

LHCC Referee Meeting

March 3, 2015

ALICE Status Report

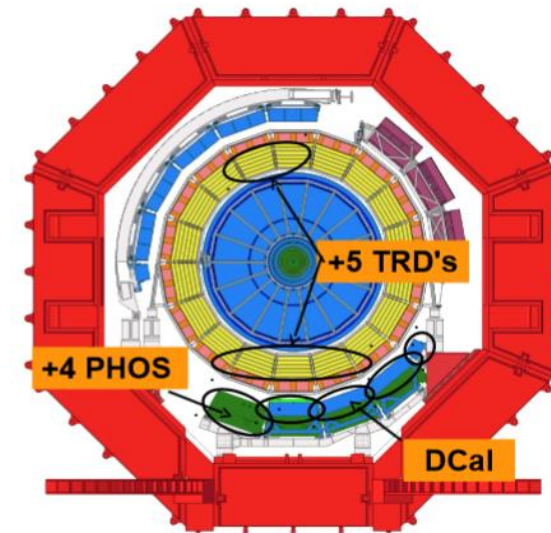
Predrag Buncic

CERN



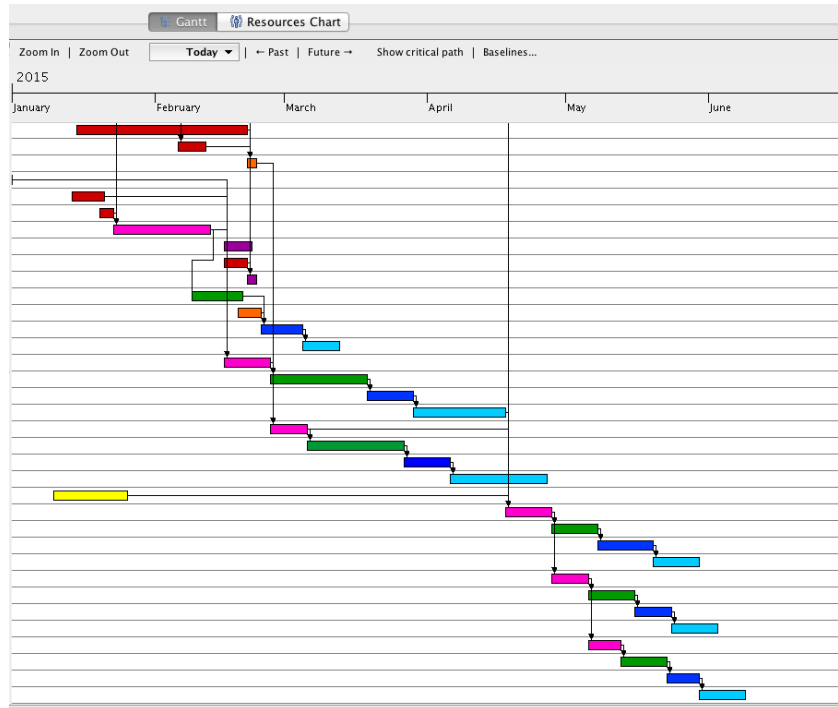
ALICE

- Expecting increased event size
 - 25% larger raw event size due to the additional detectors
 - Higher track multiplicity with increased beam energy and event pileup
- Concentrated effort to improve performance of ALICE reconstruction software
 - Improved TPC-TRD alignment
 - TRD points used in track fit in order to improve momentum resolution for high p_T tracks
 - Streamlined calibration procedure
 - Reduced memory requirements during reconstruction and calibration (~500Mb, the resident memory is below 1.6GB and the virtual - below 2.4 GB)





Re-processing & re-commissioning

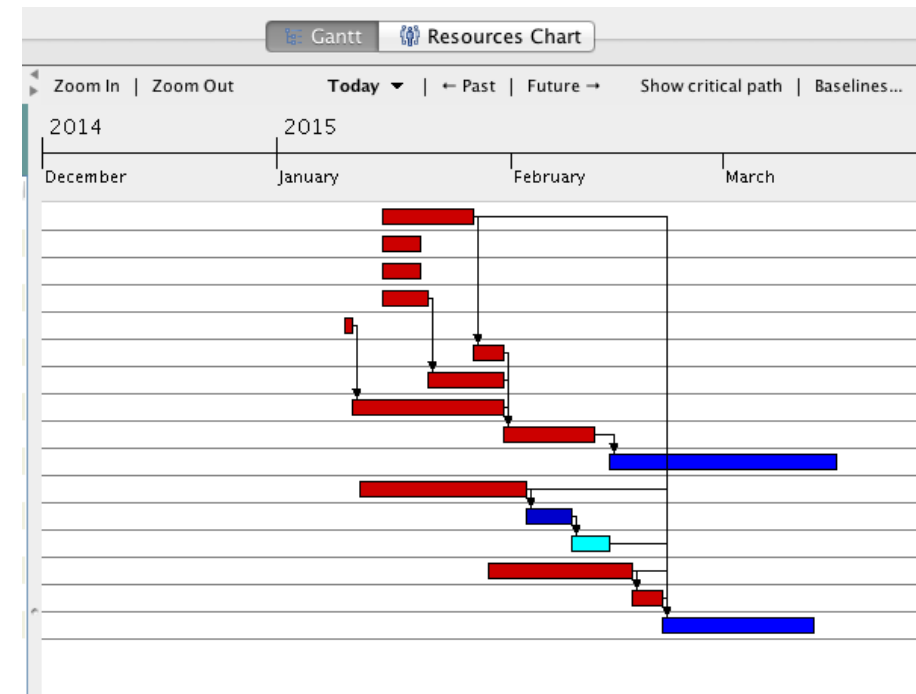


• Re-processing (Run 1)

