# Investing in Curation
## A **Shared** Path to Sustainability

[Jamie.Shiers@cern.ch](mailto:Jamie.Shiers@cern.ch)

Data Preservation in HEP (DPHEP)

# Outline

1. Pick 2 of the messages from the roadmap & comment
   - I could comment on all – but not in 10'

1. What (+ve) impact has 4C already had on us?
   - Avoiding overlap with the above
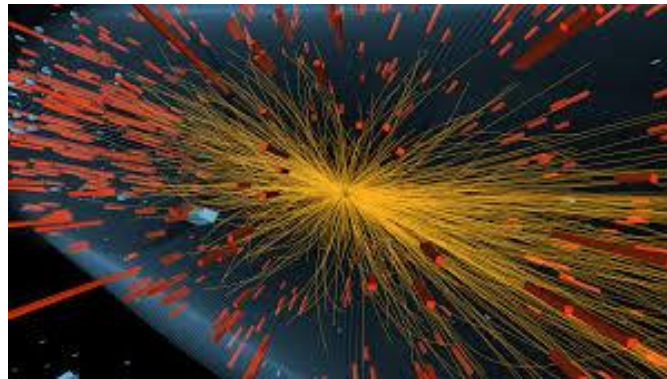
1. A point for discussion – shared responsiblity / action

# THE MESSAGES

# The 4C Roadmap Messages

1. Identify the value of digital assets and make choices

2. Demand and choose more efficient systems

3. Develop scalable services and infrastructure

4. Design digital curation as a sustainable service

5. Make funding dependent on costing digital assets across the whole lifecycle

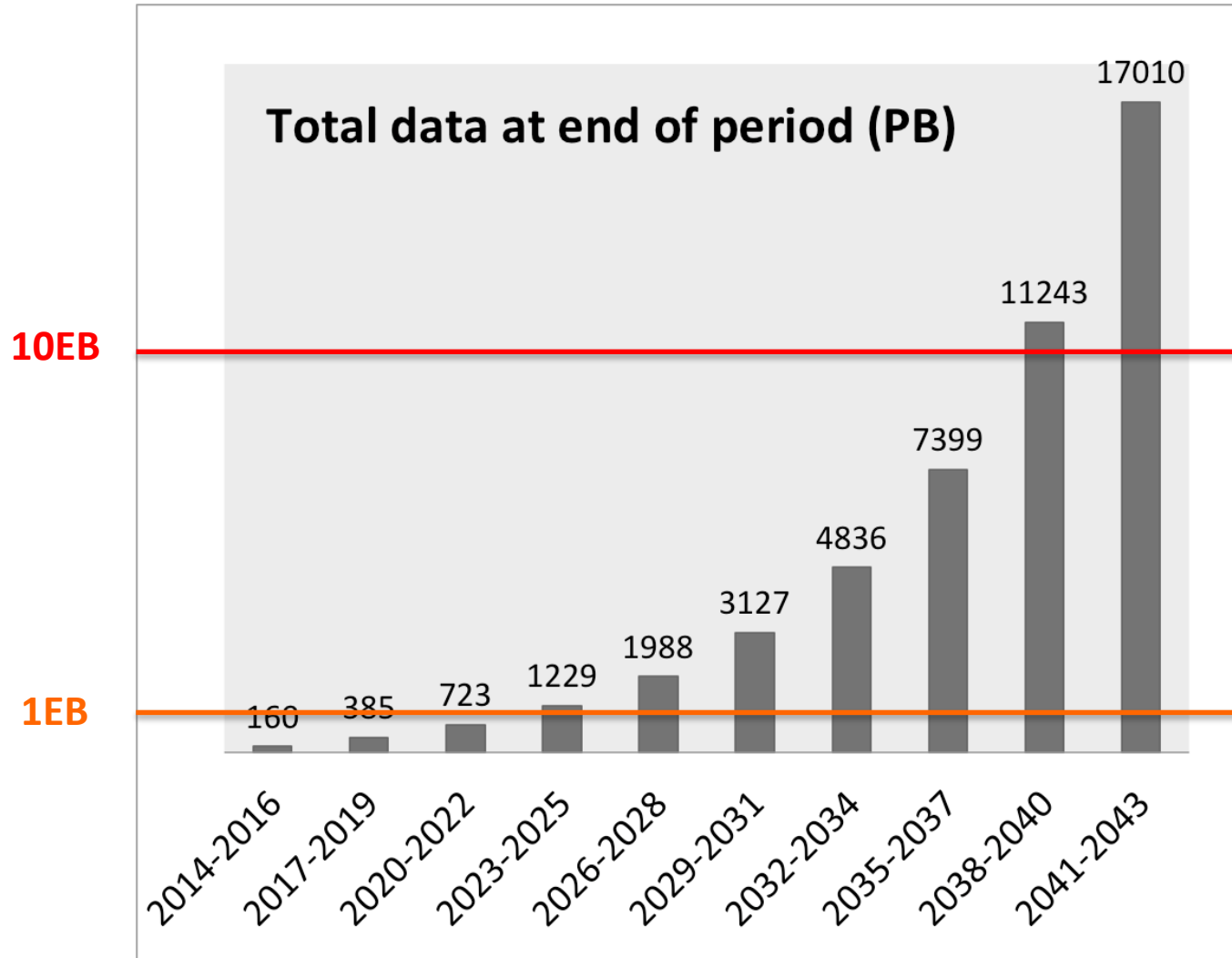6. Be collaborative and transparent to drive down costs
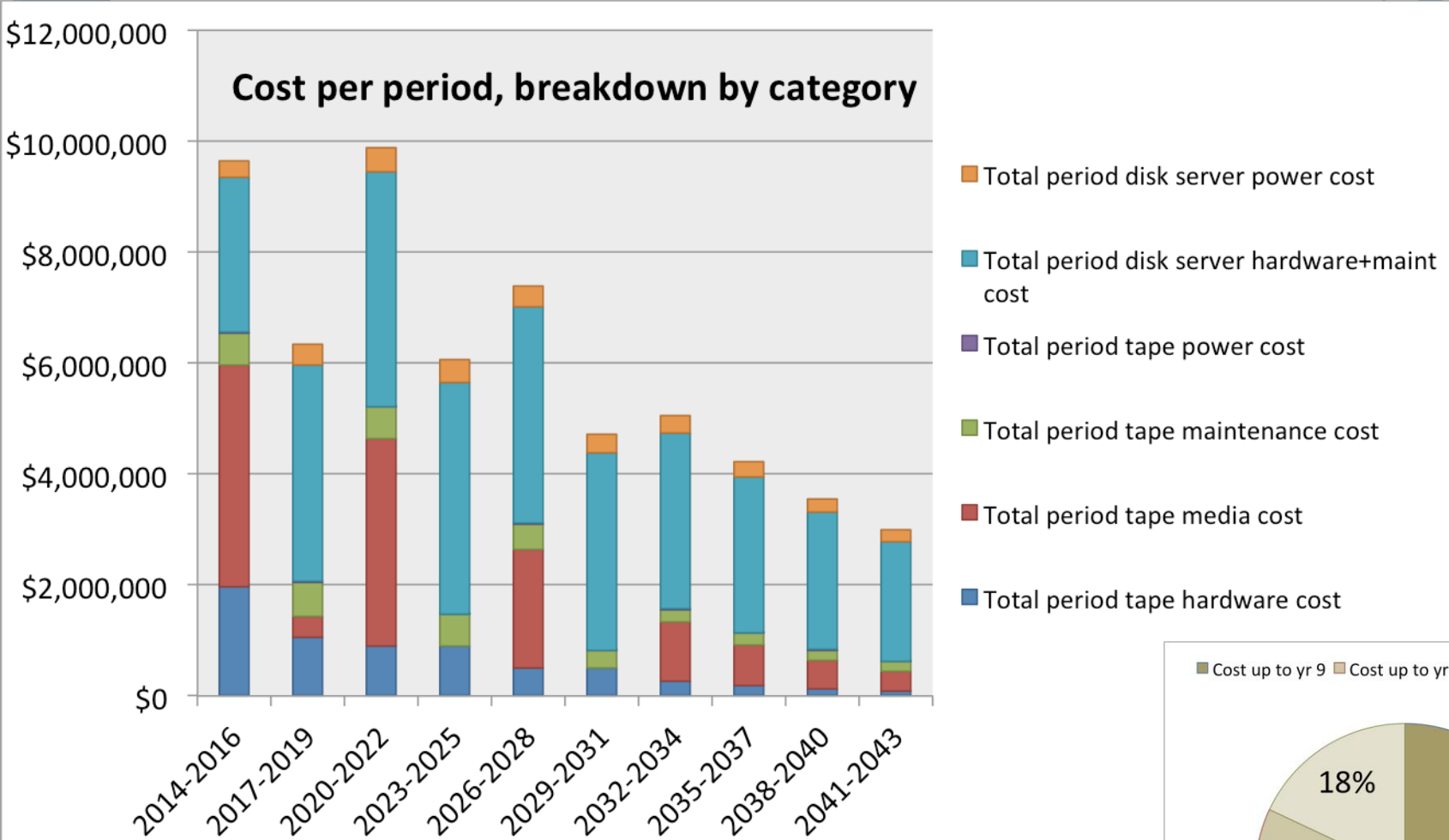
# IMPACT OF 4C ON DPHEP



