# The ALICE Grid upgrade, methods and tools for LHC Run 3 and beyond

L. Betev
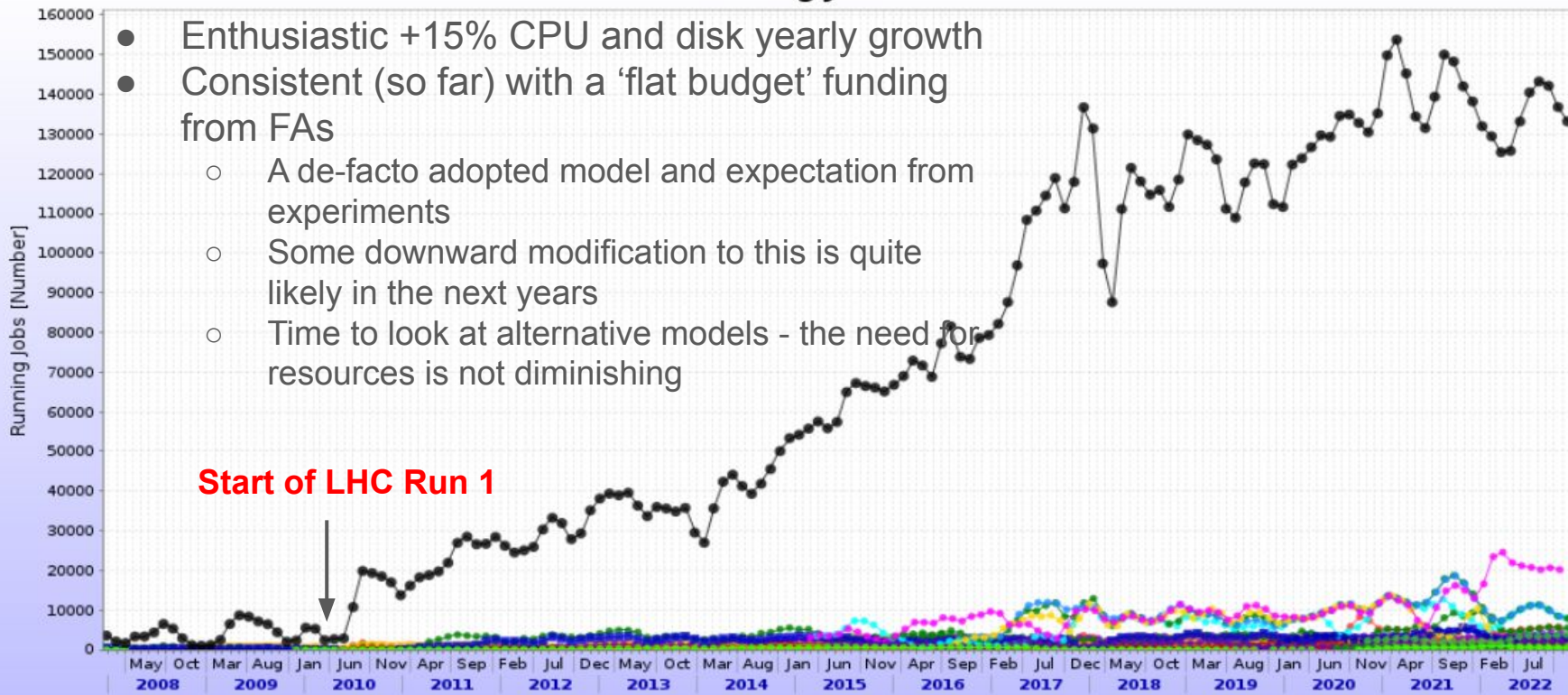
Journées LCG-France, LPNHE, 6-8 June 2023

# The ALICE Grid - individual computing centres



North America - 3 T2s

Europe - 1 T0, 7T1s, 42 T2s

Asia - 7 T2s, 1 T1

South America - 2 T2s

Africa - 1 T2

# ALICE resources evolution

**Running Jobs**

- Enthusiastic +15% CPU and disk yearly growth
- Consistent (so far) with a 'flat budget' funding from FAs
  - A de-facto adopted model and expectation from experiments
  - Some downward modification to this is quite likely in the next years
  - Time to look at alternative models - the need for resources is not diminishing

**Start of LHC Run 1**

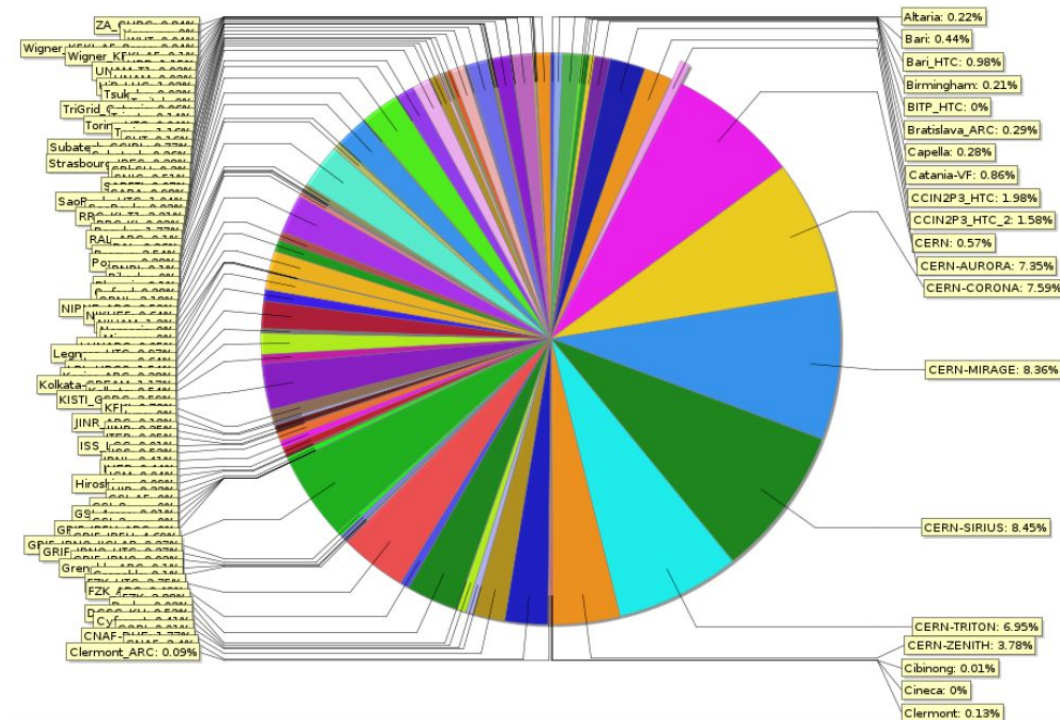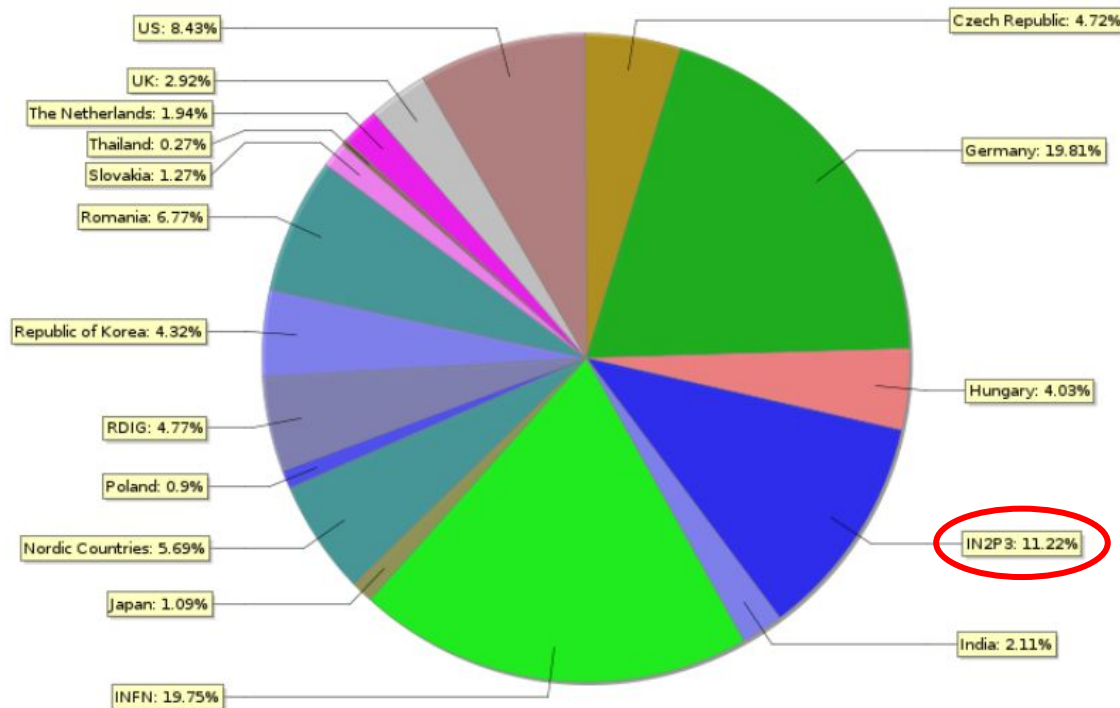# ALICE resources use per activity

