

Oracle Big Data Discovery for CERN's Control Data

Antonio Romero Marin

What's CERN

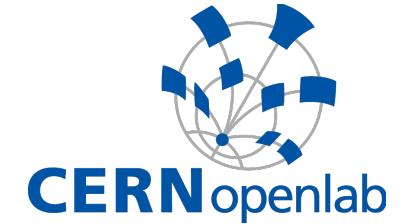


What's CERN

- European Laboratory for Particle Physics
 - Fundamental Research
- Worldwide International Collaboration
- Education & Training
- Push Frontiers of Technology



CERN openlab

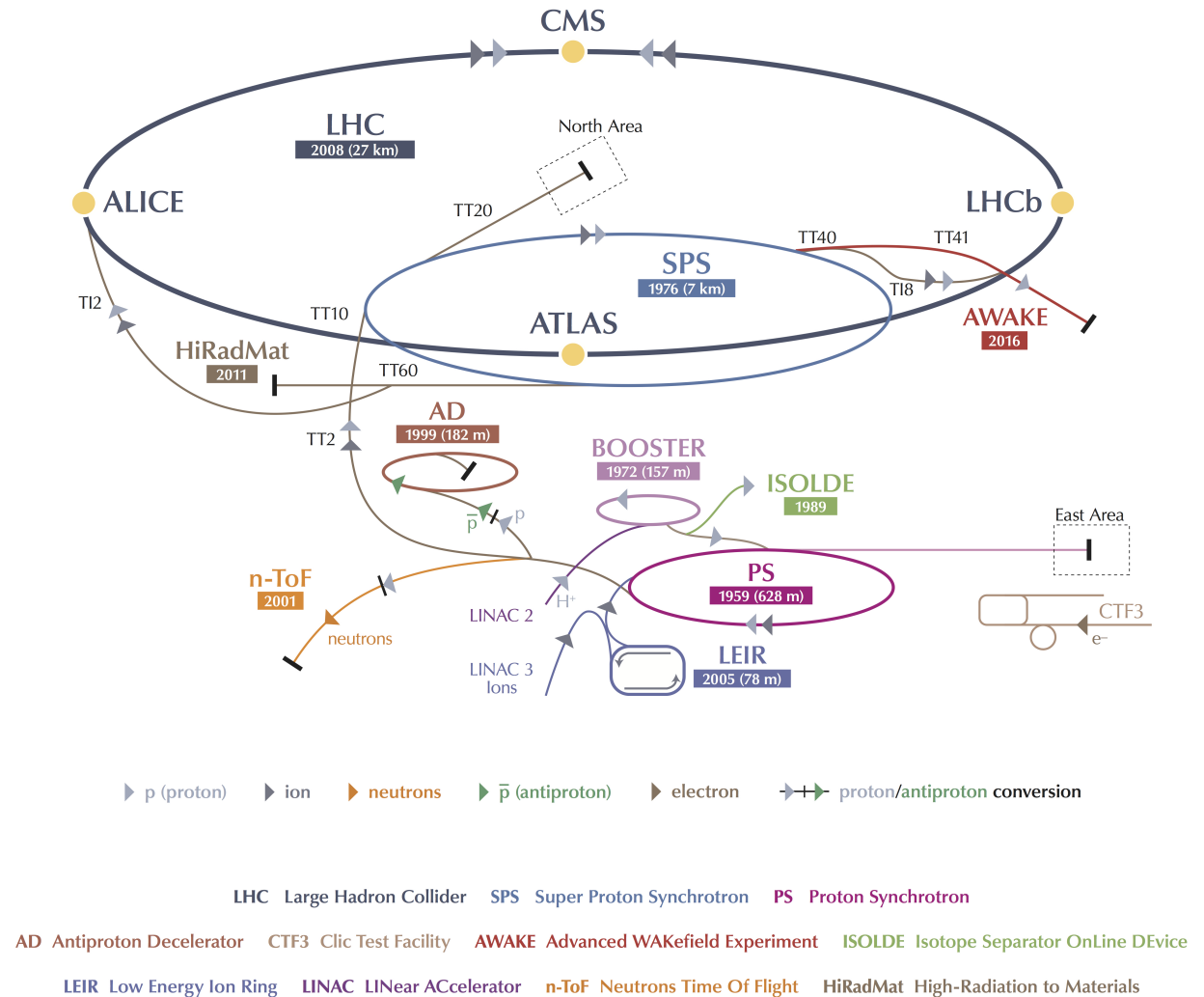


- Public-private partnership between CERN and leading ICT companies and research institutes
- Accelerate cutting-edge solutions for the worldwide LHC community and wider scientific research.
- Designed to create and disseminate knowledge
 - Publication of reports and articles
 - Workshops or seminars
 - CERN openlab Student Programme



