# R&D for Future High Throughput Computing @ GridKa
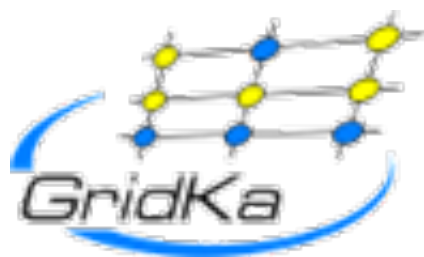
TRIUMF-Helmholtz Workshop on Scientific Computing, DESY Hamburg - 16.09.2019

**Manuel Giffels**, **René Caspart, Florian v. Cube, Tabea Fesenbecker, Max Fischer, Christoph Heidecker, Eileen Kühn, Matthias Schnepf**

Institute for Particle Physics (ETP) & Steinbuch Centre for Computing (SCC)

**German Tier-1 High Energy Physics and Astroparticle Data & Analysis Centre**
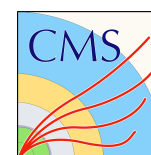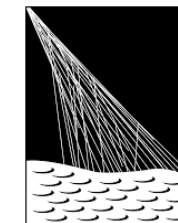
- Support all four LHC experiments
- Belle II, Pierre-Auger, several small communities

**Joint R&D** with computer science towards HL-LHC

**Resources**

- Compute: ~29k cores
- Disk: 37PB (used), Tape: 54 PB (used)
- 100 Gb/s connection to LHCONE/ LHCOPN

Among the largest and best performing T1s

Annual international **GridKa School**

- > 1800 participants since 2003
  https://gridka.school

# R&D Environment and Backgrounds

- Two collaborating HEP Computing Groups at KIT
    - SCC: GridKa Tier 1, focus on throughput and production systems
    - ETP: Institute Tier 3, focus on responsiveness and prototypes

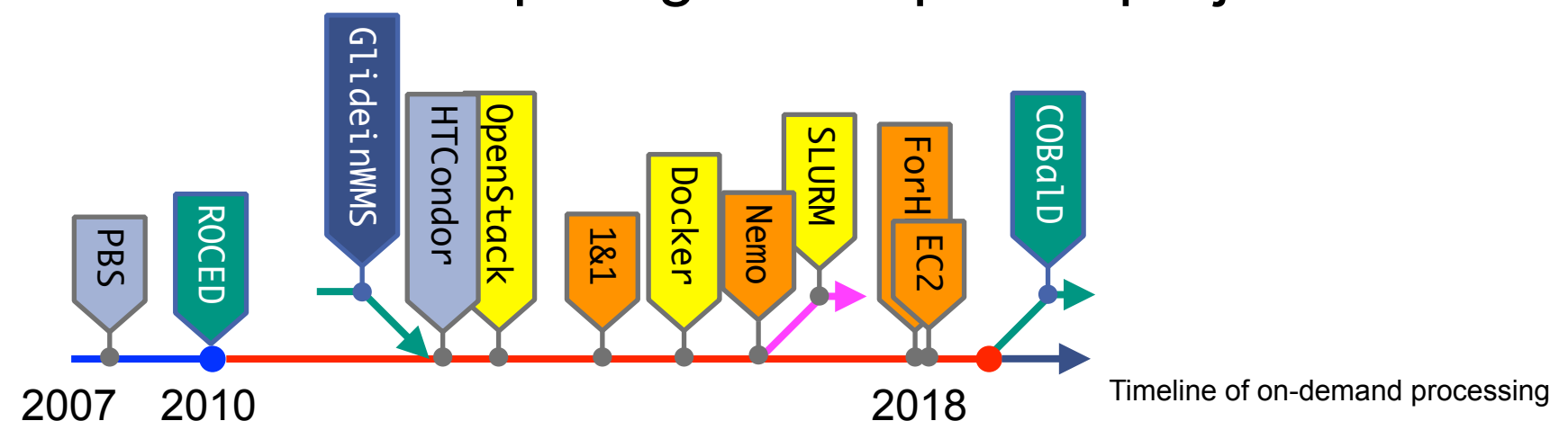Manuel Giffels

# R&D Environment and Backgrounds

- Two collaborating HEP Computing Groups at KIT
    - SCC: GridKa Tier 1, focus on throughput and production systems
    - ETP: Institute Tier 3, focus on responsiveness and prototypes
- The major topics of Research and Development
    - Dynamic on-demand processing resources via VMs/containers
      T. Hauth et al., On-demand provisioning of HEP Compute Resources on Clouds Sites and Shared HPC Centers,
      Journal of Physics **898**, 5 (2017)
    - Adaptive placement of input data via distributed coordinated caches
      M. Fischer et al., Opportunistic Data Locality for End User Data Analysis, Journal of Physics **898**, 5 (2017)

# R&D Environment and Backgrounds

- Two collaborating HEP Computing Groups at KIT
  - SCC: GridKa Tier 1, focus on throughput and production systems
  - ETP: Institute Tier 3, focus on responsiveness and prototypes
- The major topics of Research and Development
  - Dynamic on-demand processing resources via VMs/containers
    T. Hauth et al., On-demand provisioning of HEP Compute Resources on Clouds Sites and Shared HPC Centers, Journal of Physics **898**, 5 (2017)
  - Adaptive placement of input data via distributed coordinated caches
    M. Fischer et al., Opportunistic Data Locality for End User Data Analysis, Journal of Physics **898**, 5 (2017)
- Longstanding involvement in computing development projects



Timeline of on-demand processing

2007    2010                                                          2018
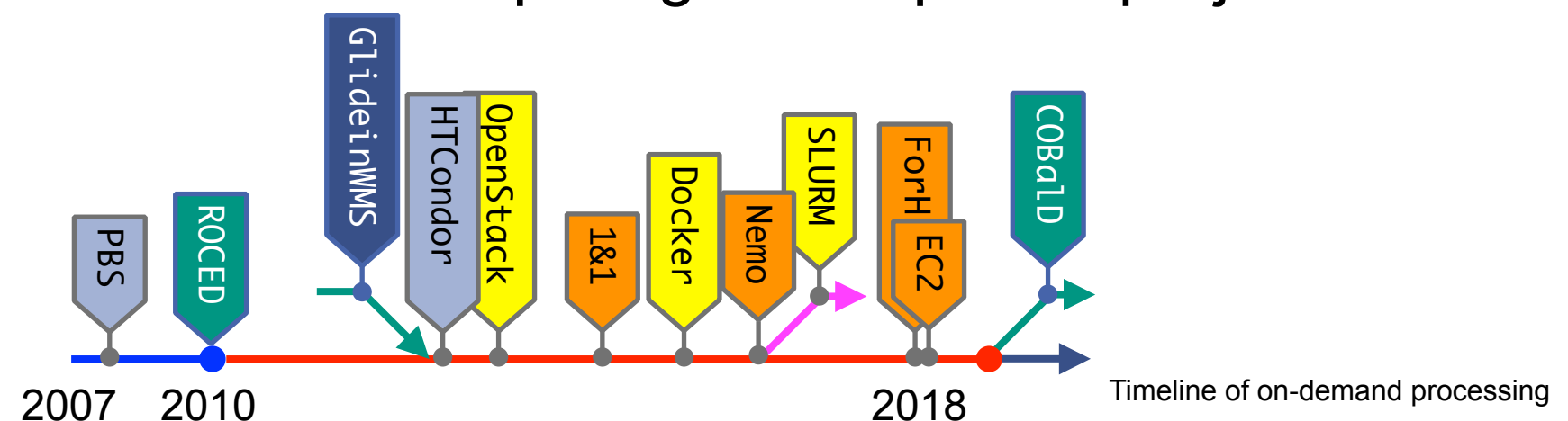
# R&D Environment and Backgrounds

- Two collaborating HEP Computing Groups at KIT
  - SCC: GridKa Tier 1, focus on throughput and production systems
  - ETP: Institute Tier 3, focus on responsiveness and prototypes
- The major topics of Research and Development
  - Dynamic on-demand processing resources via VMs/containers

    T. Hauth et al., On-demand provisioning of HEP Compute Resources on Clouds Sites and Shared HPC Centers, Journal of Physics **898**, 5 (2017)
  - Adaptive placement of input data via distributed coordinated caches

    M. Fischer et al., Opportunistic Data Locality for End User Data Analysis, Journal of Physics **898**, 5 (2017)
- Longstanding involvement in computing development projects



Timeline of on-demand processing

| | |
|---|---|
| PBS | ROCED |
| GlideinWMS | HTCondor | OpenStack |
| 1&1 | Docker | Nemo | SLURM | ForHLR2 EC2 | COBalD |

2007  2010  2018

**Opportunistic Resource**

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

# Opportunistic Resources and their Challenges

Very different to the traditional HEP environment:

# Opportunistic Resources and their Challenges

Very different to the traditional HEP environment:
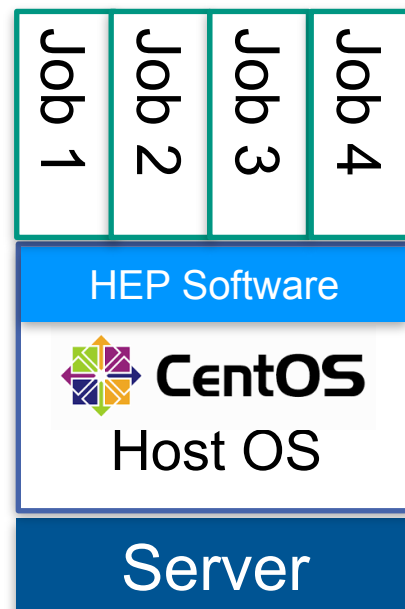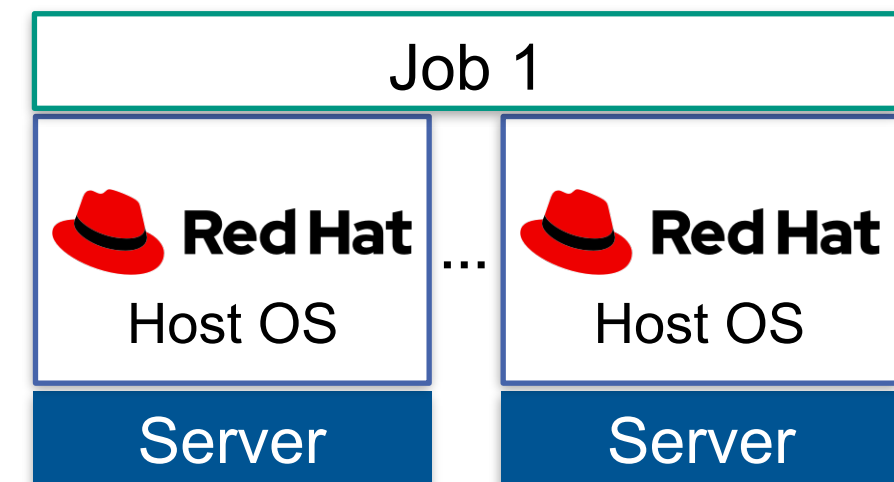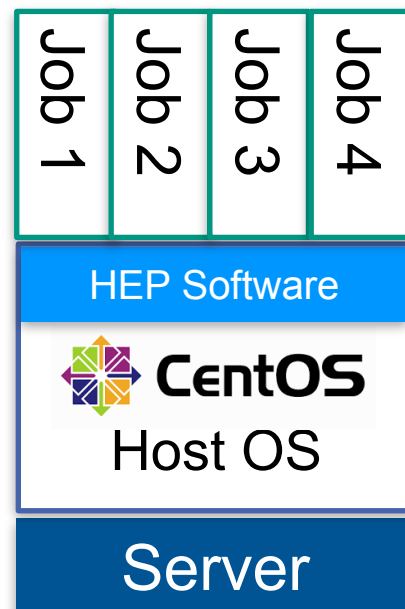


HEP-Job

# Opportunistic Resources and their Challenges

Very different to the traditional HEP environment:



HEP-Job



HPC-Job

Manuel Giffels

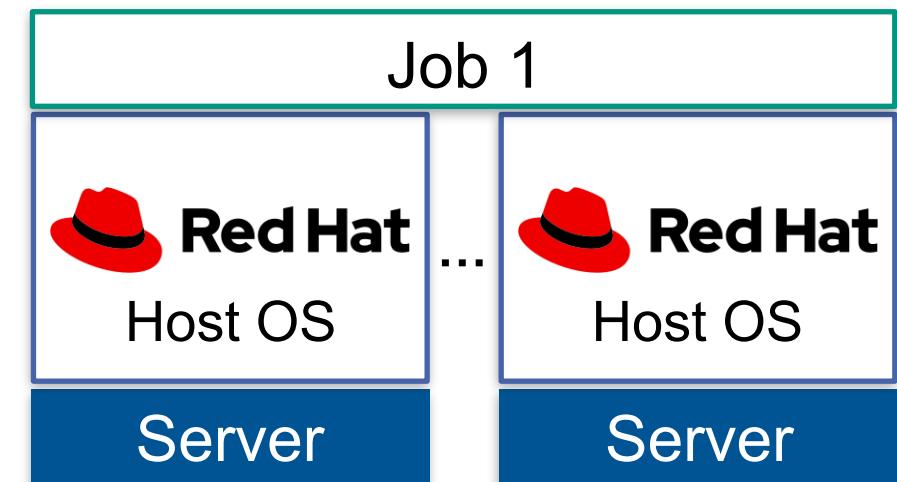# Opportunistic Resources and their Challenges

Very different to the traditional HEP environment:



HEP-Job



Virtualisation/Container



HPC-Job

# Opportunistic Resources and their Challenges
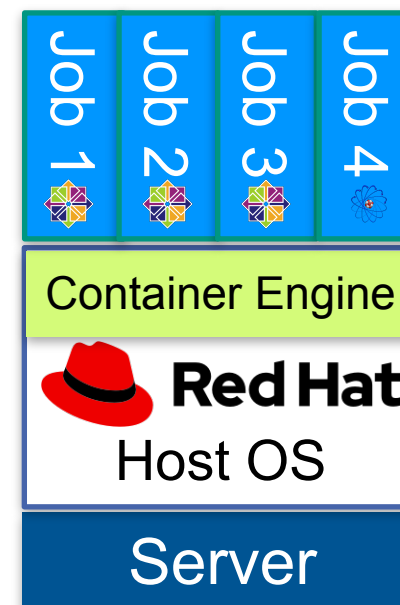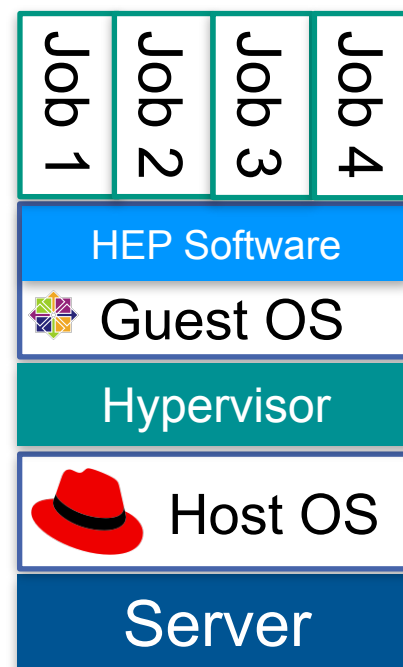
Very different to the traditional HEP environment:



HEP-Job



Virtualisation/Container



HPC-Job



Job 1 | Job 2 | Job 3 | Job 4

HEP Software

**CentOS**
Host OS

Server

# Opportunistic Resources and their Challenges
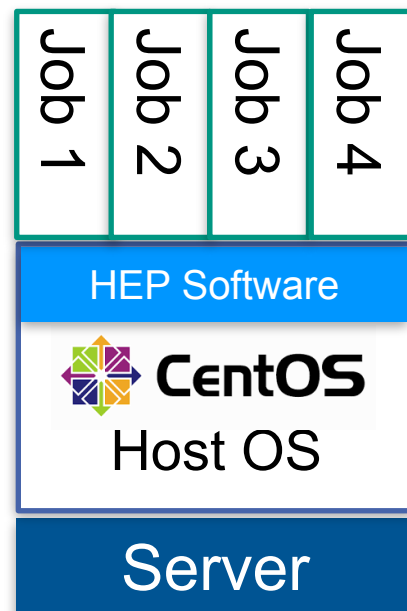
Very different to the traditional HEP environment:



HEP-Job



Virtualisation/Container



HPC-Job

Manuel Giffels

# Opportunistic Resources and their Challenges

Very different to the traditional HEP environment:



HEP-Job



Virtualisation/Container



HPC-Job

| Job 1 | Job 2 | Job 3 | Job 4 |
|---|---|---|---|

**HEP Software**

CentOS

Host OS

**Server**

| Job 1 | Job 2 | Job 3 | Job 4 |
|---|---|---|---|

HEP Software

Guest OS

Hypervisor

Host OS

**Server**

| Job 1 | Job 2 | Job 3 | Job 4 |
|---|---|---|---|

Container Engine

**Red Hat**

Host OS

**Server**

Job 1

**Red Hat**

Host OS

... **Red Hat**

Host OS

**Server**   **Server**
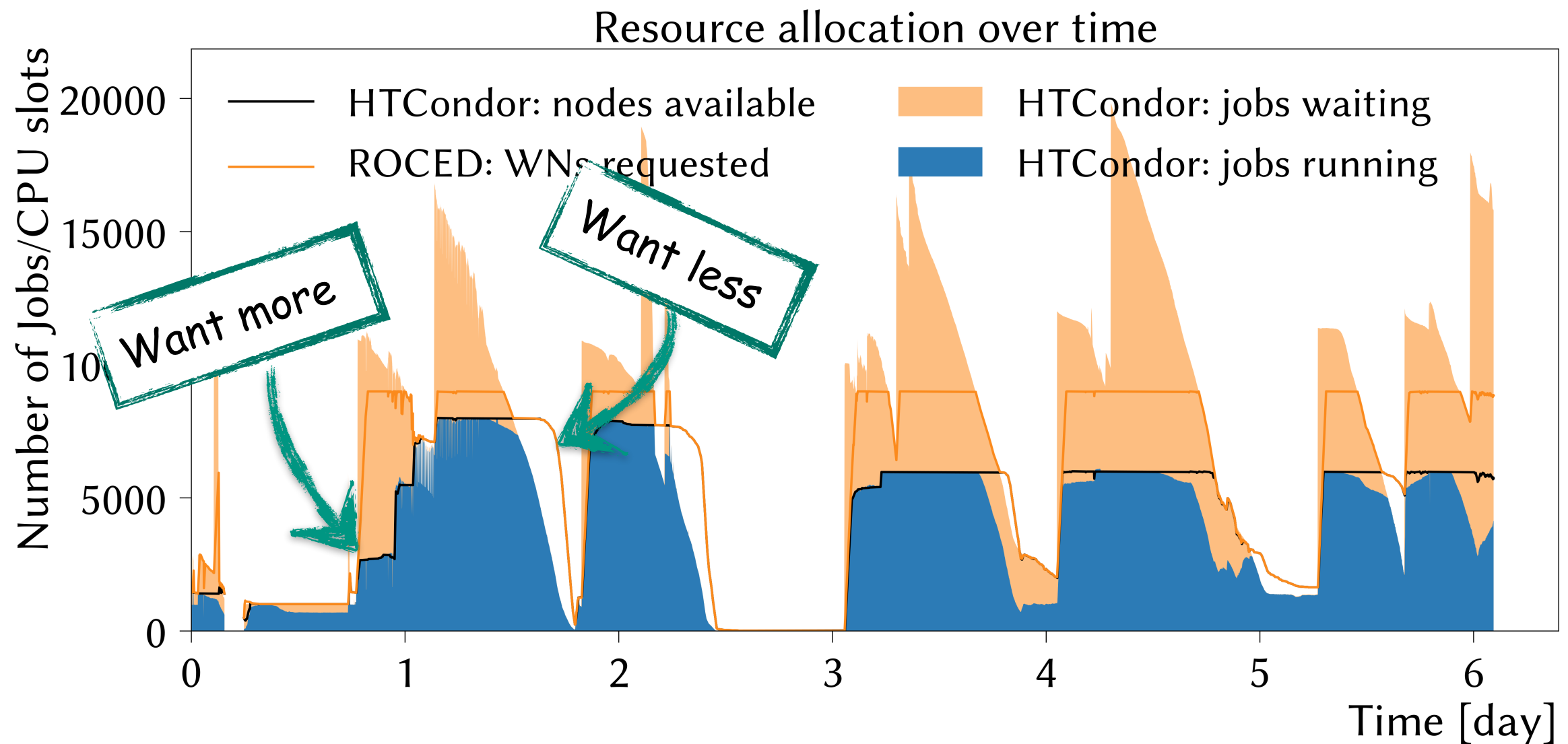
➜ Virtualisation and containerisation techniques

# Opportunistic Resources and their Challenges

- Temporary availability of opportunistic resources
- Varying demand for opportunistic resource



Resource allocation over time

Manuel Giffels

# Opportunistic Resources and their Challenges

- Temporary availability of opportunistic resources

- Varying demand for opportunistic resource

### Resource allocation over time



→ Dynamic integration and workflow management
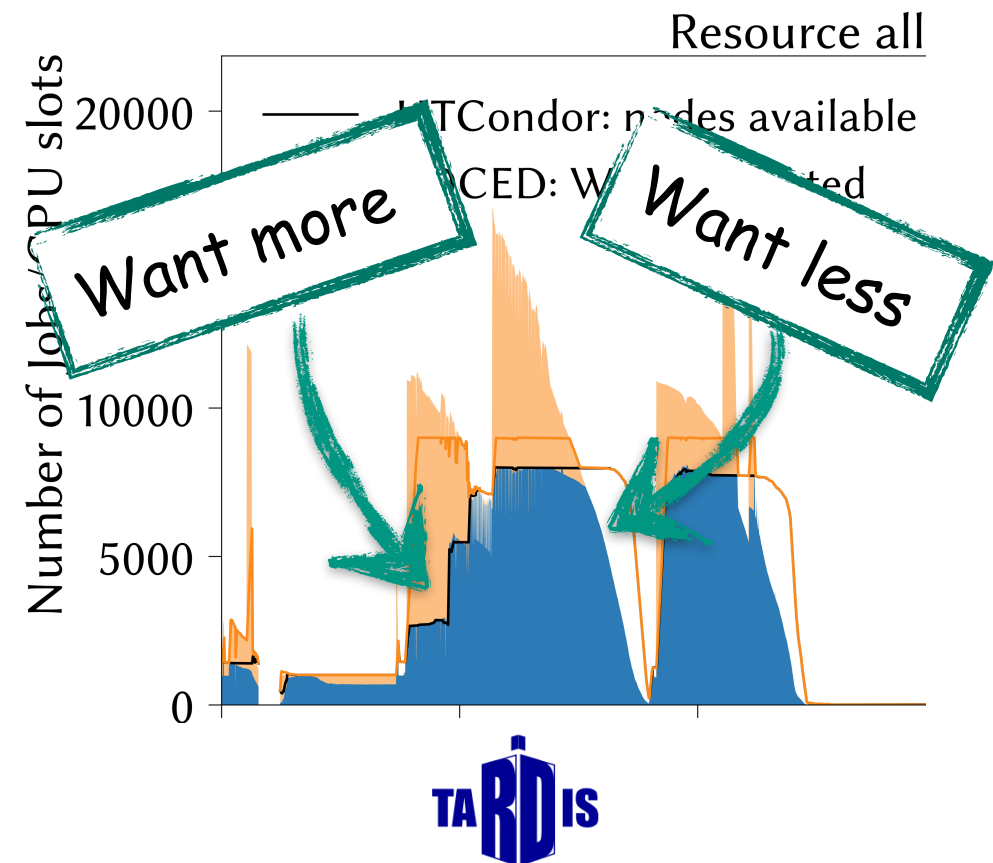
# COBalD/TARDIS Resource Manager

## Development at KIT:

- Following simplistic approach
  - Monitor resource utilisation
  - Increase well utilised resources
  - Reduce not used resources



Resource all

### COBalD

**COBalD – the Opportunistic Balancing Daemon**

COBalD is a lightweight framework to balance opportunistic resources: cloud bursting, container orchestration, allocation scaling and more. Its lightweight `model` for resources and their composition makes it easy to integrate custom resources and manage them at a large scale.

DOCUMENTATION | VERSION CONTROL | DOI

### TARDIS – Resourcemanager

**Transparent Adaptive Resource Dynamic Integration System**

Transparent Adaptive Resource Dynamic Integration System enables the dynamic integration of resources provided by different resource providers into one overlay batch system.

DOCUMENTATION | VERSION CONTROL | DOI

Manuel Giffels
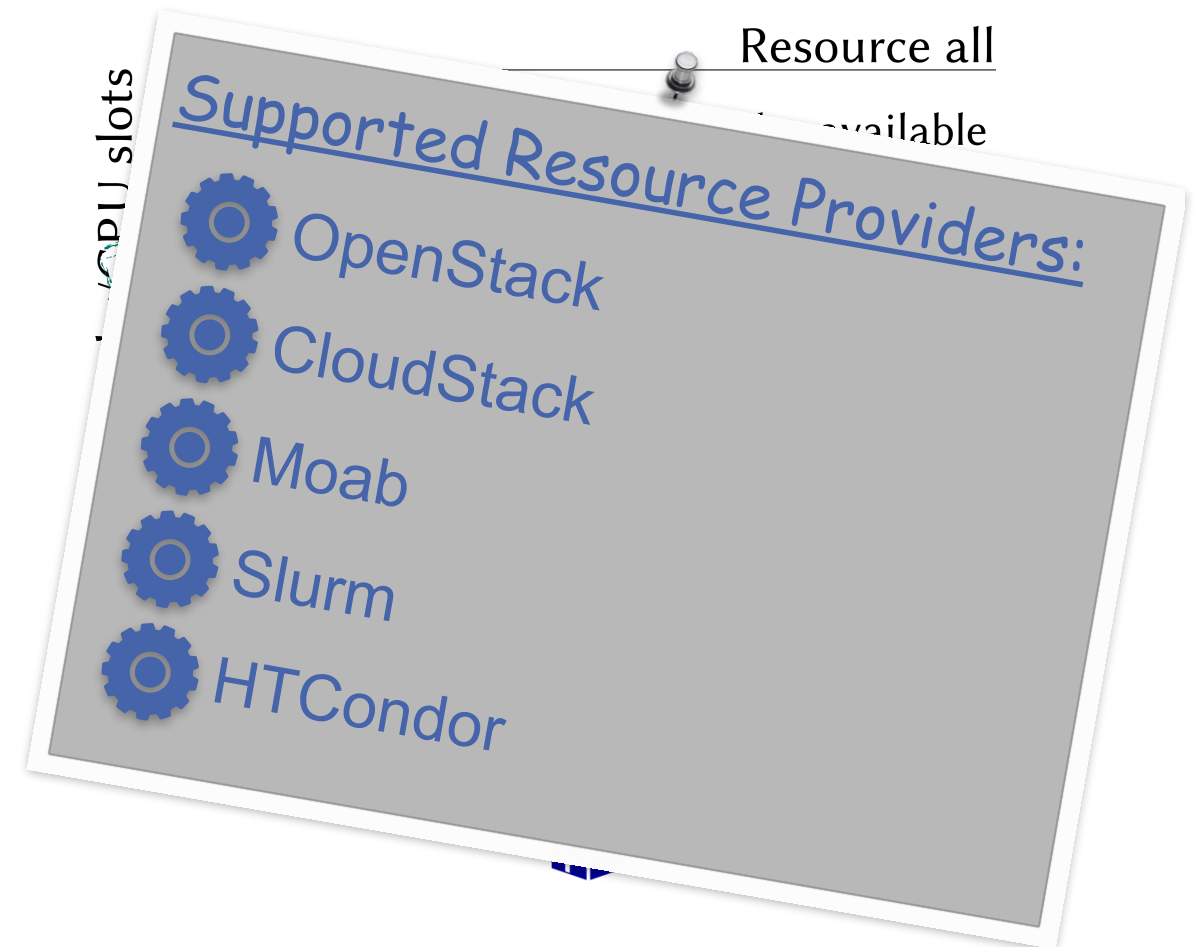
ETP / SCC

# COBalD/TARDIS Resource Manager

## Development at KIT:

- Following simplistic approach

  - Monitor resource utilisation

  - Increase well utilised resources

  - Reduce not used resources



**Resource all**
- HTCondor: nodes available
- CED: W... ...ted

*Want more*

*Want less*

## COBalD

**COBalD – the Opportunistic Balancing Daemon**

COBalD is a lightweight framework to balance opportunistic resources: cloud bursting, container orchestration, allocation scaling and more. Its lightweight `model` for resources and their composition makes it easy to integrate custom resources and manage them at a large scale.

| 📖 DOCUMENTATION | ⌥ VERSION CONTROL | ✸ DOI |

## TARDIS – Resourcemanager

**Transparent Adaptive Resource Dynamic Integration System**

Transparent Adaptive Resource Dynamic Integration System enables the dynamic integration of resources provided by different resource providers into one overlay batch system.

| 📖 DOCUMENTATION | ⌥ VERSION CONTROL | ✸ DOI |

# COBalD/TARDIS Resource Manager

## Development at KIT:

- Following simplistic approach

  - Monitor resource utilisation

  - Increase well utilised resources

  - Reduce not used resources

### Supported Resource Providers:
- OpenStack
- CloudStack
- Moab
- Slurm
- HTCondor

Resource all
available

CPU slots

## COBalD

### COBalD – the Opportunistic Balancing Daemon

COBalD is a lightweight framework to balance opportunistic resources: cloud bursting, container orchestration, allocation scaling and more. Its lightweight `model` for resources and their composition makes it easy to integrate custom resources and manage them at a large scale.

DOCUMENTATION | VERSION CONTROL | DOI

## TARDIS – Resourcemanager

### Transparent Adaptive Resource Dynamic Integration System

Transparent Adaptive Resource Dynamic Integration System enables the dynamic integration of resources provided by different resource providers into one overlay batch system.

DOCUMENTATION | VERSION CONTROL | DOI

# Opportunistic Resources and their Challenges

**Where to send my jobs?**

**Which resources are available?**

Site A

Site B

Site C

# Opportunistic Resources and their Challenges

- Each opportunistic resource is different (very heterogenous system)

**Where to send my jobs?**

**Which resources are available?**

Site A

Site B

Site C

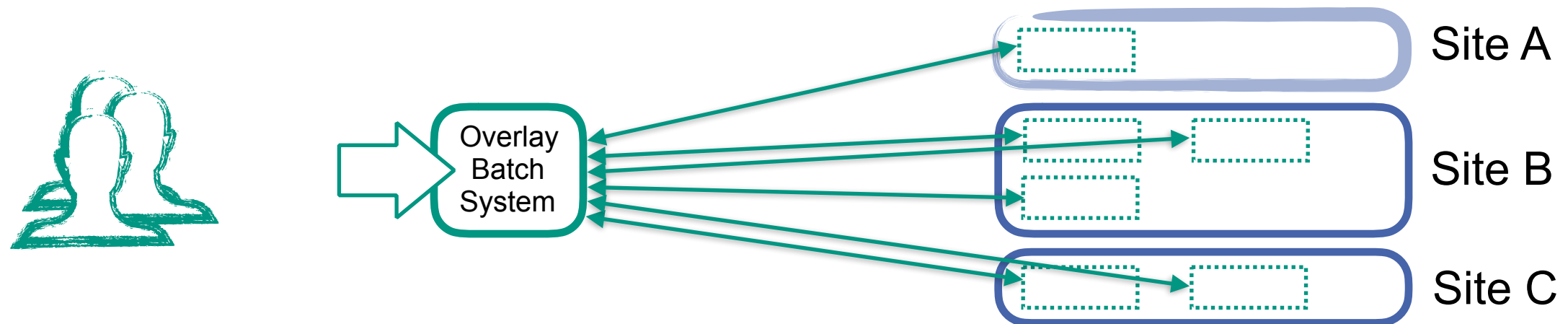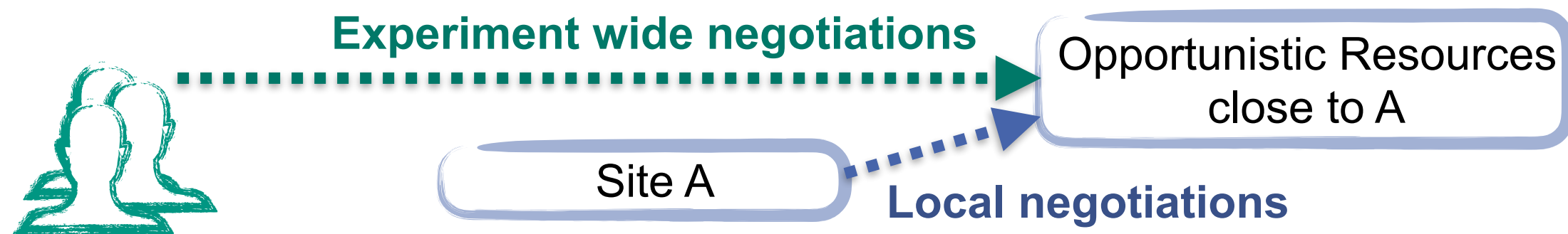# Opportunistic Resources and their Challenges