# Role of Tiers

- T0 - RAW reco + MC + analysis
- T1s - RAW reco + MC + analysis
- T2s - MC + Analysis
- Differences between tiers - custodial storage + nominal services response time
- In practice - all tiers run effectively all types of workload (except RAW reco) and availability is ~same
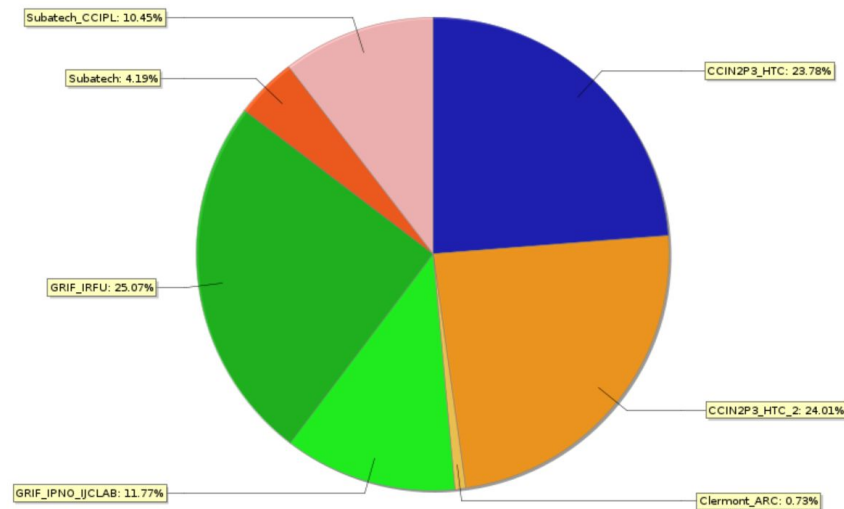- ALICE model can absorb any site size



5

# Regional contribution to ALICE computing

- ~12% FR contribution
- T1@CCIN2P3
- 4 (soon to be 3) T2s
  - Clermont
  - GRIF_IRFU
  - GRIF_IJCLAB
  - Subatech (+CCIPL)
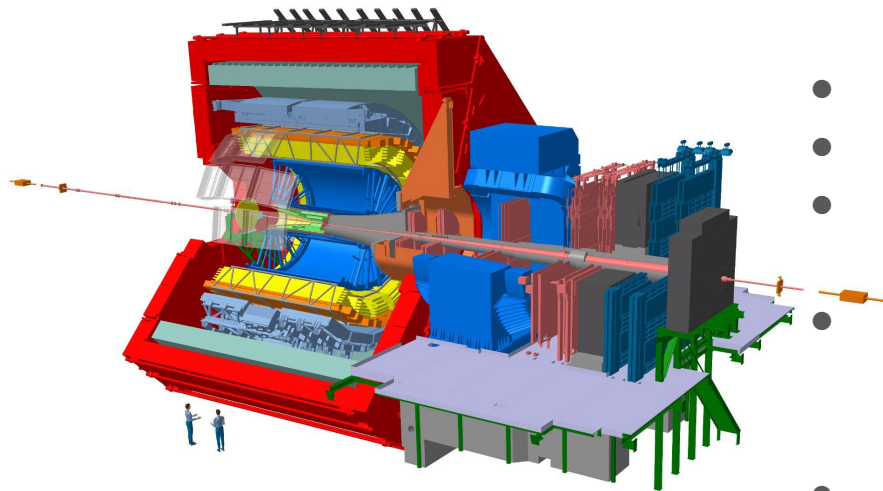- Diminishing role of T2 centres - this is an unfortunate global trend

# Repartition of resources in FR sites

- ~50% at CCIN2P3
- ~15% Subatech
- ~35% GRIF (IRFU+IJCLAB)
- ~1% Clermont
- The imminent loss of Subatech is a substantial hit to ALICE computing in France
  - In addition it is the loss of one of the oldest Grid centres in the country(!)



