# Leveraging Oracle Big Data Discovery to Master CERN's Data

Manuel Martín Márquez

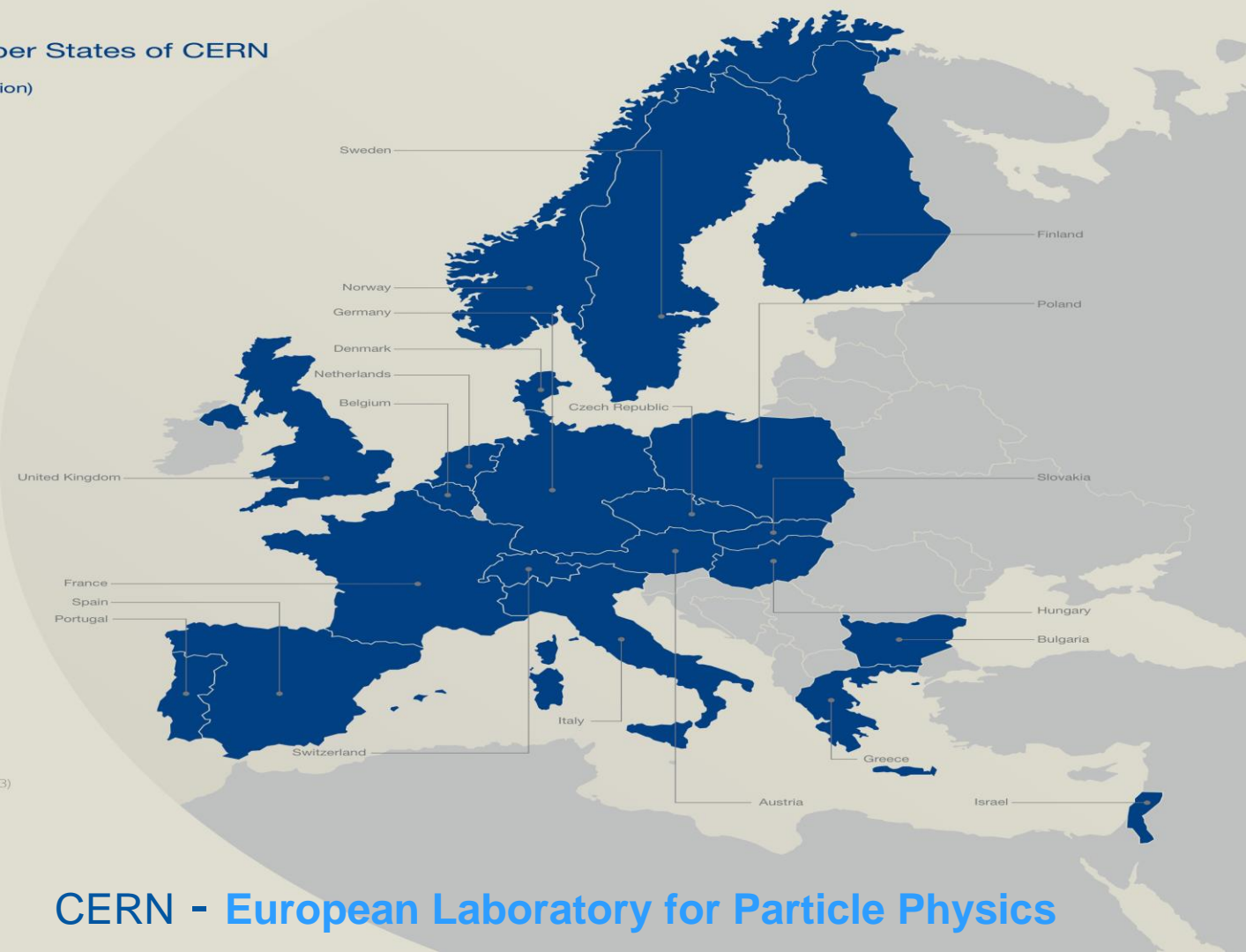**Oracle