

# Analysis Facilities

## Summary and next steps

Jamie Gooding *on behalf of the Analysis Facilities Plenary organisers*

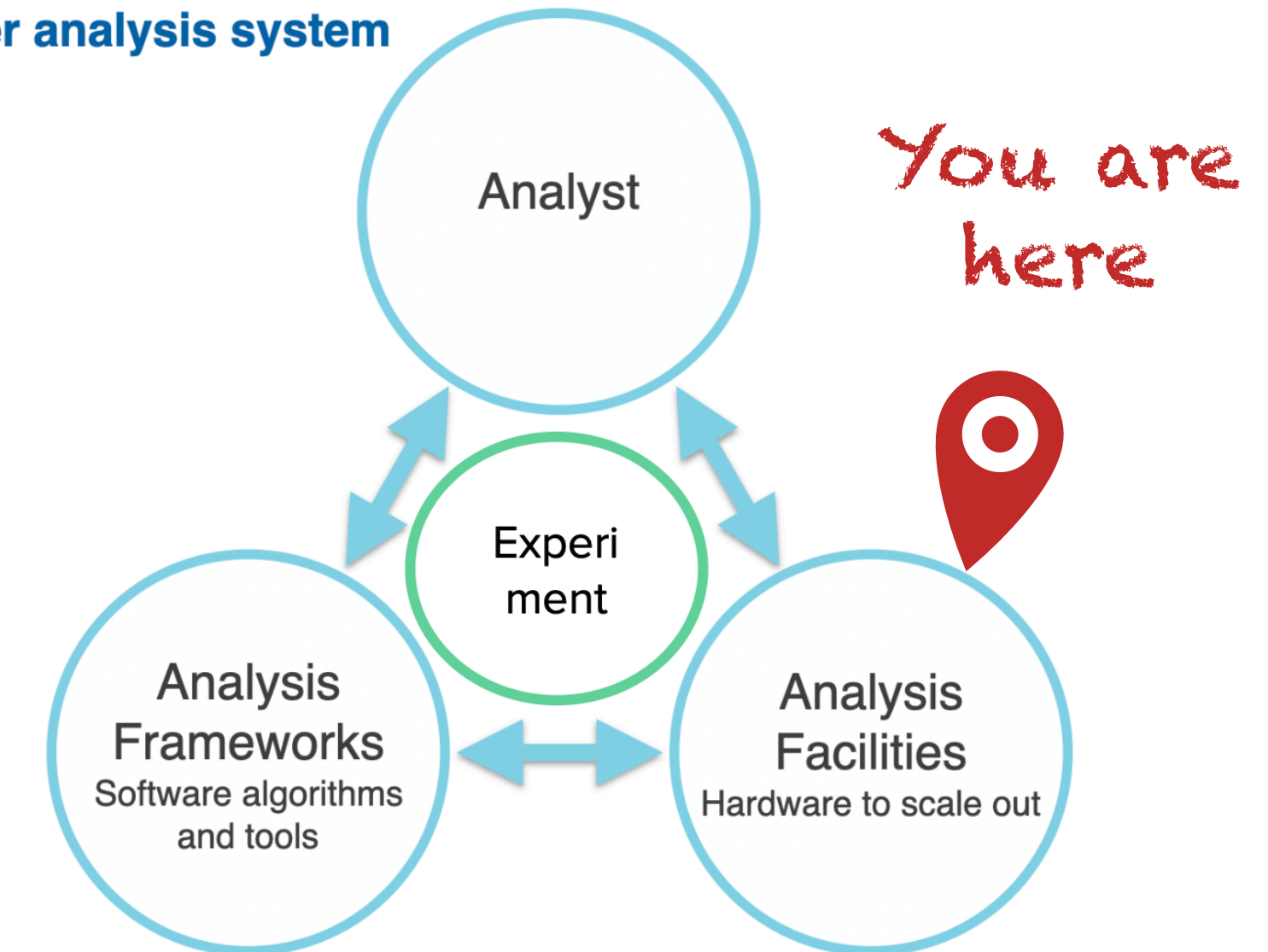
HSF/WLCG Workshop 2024  
Closing plenary, 17th May 2024



# Introduction

- Why discuss AFs at WLCG/HSF 2024?
  - Broad expertise in running existing AFs
  - Analysis in the HL-LHC era requires computing scaling
    - AFs provide a user-friendly route to deliver this
    - Need to define how this should look before moving forward
- LHCC charge: requested list of questions for experiments
  - Lively discussion yesterday, comments now implemented
- AF White Paper recently released ([arXiv:2404.02100](https://arxiv.org/abs/2404.02100))
  - Submitted to Computing and Software for Big Science

