5 December 2013, Eric Grancher

CERN, head of database services

# CERN and Oracle, a 30-year collaboration

# Outlook

- CERN

- History of using Oracle

- Current usage

- Collaboration

- Why using Oracle in our research environment?

# CERN

- European Organization for Nuclear Research
  - Founded in 1954
  - Research: Seeking and finding answers to questions about the Universe
  - Technology, International collaboration, Education



## Twenty Member States

Austria, Belgium, Bulgaria, Czech Republic, Denmark, Finland, France, Germany, Greece, Italy, Hungary, Netherlands, Norway, Poland, Portugal, Slovakia, Spain, Sweden, Switzerland, United Kingdom

## Seven Observer States

European Commission, USA, Russian Federation, India, Japan, Turkey, UNESCO

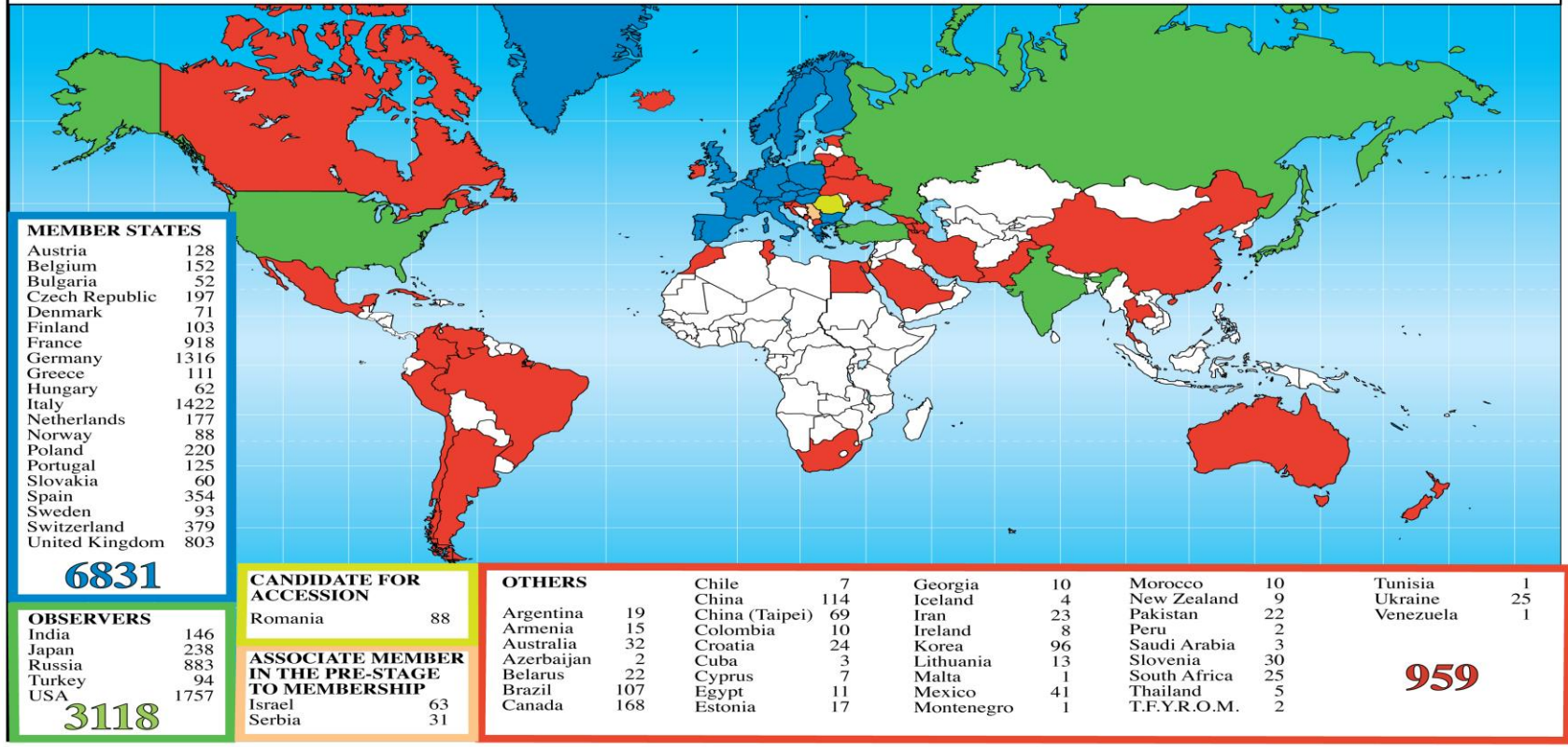## Associate Member States     Candidate State

