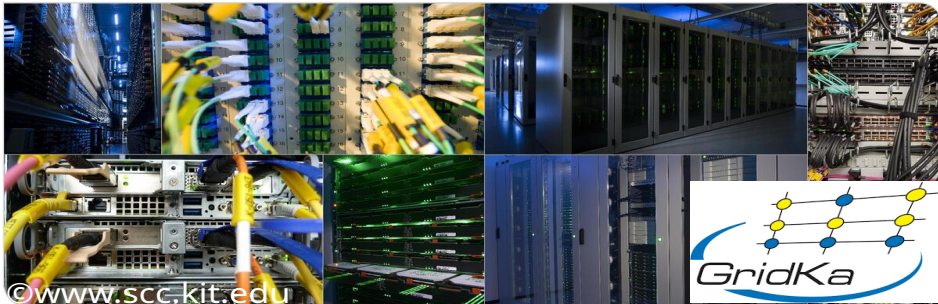# GridKa Overview Report

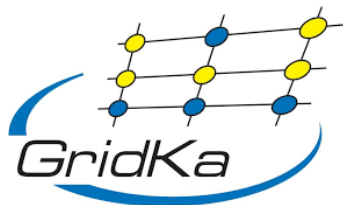**@ Annual Meeting of the ATLAS and CMS Computing Verbund**

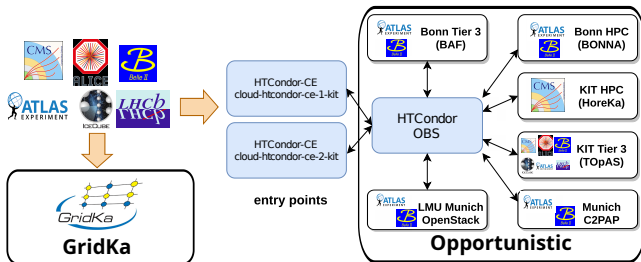**Robin Hofsaess** on behalf of the GridKa CMS and R&D team | 26.03.2024

# Outline

- GridKa Status and News
  - Pledges and Resources for CMS
  - Opportunistic Resources at GridKa
  - New NVMe cache for HPSS tape system
  - WLCG Data Challenge '24
  - ARM workers delivered
  - Update to RHEL8
  - Update of Compute Elements
  - HappyFace4 development
  - Progress in GPU Usage for CMS

# GridKa Overview



## Opportunistic Resources

Successfully integrated with COBalD/TARDIS – developed at KIT

# GridKa Status

| T1_DE_KIT | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **GGUS Tickets:** 165452 | | | | | | | | | | | | | | | | |
| **Downtimes:** | | | | | | | | | | | | | | | | |
| **SAM Status:** | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 72% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% |
| **Hammer Cloud:** | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 100% | 99% | 100% | 100% | 100% | 100% | 100% | 100% |
| **FTS Status:** | 0% | 0% | 50% | 100% | 0% | 0% | 100% | 100% | 100% | 0% | 100% | 0% | 100% | 0% | 100% | 0% |
| **Site Readiness:** | 74% | 79% | 99% | 89% | 100% | 99% | 99% | 68% | 97% | 99% | 95% | 100% | 97% | 99% | 97% | 87% |
| **Life Status:** | | | | | | | | | | | | | | | | |
| **Prod Status:** | | | | | | | | | | | | | | | | |
| **CRAB Status:** | | | | | | | | | | | | | | | | |
| | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 |
| | | | | | | | | | | Mar | | | | | | |

## Site Status

GridKa is one of the most reliable Tier 1 centers! (Readiness Report)

# Pledges and Resources for CMS

| Resource | Pledges 2024 |
|---|---|
| CPUs | 93k HEPscore23 (≈ 6900k Cores) |
| Disk | 12.2 PB |
| Tape | 38 PB |