- Steady RAW and MC activities
- Full detector re-calibration and 2 years worth of software updates
- All Run 1 RAW data processing with the **same** software

• Re-commissioning (Run 2)

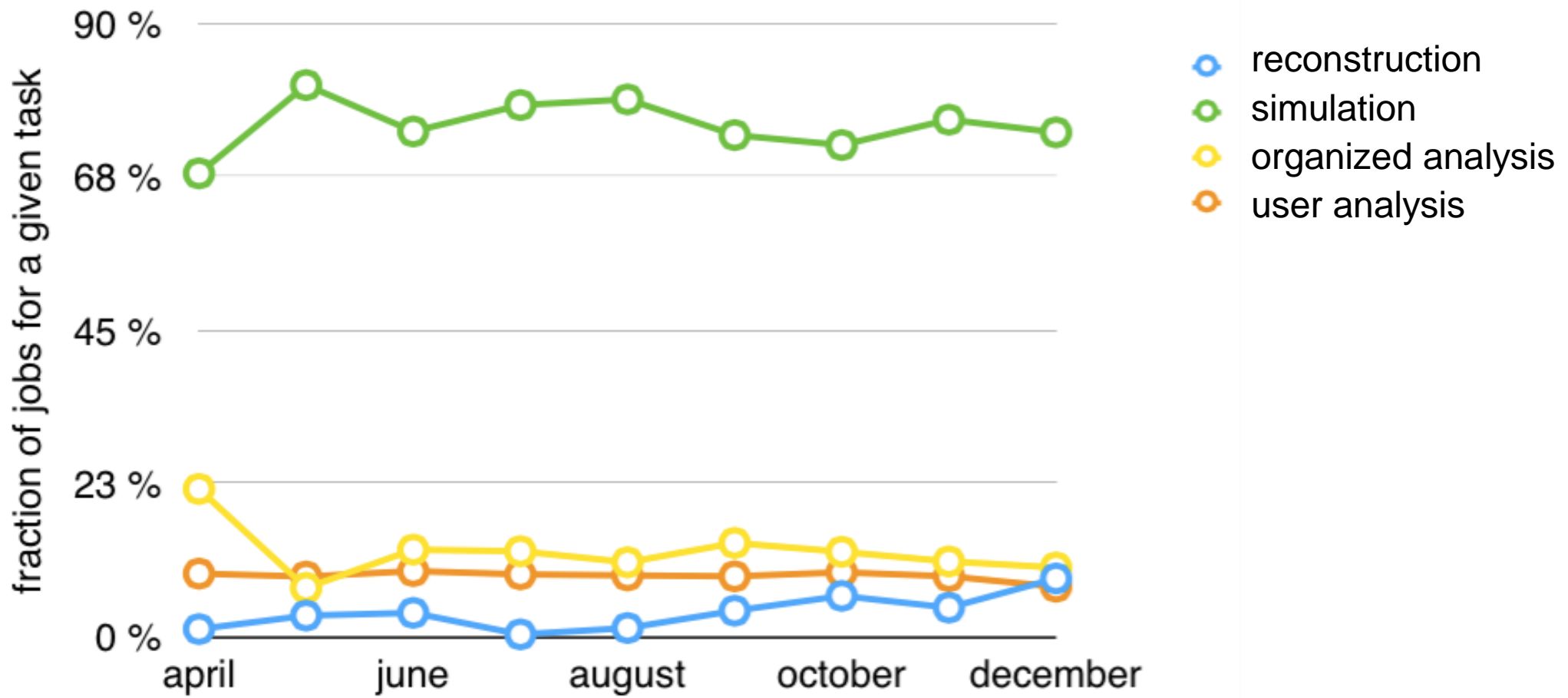
- Test of upgraded detectors readout, Trigger, DAQ, new HLT farm
- Full data recording chain, with conditions data gathering
- Cosmics trigger data taking with Offline processing