7

# ALICE upgrade general

- p-p and HI physics
- 10x integrated luminosity $L \sim 10nb^{-1}$ (B=0.5T) + $3nb^{-1}$ (B=0.2T)
- 100x event rate of Run 1/2, 10x more data
- Continuous readout
- Focus on data compression and real time (synchronous) data reconstruction
- => Reasonable rates and data volumes after compression to storage and secondary data formats
- Adherence to 'flat budget' resources funding for data processing and analysis

8

# Data flow and data rates

# The O2 facility (EPNs)



- Container-hosted computing facility located at the ALICE site, PUE<1.07
- High-throughput system, heterogeneous computing platform (CPU+GPU)
- 250 dual CPU nodes (ROME, 64 cores, 512GB RAM) with 8 AMD (MI50, 32GB) GPUs/node
- Functions
  - Data aggregation (Detector STFs to global CTF)
  - Synchronous global reconstruction
  - Calibration and data volume reduction
  - Quality control
  - Asynchronous (offline) reconstruction
- Containers house a backup EOS storage in case of network interruption to CC

# Synchronous data processing

- Goal - to compress the RAW data by about factor 35 (3.5TB/s -> 100GB/s)
- Through zero suppression, clusterization, tracking, optimized data format
  - ***Mandatory use of GPUs (40x faster than CPUs)***
  - All synchronous level software is written for GPUs for all detectors



*Pb-Pb @50 kHz IR*
*2 ms drift time*
*TPC reconstructed tracks from*
*different colour-coded events*

Unassigned clusters (noise)   Reconstructed tracks
Removed clusters   Failed fits

# O2 Software framework



- Developed by ALICE in collaboration with the FAIR group at GSI Darmstadt
- Three major parts
  - Transport layer, based on FairMQ - message passing toolkit
  - Data model - ALICE-specific object description and content
  - Data processing layer - set of data processors implicitly organized in a logical dataflow for data transformation
- Trivially parallel and integrates tools for GPU offloading
- Natural use of multicore processing and shared memory - *move to multicore*

*Users define tasks and their dependencies*

*O2 DPL builds a workflow (DAG) out of the specified tasks*

*O2 DPL builds a topology of FairMQ devices and maps tasks on it into account available resources*

12

# Analysis facilities (AFs)

AOD
MCAOD

Single access to
Input data

Analysis
train

Task 1

Task 2

Task 3

Task 4

Task N

- New element of the computing model
- Data transferred to AF from T0/T1s/T2s
- Goals
  - Provide a location with comprehensive data samples from asynchronous and MC data processing at ~10% statistics
  - Fast tuning of analysis algorithms - once ready, run on full sample on the Grid
  - First data and low statistics analysis (if compatible)
- Incorporated in the Grid framework
- Sites tuned for fast I/O between storage and CPU
  - Approximate total size 6-8k cores, 10PB storage
  - ~15MB/s/core throughput
- As of today - GSI Darmstadt and KFKI Budapest (⅔ of the AF target, looking for more suitable sites)

13

# Grid middleware development - JAliEn

File catalogue

Task queue

JCentral — Central services

JSite
Multiplexer, object caching — VO-box

JobAgent

Payload
O2

JBox
Transition between Java serialized objects and payload — Worker nodes

- Evolution of the AliEn middleware
  - Refactored and rewritten in Java
- Highly efficient and scalable communications infrastructure
- Persistent, compressed, SSL channels
- Multiplexing and object caching
- Use of Java serialized objects
- Platform independent
- Multi-core enabled, HPC ready
- Deployed gradually on the existing infrastructure - no interference with operations

# Site services evolution

JSite
Multiplexer, object caching
CE interface (either gateway -
HTCondor/ARC or local batch)

MonALISA
Cache for local monitoring
information
Object filtering and
communication with central ML
instance

VO-box

- New middleware for sites - simplification of operation
- JAliEn was installed gradually and in combination with the local CE updates
- From 5 services to 2
  - The remaining services are quite reliable and effectively do not require site manager intervention
- Automatic updates to new version of VO-box services
- Monitoring of all relevant info for the site is provided on a single page here

# JobRunner, JobAgent and JobWrapper



- Entirely new method for both resources and job control
- Fully containerized workload
- Ability to run multiple jobs within the control of the same JobRunner
  - Effective control of any set of resources provided

# Payload containers

- By default, all jobs are wrapped in a CentOS 7.9 container
- Other images are available
  - Rocky 8.6: For newer payloads and **GPUs** (special for the EPN cluster)
  - Rocky 9.0 + RHEL 9: Already certified
  - Debug containers, for example with vtune, strace
- GPUs are supported in Apptainer (formerly Singularity)
- All of the above allows for fulfilling various job requirements, independent of the underlying OS
- Allows for use of HPCs or other specialized clusters (for example EPNs)

# Job isolation and control - applying `taskset`



- Total CPU usage goes above the requested 8 cores

- CPU consumption is limited with `taskset`
- Total CPU usage is flat at 8 cores

- Applicable for sites with non-constrained resources and full node submission

# Improving job efficiency through CPU pinning



*Sample host CPU architecture*

- Various core/cache pinning configurations possible
  - Same NUMA Node and independent L1,L2 cache
  - Different NUMA Nodes and independent L1,L2 cache
  - Same NUMA Node and sharing L1,L2 cache
  - Random core assignment
  - No pinning

19

# Improving job efficiency through CPU pinning

- Most efficient configuration - same NUMA node, independent L1/L2 cache - compared to no pinning
- Only possible if full control of the CPU - whole node
- Already in production at LBNL - Lawrencium HPC

# New tools and monitoring - SiteSonar

- Tool to evaluate site capabilities and installations - probes invoked at the beginning of execution
  - Collects data from ~10K Grid nodes daily

# Storage performance - FileCrawler

- Checks storage integrity on sites by mimicking normal jobs
- Random files, proportional to the storage size
- Reporting on file health, throughput and accessibility
- Early detection of storage issues

**Status codes extracted from the crawler**

SE Name: ALICE::HIROSHIMA::EOS    Interval: Last week

| Status Type | Status Code | Status Count | Status Code Ratio | Download throughput |
|---|---|---|---|---|
| FILE_OK | S_FILE_CHECKSUM_MATCH | 26972 | 99.79 % | 21.97 Mb/s |
| | E_CATALOGUE_MD5_IS_BLANK | 2 | 0.01 % | 19.04 Mb/s |
| INTERNAL_ERROR | XRDFS_CANNOT_CONFIRM_UPLOAD | 21 | 0.08 % | |
| FILE_INACCESSIBLE | XROOTD_EXITED_WITH_CODE | 35 | 0.13 % | |
| TOTAL | | 27030 | 100 % | |

