universität freiburg

# Accounting with AUDITOR

## WLCG Operations Coordination Meeting

Michael Böhler
October 3rd 2024

# AUDITOR
???

► What is AUDITOR?

► Which technologies are used?

► Who is AUDITOR?

► What can we do with AUDITOR?

► What are existing use cases?

  ► AUDITOR stands for **A**cco**u**nting **D**ata Handl**i**ng **T**oolbox For **O**pportunistic **R**esources

► What is AUDITOR?

► Which technologies are used?

► Who is AUDITOR?

► What can we do with AUDITOR?

► What are existing use cases?



**Internal Auditor (IA)**

*[in-ˈtər-nᵉl ˈò-də-tər]*

A trained professional employed by companies to provide independent and objective evaluations of financial and operational business activities, including corporate governance.

Investopedia

► AUDITOR stands for **A**ccounting **D**ata Handling **T**oolbox For **O**pportunistic **R**esources
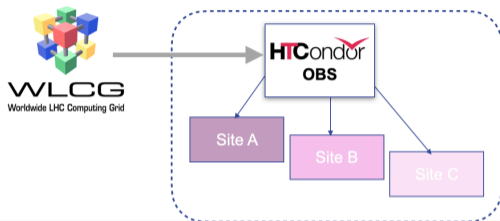
- ▶ What is AUDITOR?
- ▶ Which technologies are used?
- ▶ Who is AUDITOR?
- ▶ What can we do with AUDITOR?
- ▶ What are existing use cases?

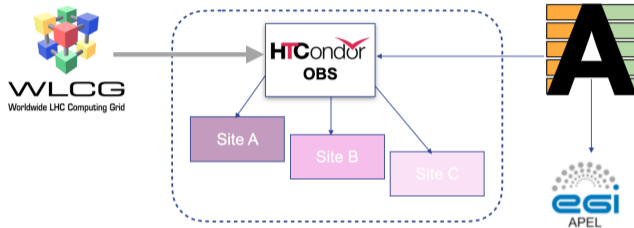- ▶ AUDITOR stands for **A**ccou**u**nting **D**ata Hand**l**ing **T**oolbox For **O**pportunistic **R**esources

- ▶ 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
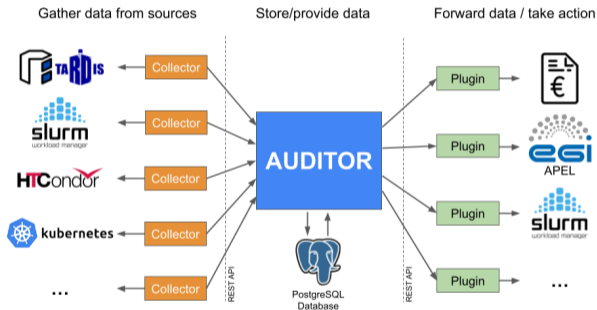
## Original Motivation
Accounting opportunistic resources



- ▶ 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
## Accounting Ecosystem



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

**Modular accounting ecosystem**

► **Collectors**
  ► Accumulate data

► **Core component**
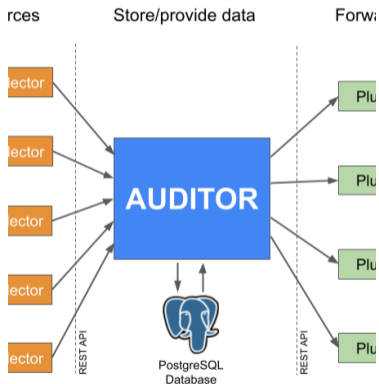  ► Accept data
  ► Store data
  ► Provide data

► **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**
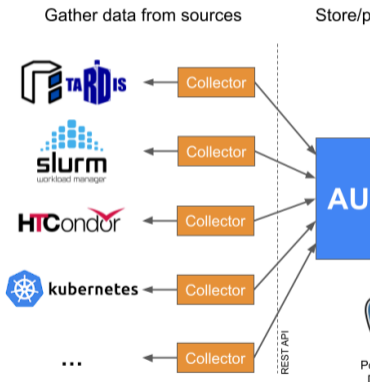- ▶ Client libraries in Rust and Python

# Record
Unit of accountable resources

- ▶ `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`

```
{
    "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
},
```

Accumulate data



Gather data from sources   Store/p...

- **TARDIS Collector**
  - Collect drone information
- **SLURM Collectors**
  - Collect information about SLURM jobs via SLURM CLI commands
- **HTCondor Collector** (developed @ KIT)
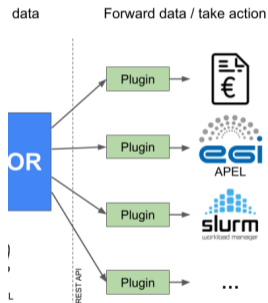  - Equivalent of SLURM collector for HTCondor
- **Kubernetes Collector** (developed @ Uni Wuppertal)
  - Collects information from Kubernetes pods

## Plugins

Take action based on stored data



data — Forward data / take action

- **Priority plugin**
  - Compute priorities from a list of records
  - Update priorities on a batch cluster
- **APEL accounting plugin**
  - Properly accounts individual sites 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
Contributors and Documentation

▶ Extensive documentation



Overview of the AUDITOR ecosystem. AUDITOR accepts records from collectors, stores them in a PostgreSQL database and offers these records to plugins which take an action based on the records.