International Collaboration for Data Preservation and
Long Term Analysis in High Energy Physics
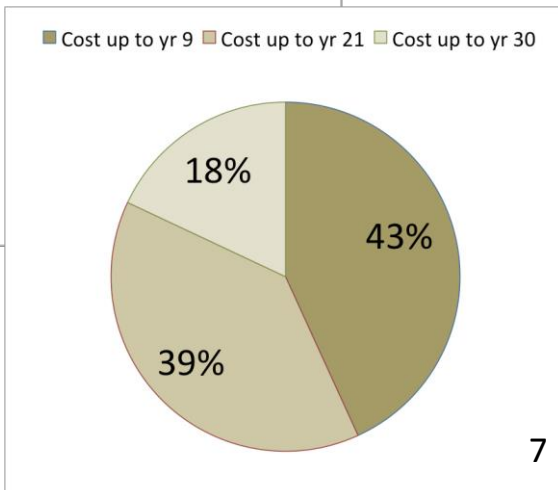
# "LHC Cost Model" (simplified)

Start with 10PB, then +50PB/year, then +50% every 3y (or +15% / year)



**Total data at end of period (PB)**

| Period | Total (PB) |
|---|---|
| 2014-2016 | 160 |
| 2017-2019 | 385 |
| 2020-2022 | 723 |
| 2023-2025 | 1229 |
| 2026-2028 | 1988 |
| 2029-2031 | 3127 |
| 2032-2034 | 4836 |
| 2035-2037 | 7399 |
| 2038-2040 | 11243 |
| 2041-2043 | 17010 |

10EB

1EB

# Case B) increasing archive growth



Cost per period, breakdown by category

- Total period disk server power cost
- Total period disk server hardware+maint cost
- Total period tape power cost
- Total period tape maintenance cost
- Total period tape media cost
- Total period tape hardware cost

Total cost: ~$59.9M
(~$2M / year)

Cost up to yr 9   Cost up to yr 21   Cost up to yr 30

43%
39%
18%

Large Hadron Collider (LHC)

Scientists accelerate two beams of protons around the 17-mile ring, smashing them together at 186,000 miles per second.