**Averaged metrics for the selected interval**

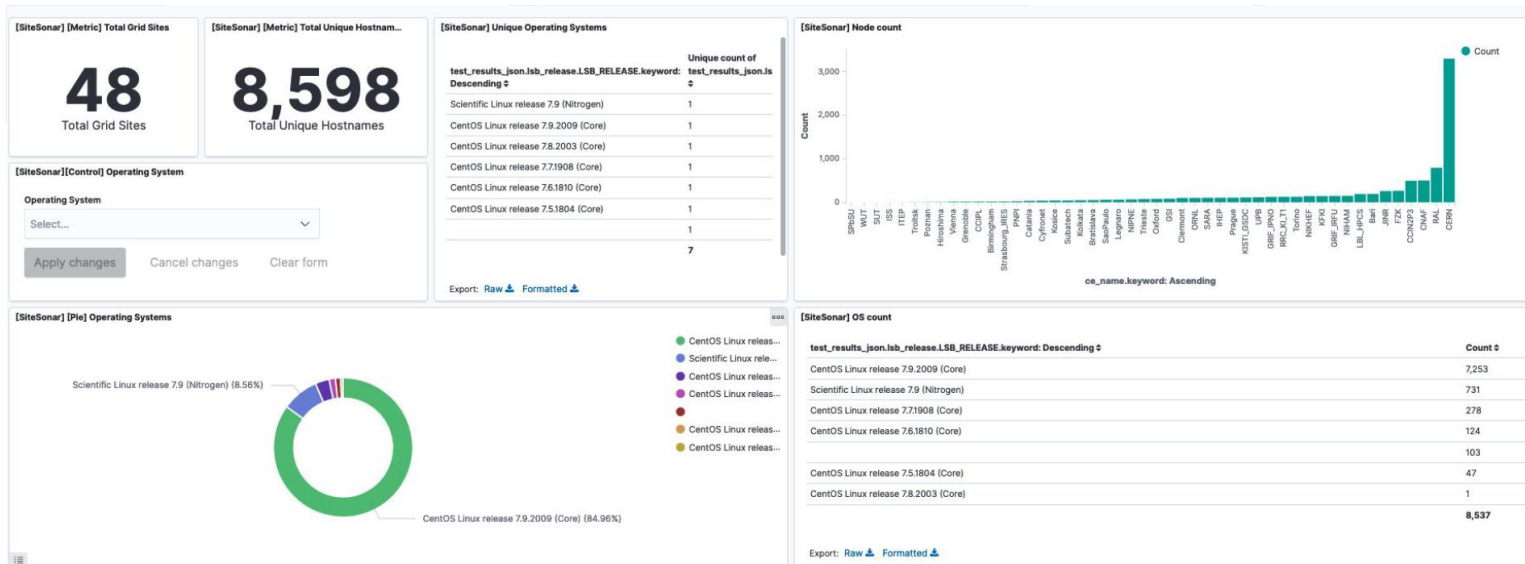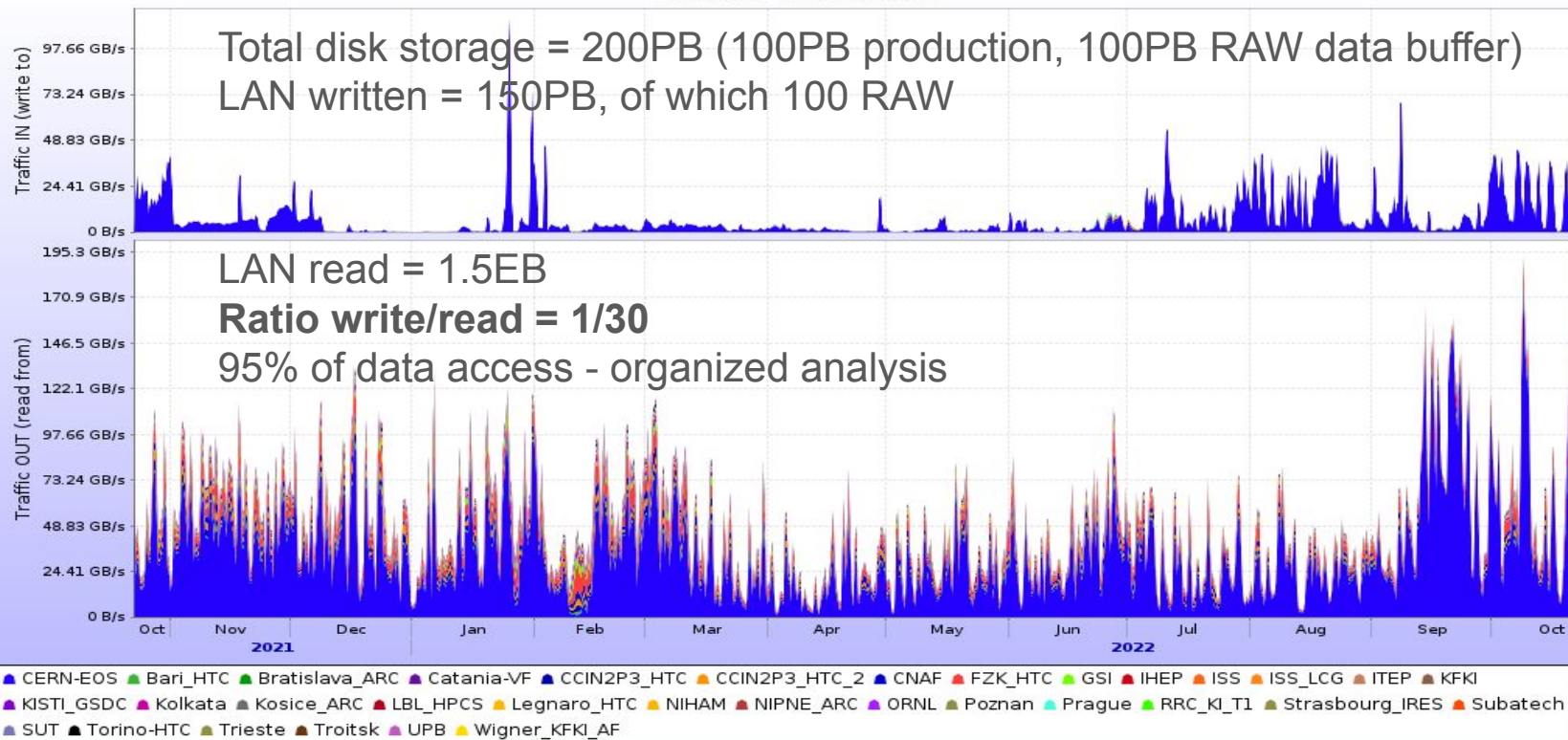| SE Name | Start | End | Success ratio ▲ | Corrupt ratio | Inaccessible ratio | Internal error ratio |
|---|---|---|---|---|---|---|
| SARA::DCACHE | 18 Oct 2022 06:08 | 17 Nov 2022 10:32 | 99.87 % | 0.09 % | 0.05 % | 0.00 % |
| Hiroshima::EOS | 18 Oct 2022 06:08 | 17 Nov 2022 10:33 | 99.73 % | 0.00 % | 0.18 % | 0.09 % |
| SNIC::DCACHE | 18 Oct 2022 06:12 | 17 Nov 2022 10:28 | 99.68 % | 0.02 % | 0.27 % | 0.03 % |
| Vienna::EOS | 18 Oct 2022 06:07 | 17 Nov 2022 10:38 | 99.60 % | 0.24 % | 0.16 % | 0.00 % |
| NIPNE::EOS | 18 Oct 2022 06:09 | 17 Nov 2022 13:03 | 99.58 % | 0.03 % | 0.37 % | 0.03 % |
| Trieste::SE | 18 Oct 2022 06:11 | 17 Nov 2022 12:11 | 99.54 % | 0.11 % | 0.35 % | 0.00 % |
| Bari::SE | 18 Oct 2022 06:04 | 17 Nov 2022 12:22 | 99.50 % | 0.08 % | 0.42 % | 0.00 % |
| IHEP::SE | 18 Oct 2022 06:07 | 17 Nov 2022 10:20 | 99.35 % | 0.11 % | 0.53 % | 0.01 % |
| Torino::SE2 | 18 Oct 2022 06:09 | 17 Nov 2022 11:07 | 99.34 % | 0.13 % | 0.53 % | 0.00 % |
| Troitsk::SE | 18 Oct 2022 06:04 | 17 Nov 2022 10:43 | 99.26 % | 0.54 % | 0.19 % | 0.01 % |
| CERN::EOS | 18 Oct 2022 06:12 | 17 Nov 2022 10:47 | 99.19 % | 0.08 % | 0.65 % | 0.07 % |
| CNAF::SE | 18 Oct 2022 06:10 | 17 Nov 2022 10:35 | 99.06 % | 0.02 % | 0.92 % | 0.00 % |
| FZK::SE | 18 Oct 2022 06:11 | 17 Nov 2022 10:33 | 98.86 % | 0.06 % | 1.07 % | 0.01 % |
| Legnaro::SE | 18 Oct 2022 06:04 | 17 Nov 2022 10:26 | 98.54 % | 0.03 % | 1.34 % | 0.09 % |
| UPB::EOS | 18 Oct 2022 06:08 | 17 Nov 2022 10:32 | 98.49 % | 0.07 % | 1.44 % | 0.00 % |
| ORNL::EOS | 18 Oct 2022 06:06 | 17 Nov 2022 10:31 | 98.18 % | 0.46 % | 1.36 % | 0.00 % |
| NDGF::DCACHE | 18 Oct 2022 06:04 | 17 Nov 2022 10:30 | 97.89 % | 0.23 % | 1.87 % | 0.00 % |
| NIHAM::EOS | 18 Oct 2022 06:08 | 17 Nov 2022 10:49 | 97.75 % | 0.12 % | 2.13 % | 0.00 % |
| GRIF::EOS | 18 Oct 2022 06:05 | 17 Nov 2022 10:31 | 97.75 % | 0.05 % | 2.20 % | 0.00 % |
| Subatech::EOS | 17 Oct 2022 17:38 | 16 Nov 2022 16:28 | 97.46 % | 0.06 % | 0.91 % | 1.57 % |
| JINR::EOS | 18 Oct 2022 06:11 | 17 Nov 2022 12:13 | 95.93 % | 0.13 % | 3.92 % | 0.03 % |
| RRC_KI_T1::EOS | 18 Oct 2022 06:06 | 17 Nov 2022 10:30 | 95.86 % | 0.09 % | 1.47 % | 2.57 % |
| KISTI_GSDC::EOS | 18 Oct 2022 06:07 | 17 Nov 2022 10:57 | 95.04 % | 3.49 % | 1.47 % | 0.01 % |
| CCIN2P3::SE | 18 Oct 2022 06:11 | 17 Nov 2022 10:37 | 94.27 % | 0.02 % | 5.69 % | 0.02 % |
| Kosice::EOS | 18 Oct 2022 06:07 | 17 Nov 2022 11:40 | 93.05 % | 0.11 % | 6.84 % | 0.00 % |
| Prague::SE | 18 Oct 2022 06:06 | 17 Nov 2022 10:44 | 90.18 % | 0.02 % | 9.79 % | 0.01 % |
| Birmingham::EOS | 18 Oct 2022 06:05 | 17 Nov 2022 10:26 | 87.70 % | 0.06 % | 12.23 % | 0.01 % |
| Strasbourg_IRES::SE2 | 18 Oct 2022 06:04 | 17 Nov 2022 12:46 | 87.68 % | 0.03 % | 12.26 % | 0.03 % |
| Catania::SE | 18 Oct 2022 06:07 | 17 Nov 2022 10:23 | 86.12 % | 0.03 % | 13.84 % | 0.00 % |
| KISTI_GSDC::SE2 | 18 Oct 2022 06:07 | 17 Nov 2022 10:41 | 86.03 % | 0.17 % | 13.80 % | 0.00 % |
| LBL_HPCS::EOS | 18 Oct 2022 06:04 | 17 Nov 2022 10:23 | 85.88 % | 1.21 % | 12.90 % | 0.00 % |
| Poznan::SE | 17 Oct 2022 23:50 | 17 Nov 2022 10:42 | 79.63 % | 0.33 % | 20.04 % | 0.00 % |
| ISS::FILE | 18 Oct 2022 06:07 | 17 Nov 2022 05:19 | 78.76 % | 0.07 % | 21.12 % | 0.04 % |
| Kolkata::EOS2 | 18 Oct 2022 06:09 | 17 Nov 2022 12:09 | 68.71 % | 0.61 % | 30.57 % | 0.12 % |

