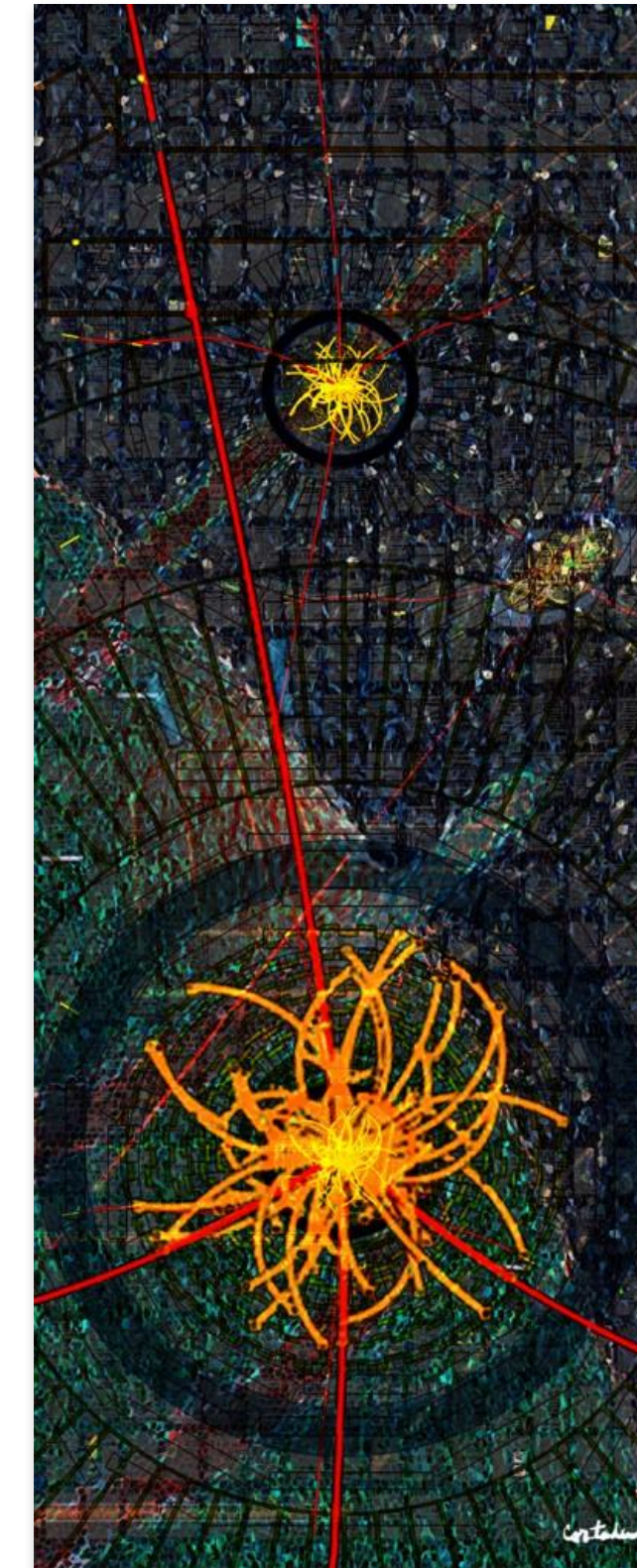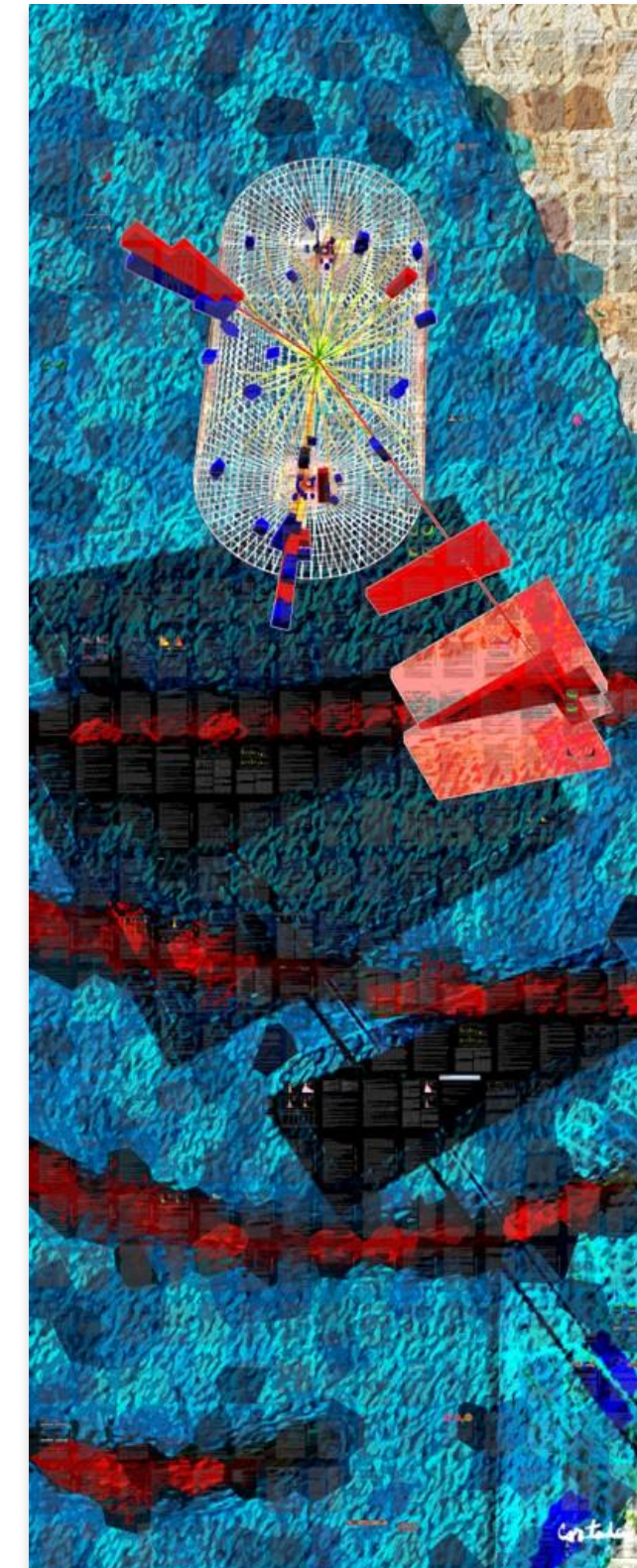