Business Analytics Innovation**
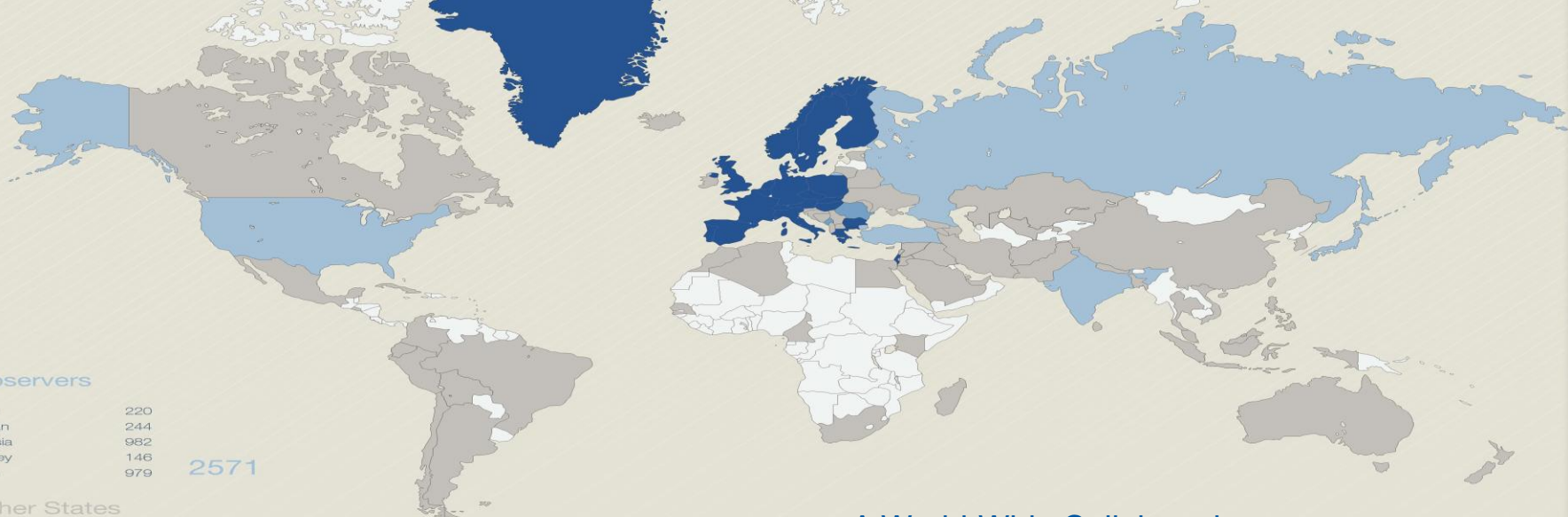
**12 October- Stockholm, Sweden**