- Each opportunistic resource is different (very heterogenous system)
  - ➡ Hide complexity from users and computing operations of experiments



➜ Transparent integration of resources needed

# Opportunistic Resources and their Challenges

- Each opportunistic resource is different (very heterogenous system)
  - ➡ Hide complexity from users and computing operations of experiments

Site A

Overlay Batch System

Site B

Site C

  - ➡ Local/regional negotiations with resource providers more promising than experiment wide negotiations

**Experiment wide negotiations**

Opportunistic Resources close to A

Site A

**Local negotiations**

➜ Transparent integration of resources needed

Manuel Giffels

# The Entire Picture

**Local Site**

Users

Job submission

Overlay Batch System (OBS)

HTCondor
High Throughput Computing

TARDIS

increase resources

utilization

decrease resources

usage monitoring

**External Site**

Access Point

schedule request and start resource

Resource Pool

VM

batch job

container

usage monitoring

integrate into OBS

jobflow

One TARDIS instance per site

# **Use case:** Backfilling of HPC Resources

Manuel Giffels

# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)
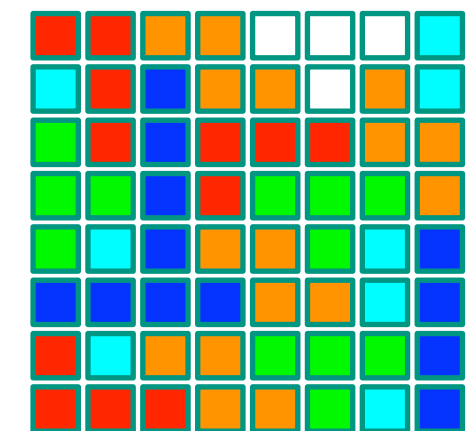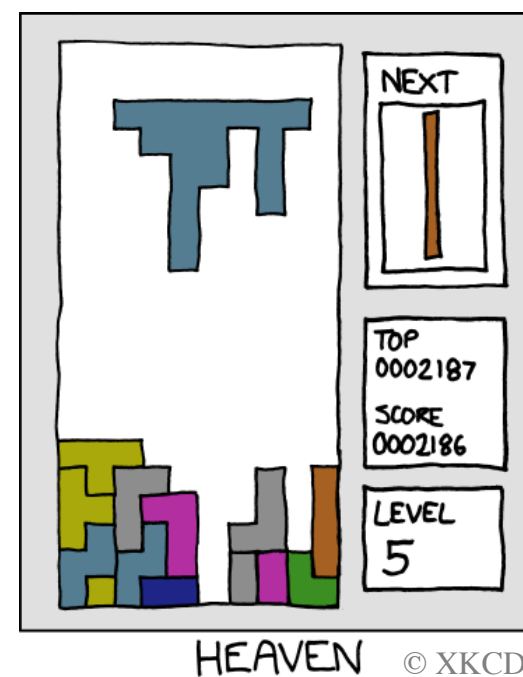
# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

Utilization=0.875

NEXT

TOP
000000
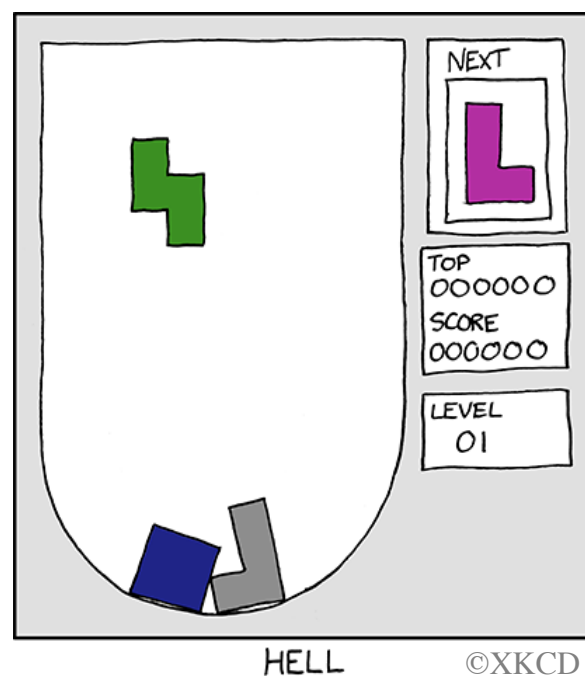
SCORE
000000

LEVEL
01

HELL ©XKCD

Manuel Giffels

ETP / SCC

# Use case: Backfilling of HPC Resources



- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization
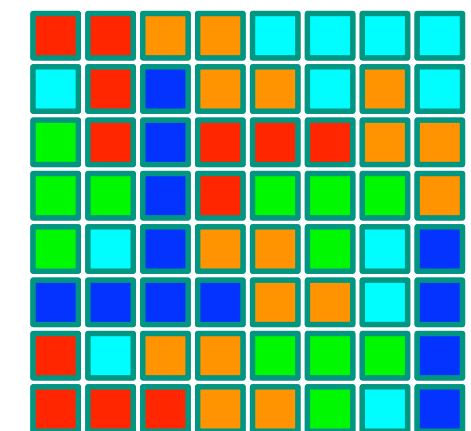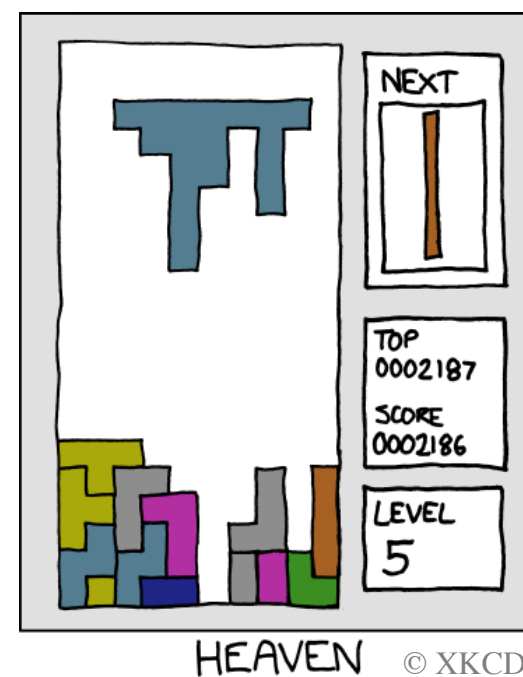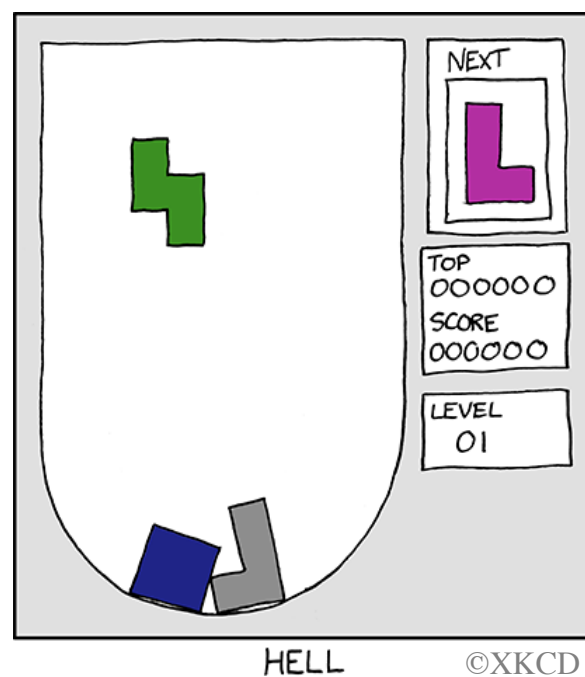
Utilization=0.875



HELL ©XKCD



HEAVEN © XKCD

HEP Job

# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization

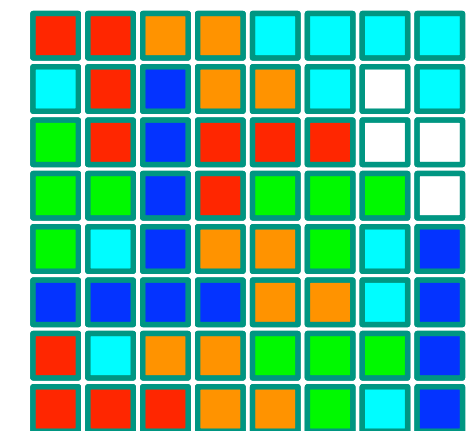Utilization=0.875
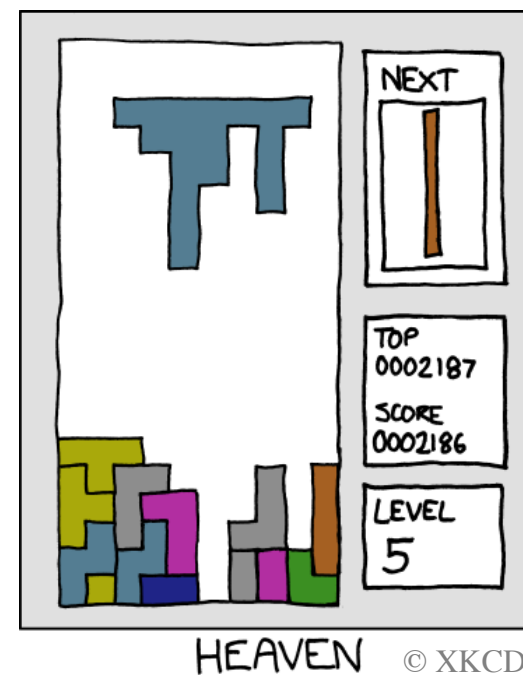
Utilization=1.0

HELL ©XKCD

HEAVEN © XKCD

☐ HEP Job

# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization



Utilization=0.875

Utilization=1.0
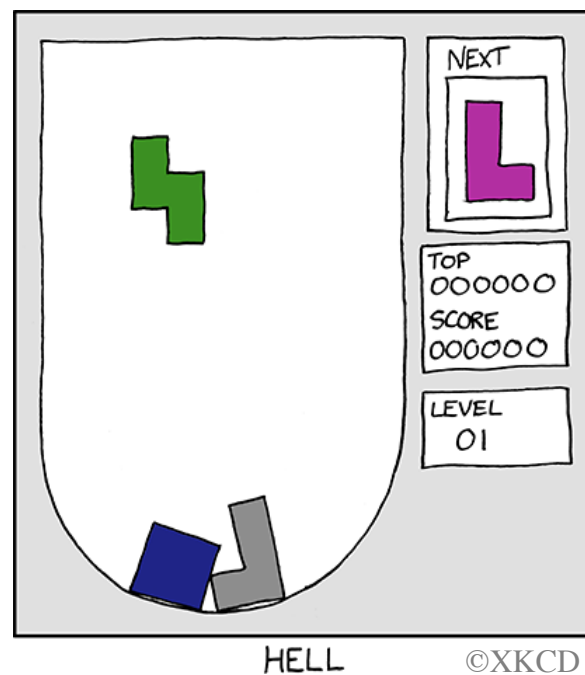
HELL ©XKCD

HEAVEN © XKCD

☐ HEP Job

Manuel Giffels

# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization



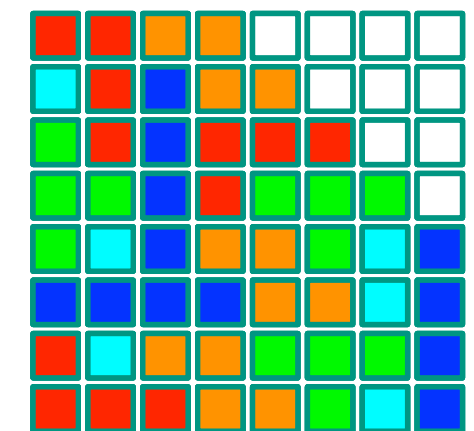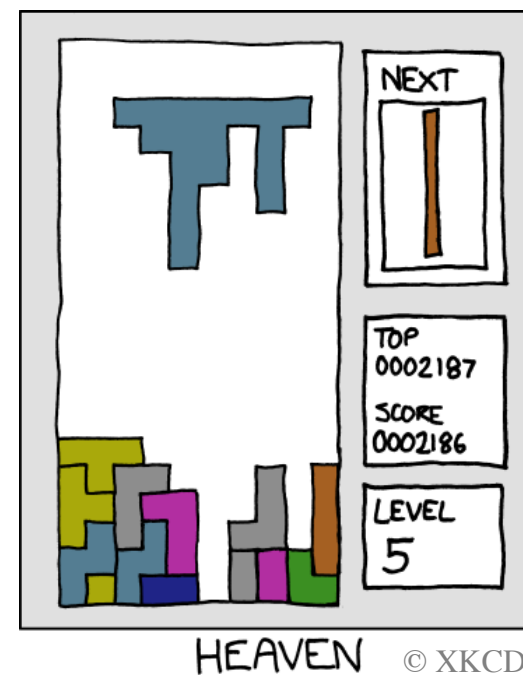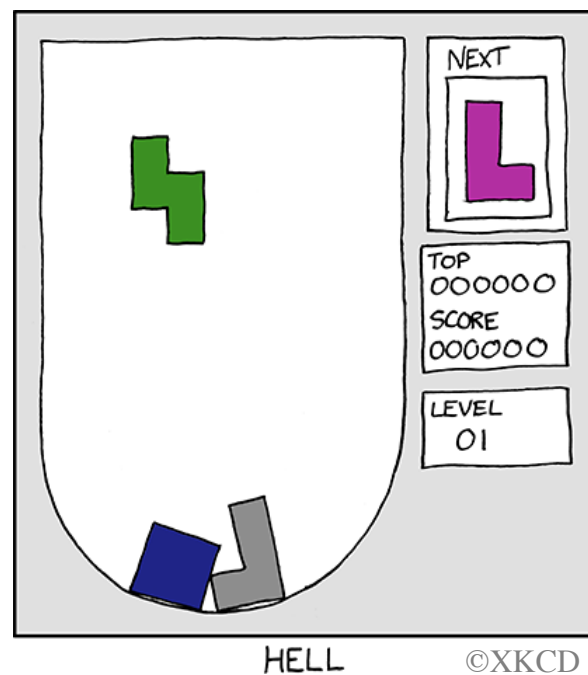Utilization=0.875

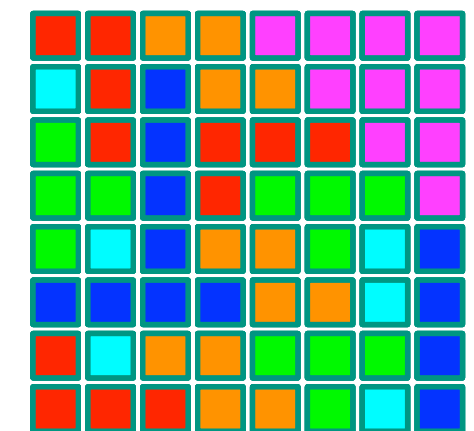Utilization=1.0

HEP Job

# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization



Utilization=0.875

Utilization=1.0

HEP Job

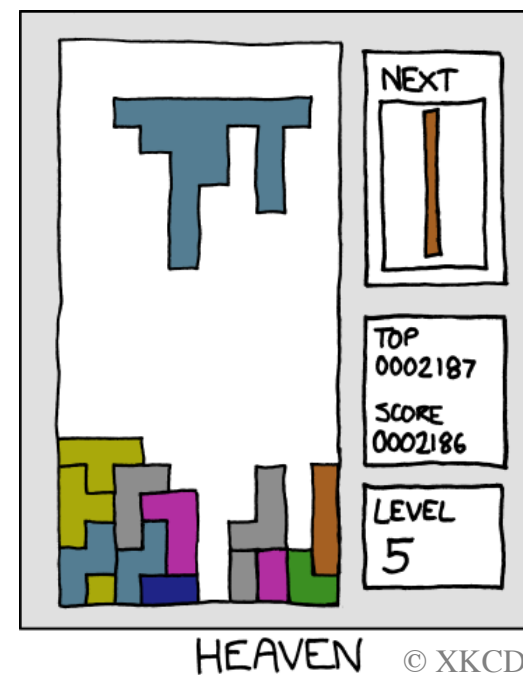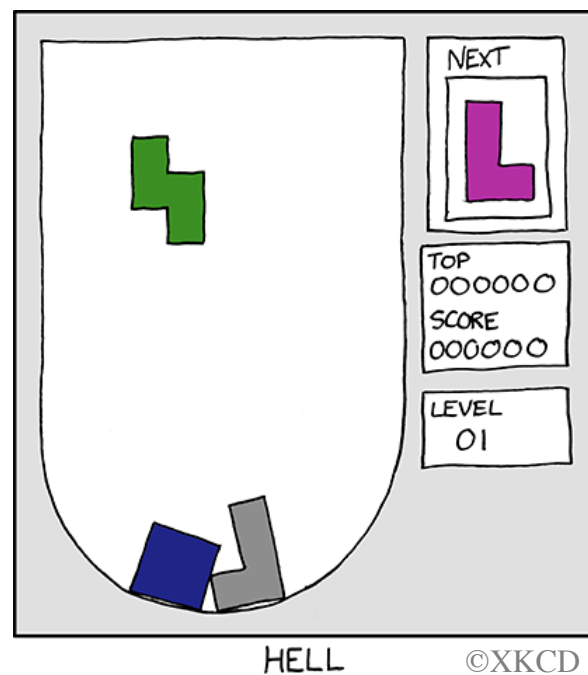HELL ©XKCD

