



CERN's Journey into Big Data and Analytics

Manuel Martin Marquez

Senior Data Scientists and Data Streaming Project Leader CERN – European Organization for Nuclear Research





2. 3. 4 5. 6

CERN Introduction

Why CERN cares about Big Data?

What does Big Data looks like? What are the benefits and risks?

The data process: jumping into Big Data

Profiting from the data investments

Making the most of it: Advance and scalable analytics

Some conclusions from our experince.



01. CERN – European Organization for Nuclear Research







01. CERN – A Worldwide Collaboration





Other States

Afghanistan	1	El Salvador 1 Pakistan		41	
Albania	2	Estonia 16		Palestine (O.T.).	4
Algeria	8	Georgia 36 P		Peru	8
Argentina	11	Gibraltar 1 Phil		Philippines	1
Armenia	25	Hong Kong 1 Saudi Arabia		Saudi Arabia	З
Australia	25	Iceland	ind 4 Senegal		1
Azerbaijan	8	Indonesia 1 Singapore		Singapore	2
Bangladesh	4	Iran 28 Sint Maarten		Sint Maarten	2
Belarus	47	Ireland 22 Sid		Slovenia	27
Bolivia	3	Jordan 2 South Afric		South Africa	16
Bosnia &		Kenya	1	Sri Lanka	5
Herzegovina	1	Korea, D.P.R.	1	Syria	2
Brazil	108	Korea Rep.	117	Thailand	12
Cameroon	1	Kuwait	1	T.F.Y.R.O.M.	1
Canada	134	Lebanon	banon 12 Tunisia		6
Cape Verde	1	Lithuania	19 Ukraine		55
Chile	12	Luxembourg	4	Uzbekistan	4
China	280	Madagascar	4	Venezuela	9
China (Tapei)	45	Malaysia	15	Viet Nam	9
Colombia	30	Mauritius	/////	Zimbabwe	//2
Croatia	35	Mexico	64		
Cuba	7	Montenegro	3		
Cyprus	16	Morocco	/ 12		
Ecuador	3	Nepal	5		1144
Egypt	19	New Zealand	1/17		1475

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

Distribution of All CERN Users by Nationality on 14 January 2014



01. Fundamental Research – The Building Blocks of the Universe





13.8 billion years



01. The Large Hadron Collider – A Marvel of Technology



ALICI

World's largest scientific instrument 27km, 6000+ superconducting magnets

Fastest racetrack on Earth

Protons circulate 11245 times/s (99.9999991% the speed of light)

Emptiest place in the solar system High vacuum inside the magnets

Hottest spot in the galaxy During Lead ion collisions create temperatures 100 000x hotter than the heart of the sun;

HCh

CERN Prévess



ATLAS

01. The Large Hadron Collider – An Extremely Complex Instrument



CERN Accelerator Complex is unique installation Therefore, we have to face unique challenges

Control and Operations

Million of sensors, large number of control devices, front-end equipment, etc. Many critical systems: Cryogenics, Vacuums, Machine Protection, etc.



01. The Experiments – ATLAS, CMS, ALICE and LHCb Detectors





150 Million of sensor Control and detection sensors

Massive 3D camera

Capturing 40+ million collisions per second



01. The Experiments - Millions of Collisions Per Second







2.

3.

4

5.

6



CERN Introduction

Why CERN cares about Big Data?

What does Big Data looks like? What are the benefits and risks?

The data process: jumping into Big Data

Profiting from the data investments

Making the most of it: Advance and scalable analytics

Some conclusions from our experince.



02. What are the Challenges Ahead - Toward High Luminosity

LHC / HL-LHC Plan





The constant to push the machine to its limits

Exiting future for scientists but a huge challenge for the engineers



CERN

02. CERN's Database Status – More Data and Faster







02. CERN's Accelerator Logging Service – Control IoT System







02. CERN's Accelerator Logging Service – Data Evolution





Storage Evolution - Size in GB / day



02. The Evolution Towards Big Data



What we do have...



- A cost-effective system
- The functionalities covers more than 90% of current requirements
- Data access with SQL, custom Java APIs and applications
- ...then why to change?
 - Better performance and unlock bigger datasets
- Leverage analytics capabilities: new data exploitation approaches
- More heterogeneous data access models



3.

4

5.

6



- **CERN Introduction**
- 2. Why CERN cares about Big Data?
 - What does Big Data looks like? What are the benefits and risks?
 - The data process: jumping into Big Data
 - Profiting from the data investments
 - Making the most of it: Advance and scalable analytics
 - Some conclusions from our experince.











The new architecture fulfils current and future requirements...