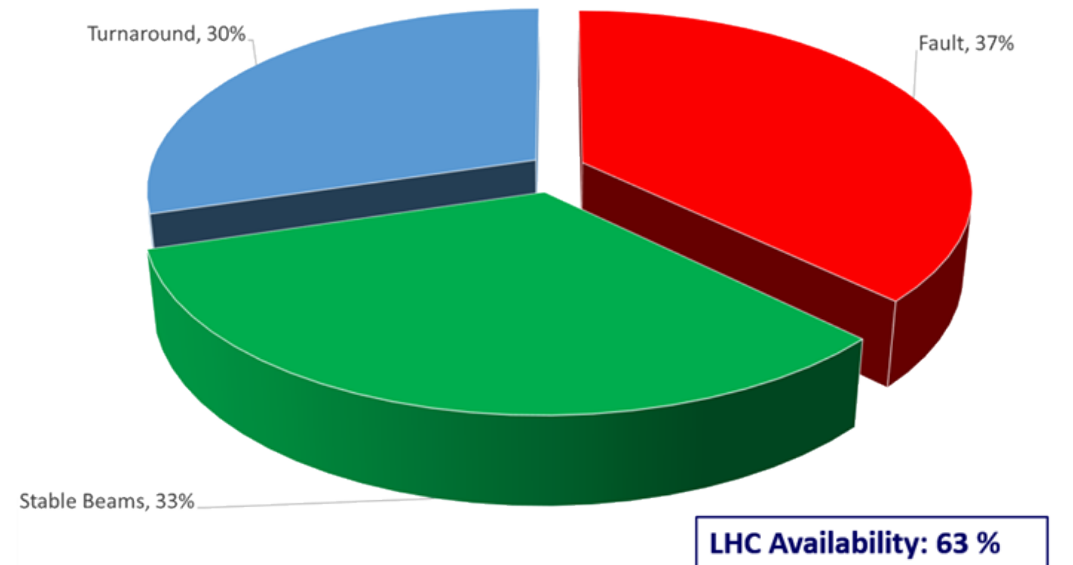
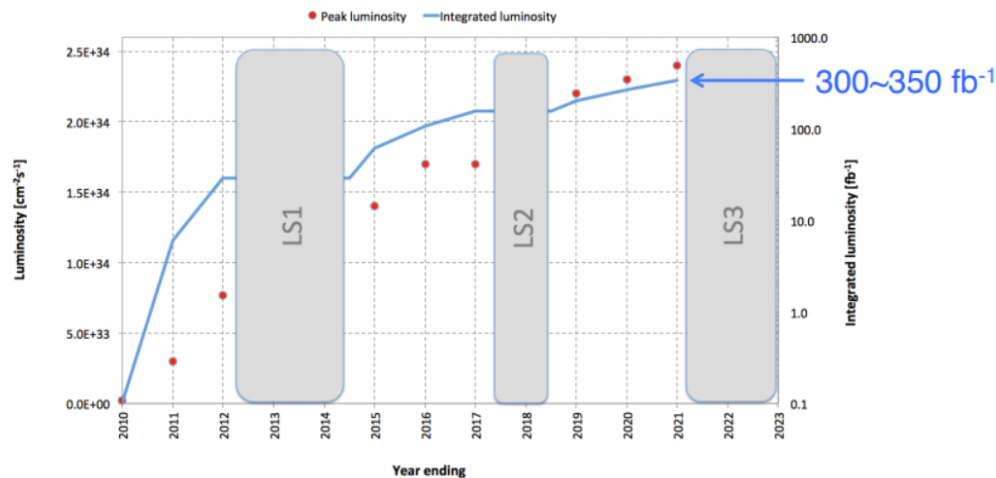
CERN Accelerator Complex

- Control and operations
 - Millions of sensors, signals
 - Large number of control devices
 - Equipment
- Monitoring and logging
- Supporting IT infrastructure
 - Databases
 - Network
 - Services

