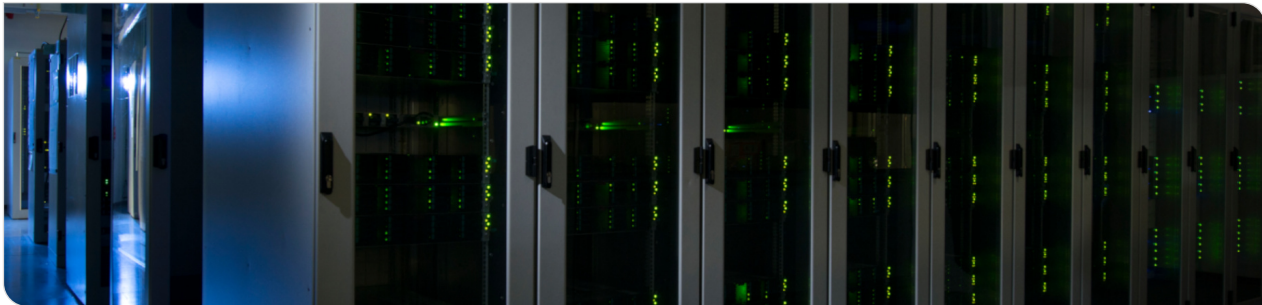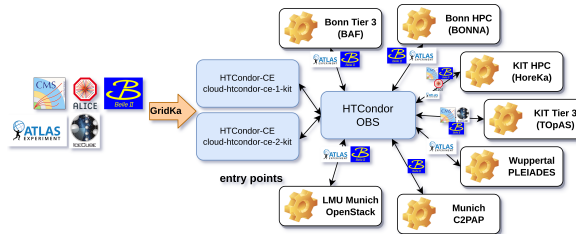# AUDITOR:
# An Accounting tool for Grid Sites and Opportunistic Resources

Matthias J. Schnepf on behalf of the AUDITOR Group | 4. November 2024
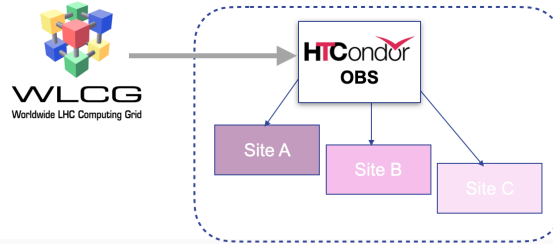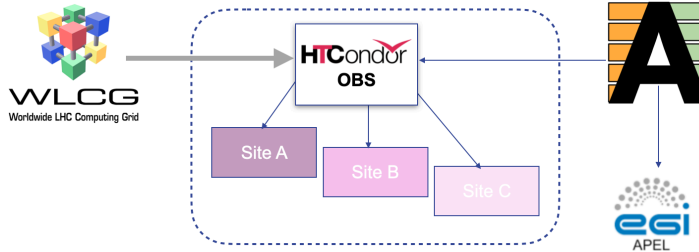
# Grid Subsites



- Grid Site provide CE and storage
- several small/non-Grid like resource provider
- CE provides Grid-like access to several computing resources
- all resources are in one resource pool/**O**verlay **B**atch **S**ystem (OBS)

# Accounting opportunistic resources



- COBalD/TARDIS allows multiple resources to be clustered in an **O**verlay **B**atch **S**ystem
  - Sub clusters **cannot be accounted individually** with existing tools
  - Requires a dedicated mechanism for accounting
- Challenges
  - Vastly different infrastructures
  - Many potential use cases
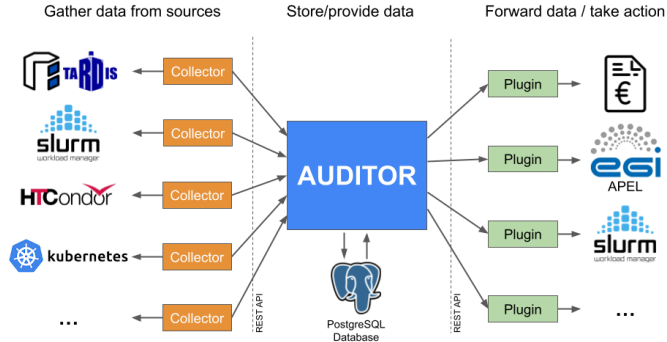
# Accounting opportunistic resources



