
ATLAS Data Management: Evolution Challenges

Vincent Garonne



Kicktipp / ADC Euro2016

Pos.	Nom	Journées									Total		
		<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>16</u>	<u>Qu</u>	<u>Se</u>	<u>Fi</u>	<u>Bon.</u>	<u>Victoires</u>	<u>Tot</u>
1.	IPERLUCA	17	16	4							0	1,33	37
2.	Tomas.Kouba	23	12	2							0	1,00	37
3.	TheCrow	13	13	2							0		28
4.	Ivan	14	10	3							0		27
5.	vincent	11	8	3							0		22
6.	tartan_army	10	6	4							0	0,33	20
7.	tjavurek	7	7	2							0		16
8.	AleDiGGi	3	6	4							0	0,33	13

ATLAS DDM: Current numbers

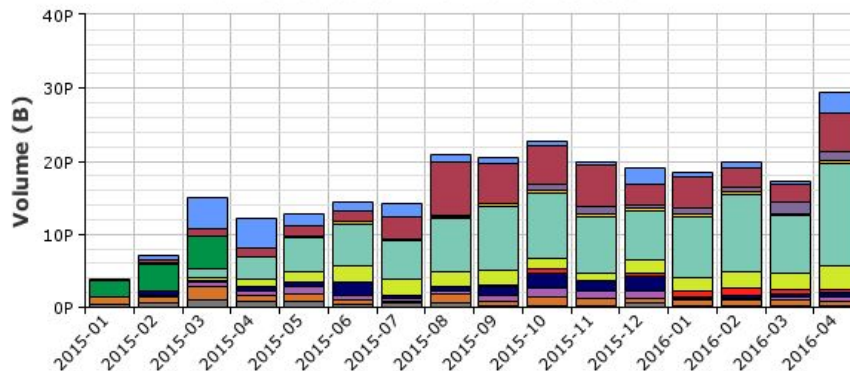
The ATLAS DDM System/Rucio has demonstrated very large scale data management:

- 2000 ATLAS users
- 200 PB on 130 sites
- 1B file replicas
- 2.5 M file transfers/Day, 1.5 PB/Day
- 6 M deleted files/Day, 2 PB/Day
- 1M jobs/day

It's a lot of data !

Average ATLAS traffic: 10GB/s

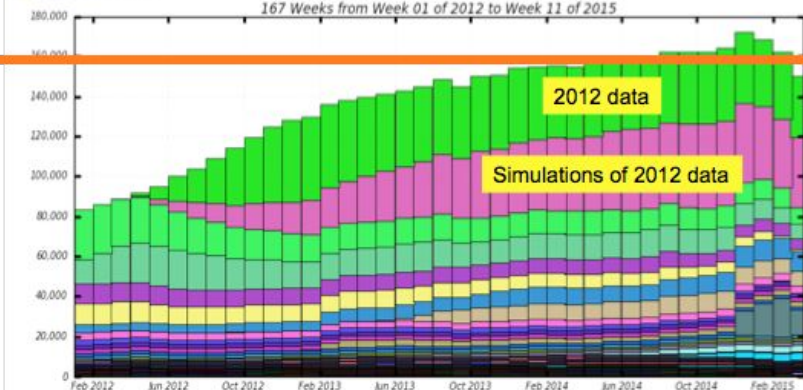
dashboard **Transfer Volume**
2015-01-01 00:00 to 2016-05-01 00:00 UTC



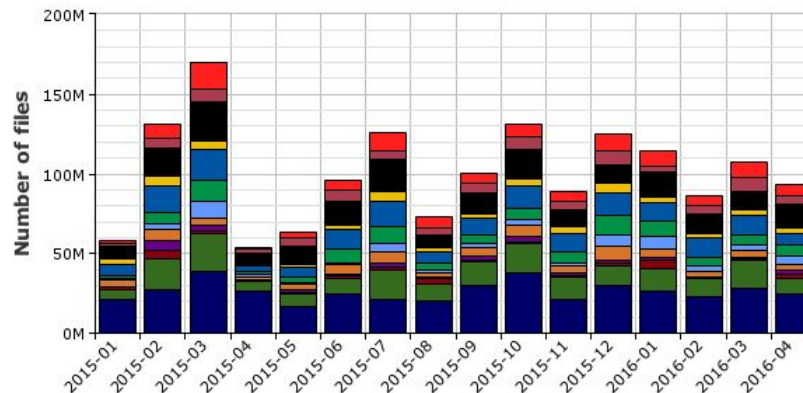
Activities



dashboard **Number of Physical Bytes (in TBs)**
167 Weeks from Week 01 of 2012 to Week 11 of 2015



