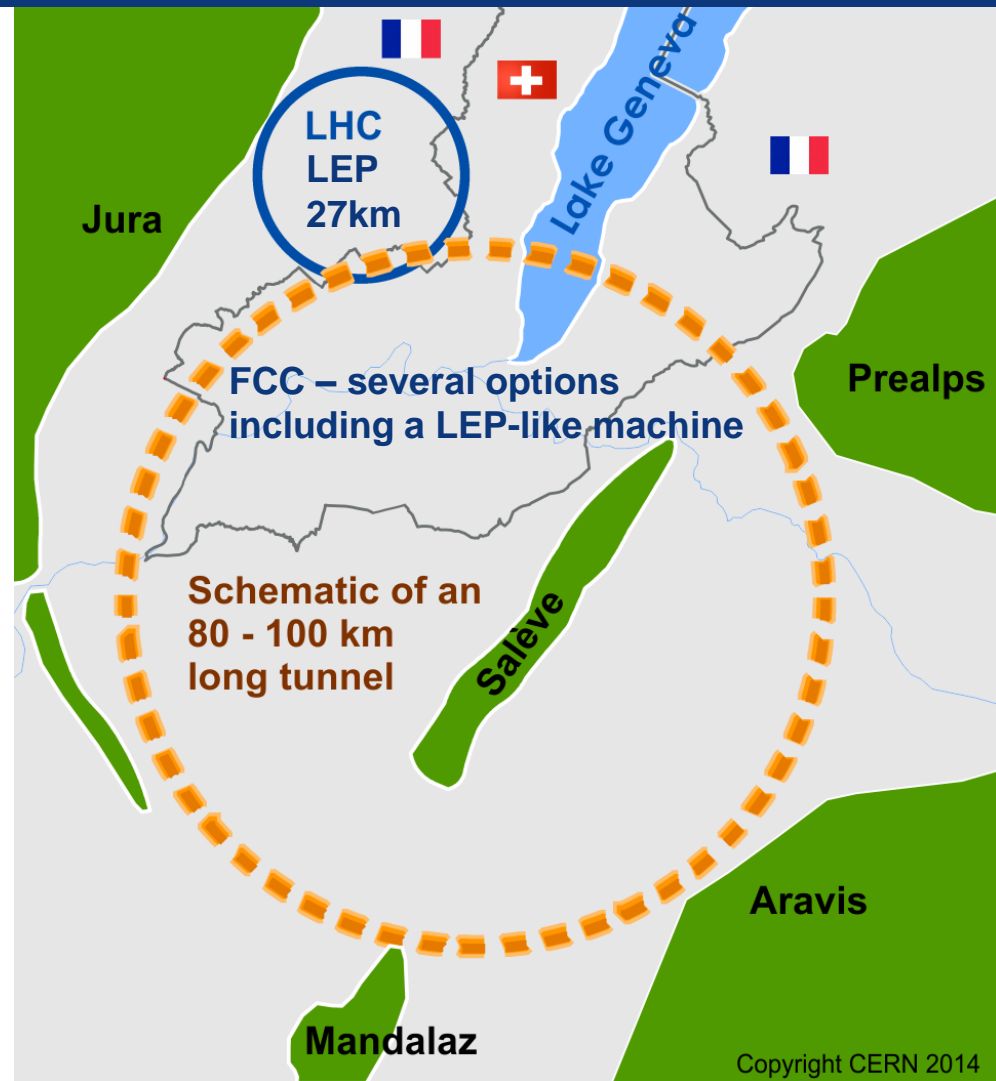
- We currently collect around 50PB of data per year (100 days running period)
- The data volume will rise from ~200PB today to ~10EB by ~2035 (end of LHC data taking)
- **How do we collect, distribute and process such data volumes?**
- **How do we turn “data into discoveries”?**
- **How do we preserve it for future re-use?**

Big Data: From LEP to the LHC to the FCC

From LEP (1989 – 2000) to the LHC (2009 – 2035) to the “FCC”

- “Big data” from hundreds of TB to hundreds of PB to (perhaps) hundreds of EB
- FCC-ee option: “repeat” LEP in just 1 day!
- FCC-hh: 7 times LHC energy, 10^{10} Higgs bosons

What is the “business case” for all this investment?



~30 years of LEP – what does it tell us?

- ▶ Today's “**Big Data**” may become tomorrow's “**peanuts**”
 - ▶ 100TB per LEP experiment: **immensely challenging** at the time; now “trivial” for both CPU and storage
 - ▶ With time, **hardware costs** tend to zero
 - ▶ O(CHF 1000) per experiment per year for archive storage
 - ▶ **Personnel costs** tend to O(1FTE) >> **CHF 1000!**
 - ▶ Perhaps as little now as 0.1 – 0.2 FTE per LEP experiment to keep data + s/w alive – no new analyses included
- ▶ “**Data**” is not just “**bits**”, but also **documentation, software + environment + “knowledge”**
 - ▶ “**Collective knowledge**” particularly hard to capture
 - ▶ Documentation “refreshed” after 20 years (1995) – now in Digital Library in PDF & PDF/A formats (was Postscript)

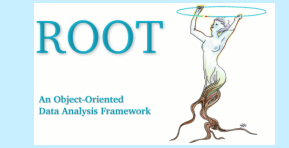
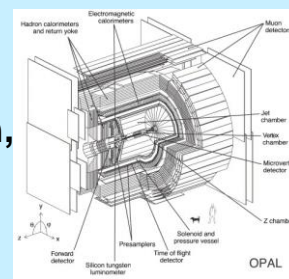
What is HEP data?