**Tier 3: (opportunistic)**

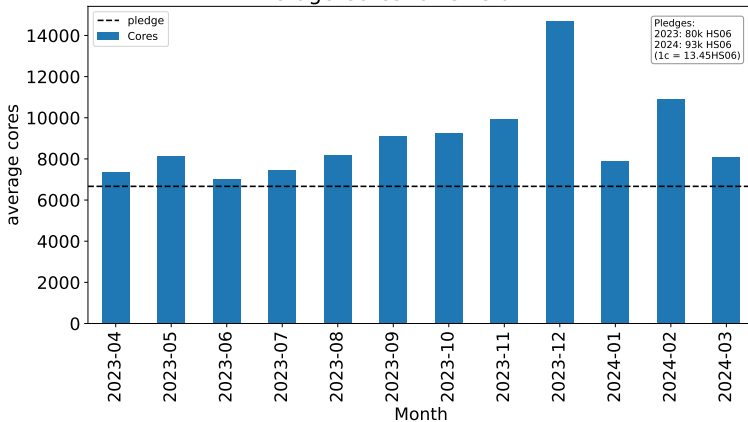| | |
|---|---|
| **dCMS** Disk | 2.8 PB |
| CPUs | 1000 cores |
| *GPUs | 56 |

*Prototype integration
of GPUs into the Grid

## Summary

- The pledged tape and disk will be fully available in April
- All pledges will be fulfilled
- Additionally: We provide 56 GPUs with our Tier 3 (accessible for CMS production and via CRAB/CMS Connect for users)

Average Cores for CMS at T1 KIT

Pledges:
2023: 80k HS06
2024: 93k HS06
(1c = 13.45HS06)

## Tier 1 Compute Resources

**Always over pledge for all VOs**

GridKa Status and News

# Integration of HoreKa



CoreHours (KIT HPC Contribution)

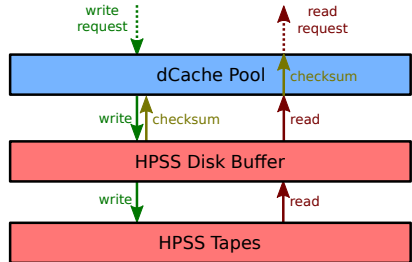## HoreKa HPC Resources

**(Opportunistic) contribution up to an average T2!**

# NVMe Cache for HPSS Tape System

- New NVMe cache **in production** since Q4 '23
  - 300TB fast NVMe
  - 4x100G network
  - To replace old HDD cache
- Currently, 95 PB on GridKa tape in total

**Milestone:**

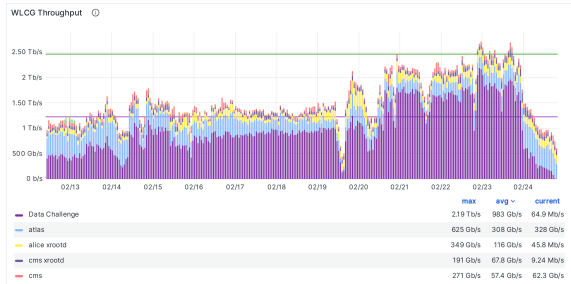Stable operation with the desired **300 MB/s - 400 MB/s** read and write rate per tape drive



More details on our tape system in backup

GridKa Status and News

# WLCG Data Challenge '24: Overview

- 12.02 to 24.02
- 1st week:
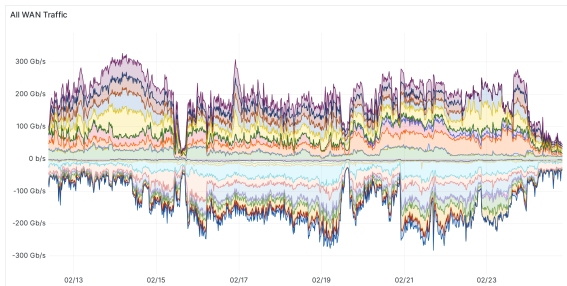  1.2 Tb/s target ✓
- 2nd week:
  2.4 Tb/s target (✓)



## Result

In general, a success and valuable lessons learned for further optimization

# WLCG Data Challenge '24: GridKa Network Perspective

- ■ LHCOPN: 300 Gb/s in each direction for DC 24
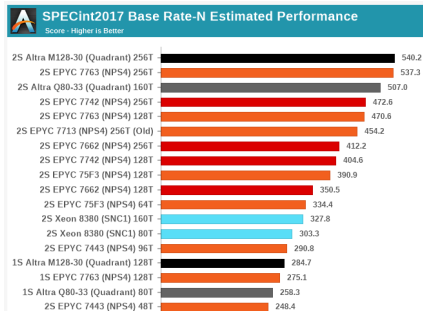- ■ LHCONE+Internet: 200 Gb/s in each direction