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

▶ 8 contributors
▶ from 3 universities
  ▶ Freiburg (main developement), KIT, Uni Wuppertal
▶ 15 releases - latest v0.6.2
▶ Continuous improvements: Commits



▶ Also on https://zenodo.org/records/13239266

# WLCG Accounting Use Case

First working prototype



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

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



▶ AUDITOR is able to provide the accounting of the DE-Tier 1
  ▶ one of the largest WLCG Tier-1s

# 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

AUDITOR

- ▶ 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

# References



AUDITOR

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

Michael Boehler
Albert-Ludwigs-Universität Freiburg
michael.boehler@physik.uni-freiburg.de

Back-Up...

# WLCG Accounting Use Case
APEL Accounting with AUDITOR



- ▶ 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

Configuration

```
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: subject
```

- ▶ `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
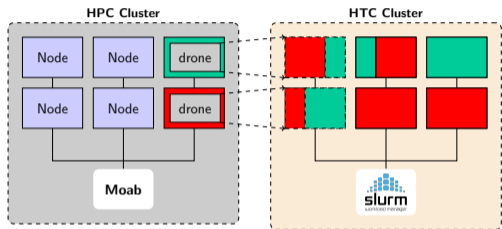
# APEL Plugin

Configuration in Release v0.6.2

```
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

Adapting priorities on HTC cluster based on provided HPC resources



- ▶ 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

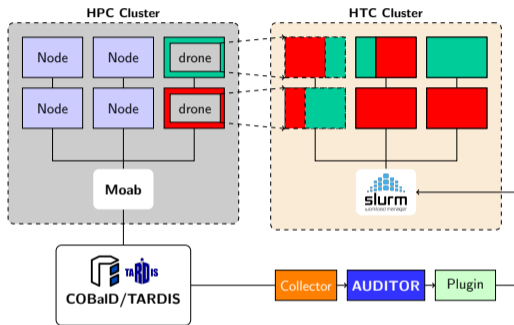universität freiburg
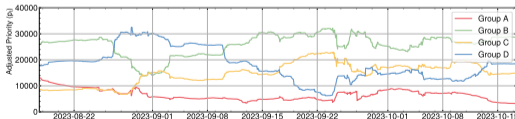
# Priority Use Case

Adapting priorities on HTC cluster based on provided HPC resources



- ▶ 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