- COBalD/TARDIS allows multiple resources to be clustered in an **O**verlay **B**atch **S**ystem
  - Sub clusters **cannot be accounted individually** with existing tools
  - Requires a dedicated mechanism for accounting
- Challenges
  - Vastly different infrastructures
  - Many potential use cases
- **AUDITOR** provides multi-purpose accounting ecosystem

# AUDITOR: Modular Accounting Ecosystem



Gather data from sources | Store/provide data | Forward data / take action

**AUDITOR: Acco**U**nting** D**ata handl**I**ng** T**oolbox for** O**pportunistic** R**esources**

- **Collectors**
  - Accumulate data
- **Core component**
  - Accept data
  - Store data
  - Provide data
- **Plugins**
  - Take action based on stored data

**Documentation and code**

→ https://github.com/ALU-Schumacher/AUDITOR

# AUDITOR: Core component



- Implemented in **Rust**
  - Access via REST interface
- Unit of accountable resources: **Record**
- Data stored in PostgreSQL
- Completely stateless
  - No dataloss
  - Suitable for high availability setups
- Provided as **RPM** or **Docker container**
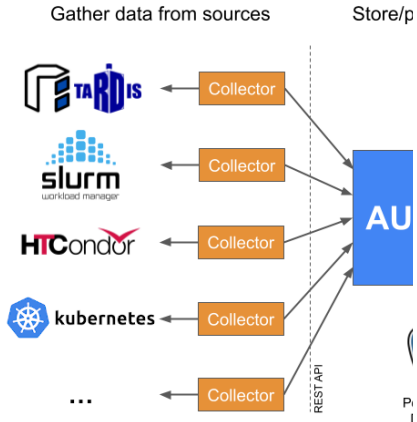- Client libraries in Rust and Python

# Record in AUDITOR

- record_id: uniquely identifies the record
- meta: multiple key value pairs of the form `String -> [String]`
- components: arbitrary number of resources that are to be accounted for (CPU, RAM, Disk, GPU, …)
  - scores: (multiple) accounting scores supported
- start_time, end_time: datetime in UTC
- runtime: calculated as `end_time - start_time`

→ meta & component fields allow for maximal flexibility

```json
{
    "record_id": "hpc-4126142",
    "meta": {
        "group_id": [ "atlpr" ],
        "site_id": [ "hpc" ],
        "user_id": [ "atlpr001" ]
    },
    "components": [
        {
            "name": "Cores",
            "amount": 8,
            "scores": [
                {
                    "name": "HEPSPEC06",
                    "value": 10.0
                },
                {
                    "name": "HEPScore23",
                    "value": 10.0
                }
            ]
        },
        {
            "name": "Memory",
            "amount": 16000,
            "scores": []
        }
    ],
    "start_time": "2023-02-24T00:27:58Z",
    "stop_time": "2023-02-24T03:41:35Z",
    "runtime": 11617
},
```
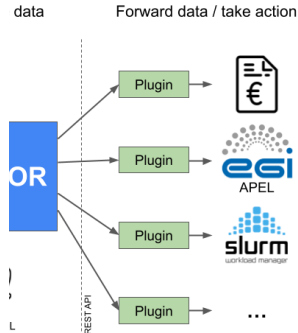
# Collers: Accumulate data

Gather data from sources    Store/p



- collect metadata about used resources
  - walltime
  - CPU usage
  - weight/score (e.g. HEPScore23 or $CO_2$ equivalent)
  - ...
- provided collectors
  - **TARDIS Collector** (developed @ Freiburg)
  - **SLURM Collectors** (developed @ Freiburg)
  - **HTCondor Collector** (developed @ KIT)
  - **Kubernetes Collector** (developed @ Uni Wuppertal)
- You miss a collector? Feel free to join

# Plugins

- **Priority plugin**
  - Compute priorities from a list of records
  - Update priorities on a batch cluster
- **APEL accounting plugin**
  - Properly accounts individual sites or subsites (e.g. behind COBalD/TARDIS)
  - Reports accounting data to the APEL accounting platform
- **Utilization report** (future project)
  - Analyse requested vs. consumed resources of a user
  - Send a weekly report with possible savings and $CO_2$ footprint