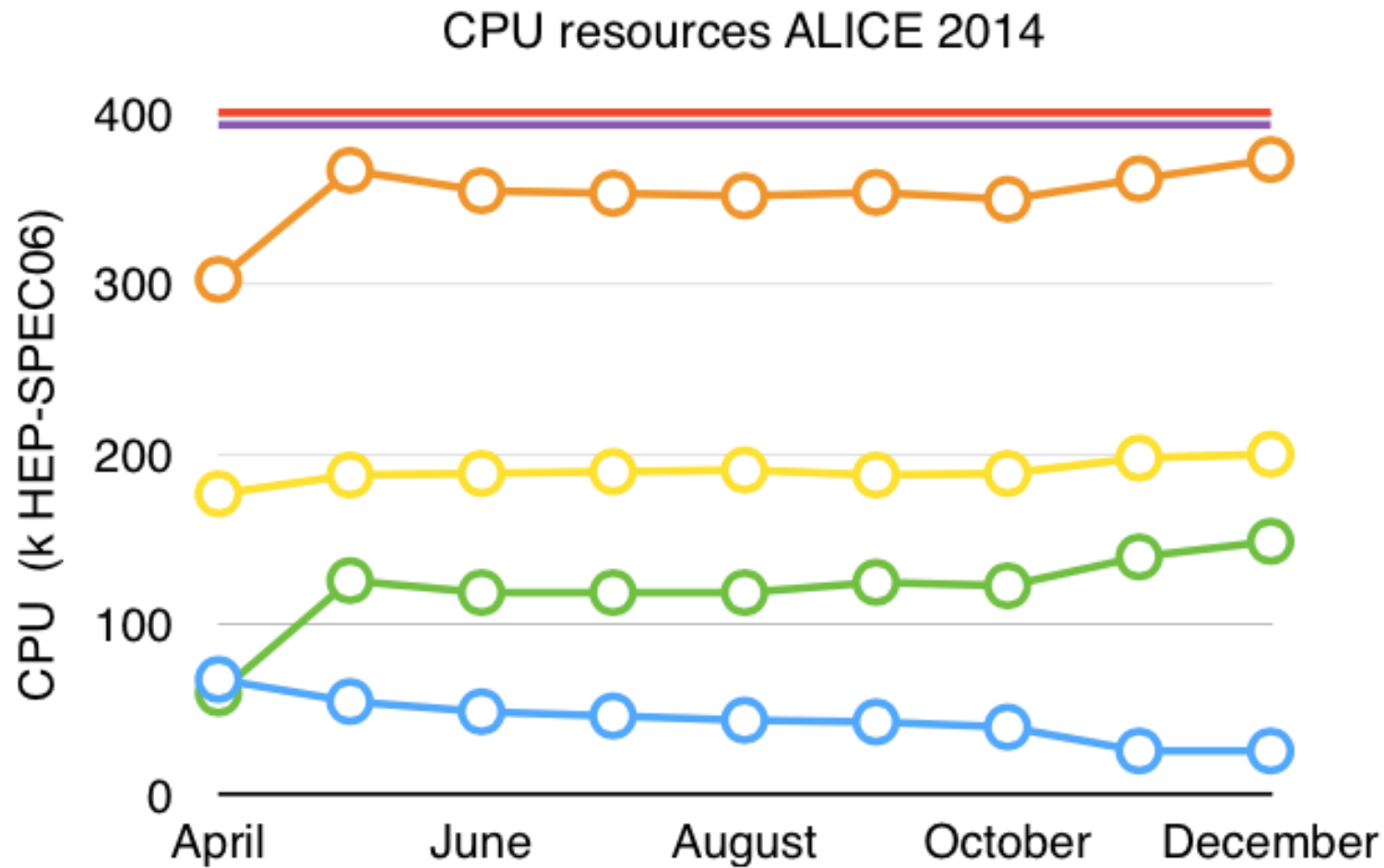
CPU Shares

CPU Usage





CPU Resources in 2014

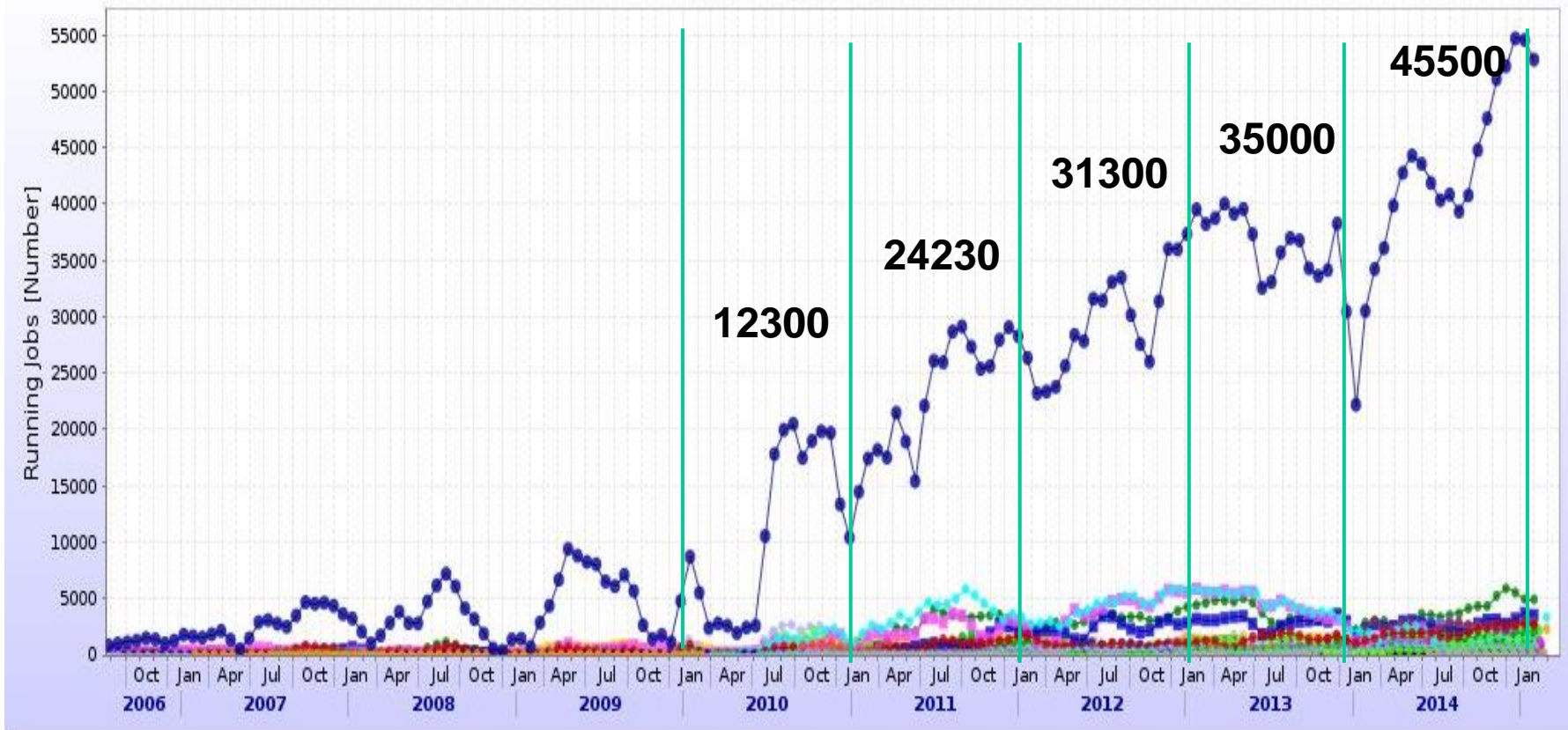


- T0
- T1
- T2
- SUM
- Required
- Pledged





Available CPU evolution



Year on year increase

+97
%

+30
%

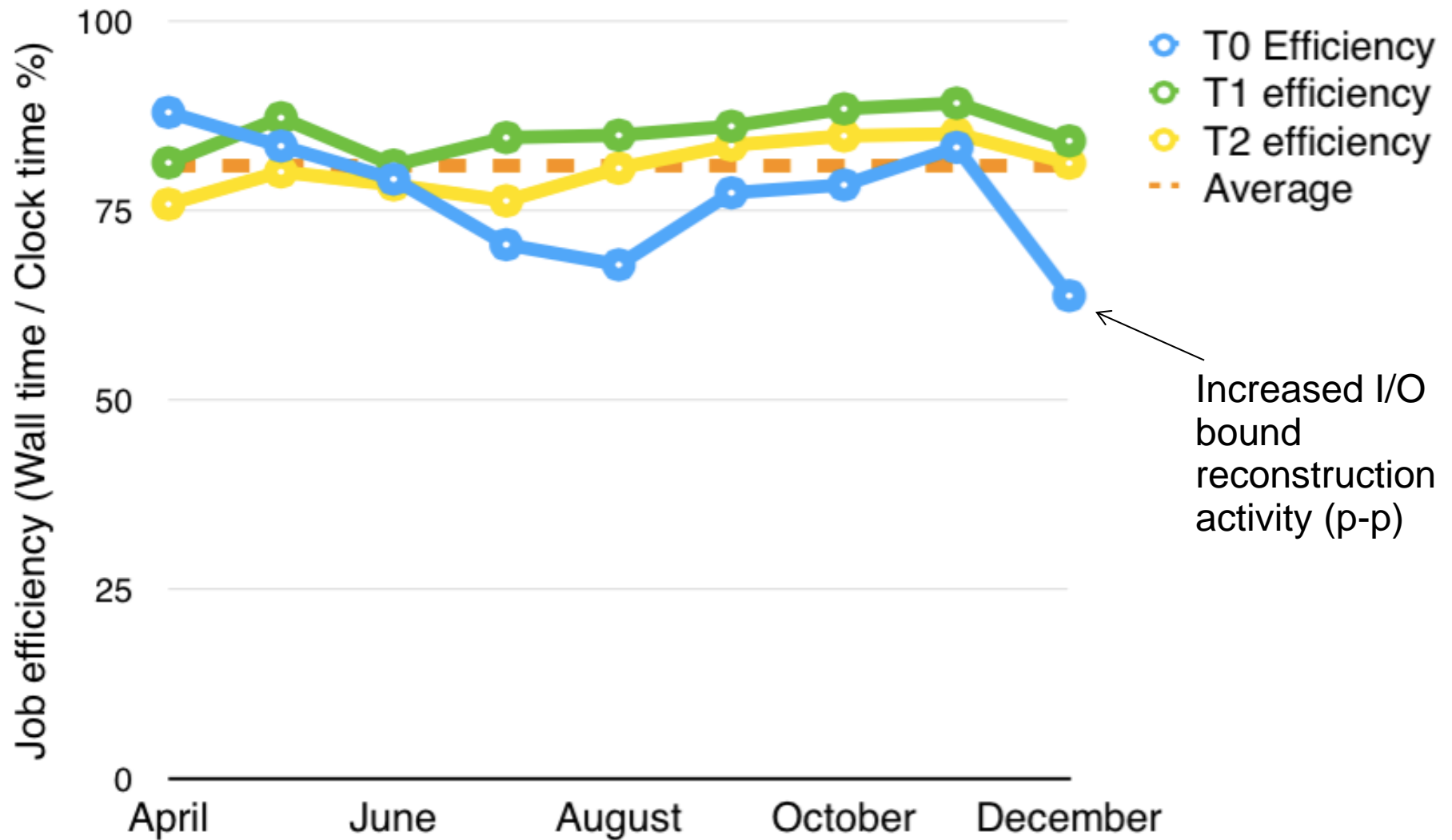
+12
%

+30
%

- 22% average per year – slightly above the WLCG projection

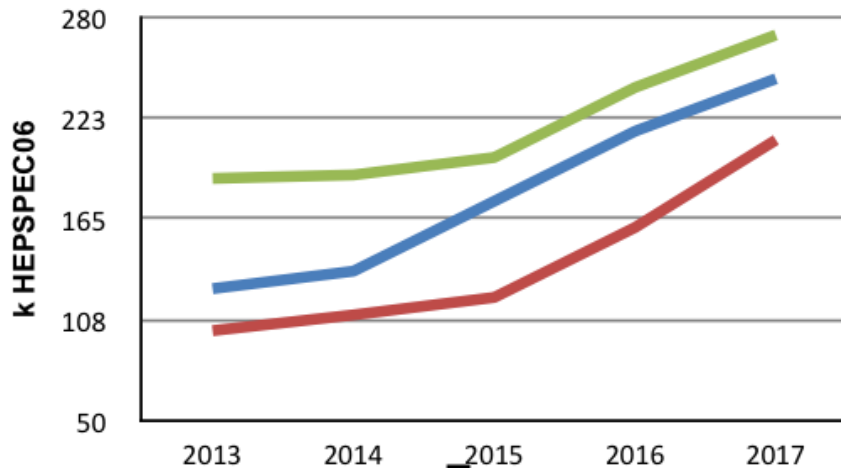


CPU Shares

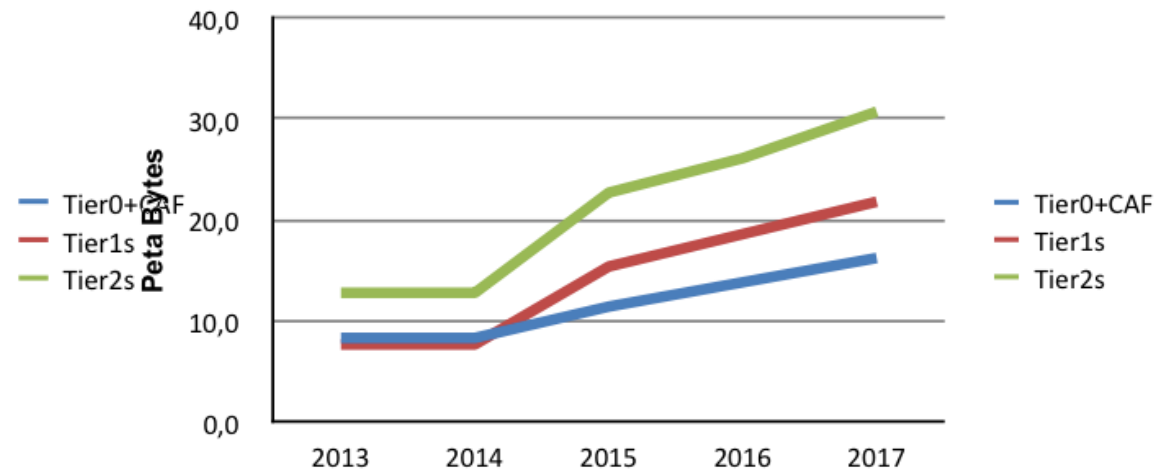


Increased I/O bound reconstruction activity (p-p)

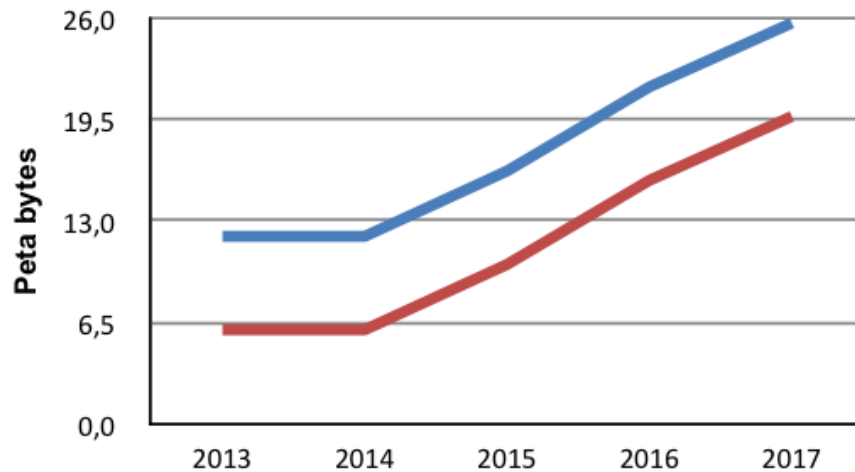
CPU



Disk



Tape



No significant changes compared to what we showed last time



ALICE

O2 TDR

ALICE UPGRADE

ALICE A Large Ion Collider Experiment | February 2015

ALICE
Technical Design Report

CERN-LHCC-2015-xxxx
ALICE-TDR-xxxx
February 18, 2015

Upgrade of the Online - Offline computing system

Technical Design Report

Technical Design Report for the Upgrade of the Online - Offline computing system | CERN-LHCC-2015-xxxx (ALICE-TDR-xxxx)

- Reduce data volume early and by a large factor
 - TPC x20, overall x14
 - Requires considerable computing capacity @ P2
- Use h/w accelerators to speed up the computation and reduce the cost
 - Requires new and flexible s/w framework that can adapt itself to different environments
- Assume that Grid capacity will continue to grow within constraints of a flat budget
 - Available resources must be used with maximum efficiency
 - 2/3 of raw data processing will happen @ P2
 - All intermediate data formats will be transient or temporary
 - Only the compressed raw data and analysis level data will be persistent