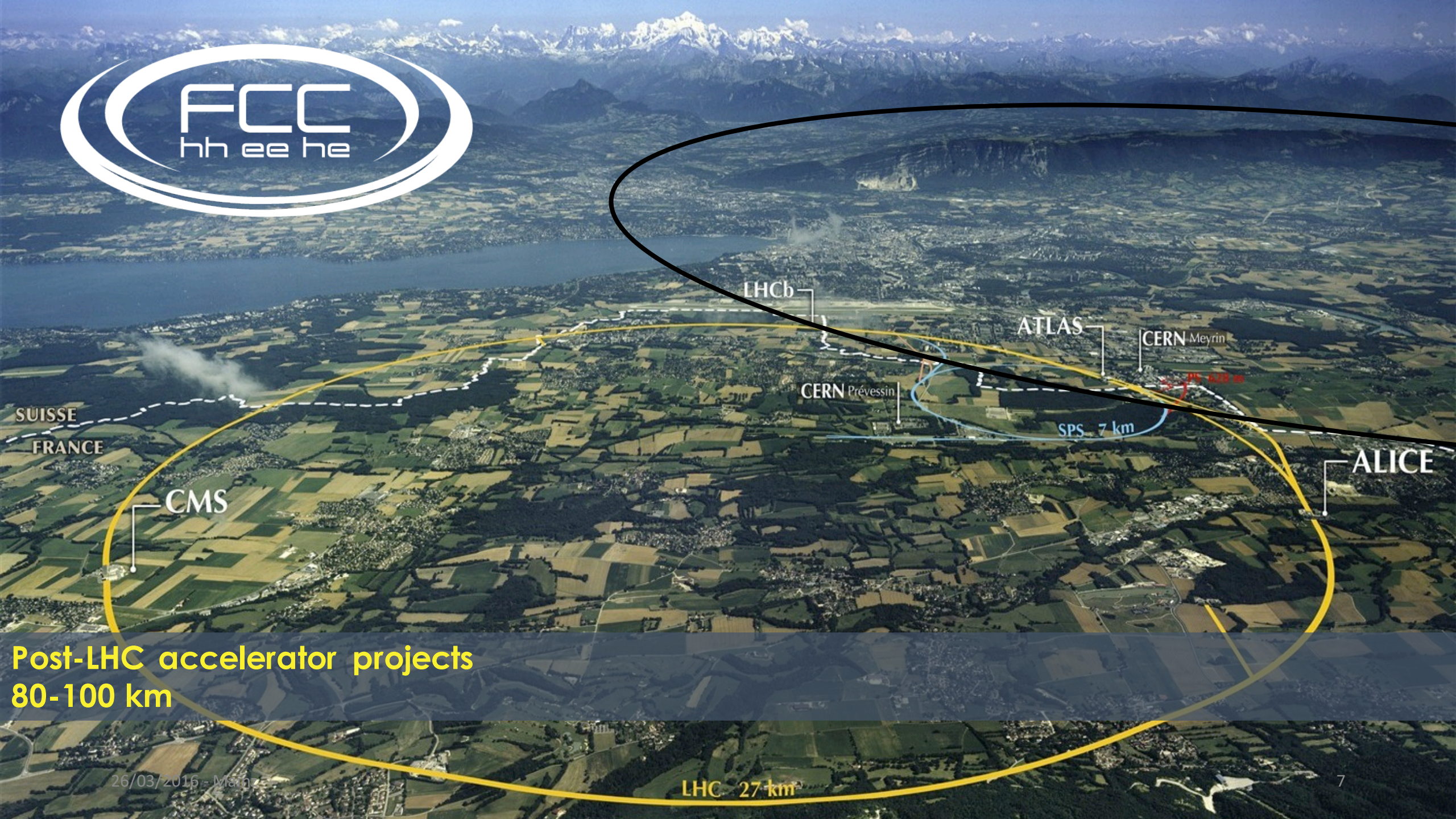


Data Analytics Challenges

- Some faults cannot be avoided
 - Decrease the availability for running physics
- Preventive maintenance is not enough
 - Does not take into account the condition of the equipment
- LHC upgrades will further increase luminosity
 - Computing resources needs will be higher
 - Data generated will increase drastically



A. Apollonio



LHCb

ATLAS

CERN Meyrin

CERN Prévessin

SPS 7 km

PS 6.18 km

SUISSE
FRANCE

CMS

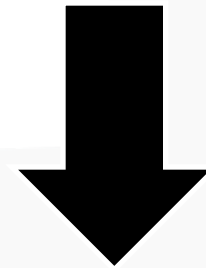
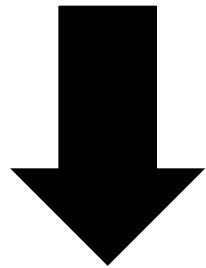
ALICE

Post-LHC accelerator projects
80-100 km

LHC 27 km

Data Analytics Objective

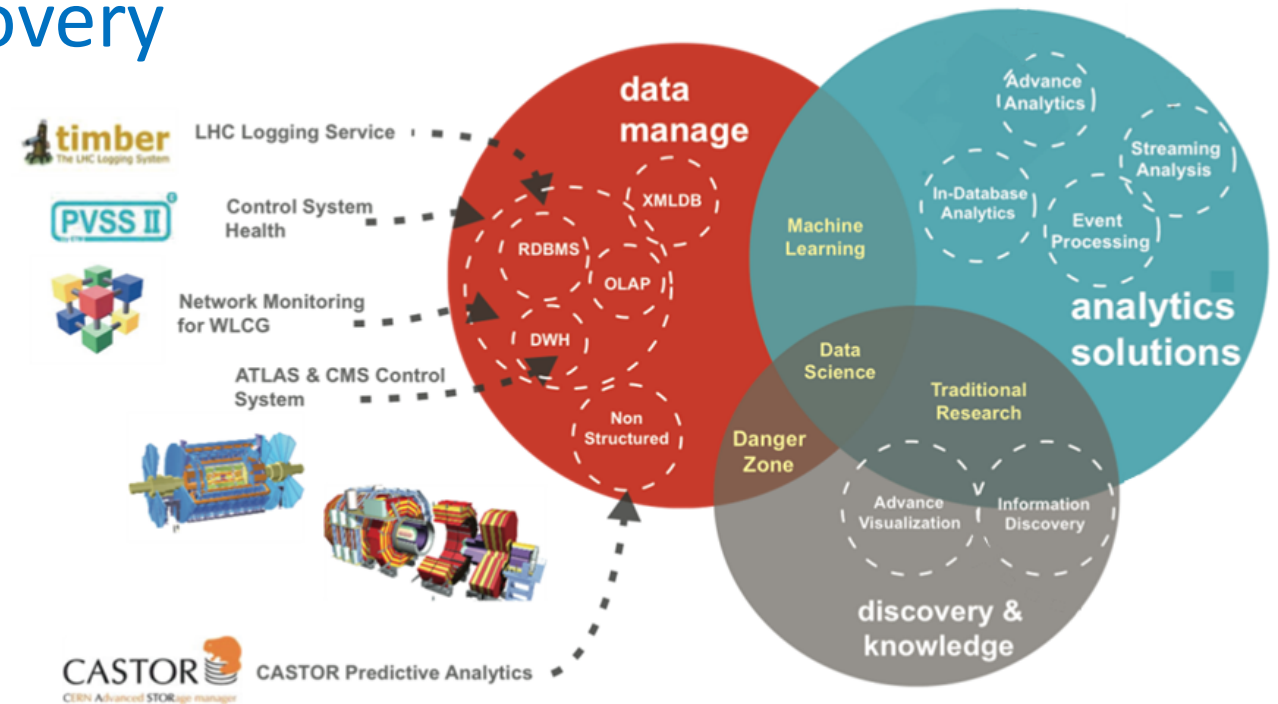
Control and Monitoring Systems



Intelligent, Predictive and Proactive Systems

Areas of investigation

- Predictive maintenance and system optimization
- Data extraction, transformation and loading (ETL)
- Data Visualization and Discovery



