

## Dynamic Federation of Heterogeneous Compute Resources in High Energy Physics & Beyond PUNCH-Lunch, 16.09.2021 Manuel Giffels & Kilian Schwarz



KIT – The Research University in the Helmholtz Association

### www.kit.edu



# The High Energy Physics Use-case

- Particle detectors record physics event data
- Each detector used by a collaboration of scientists





Slide by Max Fischer



# The High Energy Physics Use-case

- Particle detectors record physics event data
- Each detector used by a collaboration of scientists



![](_page_2_Picture_4.jpeg)

![](_page_2_Picture_5.jpeg)

![](_page_2_Picture_10.jpeg)

- storage/compute centres
- Resources shared via a worldwide computing Grid

Slide by Max Fischer

![](_page_2_Picture_15.jpeg)

# The High Energy Physics Use-case

- Particle detectors record physics event data
- Each detector used by a collaboration of scientists

![](_page_3_Picture_3.jpeg)

![](_page_3_Picture_4.jpeg)

![](_page_3_Picture_5.jpeg)

Collaborators provide storage/compute centres

![](_page_3_Picture_10.jpeg)

#### Collaborations automate 800÷ -00 common pre-processing Scientists run individual end-user analyses 820÷ 820× 820-000+ 800

## Resources shared via a worldwide computing Grid

![](_page_3_Figure_13.jpeg)

![](_page_3_Figure_15.jpeg)

![](_page_3_Figure_16.jpeg)

![](_page_3_Figure_17.jpeg)

![](_page_4_Picture_1.jpeg)

## ~20 years experience

![](_page_4_Picture_5.jpeg)

![](_page_4_Picture_7.jpeg)

![](_page_4_Picture_8.jpeg)

![](_page_5_Picture_1.jpeg)

#### **Global collaboration**

42 countries 170 computing centres Over 2 million tasks daily 1 million computer cores 1 exabyte of storage

![](_page_5_Picture_4.jpeg)

![](_page_5_Picture_5.jpeg)

### ~20 years experience

![](_page_5_Picture_9.jpeg)

![](_page_5_Figure_10.jpeg)

![](_page_5_Picture_13.jpeg)

![](_page_5_Picture_14.jpeg)

![](_page_6_Figure_1.jpeg)

## ~20 years experience

![](_page_6_Picture_5.jpeg)

![](_page_6_Figure_6.jpeg)

![](_page_6_Picture_8.jpeg)

![](_page_6_Picture_9.jpeg)

![](_page_7_Figure_1.jpeg)

## ~20 years experience

![](_page_7_Picture_5.jpeg)

#### **Global Trust Federation**

Established a trusted set of identity credential providers avoiding user registration at each entity.

![](_page_7_Picture_9.jpeg)

![](_page_7_Picture_10.jpeg)

![](_page_8_Figure_1.jpeg)

## ~20 years experience

![](_page_8_Picture_5.jpeg)

![](_page_8_Picture_6.jpeg)

#### **Global Trust Federation**

Established a trusted set of identity credential providers avoiding user registration at each entity.

![](_page_8_Picture_12.jpeg)

![](_page_9_Figure_1.jpeg)

![](_page_9_Picture_4.jpeg)

![](_page_9_Picture_5.jpeg)

#### **Global Trust Federation**

Established a trusted set of identity credential providers avoiding user registration at each entity.

## Most relevant achievement is not a technical one, it is the establishment of the Global Trust Federation.

Manuel Giffels

![](_page_9_Picture_13.jpeg)

![](_page_10_Figure_1.jpeg)

![](_page_10_Picture_4.jpeg)

Integrate resources into a globally distributed batch system and remove a significant fraction of the initial Grid middleware

Manuel Giffels

![](_page_10_Picture_8.jpeg)

![](_page_10_Picture_9.jpeg)

![](_page_11_Figure_1.jpeg)

![](_page_11_Picture_4.jpeg)

![](_page_11_Picture_7.jpeg)

![](_page_11_Picture_8.jpeg)

![](_page_12_Figure_1.jpeg)

![](_page_12_Picture_4.jpeg)

![](_page_12_Picture_7.jpeg)

![](_page_12_Picture_8.jpeg)

![](_page_13_Figure_1.jpeg)

![](_page_13_Picture_4.jpeg)