Digital information
The data themselves, volume estimates for preservation data of the order of **a few to 10 PB**

Other digital sources such as databases to also be considered

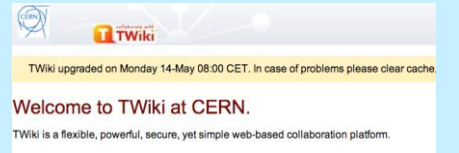
Software Simulation, reconstruction, analysis, user, in addition to any external dependencies



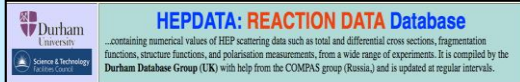
CERNLIB Access

- Access to the CERN Program Library is free of charge to all HEP users worldwide.
- Non-HEP academic and not-for-profit organizations: 1KSF/year

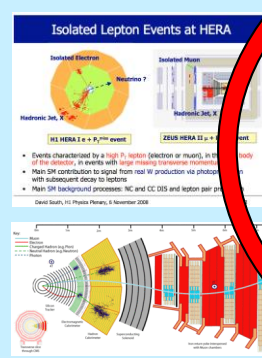
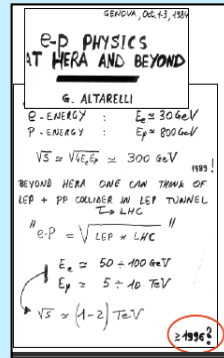
Meta information
Hyper-news, messages, wikis, user forums..



Publications **arXiv.org**



Documentation
Internal publications, notes, manuals, slides



Expertise and people



Nobel Prize in Physics 2013



The Nobel Prize in Physics 2013 was awarded jointly to François Englert and Peter W. Higgs *who had never met until July 2012 at CERN* "for the postulation of a mechanism that contributes to our understanding of the origin of mass of subatomic particles, and which recently was *confirmed through the discovery of the predicted particle, by the ATLAS and CMS experiments at CERN's Large Hadron Collider*".

From ideas of individual theoretical physicists...

...to collective innovation

BROKEN SYMMETRIES AND THE MASSES OF GAUGE BOSONS

Peter W. Higgs

Tait Institute of Mathematical Physics, University of Edinburgh, Edinburgh, Scotland
(Received 31 August 1964)

In a recent note¹ it was shown that the Goldstone theorem² about the "vacuum" solution $\phi(x) = 0$, $\phi(x) = \phi_0$ in theories in which symmetry under certain zero-mass bosons is conserved.

It is shown that the longitudinal modes of the conserved current are not zero, but that the scalar field $\phi(x)$ is not zero over in coupling tends to the relativistic limit. It is shown that the scalar field $\phi(x)$ is not zero over in coupling tends to the relativistic limit. It is shown that the scalar field $\phi(x)$ is not zero over in coupling tends to the relativistic limit.

The simplest case is a gauge theory with Goldstone fields ϕ_1, ϕ_2 and through the Lagrangian

$$L = -\frac{1}{2}(\nabla_\mu \phi_i)^2$$

where

$$\nabla_\mu = \partial_\mu + ig \frac{\tau_a}{2} \phi_a$$

e is a dimensionless constant taken simultaneous kind on ϕ_1, ϕ_2 . Let us suppose spontaneous breaking. Consider the equation treating $\Delta \phi_i, \Delta$ governing the p

*Work supported in part by the U. S. Atomic Energy Commission and in part by the Graduate School from funds supplied by the Wisconsin Alumni Research Foundation.

¹S. Feynman and M. Gell-Mann, Phys. Rev. **109**, 13 (1958).

²T. D. Lee and C. N. Yang, Phys. Rev. **119**, 1410 (1960); S. B. Treiman, Nuovo Cimento **15**, 916 (1960).

³H. Okubo and R. E. Marshak, Nuovo Cimento **25**, 56 (1955); Y. Nambu, Nuovo Cimento **21**, 922 (1963).

⁴Estimates of the rate for $K^+ \rightarrow \pi^+ + e^+ + \nu_e$ due to induced neutral currents have been calculated by several authors. For a list of previous references see Mirza A. Baq Dég, Phys. Rev. **133**, 424 (1963).

⁵M. Baker and S. Glashow, Nuovo Cimento **25**, 857

(1962). They predict a branching ratio for decay mode (1) of $\sim 10^{-6}$.