Use Case - FCC RAMS studies

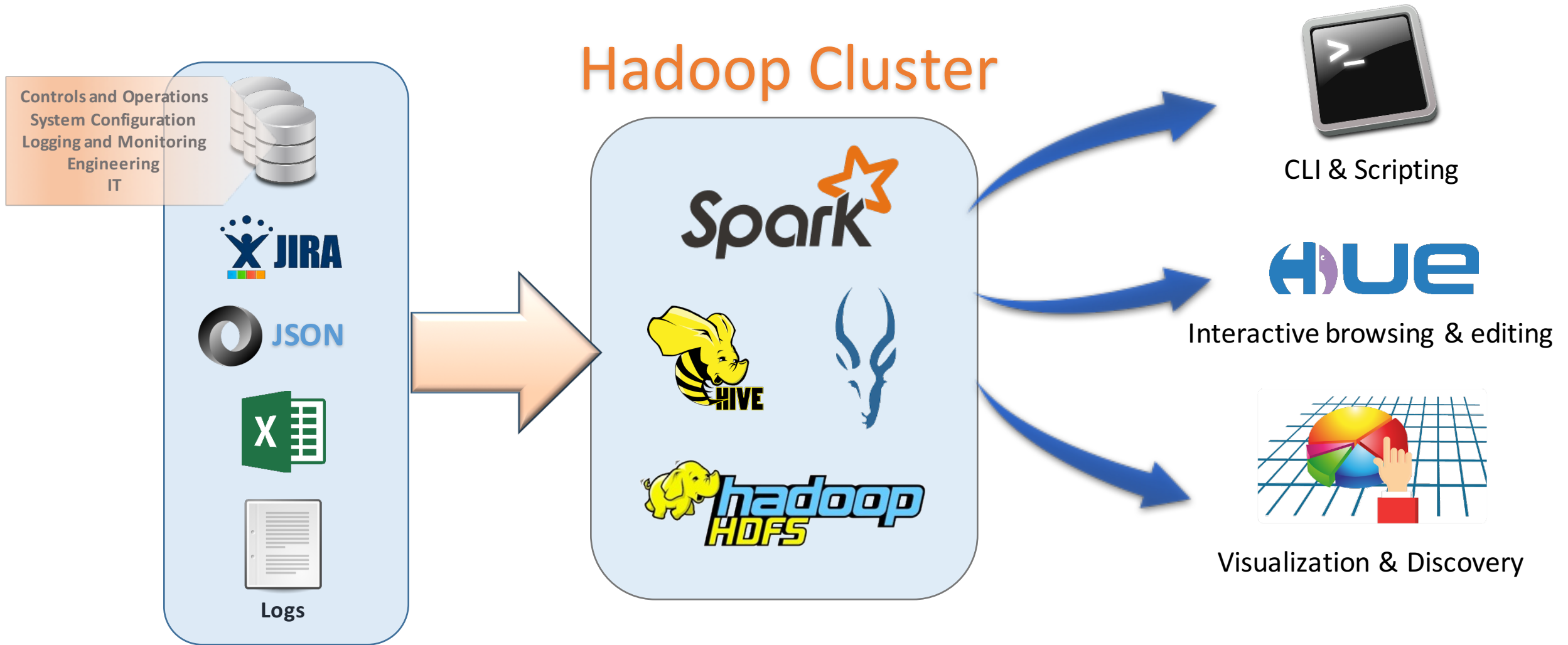
- Reliability, Availability, Maintainability and Safety (RAMS) studies for the Future Circular Collider (FCC)
- Study and increase the reliability and availability of the LHC
- Use RAMS findings to assess the feasibility of the needs of FCC
- Data distributed across multiple sources
 - Operations e-logbook
 - Accelerator Fault Tracking project
 - Accelerator logging service
 - Accelerator schedules
 - Cryogenics
 - Add more...
 - Vacuum, Power Converters, etc.



Requirements

- Flexible
 - Data => structured, semi-structured and non-structured, data editing
 - Use => Interactive, CLI & Scripting
- ETL functionalities
- Scalable
 - Data is foreseen to increase significantly (+datasets)
 - Processing
- Powerful
 - Browse data
 - Correlate Information
 - Visualization
 - Analytics

Hadoop based solution



HUE – Hadoop User Experience

- Hue is an open source suite of web-based applications for analyzing data with any Apache Hadoop
- It features:
 - SQL Editors for Hive, Impala, MySQL, PostGres, Sqlite and Oracle
 - Dynamic search dashboards for Solr
 - Spark Notebooks
 - Browsers for YARN, HDFS, Hive table Metastore, HBase, ZooKeeper
 - Pig Editor, Sqoop2, Oozie workflows Editors and Dashboards
 - Wizards to import data into Hadoop

HUE – Hadoop User Experience

The screenshot displays the HUE interface with two charts. The top chart is a bar chart showing the number of rows for various variable names. The bottom chart is a line chart showing the average value of a variable over time.

Top Chart: Bar Chart

Impala SQL:

```
select variable_name, count(*) as numrows
from lhlog_cryo
group by variable_name
order by numrows desc
limit 10
```

X-AXIS: variable_name

Y-AXIS: numrows

variable_name	numrows
QRLAB_23L1_GT943.POSST	6,887,08
QRLHA_05R4_GT938.POSST	6,000,000
QRLAB_15L6_GT947.POSST	5,800,000
QRLHA_05R4_GT937.POSST	5,000,000
QRLFE_04L8_GT930.POSST	5,000,000
QRLCC_07L4_GT947.POSST	4,800,000
ATLAS:LUMI_TOT_INST	4,600,000
QRLDE_06R8_GT931.POSST	4,500,000
QRLHA_05L4_GT935.POSST	4,400,000
QRLEA_06L2_GT931.POSST	4,300,000

Bottom Chart: Line Chart

Impala SQL:

```
select avg(value) as value, extract(utc_timestamp, "hour") as hour_
from lhlog_cryo
where utc_timestamp > "2015-06-12"
and utc_timestamp < "2015-06-13"
and variable_name = "QRLAB_23L1_GT943.POSST"
group by hour_
order by hour_
```

