

# Analysis Facilities LHCC charge experiments questions

|   |          |
|---|----------|
| <b>Analysis Facilities LHCC charge<br/>experiments questions</b>    | <b>1</b> |
| DRAFT DRAFT Analysis Facilities: Questions for the WLCG Experiments | 2        |
| ALICE   | 3        |
| ATLAS   | 7        |
| CMS   | 7        |
| LHCb  | 10       |
| Experiments Common  | 10       |
| CERN input for WLCG AF discussion (from AF@CERN WG)                 | 11       |
| Other input (Markus)  | 12       |

##### document for LHCC has been moved into this new (almost) [clean gdoc](#) . Please keep on commenting here, we will incorporate the changes in the final doc. #####

This document regards the first part of the charge below.

## **1- Establishing a list of questions**

*Analysis Facilities is a broad topic. The use cases and the expectation of the experiments, and the analysers may evolve over time. The scope and the expected content of the document to be provided by the experiments must be clearly defined. To this end, a list of questions must be defined first that seek to define the expectations from experiments for Analysis Facilities. The questions, to be answered by experiments, must be picked such that the answers are useful for sites and are representative of a broad spectrum of analyses and analysers. The list of questions might evolve in the future.*

*The LHCC charges WLCG with establishing such a list. This process will be iterative and must involve the HSF Analysis Facilities Forum, the sites and the experiments. The monthly GDB (Grid Deployment Board) meetings of WLCG could be a place for the relevant discussions to occur, but it is up to WLCG to establish the appropriate mechanism. A first list will be established by the June 2024 LHCC meeting and presented by a WLCG representative at the LHCC to the LHCC for comments.*

Experiments can list their own questions separately before we get to a common list.

## DRAFT 2.0 (16th May at 18:45)

### Analysis Facilities: Questions for the WLCG Experiments

1. How does the experiment expect the Analysis Model will evolve for Run4 (and Run5) compared to Run3, considering both evolutionary advancements and potentially disruptive revolutions?
  - a. Please briefly describe the Run3 Analysis Model and highlight the most important expected changes foreseen/planned for Run4 (and Run5)
  - b. Which main analysis workflows run in Run3, and which types of workflows do you foresee being needed in Run4 (and Run5).
  - c. How many data reduction steps? How tightly chained are they?
  - d. If reduced data formats exist (.e.g. PHYSYLITE and nanoAOD): How many analysis will be covered by them?
  - e. ATLAS-1, ATLAS-2, ATLAS-4, CMS-1, CMS-3, CMS-4, CMS-5, LHCb-1
2. Please describe the data formats used today and in 2030 for analysis?
  - a. Data volumes (per year and total per Run, data and MC)
  - b. The number of versions and the number of replicas?
  - c. How many will be centrally managed?
  - d. Will analysis need to access extra information from other resources?
  - e. ATLAS-1, ATLAS-5, CMS-1
3. How much compute power is used today for analysis?
  - a. How much is coming from pledged resources and how much from unpledged resources?
  - b. Do you have any estimate of how much from local interactive and how much from local batch?
  - c. ATLAS-1, ATLAS-5
4. What are the main pain points that users experience today in analysis and how does the experiment plan to improve them in the coming years?
  - a. What features are currently found to be most beneficial by users? What is missing for a more effective analysis?
  - b. Do you foresee any technological or infrastructure evolution/revolution that would help in improving the analysis experience?
  - c. ATLAS-6, ATLAS-7, ATLAS-8
5. What is an Analysis Facility from the experiment point of view?
  - a. Please briefly describe the present status of AFs in the experiment.
  - b. Which functionalities should be included in an AF, and which would not make sense to include there (please refer also to the [HSF AF](#))?
  - c. Is an AF using resources locally or is it distributed or a mix?
  - d. For the local resources, would AF cover interactive, batch or both?
  - e. Using unpledged or pledged compute power (or both)?
  - f. Tailored for certain analysis workflows and/or specific working groups?
  - g. Should this support all the users of an experiment, or only some of them?
  - h. ALICE-1, ALICE-7, ALICE-8, ATLAS-1, ATLAS-3, CMS-2, LHCb-1, LHCb-2