Ultimately, scientists hope to find in the collisions proof of the "God particle", the Higgs boson, which is thought to give mass to matter.

# 1. Identify the value of digital assets and make choices

- Today, significant volumes of HEP data are thrown away "at birth" – i.e. via very strict filters (aka triggers) B4 writing to storage

> **To 1ˢᵗ approximation ALL remaining data needs to be kept for a few decades**

- "Value" can be measured in a number of ways:
  - Scientific publications / results;
  - Educational / cultural impact;
  - "Spin-offs" – e.g. superconductivity, ICT, vacuum technology.

# Why build an LHC?



## THE STANDARD MODEL



**BEFORE!**

# 1 – Long Tail of Papers



Nov. 2000:
End of Data Taking

# 2 – New Theoretical Insights



# 3 – "Discovery" to "Precision"



**possible long-term time line**

Alain Blondel TLEP design study r-ECFA 2013-07-20

# Use Case Summary

1. Keep data usable for ~1 decade

2. Keep data usable for ~2 decades

3. Keep data usable for ~3 decades

**Volume: 100PB + ~50PB/year (+400PB/year from 2020)**

# Balance sheet – Tevatron@FNAL

- 20 year investment in Tevatron          **~ $4B**
- Students                    $4B
- Magnets and MRI          $5-10B   **}   ~ $50B total**
- **Computing                    $40B**

*Very rough calculation – but confirms our gut feeling that investment in fundamental science pays off*

I think there is an opportunity for someone to repeat this exercise more rigorously

cf. STFC study of SRS Impact

http://www.stfc.ac.uk/2428.aspx

Science & Technology
Facilities Council

# 2. Demand and choose more efficient systems

- From my point of view, we can divide the services into two areas:

    - **Generic, potentially "shared" infrastructure**
    - **Discipline oriented services, addressing specific Use Cases of communities, using common components (VREs)**

- "Bit preservation" falls into the former – and is harder than many people thing – particularly at scale (PB, EB, ZB, YB, …)

Suppose these guys can build / share the most cost effective, scalable and reliable federated storage services, e.g. for peta- / exa- / zetta- scale bit preservation?
Can we ignore them?

The answer today is often yes – Something I am convinced we need to overcome, if we really care about:
1. Efficiency;
2. Sustainability;
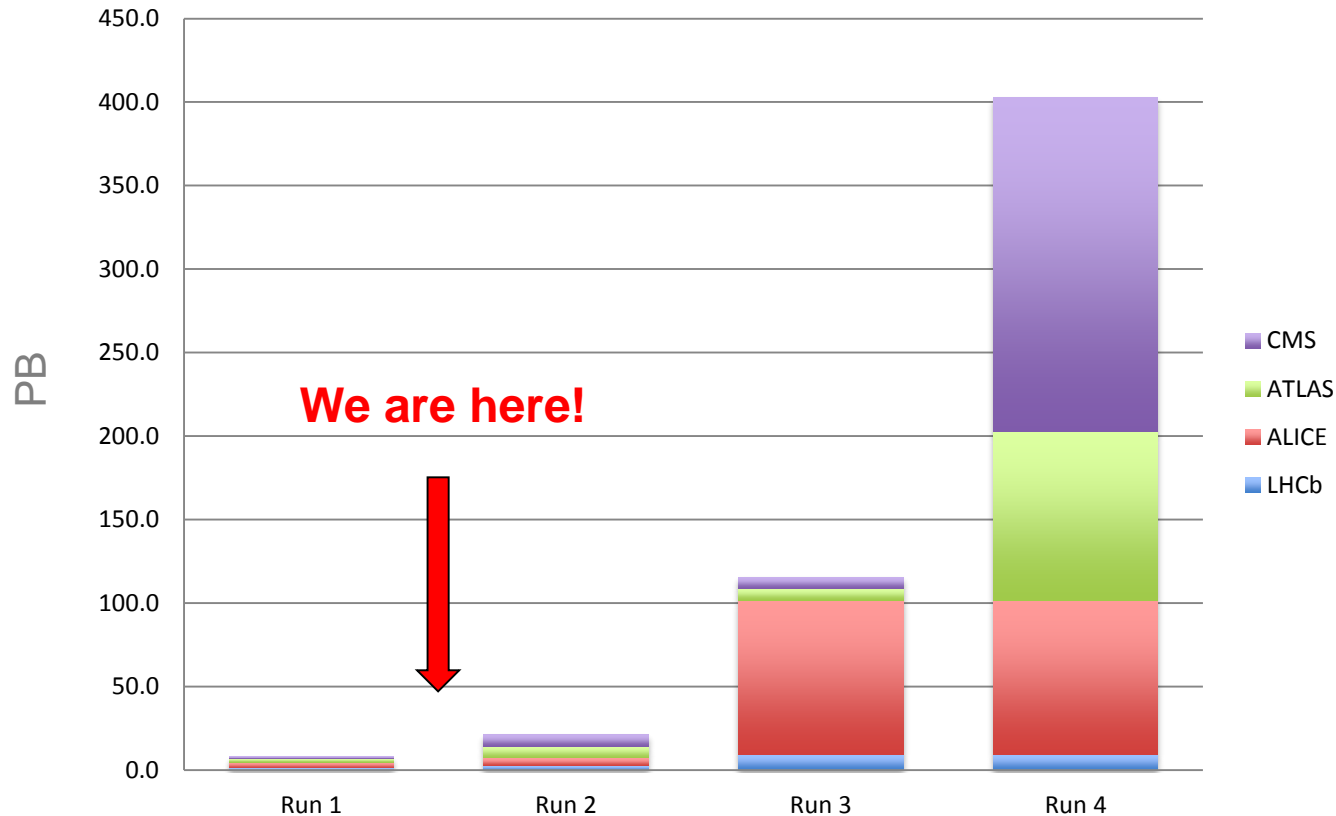3. Cost effectiveness

# 3. Develop scalable services and infrastructure

- What do we mean by "scalable"?