X-AXIS: hour_

Y-AXIS: value

hour_	value
0	12.33722
1	12.337
2	12.325
3	12.326
4	12.337
5	12.318
6	12.313
7	12.321
8	12.324
9	12.320
10	12.320
11	12.327
12	12.310
13	12.318
14	12.322
15	12.323
16	12.326
17	12.323
18	12.310
19	12.316
20	12.321
21	12.328
22	12.327
23	12.322

HUE – Hadoop User Experience

- One tool to use multiple Hadoop components
- Easy to use
- Compatible with multiple versions, open source
- Extensible
- But
 - Requires language knowledge to explore and transform the data
 - Limited for Data Discovery

Data Discovery

- Interactive and visual analytics
 - Find hidden patterns
 - Get **new insights**
- Intended to be used by the end users
 - Enabling them to use their intuition and knowledge of the data
- Powerful customization of dashboards and visualizations
 - Without intervention of IT
- Integrate multiple data sources
 - Analyze information of any type and any source

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

Components

- BDD Data Processing (Spark on YARN)
 - Hive Table Detector
 - Profiling and sampling
 - Transformations and enrichments
 - Refresh/incremental update datasets automatically or manually
- Dgraph (In-Memory Discovery Indexes)
 - In-memory, columnar, multi-core architecture
- Web Studio
 - Catalog, explore, transform and discover UI's

Architecture overview

cloudera®

CDH 5.5.1
16 nodes, 24 GB ram
Intel Xeon L5520 @ 2.27GHz
165 TB HDFS

Oracle Big Data Discovery
Libraries + Hive table detector



Resource Management (YARN)

Data Storage



Data Integration



C
o
o
r
d
i
n
a
t
i
o
n

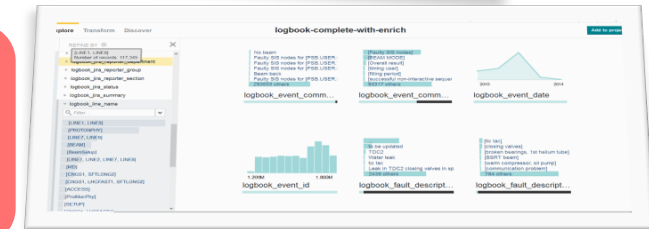
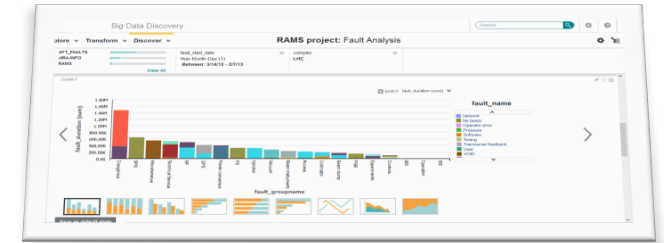