⁶N. P. Samios, Phys. Rev. **121**, 275 (1961).

⁷The best previously reported estimate comes from the limit on $K^+ \rightarrow \pi^+ + \mu^+ + \nu_\mu$. The 90% confidence level is $|g_{\mu\mu}|^2 < 10^{-2} |g_{ee}|^2$; M. Bartos, K. Lande, L. M. Lederman, and William Chinowsky, Ann. Phys. (N.Y.) **5**, 156 (1958).

The absence of the decay mode $\mu^+ \rightarrow e^+ + \nu_e + \nu_\mu$ is not a good test for the existence of neutral currents since this decay mode may be absolutely forbidden by conservation of mass number. G. Feenberg and L. M. Lederman, Ann. Rev. Nucl. Sci. **12**, 445 (1963).

⁸S. N. Biswas and S. K. Bose, Phys. Rev. Letters **12**, 176 (1964).

BROKEN SYMMETRY AND THE MASS OF GAUGE VECTOR MESONS*

F. Englert and R. Brout

Faculté des Sciences, Université Libre de Bruxelles, Bruxelles, Belgium
(Received 26 June 1964)

It is of interest to inquire whether gauge vector mesons acquire mass through interaction; by a gauge vector meson we mean a Yang-Mills field¹ associated with the extension of a Lie group from global to local symmetry. The importance of this problem resides in the possibility that strong-interaction physics originates from massive gauge fields related to a system of conserved currents.² In this note, we shall show that in certain cases vector mesons do indeed acquire mass when the vacuum is degenerate with respect to a compact Lie group.

Theories with degenerate vacuum (broken symmetry) have been the subject of intensive study since their inception by Nambu.^{3,4} A characteristic feature of such theories is the possible existence of zero-mass bosons which tend to restore the symmetry.^{5,6} We shall show that it is precisely these singularities which maintain the gauge invariance of the theory, despite the fact that the vector meson acquires mass.

We shall first treat the case where the original fields are a set of bosons ϕ_A which transform as a basis for a representation of a compact Lie group. This example should be considered as a rather general phenomenological model. As such, we shall not study the particular mechanism by which the symmetry is broken but simply assume that such a mechanism exists. A calculation performed in lowest order perturbation theory indicates that

these vector mesons which are coupled to currents that "rotate" the original vacuum are the ones which acquire mass [see Eq. (6)].

We shall then examine a particular model based on chirality invariance which may have a more fundamental significance. Here we begin with a chirality-invariant Lagrangian and introduce both vector and pseudovector gauge fields, thereby guaranteeing invariance under both local phase and local γ_5 -phase transformations. In this model the gauge fields themselves may break the γ_5 invariance leading to a mass for the original Fermi field. We shall show in this case that the pseudovector field acquires mass.

In the last paragraph we sketch a simple argument which renders these results reasonable.

(1) Lest the simplicity of the argument be shrouded in a cloud of indices, we first consider a one-parameter Abelian group, representing, for example, the phase transformation of a charged boson; we then present the generalization to an arbitrary compact Lie group. The interaction between the ϕ and the A_μ fields is

$$H_{int} = ie A_\mu \phi^* \nabla^\mu \phi - e^2 \phi^* \phi A_\mu A^\mu \quad (1)$$

where $\phi = (\phi_1 + i\phi_2)/\sqrt{2}$. We shall break the symmetry by fixing $\langle \phi \rangle \neq 0$ in the vacuum, with the phase chosen for convenience such that $\langle \phi \rangle = \langle \phi^* \rangle = \langle \phi \rangle / \sqrt{2}$.

We shall assume that the application of the



WLCG

1964

2012



The Worldwide LHC Computing

Grid

October 2016:

-63 MoU's

-167 sites; 42 countries

-Tier0, Tier1s & Tier2s

-O(1), O(10), O(100)