Israel, Serbia                                Romania

## People

~2400 Staff, ~900 Students, post-docs and undergraduates, ~9000 Users, ~2000 Contractors
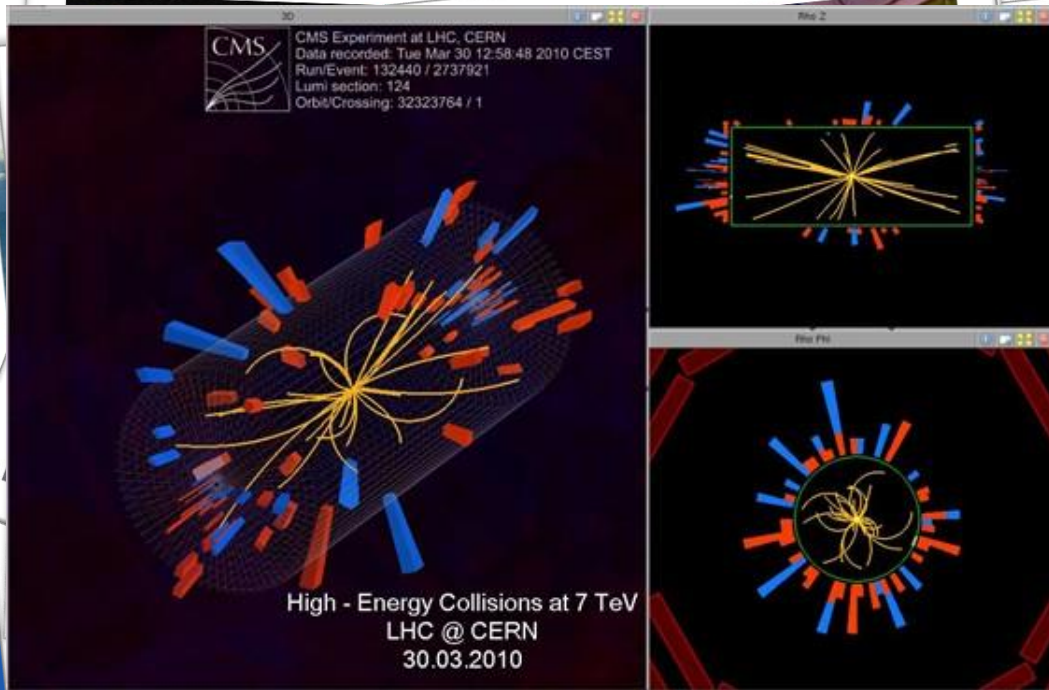
# A European Laboratory with Global reach



## Distribution of All CERN Users by Location of Institute on 14 January 2013

**MEMBER STATES**

| | |
|---|---|
| Austria | 128 |
| Belgium | 152 |
| Bulgaria | 52 |
| Czech Republic | 197 |
| Denmark | 71 |
| Finland | 103 |
| France | 918 |
| Germany | 1316 |
| Greece | 111 |
| Hungary | 62 |
| Italy | 1422 |
| Netherlands | 177 |
| Norway | 88 |
| Poland | 220 |
| Portugal | 125 |
| Slovakia | 60 |
| Spain | 354 |
| Sweden | 93 |
| Switzerland | 379 |
| United Kingdom | 803 |

**6831**

**OBSERVERS**

| | |
|---|---|
| India | 146 |
| Japan | 238 |
| Russia | 883 |
| Turkey | 94 |
| USA | 1757 |