HEAVEN © XKCD

# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization

- Enable pre-emption to free resources if needed



Utilization=0.875

Utilization=1.0


HELL ©XKCD


HEAVEN © XKCD

HEP Job

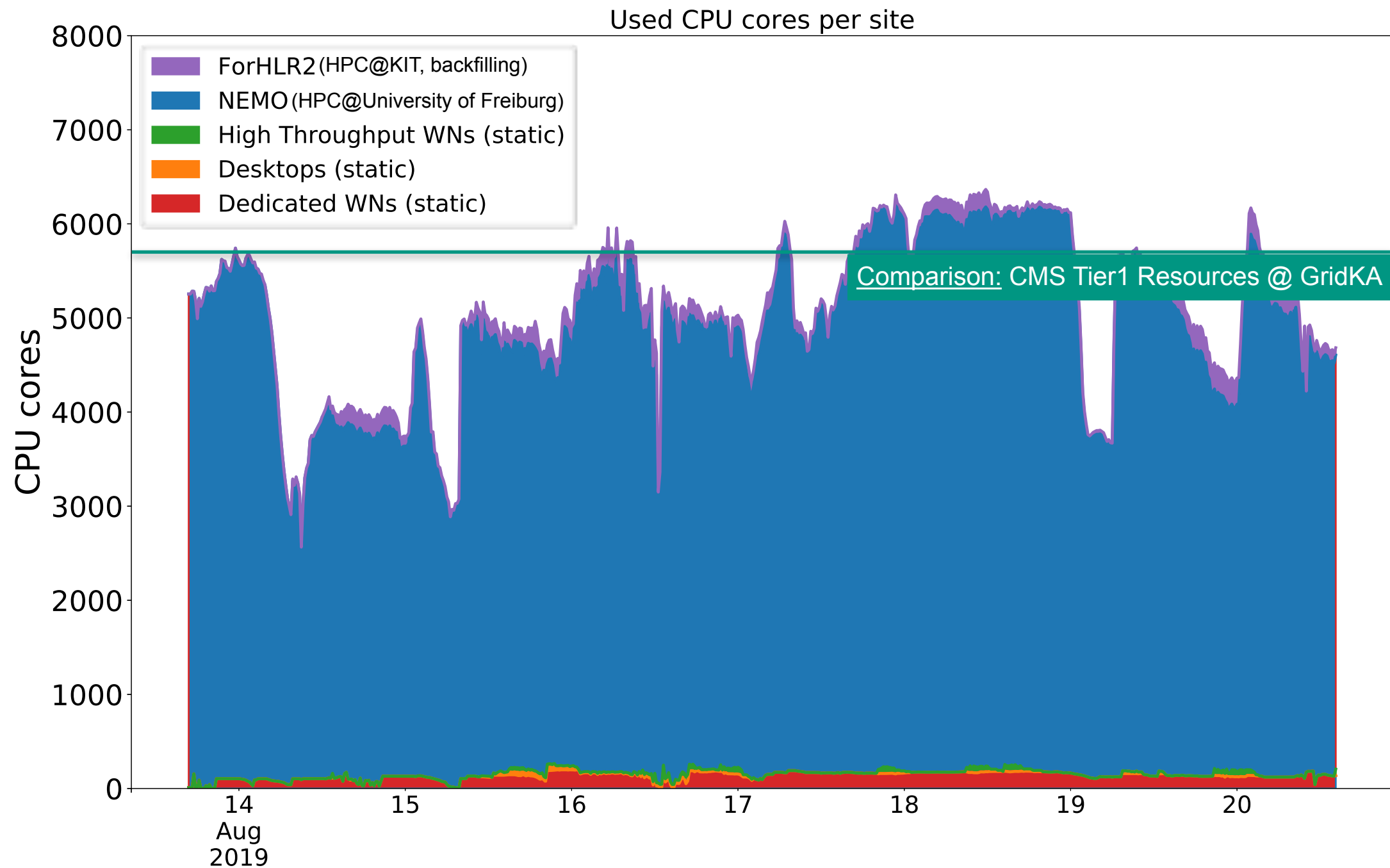# Use case: Backfilling of HPC Resources

- HPC schedulers are optimized to schedule large scientific simulation and calculation workflows (Many cores/nodes)

- Leads inevitably to unused resources due to draining of machines and fragmentation

- Backfilling with small and short running HEP jobs increases the overall resource utilization

- Enable pre-emption to free resources if needed



Utilization=0.875
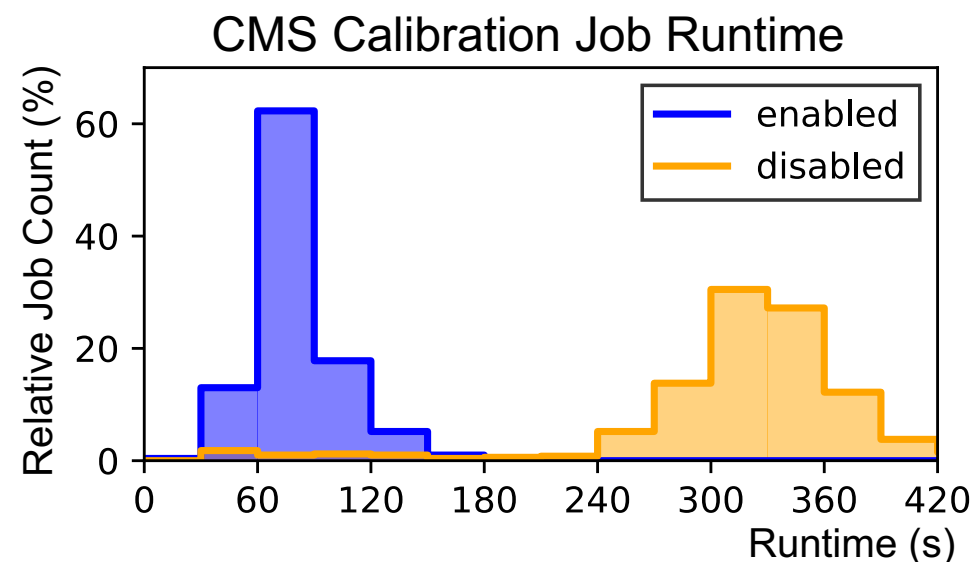
Utilization=1.0

HEP Job

NEXT

TOP
000000

SCORE
000000

LEVEL
01

HELL ©XKCD

NEXT

TOP
0002187

SCORE
0002186

LEVEL
5

HEAVEN © XKCD

# Opportunistic "TIER 1" for a Day



**Used CPU cores per site**

Legend:
- ForHLR2 (HPC@KIT, backfilling)
- NEMO (HPC@University of Freiburg)
- High Throughput WNs (static)
- Desktops (static)
- Dedicated WNs (static)

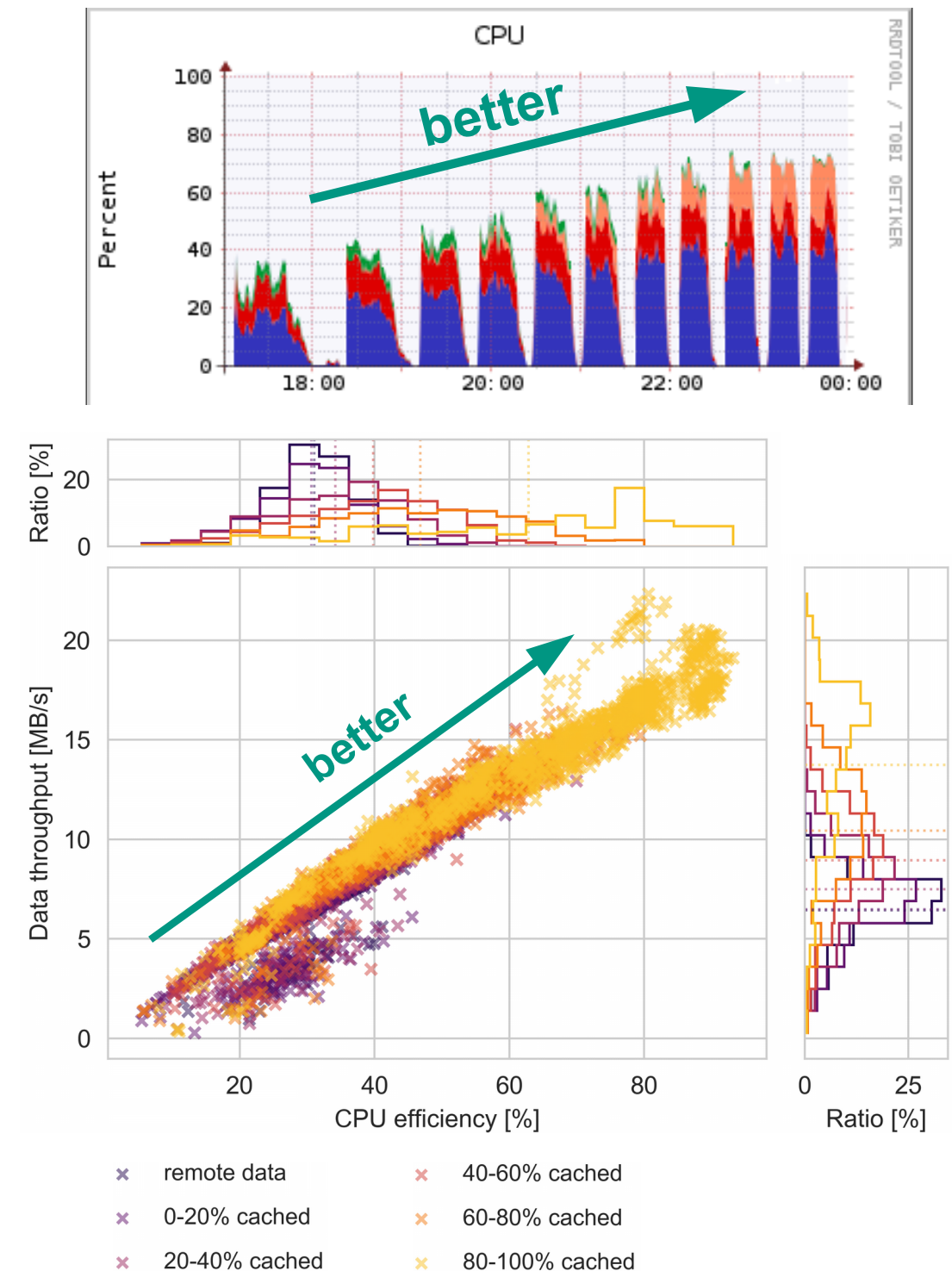Comparison: CMS Tier1 Resources @ GridKA

➜ Production ready software at scale in stable operation!

# Data Caches

- CPU efficiency/runtime strongly depends on available I/O bandwidth
- Opportunistic sites have potentially slower WAN connections
- Many opportunistic sites offering access to fast storage systems (HPC, S3, etc.)
- Utilise data caches to enable opportunistic sites for recurrent I/O-intense workflows
- Transparent data access was also a hot topic in the Helix Nebula Science Cloud



CMS User Analysis



CMS Calibration Job Runtime
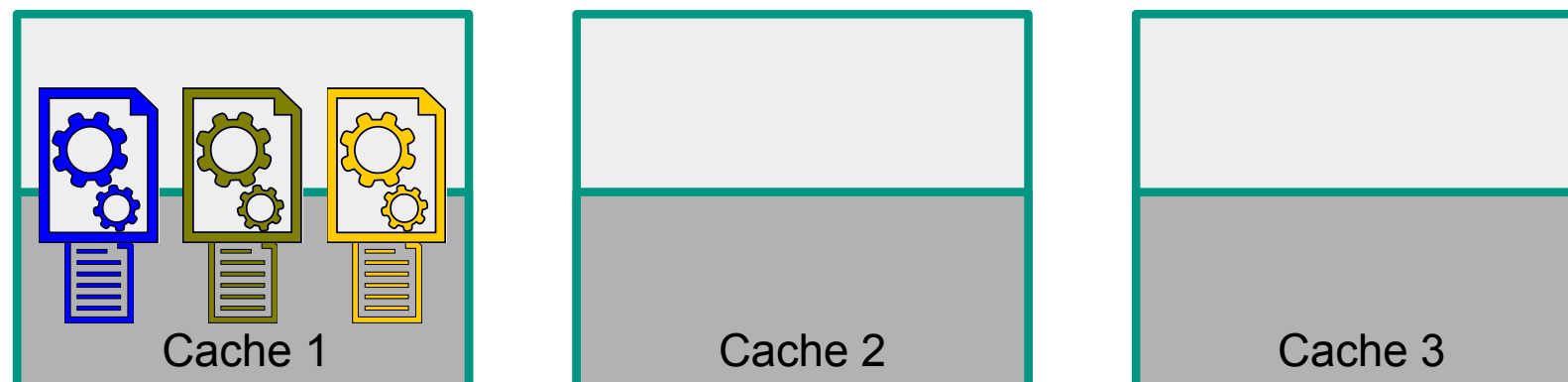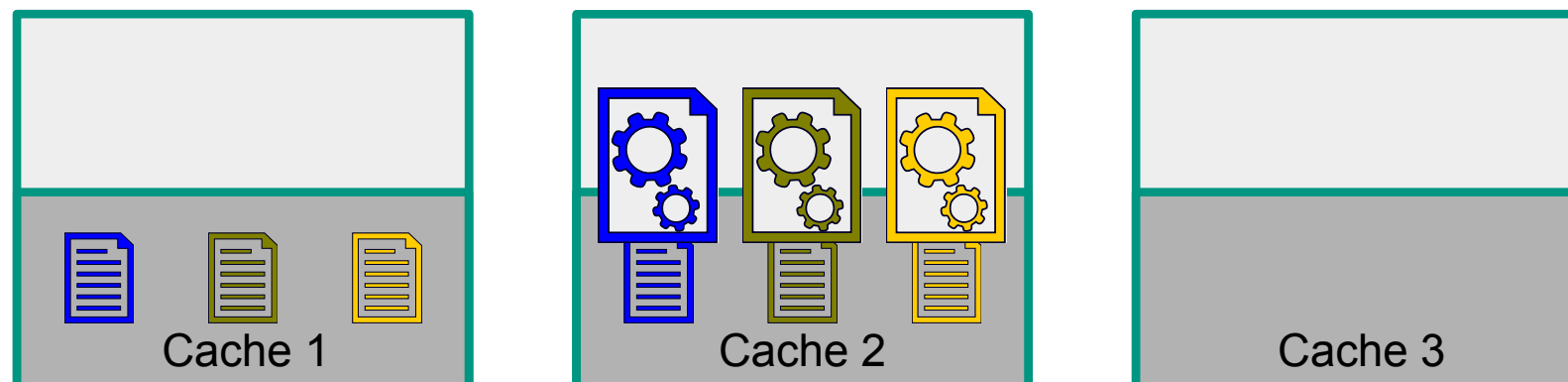
Manuel Giffels

ETP / SCC

# Caching

- Common solution for repeated access to the same data
- Cache data as close as possible to the CPU

Suitable for:

- HEP workflows that process the same data frequently
- CPU resources without permanent storage