# AUDITOR

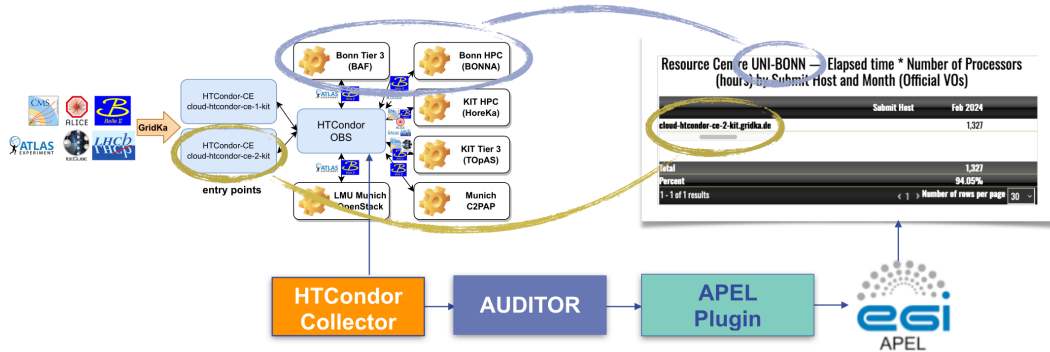- Extensive documentation



Some numbers:

- 8 contributors
- from 3 universities
    - Freiburg (main developement), KIT, Uni Wuppertal
- 16 releases - latest v0.6.3
- Continuous improvements: Commits
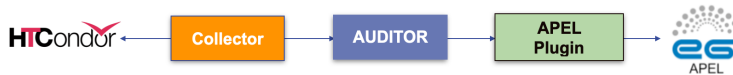




https://doi.org/10.21203/rs.3.rs-4741479/v1

# WLCG Subsite Accounting



- Grid infrastructure hosted and maintained in Karlsruhe, resources provided by Bonn
- **AUDITOR** accounting pipeline allows to account for sub-clusters individually

# WLCG Grid Site Accounting

- KIT replaced accounting of the APEL client by AUDITOR pipeline in May 2024



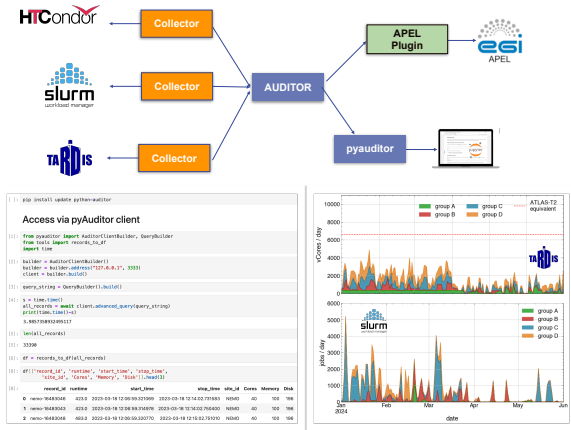Infrastructure High Throughput Compute Accounting

The EGI Accounting Portal is an EGI service provided by CESGA

This work is co-funded by the EOSC-hub project (Horizon 2020) under Grant number 777536.

Resource Centre FZK-LCG2 — Elapsed time * Number of Processors (hours) by Submit Host and Month (Official VOs)