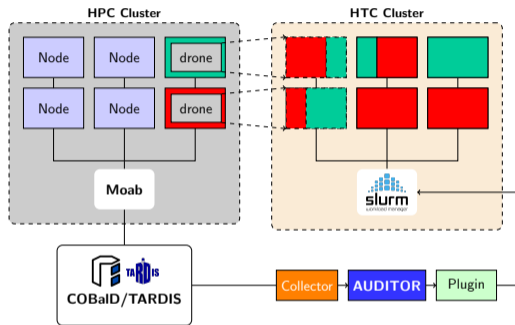
## Priority Use Case

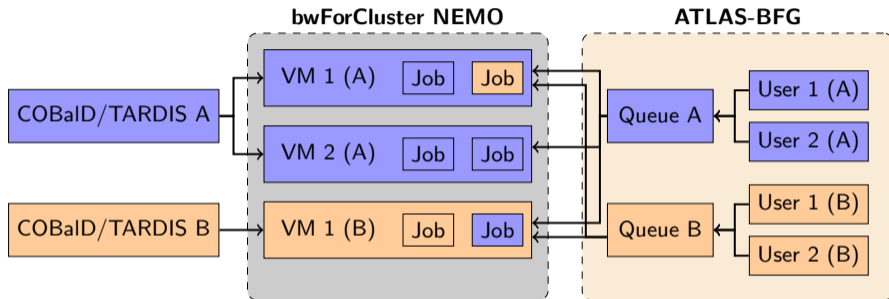Adapting priorities on HTC cluster based on provided HPC resources



- ▶ 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

- ▶ 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
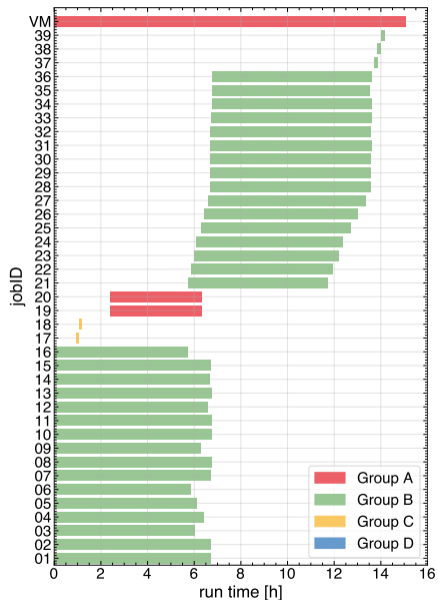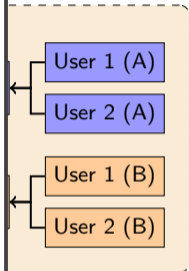
# Priority Use Case

Motivation



COBalD/TARDIS

COBalD/TARDIS

**LAS-BFG**

User 1 (A)

User 2 (A)

User 1 (B)

User 2 (B)

- ▶ Four local HEP rese...
- ▶ Each served with its...
- ▶ Each has its own SL...
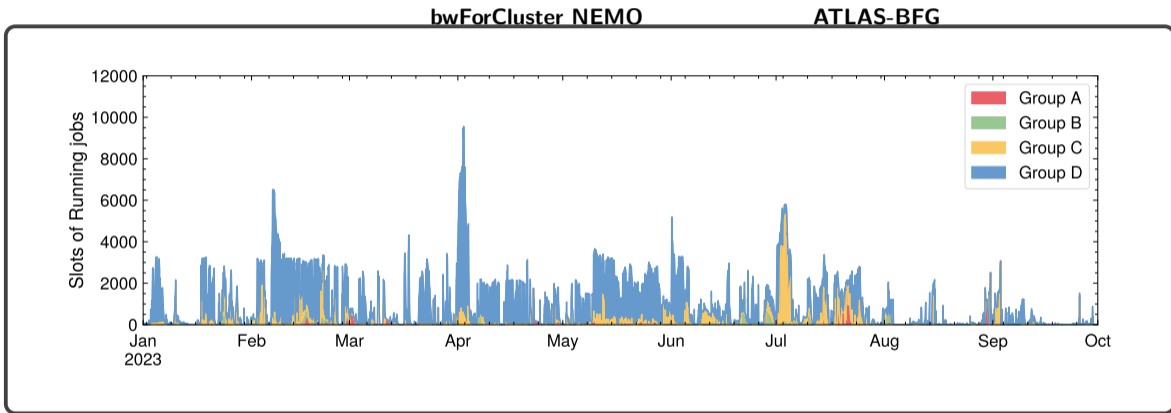- ▶ Efficient use of resou...

universität freiburg
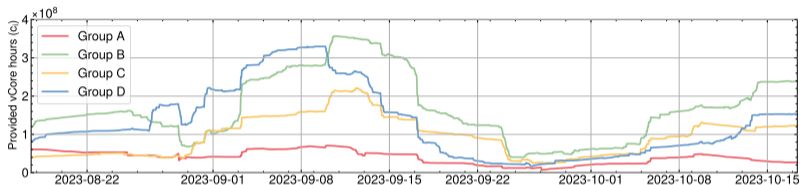
# Priority Use Case

Motivation



- ▶ Each has its own SLURM partition (job queue)
- ▶ Efficient use of resources due to sharing VMs across HEP groups

# Priority Use Case
Results from HEP groups @ University of Freiburg
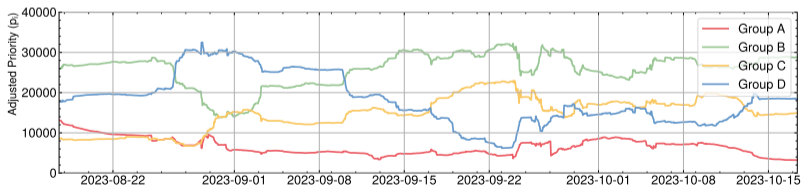
▶ Provided resources of the four local HEP groups



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

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

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

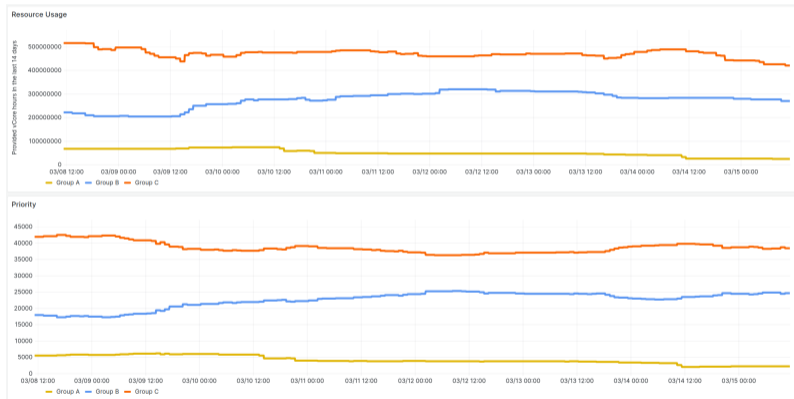▶ Priority is adjusted according to the provided resources



Priority $p_i$ is defined as:

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

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

# Real-time monitoring

e.g. Priority Use Case, but also AUDITOR stats



- ▶ Resource usage and priority can be made available for Prometheus
- ▶ Real-time monitoring of priority adjustments (with e.g. Grafana)