## A better analysis system



## Questions for the LHCC on Analysis Facilities, Alessandra Forti

### Analysis Facilities White Paper

D. Ciangottini<sup>1,b</sup>, A. Forti<sup>2,b</sup>, L. Heinrich<sup>3,b</sup>, N. Skidmore<sup>4,b</sup>,  
C. Alpigiani<sup>5</sup>, M. Aly<sup>2</sup>, D. Benjamin<sup>6</sup>, B. Bockelman<sup>7</sup>, L. Bryant<sup>8</sup>, J. Catmore<sup>9</sup>, M. D'Alfonso<sup>35</sup>, A. Delgado Peris<sup>17</sup>, C. Doglioni<sup>2</sup>, G. Duckeck<sup>10</sup>, P. Elmer<sup>11</sup>, J. Eschle<sup>12</sup>, M. Feickert<sup>13</sup>, J. Frost<sup>14</sup>, R. Gardner<sup>8</sup>, V. Garonne<sup>6</sup>, M. Giffels<sup>36</sup>, J. Gooding<sup>15</sup>, E. Gramstad<sup>9</sup>, L. Gray<sup>16</sup>, B. Hegner<sup>31</sup>, A. Held<sup>13</sup>, J. Hernández<sup>17</sup>, B. Holzman<sup>16</sup>, F. Hu<sup>8</sup>, B. K. Jashal<sup>18,19</sup>, D. Kondratyev<sup>20</sup>, E. Kourlitis<sup>3</sup>, L. Kreczko<sup>21</sup>, I. Krommydas<sup>22</sup>, T. Kuhr<sup>10</sup>, E. Lancon<sup>6</sup>, C. Lange<sup>23</sup>, D. Lange<sup>11</sup>, J. Lange<sup>34</sup>, P. Lenzi<sup>1</sup>, T. Linden<sup>24</sup>, V. Martinez Outschoorn<sup>27</sup>, S. McKee<sup>25</sup>, J. F. Molina<sup>17</sup>, M. Neubauer<sup>26</sup>, A. Novak<sup>35</sup>, I. Osborne<sup>11</sup>, F. Ould-Saada<sup>9</sup>, A. P. Pages<sup>28</sup>, K. Pedro<sup>16</sup>, A. Perez-Calero Yzquierdo<sup>17</sup>, S. Piperov<sup>20</sup>, J. Pivarski<sup>11</sup>, E. Rodrigues<sup>29</sup>, N. Sahoo<sup>30</sup>, A. Sciaba<sup>31</sup>, M. Schulz<sup>31</sup>, L. Sexton-Kennedy<sup>16</sup>, O. Shadura<sup>32</sup>, T. Šimko<sup>31</sup>, N. Smith<sup>16</sup>, D. Spiga<sup>1</sup>, G. Stark<sup>33</sup>, G. Stewart<sup>31</sup>, I. Vukotic<sup>8</sup>, G. Watts<sup>5</sup>,

# Updates to LHCC questions

# LHCC questions 1-3

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?

- Please briefly describe the Run3 Analysis Model and highlight the most important expected changes foreseen/planned for Run4 (and Run5)
- Which main analysis workflows run in Run3, and which types of workflows do you foresee being needed in Run4 (and Run5).
- How many data reduction steps? How tightly chained are they?
- If reduced data formats exist (.e.g. PHYSLITE and nanoAOD): How many analysis will be covered by them?
- 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?

- Data volumes (per year and total per Run, data and MC)
- The number of versions and the number of replicas?
- How many will be centrally managed?
- Will analysis need to access extra information from other resources?
- ATLAS-1, ATLAS-5, CMS-1

3. How much compute power is used today for analysis?

- How much is coming from pledged resources and how much from unpledged resources?
- Do you have any estimate of how much from local interactive and how much from local batch?
- ATLAS-1, ATLAS-5

Analysis Facilities: Questions for the WLCG Experiments

New sub-question

Rewording

## Analysis Facilities: Questions for the WLCG Experiments

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 for what concerns AF in the experiment.
  - b. Which functionalities refer, or could refer, to an AF, and which would not make sense to include there (please refer also to the [HSF AF](#))?
  - c. Is an AF using resources local or is it distributed or a mix?
  - d. For the local resources, would AF cover interactive, batch or both?
  - e. Using unpledged vs pledged compute power?
  - f. Tailored for certain analysis workflows and/or specific working groups?
  - g. Should this support all the users of an experiment, or only a part of it?
  - h. ALICE-1, ALICE-7, ALICE-8, ATLAS-1, ATLAS-3, CMS-2, LHCb-1, LHCb-2