![](_page_13_Picture_7.jpeg)

![](_page_13_Picture_8.jpeg)

![](_page_14_Figure_1.jpeg)

![](_page_14_Picture_4.jpeg)

![](_page_14_Picture_7.jpeg)

![](_page_14_Picture_8.jpeg)

![](_page_15_Figure_1.jpeg)

![](_page_15_Picture_4.jpeg)

![](_page_15_Picture_7.jpeg)

![](_page_15_Picture_8.jpeg)

![](_page_16_Figure_1.jpeg)

![](_page_16_Picture_4.jpeg)

![](_page_16_Picture_7.jpeg)

![](_page_16_Picture_8.jpeg)

# **Upcoming Computing Challenges in HEP**

- HL-LHC poses unprecedented challenges to HEP computing
- Assuming flat budget and 10-20% technology advance per year
- CPU shortcoming of factor 2 3 estimated in 2027
- Critical conditions for HL-LHC (Run 4)

![](_page_17_Picture_8.jpeg)

![](_page_17_Figure_9.jpeg)

![](_page_17_Picture_11.jpeg)

![](_page_17_Picture_12.jpeg)

# **Upcoming Computing Challenges in HEP**

- HL-LHC poses unprecedented challenges to HEP computing
- Assuming flat budget and 10-20% technology advance per year
- CPU shortcoming of factor 2 3 estimated in 2027
- Critical conditions for HL-LHC (Run 4)

## **Potential Measures:**

Software & Computing Evolution (R&D)

Add additional resources (opportunistic resources)

![](_page_18_Picture_11.jpeg)

![](_page_18_Figure_12.jpeg)

Manuel Giffels

![](_page_18_Picture_15.jpeg)

**Opportunistic Resources** Any resources not permanently dedicated to but temporarily available for a specific task, user or group.

![](_page_19_Picture_4.jpeg)

![](_page_19_Picture_5.jpeg)

![](_page_19_Picture_7.jpeg)

![](_page_19_Picture_8.jpeg)

![](_page_19_Picture_9.jpeg)

**Opportunistic Resources** Any resources not permanently dedicated to but temporarily available for a specific task, user or group.

## **Challenges:**

Different OS & Software availability

![](_page_20_Picture_4.jpeg)

**Container Technology** 

![](_page_20_Picture_8.jpeg)

![](_page_20_Picture_9.jpeg)

![](_page_20_Picture_11.jpeg)

![](_page_20_Picture_12.jpeg)

![](_page_20_Picture_13.jpeg)

**Opportunistic Resources** Any resources not permanently dedicated to but temporarily available for a specific task, user or group.

## **Challenges:**

- Different OS & Software availability
- Very heterogenous systems, not all resources are suited for all tasks
- Varying availability of and demand for those resources

![](_page_21_Picture_6.jpeg)

![](_page_21_Picture_9.jpeg)

![](_page_21_Picture_10.jpeg)

![](_page_21_Picture_14.jpeg)

![](_page_21_Picture_15.jpeg)

**Opportunistic Resources** Any resources not permanently dedicated to but temporarily available for a specific task, user or group.

## **Challenges:**

- Different OS & Software availability
- Very heterogenous systems, not all resources are suited for all tasks
- Varying availability of and demand for those resources
- No global trust federation/Grid entry point available

![](_page_22_Picture_7.jpeg)

![](_page_22_Picture_8.jpeg)

![](_page_22_Picture_9.jpeg)

![](_page_22_Picture_13.jpeg)

![](_page_22_Picture_14.jpeg)

![](_page_22_Picture_16.jpeg)

![](_page_22_Picture_17.jpeg)

![](_page_22_Picture_18.jpeg)

**Opportunistic Resources** Any resources not permanently dedicated to but temporarily available for a specific task, user or group.