➢ **Linear scaling with "capacity" is not a realistic option: increased "capacity" with constant {budget; effort} over a long period (decades) is a more "realistic" target**

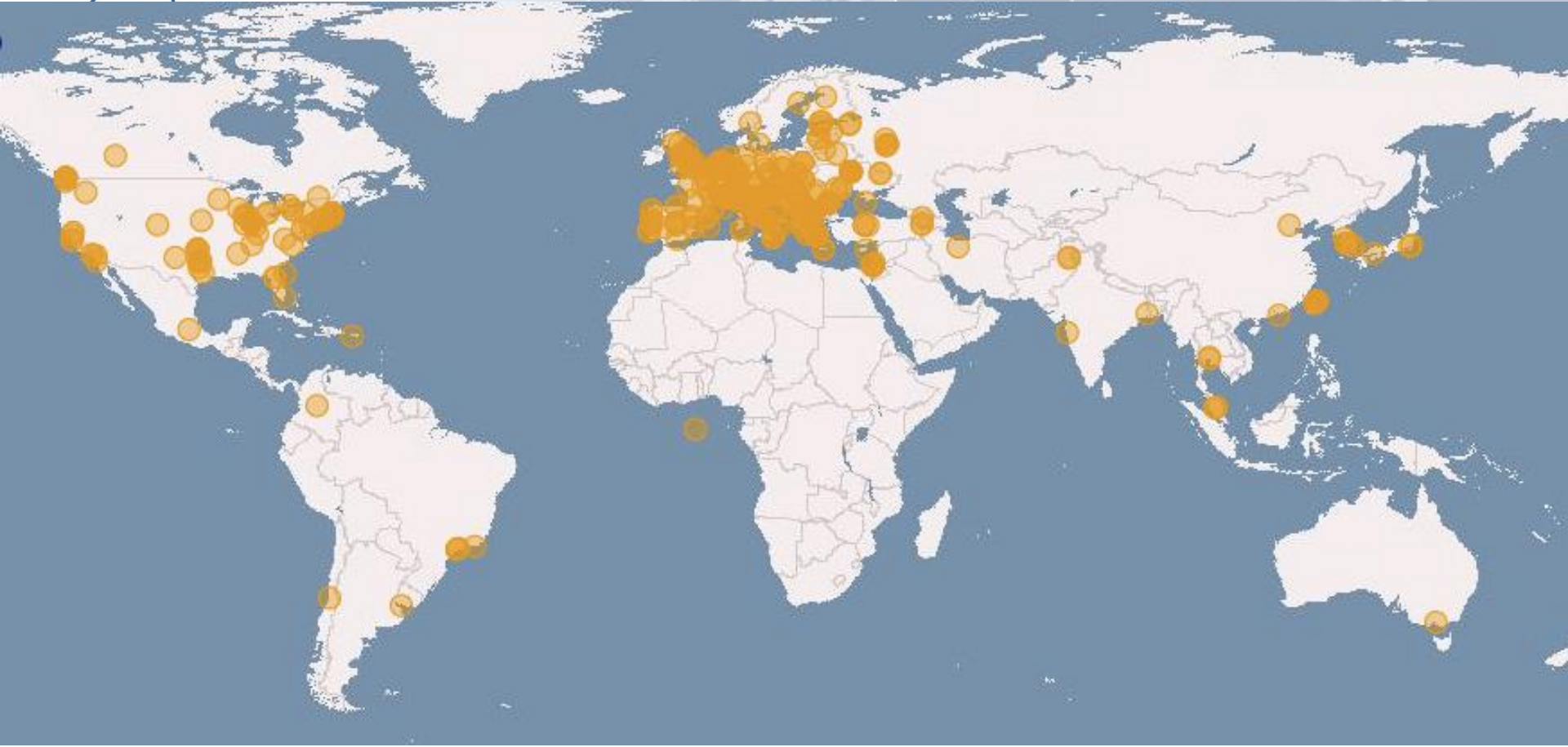- And – the WLCG experience shows – constant service improvement is also possible