## Conclusion

From the network perspective, GridKa was not at its limits

# ARM Worker Nodes Delivered

- 15 nodes ordered with 2x AMPERE_Altra Max (2x128 cores)

- Test machines are promising (80 Cores 3,3 GHz)

| AMD (64c) | ARM (80c) |
|-----------|-----------|
| 2828 HS23 | 1606 HS23 |
| 600W | 380W |
| 4.71 HS/W | **5.74** HS/W |



SPECint2017 Base Rate-N Estimated Performance
Score · Higher is Better

| | |
|---|---|
| 2S Altra M128-30 (Quadrant) 256T | 540.2 |
| 2S EPYC 7763 (NPS4) 256T | 537.3 |
| 2S Altra Q80-33 (Quadrant) 160T | 507.0 |
| 2S EPYC 7742 (NPS4) 256T | 472.6 |
| 2S EPYC 7763 (NPS4) 128T | 470.6 |
| 2S EPYC 7713 (NPS4) 256T (Old) | 454.2 |
| 2S EPYC 7662 (NPS4) 256T | 412.2 |
| 2S EPYC 7742 (NPS4) 128T | 404.6 |
| 2S EPYC 75F3 (NPS4) 128T | 390.9 |
| 2S EPYC 7662 (NPS4) 128T | 350.5 |
| 2S EPYC 75F3 (NPS4) 64T | 334.4 |
| 2S Xeon 8380 (SNC1) 160T | 327.8 |
| 2S Xeon 8380 (SNC1) 80T | 303.3 |
| 2S EPYC 7443 (NPS4) 96T | 290.8 |
| 1S Altra M128-30 (Quadrant) 128T | 284.7 |
| 1S EPYC 7763 (NPS4) 128T | 275.1 |
| 1S Altra Q80-33 (Quadrant) 80T | 258.3 |
| 2S EPYC 7443 (NPS4) 48T | 248.4 |

Source: anandtech.com

## Provisioning: Q2 '24 (expected)

Will accept jobs from our CEs
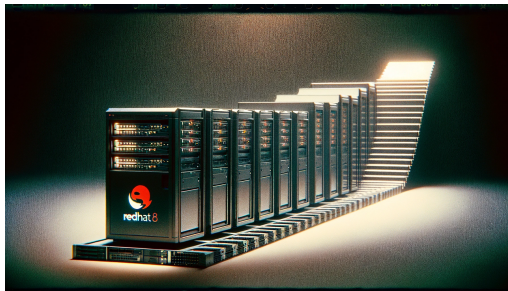
More info: CHEP poster on energy efficiency of ARM

GridKa Status and News
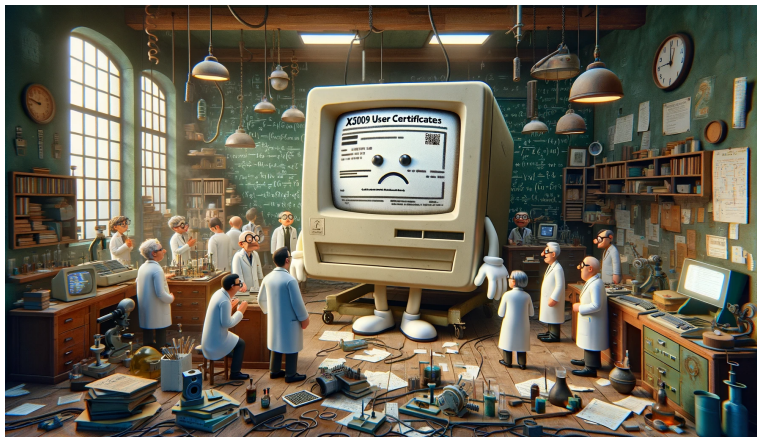○○○○○○○○●○○○○○○○○○

# ARM Worker Nodes Delivered

# Update to RHEL8



- CentOS 7 will reach EOL end of June
- All machines will be updated to RHEL8 in the near future
  - Includes also the Tier 3 resources