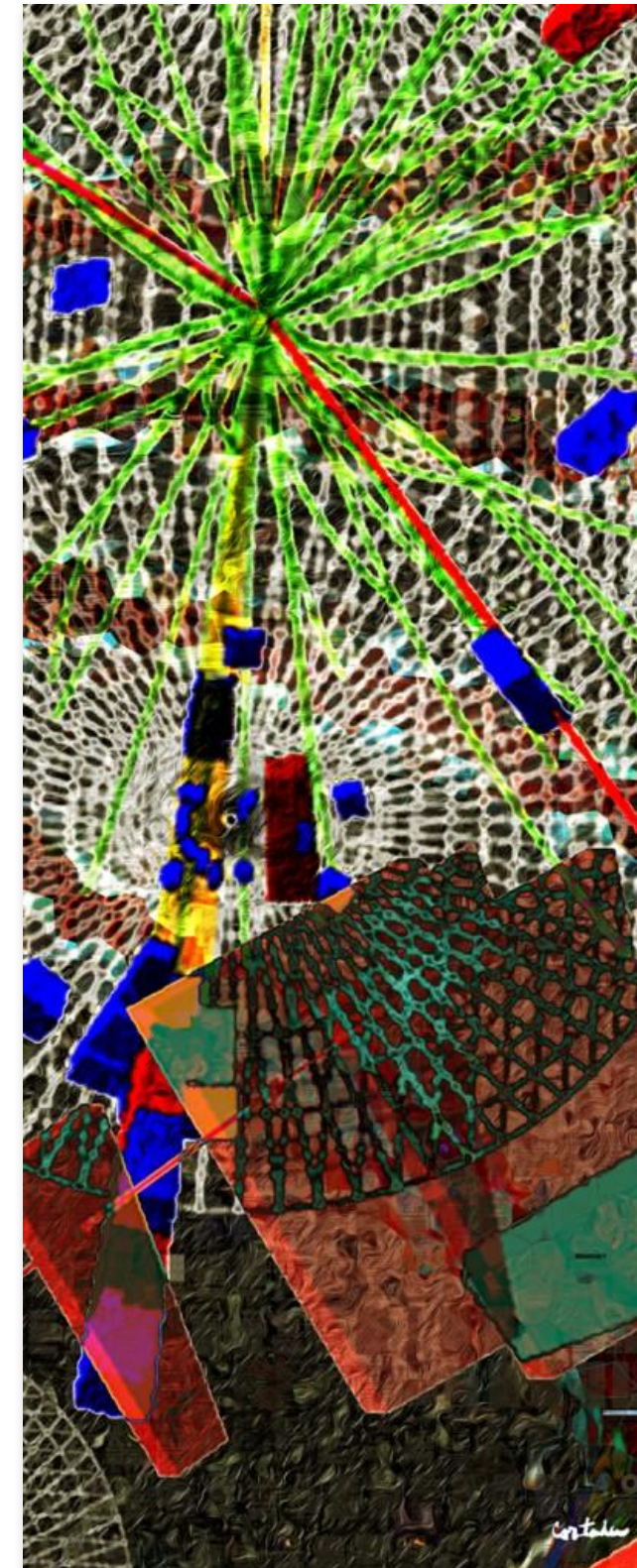