Heterogeneity of clients:

- Not limiting the accessing technologies
- Different expertise



More and better ways to exploit the data

- Bigger datasets
- Real-time or streaming analytics,
- More complete and informative decisions



Leverage Collaboration

- Across different locations
- Across academic and processional backgrounds





...but what are the costs and risks?

- Human and expertise
- Economical costs underestimation
- Overestimation of expectation
- Lack of maturity of the technologies to run critical services
 - Evolution and changing pace of the technologies





...how minimize them?

- Academic programs have quickly adapted to the expertise need
- Avoid attachments to technologies by decupling services
- Strategies to find the balance between productivity and exploration
- Same old problems will show up use your expertise
- Cloud, Cloud and Cloud:
 - Fast deployments
 - Focus on solutions and not architecture
 - Minimize the costs



3.

4

5.

6



- **CERN Introduction**
- 2 Why CERN cares about Big Data?
 - What does Big Data looks like? What are the benefits and risks?
 - The data process: jumping into Big Data
 - Profiting from the data investments
 - Making the most of it: Advance and scalable analytics
 - Some conclusions from our experince.



04. Future Circular Collider: Enhance Collaboration with Big Data







04. The Data Process: Making Data Accessible







04. The Process: Offering the Right Tools – Data to Value



25

Data Analysis, visualization and discovery involves several steps:

- Data integration, transformation, exploration, and discovery
- Agile data governance

Reliability Availability & Maintainability for CERN's Accelerator Complex





3.

4

5.

6



- **CERN Introduction**
- 2. Why CERN cares about Big Data?
 - What does Big Data looks like? What are the benefits and risks?
 - The data process: jumping into Big Data
 - Profiting from the data investments
 - Making the most of it: Advance and scalable analytics
 - Some conclusions from our experince.



05. The Data Process - Data to Value how to...

CERN

The goal...



- Decentralize the data process
- Push the responsibility of the process to the people who better knows it

...but is that possible?



- Use the right tools flexible, intuitive and almost-zero learning curve
- Hide the technology complexity
- Help to focus on the data that matters



Allow a collaborative environment



05. The Data Process - Hiding Technology Complexity



÷



Unlock the exploitation of data from different perspectives Avoid technological barrier

Add an abstraction layer that minimize the risks

ORACLE Big Data Discovery

re 🗸 Transform 🗸 Discover 🗸		Lo	gbook BDI): logbook				¢ '	* = =
OGBOOK	logbook_complex TESTS	0							
REFINE BY 🌣 🗙	Transformation Editor								
logbook_complex	Use refinement state as a condition	al statement			Functions	Attribute	Apply to "logbook_logbook"		
Q Filter ▼	1					>	Create New Attribute		
LHC	1			Double-click or drag to inser	ta				
PS Complex				function					1
SPS	1						Lintor toxt		
ISOLDE	1			Q Filter			Data Type		
HISTORY	1			abs(long)	^		String	~	
BE-BI				acos			Single Assign		
UA9				addTime(date,integer,string)					
TE-ABT	1			asin(double)					
TE				atan(double)			Cancel Preview	Add to Sc	rip
	1			atan2(double double)					
TE-VSC	725K of 725K 525K	41	🔶 FAVORIT					-	
R2E	Records Sampled Records Viewe	ad Attributes	~	ardinality(object)				-0	
EN-ICE	All Attributes			cbrt(double)			Nort: By name 🗸		=
AT-MCS	≡ II logbook line name	E Iogbook_logbook		ceil(double)		record id =	■ Iogbook shift end	≡ # lo	gb
Select All			-	concat(string[])		_			-
Select All	IPROTONPHYI	LHC OR		concatWithToken(string,string[])			2012-08-09 23:00:00 LITC	1 046	. 70
logbook_devices_accelerator	[PROTONPHY]			contains(string,string)			2012-08-09 23:00:00 UTC	1.046	70
logbook_devices_class_description				copySian(double.double)			2013-02-14 07:00:00 UTC	1.051	96
logbook_devices_class_version	[PBOTONPHY]	LHC OP		copySign(float float)			2012-08-10 07:00:00 UTC	1,046	.70
logbook_devices_classname	[PBOTONPHY]						2013-02-14 07:00:00 UTC	1.051	.96
logbook_devices_description	[PBOTONPHY]	LHC OP		E cos(double)	~		2012-08-10 07:00:00 UTC	1.046	, 20
logbook_devices_devicename		LHC OP		■ PAGE 1 OF 3 ▶ ▶			2012-08-13 07:00:00 UTC	1 046	78
logbook_devices_fecname	[PBOTONPHY]						2012-07-12 07:00:00 UTC	1.045	.80
logbook_devices_implementation	[PBOTONPHY]	LHC OP		Unselected line	1,717,723		2013-02-14 07:00:00 UTC	1.051	.96
logbook event comment	[PROTONPHY]	LHC OP		STABLE BEAMS	1,538,448		2012-11-25 07:00:00 UTC	1.049	.98
logbook event date	[PROTONPHY]	LHC OP		Unselected line	2,176,997		2012-07-27 23:00:00 UTC	1,046	.29
logbook event id	[PROTONPHY]	LHC OP		Unselected line	1,541,296		2012-08-10 15:00:00 UTC	1.046	.71