22

# Network and data processing

- Jobs are dispatched to the Grid sites that already have the data
  - Minimizes WAN traffic and RTT efficiency penalty
- Grid site local file access (95%), remote (5%)
  - Remote access due to local SE issues, usually temporary
- Multiple replicas sorted topologically: apps first access local replica, then the next closest
  - Sorting by network topology, availability, network quality, geo-location and other metrics
- Storing multiple replicas
  - One replica is written to the local storage element
  - The other replicas are written to the remote (but close) storage elements
  - Remote writes might go through LHCOPN / LHCONE

# Data access - LAN

**LAN server traffic**

Total disk storage = 200PB (100PB production, 100PB RAW data buffer)
LAN written = 150PB, of which 100 RAW

LAN read = 1.5EB
**Ratio write/read = 1/30**
95% of data access - organized analysis



● CERN-EOS ● Bari_HTC ● Bratislava_ARC ● Catania-VF ● CCIN2P3_HTC ● CCIN2P3_HTC_2 ● CNAF ● FZK_HTC ● GSI ● IHEP ● ISS ● ISS_LCG ● ITEP ● KFKI
● KISTI_GSDC ● Kolkata ● Kosice_ARC ● LBL_HPCS ● Legnaro_HTC ● NIHAM ● NIPNE_ARC ● ORNL ● Poznan ● Prague ● RRC_KI_T1 ● Strasbourg_IRES ● Subatech
● SUT ● Torino-HTC ● Trieste ● Troitsk ● UPB ● Wigner_KFKI_AF

24

# Data access - WAN (LHCONE/LHCOPN)

**WAN server traffic**

WAN written = 6PB

WAN read = 60PB, 4% of LAN
Remote write/read due to local SE issues, usually temporary

Traffic IN (write to): 1.953 GB/s, 1.465 GB/s, 1000 MB/s, 500 MB/s, 0 B/s

Traffic OUT (read from): 9.766 GB/s, 9.277 GB/s, 8.789 GB/s, 8.301 GB/s, 7.813 GB/s, 7.324 GB/s, 6.836 GB/s, 6.348 GB/s, 5.859 GB/s, 5.371 GB/s, 4.883 GB/s, 4.395 GB/s, 3.906 GB/s, 3.418 GB/s, 2.93 GB/s, 2.441 GB/s, 1.953 GB/s, 1.465 GB/s, 1000 MB/s, 500 MB/s, 0 B/s