# Data: Outlook for HL-LHC



- Very rough estimate of a new RAW data per year of running using a simple extrapolation of current data volume scaled by the output rates.
  - To be added: derived data (ESD, AOD), simulation, user data…
- ➤ **At least 0.5 EB / year (x 10 years of data taking)**
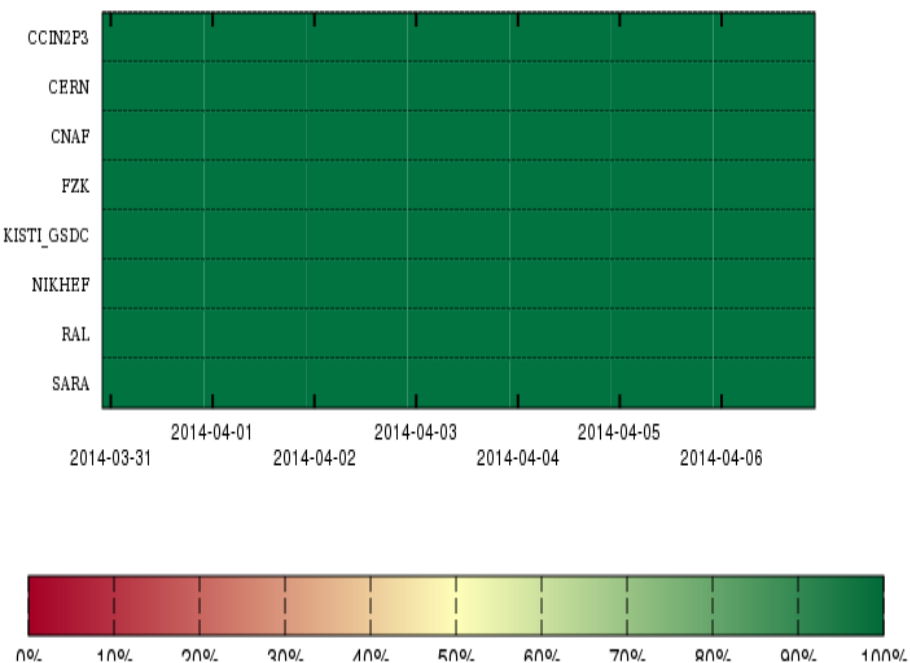
# WLCG Collaboration Today



- Distributed infrastructure of 150 computing centers in 40 countries
- 300+ k CPU cores (~ 2M HEP-SPEC-06)
- The biggest site with ~50k CPU cores, 12 T2 with 2-30k CPU cores
- Distributed data, services and operation infrastructure

# WLCG Collaboration Tomorrow



- How will this evolve to HL-LHC needs?
- To what extent is it applicable to other comparable scale projects?
- **Already evolving, most significantly during Long Shutdowns, but also during data taking!**
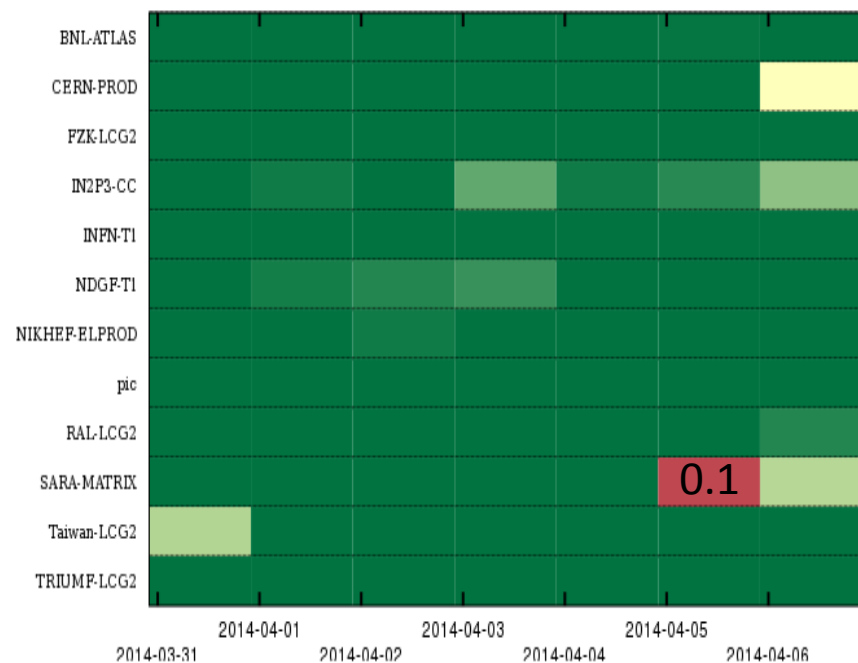
## Site reliability using ALICE_CRITICAL
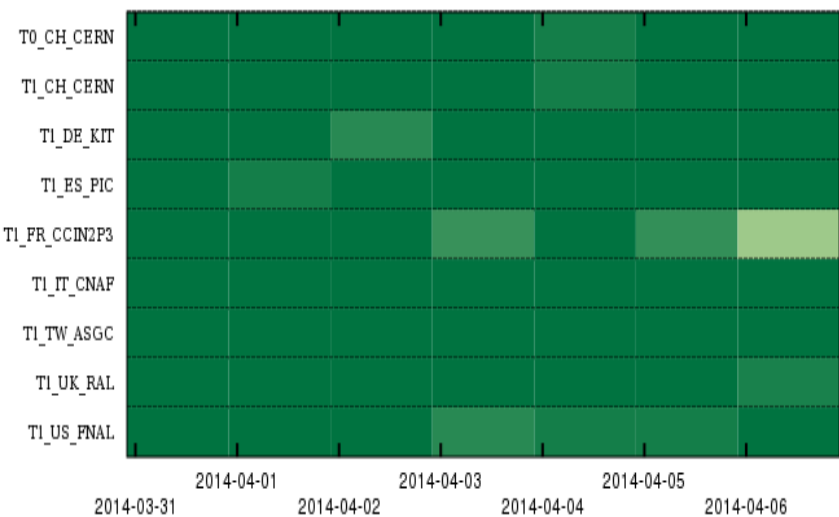168 hours from 2014-03-31 00:00 to 2014-04-07 00:00