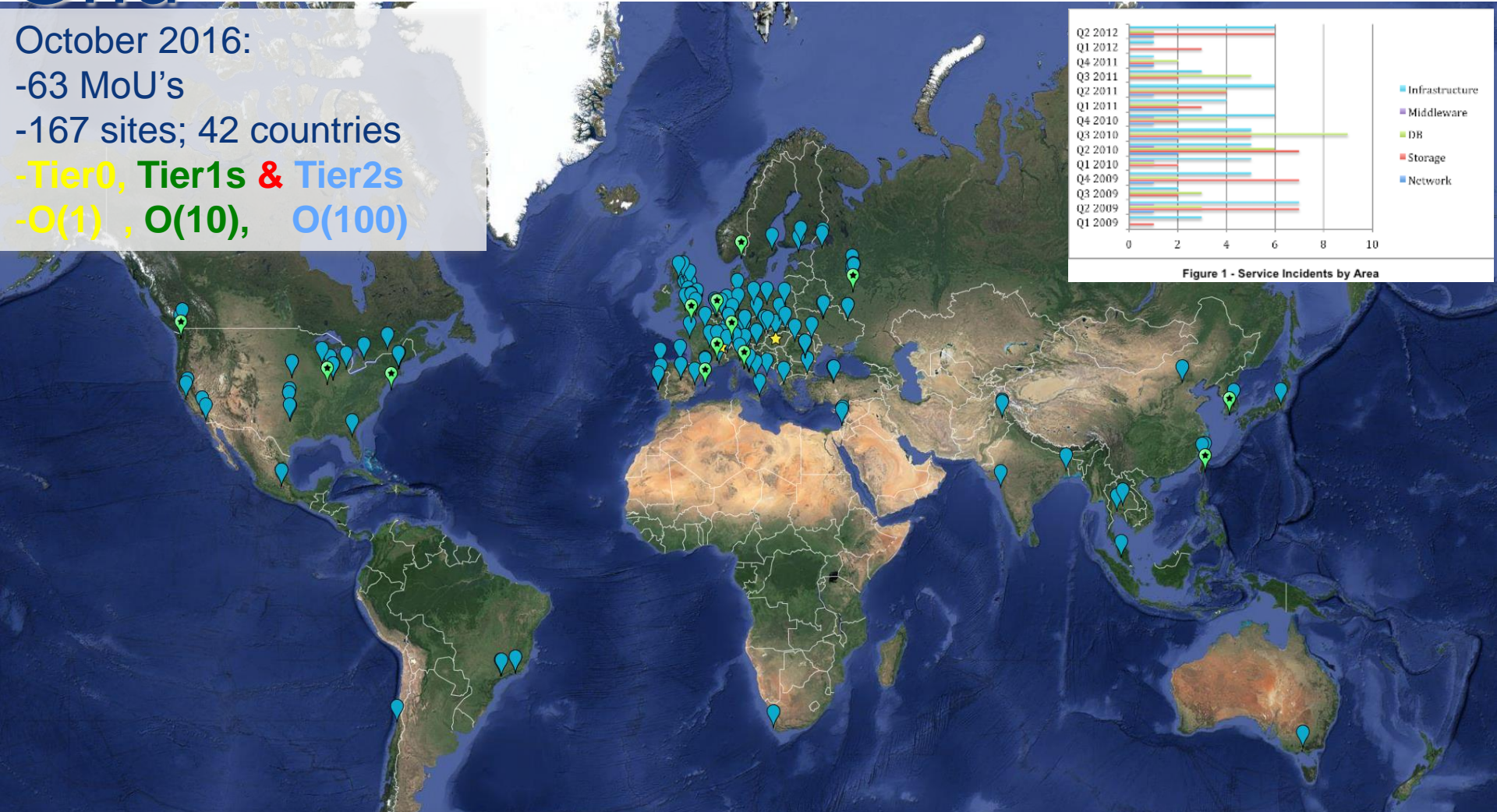


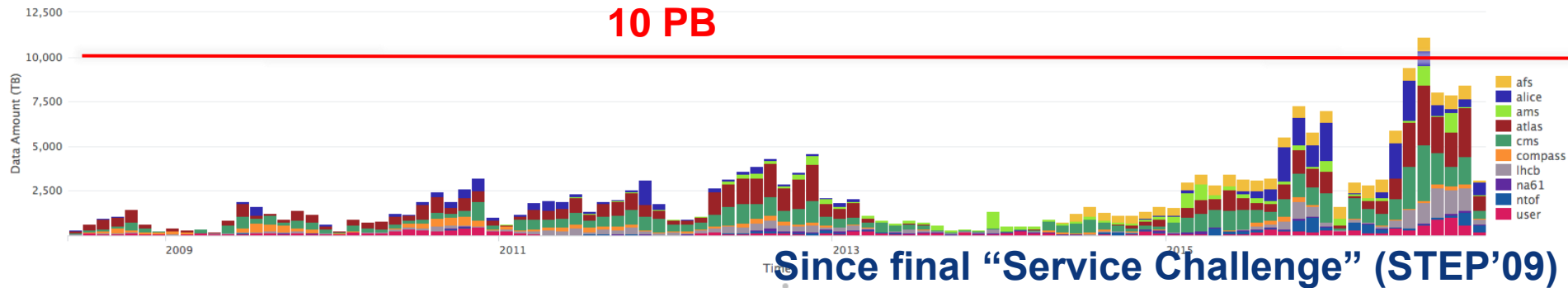
Figure 1 - Service Incidents by Area

- CPU: 3.8 M HepSpec06
 - If today's fastest cores: ~ 350,000 cores
- Disk 310 PB
- Tape 390 PB