| Submit Host | Apr 2024 | May 2024 | Jun 2024 | Jul 2024 | Aug 2024 | Total | Percent |
|---|---|---|---|---|---|---|---|
| htcondor-ce-1-kit.gridka.de:9619/htcondor-ce-1-kit.gridka.de-condor | 10,201,871 | 11,568,000 | 1,392,046 | 0 | 0 | 23,161,918 | 13.72% |
| htcondor-ce-2-kit.gridka.de:9619/htcondor-ce-2-kit.gridka.de-condor | 8,276,691 | 6,245,383 | 0 | 0 | 0 | 14,524,074 | 8.6% |
| htcondor-ce-3-kit.gridka.de:9619/htcondor-ce-3-kit.gridka.de-condor | 7,815,697 | 4,672,054 | 0 | 0 | 0 | 12,487,752 | 7.39% |
| htcondor-ce-4-kit.gridka.de:9619/htcondor-ce-4-kit.gridka.de-condor | 7,477,659 | 708,594 | 0 | 0 | 0 | 8,186,253 | 4.85% |
| cloud-htcondor-ce-2-kit.gridka.de | 216,816 | 179,685 | 0 | 0 | 0 | 396,501 | 0.23% |
| pps-htcondor-ce.gridka.de:9619/pps-htcondor-ce.gridka.de-condor | 640 | 0 | 0 | 0 | 0 | 640 | 0% |
| htcondor-ce-1-kit.gridka.de | 0 | 0 | 3,083,567 | 9,665,501 | 10,043,489 | 22,792,557 | 13.5% |
| htcondor-ce-2-kit.gridka.de | 0 | 1,920,223 | 9,015,181 | 8,398,234 | 8,866,979 | 28,200,616 | 16.7% |
| htcondor-ce-3-kit.gridka.de | 0 | 2,340,311 | 9,307,988 | 8,008,737 | 9,126,696 | 28,783,731 | 17.04% |
| htcondor-ce-4-kit.gridka.de | 0 | 7,991,571 | 5,739,147 | 7,842,016 | 8,768,026 | 30,340,760 | 17.97% |
| pps-htcondor-ce.gridka.de | 0 | 438 | 1,639 | 1,619 | 1,050 | 4,745 | 0% |
| **Total** | **33,991,374** | **35,626,760** | **28,539,568** | **33,916,107** | **36,806,240** | **168,879,549** | |
| **Percent** | **20.13%** | **21.10%** | **16.90%** | **20.08%** | **21.79%** | | |
| 1 - 11 of 11 results | | | | | | ‹ 1 › Number of rows per page 30 ⌄ | |

- AUDITOR is able to provide the accounting of the DE-Tier 1
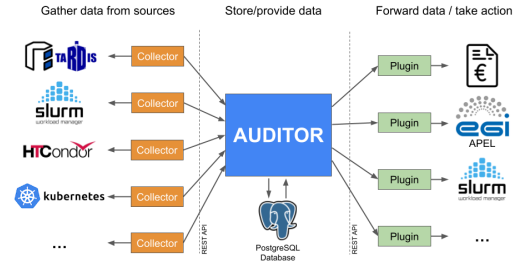  - largest WLCG Tier-1 that supports allfour WLCG experiments

# Collecting Accounting Info with AUDITOR



- accounting data can be collected in one or more AUDITOR instances from multiple sources
- APEL plugin can report for one or more queues
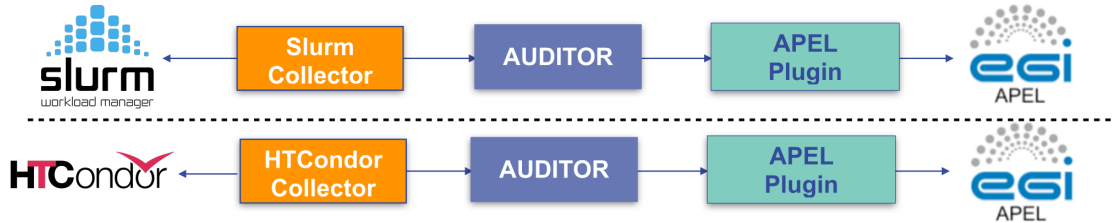- pyauditor allows to integrate AUDITOR client into python env

# Conclusion

- Provides an accounting ecosystem for various use cases
- Allows to account for different resources shared by one overlay batch system
- Allows to collect accounting data from multiple sources
- provision via containers independent of the OS
- Flexible structure of records and ecosystem allows to quickly adapt to future use cases
- Website: https://alu-schumacher.github.io/AUDITOR/
- GitHub: https://github.com/ALU-Schumacher/AUDITOR/
- FIDIUM: https://fidium.erumdatahub.de
- Email: auditor@physik.uni-freiburg.de

# Backup

# WLCG Accounting Use Case



- Collect accounting data from SLURM or HTCondor
- Store data as records in AUDITOR DB
- APEL plugin retrieves records from AUDITOR
  - creates APEL job summary from records
  - sends summary to defined APEL server
- Sites planing to use AUDITOR for accounting:
  - DESY-HH, Uni Wuppertal, ... ← ATLAS DE T1 (GridKa) moved reporting to AUDITOR