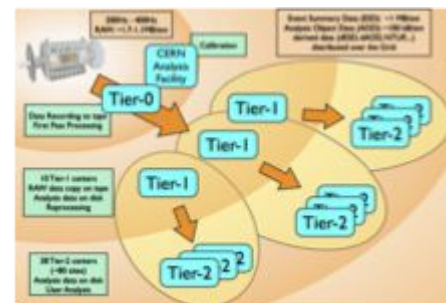
dashboard **Deletion Successes**
2015-01-01 00:00 to 2016-05-01 00:00 UTC



Destinations

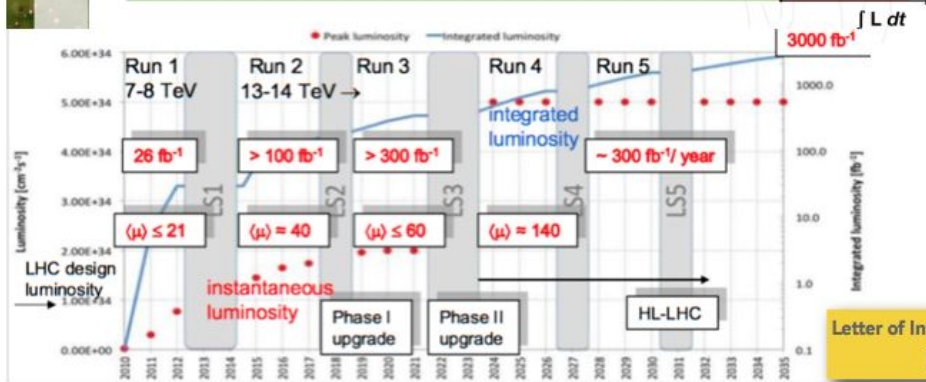
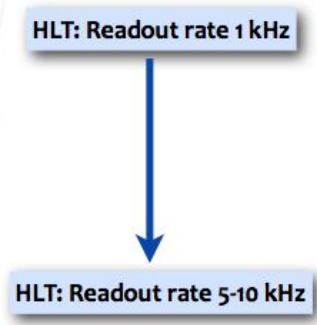
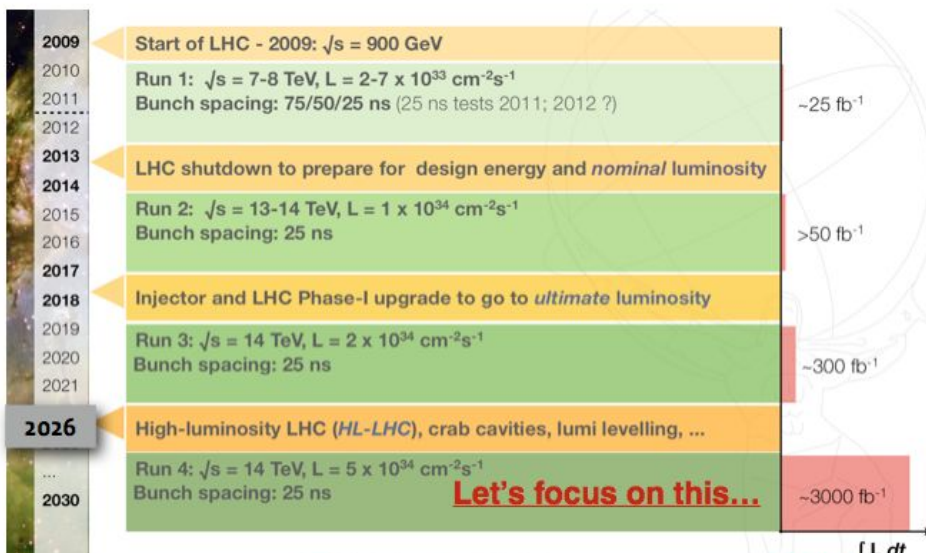