The ALICE O2 Project: Architecture



Detector

Synchronous readout, aggregation compression

- FLP 2.5x
- EPN 8x

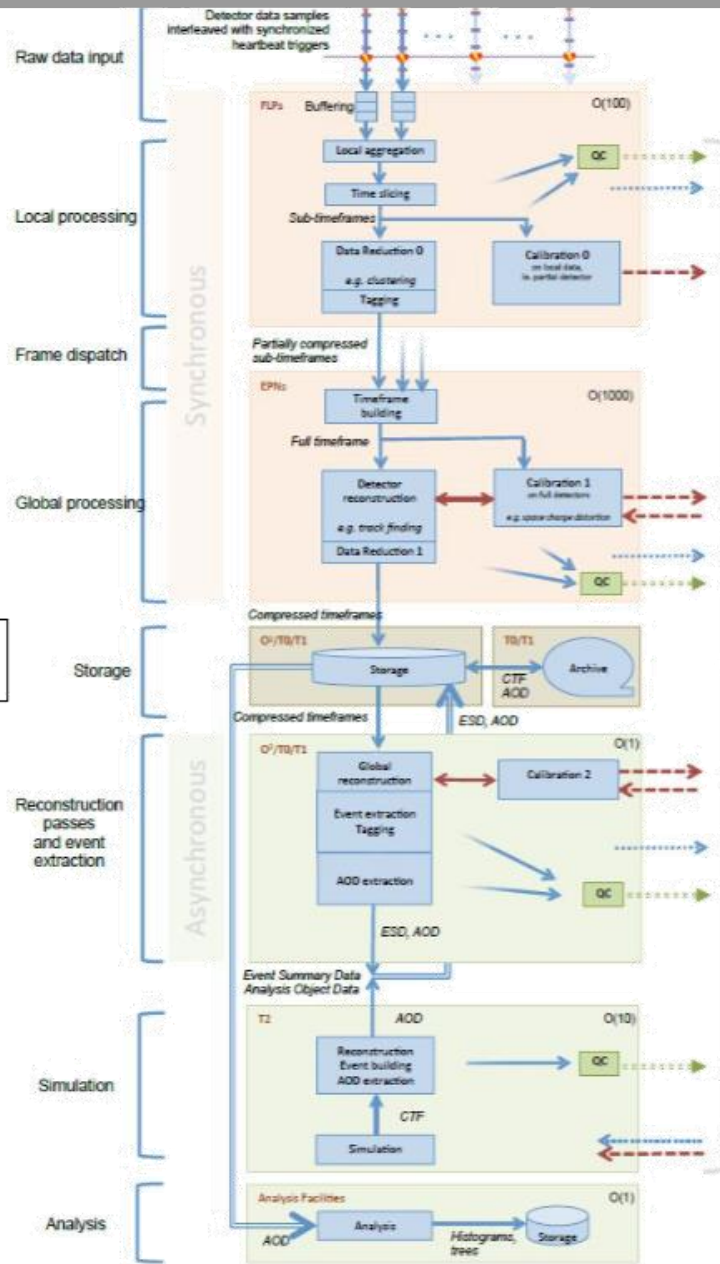
Asynchronous reconstruction

75 % EPN / 25 % T0-T1

Storage

Simulation

Analysis

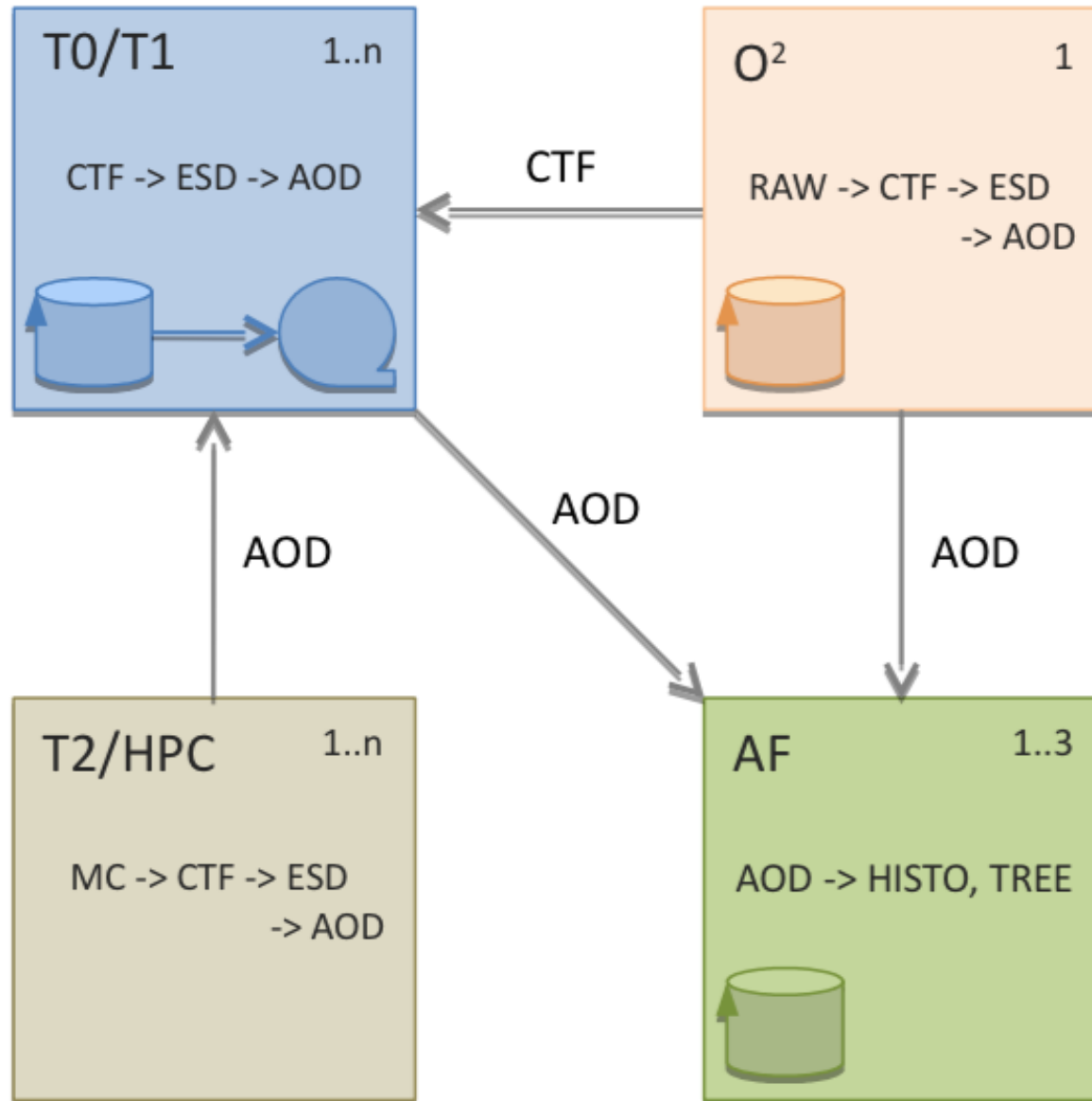


Size of O2 facility:

- ~100'000 CPU cores
- 5000 GPUs
- 50 PB of disk



Roles of Tiers



- **Motivation**
 - Analysis is the least efficient of all workloads that we run on the Grid
 - I/O bound in spite of attempts to make it more efficient by using the analysis trains
 - Increased data volume will only magnify the problem
- **Solution**
 - Collect AODs on a few dedicated sites that are capable of locally processing quickly large data volume
 - Typically (a fraction of) HPC facility (20-30'000 cores) and 5-10 PB of disk on very performant file system
 - Run organized analysis on local data like we do today on the Grid

- ALICE expects grid resources to evolve and grow at 20% per year rate which is consistent with a flat funding
- We expect 20Gb/s share of network connectivity between CERN and T1s in order to be able to export 1/3 of raw data to T1s
- On T1s data will need to be archived on tape and subsequently processed (calibration & reconstruction)
- Since T2s will be used almost exclusively for simulation jobs (no input) and resulting AODs will be exported to T1s/AFs, we expect to significantly lower the future needs for storage on T2s and would like to use available funding to buy more CPUs
- While in this model the sites will be mostly specialized for a given task, we still want to retain ability to run any kind of job on any available resource
- Data management is going to be the biggest issue, we need a uniform solution

- **ALFA & AliceO2**

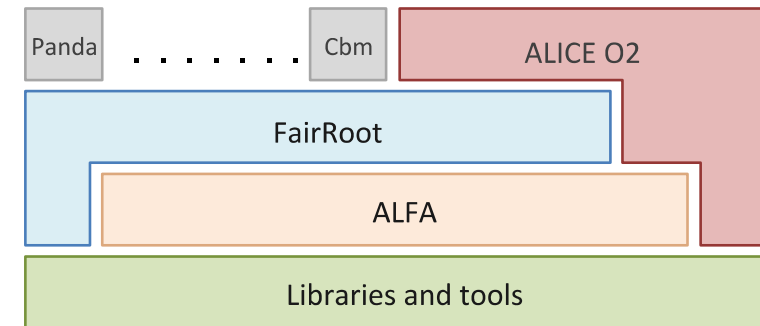
- Common software layer between FairRoot and AliceO2
- Git repository, Cdash/Ctest infrastructure in place
- Template for implementation of a detector code in a new frameworks
- DDS – Dynamic Deployment System v0.8 released
- Interface to existing HLT modules in place allowing to test existing algorithms (TPC reconstruction on GPUs)

- **AliEn & PanDA**

- Collaboration between ATLAS and ALICE to allow Grid like interface to HPC resources
- Architecture agreed and manpower finally identified

- **ALICE & OpenLab**

- Collaboration on productizing EOS storage (COMTRADE)
- Rack-scale computing (Intel, Cisco..)



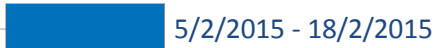


O2 TDR schedule

2015



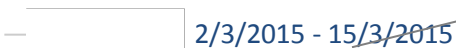
Comments TDR by the O² members



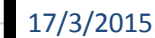
TDR editing



Comments TDR by the ALICE Coll.



ALICE internal review



Modification by the authors



Final editing of the TDR and UCG



- Borut Paul Kersevan, IJS, ATLAS
- Tonko Lubicic, BNL, STAR
- Niko Neufeld, CERN, LHCb

LHCC review



LHC Committee





Conclusions

- Resources for Run 2 are sufficient and we are showing that we can use them efficiently
- The primary goal of the O2 facility is data compression
- In a new computing model we try to minimize the amount of data moving between Tiers and carry out most of the processing on local datasets
- We expect the Grid to grow by 20% per year
- We expect 1/3 of raw data processing to be done on T1s
- Since T2s will be used mostly for simulation and we can rebalance the CPU/disk ratio and buy more CPUs
- Dedicated Analysis Facilities are a new element that should improve analysis efficiency
- TDR is on schedule