# APEL Plugin

```
log_level: INFO
time_json_path: /etc/auditor_apel_plugin/time.json
report_interval: 86400

site:
  publish_since: 2023-01-01 13:37:42+00:00
  sites_to_report:
    SITE_A: ["site_id_1", "site_id_2"]
    SITE_B: ["site_id_3"]

  benchmark_type: hepscore23

auditor:

  benchmark_name: hepscore23
  cores_name: Cores
  cpu_time_name: TotalCPU
  cpu_time_unit: milliseconds
  nnodes_name: NNodes
  meta_key_site: site_id
  meta_key_submithost: headnode
  meta_key_voms: voms
  meta_key_username: username
```

- **block 1**: configure serivce
  - file to store current state
  - time in seconds between reports

- **block 2**: configure site(s) to be reported
  - sites_to_report:
    keys: names of the sites in the GOCDB,
    values: corresponding site names in AUDITOR records

- **block 3**: configure metrics to be reported
  - meta_key_voms:
    key in meta field to be used as voms

https://alu-schumacher.github.io/AUDITOR//v0.5.0/#apel-plugin
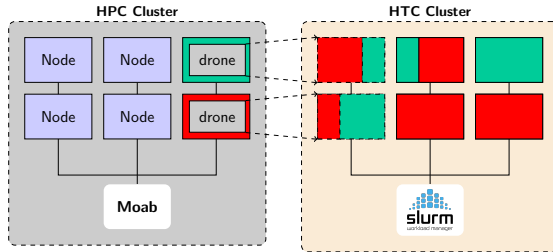
# APEL Plugin

```
optional:
  GlobalUserName: !MetaField
    name: subject
    datatype_in_message: TEXT
  VO: !MetaField
    name: voms
    datatype_in_message: TEXT
    regex: (?<=%2F).*?\S(?=%2F)
  VOGroup: !MetaField
    name: voms
    datatype_in_message: TEXT
    regex: (?=%2F).*?\S(?=%2F)
  VORole: !MetaField
    name: voms
    datatype_in_message: TEXT
    regex: (?=Role).*
  SubmitHost: !MetaField
    name: headnode
    datatype_in_message: TEXT
```
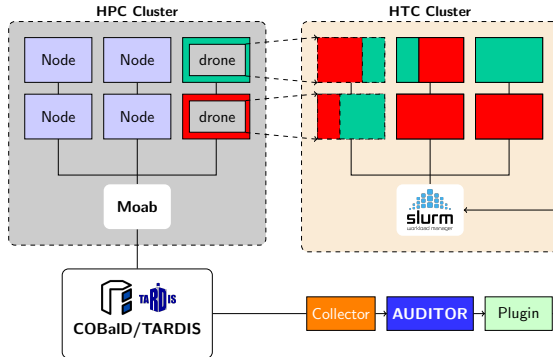
- dynamic mapping of any MetaField via regex
  - this allows to report accounting data for different VOs submitted with **tokens**
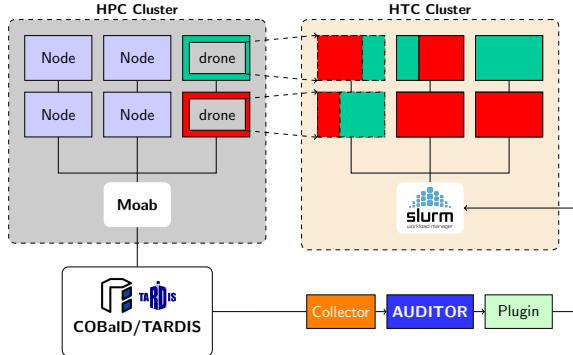- plugin configuration a bit more complicated, but much more flexible
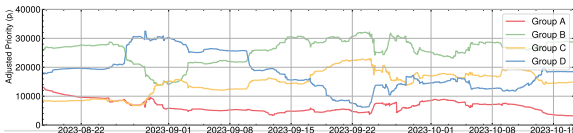
## Priority Use Case



- HPC resources integrated with COBalD/TARDIS
- Several HEP groups provide HPC resources
- Resources shared among HEP groups
- How to guarantee fair share on HTC cluster?