**3118**

**CANDIDATE FOR ACCESSION**

| | |
|---|---|
| Romania | 88 |

**ASSOCIATE MEMBER IN THE PRE-STAGE TO MEMBERSHIP**

| | |
|---|---|
| Israel | 63 |
| Serbia | 31 |

**OTHERS**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | Chile | 7 | Georgia | 10 | Morocco | 10 | Tunisia | 1 |
| | | China | 114 | Iceland | 4 | New Zealand | 9 | Ukraine | 25 |
| Argentina | 19 | China (Taipei) | 69 | Iran | 23 | Pakistan | 22 | Venezuela | 1 |
| Armenia | 15 | Colombia | 10 | Ireland | 8 | Peru | 2 | | |
| Australia | 32 | Croatia | 24 | Korea | 96 | Saudi Arabia | 3 | | |
| Azerbaijan | 2 | Cuba | 3 | Lithuania | 13 | Slovenia | 30 | | |
| Belarus | 22 | Cyprus | 7 | Malta | 1 | South Africa | 25 | | |
| Brazil | 107 | Egypt | 11 | Mexico | 41 | Thailand | 5 | | |
| Canada | 168 | Estonia | 17 | Montenegro | 1 | T.F.Y.R.O.M. | 2 | | |

**959**

5

# LHC

- The largest particle accelerator & detectors



17 miles (27km) long tunnel

Thousands of superconducting magnets

Coldest place in the Universe: *1.9 degrees kelvin*

Ultra vacuum: 10x emptier than on the Moon

600 million collisions per second / analysis is like finding a needle in 20 million haystacks
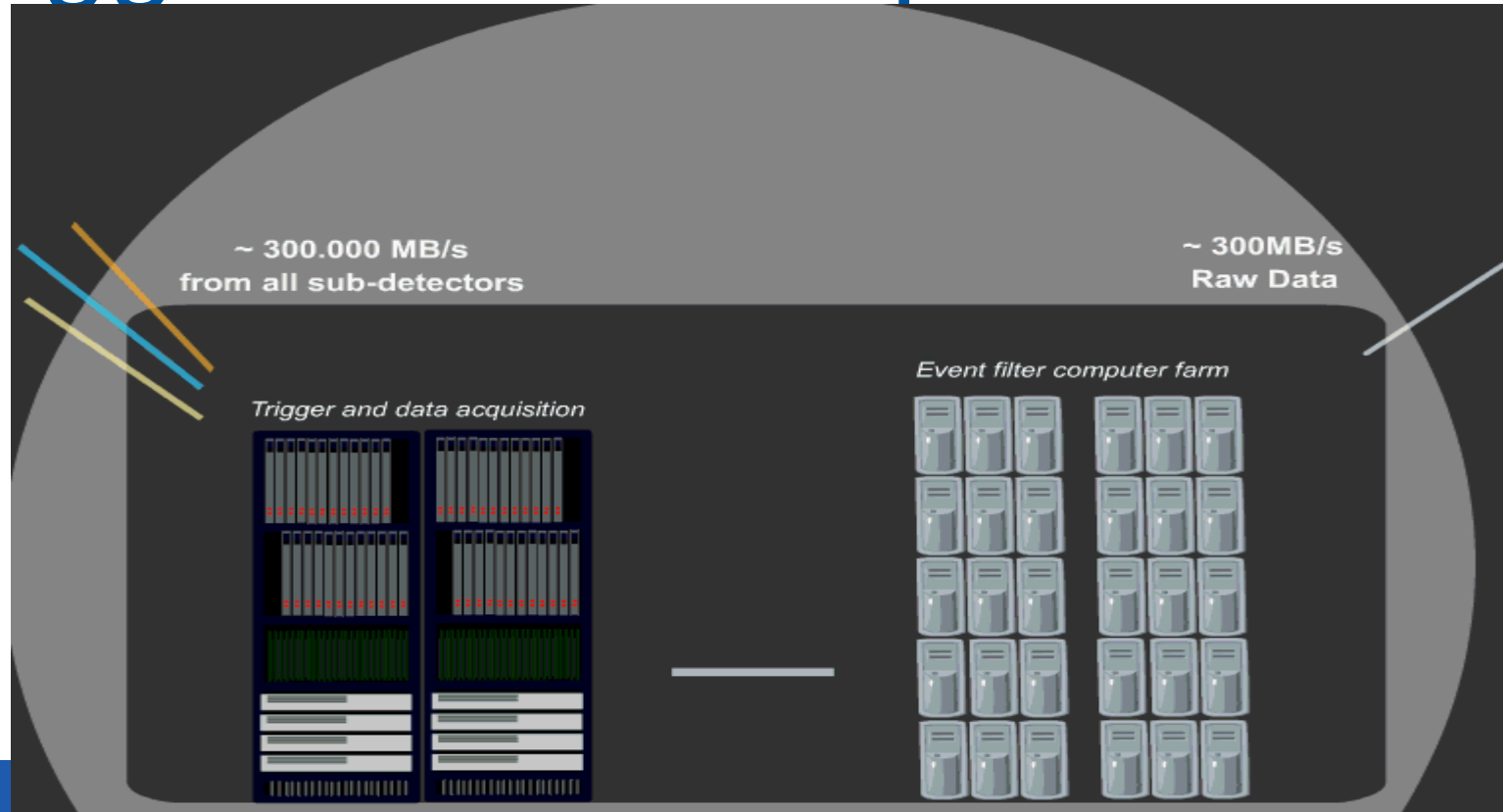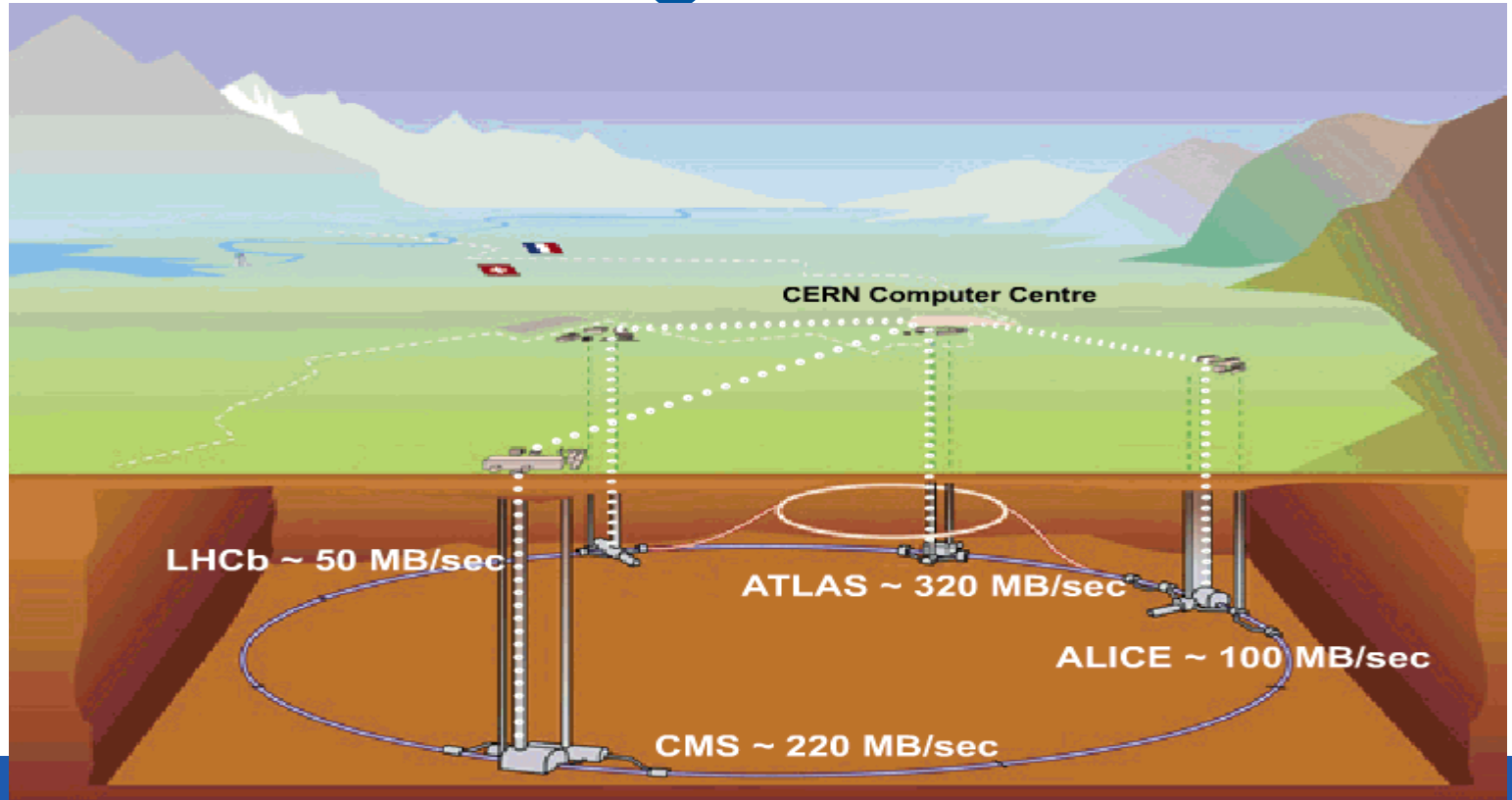
# Events at LHC