The twenty one Member States of CERN

Member States (Dates of accession)

- Austria (1959)
- Belgium (1953)
- Bulgaria (1999)
- Czech Republic (1993)
- Denmark (1953)
- Finland (1991)
- France (1953)
- Germany (1953)
- Greece (1953)
- Hungary (1992)
- Israel (2014)
- Italy (1953)
- Netherlands (1953)
- Norway (1953)
- Poland (1991)
- Portugal (1986)
- Slovakia (1993)
- Spain (1/1961-12/1968-1/1983)
- Sweden (1953)
- Switzerland (1953)
- United Kingdom (1953)

CERN - European Laboratory for Particle Physics

## A World-Wide Collaboration

### Observers

| | |
|---|---|
| India | 220 |
| Japan | 244 |
| Russia | 982 |
| Turkey | 146 |
| USA | 979 |

**2571**

### Other States

| | | | | | |
|---|---|---|---|---|---|
| Afghanistan | 1 | El Salvador | 1 | Pakistan | 41 |
| Albania | 2 | Estonia | 16 | Palestine (O.T.) | 4 |
| Algeria | 8 | Georgia | 36 | Peru | 8 |
| Argentina | 11 | Gibraltar | 1 | Philippines | 1 |
| Armenia | 25 | Hong Kong | 1 | Saudi Arabia | 3 |
| Australia | 25 | Iceland | 4 | Senegal | 1 |
| Azerbaijan | 8 | Indonesia | 1 | Singapore | 2 |
| Bangladesh | 4 | Iran | 28 | Sint Maarten | 2 |
| Belarus | 47 | Ireland | 22 | Slovenia | 27 |
| Bolivia | 3 | Jordan | 2 | South Africa | 16 |
| Bosnia & Herzegovina | 1 | Kenya | 1 | Sri Lanka | 5 |
| Brazil | 108 | Korea, D.P.R. | 1 | Syria | 2 |
| Cameroon | 1 | Korea Rep. | 117 | Thailand | 12 |
| Canada | 134 | Kuwait | 1 | T.F.Y.R.O.M. | 1 |
| Cape Verde | 1 | Lebanon | 12 | Tunisia | 6 |
| Chile | 12 | Lithuania | 19 | Ukraine | 55 |
| China | 280 | Luxembourg | 4 | Uzbekistan | 4 |
| China (Tapei) | 45 | Madagascar | 4 | Venezuela | 9 |
| Colombia | 30 | Malaysia | 15 | Viet Nam | 9 |
| Croatia | 35 | Mauritius | 1 | Zimbabwe | 2 |
| Cuba | 7 | Mexico | 64 | | |
| Cyprus | 16 | Montenegro | 3 | | |
| Ecuador | 3 | Morocco | 12 | | |
| Egypt | 19 | Nepal | 5 | | |
| | | New Zealand | 7 | | |

**1415**

### Member States

| | | | | | |
|---|---|---|---|---|---|
| Austria | 99 | Greece | 152 | Slovakia | 88 |
| Belgium | 106 | Hungary | 68 | Spain | 337 |
| Bulgaria | 75 | Israel | 51 | Sweden | 75 |
| Czech Republic | 202 | Italy | 1686 | Switzerland | 180 |
| Denmark | 53 | Netherlands | 153 | United Kingdom | 640 |
| Finland | 87 | Norway | 61 | | |
| France | 751 | Poland | 229 | | |
| Germany | 1150 | Portugal | 109 | | |