6. What are the capabilities and support model that you would need from the AF for them to be effective?
  - a. What technical capabilities?
  - b. Which specialised HW?
    - i. GPU
    - ii. High Memory: do you need special nodes?
    - iii. Disk with high IOPS? NVME?
    - iv. Caches?
    - v. Specific type of storage?
    - vi. What bandwidth?
  - c. Personpower support
    - i. installation, quick answering to users needs, documentation
  - d. ALICE-4, ALICE-5, ATLAS-9, LHCb-4, LHCb-5
7. What is the motivation for deploying AF?
  - a. Based on the above answers, please describe your strategy towards deployment and integration of AF.
  - b. Please detail also whether analyses would be centrally organised on the AF, and if so how.
  - c. Please estimate how many AF would be needed for each type defined in Q5.
  - d. ALICE-1, ALICE-6, ALICE-3
8. Are you already able to identify a few use cases of analysis which could be useful to benchmark an AF of the type(s) defined in Q5?
9. List R&D that is being done to help address your concerns in questions Q1 and Q4?
  - a. CMS-7

→ uncovered: ALICE-2, CMS-6, LHCb-3

===== END DRAFT =====

## ALICE

### 1. What is the motivation for deploying Analysis Facilities (AF) in ALICE?

The ALICE Run3/4 upgrade Technical Design Report (TDR) for computing [1] outlines several key elements designed to enhance analysis efficiency and turnaround time. These elements include:

- Simplified and streamlined analysis data model: Facilitates faster internal navigation of the persistent data, enabling multicore and parallel processing capabilities. It also ensures **universality** for all major physics analyses.
- Enhanced organized analysis tools for parallel processing: The focus here is on the development particularly in two aspects:
  - Separation of data streaming and analysis tasks **to improve efficiency**.
  - Combining tasks from different physics analyses **to minimize I/O requirements**.
- Dedicated computing facility: This facility will enable rapid turnaround for analysis of the **entire data set residing on-site**, with a target of a **maximum of 12 hours**. The

datasets consist of selection of statistics, up to 10% of the entire ALICE AO2D sample

These elements have been successfully implemented by the start of Run3 in 2022 and are used in production. Some details follow below.

## 2. What is the new ALICE Analysis AOD design?

The new ALICE analysis framework has been presented at two consecutive CHEP conferences in 2019 and 2021. The following description is extracted from the published articles [2,3].

The analysis data model in Run 3 is a collection of flat tables, arranged in a relational database-like structure using index connections. Individual tables are represented as distinct C++ types, defined by their list of columns and backed by Arrow ChunkedArrays in memory. Any self-contained cluster of data, like the one containing the information about collisions or the tracks, is represented as a single table. The information stored in separate tables can be combined by performing join operations (similar to that in SQL). In addition, using flat data structures avoids the deserialisation overhead and significantly improves the IO speed. It also makes possible the use of bulk IO interface in ROOT, which is much more efficient compared to the per-event one.

The AO2D data format provides a significant reduction in required storage in a number of ways. It is based on the ROOT tree format that provides a good compression ratio. Additionally, floating point values are truncated to their respective uncertainties. Finally, only a minimal set of values is stored, while the others are calculated on-the-fly based on the analysis needs. This leads to reduction in the AOD size in terms of disk space by a factor of five or better, as measured with converted Run 1/2 data. The change in the data model, the reduced data size and the bulk IO improved by one order of magnitude the analysis throughput (measured as events per second).

## 3. How to organise the analysis on the AFs?

ALICE analysis tasks mostly deal with large datasets, on the order of several PB per dataset. Since LHC Run2, the bulk of the analysis was done through the organised analysis system 'Lego', which combined multiple physics tasks into a common workflow, running on the same data and on the Grid infrastructure. This system brought two main benefits - minimising the I/O requirement of the analysis and streamlining the analysis process by assuring predictable turnaround time, as opposed to individual user analysis competing for the same resources.

In Run3, 'Lego' was rewritten and considerably improved into the new system 'Hyperloop', adhering to the same basic principles, however leveraging the increased performance of the new AO2Ds and the multicore queues to increase the analysis efficiency. In addition to the data processed to AO2Ds from Run3 onward, the AODs from Run1 and Run2 were converted to AO2Ds, making the Hyperloop a universal tool for all ALICE physics analysis on a large scale. Since Hyperloop operates on the concept of dataset, running on the AFs

does not require any modification of the code, but just a pointer to the AF, where the analysis should be performed. The tool is also coupled with the ALICE data transfer utilities, which pre-place the desired data on the AF storage, before launching the Hyperloop run.