Current Challenges



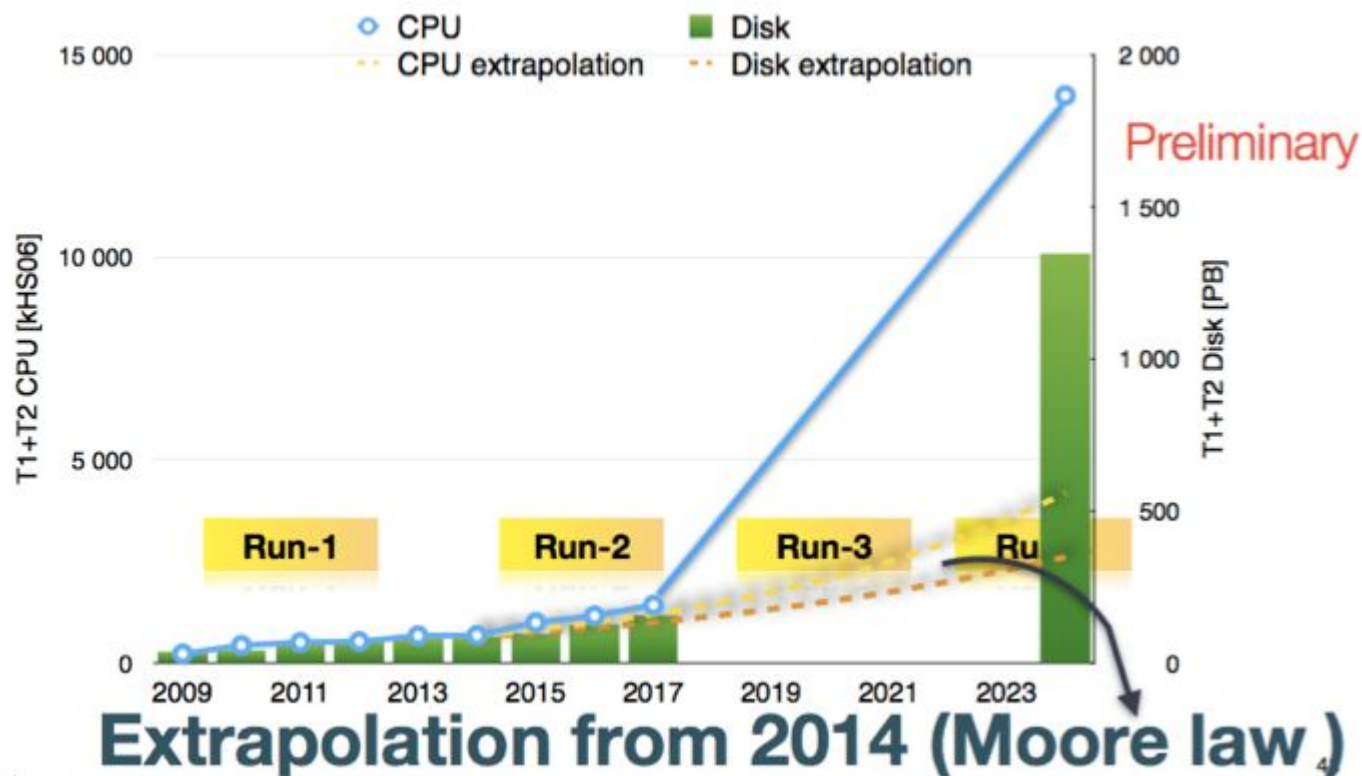
- New trends in data management
 - Original model was based on network being the weak point
 - But network has proven to be cheaper and better than expected
 - Break the rigid hierarchical model of data flow and sending jobs to data
 - Dynamic data placement
 - Remote data access over wide area network
- Event-level workflow instead of file-level
- Need more CPU and disk but with flat budget → opportunistic resources
 - High Performance Computing (supercomputers)
 - Volunteer Computing (general public)
 - Cache

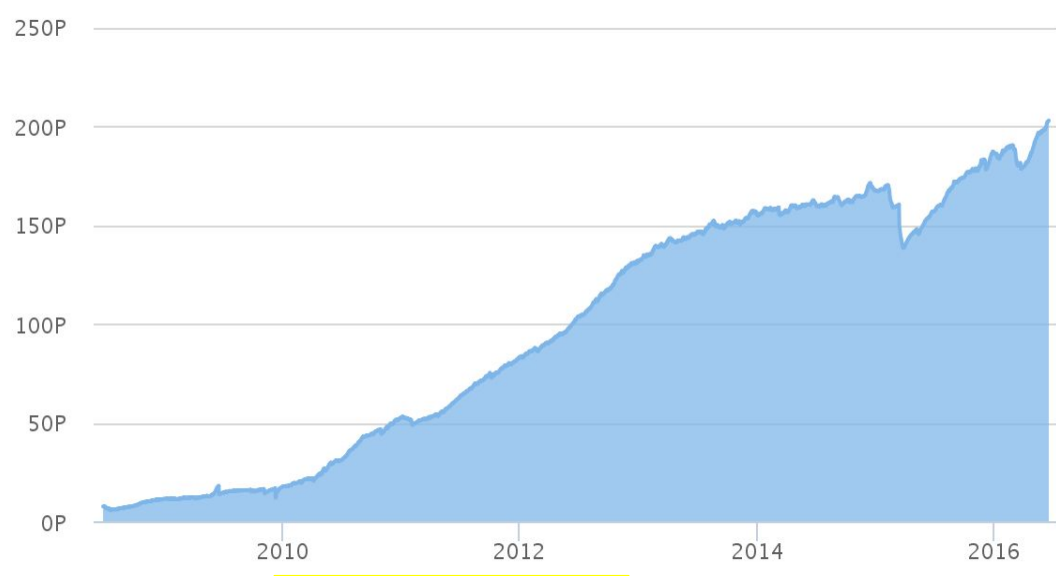
LHC Upgrade Timeline



Letter of Intent for the Phase-II Upgrade of the ATLAS Experiment
<https://cds.cern.ch/record/1502664>