05. The Data Process - Fishing into the (Data) Lake



Help to focus on the data that mattersWhat do I really have available?

ORACLE Big Data Discovery Search Q, Logbook BDD: logbook Explore v Transform v Discover v LOGBOOK logbook_complex 0 TESTS Clear All 0 Ģ REFINE BY \times 725K of 725K 525K 41 \boxtimes 0 NAME V + FAVORITES DATA TYPE V **M HIDDEN Records Sampled Records Viewed** Attributes logbook_complex \mathbf{x} DISPLAY: Favorites @ Sort: Relationship to logbook_complex V Q Filter \sim LHC **PS** Complex SPS LHC OP Herve Michel Genoud LHC BE-OP-PS ISOLDE PS Complex SPS Richard Scrivens BE-OP-PSB HISTORY SPS PS BE-OP-SPS Fabrice Chapuis ISOLDE BE-ABP-HSL PSB Bettina Mikulec BE-BI BE-BE-EB HISTORY CTE Bodolphe Maillet UA9 BE-BI ADE Jose-Luis Sanchez Alvarez BE-OP-AD TE-ABT 9 others 75 others 85 others 19 others TE logbook_complex logbook_logbook logbook_jira_reporter logbook_jira_reporter_... BF SHUTDOWN TE-VSC R2E EN-ICE AT-MCS BE-CO-APS CPS Resolved BE Select All BE-CO-FE PS Closed GS BE-CO-HT PSB Done **BE-OP-AD** TE loobook devices accelerator BF Open BE-CO-DO OP EN Iogbook_devices_class_description In Progress BE-OP-PS Cryogenics PH 20 others 61 others Reopened Iogbook devices classname logbook_jira_assignee... logbook_fault_groupn... logbook_jira_assignee... logbook_jira_status Ioabook devices description Iogbook devices devicename Iogbook_devices_fecname Iogbook devices implementation logbook event comment Unselected line [PROTONPHY] logbook_event_date BUN [BEAM] logbook_event_id BEAM IN [BeamSetup] SETUP [MD] logbook_fault_description ICNGS1 SETLONG2 SETTING LIP logbook +



05. The Data Process - Sharing the Findings



Use the right tools – flexible, intuitive and almost-zero learning curve

Push the responsibility to explore the data to the people who better knows it







So we already have achieved...



Making the data available and integration of data sources

The data is explored by the people who better knows it



We can easily share findings and explorations

...but all that is part of the advance analytics process





Allow more effective advance analytics prototypes

Row Data -> Integration -> Transformation -> Focus -> Create Prototype



3.

4

5.

6



- **CERN Introduction**
- 2. Why CERN cares about Big Data?
 - What does Big Data looks like? What are the benefits and risks?
 - The data process: jumping into Big Data
 - Profiting from the data investments

Making the most of it: Advance and scalable analytics

Some conclusions from our experince.





Lets apply predictive maintenance to the LHC control system..

Aperture order

Aperture measured

Difference = measured - order



Row Data -> Integration ->Transformation -> Focus -> Create Prototype





06. Making the most of it: Scalable Advance Analytics

CERN

...nice prototype but now I need to scale



Need for a new set of skills and therefore profiles





3.

4

5.

6



- **CERN Introduction**
- 2. Why CERN cares about Big Data?
 - What does Big Data looks like? What are the benefits and risks?
 - The data process: jumping into Big Data
 - Profiting from the data investments
 - Making the most of it: Advance and scalable analytics
 - Some conclusions from our experince.



05. Some Conclusions



- Embark into Big Data is a risk but the risk of not doing it is much higher
- Unlocks new ways to exploit your investment on Data
 - Overcomes technical limitations
 - Allows more heterogeneous data access



- But making the data is available is just half of the way
 - Profit from data is a full process
 - Self-service tools enhance collaboration
 - Reduce IT intervention improve productivity
 - Scale your advance analytics



- Cloud is essential
 - Fast and flexible deployments
 - Focus on solutions and not architecture
 - Minimize the risks and costs



And please do not underestimate costs or overestimate expectations