#### **4. What is the general vision for AF and what are the considerations, structure and operation?**

Design and setup of an AF is based on several key considerations:

1. **I/O Performance:** The storage element (SE) to CPU bandwidth needs to be at least 10MB/s/core. This translates to specific requirements for the combined site structure: global throughput of the SE and local area network (LAN), and individual worker node throughput based on the number of cores per worker node.
2. **WAN connectivity:** Since AFs are primarily data consumers, they require adequate wide area network (WAN) connection to the Grid for data transfers from Grid SEs.
3. **Software reuse:** Existing Grid data transfer and management tools should be used also to manage the data on the AFs.
4. **Workflow compatibility:** AFs should use the same access methods, workflows and software packages as on the Grid, with minimal modifications.
5. **Minimum amount of resources:** Sites must satisfy a minimum threshold of storage capacity and CPU cores to be suitable for AFs.

#### **5. What are the possible target sites for AF deployment and specific implementation details?**

Many existing Grid sites meet the I/O and WAN requirements. ALICE prioritises using or adapting existing Tier-2 sites due to their flexibility and to have minimal disruption to high-priority T0/T1 operations, including interference with other VOs running on these.

Minor modifications were made to the ALICE Grid software jAliEn to fulfil the third and fourth considerations as outlined above. These modifications included:

- Possibility to define specific data transfer paths for the AFs.
- Introducing the concept of data sets with percentages of production statistics and transfer these upon request by operators running the Hyperloop analysis trains.
- Segregating AF storage for data bookkeeping, dataset lifetime management, transfer prioritisation, and space quotas.

AFs are integrated into the Grid infrastructure as Tier-2 sites. Users access them through standard Grid mechanisms and submit payload through the jAliEn workload management system with the option to target workflows specifically at the AFs. The AF data access remains the same as for regular Grid payloads.

#### **6. How many AFs are needed and what are the operational benefits?**

The system allows for multiple coexisting AFs with non-overlapping data sets at each one while the minimal size of an AF avoids significant resource fragmentation. This approach

leverages the already developed Grid and analysis software, minimizing the effort to incorporate the AF concept. It offers several advantages:

- Faster analysis turnaround: This benefit is crucial for high-priority analysis and iterative software testing, which later can be run on large datasets on the Grid.
- Reduced site burden: Volunteer Tier-2 sites can be assigned as AFs without needing additional support functions, procedures and personnel.
- The total data sample at the AFs is approximately 10% of the total AO2D and each AF holds a non-overlapping portion of this sample.

With respect to the total amount of expected AO2D volume in Run3 and Run4, the aggregated AF storage capacity to be able to accommodate 10% of AO2Ds should be approximately 12PB and to process this data in the required timeframe, 18K CPU cores. For operational reasons, a single AF should be aimed to cover at least 20% of the required total CPU and storage capacity.

## 7. What are the plans for interactive AF use?

During Run1/2 ALICE has extensively used interactive AF for fast prototyping and validation of analysis code. The old approach was based on Proof clusters. While it provided speed-up and fast feedback, it also showed significant limitations of such interactive analysis. These limitations can be listed in the following way: inefficient use of resources, non-scalable solutions, lack of robustness, significant maintenance burden. Based on these observations and taking into account that Proof is obsolete, the ALICE collaboration decided to keep the interactive analysis on the desktops and laptops of the users, while the AF are used for organised analysis as explained above.

## 8. What is the current state of affairs and future improvement goals?

ALICE currently operates three AFs: GSI Darmstadt (fully dedicated), Wigner Budapest (partial resources), and LBNL Berkeley (partial resources). These AFs offer a total capacity of approximately 10PB storage and 12,000 CPU cores, representing sufficient resources for the data sample expected for Run3. These will naturally grow toward Run4 to cover the 10% target data sample for the combined Run3 and Run4 dataset. In the past two years, hyperloop trains have run over 10,000 times on AF datasets, consistently achieving an average turnaround time within the specified 12-hour limit.