ATLAS Resource Needs at T1s & T2s





Run-1

Run-2

— PMS

300P

250P

200P

150P

100P

50P

0P

2010

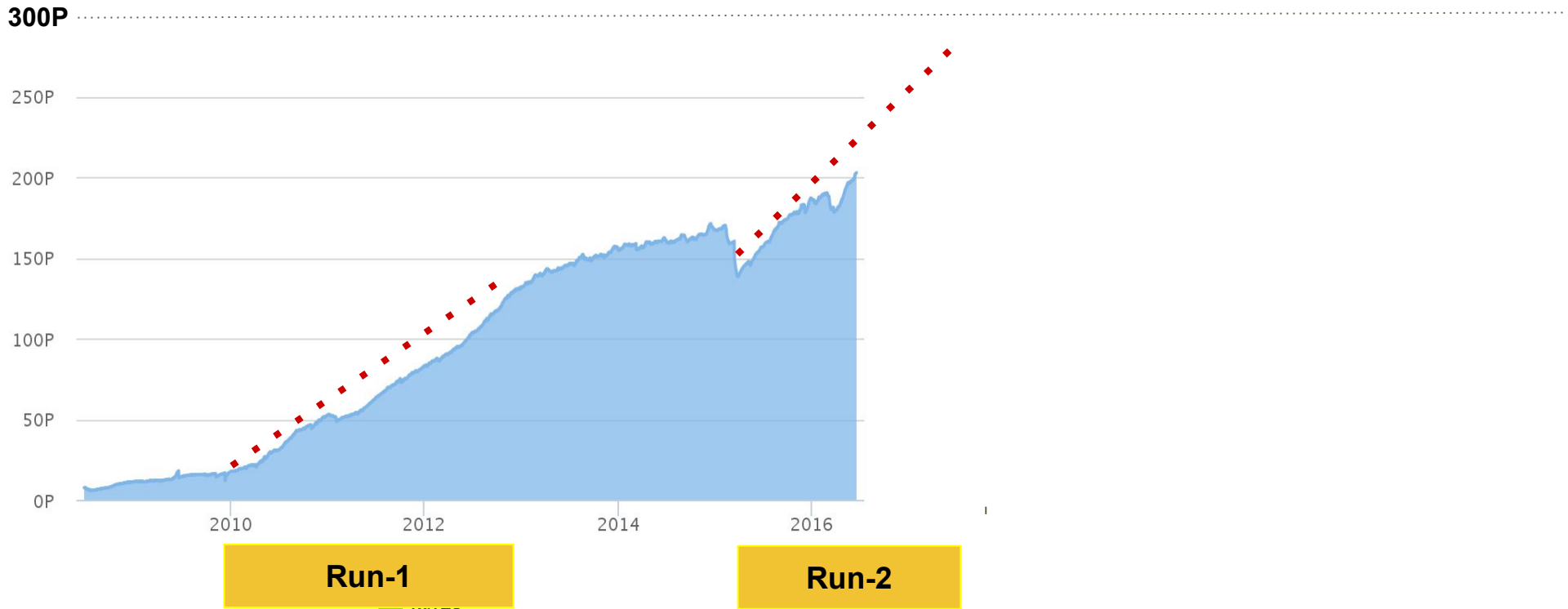
2012

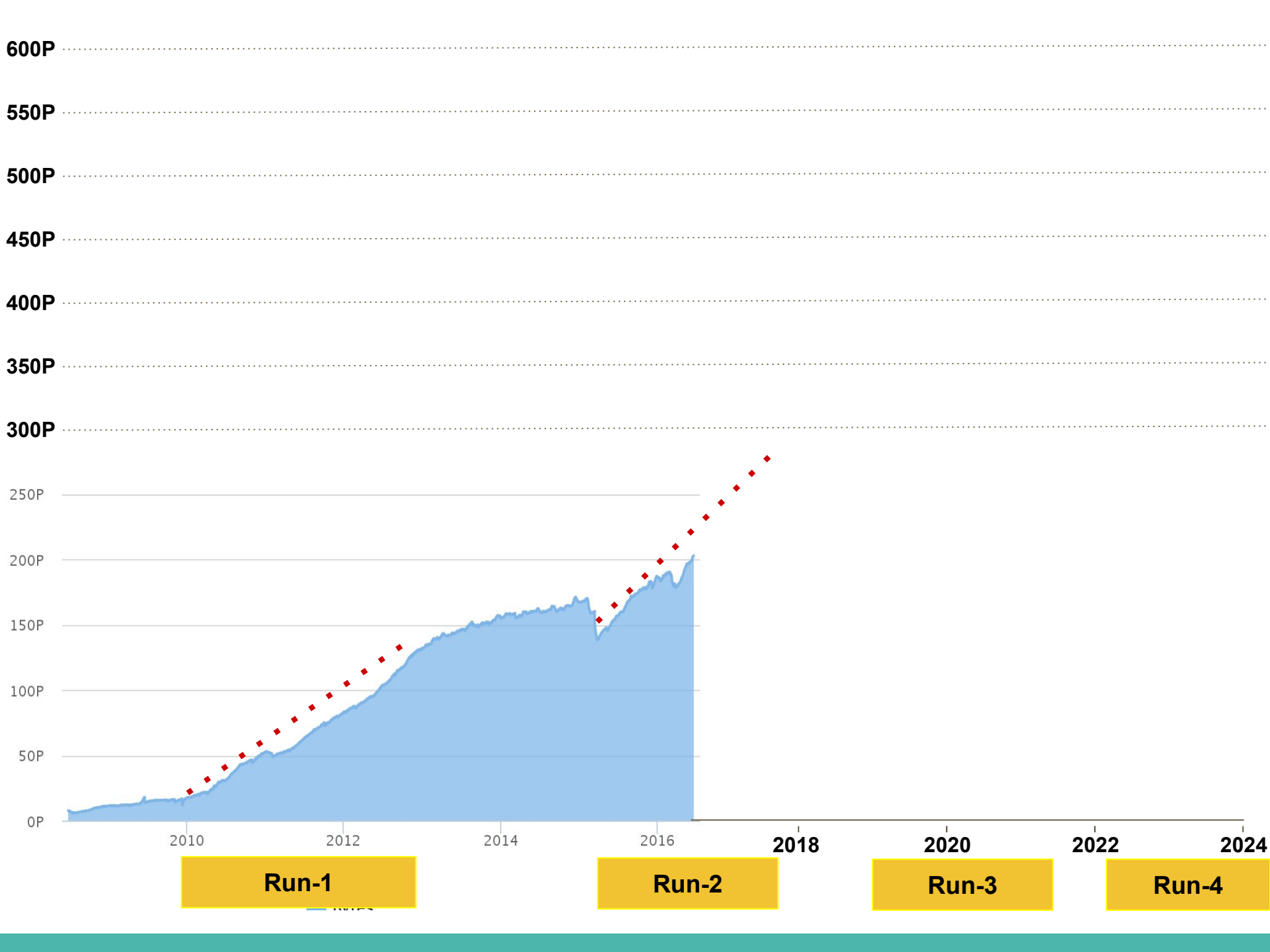
2014

2016