**6352**

### Candidate for Accession

| | |
|---|---|
| Romania | 118 |

### Associate Members in the Pre-stage to Membership

| | |
|---|---|
| Serbia | 41 |

# Higgs Boson Discovery
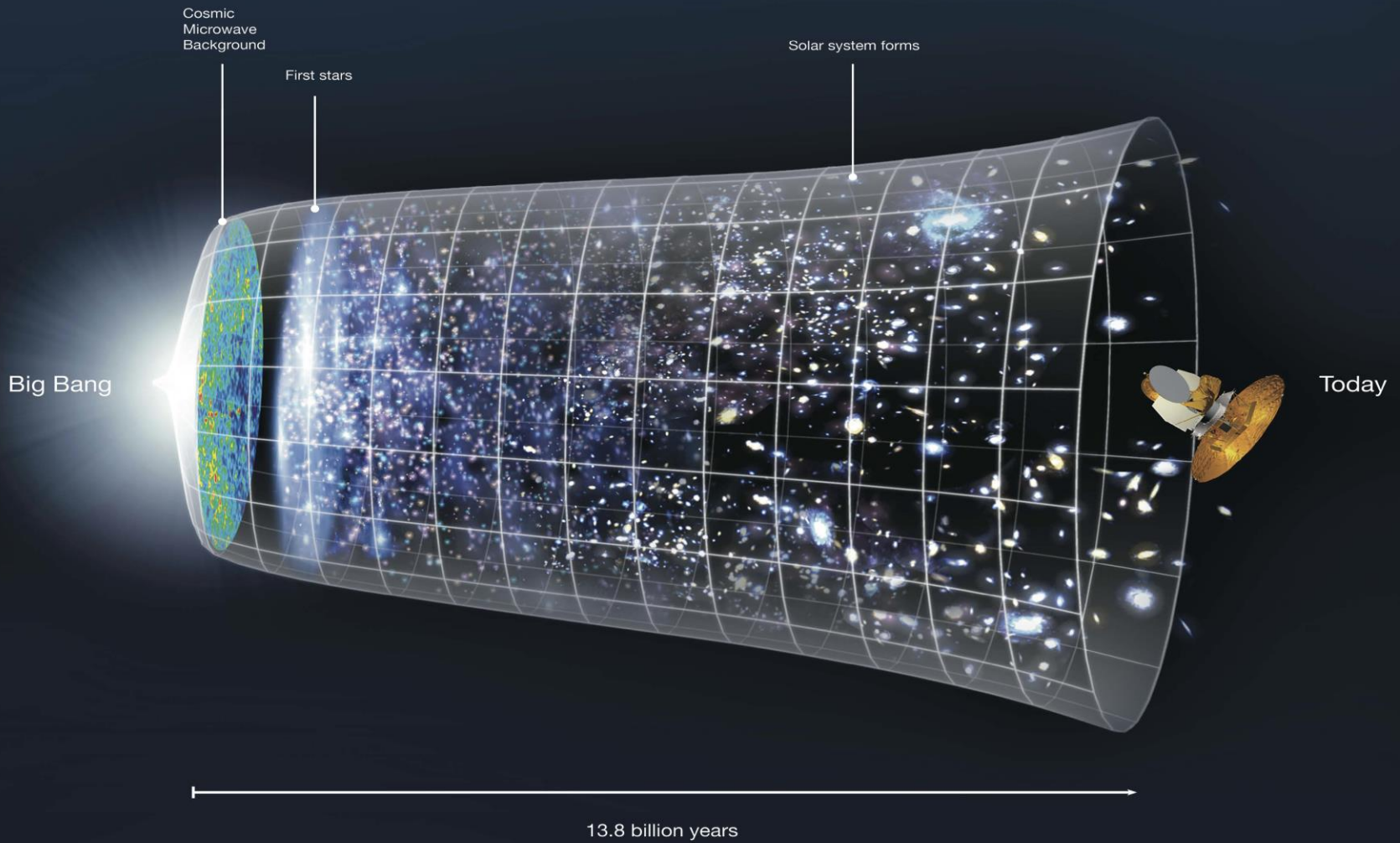