Problematic on distributed resources with multiple caches



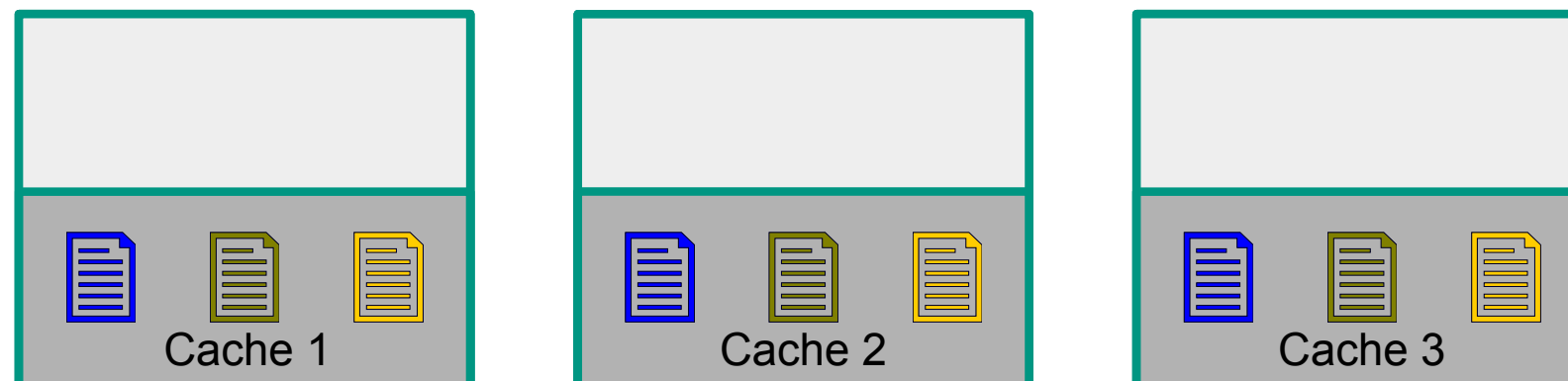Cache 1          Cache 2          Cache 3

# Caching

- Common solution for repeated access to the same data
- Cache data as close as possible to the CPU

Suitable for:

- HEP workflows that process the same data frequently
- CPU resources without permanent storage

Problematic on distributed resources with multiple caches



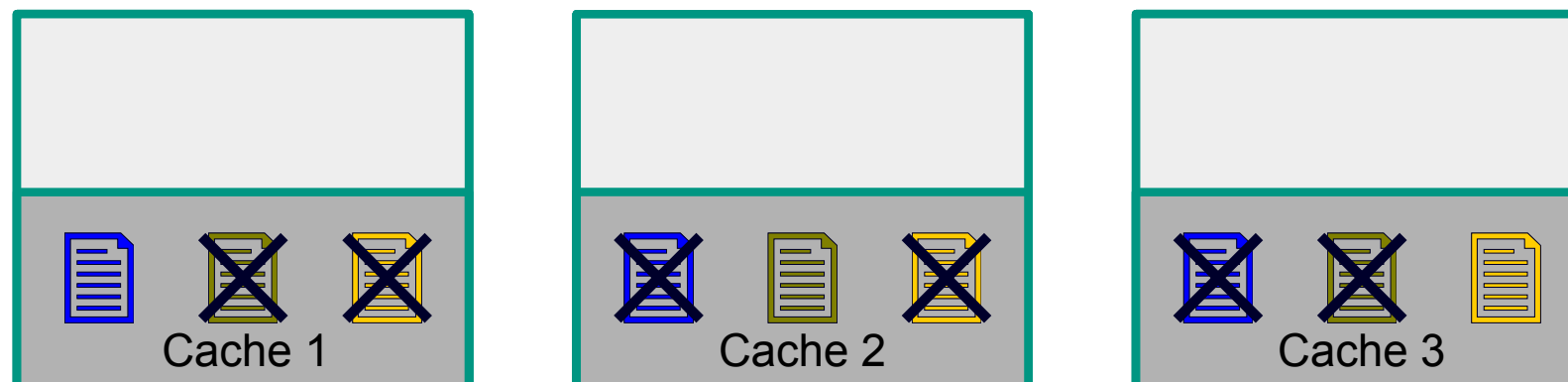Cache 1     Cache 2     Cache 3

# Caching

- Common solution for repeated access to the same data
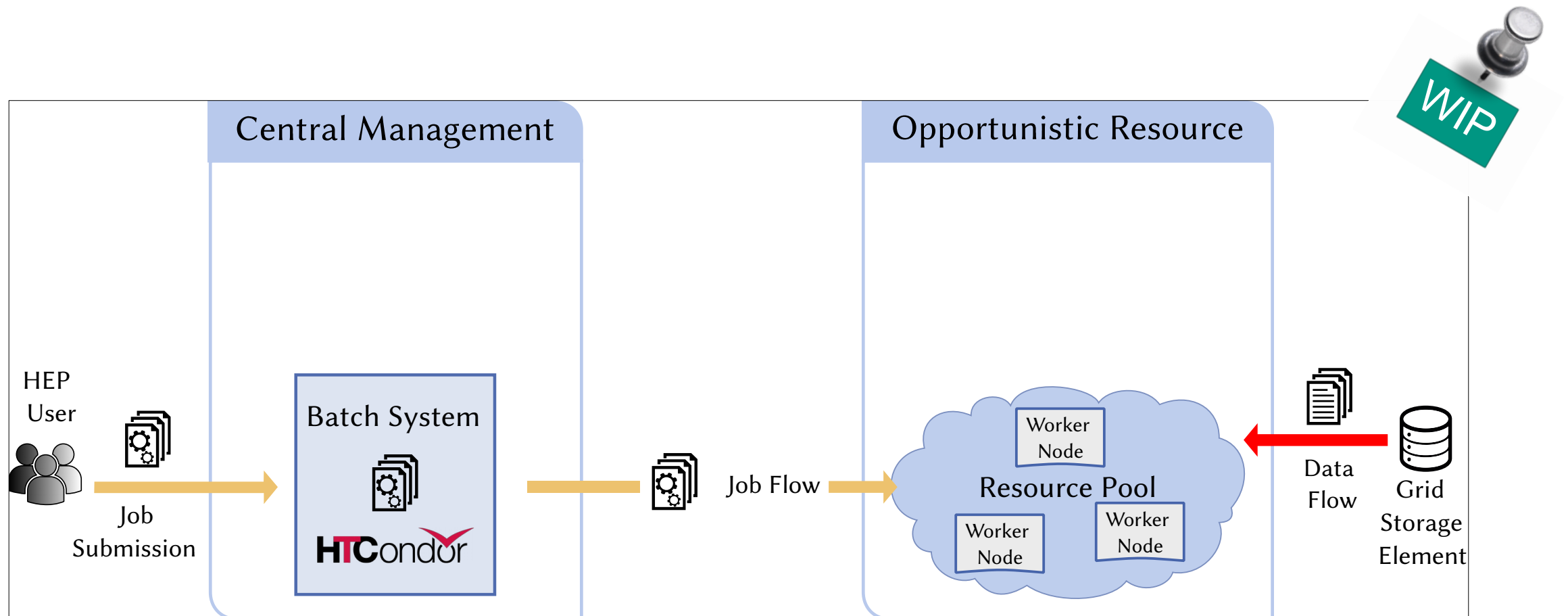- Cache data as close as possible to the CPU

Suitable for:

- HEP workflows that process the same data frequently
- CPU resources without permanent storage

Problematic on distributed resources with multiple caches

# Caching

- Common solution for repeated access to the same data
- Cache data as close as possible to the CPU

Suitable for:

- HEP workflows that process the same data frequently
- CPU resources without permanent storage

Problematic on distributed resources with multiple caches



Cache 1    Cache 2    Cache 3

Waste of storage capacity due to replication of data!
Caches must be coordinated!

# Coordinated Caching



Central Management

Opportunistic Resource

WIP

HEP User
Job Submission

Batch System
HTCondor

Job Flow

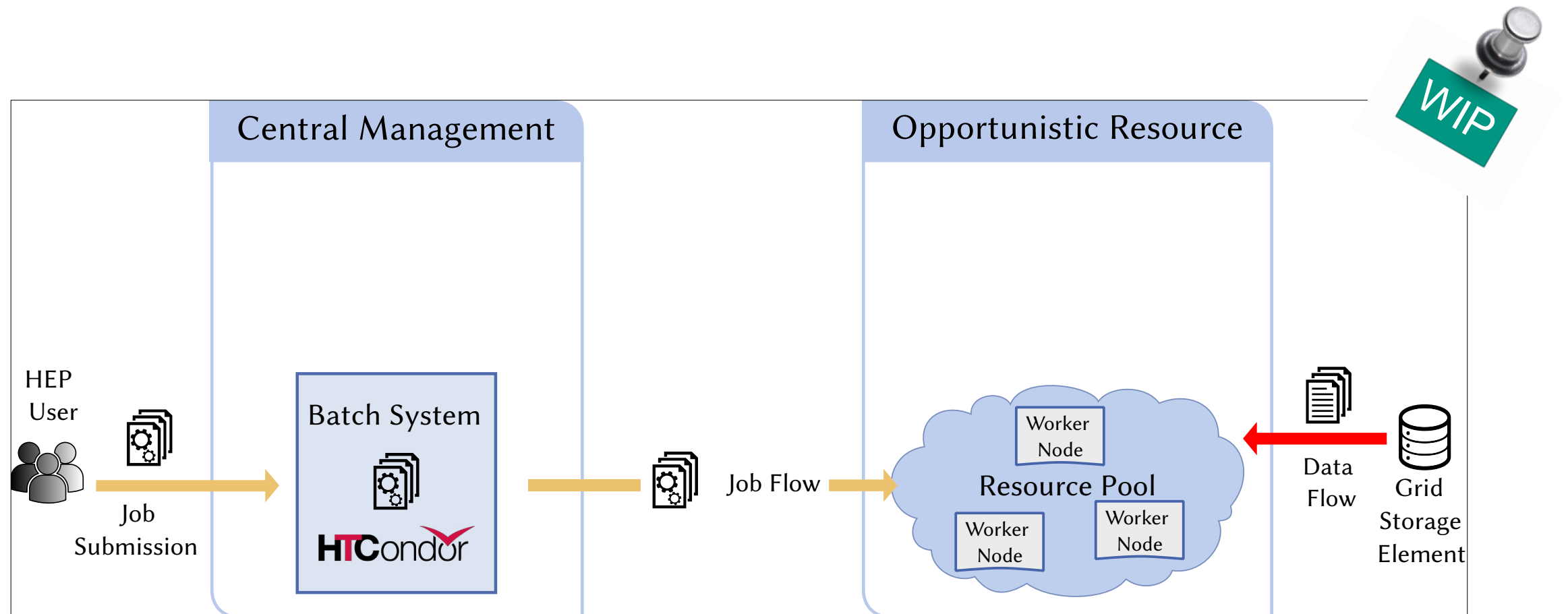Resource Pool
Worker Node
Worker Node
Worker Node

Data Flow

Grid Storage Element

Longterm experiences @ KIT:
"Data Locality via Coordinated Caching for Distributed Processing",
M. Fischer et al., J. Phys.: Conf. Ser.762 012011 (2016)
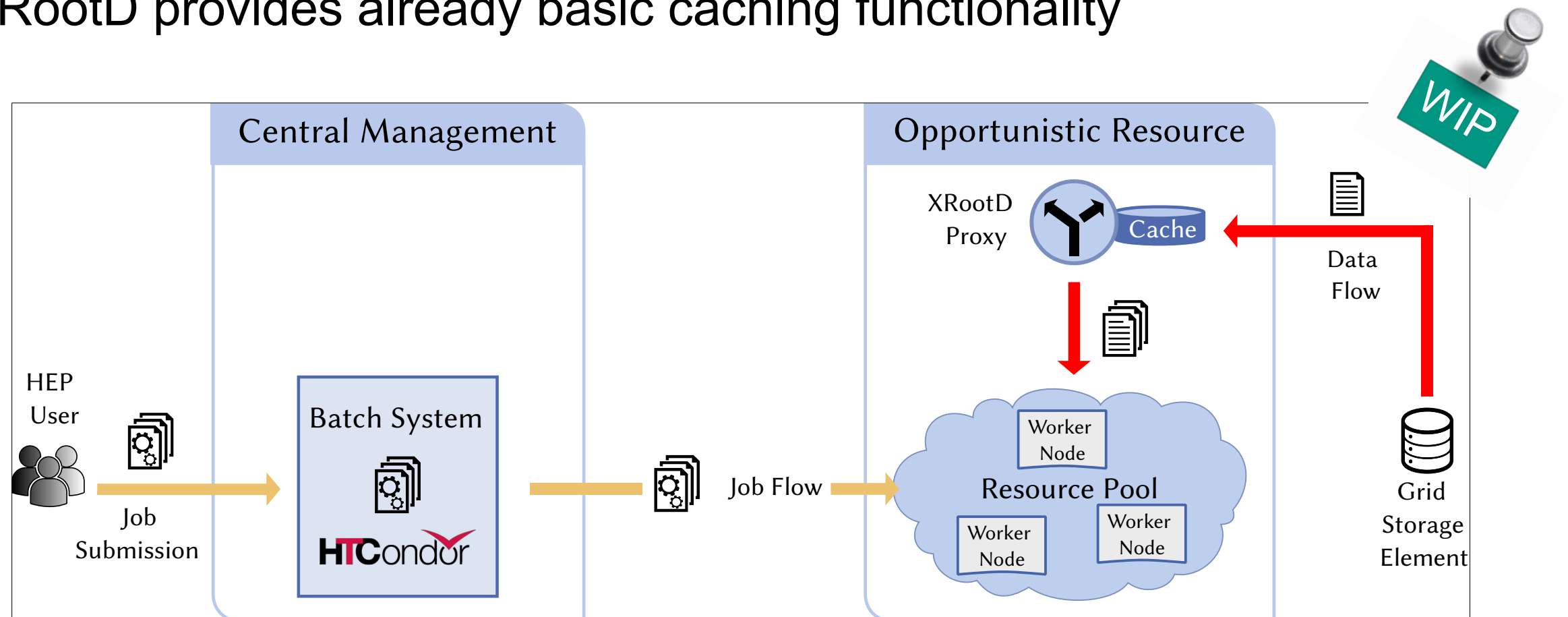
# Coordinated Caching

■ HTCondor schedules jobs to resources



Longterm experiences @ KIT:
"Data Locality via Coordinated Caching for Distributed Processing",
M. Fischer et al., J. Phys.: Conf. Ser.762 012011 (2016)
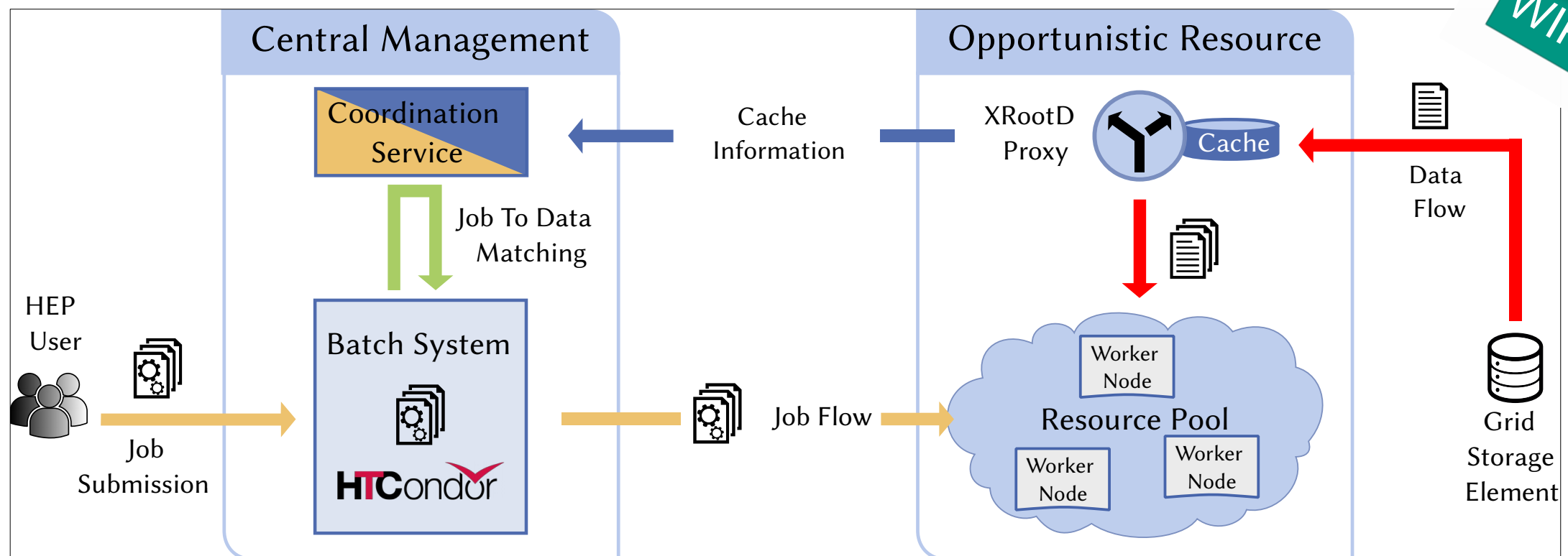
# Coordinated Caching