## Site reliability using ATLAS_CRITICAL
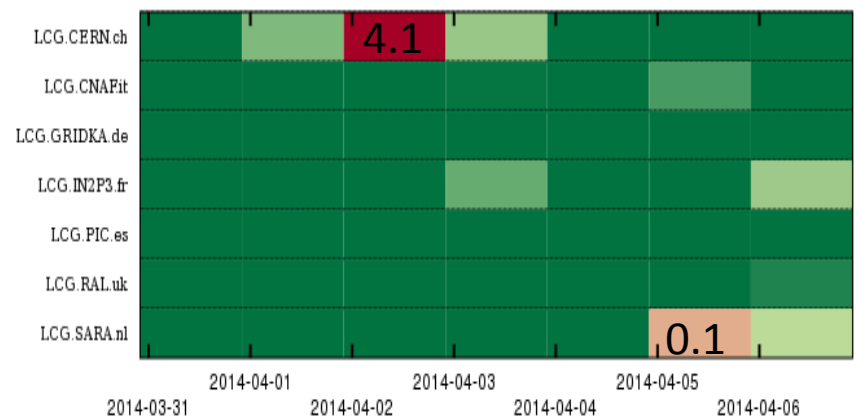168 hours from 2014-03-31 00:00 to 2014-04-07 00:00

## Site reliability using CMS_CRITICAL_FULL
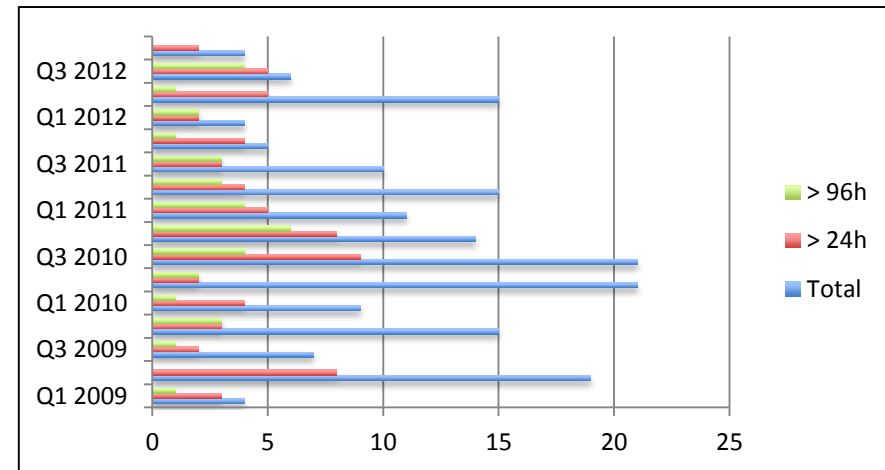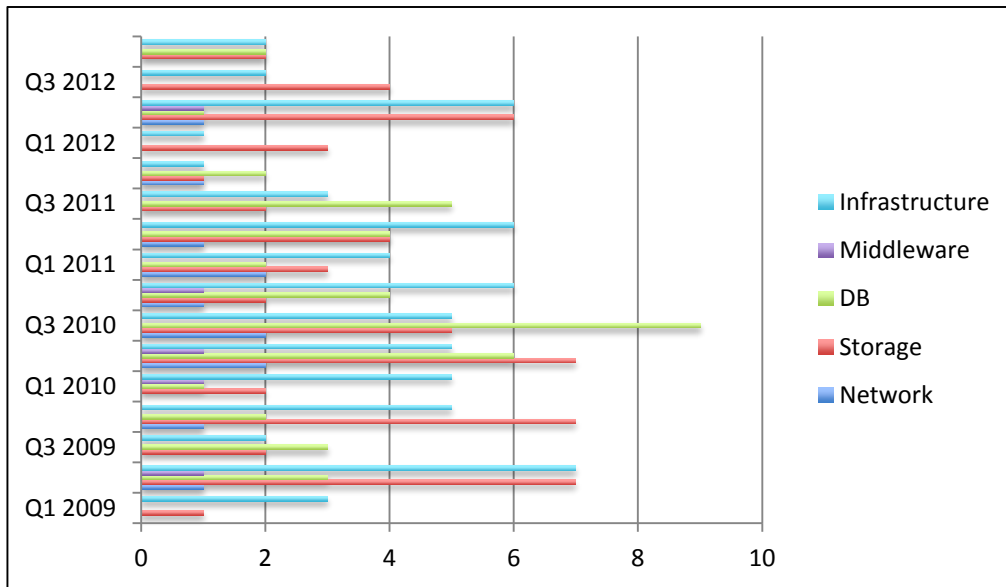168 hours from 2014-03-31 00:00 to 2014-04-07 00:00

## Site reliability using LHCb_CRITICAL
168 hours from 2014-03-31 00:00 to 2014-04-07 00:00

# WLCG Service Incidents

- Aka "post-mortems"