## 2012

Higgs to 4μ candidate event

ATLAS
EXPERIMENT
http://atlas.ch

Run:      204769
Event:    71902630
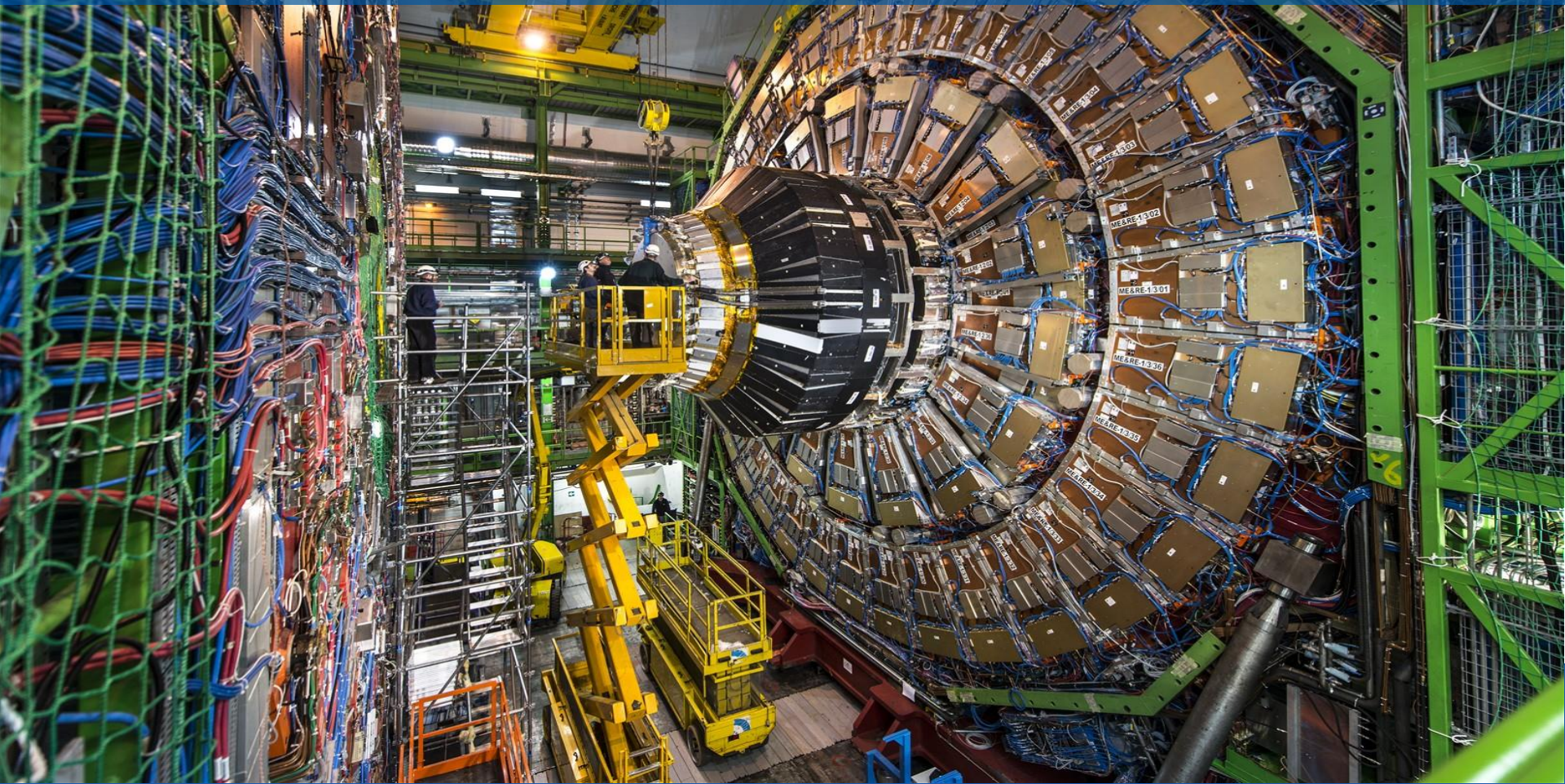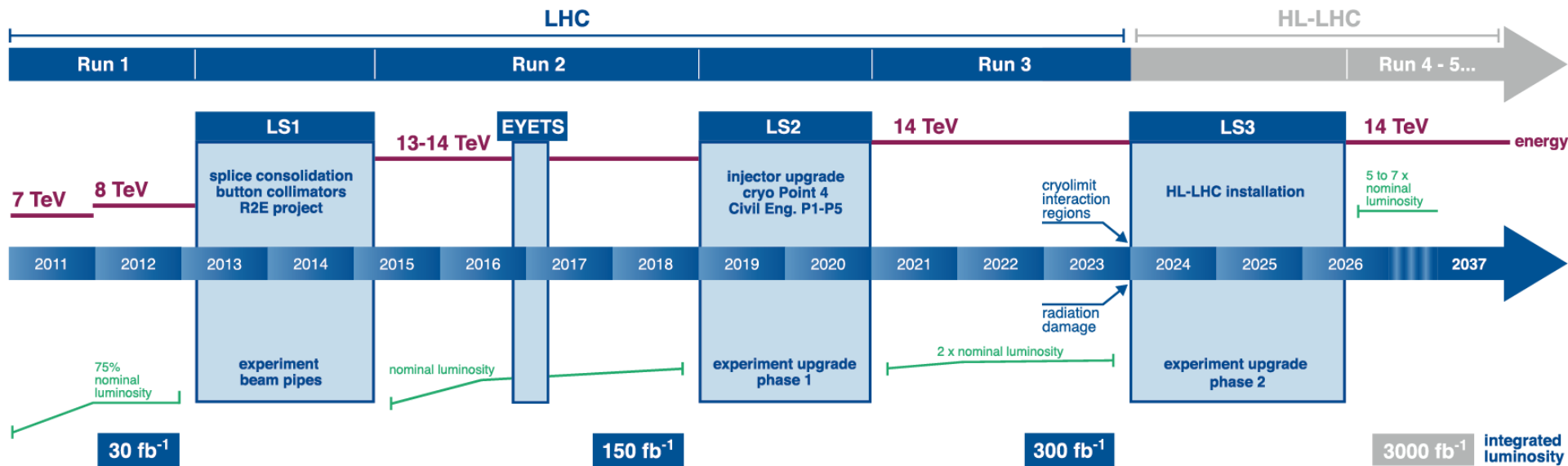Date:     2012-06-10
Time:     13:24:31 CEST