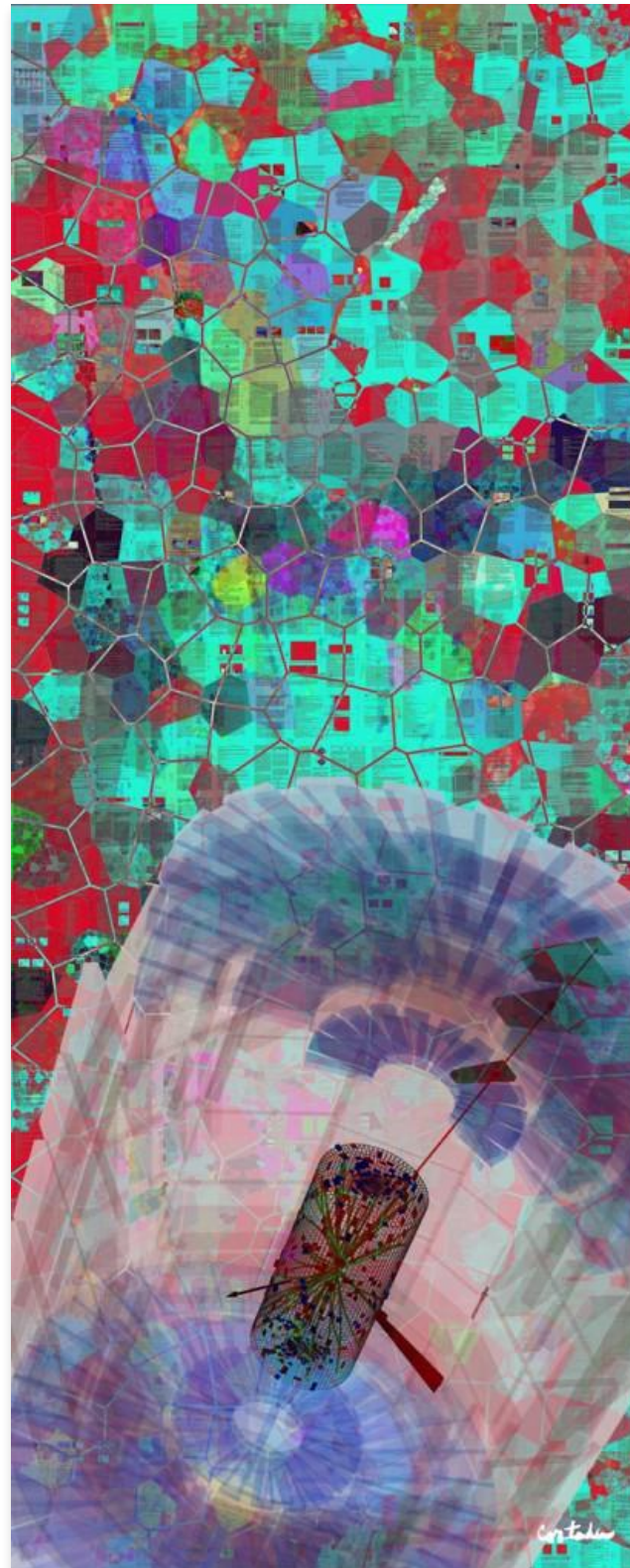# HSF/WLCG meeting debrief