# Resolution of Incidents



**Incidents**

**Data taking**

Legend:
- > 96h
- > 24h
- Total

# 4. Design digital curation as a sustainable service

# 4. Design digital curation as a sustainable service



The last years have seen the end of several experiments

LEP, 2 November 2000

HERA, 30 June 2007

PEP-II, 7 April 2008

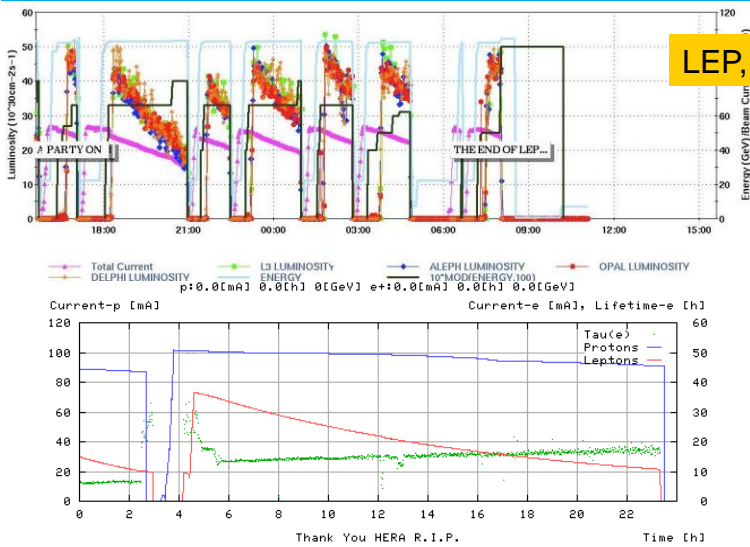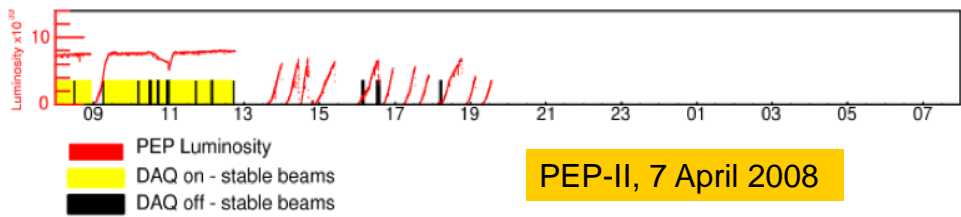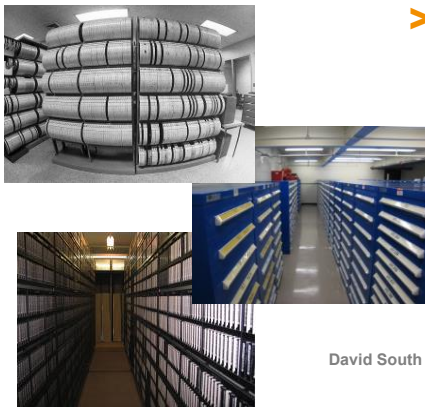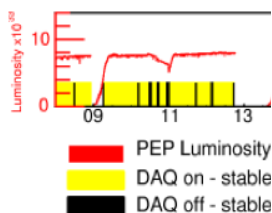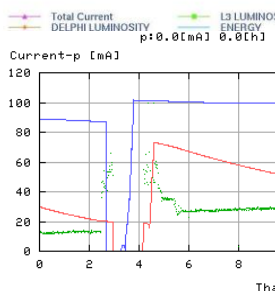# 4. Design digital curation as a sustainable service

## The last years have seen the end of several experiments

### After the collisions have stopped

> Finish the analyses! But then what do you do with the data?

  ▪ Until recently, there was no clear policy on this in the HEP community

  ▪ It's possible that older HEP experiments have in fact simply lost the data

> Data preservation, including long term access, is generally not part of the planning, software design or budget of an experiment

  ▪ So far, HEP data preservation initiatives have been in the main not planned by the original collaborations, but rather the effort a few knowledgeable people

> The conservation of tapes is not equivalent to data preservation!

  ▪ *"We cannot ensure data is stored in file formats appropriate for long term preservation"*

  ▪ *"The software for exploiting the data is under the control of the experiments"*

  ▪ *"We are sure most of the data are not easily accessible!"*

# 4. Design digital curation as a sustainable service

CERN-Council-S/106
Original: English
7 May 2013

ORGANISATION EUROPEENNE POUR LA RECHERCHE NUCLEAIRE

**CERN** EUROPEAN ORGANIZATION FOR NUCLEAR RESEARCH

| Action to be taken | | Voting Procedure |
|---|---|---|
| For Approval | **EUROPEAN STRATEGY SESSION OF COUNCIL** 16th Session - 30 May 2013 **European Commission Berlaymont Building - Brussels** | Simple Majority of Member States represented and voting |

The European Strategy for Particle Physics
Update 2013

...ments

...do with the data?

...the HEP community

...fact simply lost the data

...ess, is generally not part of ...an experiment

...en in the main not planned by the ...knowledgeable people

...tapes is not equivalent to

*...a is stored in file formats appropriate for ...”*

*...piting the data is under the control of the*

▪ *"We are sure most of the data are not easily accessible!"*

4.



i) The success of particle physics experiments, such as those required for the high-luminosity LHC, relies on innovative instrumentation, state-of-the-art infrastructures and large-scale data-intensive computing. *Detector R&D programmes should be supported strongly at CERN, national institutes, laboratories and universities. Infrastructure and engineering capabilities for the R&D programme and construction of large detectors, as well as infrastructures for data analysis, data preservation and distributed data-intensive computing should be maintained and further developed.*

5. Make funding dependent on costing digital assets across the whole lifecycle

# 5.   Make funding dependent on costing digital assets across the whole lifecycle

[http://science.energy.gov/funding-opportunities/digital-data-management/](http://science.energy.gov/funding-opportunities/digital-data-management/)

- *"The focus of this statement is sharing and preservation of digital research data"*

- **All proposals submitted to the Office of Science (after 1 October 2014) for research funding must include a Data Management Plan (DMP) that addresses the following requirements:**

1. **DMPs should describe whether and how data generated in the course of the proposed research will be shared and preserved.**

    If the plan is not to share and/or preserve certain data, then the plan must explain the basis of the decision (for example, cost/benefit considerations, other parameters of feasibility, scientific appropriateness, or limitations discussed in #4).