# Update of CEs: Farewell to X509!



Within the next months, our **CEs** will be updated and X509 certificates will be deprecated.
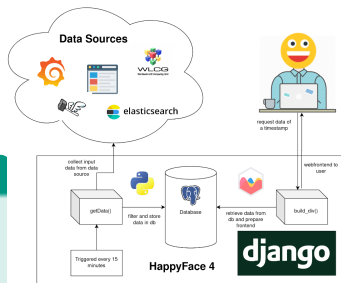
# HappyFace4

- What is HappyFace?
  - Meta monitoring tool for computing sites
  - Gives user a fast overview of the site status
  - GridKa production instance: happyface@ETP



## Purpose

- Advanced (meta) monitoring for our Tier 1 center
- All necessary information for reliable operations gathered in one place

→ Extremely helpful for shifters to detect and report problems fast and in detail!
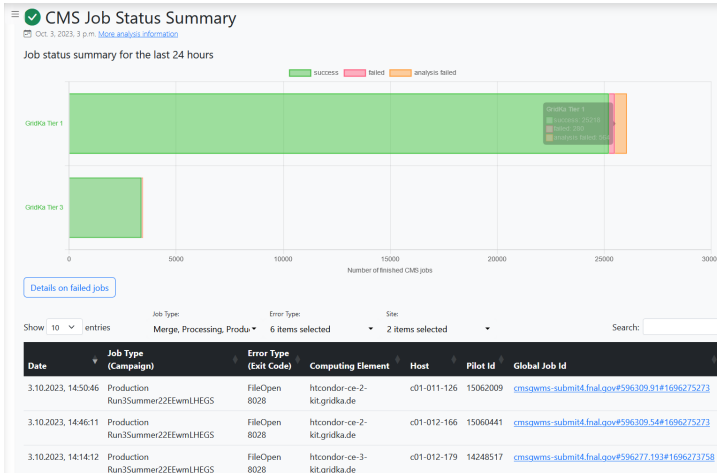
# HappyFace4: Example



## Don't panic

- **Red** $\neq$ GridKa broken!

- **All** red and yellow issues are regularly followed up by our operations team to be understood :-)

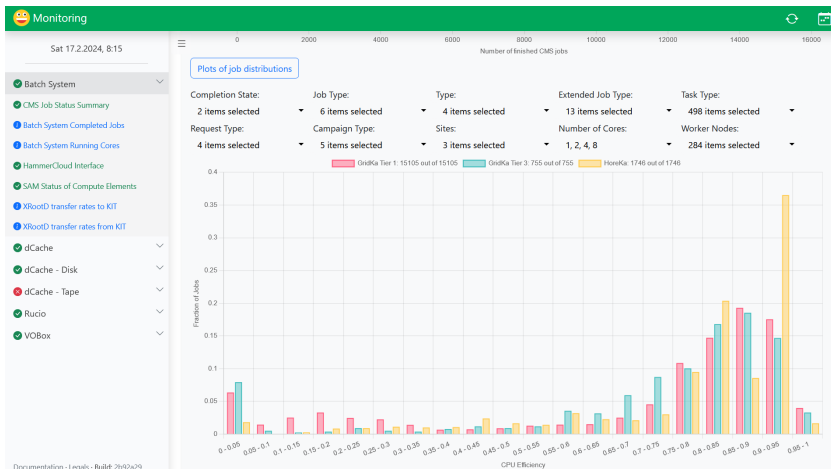# HappyFace4: Details

# HappyFace4: Recent Development



## New Features                              (credits to Artur Gottmann!)

CPU efficiency monitoring per job grouped by site, job type, cores, ...

# HappyFace4: Recent Development

# Progress in GPU Usage for CMS

*GridKa provides GPUs as a (opportunistic) prototype service via the Tier 3!*

- **Current Development and Challenges:**
    - Efficient grid integration of GPU resources
    - Scheduling optimization of small scale GPU jobs
    - Contribution to the future of GPU in the HEPScore benchmark
    - Energy efficiency benchmarks of HEP GPU applications

- **Currently used by CMS DeepTau group**

GridKa Status and News

# Summary

## GridKa Status

- **No problems** to report
- One of the most **reliable** Tier 1s
- All pledges will be **fulfilled**

## Hardware

- New **NVME cache** is fully operational
- **ARM** machines delivered are about to be provisioned