Oct Nov Dec Jan Feb Mar Apr May Jun Jul Aug Sep Oct
2021                2022

CERN-EOS ▲ Bari_HTC ▲ Bratislava_ARC ▲ Catania-VF ▲ CCIN2P3_HTC ▲ CCIN2P3_HTC_2 ▲ CNAF ▲ FZK_HTC ▲ GSI ▲ IHEP ▲ ISS ISS_LCG ▲ ITEP ▲ KFKI
▲ KISTI_GSDC ▲ Kolkata ▲ Kosice_ARC ▲ LBL_HPCS ▲ Legnaro_HTC ▲ NIHAM ▲ NIPNE_ARC ▲ ORNL ▲ Poznan ▲ Prague ▲ RRC_KI_T1 ▲ Strasbourg_IRES ▲ Subatech
▲ SUT ▲ Torino-HTC ▲ Trieste ▲ Troitsk ▲ UPB ▲ Wigner_KFKI_AF
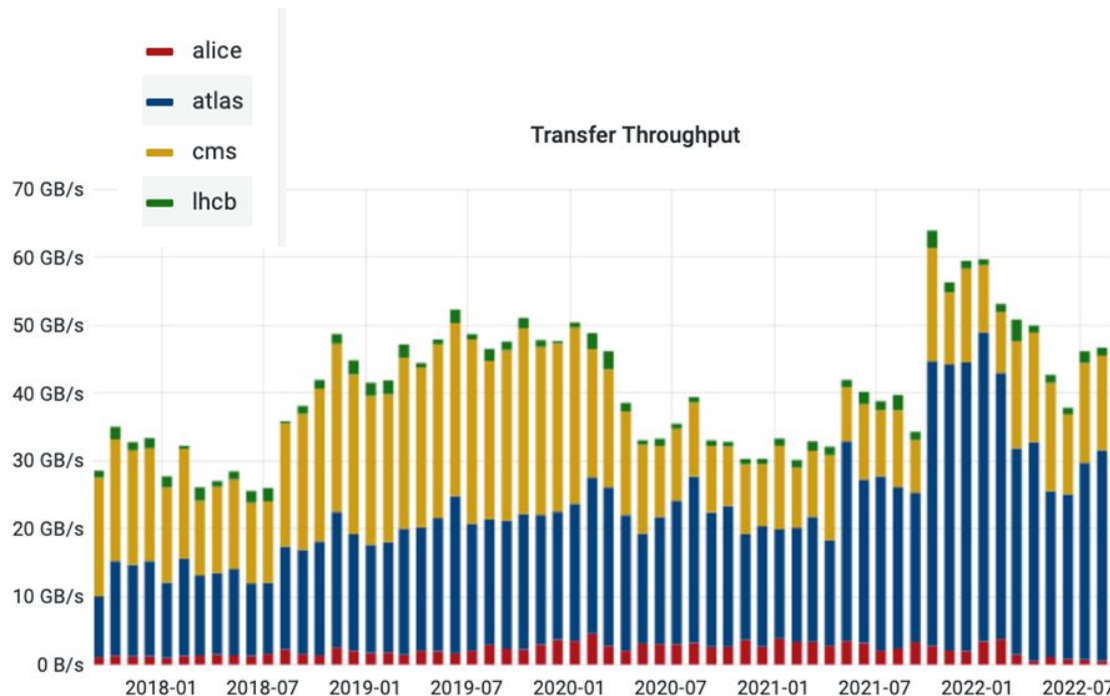
25

# ALICE in the big picture - WLCG data transfers

- Includes RAW data distribution and other LHCONE/LHCOPN transfers
- ALICE computing model and network use is beneficial to remote sites
- Network requirements are mild and well within the capabilities of regional T2s
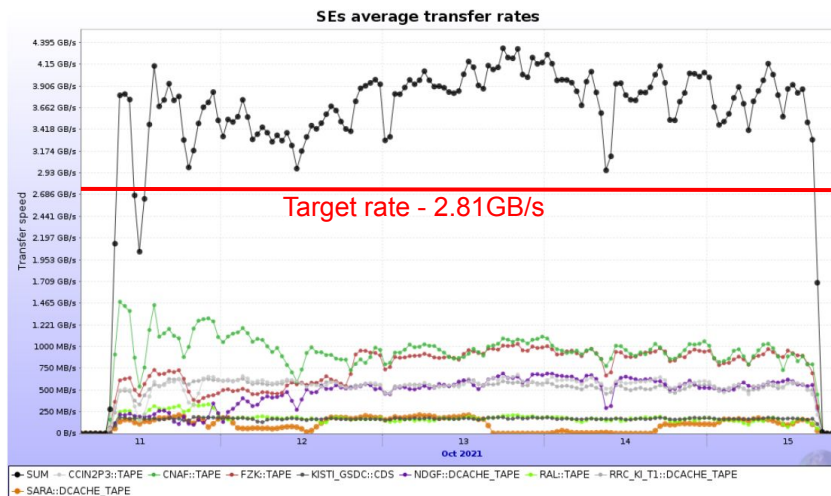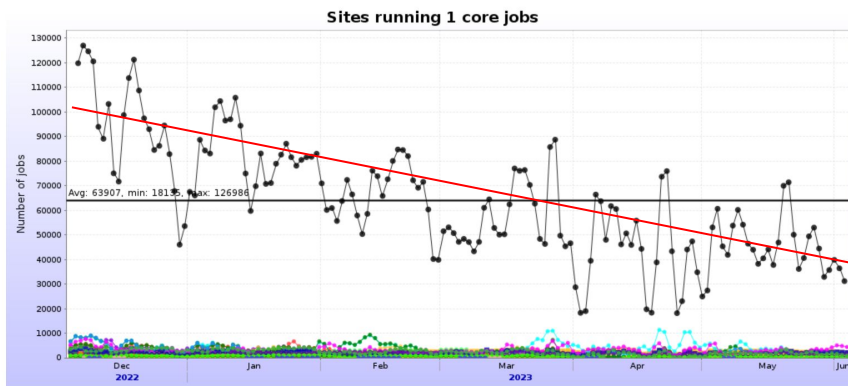


26

# Expected data rates in Run3 - replication of RAW

| T1 Centre | Target rate GB/s | Achieved rate GB/s |
|-----------|------------------|--------------------|
| CNAF | 0.8 | 0.94 (116%) |
| IN2P3 | 0.4 | 0.54 (130%) |
| KISTI | 0.15 | 0.16 (106%) |
| GridKA | 0.6 | 0.76 (123%) |
| NDGF | 0.3 | 0.47 (144%) |
| NL-T1 | 0.08 | 0.1 (122%) |
| RRC-KI | 0.4 | 0.53 (128%) |
| RAL | 0.08 | 0.17 (172%) |

Sum 2.81GB/s

- Full traffic simulated during data challenge
- Channels tuned to slightly above the target rate, within reasonable limit
- The bulk of the bandwidth will be used after the Pb-Pb data taking period, for ~3 months
  - Since there is no Pb-Pb this year, we remain at the level of data challenges