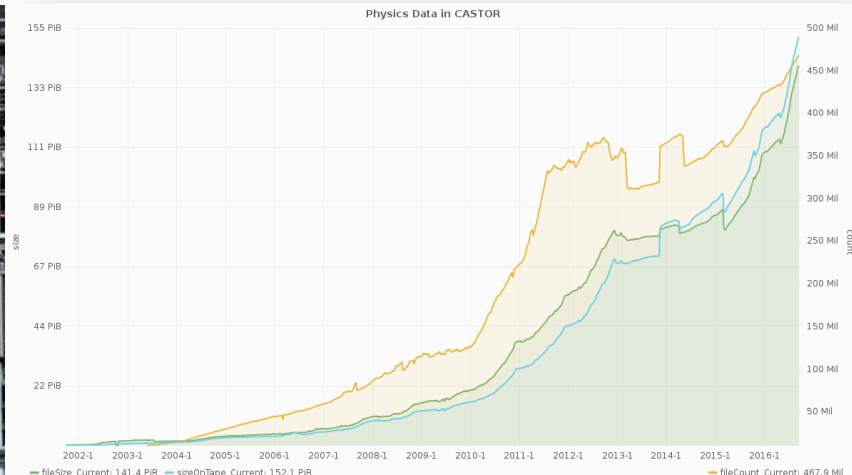
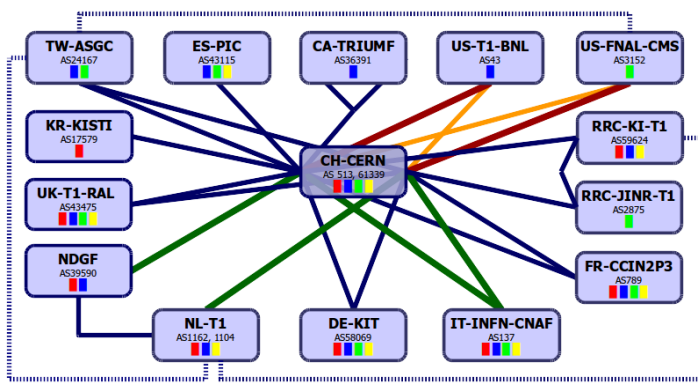
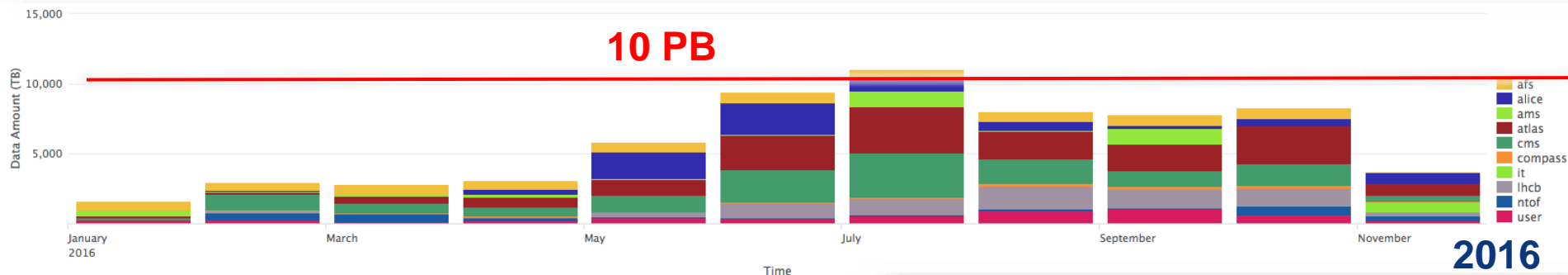
Running jobs: 441,353
Active cores: 630,003
Transfer rate: 35.32 GiB/sec