## **Challenges:**

- Different OS & Software availability
- Very heterogenous systems, not all resources are suited for
- Varying availability of and demand for those resources
- No global trust federation/Grid entry point available

![](_page_23_Picture_7.jpeg)

![](_page_23_Picture_8.jpeg)

![](_page_23_Picture_9.jpeg)

![](_page_23_Picture_13.jpeg)

![](_page_23_Picture_14.jpeg)

OBS + Pilots in more generalized way

Manuel Giffels

![](_page_23_Picture_18.jpeg)

![](_page_23_Picture_19.jpeg)

![](_page_23_Picture_20.jpeg)

![](_page_23_Picture_21.jpeg)

![](_page_23_Picture_22.jpeg)

![](_page_23_Picture_23.jpeg)

![](_page_23_Picture_24.jpeg)

# The COBaID View of Resource Scheduling

[COBalD - the **O**pportunistic **Bal**ancing **D**aemon]

![](_page_24_Picture_4.jpeg)

Based on a slide by Max Fischer

Manuel Giffels

![](_page_24_Picture_8.jpeg)

![](_page_24_Picture_9.jpeg)

![](_page_24_Picture_10.jpeg)

# The COBalD View of Resource Scheduling

[COBalD - the **O**pportunistic **Balancing D**aemon]

## Resource Meta-Scheduling for Job Scheduler is a "hard" problem

Usually based on predictions of the future resource availability and future mixture of job classes (CPU intense, I/O intense, ...)

![](_page_25_Picture_6.jpeg)

Based on a slide by Max Fischer

![](_page_25_Picture_10.jpeg)

![](_page_25_Picture_11.jpeg)

# The COBalD View of Resource Scheduling

[COBalD - the **O**pportunistic **Balancing Daemon**]

## Resource Meta-Scheduling for Job Scheduler is a "hard

Usually based on predictions of the future resource avail future mixture of job classes (CPU intense, I/O intense, ...)

Works perfectly fine in homogenous environments.

aliu

Based on a slide by Max Fischer

![](_page_26_Picture_10.jpeg)

![](_page_26_Figure_11.jpeg)

![](_page_26_Picture_12.jpeg)

# The COBaID View of Resource Scheduling

[COBalD - the **O**pportunistic **Balancing Daemon**]

## Resource Meta-Scheduling for Job Scheduler is a "hard

Usually based on predictions of the future resource avail future mixture of job classes (CPU intense, I/O intense, T.)

## However: We usually care only about a simpler problem!

![](_page_27_Figure_5.jpeg)

### Resource allocation over time

Works perfectly fine in homogenous environments.

y anu

Based on a slide by Max Fischer

![](_page_27_Picture_14.jpeg)

![](_page_27_Figure_15.jpeg)

![](_page_27_Picture_16.jpeg)

# The COBalD View of Resource Scheduling

[COBalD - the **O**pportunistic **Balancing Daemon**]

## Resource Meta-Scheduling for Job Scheduler is a "hard

- Usually based on predictions of the future resource avail future mixture of job classes (CPU intense, I/O intense, ...)
- However: We usually care only about a simpler problem!

![](_page_28_Figure_5.jpeg)

### Resource allocation over time

Works perfectly fine in homogenous environments.

y anu

Based on a slide by Max Fischer

![](_page_28_Picture_15.jpeg)

![](_page_28_Figure_16.jpeg)

![](_page_28_Picture_17.jpeg)

# The COBalD View of Resource Scheduling

[COBalD - the **O**pportunistic **Balancing D**aemon]

- Resource Meta-Scheduling for Job Scheduler is a "hard"
- COBaID cares only about resources, not jobs
  - Observe how much and how well each resource is used
  - Increase well-used resources, decrease unused resources

![](_page_29_Picture_9.jpeg)

![](_page_29_Picture_10.jpeg)

Slide by Max Fischer

![](_page_29_Picture_13.jpeg)

![](_page_30_Figure_0.jpeg)

![](_page_30_Picture_4.jpeg)

![](_page_30_Figure_7.jpeg)

![](_page_31_Figure_0.jpeg)

![](_page_31_Picture_6.jpeg)

# **The COBaID/TARDIS - Resource Scheduler**

[COBaID - the **O**pportunistic **Balancing Daemon**] [**T**ransparent **A**daptive **R**esource **D**ynamic Integration **S**ystem]

## COBalD:

- Look at what is used, not what is requested
  - Simple logic: more used, less unused resources
  - COBalD acquires/releases resources
  - Batch system scheduler handles jobs
- TARDIS (a COBalD plugin):
- Defines VM/Container/Job(Script) as resource
- Provides access to resource provider APIs OpenStack, CloudStack, HTCondor, Slurm, Moab and K8S
- Integrates resources into Overlay Batch System HTCondor and Slurm are supported
- Manages resource life cycle