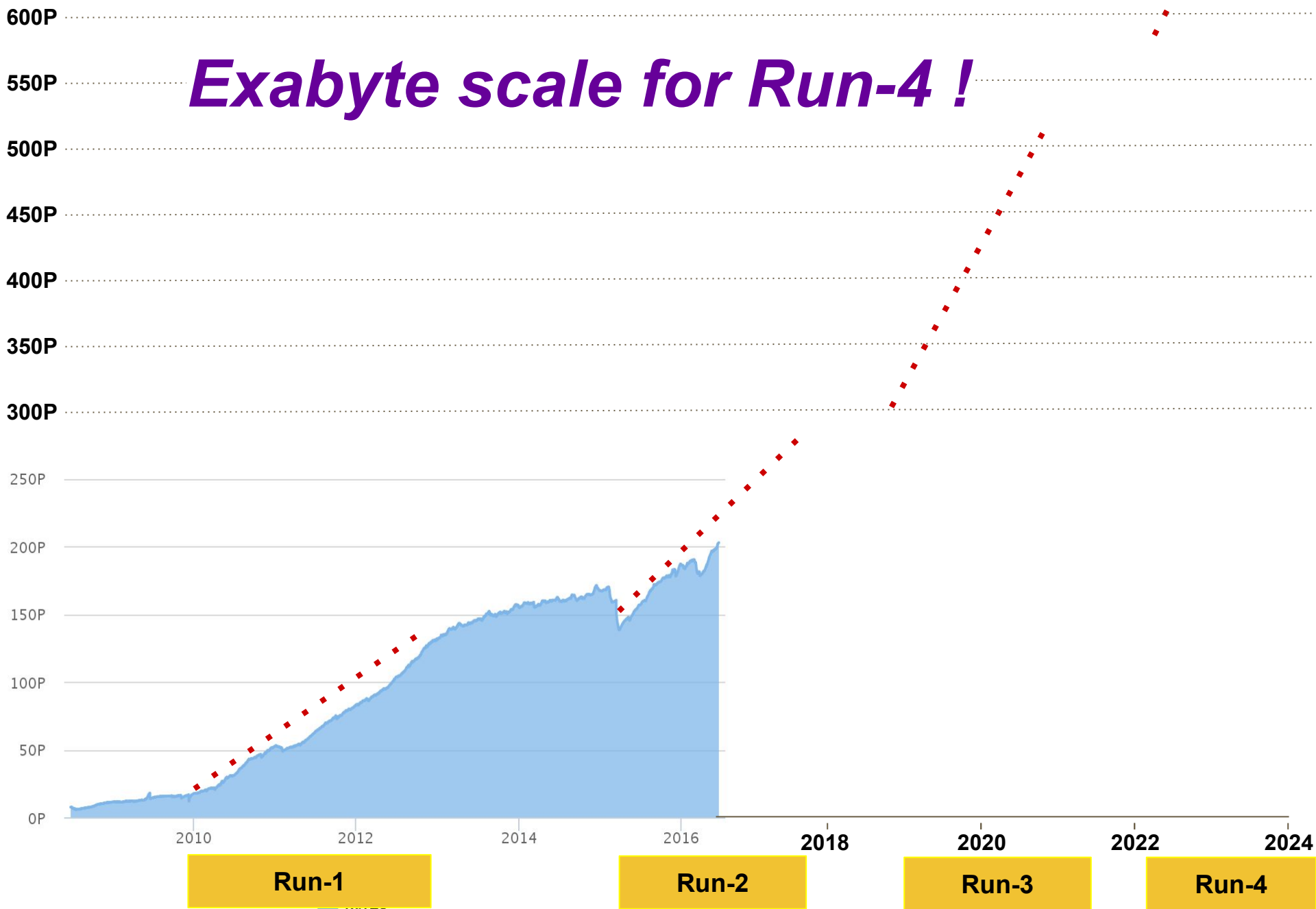
Run-1

Run-2





Exabyte scale for Run-4 !



ATLAS DDM: Current Numbers (2016)

The ATLAS DDM System/Rucio has demonstrated very large scale data management:

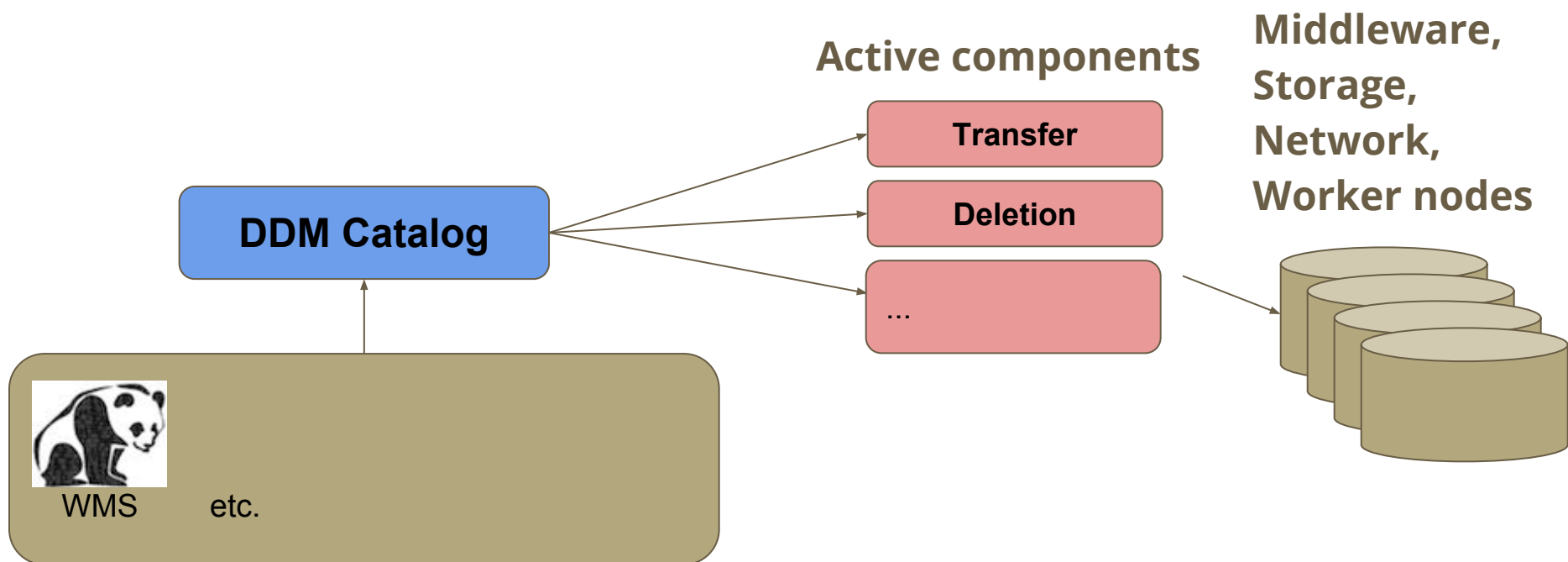
- 2000 ATLAS users
- 200 PB on 130 sites
- 1B file replicas
- 2.5 M file transfers/Day, 1.5 PB/Day
- 6 M deleted files/Day, 2 PB/Day
- 1M jobs/day

ATLAS DDM: Future Numbers (2026)?