As a future improvement of the AF functionality it can be considered to add accelerators to the existing or to add new hardware, which will enable the data processing using fast and efficient machine learning algorithms. These can be added on specific AF, in which case the appropriate datasets must be placed there, or on all. The addition of accelerators will not require changing the workload and data management tools.

[1] ALICE Collaboration. *Upgrade of the Online-Offline computing system, technical design report*. ALICE-TDR-019 ed., CERN-LHCC-2015-006, 2015.

[2] <https://doi.org/10.1051/epjconf/202024506032>

[3] <https://doi.org/10.1051/epjconf/202125103063>

## ATLAS

1. What is the standard analysis model for the majority of analyses performed in the experiment, and how do you expect it to evolve over the next 5 years, in terms of:
  - a. Input data per analysis (in bytes, events, files, and datasets)
  - b. Data reduction steps (number, resources on which they run)
2. What fraction of analyses that the experiment performs do you expect to evolve to use the standard analysis model?
3. For the analyses that will not evolve in the above way, what expectations do you have? What sorts of analyses will be included in this fraction, and how might their resource requirements differ from the standard analyses?
4. Are there any opportunities for disruptive technology changes / analysis model changes in the next 5 years? If so, what might they be, and how might they be rolled out to the experiment?
5. Based on the above model and related considerations, what resources do you expect to be used for data analysis in 2030?
6. Looking at today's global analysis infrastructure (grid, cloud, local cluster, national analysis facility, etc), what technologies are missing that should be provided in order to make analysis more efficient in 2030?
7. Looking at today's global analysis infrastructure, what (r)evolutionary steps are needed to improve the analysis experience in 2030?
8. Based on the most popular sites for (interactive) analysis today, what features do you believe drive users to those sites? Is it reasonable to expect that if other sites have similar features, they would be equally popular?
9. How do you see the expected support model for Analysis Facilities when compared to WLCG Grid Sites, in terms of required personpower for the installation and deployment, as well as the required monitoring and daily operations?

## CMS

Analysis facilities (AF) in the sense of infrastructures and services that provide platforms for analysing physics data with the necessary software and access to computational resources have existed since the early days of CMS, examples being the LPC at Fermilab and the CAF at CERN. Note that interactive computing resources, which some existing AFs are partially built upon, are not pledged.

The LHCC has charged the experiments to “define the expectations for Analysis Facilities” by writing a list of questions to be answered by the experiments periodically. This should result in a list of requirements, e.g. features and use cases to be supported. In order to gauge whether sufficient progress is being made in fulfilling these requirements, the experiments are asked to periodically report on whether the desired use cases and features are supported by the R&D initiatives that the experiments participate in and, if not, whether there are concrete milestones to do so. Use cases and features could be absolute requirements or “nice-to-have” options.

1. What are the experiment-specific elements of the Run 4 computing model which have a bearing on the quantitative requirements for analysis? In particular:

- a. What input data formats will be used for analysis and by what fractions of analyses?
  - b. What are the data volumes for each format?
  - c. What is the average fraction of the overall dataset needed for a typical analysis?
  - d. What is the expected turnaround time for an analysis using a given data format?
  - e. How many centrally-managed replicas of data in each format will be available for analysis?
  - f. Will the experiment use caches, or is this up to the sites?
  - g. Will the access interface be POSIX, through object stores, or both?
  - h. What is the role of skims in data reduction? And should they be restricted to central skims?
  - i. What are the requirements for bandwidth to and from tape, disk storage, and WAN connectivity, etc. to serve the data?
2. Does the experiment require the particular end-user analysis use cases enumerated in the HSF White Paper on Analysis Facilities [1], in particular:
- a. Ability to perform fast research iterations on large datasets interactively
  - b. Ability to convert interactive to batch-schedulable workloads
  - c. Ability to interact with the WLCG and scale outside of the facility on occasion
  - d. Ability to efficiently train machine learning models for HEP
  - e. Ability to reproducibly instantiate desired software stack (at one or more sites).
  - f. Ability to collaborate in a multi-organisational team on a single resource
  - g. Ability to move analyses to new facilities (including forward compatibility)
  - h. Ability to efficiently access collaboration data as well as make intermediate data products available to the team
  - i. Ability to express interdependent distributed computations at small and large scales
3. What additional analysis use cases does the experiment foresee supporting?  
Examples could include:
- a. Availability to all experiment members, or members of multiple VOs
  - b. Ability to re-run only the computation affected by a change rather than the entire analysis
  - c. Availability of a declarative analysis description language which is independent of the execution backend
  - d. Availability of tools for seamless handling of systematic uncertainties, multidimensional histogramming, etc.
  - e. Ability to add columns to a data sample easily, e.g. from larger data formats, and join columns, e.g. from more than one input file.
  - f. Support for different modes of execution, e.g. single- and multi-core, local vs. distributed, quasi-interactive vs. batch) with seamless transition and low latency
  - g. Access to accelerators
  - h. Ability to be forward compatible with a new technology (a new Dask ++ tool, or anything of this sort)