Dirk Hufnagel (FNAL), with help from Andrew Melo for the note-taking
**CREST Meeting**
May 24, 2024

# HSF/WLCG Workshop When/Where/What

- Periodical WLCG Workshop traditionally done just before every CHEP, in recent years often combined with HSF

- This one done at DESY, Hamburg, there won't be one before CHEP in October

- Mix of plenary and parallel session (the latter often split into WLCG and HSF sessions)

- Caveat: The WLCG sessions were more relevant to me and while I tried to cover both (with Andrews help), we didn't attend / have notes for all the HSF sessions



**WLCG/HSF Workshop 2024**

13–17 May 2024
DESY
Europe/Zurich timezone

Enter your search term

Overview
Timetable
My Conference
    My Contributions
Community Software call for Abstracts
Photograph
Registration info
Visas & letters
Zoom Connections

About DESY
How to get to DESY
Onsite Information
Accommodation
Social Programme
Restaurant List
Craft Beers
Things to do in Hamburg
Running Route

Privacy Notice
Code of Conduct
Organising Committee

Contact
✉ wlcg-hsf-2024-contact@…

Welcome to the WLCG/HSF Workshop at DESY Hamburg, May 13-17 2024

Limited participation will be possible from remote, please register using the dedicated registration form if you wish to participate remotely. You will need to register to have access to the zoom connection(s).

**Starts** 13 May 2024, 12:00
**Ends** 17 May 2024, 13:00
Europe/Zurich

DESY
Deutsches Elektronen-Synchrotron (DESY)
Notkestraße 85
22607 Hamburg
Germany

# Relevant Links

- ## Workshop indico with agenda
  - https://indico.cern.ch/event/1369601/overview

- ## Google Doc with Minutes (14 pages, more details than this talk)
  - https://docs.google.com/document/d/1UHV3st8UE_sKxoToGvsT0eM7_iAvbUC4I-lkwi7bmrw/edit?usp=sharing

- ## Workshop internal recaps (in the Friday Plenary)
  - WLCG
    - https://indico.cern.ch/event/1369601/contributions/5908514/
  - HSF (not really a recap at all unfortunately, due to HSF sessions being more scattered in topic matters)
    - https://indico.cern.ch/event/1369601/contributions/5908515/
  - Analysis Facilities
    - https://indico.cern.ch/event/1369601/contributions/5908516/
  - DC24
    - https://indico.cern.ch/event/1369601/contributions/5908522/

# Introduction/Goals

- ## WLCG
  - Want to complete the strategic plan, nothing new here, just getting more urgent
    - What role (if any) does WLCG play in AF
    - Evolution of pledging (especially looking at special resources like HPC and Cloud)
    - GPU
  - DC24 Postmortem
  - Analysis Facilities
  - Operations and Facilities