SEs average transfer rates

Target rate - 2.81GB/s

27

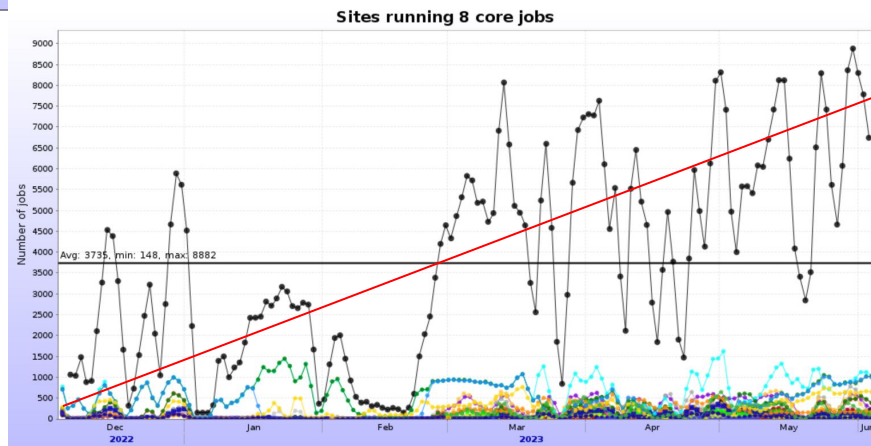# Multicore use



Sites running 1 core jobs

Steady decline of single-core payloads
- Legacy analysis and MC

Proportional increase of multicore
- 2022 and 2023 data processing
- MC
- New organized analysis (Hyperloop)
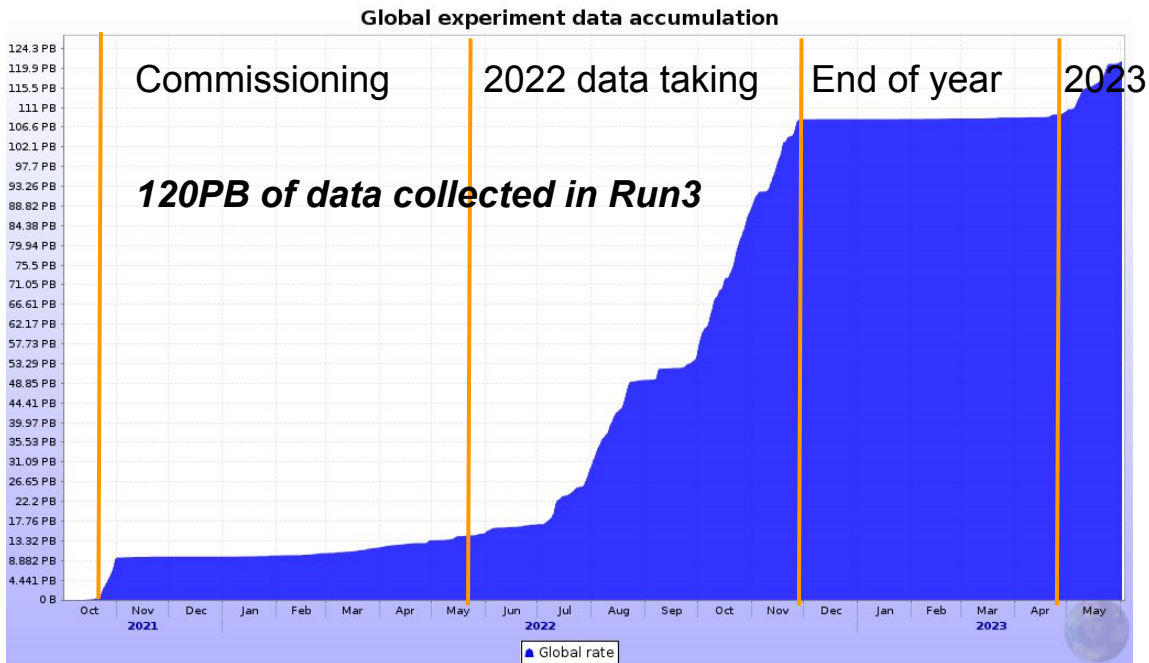


Sites running 8 core jobs

# Site upgrade to 8-core

- Site conversion to 8 cores ongoing
  - ~90% of Grid capacity already there
  - 100% of FR sites on 8 core queues
- Good experience with whole node submission
  - Steady running (LBNL Lawrencium, Perlmuter) + ORNL + GridKA (in progress) + KISTI
  - All HPC resources are whole-node, use will expand
  - Possibility to improve job performance
    - ~8% reduction in execution time through optimal NUMA assignment
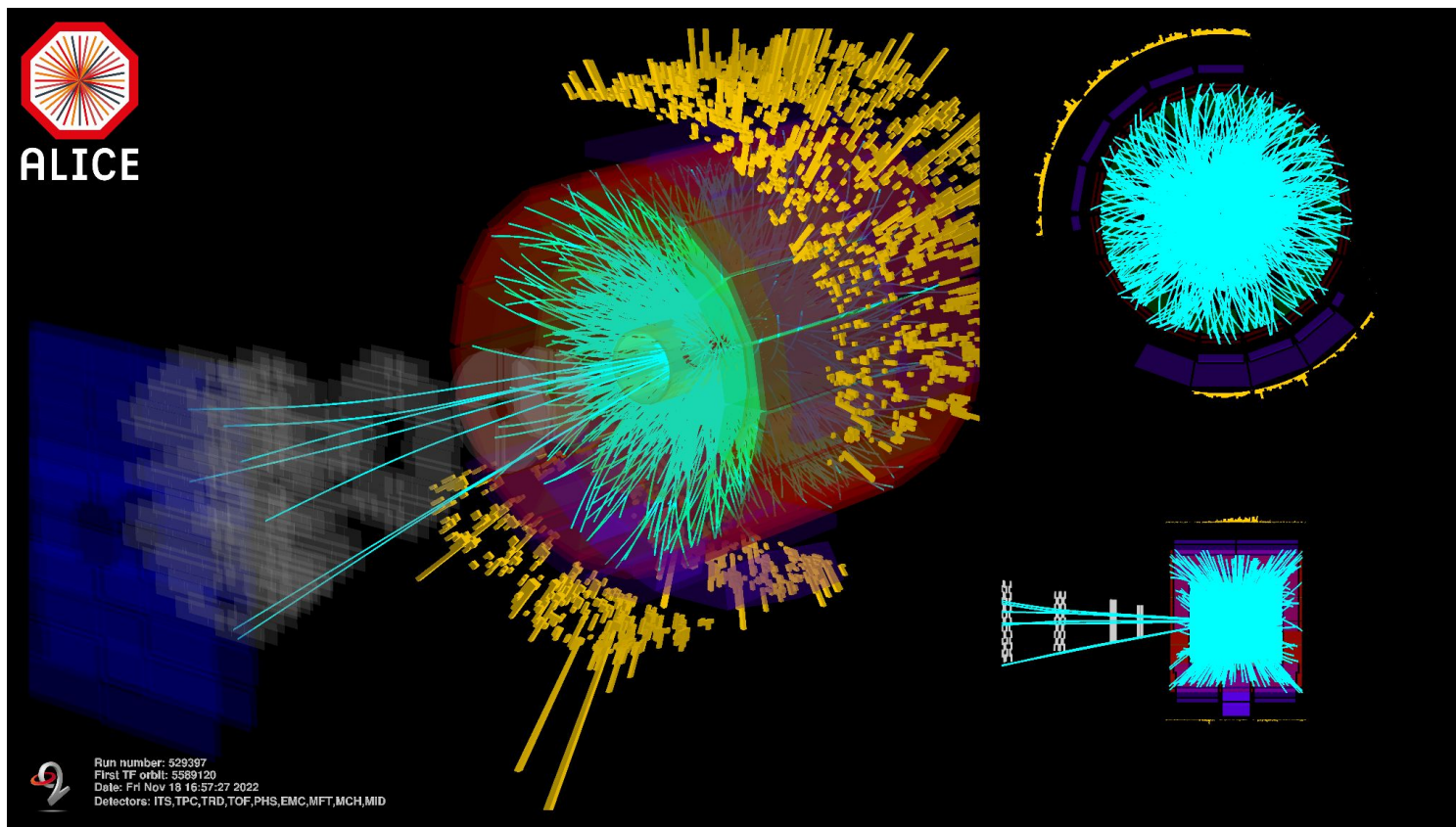    - More flexibility with CPU vs. I/O intensive tasks