- HTCondor schedules jobs to resources
- XRootD provides already basic caching functionality



Longterm experiences @ KIT:
"Data Locality via Coordinated Caching for Distributed Processing",
M. Fischer et al., J. Phys.: Conf. Ser.762 012011 (2016)

# Coordinated Caching

- HTCondor schedules jobs to resources

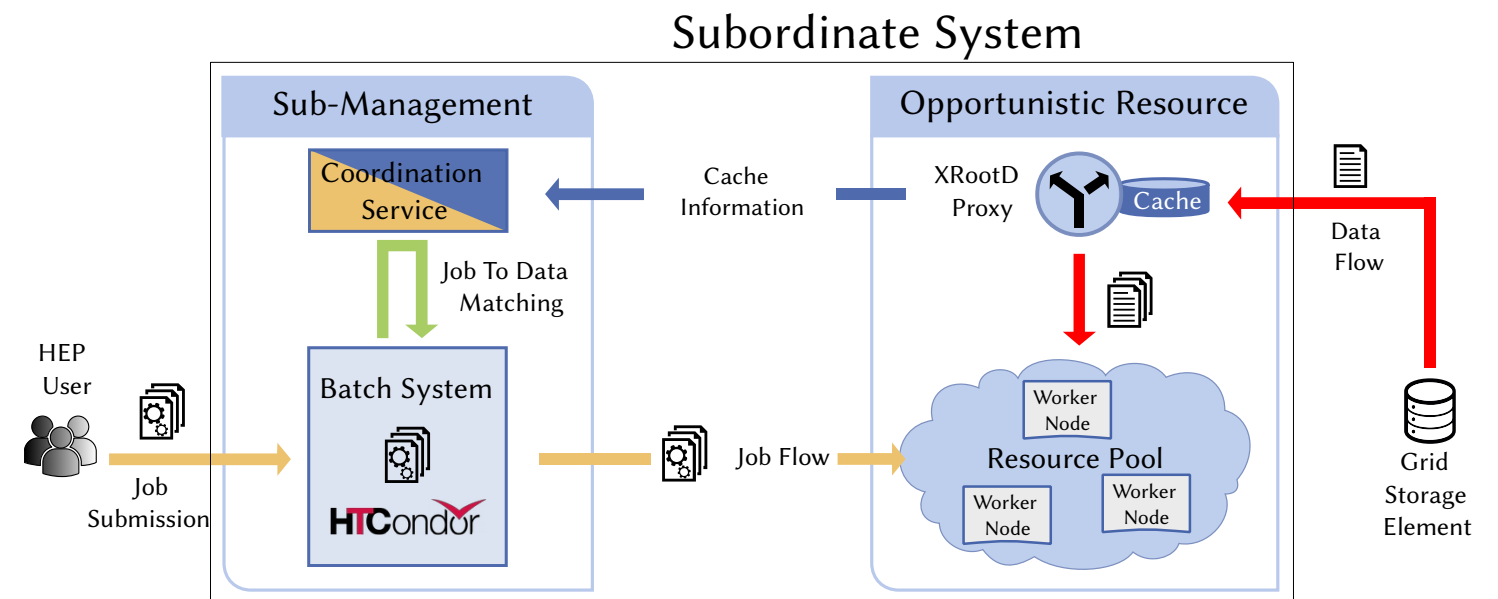- XRootD provides already basic caching functionality



- NaviX coordination service currently under development at KIT
  - Implicit data placement via job scheduling
  - Schedule jobs to cached data

Longterm experiences @ KIT:
"Data Locality via Coordinated Caching for Distributed Processing",
M. Fischer et al., J. Phys.: Conf. Ser.762 012011 (2016)
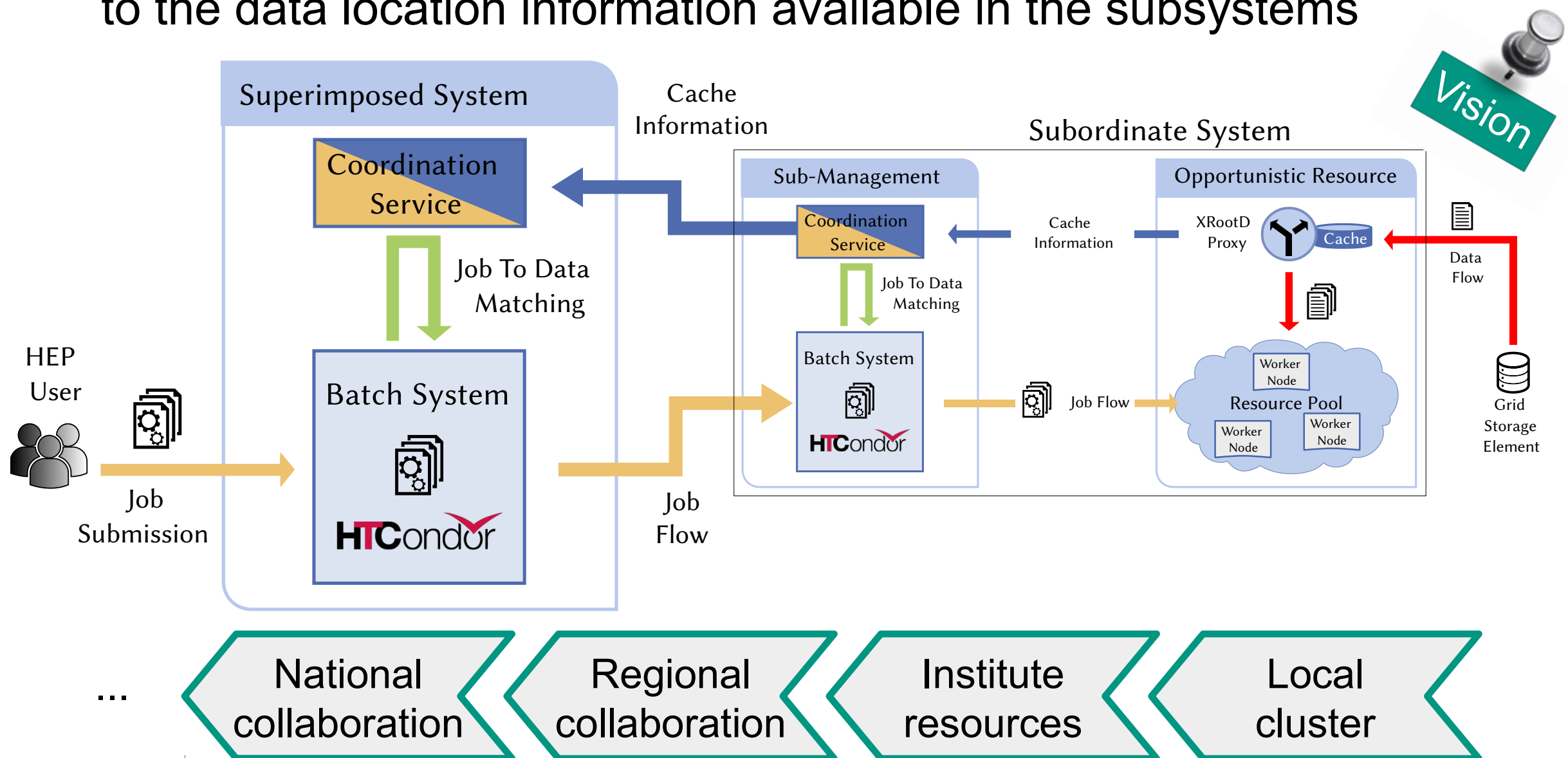
# Coordinated Caching

<u>Vision:</u> Build a hierarchical system of local, site and national caches

# Scalability by Design

Vision: Build a hierarchical system of local, site and national caches

- XRootD and HTCondor take care of hierarchical upscaling
- Job-to-cache coordination can be performed on all levels with regard to the data location information available in the subsystems
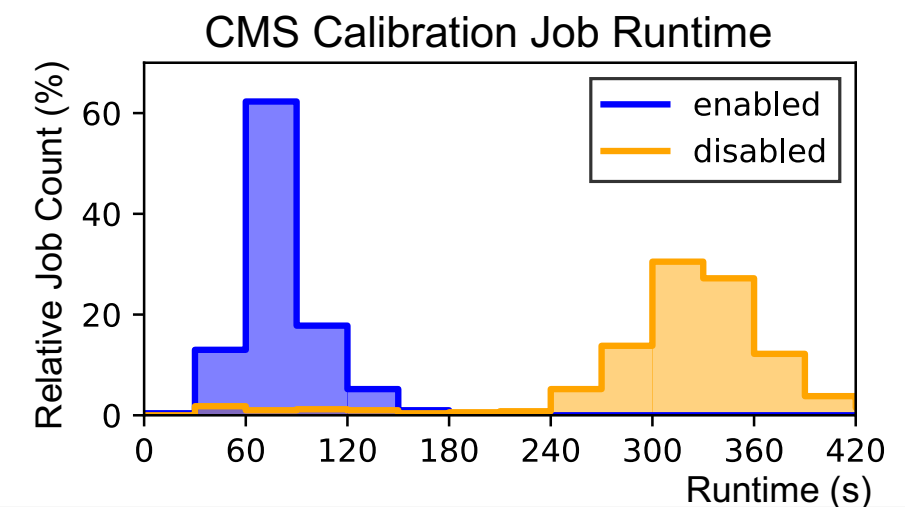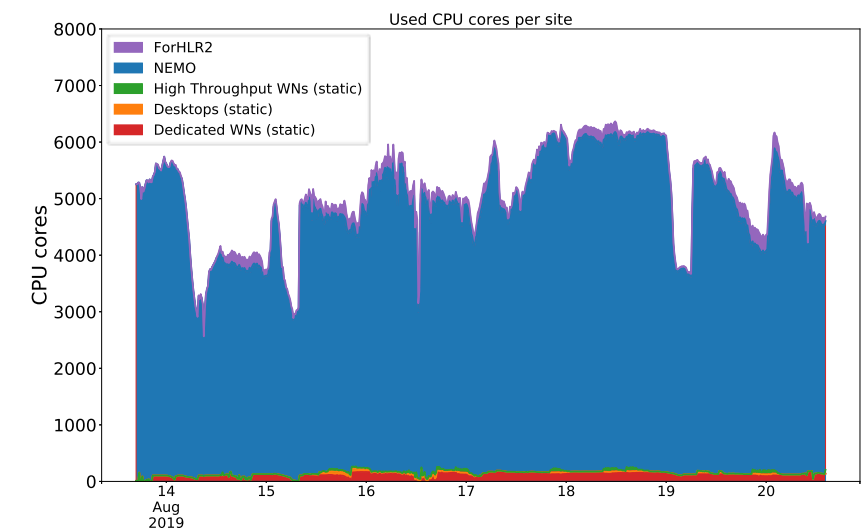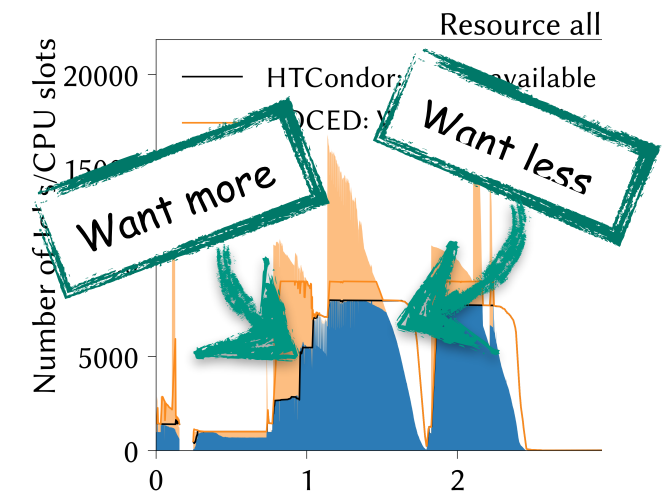
# Conclusion

**Dynamic on-demand provisioning of resources**

- COBalD/TARDIS resource manager developed at KIT

- Enables transparent and dynamic on-demand provisioning of opportunistic resources

- Enables backfilling of HPC resources
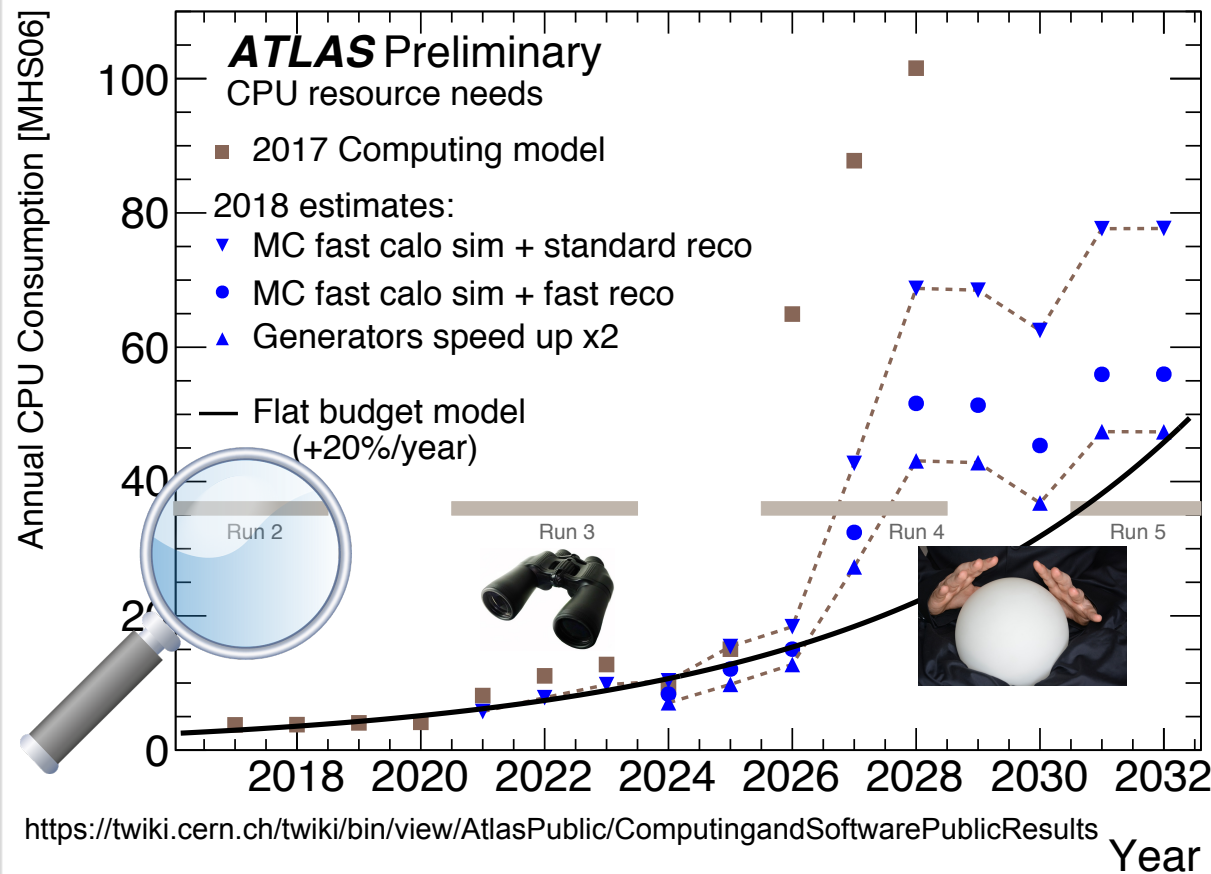
- Production ready software at scale

COBalD

TaRDIS

**Distributed coordinated caching**

- NaviX coordination service in development at KIT

- Working towards a scalable caching solution based upon HTCondor and XRootD
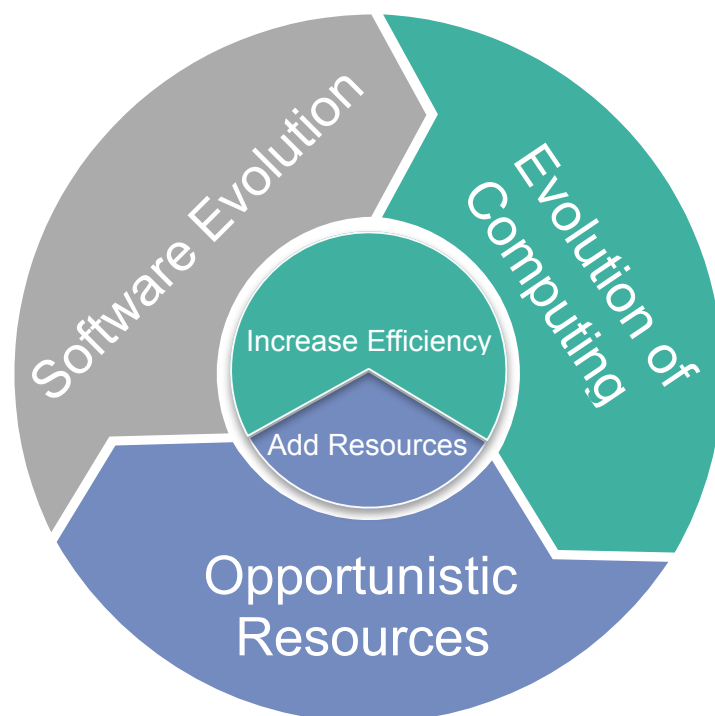
- First prototype available