Luminosity :
$10^{34} cm^{-2} s^{-1}$

40 MHz – every 25 ns

20 events overlaying

# Trigger & Data Acquisition



~ 300.000 MB/s
from all sub-detectors

~ 300MB/s
Raw Data

Event filter computer farm

Trigger and data acquisition

# Data Recording



CERN Computer Centre

LHCb ~ 50 MB/sec

ATLAS ~ 320 MB/sec

ALICE ~ 100 MB/sec

CMS ~ 220 MB/sec

# World's largest computing grid - WLCG



- 1 PB raw data per second before filtering / >20 PB of new data annually
- 68,889 physical CPUs / 305,935 logical CPUS
- 157 computer centres around the world

# Oracle at CERN, 1982 accelerator control

http://cds.cern.ch/record/443114?ln=en

## ORACLE - the data base management system for LEP

J.Schinzel

Following the decision that an efficient data base system is required for the LEP project and that the systems at present in use at CERN are not adequate, an enquiry into possible data base management systems on the market was launched early this year.

# Oracle at CERN, version 2.3



Credit: N. Segura
Chinchilla

# Accelerator logging



Credit:
C. Roderick

15

# Accelerator logging



Credit:
C. Roderick

16

# Accelerator logging

50TB/year, rate to increase to 100 – 150 TB in 2014

(Quench Protection System)

**INTEGRATED LOGGED RECORDS**

5.3 trillion records logged
175 TB

LHC Operation

LHC Comissioning

1 TB Logged

Logged Records

17

# Administrative systems



- AIS has standardized on Oracle as database and uses the Oracle database as interface between the tools

- Java EE and Apex, deployment with Weblogic

- Oracle E-Business HR, managed as Oracle Managed Cloud Services ("OnDemand")

# Engineering applications

- An integrated PLM platform based on commercial tools.
- Simplified web interfaces for precise tasks

3D CAD
CATIA

Design data
management

Manufacturing
follow-up

Installation
follow-up

Maintenance
management

Data publishing

Workflow actions

**CERN Engineering & Equipment Data Management System**

**Smarteam**
CAD data, Items,
Drawings & Documents

**Agile PLM**
Design data
Drawings & Documents

**Infor EAM**
Asset tracking
Work management

**Other DBs**
Layout databases
Production databases

# Design data management

## Design baseline with full configuration management

- Workflows, versioning rules and access control based on project dependent contexts.
- Fully web-based and distributed approval processes.

# Manufacturing follow-up

Credit:
D. Widegren

**Follow-up of each manufactured component**

- Manufacturing & test data captured at manufacturing sites.
- Predefined manufacturing workflows for each equipment type.

# Installation follow-up

Credit:
D. Widegren

## Detailed logging of Installation & Commissioning tasks

- Over 150.000 jobs logged – allows detailed progress reporting
- Resolution of non-conformities with distributed approval processes.