    **At a minimum, DMPs must describe how data sharing and preservation will enable validation of results, or how results could be validated if data are not shared or preserved.**

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# ALICE Data Preservation plans

- The ALICE collaboration is committed to develop a long term program for Data Preservation to serve the triple purpose of
    i. preserving data, software and know-how inside the Collaboration,
    ii. sharing data and associated software and documentation with the larger scientific community, and
    iii. give access to reduced data sets and associated software and documentation to the general public for educational and outreach activities.
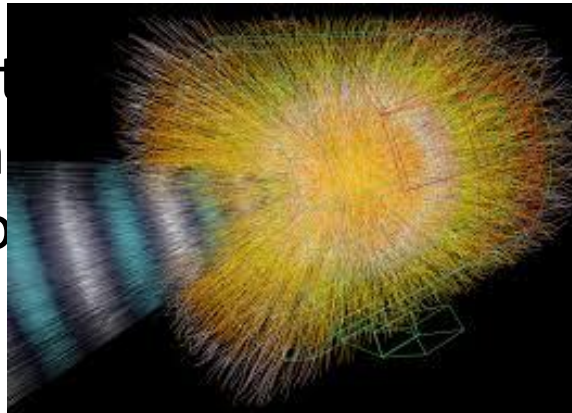
# ALICE Data Preservation plans

- The ALICE collaboration is committed to develop a long term program for Data Preservation to serve the triple purpose of

**The goal is to require the reproducibility of analysis in such virtualized environments as a prerequisite for publishing results.**

and

iii. give access t_____ and associated software an_____ the general public for educatio_____ tivities.

# 6.   Be collaborative and transparent to drive down costs

# 6.   Be collaborative and transparent to drive down costs

## 2020 Vision for LT DP in HEP

- *Long-term – e.g. FCC timescales: disruptive change*

  - By 2020, all **archived data** – e.g. that described in DPHEP Blueprint, including LHC data – easily **findable**, fully **usable** by **designated communities** with clear (Open) access policies and possibilities to annotate further

  - Best practices, tools and services well run-in, fully documented and sustainable; built in common with **other disciplines**, based on standards

  - **DPHEP portal**, through which data / tools accessed
    - ➢ **"HEP FAIRport": Findable, Accessible, Interoperable, Re-usable**

- ➢ **Agree with Funding Agencies clear targets & metrics**

# 6. Be collaborative and transparent to drive down costs

## Collaboration – Benefits

- In terms of 2020 vision, collaboration with other projects has arguably advanced us (in terms of implementation of the vision) by several years

- **I typically quote 3-5 years and don't think that I am exaggerating**

- **Concrete examples include "Full Costs of Curation", as well as proposed "Data Seal of Approval+"**

- With or without project funding, we should continue – and even strengthen – this collaboration
    - APA events, iDCC, iPRES etc. + joint workshops around RDA

- **The HEP "gene pool" is closed and actually quite small – we tend to recycle the same ideas and "new ones" sometimes needed**

6

# 6. Be collaborative and transparent to drive down costs

## APARSEN Training & Knowledge Base

- **_Long-_**

  – By
     inc
     **co**
     an

  – Be
     su:
     sta

  – **DP**
     ➢

➢ **Agre**

# POINT FOR DISCUSSION

# http://science.energy.gov/funding-opportunities/digital-data-management/

- *"The focus of this statement is sharing and preservation of digital research data"*

- **All proposals submitted to the Office of Science (after 1 October 2014) for research funding must include a Data Management Plan (DMP) that addresses the following requirements:**

1. **DMPs should describe whether and how data generated in the course of the proposed research will be shared and preserved.**

   If the plan is not to share and/or preserve certain data, then the plan must explain the basis of the decision (for example, cost/benefit considerations, other parameters of feasibility, scientific appropriateness, or limitations discussed in #4).

   **At a minimum, DMPs must describe how data sharing and preservation will enable validation of results, or how results could be validated if data are not shared or preserved.**

U.S. DEPARTMENT OF **ENERGY** | Office of Science

- *"Th............................................................................................earch dat..................*

- **All proposals submitted to the Office of Science (after 1 October 2014) for research funding must include a Data Management Plan (DMP) that addresses the following requirements**

1. **DMPs should describe whether and how data generated in the course of the proposed research will be shared and preserved.**

   If the plan is not to share and/or preserve certain data, then the plan must explain the basis of the decision (for example cost/benefit considerations, other parameters of feasibility, scientific appropriateness, or limitations discussed in #4)

**User Communities & Virtual Research Environments**

**(International) Funding Agencies**

**Service Providers (e-infrastructures)**

...escribe h...
...lts, or ho...

ENERGY | Science

38

- *"Th... ...earch dat...*

**User Communities &
Virtual Research Environments**

- **All proposals submitted to the Office of Science (after 1 October 2014) for research funding must include a Data Management Plan (DMP) that addresses the following requirements.**

1. **DMPs should describe wh... ...ata generated in the course of the proposed research wi... ...preserved.**

   If the plan is not to shar... ...ain data, then the plan must explain the basis of th... ...st/benefit considerations, other parameters of... ...ateness, or limitations discussed in #4)

**(International)
Funding Agencies**

escribe h...
lts, or h...

**Service Providers
(e-infrastructures)**

Science

# Summary

- Thanks to stimulation from and interactions with the 4C project, we have developed a simple "cost model" for LHC data ("bit preservation")

- This has been input to the CERN Resource Review Board (LHC machine, experiments, computing) – a sustainable funding scheme for decades

- A portal will help us inform people of the current state, including the implementation of DMPs and access to the preserved data (HEP FAIRport)