- ## HSF
  - 10 year anniversary
  - Original plan was for HSF to develop and spin off software, that hasn't really happened
  - More of a think tank about software directions/requirements, others do the development
  - Lots of discussion on the human element, how to retain manpower for HEP software development, how to make sure there are career paths etc (also nothing truly new here)
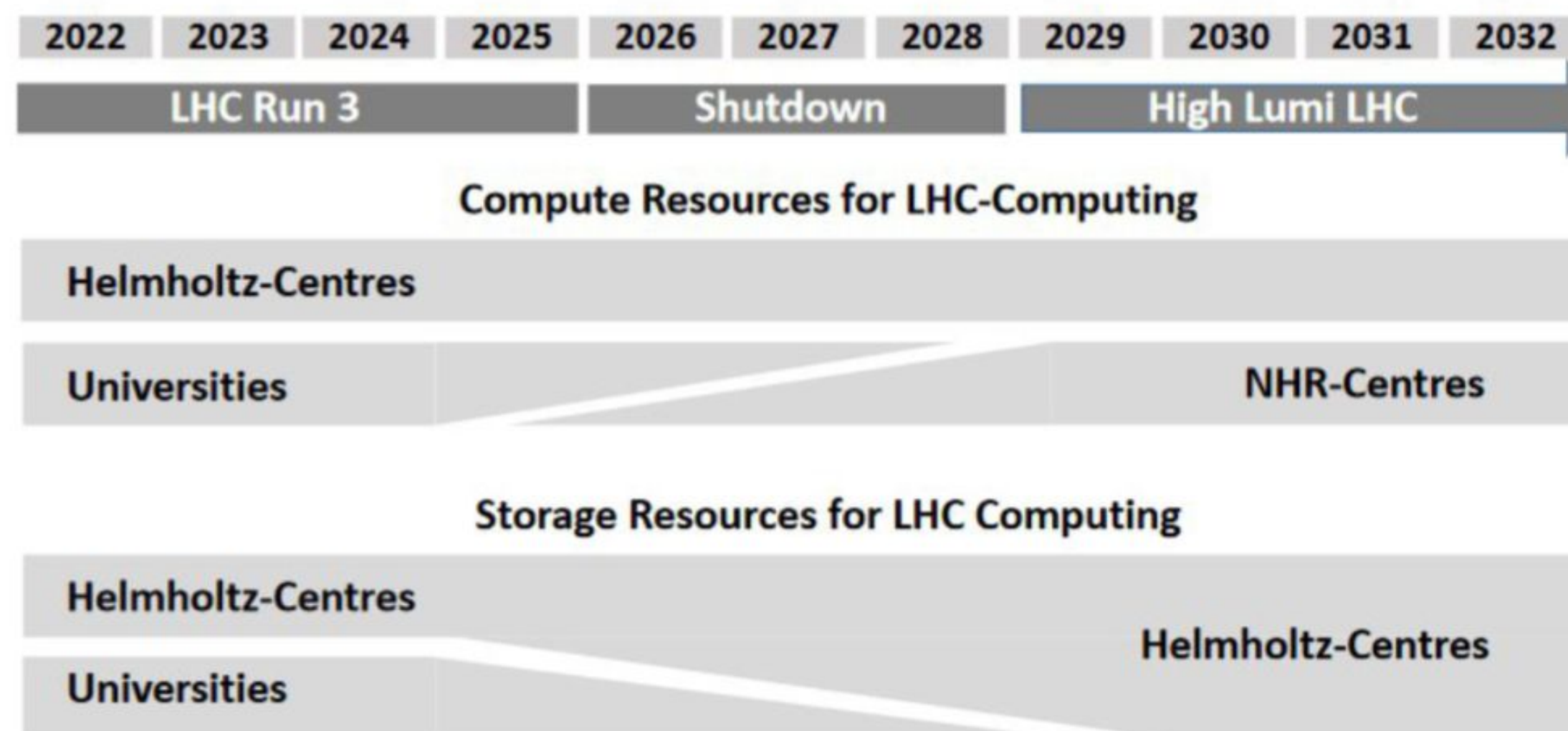
# German T2 Infrastructure

- Workshop was at Desy, so attendance from German institutes was high and there were quite a few talks from them. One interesting on was about the evolution of the German T2 infrastructure, which will be **fully** transitioned to HPC (for compute) and other shared infrastructure (for storage) within the next few years.

**Gradual transition** from university-based Tier-2 centres **to NHR (CPU) and Helmholtz-Centres (mass storage)** towards beginning of HL-LHC, i.e. 20% per year.

**Local ATLAS/CMS groups** keep supervising the NHR resources and apply for funding from federal government for **dedicated personnel**:

- ATLAS:
    - NHR@KIT - Freiburg group
    - NHR@Göttingen - Göttingen group
- CMS:
    - NHR@KIT - Karlsruhe group
    - NHR@Aachen - Aachen group

| 2022 | 2023 | 2024 | 2025 | 2026 | 2027 | 2028 | 2029 | 2030 | 2031 | 2032 |
|------|------|------|------|------|------|------|------|------|------|------|
| LHC Run 3 | | | | Shutdown | | | High Lumi LHC | | | |

**Compute Resources for LHC-Computing**

Helmholtz-Centres

Universities → NHR-Centres

**Storage Resources for LHC Computing**

Helmholtz-Centres

Universities → Helmholtz-Centres

*Strategy paper by KET from 2022:*
https://www.ketweb.de/sites/site_ketweb/content/e199639/e312771/KET-Computing-Strategie-HL-LHC-final.pdf

# Technology Evolution

- **Always a highly anticipated talk from Bernd Panzer, looking at technology trends in the market, under the aspect of what will be available to buy and at what price point**
  - [Slides](#)