CERN Aerial View
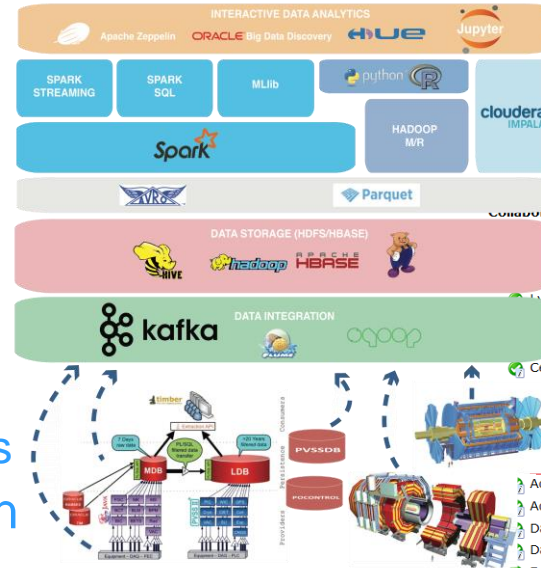
# Hadoop and Analytics – IT-DB-SAS

- **New scalable data services**
  - Scalable databases
  - Hadoop ecosystem
  - Time Series databases
- **Big Data Analytics**
- **Activities and objectives**
  - Support of Hadoop Components
  - Further value of Analytics solutions
  - Define scalable platform evolution
- **Hadoop Production Service**

# CERN Accelerator Logging Service



- +800 extraction clients
- +5 million extraction requests per day
- 130 custom applications

~ 1 million signals
~ 300 data loading processes
~ 4 billion records per day
~ 160 GB / day
→ 52 TB / year *stored*

~ 250'000 Signals
~ 50 data loading processes
~ 5.5 billion records per day
~ 275 GB / day
→ 100 TB / year throughput

Credit: BE-CO-DS

13

# CERN Accelerator Logging Service

- New Landscape bring new challenges
  - Better Performance on bigger datasets
    - Big Data queries: Impala, Spark SQL
  - Leverage analytics capabilities
    - Spark Analytics: Python, ML, R
  - More heterogeneous data access models

Credit: BE-CO-DS

**Storage Evolution - Size in GB / day**

**QPS**

# CERN Accelerator Logging Service



Credit: BE-CO-DS

# Accelerator Postmortem Analysis

- Postmortem Analysis
  - Diagnostic on failures
    - Continue operations safely
    - Intervention Required
- Designed for CERN LHC
  - Extended to injectors complex (SPS)
  - External Post Operational Checks
  - Injection Quality Checks

# Accelerator Postmortem Analysis

- Challenges:
  - Stringent Timing Constraint
  - Better scalability
    - data storage
    - IO throughput
  - Big Data Streaming Analytics