- i. ... list could be augmented by other experiments ...
4. Does the experiment require other features that were described in the HSF White Paper on Analysis Facilities, in particular:
  - a. A secure, unified, federated identity for authentication to storage and computing resources.
  - b. Common namespace for storage
  - c. A sustainable support model (both for user support and infrastructures).
  - d. Standard exercises and end user documentation
  - e. Monitoring and metrics, including benchmarking applications to evaluate overall throughput in light of the experiment-specific assumptions about data formats and volumes.
  - f. CVMFS
  - g. Conda
  - h. Standard Linux containers
5. Does the experiment require any additional features or services? Examples could include:
  - a. Frontier/Squid
  - b. CI/CD pipelines
  - c. Common environments and user interfaces
  - d. Scale out infrastructures like OpenShift or Kubernetes
  - e. ... list could be augmented by other experiments ...
6. Since the last review are there new use cases or desired features that should be added to the lists in questions 3 and 5?
7. Questions about R&D initiatives or proto-facilities:
  - a. Which R&D initiatives or proto-AFs are experiment members participating in?
  - b. Do these individual initiatives support the desired use cases and features or not, or do they have milestones to support them in the future? What milestones were accomplished, dropped, or modified since the last review?
  - c. Are there any initiatives that support all of the desired use cases and features?
  - d. Are there any desired use cases or features that are not supported and not planned to be supported by any of the initiatives?
  - e. Have there been any benchmarking tests to evaluate the ability of the initiatives to scale up to experiment-specific requirements, or are there milestones to perform such tests in the future?
  - f. Are there published results about the initiatives or physics analyses that leveraged the proto-facilities?

[1] <https://arxiv.org/abs/2404.02100>

## LHCb

The points raised by the Analysis Facilities whitepaper[1] are very valid, and we agree with the way CMS translated them to questions. Among those, we find important to clarify:

1. whether Analysis Facilities should provide built-in support for workflow management systems (such as REANA) ? Who will contribute to the upkeep of HEP-specific integrations for these systems in this case?
2. How will experimental fair use policies for distributed computing look within the framework of an AF, will they cover usage from both direct job submission and worker-based modes of use?
3. will Facilities entirely supersede existing centralised resources at CERN (e.g., CERN Condor pools)?
4. what technologies/packages/services will be supported for users to interface with AFs (e.g., HTCondor, Dask)?
5. the criteria that “new” technologies/packages/services will have to meet to be officially supported by AFs?

Furthermore, questions narrowing the needs in terms of:

1. documentation/training
2. monitoring of the use of the infrastructure (to know which technical features, how much data is transferred and how...)

will be crucial to encourage their use and to allow improving them in a way that is efficient for the sites and the experiments.

## Experiments Common

It would be very Interesting to know if the sites' R&D activities are aligned with experiment needs.

Discussion about the ALICE analysis facility: a consolidated concept in use already. Is an ALICE analysis facility something that would suit the other experiments.

Do the services provided today by Grid sites fit the purpose of an analysis facility ? (see example of ALICE).

Which hardware type do you need in an analysis facility ? E.g. I/O specification, GPUs, ...

How much does the Analysis Facility model of the 4 experiments need to be compatible? (a question more for the sites)

Does analysis have to be portable from one site to the other? What does that imply?