Rewording



# LHCC questions 6-9

6. What are the capabilities and support model that you would need from the AF for you to have them 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

## Analysis Facilities: Questions for the WLCG Experiments

← Rewording

← New sub-question

← New question

## Last call for comments!

- Document **will be frozen to major comments after today**
- Minor comments (small language/grammar changes) still welcome

[Analysis Facilities: Questions for the WLCG Experiments](#)

[Analysis Facilities: Questions for the WLCG Experiments](#)



# Analysis Facilities Plenary talks



# 200 Gb/s analysis demonstrator

- Challenge of demonstrating analysis at 25% HL-LHC scale
  - ATLAS data implementation in ServiceX running at UChicago
  - CMS data implementation in Uproot running at Nebraska T2
- Many lessons learned:
  - Understanding memory usage in Python vital
  - Bugs caught (and fixed!), e.g., new Uproot version with XRootD fixes
  - Role of network imbalances
  - Best practices in XCache tuning

## Derived Values – Example CMS ‘napkin math’

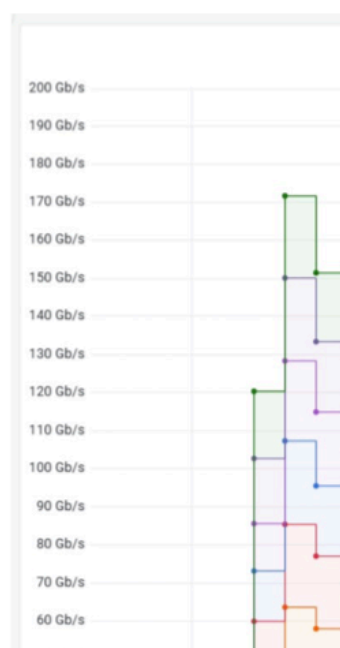
- ▶ Start with 200TB read in 30 minutes. => ~900Gbps sustained.
- ▶ 25% scale => 200Gbps sustained. Hence, **200Gbps challenge**.

## Uproot Results

- ▶ Highest data-rate configuration (TaskVine):
  - ▶ Data read (compressed): 58.33TB
  - ▶ Average data rate: 221Gbps
  - ▶ Peak data rate: 240Gbps
  - ▶ Files processed: 63,762 (17 failed)
- ▶ Highest event-rate configuration (Dask):
  - ▶ total event rate : 32,256 kHz
  - ▶ Processed 40,276,003,047 events total
  - ▶ Per-core event rate : 27.66 kHz

## ServiceX Results

- ▶ To reduce the overhead of small datasets, we ran on a subset that consisted of the bulk of the data.
- ▶ Highlight run:
  - ▶ 4 Datasets
  - ▶ 146TB total
  - ▶ 19,074,862,754 Events
  - ▶ 170Gbps
  - ▶ Limited to 1,000 pods.
  - ▶ Time: 32:28
  - ▶ Event Rate: 9,787 kHz



200 Gb/s analysis demonstrator, *Brian Bockelman*

# User experience discussion

- Discussion around requested (expected) UX at AFs
- A sample of the discussion:
  - Where to draw line between “installed” tools and what users bring with them?
  - How can users be supported without prescribing tools?
  - Where exactly do our current pain points lie?
  - How do ML tools fit into AF ecosystem?
  - Tools in AFs to ease use of Grid by beginners?

## Questions...

- What is impossible?
  - When should a job be done on the GRID rather than an AF?
- Can we (the experiments) collaborate on building these?
- What work gives us the most users?
- Are there technical decision we make that the user will care about?
  - But are unaware of?



I am going to go very fast through this talk!

## Tools & Environments

Installed Tools: notebooks, ssh, etc.?

What should the user bring with them?

Method to share container environments in an analysis group  
(e.g. docker, [dev-containers](#), binderhub)

Scaling Impact?

Current Analysis Facilities tend to specialize – is this the right way forward?

- Services provided – ServiceX, REANA, etc.
- Locally installed tools like snakemake...
- Should coffea be there?

Machine Learning



Toolset and workflows – plethora of workflows

- Access to GPU's (efficient!)
- Workflow support: don't make a choice?