Dgraph
In-Memory Columnar
Database

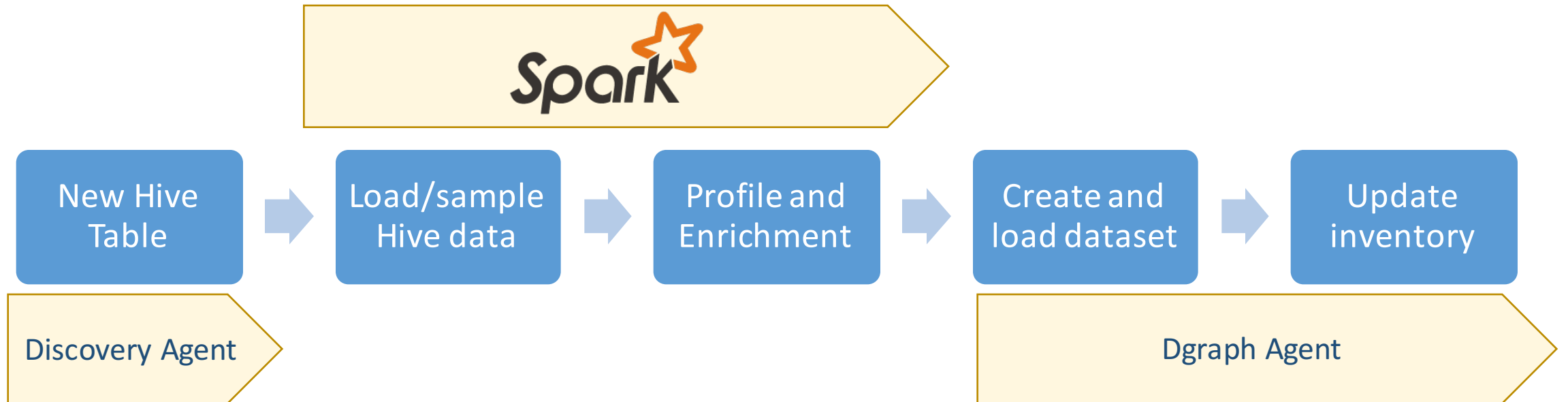
ORACLE®
EXALYTICS



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



Data Processing Workflow

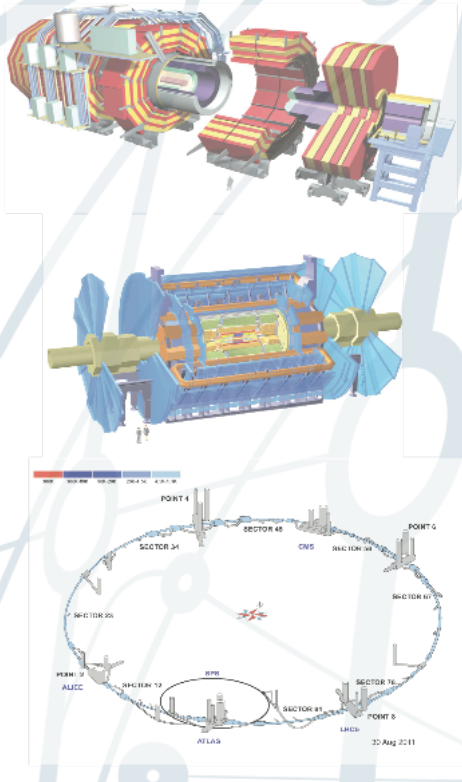


Important Technical Features

- Supports HDP and CDH Hadoop distributions
- Kerberos support
- Spark on YARN
- Data set and project level security in Studio
- Incremental updates and refresh datasets
- Applications and dashboards
- Ability to publish and share transformation scripts
- Custom visualization support (Javascript, D3, EQL)

Use Case: FCC RAMS

Scenario




- Power Converters
- Cryogenics
- Machine Protection
- Accelerator Major Events
- Accelerator Fault Tracking
- Accelerator Logging
- Operations logbook

•NoSQL
•XML
•JSON
•Text
•RDBMS






Datasets Catalog

 **7 Projects**
[View all](#)




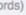
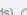


 **11 Data Sets**
[View all](#)

 **Add Data Set**

Recently Viewed Data Sets


aft_faults Data Source: fcc_rams.aft_faults (1.8k records)  Preview	op_logbook Data Source: fcc_rams.op_logbook (1.1M records)  Preview	cms_lhcfills Data Source: cms_lhcfills.csv (621 records)  Preview
--	--	---

Most Popular Data Sets [View More](#)


op_logbook Data Source: fcc_rams.op_logbook (1.1M records)  Preview	apache_mwctl_prod_db... Data Source: apache-mwctl-test-db-a0... (183 records) Preview	apache_mwctl_prod_db... Data Source: apache-mwctl-test-db-a0... (745 records) Preview	aft_faults Data Source: fcc_rams.aft_faults (1.8k records)  Preview	cms_lhcfills Data Source: cms_lhcfills.csv (621 records)  Preview	aft_cardiogram Data Source: fcc_rams.aft_cardiogram (1.2M records)  Preview	cms_runs_for_fill Data Source: cms_runs_for_fill.csv (2.5k records) Preview
cms_tmb_rates Data Source: cms_tmb_rates.csv (820.6k records)  Preview	high_level_summary_sc... Data Source: fcc_rams.high_level_sum... (63 records)  Preview	lhlog_cryo Data Source: fcc_rams.lhlog_cryo (1.1B records)  Preview New	naive_ml_prediction Data Source: naive_ml_prediction.csv (1.2M records) Preview			

lhlog_cryo (1,078,271,323 records)

Data Set Info **Used in Projects (1)** Related Data sets (4)

 **Tags** *To add tags, click the Tags button at left.*

Attributes


Data source: fcc_rams.lhlog_cryo
Data source type: Hive
Hive Table name: fcc_rams.lhlog_cryo
Created on: 1/17/2016 4:52:14 PM (UTC)
Access: Public [edit](#)
Data set key: edp_cli_edp_bfb068c8-3280-4bcf-a107-87a3d4ef7e04

Actions

- [Explore](#)
- [Add to project](#)
- [Edit tags](#)
- [Reload data set](#)
- [Delete](#)

Summary

0 Views
Last Updated
1/17/2016 4:52:16 PM (UTC)

Newly Added Data Sets [View More](#)

lhlog_cryo Data Source: fcc_rams.lhlog_cryo (1.1B records)  Preview New	op_logbook Data Source: fcc_rams.op_logbook (1.1M records)  Preview	apache_mwctl_prod_db... Data Source: apache-mwctl-test-db-a0... (745 records) Preview	apache_mwctl_prod_db... Data Source: apache-mwctl-test-db-a0... (183 records) Preview	cms_tmb_rates Data Source: cms_tmb_rates.csv (820.6k records)  Preview	cms_runs_for_fill Data Source: cms_runs_for_fill.csv (2.5k records) Preview	cms_lhcfills Data Source: cms_lhcfills.csv (621 records)  Preview
---	--	---	---	---	---	--

Data Transformation UI - ETL

Transformation Editor

Use refinement state as a conditional statement

Enable automatic typeahead

+ Functions

+ Attribute

1 diffDates(op_end_time, creation_time_utc, SECONDS)

Configure Output Settings

Apply to "op_duration"

Create New Attribute

New Attribute Name

op_duration_seconds

Data Type

Double

Single Assign

Cancel Preview Add to Script

TRANSFORM SCRIPT

- elogbook_fault_id - Transform
- fault_id - Create
- fault_classification_id - Transform