Data transfers / acquisition

Transferred Data Amount per Virtual Organization for WRITE Requests



Transferred Data Amount per Virtual Organization for WRITE Requests



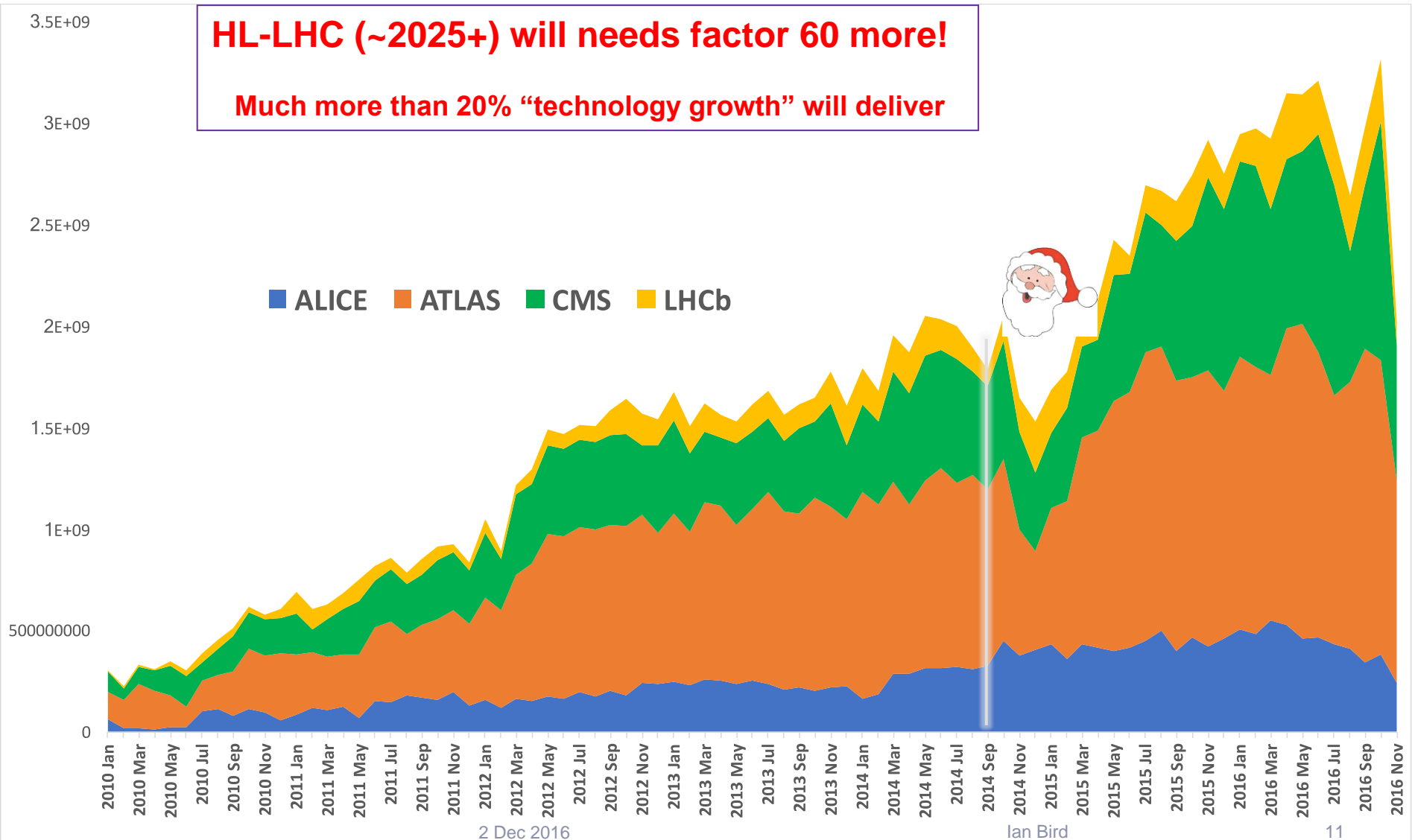
- TO-T1 and T1-T1 traffic
 - T1-T1 traffic only
 - Alice - Atlas - CMS - LHCb
 - edoardo.martelli@cern.ch 20160912

CPU delivered

HL-LHC (~2025+) will need factor 60 more!

Much more than 20% "technology growth" will deliver

■ ALICE ■ ATLAS ■ CMS ■ LHCb



Why Build a Grid? (and not Cloud*)

- Much R&D into Computing for the LHC was done in the mid-late 1990s
 - LEP had already moved to distributed computing; Unix was the main O/S; Intel + Linux not then dominant
- **In 2000, decisions (and funding) needed**
- Several rounds of EU-funded projects
 - EDG, EGEE I, II, III, EGI, ... + others elsewhere
- **WLCG Service Challenges to “harden” Grid**

WLCG Service Challenges

- As much about people and collaboration as about technology
- **Getting people to provide a 24 x 7 service for a machine on the other side of the planet for no clear reason was going to be hard!**
- Regional workshops – both motivational as well as technical – plus daily Operations Calls
- In a grid, **something** is broken all of the time!
- Clear KPIs, “critical services” & response targets: **measurable improvement in service quality despite ever increasing demands**





Distributed Computing = Distributed Spending

**Much more attractive to funding agencies;
Many other benefits to Universities and Institutes.**



“Higgs Discovery Day”

- To find the Higgs, you need the accelerator, the experiments and *the Grid*
- Rolf Dieter Heuer, CERN DG