![](_page_32_Figure_14.jpeg)

![](_page_32_Picture_17.jpeg)

# **Opportunistic Resources & WLCG in Practice**

![](_page_33_Figure_1.jpeg)

- Dynamic, transparent and on-demand integration using **COBaID/TARDIS** resource manager developed at KIT
- Single point(s) of entry to opportunistic resources (HPCs, local clusters and Clouds) using traditional Grid Compute Elements
- HTCondor overlay batch system to federate a variety of compute resources

![](_page_33_Picture_7.jpeg)

![](_page_33_Picture_8.jpeg)

![](_page_33_Picture_9.jpeg)

![](_page_33_Figure_11.jpeg)

![](_page_33_Figure_12.jpeg)

![](_page_33_Figure_13.jpeg)

Establish a federated heterogeneous compute infrastructure for PUNCH Integrate data storages, archives and opportunistic caches

![](_page_34_Figure_2.jpeg)

Introduce data-locality aware scheduling Benefit from experiences, concepts and tools available in HEP community

### **Opportunistic Resources**

#### Kilian Schwarz

#### **Infrastructure Evolution**

![](_page_36_Figure_1.jpeg)

Caches require less operational cost while minimising external I/O, and providing high data locality

# Dynamic disk caches have been developed in ErUM-Data-Pilot

- XrootD-plug-in based ingredients
  - ProxyPrefix (local clients are redirected to the local XrootD forward proxy)
  - RedirLocal (local clients read from local file systems if data is available)
  - In a classic XrootD based Grid SE clients would always read via XrootD redirector and data server

![](_page_37_Figure_5.jpeg)

# Disk Caching Proxy Setup (Disk Cache on the fly)

- · Clients access data through the redirector
- If data exists redirector redirects clients to local file system
- Otherwise redirector redirects to disk caching forward proxy
- Disk caching proxy forwards request to external site and retrieves the data
- · Data are being cached on local file system for later use

![](_page_38_Figure_6.jpeg)

### **Cache Scenarios for Benchmarks**

![](_page_39_Figure_1.jpeg)

### Cache performance test setup

- Jobs submitted via Slurm to LOEWE CSC (Frankfurt) running in Singularity container
- Sample analysis reads events from branch and does nothing further
- Dynamic Disk Cache setup at LOEWE CSC
- Data located at GSI Lustre Cluster served via XrootD data server
- Bandwidth connection GSI LOEWE CSC 10 Gb/s

### Cache performance results

#### ALICE ESD Data in Singularity

![](_page_41_Figure_2.jpeg)

## Modelling of heterogeneous distributed systems with coordinated caches

- Fixing workload and scheduler: How do we achieve the most efficient data processing?
  - How do we design the network?
  - How many caches do we need to place? Where?
  - What is the optimal size and caching logic of these caches?
  - Can we replace managed storage with caches without losing momentum?
  - Which figure of merit do we want to optimize? User walltime? Monetary costs? WAN bandwidth?
  - Do we cache for all workflow types? If not, for which?

![](_page_42_Picture_8.jpeg)

# Questions can be answered by simulation

![](_page_43_Picture_1.jpeg)

#### **Simulator LAPIS**

![](_page_43_Figure_3.jpeg)

T. Feßenbecker, Modeling of distributed coordinated caching for LHC data analyses

Central metric for job characterization: job walltime

$$t_{ ext{wall}} = \max\left(t_{ ext{calculation}}, t_{ ext{transfer}}
ight) = \max\left(rac{t_{ ext{CPU}}}{arepsilon_{ ext{instr}}}, rac{V \cdot (1-h)}{b_{ ext{remote}}}, rac{V \cdot h}{b_{ ext{cache}}}
ight)$$

### Summary and outlook

- Dynamic disk caches have the potential to increase throughput significantly
- Dynamic disk caches are an efficient tool for integrating opportunistic resources
- DCOTF and XCache DCA have a similar performance
- XCache creates a bottleneck in terms of CPU and bandwidth (number of threads on caching node)
- No cache is limited by intersite network bandwidth
- Containerisation does not affect performance
- Strategic placement of dynamic disk caches can be optimised via KIT cache simulation