FCC
hh ee he

**Post-LHC accelerator projects (80-100 km)**

SUISSE
FRANCE

CMS

LHCb

ATLAS

CERN Meyrin

CERN Prévessin

SPS 7 km

ALICE

LHC 27 km

# Architecture overview

CDH 5.7.1
16 nodes, 24 GB ram
Intel Xeon L5520 @ 2.27GHz
165 TB HDFS

**cloudera**®

Coordination

Oracle Big Data Discovery
Libraries + Hive table detector

HIVE

Spark

Resource Management (YARN)

Data Storage    hadoop HDFS

Data Integration    FLUME

Big Data Discovery
v1.2.2
Dgraph & Studio

**ORACLE**®

**EXALYTICS**

4x Xeon E7-8895 v2 (15 cores each)
2 TB RAM
4.8 TB Flash + 6 x 1.2 TB 10K HDD

CERN

# Oracle Big Data Discovery Overview

- Data Exploration & Discovery
  - Interactive catalog of all data
  - Assess attribute statistics, data quality and outliers
  - Quick data exploration or create dashboards and applications
- Data Transformation with Spark in Hadoop
  - Apply built-in transformations or write your own scripts
  - Data Enrichment
    - Text: Entity extraction, relevant terms, sentiment, language detection
    - Geographical information: address, IP, reverse
  - Preview results, undo, commit and replay transforms
- Collaborative environment
  - Share and bookmarks
  - Create and share transformed datasets

# Data Transformation UI - ETL

# Discovery Applications

# Advance Analytics - Notebooks

- Easy to create and share documents that contain live code
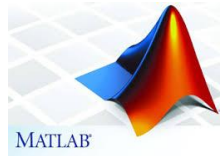- Step by step execution reproduce the analysis, charts, etc.
- Support for multiple languages/kernels
- Multiple notebook software available
  - Jupyter/IPython
  - BDD provides notebook from version 1.2.0 (BDD Shell)
    - Can be used with Jupyter/IPython
  - HUE notebooks
  - Apache Zeppelin
  - More…

# Scalable Analytics



Figure 6. Absolute difference (%) between request and feedback for QSCB_6_2CV120

- Reliability of degrading components of valves in the cryogenic system of the LHC (University of Delft)

  - BDD -> Data Extraction -> Refine Calculations

    

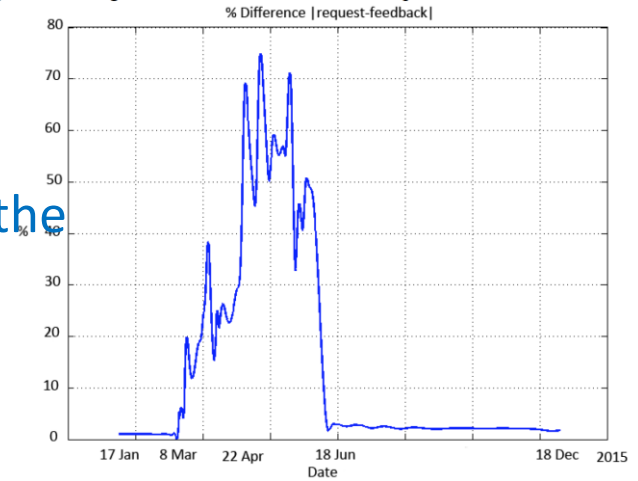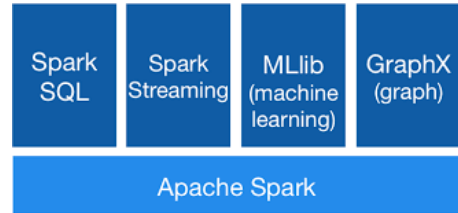  - Scalable solutions apply to all the cryogenics valves

    

# Conclusions

- Hadoop is not the solution for all your problems but..
- Unlock new ways to exploit your investment on data
  - overcome technical limitations for several CERN use cases
- Allows heterogeneous data access
  - not only SQL or custom java APIs
- Once the data is in Hadoop only half of the way is done
  - Data visualization and discovery
  - Notebooks are easy to use and powerful for advanced analytics
  - Self-service tools improve productivity
    - Users should be able to do what they need without IT intervention

www.cern.ch