The ATLAS DDM System/Rucio has **will** demonstrated very large scale data management:

- 2000 ATLAS users
- 200**0** PB on 130 sites
- 1**0** B file replicas
- 2**5** M file transfers/Day, 1**5** PB/Day
- 6**0** M deleted files/Day, 2**0** PB/Day
- 1**0** M jobs/day

DDM Evolution: Current Logical Overview



Future of Catalog ?

- Most of DDM implementations for the LHC are based on catalog
 - It's convenient to have one global and fast index for job scheduling
 - Easier to manage, few misses and availability > 99.99 %
- DDM has to scale with the (cumulative) number of data objects and operations
 - Data object(s) can be containers(s), dataset(s), file(s), event(s)
 - Most of the operations are generated by the workload management system, i.e., PanDA

Catalog Scalability & Database Technologies

- For the scalability, it'll follow the advances in databases, open and standard technologies
- Rucio has a flexible design with no dependence on particular RDBMS implementation

- Relational database management system
 - Use cases: Real-time data, transactional

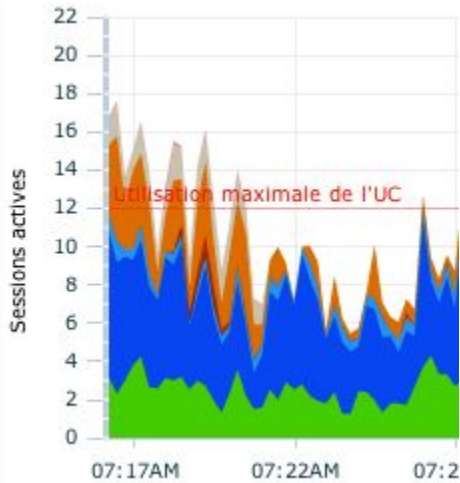


- Non relational structured storage
 - Use cases: Popularity, accounting, log analysis, ML