Commit to Project

1.8k Records 1.8k Filtered Records 21 Attributes

FAVORITES DATA TYPE NAME HIDDEN

All Attributes

Sort: By preview order

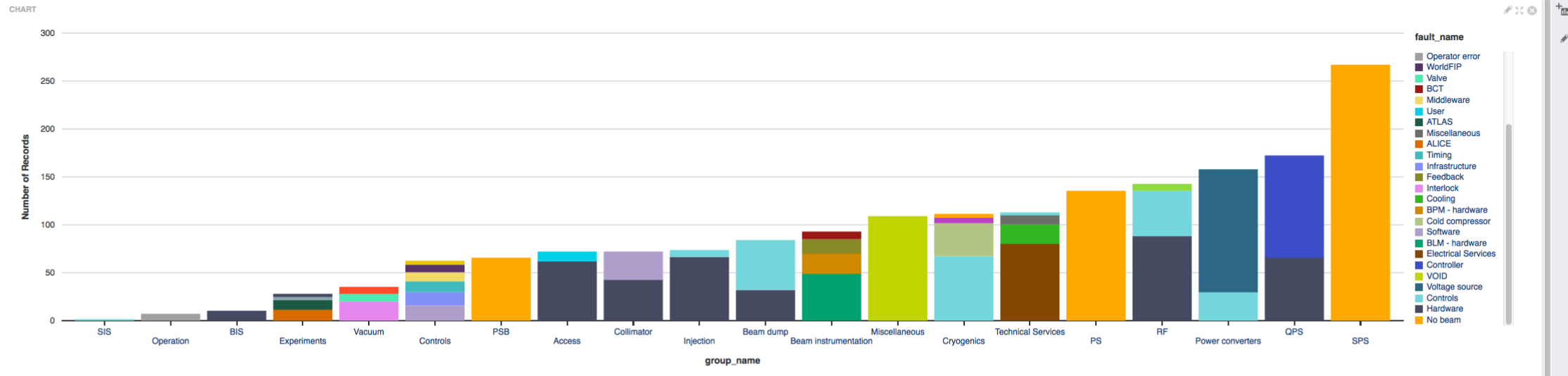
fault_id	fault_name	fault_state	group_name	is_root_cause	# op_duration	# op_duration_seconds	op_end_time	prevents_beam_op	prevents_injection	start_time
15160		OP_ENDED		N	4,266	10,932	2015-04-14 08:28:38 UTC	N	Y	2015-04-14 07:17:32 UTC
15223		CANCELLED		N				N	N	2015-04-15 06:30:00 UTC
secto... 3761		OP_ENDED		N	21,588	47,292	2014-11-24 17:30:38 UTC	N	N	2014-11-24 11:30:50 UTC
f 3765		OP_ENDED		N	99,735	66,405	2014-11-25 15:17:19 UTC	N	N	2014-11-24 11:35:04 UTC
roup.... 1909	Feedback	OP_ENDED	Beam instrumentation	N	66	56,455,621	2013-02-03 04:20:39 UTC	N	N	2013-02-03 04:19:33 UTC
condit... 3781		OP_ENDED		N	78,786	-9,846	2014-11-25 09:23:53 UTC	N	N	2014-11-24 11:30:47 UTC
roup.... 1908	No beam	OP_ENDED	SPS	N	20,870	56,457,295	2013-02-03 03:52:45 UTC	N	N	2013-02-02 22:04:55 UTC
roup.... 1924	No beam	OP_ENDED	SPS	N	147	55,550,285	2013-02-13 15:49:35 UTC	N	N	2013-02-13 15:47:08 UTC
roup.... 1907	Infrastructure	OP_ENDED	Controls	N	6,760	56,538,706	2013-02-02 05:15:54 UTC	N	N	2013-02-02 03:23:14 UTC
roup.... 1471	VOID	OP_ENDED	Miscellaneous	N	2,614	73,416,728	2012-07-21 20:55:32 UTC	N	N	2012-07-21 20:11:58 UTC
roup.... 1923	Hardware	OP_ENDED	Access	N	3,985	55,624,990	2013-02-12 19:04:30 UTC	N	N	2013-02-12 17:58:05 UTC
roup.... 1889	No beam	OP_ENDED	SPS	N	4,292	56,956,446	2013-01-28 09:13:34 UTC	N	N	2013-01-28 08:02:02 UTC
roup.... 1906	No beam	OP_ENDED	PSB	N	8,765	56,548,197	2013-02-02 02:37:43 UTC	N	N	2013-02-02 00:11:38 UTC
roup.... 117	Hardware	OP_ENDED	RF	N	3,750	137,030,838	2010-07-16 14:20:22 UTC	N	N	2010-07-16 13:17:52 UTC
roup.... 1470	No beam	OP_ENDED	PSB	N	1,322	73,419,672	2012-07-21 20:06:28 UTC	N	N	2012-07-21 19:44:26 UTC
roup.... 1922	Hardware	OP_ENDED	Injection	N	14,065	55,629,441	2013-02-12 17:50:19 UTC	N	N	2013-02-12 13:55:54 UTC
roup.... 1888	No beam	OP_ENDED	PSB	N	7,090	56,970,250	2013-01-28 05:23:30 UTC	N	N	2013-01-28 03:25:20 UTC
roup.... 1905	Controls	OP_ENDED	Beam dump	N	2,773	56,556,800	2013-02-02 00:14:20 UTC	N	N	2013-02-01 23:28:07 UTC
roup.... 1774	Cold compressor	OP_ENDED	Cryogenics	N	3,307	63,595,623	2012-11-12 13:00:37 UTC	N	N	2012-11-12 12:05:30 UTC
roup.... 116	User	OP_ENDED	Access	N	5,859	137,044,782	2010-07-16 10:27:58 UTC	N	N	2010-07-16 08:50:19 UTC
roup.... 1469	Controls	OP_ENDED	Beam dump	N	4,625	73,421,636	2012-07-21 19:33:44 UTC	N	N	2012-07-21 18:16:39 UTC
roup.... 1921	Hardware	OP_ENDED	RF	N	20,968	55,624,990	2013-02-12 19:04:30 UTC	N	N	2013-02-12 13:15:02 UTC
roup.... 1887	No beam	OP_ENDED	SPS	N	2,107	57,082,801	2013-01-26 22:07:39 UTC	N	N	2013-01-26 21:32:32 UTC
roup.... 1633	Hardware	OP_ENDED	RF	N	6,148	67,696,057	2012-09-26 02:00:03 UTC	N	N	2012-09-26 00:17:35 UTC
roup.... 1904	No beam	OP_ENDED	SPS	N	318	56,594,555	2013-02-01 13:45:05 UTC	N	N	2013-02-01 13:39:47 UTC
roup.... 1773	Hardware	OP_ENDED	RF	N	10,413	63,601,399	2012-11-12 11:24:21 UTC	N	N	2012-11-12 08:30:48 UTC
roup.... 115	Controls	OP_ENDED	Power converters	N	3,815	137,052,743	2010-07-16 08:15:17 UTC	N	N	2010-07-16 07:11:42 UTC
roup.... 1468	No beam	OP_ENDED	SPS	N	2,733	73,441,164	2012-07-21 14:08:16 UTC	N	N	2012-07-21 13:22:43 UTC