CERN input for WLCG AF discussion (from AF@CERN WG)

- Software provisioning
  - Are LCG stacks enough?
  - Is the installation of custom software packages required? How is this installation expected to be done?
  - Any other needs?
- Interface
  - What elements are expected in an AF interface?
  - What is now provided by the CERN AF Pilot:
    - JupyterLab interface
    - Notebooks for interactive computing
    - Terminal for command line access
    - Dask extensions to reserve resources on condor and to monitor progress of execution of distributed analysis
  - Is anything else expected (e.g. IDE)?
- User experience / interactivity
  - What is expected in terms of user experience when running an interactive analysis distributedly?
  - The CERN AF Pilot supports distributed interactive analysis (computation results obtained in seconds, minutes) and semi-interactive (a few hours). It also allows to reconnect to an already running analysis (e.g. close laptop, attend a meeting, reopen browser to reconnect).
- Resource requirements, allocation
  - The CERN AF pilot provides access to HTCondor resources at CERN for running interactive analysis computations distributedly.
  - Are resources for analysis expected to be obtained progressively, for example by combining a layer of resources dedicated to interactive jobs plus the common pool subject to experiment quotas? This assumes that distributed analysis computations start as soon as some resource is available.
  - Are worker nodes (where distributed analysis computations run) expected to have the same HW requirements as current batch nodes?
  - Are GPUs required?
- Transition from interactive to batch
  - How is this transition expected to happen?
  - What help is expected from the AF to make this easier?
- Storage / data management
  - Is EOS user + EOS experiments enough?
  - Do I/O performance expectations differ from present local batch workloads? Does storage need to anticipate new workloads (AI/ML) with very different characteristics?
  - What integrations with sync'n'share (e.g. CERNBox) systems are expected?
    - What's the data sharing/collaboration model?
  - Is remote access to external datasets (not at CERN) required?
  - What integrations with DDM systems are expected? Rucio? FTS?
  - Is an extension for Rucio desirable? There is a JupyterLab Rucio extension to help with the data access and discovery for Rucio datasets.
  - Are proxies (XCache) to speed up reads required?
  - Do you need local storage for large outputs?

- Interoperability between AFs
  - Are AFs expected to interoperate? In what sense?
  - Will this be achieved by allowing remote access to data?
  - Is analysis code expected to be easily movable between AFs?
- Machine Learning
  - What are the HW requirements for Machine Learning in (semi)interactive analysis, in particular for training?
  - What are the SW requirements for Machine Learning in (semi)interactive analysis?
  - Do you need support for DAGS (workflows, pipelines)?

## Other input (Markus)

What are the improvements you expect to gain from Analysis Facilities?

- 1) Reduction of development time. (Concept2Production)
  - a) Improved environment for developing and probing new analyses.
  - b) This could be met with interactive access to fast responding resources, of suitable characteristics, mainly site local storage.
  - c) This includes the training and verification of ML models.
- 2) Reduction of time to become productive for new team members. (Lowering the access barrier) (Time2Productivity)
  - a) By using programming languages that students are familiar with
  - b) By using environments (notebooks) that are easy to get started with
  - c) This comes for the development process with similar requirements as for 1), but in addition the scaleout/transition to complex batch workflows, link to the storage management systems, sharing produced data etc. becomes more important..
  - d) When limited to medium sized data, this is close to what REANA tries to provide
- 3) Reduction of time to results/improved efficiency for analyses (Production2Results)
  - a) Here the focus would be on availability of suitable hardware with sufficiently high priority
  - b) Means to improve efficiency, like the ALICE Analysis Trains where data is read once and multiple analyses can share the same I/O costs.
  - c) Improved access to data, local and remote ( caches (XCache)), faster disks, closer integration with the tape system ( tape carousel )..
  - d) Interactive access can be helpful to sort out problems
  - e) Integration with complex, multistage workflows matters (DAGs), integration with the experiments system to prioritise work.
  - f) This could also be something as simple as a well connected batch system dedicated to analysis, with a few tweaks, like a more responsive scheduler ( NAF at DESY can serve as an example)
- 4) To reduce the time for the very last step in the analysis before publication. ( Production2Publication)

- a) Quick turnaround on highly pre-processed data
  - b) Interactive access via workbooks to react quickly and fine tune the analysis
  - c) Sharing the results for review and discussion within the team
  - d) Requirements very similar to 1 and 2, but less development of code/ ML models that transition to production
- 5) All of the above. AF as an alternative way to do all we are currently doing.  
(GrandUnifiedAF)