| | AliEn proxy | | LDAP | | CVMFS | |
| Service | Status | Time left | Status | Cores ▲ | Status | Revision |
|---|---|---|---|---|---|---|
| 48. Perlmutter | - | - | | 64 | | 15557 |
| 20. EPN | - | - | | 16 | | 15557 |
| 74. Wigner_KFKI_AF_8core | | 1d 23:23 | | 8 | | 15557 |
| 72. Vienna | | 1d 23:13 | | 8 | | 15557 |
| 71. UPB | | 1d 23:31 | | 8 | | 15557 |
| 70. UNAM | | 1d 23:52 | | 8 | | 15557 |
| 69. UIB_LHC | | | | 8 | | 15557 |
| 66. Torino-HTC | | 1d 23:18 | | 8 | | 15557 |
| 64. Subatech_CCIPL | - | | | 8 | | 15557 |
| 63. Subatech | | 1d 23:26 | | 8 | | 15557 |
| 60. SARFTI | | 1d 23:25 | | 8 | | 15557 |
| 59. SARA | | 1d 23:37 | | 8 | | 15557 |
| 58. SaoPaulo_HTC | | 1d 23:01 | | 8 | | 15557 |
| 57. RRC_KI_T1 | | 1d 23:12 | | 8 | | 15557 |
| 55. RAL | | 1d 23:15 | | 8 | | 15557 |
| 54. Prague | | 1d 23:44 | | 8 | | 15557 |
| 52. Polaris | | 1d 23:17 | | 8 | | 15557 |
| 51. PNPI | | 1d 23:27 | | 8 | | 15557 |
| 45. NIPNE_ARC | | 1d 23:29 | | 8 | | 15557 |
| 44. NIKHEF | | 1d 23:01 | | 8 | | 15557 |
| 42. Nemesis | | 1d 23:48 | | 8 | | 15557 |
| 40. Legnaro_HTC | | 1d 23:55 | | 8 | | 15557 |
| 35. KFKI | | 1d 23:14 | | 8 | | 15557 |
| 34. JINR_ARC | | 1d 23:31 | | 8 | | 15557 |
| 31. IHEP | | 1d 23:00 | | 8 | | 15557 |
| 29. Hiroshima | | 1d 23:39 | | 8 | | 15557 |
| 27. GSI_8core | - | | | 8 | | 15557 |
| 25. GRIF_IRFU | | 1d 23:00 | | 8 | | 15557 |
| 24. GRIF_IPNO_IJCLAB | | 1d 23:22 | | 8 | | 15557 |
| 22. FZK_HTC | | 1d 23:19 | | 8 | | 15557 |
| 21. FZK | | 1d 23:56 | | 8 | | 15557 |
| 19. DCSC_KU | | 1d 23:59 | | 8 | | 15557 |
| 17. CNAF-DUE | | 1d 23:30 | | 8 | | 15557 |
| 16. CNAF | | 1d 23:46 | | 8 | | 15557 |
| 15. Clermont_ARC | | 1d 23:13 | | 8 | | 15557 |
| 14. CERN-ZENITH | | 1d 23:53 | | 8 | | 15557 |
| 13. CERN-TRITON | | 1d 23:07 | | 8 | | 15557 |
| 12. CERN-SIRIUS | | 1d 23:54 | | 8 | | 15557 |
| 11. CERN-MIRAGE | | 1d 23:09 | | 8 | | 15557 |
| 10. CERN-CORONA | | 1d 23:10 | | 8 | | 15557 |
| 8. CCIN2P3_HTC_2 | | 1d 23:44 | | 8 | | 15557 |
| 7. CCIN2P3_HTC | | 1d 23:28 | | 8 | | 15557 |
| 2. Bari_HTC | | 1d 23:51 | | 8 | | 15557 |
| 1. Altaria | | 1d 23:44 | | 8 | | 15557 |
| 26. GSI_4core | - | | | 4 | | 15557 |
| 76. Yerevan | | 1d 23:02 | | 1 | | 15557 |
| 75. WUT | | 1d 23:59 | | 1 | | 15557 |
| 73. Wigner_KFKI_AF | | 1d 23:27 | | 1 | | 15557 |
| 68. Troitsk | | 1d 23:14 | | 1 | | 15557 |
| 67. Trieste | | 1d 23:16 | | 1 | | 15557 |
| 65. SUT | | 1d 23:43 | | 1 | | 15557 |

29

# 2022 + 2023 data collection and processing

- After a period of detector commissioning in 2021-2022
- Steady p-p data taking in 2022
- Dec 2022 to Apr 2023 - calibration and processing, followed by skimming
- 2023 has started well and is the first year with Pb-Pb beam (looking forward to it)

**Global experiment data accumulation**

Commissioning          2022 data taking          End of year          2023

*120PB of data collected in Run3*

Global rate

30

# Event from 18 Nov 2022, low IR Pb-Pb@5.36TeV

# Summary

- During the LHC shutdown ALICE upgraded the detector and entire software stack
- Change of physics focus, triggerless readout, up to 100x event rate
  - Requires online compression and offline filtering
  - To stay within the resources envelope of WLCG
  - Steady data taking and processing since start of Run3 (2022 onward)
- New O2 software - multiprocess/shared memory, multicore
- New central GRID software
  - Entirely rewritten in Java
- JAliEn keeps the logic and functionality of AliEn
  - It is faster and simpler to deploy and operate
  - Much easier to maintain and to add new features
  - Incorporated support for GPUs and HPCs
- Site operation is simplified and made more reliable
  - Lowers the threshold for new sites joining (not the current trend…)

32

# Thank you!

- To the entire French GRID community for the excellent and steady support

- ***Farewell to Subatech as GRID site***
- Special thanks to the team and especially to ***Jean-Michel Barbet***

- In the ALICE Grid since the start of operation many years ago…
- One of the first sites to adopt, test and debug new software or principles of operation
- We are very sorry to see it go… and wish JM Happy Trails!



33