# BACKUP

Manuel Giffels

# HEP Computing Challenges Ahead



https://twiki.cern.ch/twiki/bin/view/AtlasPublic/ComputingandSoftwarePublicResults

- ATLAS/CMS CPU resource estimates
- Assuming flat budget and 20% technology advance per year
- CPU and storage resource shortfall between needs and technology in 2027
- ⊕ Politically endorsed integration of cloud and high performance computing resources (US)
- → (R)evolution of computing model is required to master future challenges



New interesting research topics:

- Exploitation of modern technologies
- Improvement of algorithms and utilization of ML
- Dynamic integration of opportunistic resources (HPC, cloud, volunteer computing)
- Data lakes and data caching technologies

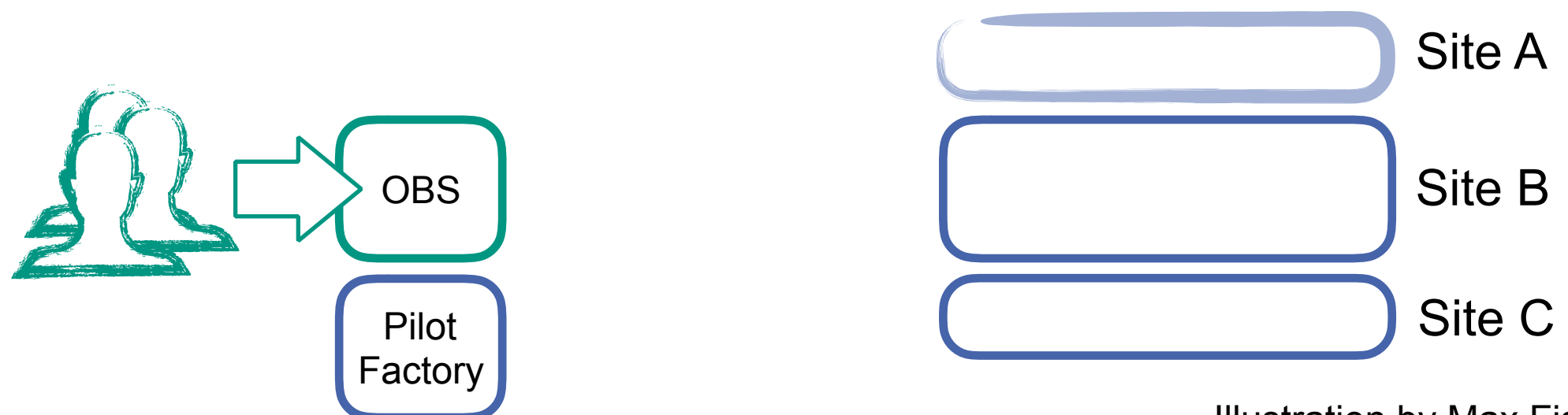# Transparent Integration: Overlay Batch System (OBS)

OBS

Pilot
Factory

Site A

Site B

Site C

Illustration by Max Fischer (KIT)

Manuel Giffels

# Transparent Integration: Overlay Batch System (OBS)

Or how the Grid is used today:



OBS

Pilot
Factory

Site A

Site B

Site C

Illustration by Max Fischer (KIT)

Or how the Grid is used today:

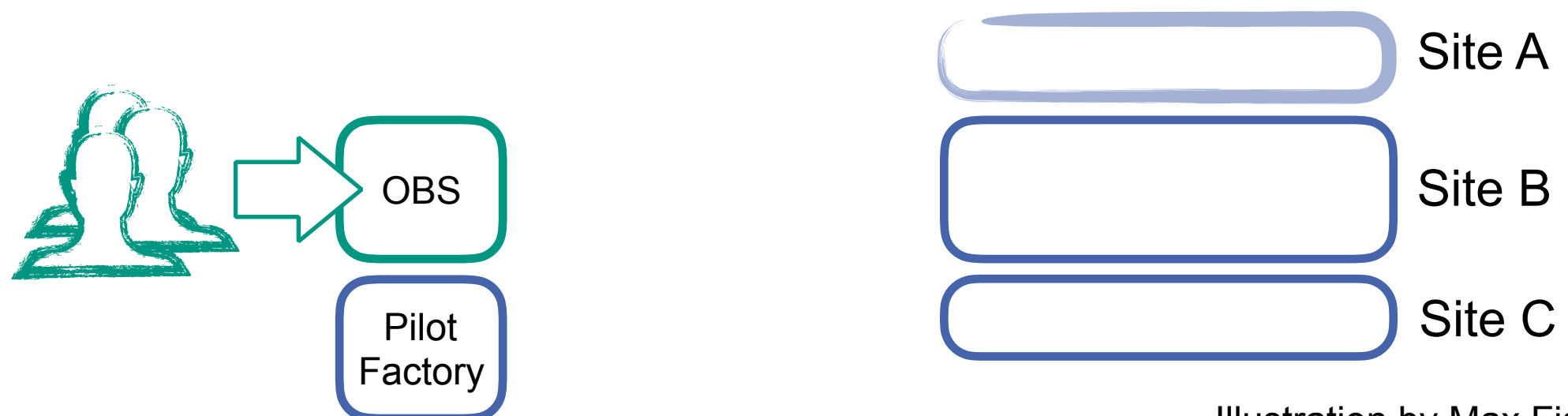- Pilot factory submits placeholder jobs (pilots) to different sites

Illustration by Max Fischer (KIT)

# Transparent Integration: Overlay Batch System (OBS)

Or how the Grid is used today:

- ■ Pilot factory submits placeholder jobs (pilots) to different sites
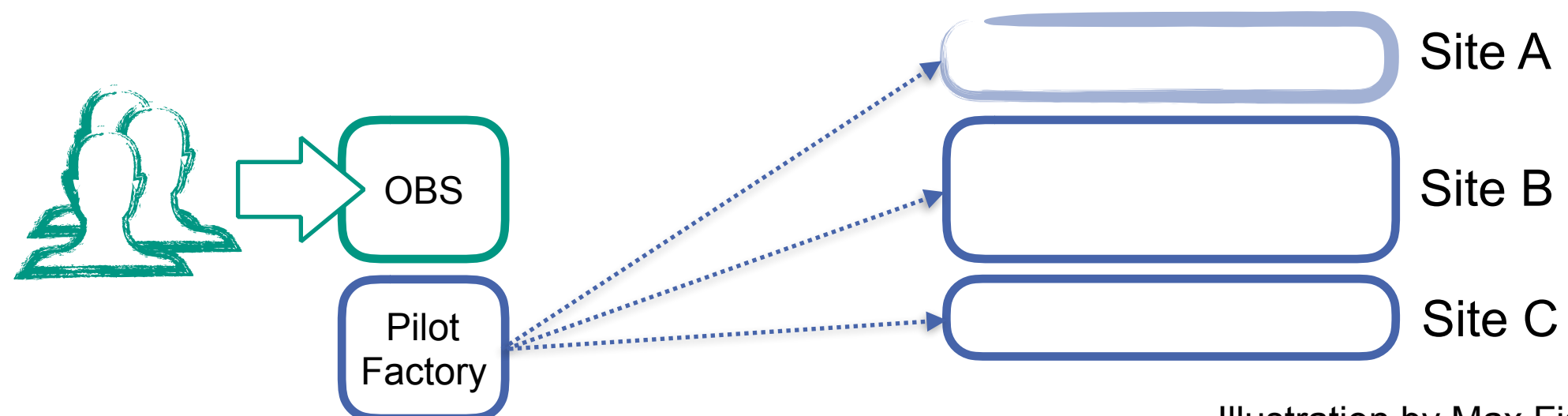- ■ Pilot allocates resources at the site



OBS

Pilot Factory

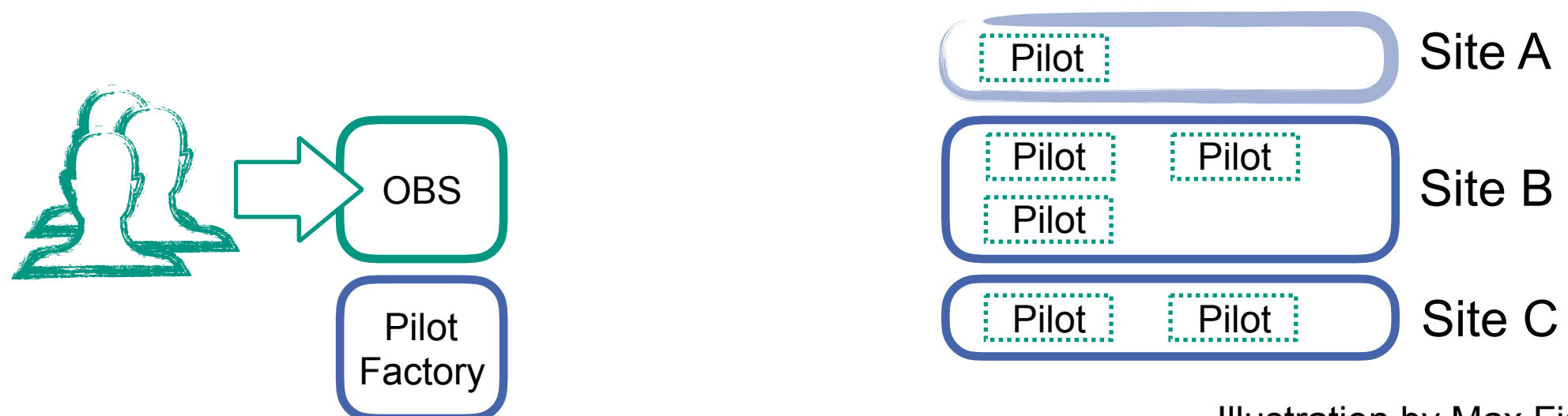| Pilot | Site A |
| Pilot | Pilot |
| Pilot | Site B |
| Pilot | Pilot | Site C |

Illustration by Max Fischer (KIT)

# Transparent Integration: Overlay Batch System (OBS)

Or how the Grid is used today:

- Pilot factory submits placeholder jobs (pilots) to different sites

- Pilot allocates resources at the site

- Resources are integrated into the OBS

- Workload is pulled from the OBS

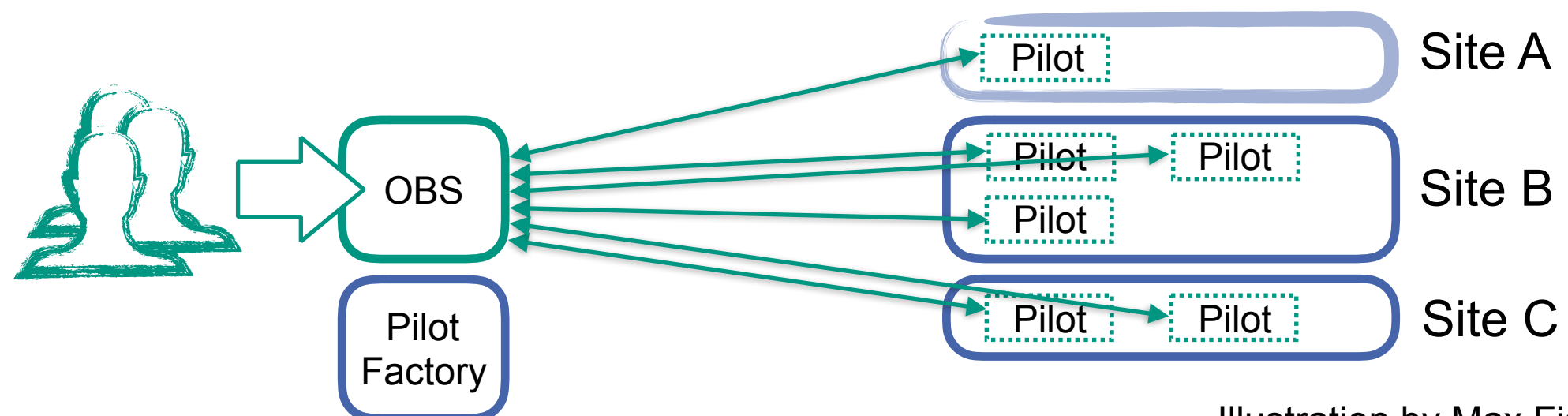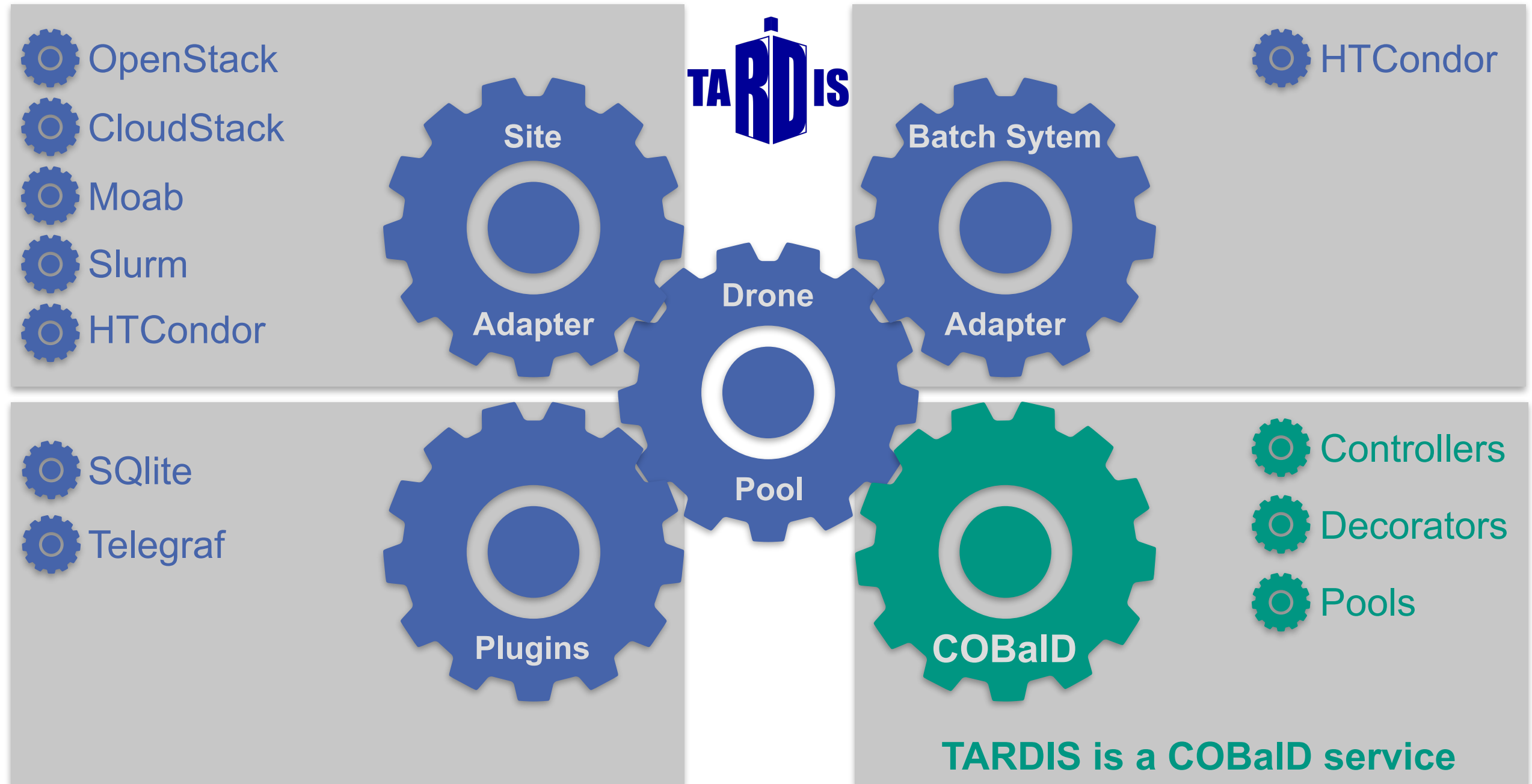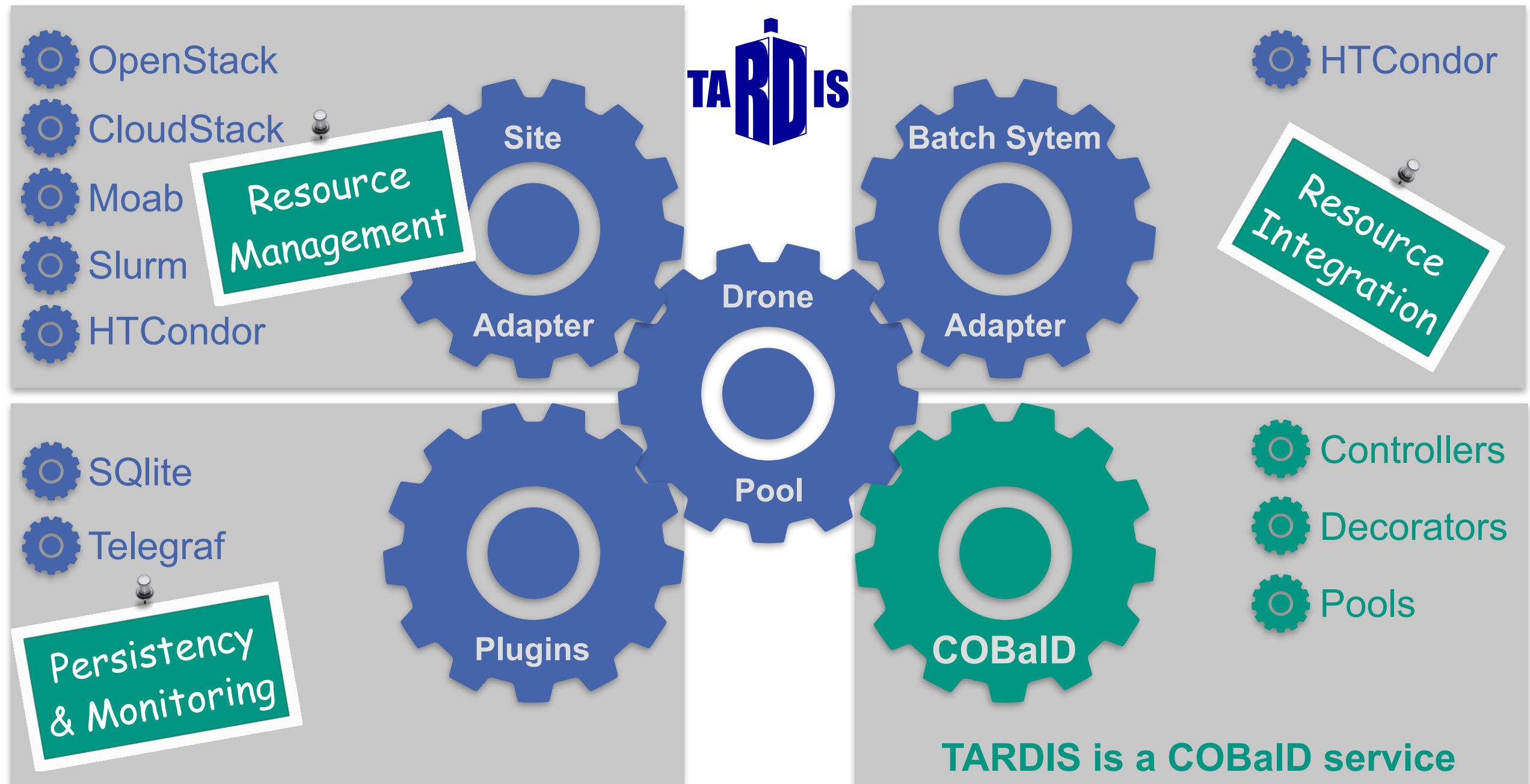➡ Users interact only with one single-point-of-entry the OBS