Discovery Applications

Explore ▾ Transform ▾ Discover ▾

FCC RAMS v1: Faults

- REFINE BY
- ▶ AFT_CARDIOGRAM
 - ▼ AFT_FAULTS
 - ▶ creation_time_utc
 - ▶ description
 - ▶ element
 - ▶ elogbook_fault_id
 - ▶ fault_classification_id
 - ▶ fault_classification_name
 - ▶ fault_creation_source_name
 - ▶ fault_description
 - ▶ fault_id
 - ▶ fault_name
 - ▶ fault_state
 - ▶ group_name
 - ▶ is_root_cause
 - ▶ op_duration
 - ▶ op_end_time
 - ▶ prevents_beam_op
 - ▶ prevents_injection
 - ▶ start_time
 - ▶ system_name
 - ▼ Other
 - ▶ expert_duration
 - ▶ expert_end_time



RESULTS TABLE

General ▾ 0 RECORDS SELECTED VIEW OPTIONS ▾ ACTIONS ▾

	group_name	op_end_time (Year-Mo-...)	element	prevents_beam_op	fault_creation_source_...	elogbook_fault_id	fault_description	fault_name	creation_time_utc (Yea...	fault_classification_id	is_root_cause	fault_state	description
<input type="checkbox"/>		4/14/15		N	LHC Logbook	1041380	null		4/14/15	4	N	OP_ENDED	
<input type="checkbox"/>				N	AFT Web Application		test from Ben		4/15/15	8	N	CANCELLED	
<input type="checkbox"/>		11/24/14		N	AFT Web Application		PLC problem in Pt 4, all ...		11/25/14	5	N	OP_ENDED	
<input type="checkbox"/>		11/25/14		N	AFT Web Application		all patrols lost in point 4		11/26/14	5	N	OP_ENDED	
<input type="checkbox"/>	Beam instrumentation	2/3/13	OFSU crash	N	LHC Logbook	1039013	{\"fault_description\": {\"gro...	Feedback	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>		11/25/14		N	AFT Web Application		interlock access during c...		11/25/14	1	N	OP_ENDED	
<input type="checkbox"/>	SPS	2/3/13		N	LHC Logbook	1039002	{\"fault_description\": {\"gro...	No beam	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	SPS	2/13/13		N	LHC Logbook	1039123	{\"fault_description\": {\"gro...	No beam	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	Controls	2/2/13	logging	N	LHC Logbook	1038991	{\"fault_description\": {\"gro...	Infrastructure	11/18/14	1	N	OP_ENDED	Connection to LHC-OP-L
<input type="checkbox"/>	Miscellaneous	7/21/12	ofc down	N	LHC Logbook	1034471	{\"fault_description\": {\"gro...	VOID	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	Access	2/12/13	patrollost in PM45	N	LHC Logbook	1039117	{\"fault_description\": {\"gro...	Hardware	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	SPS	1/28/13		N	LHC Logbook	1038905	{\"fault_description\": {\"gro...	No beam	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	PSB	2/2/13		N	LHC Logbook	1038986	{\"fault_description\": {\"gro...	No beam	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	RF	7/16/10	Ch 8 on CIB.UA47.R4B1	N	LHC Logbook	1021242	{\"fault_description\": {\"gro...	Hardware	11/18/14	1	N	OP_ENDED	RF interlock on BIS
<input type="checkbox"/>	PSB	7/21/12		N	LHC Logbook	1034469	{\"fault_description\": {\"gro...	No beam	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	Injection	2/12/13	tune kicker	N	LHC Logbook	1039116	{\"fault_description\": {\"gro...	Hardware	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	PSB	1/28/13	cooling	N	LHC Logbook	1038904	{\"fault_description\": {\"gro...	No beam	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	Beam dump	2/2/13	Trigger/Retrigger fault	N	LHC Logbook	1038981	{\"fault_description\": {\"gro...	Controls	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	Cryogenics	11/12/12	Cold Compressor S34	N	LHC Logbook	1037943	{\"fault_description\": {\"gro...	Cold compressor	11/18/14	1	N	OP_ENDED	
<input type="checkbox"/>	Access	7/16/10		N	LHC Logbook	1021240	{\"fault_description\": {\"gro...	User	11/18/14	1	N	OP_ENDED	

Next steps

- Check new upcoming features presented during OOW 15
 - Scalability improvements
 - Streaming integration (Kafka)
 - Add more advanced charts and visualizations
 - Machine Learning integration
- Evaluate Big Data Discovery cloud
- Extend to more CERN use cases
 - Controls and Operations
 - Accelerator Fault Tracking
 - Diagnostics and Monitoring
 - IT Infrastructure Monitoring
 - Server logs analysis
 - Database latency
 - Human Resources



Conclusions

- Data visualization and discovery is an important area in data analytics
 - Facilitates users to visualize and explore their data
 - Find correlations, extract insight and useful information
- Important points
 - Flexible and user-friendly platform
 - Advanced data visualization and exploration
 - Collaborative
- Application to different domains
 - Controls and Operations
 - IT Infrastructure Monitoring
 - Human Resources



www.cern.ch