## Data Challenge 24

- GridKa participated successfully
- Within **Top-3** of sites
- Network has further potential

## Opportunistic Resources

- Our opportunistic resources successfully provide up to **several 1000** additional cores

## Upcoming Updates

- All machines will be upgraded to **RHEL8**
- With updating the Compute Elements, **X509 certificates** will be deprecated

## Recent Development

- **HappyFace4**, our multi-facet observation tool for the Tier 1 operation, is constantly improved
- **GPU** integration and optimization

GridKa Status and News

# BACKUP

# The GridKa CMS and R&D Team!

- Prof. Dr. Günter Quast
- Dr. Manuel Giffels
- Dr. Artur Gottmann
- Dr. Matthias Schnepf (BELLE, T3)
- Dr. Max Fischer (ALICE, T3)
- Dr. Maximilian Horzela
- Dr. Sebastian Brommer

- Robin Hofsaess
- Tim Voigtlaender
- Jonas Eppelt
- Lars Sowa
- Cedric Verstege
- Jost von den Driesch
- Christian Winter

# Comparison of T1s

monit-grafana.cern.ch: cms-tier-1-utilization
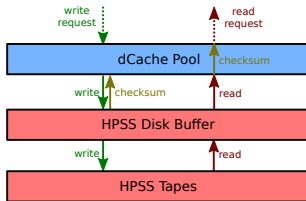
# Tape Storage at GridKa

- Since March 2022
  - Larger tape capacity: 8 → 20 TB
  - Higher tape drive speed:
    150 → 400 MB/s
  - 1/2 PB SSD+NVME buffer as
    part of the system
  - In full operation for CMS, Belle 2,
    and LHCb
- Data from the old system fully
  migrated for ATLAS, CMS, Belle 2,
  and LHCb
- Planning to finish migration for
  ALICE in summer 2024



top: tape cartridge & drive, bottom: tape library at KIT (current total capacity: 150 PB)

# Schematic overview of GridKa tape system

- Write request:
    - Incoming file transfer at dCache disk pool
    - Written from dCache to HPSS disk buffer
    - Read back for checksum consistency test
    - Within HPSS, writing to tapes initiated afterwards in **file aggregates**
- Read request:
    - File read requests appear at dCache pool
    - Requests grouped by tape & aggregate
    - **Entire aggregates** read from tapes to HPSS disk buffer
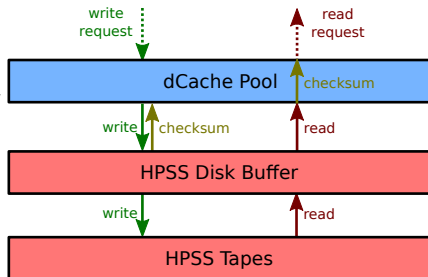    - Files read from HPSS disk to dCache pool



Files in the same directory collected
into aggregates of up to 300 GB

Important fraction of in-house written interface done by **ATLAS & CMS** representatives
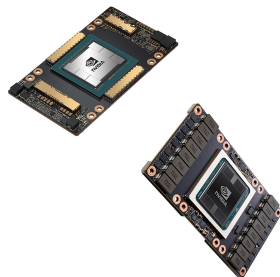
# Details on the NVMe Cache

- Current Setup:
  - AMD EPYC 9554P 64-Core
  - 300TB NVMe in XINNOR softraid ($\rightarrow$ 500TB ssd cache in total)
  - 700k IOPS for r+w
- Optimized for low latency for many clients to
- Performance benchmark:
  - 10 or 5 times 10 Clients with each sequential reads (2/3) and writes (1/3) of a 5GB file
  - Throughput: around 50-70GB/s
  - $\rightarrow$ constant 300 MB/s to 400 MB/s write speed per tape drive

**Final setup TBD (potential alternatives: GRAID or all-flash ararys)**



Backup
○○○○○●○

# GPUs at GridKa Made Available for CMS

- Several GPU's deployed at GridKa TOpAS cluster and provided to entire CMS through the grid
  - 24 × Nvidia A100
  - 24 × Nvidia V100
  - 8 × Nvidia V100S
- GPU workflows sent by CMS were successfully completed
  - High Level Trigger Test Workflow
  - Release Validation Workflow

## Conclusion

We are well prepared for heterogeneous computing era!

Backup