# Priority Use Case



- HPC resources integrated with COBalD/TARDIS
- Several HEP groups provide HPC resources
- Resources shared among HEP groups
- How to guarantee fair share on HTC cluster?

- **TARDIS collector** retrieves info of provided resources on the NEMO cluster
- **AUDITOR** accounts for provided resources of individual groups [**A** and **B**]
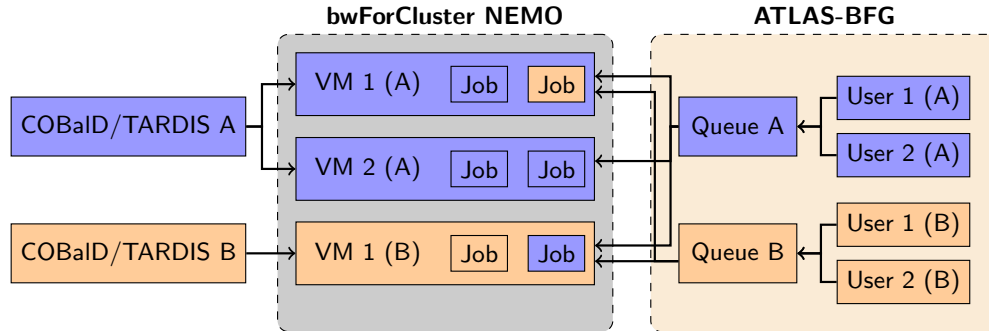- **Priority plugin** adjusts priorities on HTC cluster

# Priority Use Case

- HPC resources integrated with COBalD/TARDIS
- Several HEP groups provide HPC resources
- Resources shared among HEP groups
- How to guarantee fair share on HTC cluster?

- **TARDIS collector** retrieves info of provided resources on the NEMO cluster
- **AUDITOR** accounts for provided resources of individual groups [**A** and **B**]
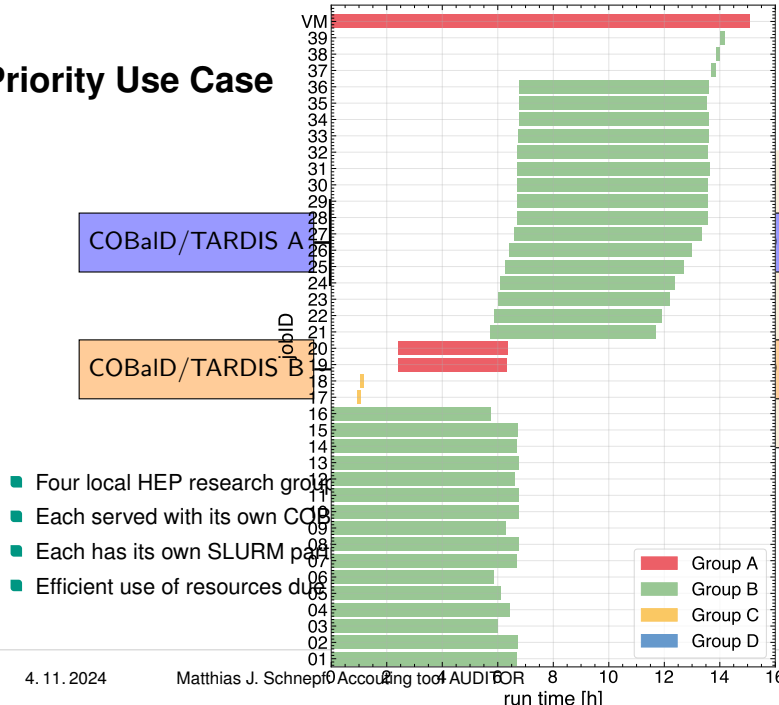- **Priority plugin** adjusts priorities on HTC cluster