- **Highlights**
  - CPU Evolution slowing down, can get EITHER more speed OR more power efficiency
    - Average performance gains over last 5 years still at ~20%, but slowing down over last 2 years
  - Cost of servers not going down right now
    - RAM prices going up, SSD prices going up (also mentioned copper prices going up)
  - Power and Cooling improvements slowing down (see comment on CPU)
  - New disk storage technology HAMR (Seagate) entering market
    - Not really available right now as single units, sold mostly as complete solutions / to hyper scalars
    - Competing WD MAMR technology even later and Toshiba even further behind
  - CERN is worrying about spinning disk IO performance, especially with disks getting larger (50TB)
    - Fear is that they might have to overprovision, just to get the needed IO rates
    - At which point SSD might become interesting
  - SSD not seen as replacing HDD anytime soon (price / manufacturing capacity reasons)
  - Tape media getting more expensive this year, 40%(!!!) for LTO-8, 15% for LTO-9
  - Lots of uncertainties (both cost evolution and experiment requirements, at least factor 2 combined)
    - CERN will try to save as much money as they can now and stockpile to account for this
    - Switching to Warranty + 2 year (7 years total)

# DC24 (Tue Morning Parallel)

- Didn't go myself, but Andrew went and took detailed notes, so if you want to know the details please look there, this only lists some highlights I picked out

- Talks from LHC (combined), DUNE and Belle-2
  - No real problems with network and storage, problems more on the "infrastructure" side (FTS/Rucio)
- Talks about European and US site perspective, respectively
- Talk from FTS
  - First slide: we were doing great. Later in first slide: probably didn't reach flex goals because of FTS.
    - Triggered some discussion apparently
- Talk from Rucio
- Talk from IAM (token provider)
  - Saw some overloads
- Talk about network
  - PerfSonar was very useful for debugging problems/misconfigurations

# HSF Community Software (Tue Morning Parallel)

- **ROOT Talk**
  - Status update
  - POW (Program of Work / Project Management)
    - Strong focus on reducing number of open issues
    - Want to report to experiments at least quarterly (stakeholders meetings)
  - Release schedule 2024: May (long terms support) and November (short term support)
  - ROOT ML philosophy: integrate with industry tools, don't compete with them
- **SciKit-HEP**
  - Python / pythonic physics analysis
  - Engagement with broader scientific python community
  - Really more a history / timeline talk, not much technical
- **Key4HEP**
  - Trying to become **the** software stack for future collider experiments
    - Already diverging, Key4HEP uses Gaudi as framework, one experiment already decided to change that
- **HSF-India**
  - 5 year program
  - Training, options for Indian researchers to visit US for some time (and vice versa)
  - Goal is built long term relationships and collaborations
- **Pythia8**
  - Status updates, lots of technical details
  - No dedicated funding for a few years now, worries about continual development of generators
- **Phoenix Event Display**
  - 2017 HSF white paper calling for common event format and common tools to visualize events
  - Used by ATLAS, FCC, LHCb, Belle-2

# WLCG Operations (Tue Afternon Parallel)

- **WLCG Infrastructure Evolution**
  - Dropping GridFTP in FTS support together with SL7 (no more GridFTP transfers!)
  - IPv6 also on all compute nodes by June 2026
    - Already clear not all sites will make that deadline
  - ARM survey: 10% of sites have some ARM resources, 33% want to by some, don't see large scale adoption (yet)
  - Retiring hardware that only supports x86_64-v1 (CMS already stopped supporting this in new CMSSW builds)

- **LHCOne Security and Jumbo Frames**
  - Adding more collaborations and sites to LHCOne, worry about diluting trust (which is basic principle behind LHCOne)
    - MULTIOne project: add tags that identify collaborations served by site (pro: simple, con: not 100% secure)
  - Conflicting reports on benefit of Jumbo Frames
    - CERN (and other large sites) will do tests (with Jumbo Frames all the way down to the compute node)

- **Token Transition Update**
  - VOMS-Admin to IAM transition happening in June (hard deadline of end of June due to SL7 EOL)
  - DC24: learned how to use and not use tokens
  - Auxiliary Services: htgettoken + HashiCorp Vault used at FNAL

# WLCG Operations (Wed Morning Plenary)

- **WLCG Helpdesk (GGUS)**
  - GGUS being retired, will run old/new in parallel for a while, old GGUS gone by end of year
    - Production ready in Oct, gradual migration of SU in Oct and Nov
    - No copying of tickets!
  - New system based on Zammad (open source)
  - REST API based (Web UI just a layer on top)
  - Concerns from audience about transition period and tickets that are open a long time
    - Might have to manually copy some tickets if needed in worst case
- **Handling of high-memory jobs. Whole node scheduling. Standard 16-core slots instead of 8-core slots. VO and site perspective.**
  - ATLAS: suggest 16-core slots, can't easily use whole nodes
  - CMS: can use whatever, up to whole nodes
  - ALICE: similar to CMS
  - Sites: varied responses, but some worries about 16-core slots or whole nodes
    - Suggestions is to resurrect the WLCG "Multicore TF" to figure this out