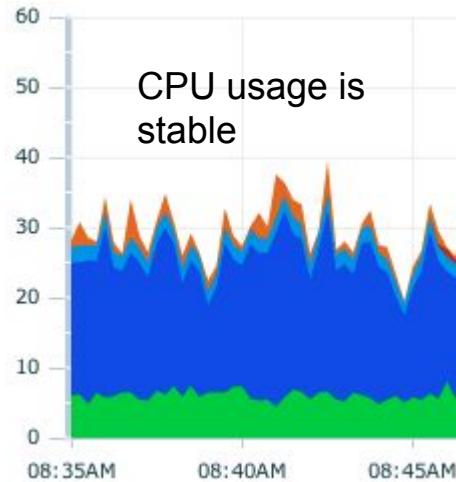


DB: Scaling tests - Load Increase

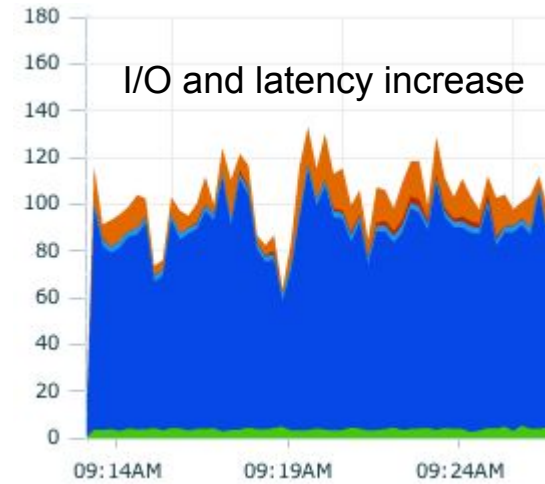
Load X1



X2



X4



Latency



Throughput

900 I/O ops.s-1
20 MB.s-1

2000 I/O ops.s-1
30 MB.s-1

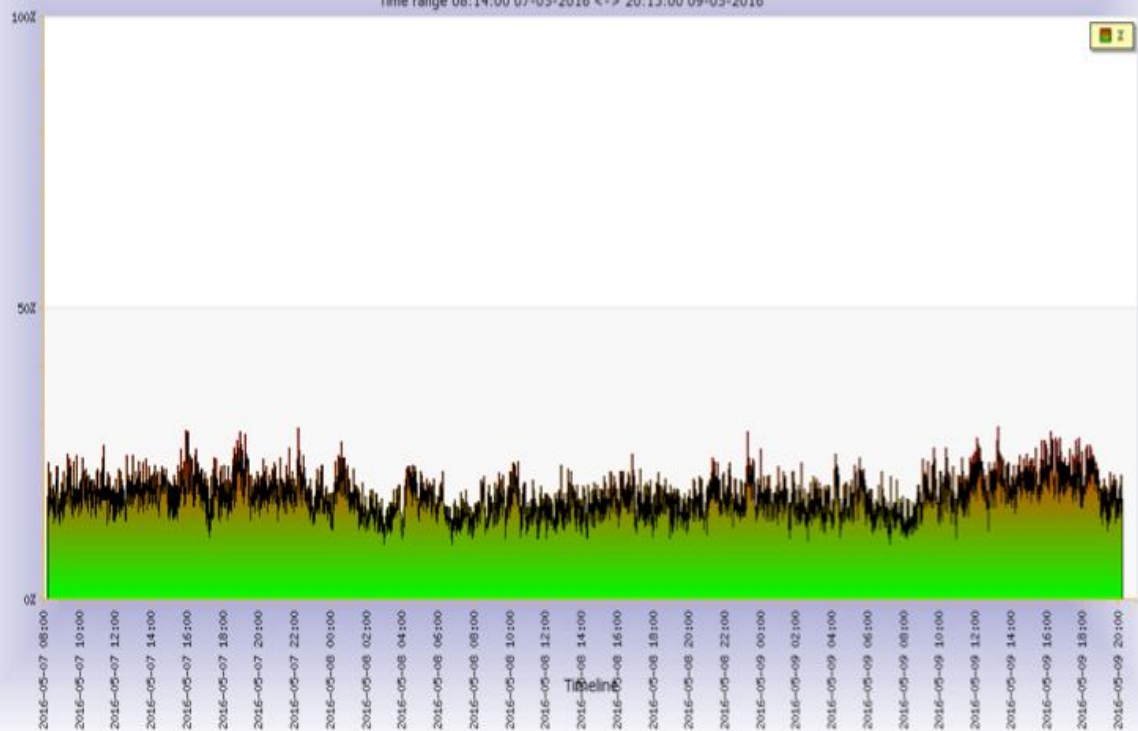
3200 I/O ops.s-1
40 MB.s-1

Oracle: Current Status

Thanks to various improvements in the Rucio code and in the DB objects, Rucio load on the database is not high.

CPU Utilization
3@ADCR.CERN.CH

Generated on Monday 09th of May 2016 08:14:52 PM
Time range 08:14:00 07-05-2016 <-> 20:13:00 09-05-2016



DBAs remarks:

- The Rucio load on the DB is acceptable. Keeping it low is important.
- A factor 4 more load is possible.
- New HW for the ADCR database is expected to be in place in Q1 of 2017 to serve the ADC applications from 2017 to 2020.
 - Factor 10 doable ?

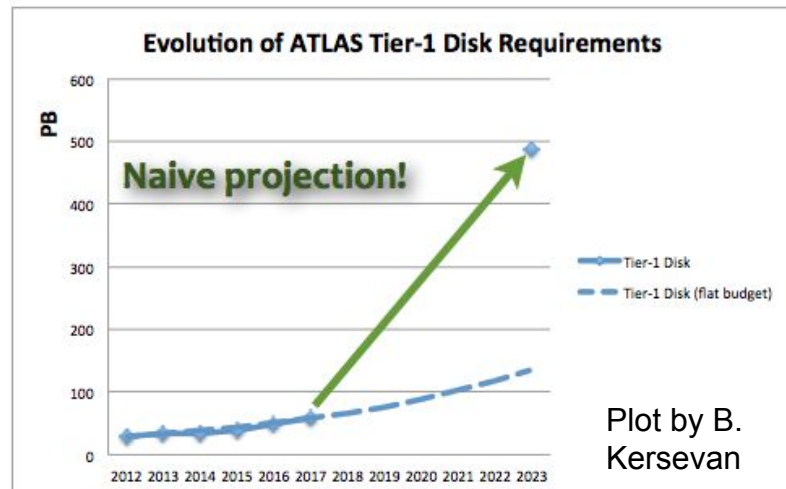
Personal Comments

- I have observed a general trend to have more and more (physics) metadata in DDM to facilitate data selection, discovery and analytics
 - Should DDM and the metadata part be strongly coupled ?
- Production, analysis and workload management workflows have a direct influence on the system' scalability
 - One scheduling change can have butterfly effect on DDM
 - Is the current strong separation between DM and WM an affordable model for run-4 ?

Storage & Network

- Resource increase w.r.t. 'flat budget':
- In ten years from now, the terabyte price on disk might be divided by 3
- Tapes foreseen to be viable for both capacity and cost
- Looking at the trends, the world gained an order of magnitude of storage over the last five years...
- We have to keep up to date with data storage/transfer technology and evolution
- The biggest predictable gain will come from network and will strongly influence the experiments computing models

⇒ **Data storage of HL-LHC experiments will NOT become a 'trivial' problem ten years from now**



Summary

- Many questions, no immediate answers