# PLM @ CERN in numbers

**Document & Drawings (incl. CAD):**

~1,500.000 documents & drawings

~7,000 new documents & drawings created per month

**Components:**

~1,300,000 registered individually followed equipment

~3,000,000 equipment interventions/jobs logged

~ 15,000 equipment interventions/jobs logged per month

# CASTOR and Oracle, tapes

- Home made mass storage system, relies on Oracle databases for name server, request handling and staging

- 4 libraries, SL8500
- 10088x4 = 40K slots (4500 free)
- Occupancy: 65PB worth of data
- Drives: 20 T10KB legacy drives; 40 T10KC drives (to be replaced by T10KD's)

# CASTOR Archive in Numbers

Credit: German Cancio Melia

**Data:**

- ~90PB of data on tape; 250M files
- Up to 4.5 PB new data per month
- Over 10GB/s (R+W) peaks



Data written to CASTOR, in TB

# Tape Writing

- Writing to tape is well managed and understood
  - Detector data can't wait for us!
- High *aggregate* transfer rates are not a problem
  - System designed to split the write stream onto several tape drives

- However, *low per-drive performance*
  - "Tape Mark" overhead per tape file up to 2s
  - CASTOR has 3 tape files per payload file
  - No significant overhead 14y ago on 20MB/s drives…!
  - Today with avg files of 350 MB: < 50% efficiency!

**Tape write performance, MB/s**
(LTO-5)

Credit: German
Cancio Melia

- Changed CASTOR tape layer to write "logical" (or "buffered") tape marks between files (via CERN-prototyped Linux kernel extension)
- Now, achieving near-native drive speed by writing data in effective chunks of ~32GB
  - Significant savings in tape drive purchases

# Data Write Rate Increase



Boosted per-drive performance: from ~55 MB/s to ~150 MB/s

Credit: German
Cancio Melia

# Experiment online systems

- Experiments rely on a SCADA system for their control

- Up to 150,000 changes / second stored in Oracle databases

# Experiment offline systems

Credit:
Vakho Tsulaia

- **Geometry DB**

  - Relational database of **Primary Numbers** for the ATLAS Detector Description
    - All data for building GeoModel description in single place
    - Contains pointers to external files
      - Identifier dictionaries
      - Magnetic field maps (becoming obsolete)
      - All such files are shipped with the s/w release, no extra steps needed for getting them

- **Conditions DB**

  - Large relational database containing information about **Detector Status**, **Data-Taking Conditions**, **Calibrations**, **Alignment ...**

  - ATLAS Conditions DB is a **COOL Database**
    - COOL: one of 3 components of the LCG Persistency Framework (other two: POOL, CORAL)

# Oracle at CERN

- From accelerator control to
  - accelerator logging,
  - administration,
  - engineering systems,
  - access control,
  - laboratory infrastructure (cabling, network configuration, etc.),
  - mass storage system,
  - experiment online systems,
  - experiment offline systems,
  - Etc.

# openlab (1/3)

- Public-private partnership between CERN and leading ICT companies, currently in fourth phase (started in 2003)

- Its mission is to accelerate the development of cutting-edge solutions to be used by the worldwide LHC community

- Innovative ideas aligned between CERN and the partners, for products "you make it, we break it"

# openlab (2/3)

- Many successes:
  - RAC on Linux x86 (9.2 PoC and 10.1 production with ASM),
  - Additional required functionality (IEEE numbers, OCCI, instant client, etc.),
  - PVSS and RAC scalability,
  - Monitoring with Grid Control,
  - Streams world wide distribution,
  - Active DG, GoldenGate,
  - Analytics for accelerator, experiment and IT,
  - Etc.
- Regular feedback with joint selection of topics, some of the projects are common with more than one partner

# openlab (3/3)

- Publications (web, paper) and presentations of results, visitors

- Maaike Limper, best poster award at The International Conference on Computing in High Energy and Nuclear Physics 2013

# Oracle in our research environment

- Even if computing is critical for HEP, it is not the goal, there is a lot to do using solutions from commercial vendors which are industry supported and scalable

- Oracle has provided solutions along the years

- We have worked with Oracle to improve the tools to our (and others') needs with success

- Good for staff to work on industry standards for their future career

# Conclusion

- Not every day you build a 30+ years collaboration
- A long way since 1982, now very wide usage with applications, tape and database
- Oracle has proven to be reliable partner who cares and supports research
- Provide feedback and ideas for enhancements
- Helps focus on our core challenges
- A collaboration which works!

www.cern.ch