- ROOT RNTuple
  - Status update and plans
  - Default compression ZStd
  - Seamless integration (TTree and RNTuple interchangeable for reads)
- Baler: ML data compression in real time
  - Open CMS Data: 58% compression
  - CFD: Output is 0.5% of input
  - Offline use case: train on actual data we want to compress
  - Online use case: train on reference data, then use for other live data
- Analysis Grand Challange benchmarking
  - AGC developed by IRIS-HEP to test feature completeness of SciKit-HEP toolkit
  - Showed results for various scenarios for T2 Munich analysis (also with/without xcache)
- Cloud Data Lake Technologies
  - Kubernetes, Parquet File Format, Object store, Apache Iceberg, Nessie, Trino Distributed SQL
  - Proof of concept with CMS Data

# WLCG Strategy (Wed Afternoon Plenary)

- **Technical Evolution**
  - 15+ years of continuous operations, not unreasonable to rethink how we do things
  - Problems to tackle
    - Keep pace with changed in the technical landscape (we aren't really cutting edge anymore)
    - Increased scrutiny and accountability (growing competition for funding)
    - Maintain relevance and engagement (in a post-Higgs world)
  - Grid Deployment Board => Technical Evolution Board
    - Not just a name change, some tentative ideas about what else to change
    - Align aspirations with capabilities (don't fragment too much)
    - How does WLCG relate to existing activities (and HEPIX)
  - LOTS of discussion
    - Many folks at the workshop won't be active during Run5, need to bring in new personnel
    - How to involve non-WLCG communities on technical discussions that overlap
- **Environmental Sustainability**
  - Strategic importance to many federations and funding agencies
  - Agree on metrics and framework to collect energy efficiency information
  - Enable use of more energy efficient hardware
  - Develop and promote and sustainability plan

# WLCG Strategy (Thu Morning Plenary)

- **Infrastructure**
  - ○ Role of WLCG, distinction between T0, T1, T2
    - ■ Actual capabilities can be quite different from what the label implies
    - ■ Goal: somehow account for this, to record actual capabilities and to give credit
      - ● And eventually doing something with this information (like job routing etc)
  - ○ Cloud integration
    - ■ Clouds are in use, strategic topic for experiments (not just commercial, but academic/research/national clouds)
    - ■ Goal: document technical solutions (internal)
    - ■ Goal: follow progress of other communities (external)
    - ■ Goal: define commercial cloud storage provisioning policies
      - ● Software infrastructure (Rucio/FTS) should support clouds (independent of experiment choices)
      - ● Cost comparison should include expertise/services you get (and not needing it in-house anymore)
      - ● ATLAS Cloud TCO should be public now ! (haven't looked myself yet)
      - ● Storage policy question not really cloud specific, more for all external storage providers
  - ○ HPC Integration
    - ■ HEP doesn't need HPC, but we can use it
    - ■ Elephant in room: might have "no other option", depends on FA
    - ■ Goal: document technical solutions (internal)
    - ■ Goal: construct dialogue with relevant global bodies to drive future allocation policies
    - ■ Discussions:
      - ● EU: need to talk to agencies for HEP to be considered a use case
      - ● UK: trying to talk to FA to influence HPC allocation would backfire
      - ● Mentioned the HPC Strategy Blueprint process (had US meeting with WLCG/USATLAS/USCMS in New York, organizing a EU workshop)
  - ○ Heterogeneous Grid Infrastructure
    - ■ Focus on software aspects: technology tracking, benchmark and accounting, Rest is on the experiments (WLCG doesn't write software).
  - ○ Evolution of services and standards
    - ■ Rely on lots of services and standards that are mostly developed in-house
    - ■ In many case non-HEP standard alternatives exist now (with larger use cases and better support)
    - ■ OTOH moving is costly and has risks too (external tools stop being supported)
    - ■ Goal: establish process for adopting non-HEP standards and services (and retiring legacy ones). Should include risk management for external dependencies.
    - ■ Discussion:
      - ● Not really a WLCG decision, up to experiments
      - ● Security cannot be an afterthought in this
      - ● Spirit of this is not to setup a new board, but to allow new things to be tested