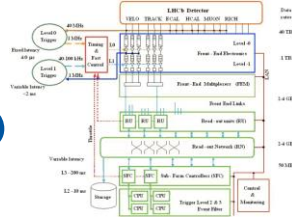


The first time “computing” mentioned at the same level as machine & experiments

“Data” Preservation in HEP

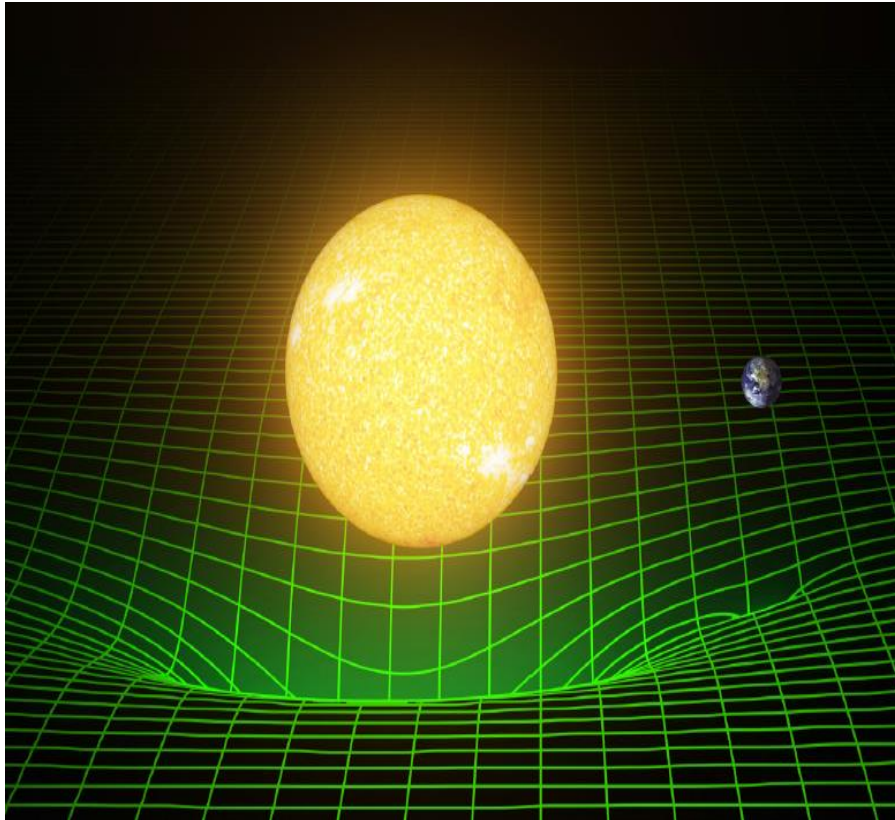
- The data from the world’s particle accelerators and colliders (HEP data) is both **costly** and **time consuming** to produce
 - HEP data contains a wealth of **scientific potential**, plus high value for **educational outreach**.
 - Many data samples **are unique**, it is essential to preserve not only the data but also the full capability to reproduce past analyses and perform new ones.
- **This means preserving data, documentation, software and "knowledge".**

What Makes HEP Different?

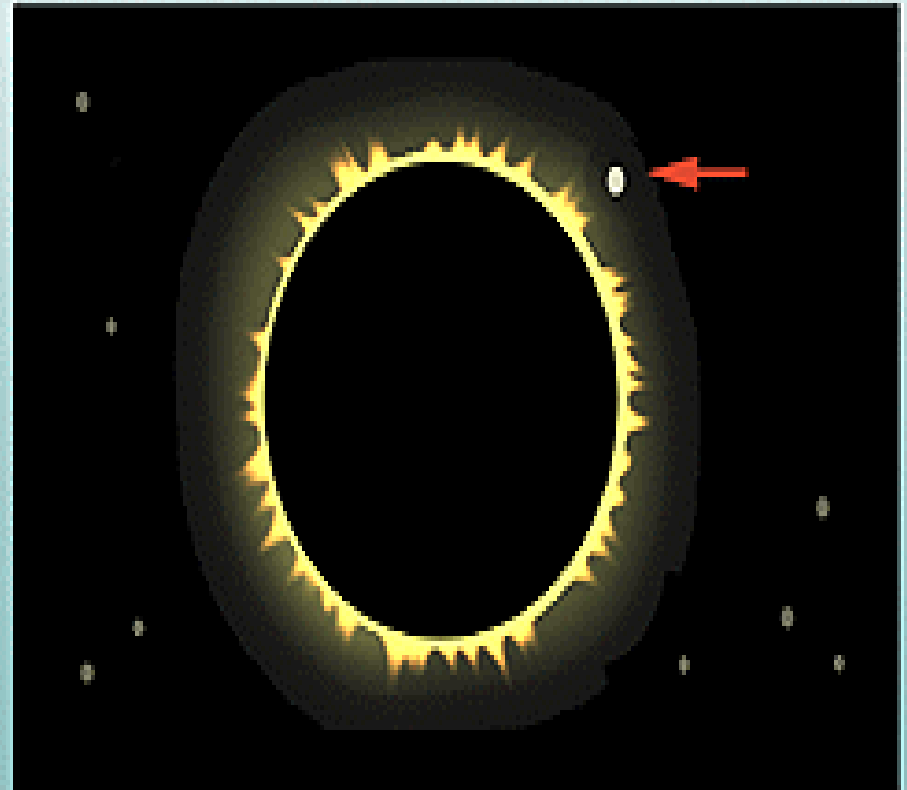


- We **throw away** most of our data before it is even recorded – “triggers”
- Our detectors are **relatively stable** over long periods of time (years) – not “doubling every 6 or 18 months”
- We make **“measurements”** – not **“observations”**
- Our projects typically last for **decades** – we **need** to keep data usable during at least this length of time
- We have **shared** “data behind publications” for more than 30 years... (HEPData)

An OBSERVATION...



BENDING LIGHT



1.3 Billion Years Ago

And another... (Black holes merging...)