[User Experience at AF, Gordon Watts](#)

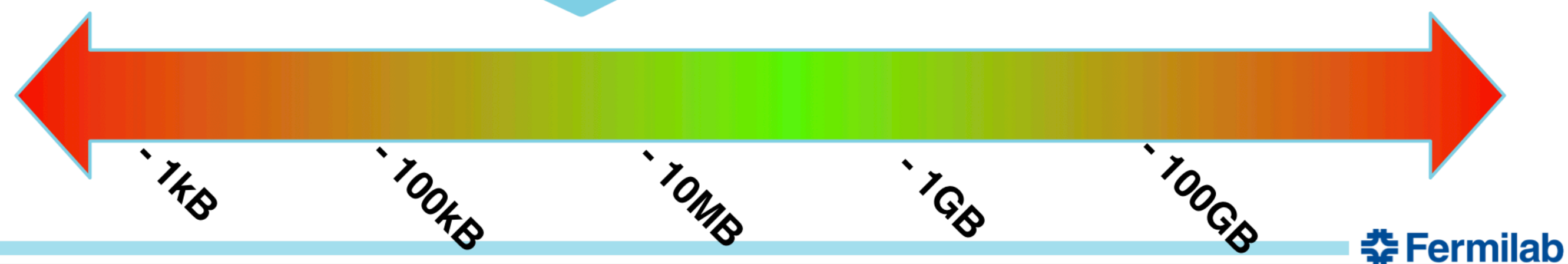
# Data access discussion

- Discussion of need to align data access with physics:
  - ~1 MB objects align well to physics content
- A sample of the discussion:
  - HEP applications generally built for POSIX
  - Community experience with POSIX → how does the path look for user adoption of object store?
  - How can/should token-based auth. be implemented at an AF?

## Can we align unit of data access to unit of physics content?

- Dataset = list of 2-4 GB files, totaling 10GB-1PB. Why?

- Sweet spot for access ~ 1MB
  - Few ragged columns for O(10k) events?
  - Many columns for O(1k) events?
  - Do we want small # events per unit?
- Whole-unit cache → off-shelf solutions
- Catalog challenge: need indirection



12 May 16, 2024 Storage for Analysis Facilities

Storage for Analysis Facilities, Dirk Hufnagel & Nick Smith

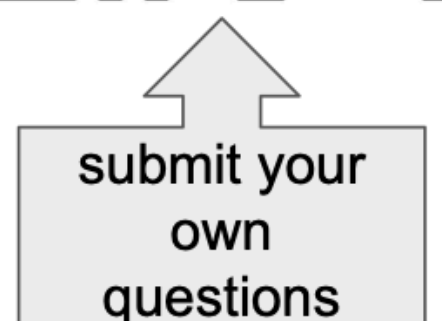
# User monitoring discussion

- Discussion of how best to monitor use of AFs
- Already many examples of monitoring best practice
- A sample of the discussion:
  - What is most useful for users/for AFs to see?
  - How to communicate to users key information (e.g., wait/fail reasons)
  - How to tell if users are suffering in silence?

## What are the questions, and from whom?

- Users (experience) <sup>1</sup>
  - What resources are available?
  - How long will my jobs wait in queue? Why do they run so slow? Why is my notebook hanging? Why did my last few jobs not finish? Why are my jobs held? Why did my jobs fail? Why are they being held?
  - How do I access my data? Is it local? How do I get X software installed? How do I run with my container?
- Resource providers (trends, performance, facility metrics) <sup>2</sup> <sup>3</sup> <sup>4</sup>
  - What resources (cpu, disk-capacity, disk-fast, network, gpu) are under-provisioned?
  - What are the performance bottlenecks?
  - What are the (unexpressed) requirements?
  - Managing the storage - scratch, precious, freeing up space, group storage
  - Scheduling bursty workflows & precious resources (GPUs, fast storage)
- A fifth category: metrics for **framework & platform developers**
  - Which data formats are physicists most often using and by which frameworks?
  - Are performance targets met? (e.g. X TB / Y minutes)
  - Where are the inefficiencies and user pain points?
  - What capabilities are missing?

[link](#)



3

[Analysis Facilities Monitoring Discussion,](#)

[Rob Gardner et al](#)

# Summary

# Summary and next steps

- Questions for experiments to be presented to LHCC by Alessandra in June 2024:
  - Comments and discussion yesterday greatly appreciated
  - Any major comments should be made today
  - Minor comments still welcome
- Plenary talks covered many aspects of AFs:
  - Discussion covered many new points, provided useful insight on current best-practice
  - Outcomes will prove valuable in reaching practical solutions to problems highlighted in AF White Paper
    - See [notes taken live](#)

Thank you, any questions?