# WLCG Strategy (Thu Morning Plenary)

- **Financial Sustainability**
  - Flat cash assumption for many years
    - Also covers many things not actually flat (infrastructure, energy)
    - Don't see this changing
  - Draft actions
    - Keep monitoring hardware and market trends
    - Multi-year resource planning
      - Some FA might allow to shape the budget over longer term
      - Already happening to some extent within experiments
    - Monitor pledges vs authorship
      - Lots of negative feedback: up to experiments, none of WLCG's business (even if the information is in principle public)
    - Fine grained pledges
      - Quartely pledges exist, have never been used (now hidden, but still available)
      - Lots of discussion
        - Experiments decide what is and isn't useful for them
        - To first order pledging over longer time period and averaging is fine (and already done)
          - Some disagreements on this but it's probably ok as long as the utilization reports later match
        - Pledges can't be 100% dynamic though, some need for a base that's always there
        - Currently comfortable over pledge, might not always be the case, would require better planning if that happens
        - Suggestion: allow a better way to pledge time-integrate resources (and not just averaging them out over the pledge time window and pledging that as always available)
  - Other federation commitments
    - Find ways to recognize commitments from federations to middleware development
    - Goal: commitment made more formal
    - Discussions:
      - If this recognition can trigger funding, have to be careful about what we want to recognize and we really want/need

- **LHCC Charge Discussion**
  - WLCG LHCC Referees wants to review the experiments AF plans and activities, so they asked the experiments to come up with a list of (common between experiments) that the LHCC will then ask the experiments in a yearly review. In principle these questions can be refined year by year, in this session we went through the initial list of questions that will be sent to the LHCC (and which they are free to modify as they see fit).
    - Just because the questions are the same, the answers from each experiment will differ
    - Before this workshop question suggestions were collected from each experiment and then they were aggregated and compiled into a suggested list of common questions, these were discussed during the workshop. Not going into this here, please see minutes for the details.
    - Final list of questions (cleaned up a bit from the draft used for discussions at the workshop are [here](#)

## 200 Gb/s analysis demonstrator

HSF

- 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*

Jamie Gooding | Analysis Facilities: Summary and next steps | HSF/WLCG Workshop 2024 | 17th May 2024 | 9

# Analysis Facilities (Thu Afternoon Plenary)



## 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***

Jamie Gooding    Analysis Facilities: Summary and next steps    HSF/WLCG Workshop 2024    17th May 2024    10
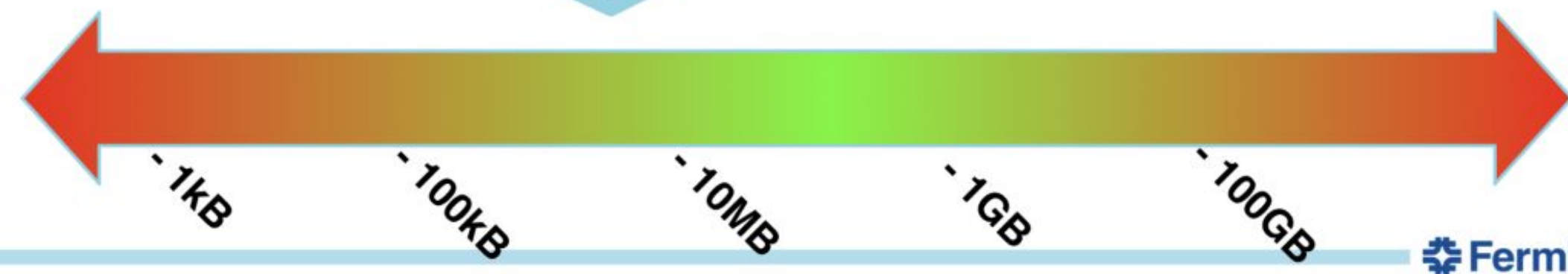
## Data access discussion

- 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

~1kB    ~100kB    ~10MB    ~1GB    ~100GB

✦ Fermilab

12    May 16, 2024    Storage for Analysis Facilities

**Storage for Analysis Facilities, *Dirk Hufnagel & Nick Smith***

Jamie Gooding    Analysis Facilities: Summary and next steps    HSF/WLCG Workshop 2024    17th May 2024    11

# 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?

link

- Users (experience) ①
  - 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) ② ③ ④
  - 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 fith 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?

submit your own questions

3

Analysis Facilities Monitoring Discussion, *Rob Gardner et al*