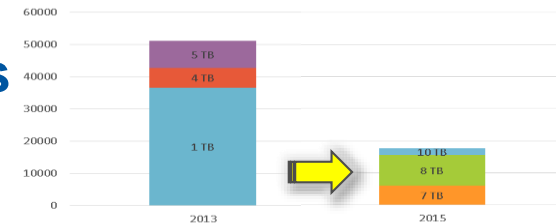
CERN Services for LTDP

- 1.State-of-the art "**bit preservation**", implementing practices that conform to the ISO 16363 standard
 - 2."**Software preservation**" - a key challenge in HEP where the software stacks are both large and complex (and dynamic)
 - 3.Analysis **capture and preservation**, corresponding to a set of agreed Use Cases
 - 4.Access to **data behind physics publications** - the HEPData portal
 - 5.An **Open Data portal** for released subsets of the (currently) LHC data
 - 6.A **DPHEP portal** that links also to data preservation efforts at other HEP institutes worldwide.
- **These run in production at CERN and elsewhere and are being prototyped (in generic equivalents) in the EOSC Pilot**



Bit Preservation: Steps Include

- Controlled media **lifecycle**
 - **Media kept for 2 max. 2 drive generations**
 - Regular media **verification**
 - When tape written, filled, every 2 years...
 - **Reducing** tape mounts
 - Reduces media wear-out & increases efficiency
 - Data **Redundancy**
 - For “smaller” communities, a 2nd copy can be created: separate library in a different building (e.g. LEP – **3 copies at CERN!**)
 - **Protecting** the physical link
 - Between disk caches and tape servers
 - Protecting the **environment**
 - Dust sensors! (Don't let users touch tapes)



Constant improvement: reduction in bit-loss rate: 5×10^{-16}

Collaboration with others

1. The elaboration of a clear "**business case**" for long-term data preservation
2. The development of an associated "**cost model**"
3. A common view of the **Use Cases** driving the need for data preservation
4. Understanding how to address Funding Agencies requirements for **Data Management Plans**
5. Preparing for **Certification** of HEP digital repositories and their long-term future.

How Much Data?

- **100TB** / LEP experiment: **3 copies** at CERN
- **1-10PB** for experiments at the HERA collider at DESY, the TEVATRON at Fermilab or the BaBar experiment at SLAC.
- The LHC experiments is already in the multi-hundred PB range (**x00PB**)
- **10EB** or more including the High Luminosity upgrade of the LHC (HL-LHC)
- At least 10 times more at FCC (**100EB-1ZB**)

The Business Case

- **For Data Preservation:**

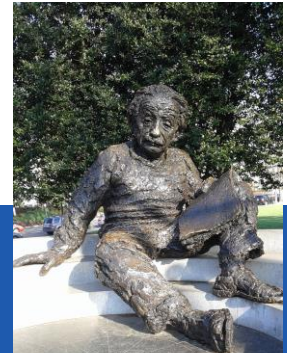
- Data continues to be analysed well after end of data taking: papers continue to be written, PhDs awarded and (sometimes) new discoveries
- ~10% of the scientific output for (<) <1 per mil of the cost

- **For LEP / LHC / FCC:**

- Studies (e.g. STFC, OECD) on “value” of CERN and other labs show ~“cost neutral” based on scientific output
 - Using accepted value of PhDs etc – no spin-offs

- **Including spin-offs (advances in superconductivity, distributed computing, physics for medicine etc.) factor of 10 – 40 ROI!**

- **Unforeseen benefits, e.g. Michelson & Morley experiment to “find ether” led indirectly to Special Relativity;**
- **Theory of “stimulated emission” eventually led to lasers – multi \$BN industry today**



Hardware costs can be significant initially but tend to zero



LTDP Conclusions

- As is well known, Data Preservation is a **Journey** and not a **destination**.
- Can we capture sufficient “**knowledge**” to keep the data usable **beyond** the lifetime of the original collaboration?
- Can we prepare for **major migrations**, similar to those that happened in the past? (Or will x86 and Linux last “forever”)
- For the HL-LHC, we may have **neither** the storage resources to keep all (intermediate) data, **nor** the computational resources to re-compute them!
- You can't **share** or re-use data, nor **reproduce** results, if you haven't first preserved it (data, software, documentation, knowledge)

80 years of “Big Data” ...



- 40 years from first **LEP** + **LHC** proposals
- ~40 years to start of **FCC** (perhaps)...
- **100 years of CERN in 2054!**

- Many studies for LEP, LHC and now FCC
 - Predictions have generally been (wildly) wrong
 - Many things – e.g, networks – have been far better than predicted, e.g. LHC OPN
 - LHC “availability” twice that of LEP!
 - (Much) **more** with the **same** (budget) or **less** (staff) – 3 orders of magnitude in scale between projects!



Possible Questions

1. Projects like the LHC, the Square Kilometre Array etc cost a significant amount of money. What steps are you taking to collaborate to reduce overall costs?
2. Petabytes and exabytes cost real money to keep. How do you decide what to keep and what to throw away? Or recalculate?
3. How can you involve – or benefit – industry in what you do?