# Priority Use Case



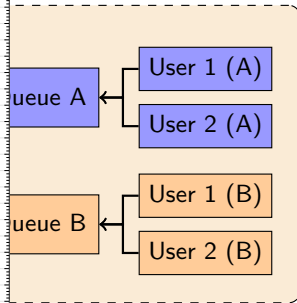**bwForCluster NEMO**     **ATLAS-BFG**

- Four local HEP research groups (A to D) with a share in NEMO
- Each served with its own COBalD/TARDIS instance
- Each has its own SLURM partition (job queue)
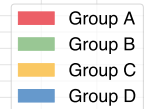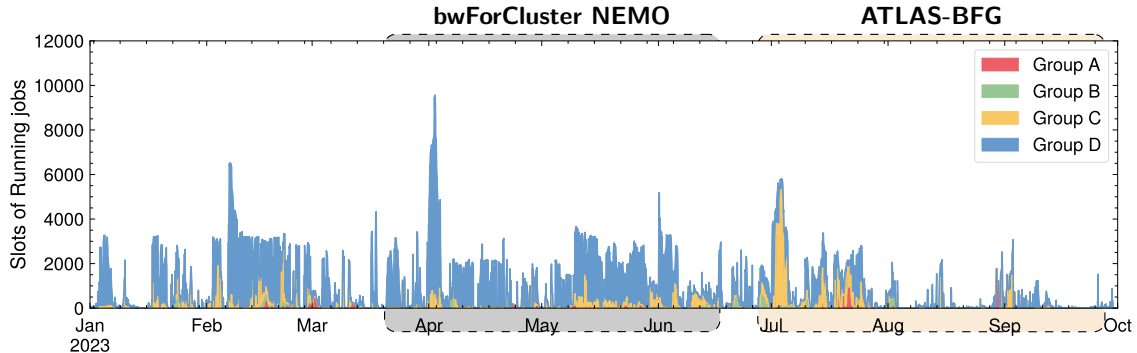- Efficient use of resources due to sharing VMs across HEP groups

# Priority Use Case

COBalD/TARDIS A

COBalD/TARDIS B

ATLAS-BFG

- Four local HEP research group
- Each served with its own COB
- Each has its own SLURM part
- Efficient use of resources due

Queue A — User 1 (A), User 2 (A)

Queue B — User 1 (B), User 2 (B)

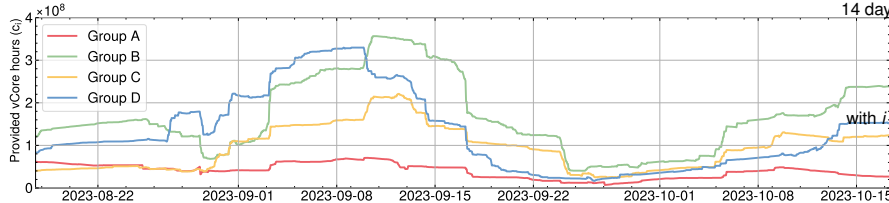Group A
Group B
Group C
Group D

jobID

run time [h]

# Priority Use Case



- Four local HEP research groups (A to D) with a share in NEMO
- Each served with its own COBalD/TARDIS instance
- Each has its own SLURM partition (job queue)
- Efficient use of resources due to sharing VMs across HEP groups

# Priority Use Case

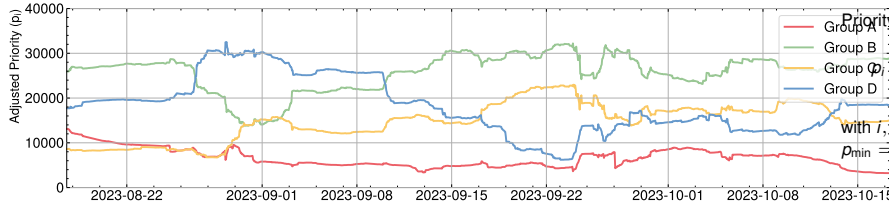- Provided resources of the four local HEP groups

Integral of provided vCore hours over last 14 days for each group:



$$c_i = \int_{t_{\text{now}} - 14\,\text{d}}^{t_{\text{now}}} N_i(t)\,\mathrm{d}t$$

with $i \in A, B, C, D$

- Priority is adjusted according to the provided resources



Priority $p_i$ is defined as:

$$p_i = \frac{c_i}{\sum_j c_j} (p_{\max} - p_{\min}) + p_{\min}$$

with $i, j \in A, B, C, D$;
$p_{\min} = 1$; $p_{\max} = 65535$