ALICE

Conclusions

Backup slides



Running scenario

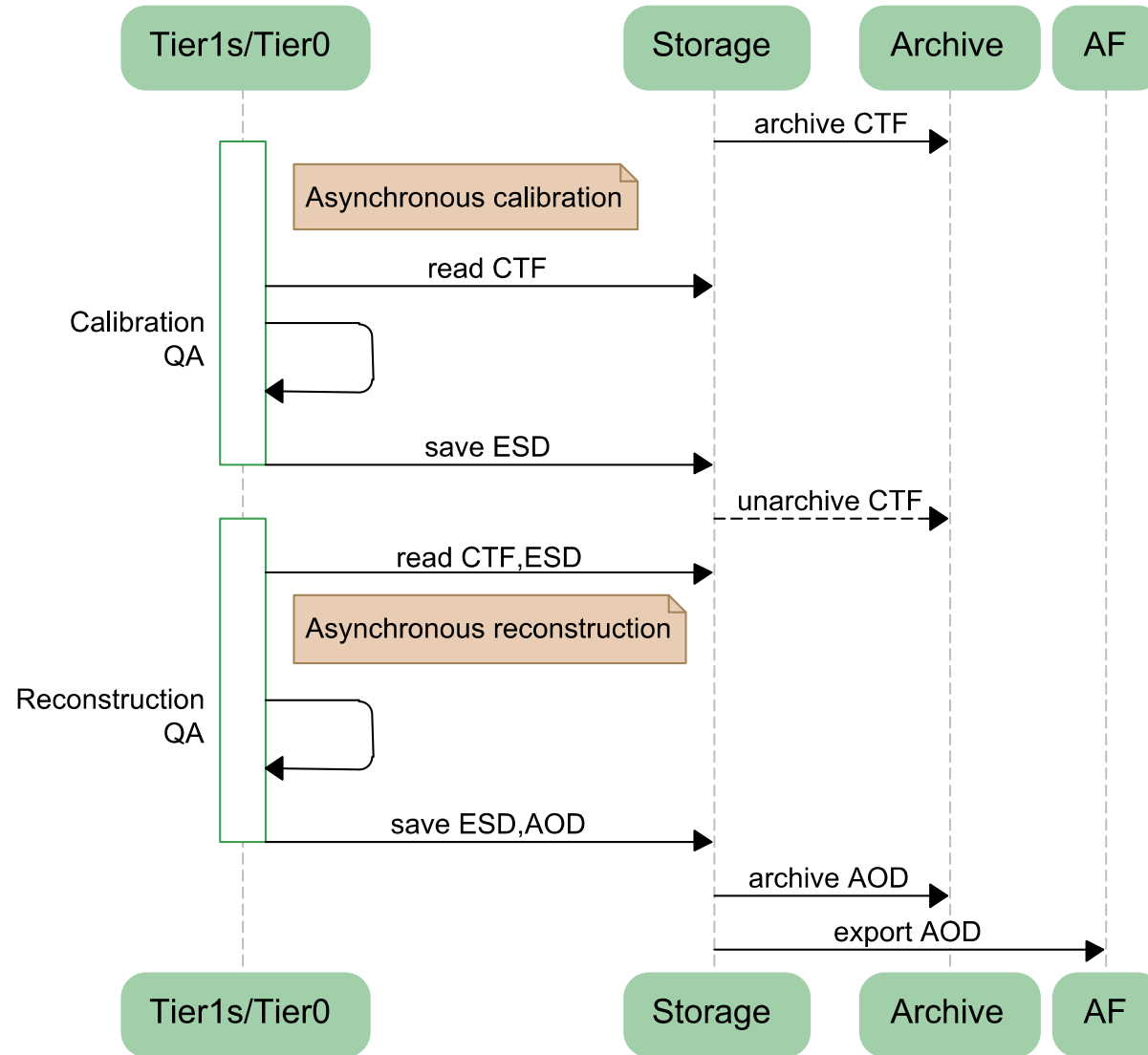
Year	System	$\sqrt{s_{NN}}$	L_{int}	$N_{collisions}$
2020	pp	14 TeV	6 pb ⁻¹	4 · 10 ¹¹
	Pb–Pb	5.5 TeV	2.85 nb ⁻¹	2.3 · 10 ¹⁰
2021	pp	14 TeV	4 pb ⁻¹	2.7 · 10 ¹¹
	Pb–Pb	5.5 TeV	2.85 nb ⁻¹	2.3 · 10 ¹⁰
2022	pp	14 TeV	4 pb ⁻¹	2.7 · 10 ¹¹
	pp	5.5 TeV	6 pb ⁻¹	4 · 10 ¹¹
2025	pp	14 TeV	4 pb ⁻¹	2.7 · 10 ¹¹
	Pb–Pb	5.5 TeV	2.85 nb ⁻¹	2.3 · 10 ¹⁰
2026	pp	14 TeV	4 pb ⁻¹	2.7 · 10 ¹¹
	Pb–Pb	5.5 TeV	1.4 nb ⁻¹	1.1 · 10 ¹⁰
	p–Pb	8.8 TeV	50 nb ⁻¹	10 ¹¹
2027	pp	14 TeV	4 pb ⁻¹	2.7 · 10 ¹¹
	Pb–Pb	5.5 TeV	2.85 nb ⁻¹	2.3 · 10 ¹⁰

Data types

Acronym	Description	Persistency
RAW	Raw data as it comes from the detector	Transient
CTF	Compressed Time Frame containing processed raw data of for a period of time ≈ 100 ms. In the case of TPC clusters not belonging to tracks are rejected and the remaining information is compressed to the maximum. Once written, CTF becomes read only data.	Persistent
ESD	Event Summary Data. Auxiliary data to CTF containing the output of the reconstruction process that assigns tracks to vertices and identifies the individual collisions.	Temporary
MC	Simulated energy deposits in sensitive detectors. Removed once the reconstruction of MC data is completed.	Transient
AOD	Analysis Object Data containing the final track parameters in a given vertex and for a given physics event. AODs are collected on dedicated facilities for subsequent analysis.	Persistent
MCAOD	Analysis Object Data for a given simulated physics event. Same as AOD with addition of kinematic information that allows comparison to MC. MCAODs are collected on dedicated facilities for subsequent analysis.	Persistent
HISTO	The subset of AOD information specific for a given analysis. Can be generated during analysis but needs to be offloaded from the Grid.	Temporary



Tier 0/1





Tier 2