Site A — Pilot

Site B — Pilot, Pilot, Pilot

Site C — Pilot, Pilot

OBS

Pilot Factory

Illustration by Max Fischer (KIT)

# Overview about COBalD/TARDIS

OpenStack

CloudStack

Moab

Slurm

HTCondor

**Site Adapter**

**Drone Pool**

**Batch Sytem Adapter**

HTCondor

SQlite

Telegraf

**Plugins**

**COBalD**

Controllers

Decorators

Pools

**TARDIS is a COBalD service**

# Overview about COBaID/TARDIS

OpenStack
CloudStack
Moab
Slurm
HTCondor

**Resource Management**

**Site Adapter**

TARDIS

**Batch Sytem Adapter**

**Drone Pool**

HTCondor

**Resource Integration**

SQlite
Telegraf

**Persistency & Monitoring**

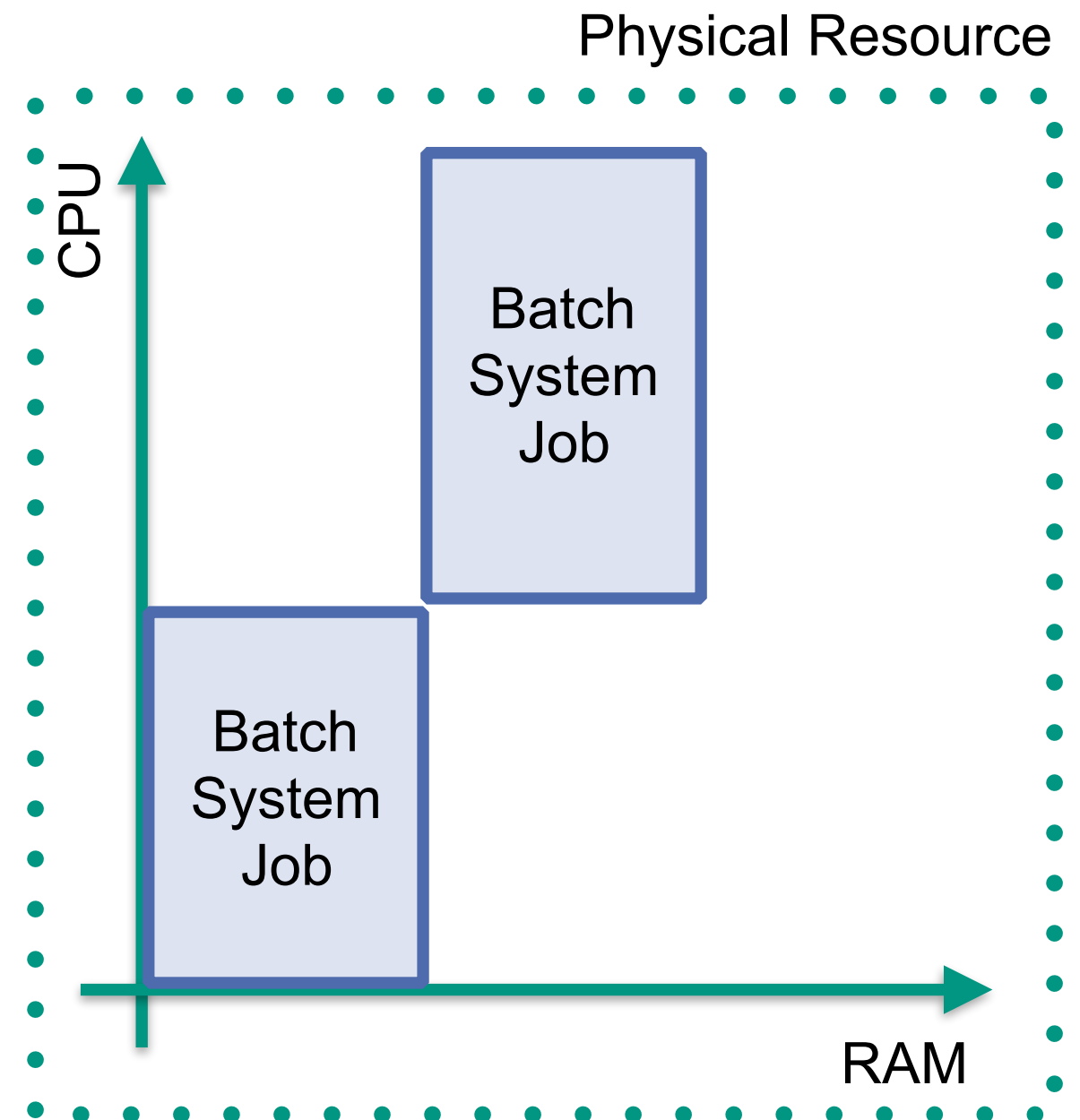**Plugins**

**COBaID**

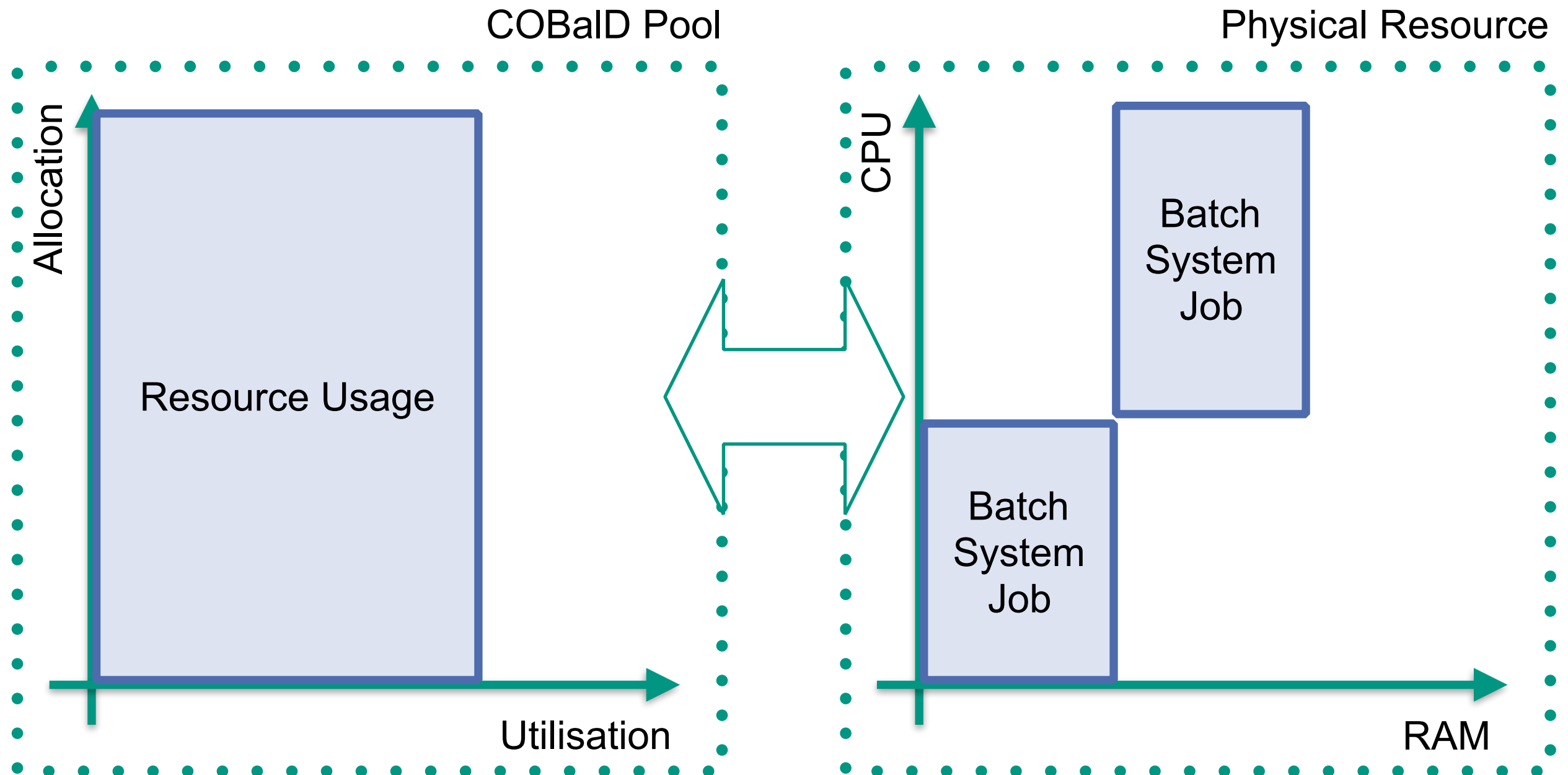Controllers
Decorators
Pools

**TARDIS is a COBaID service**

➜ Easily extendable by design through its modular structure

# COBalD Resource Pool Model

# COBalD Resource Pool Model



COBalD Pool

Physical Resource

Allocation

Utilisation

Resource Usage

CPU

RAM

Batch System Job

Batch System Job

# COBalD Resource Pool Model
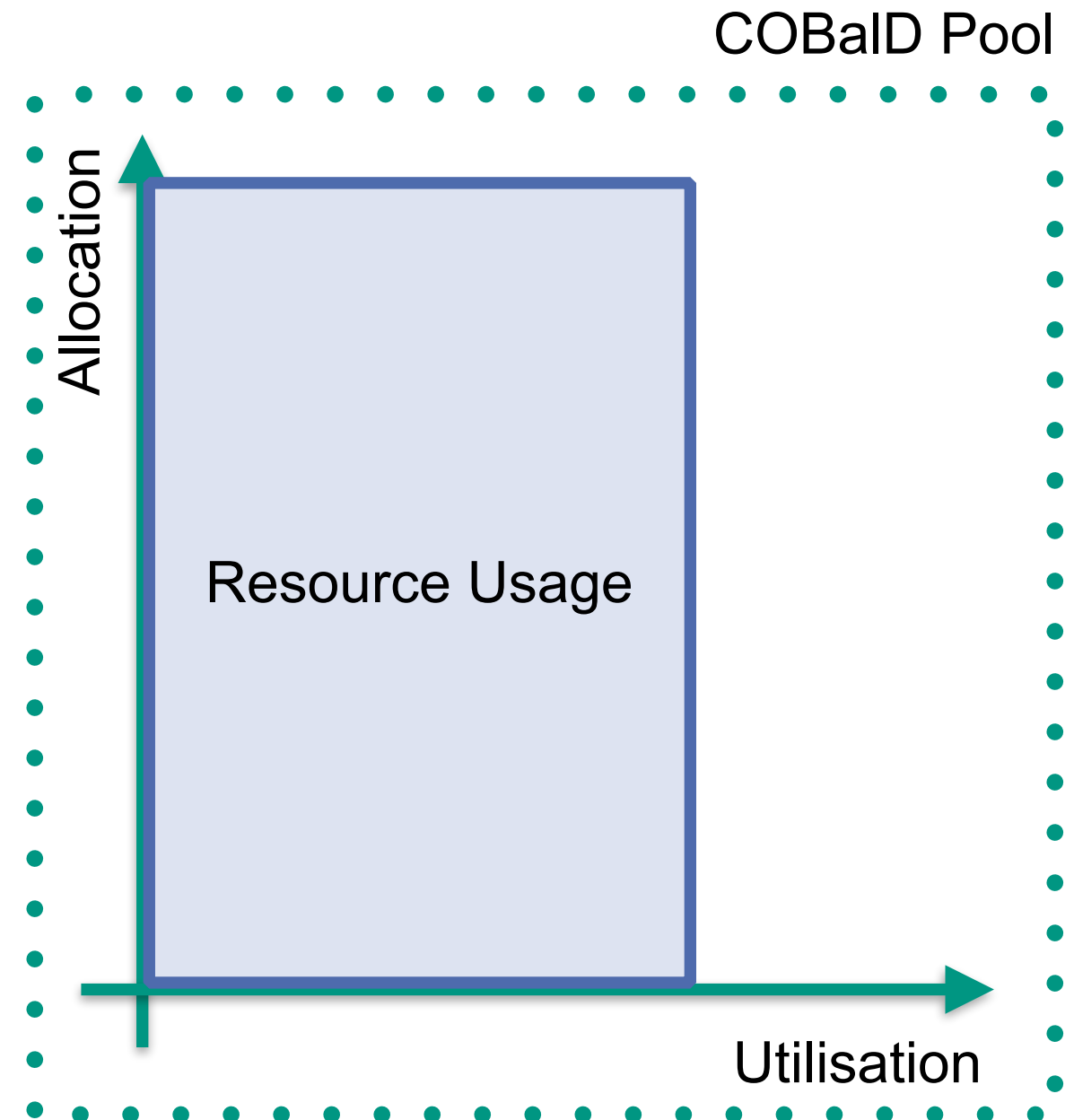


Lightweight Opportunistic Resource Management

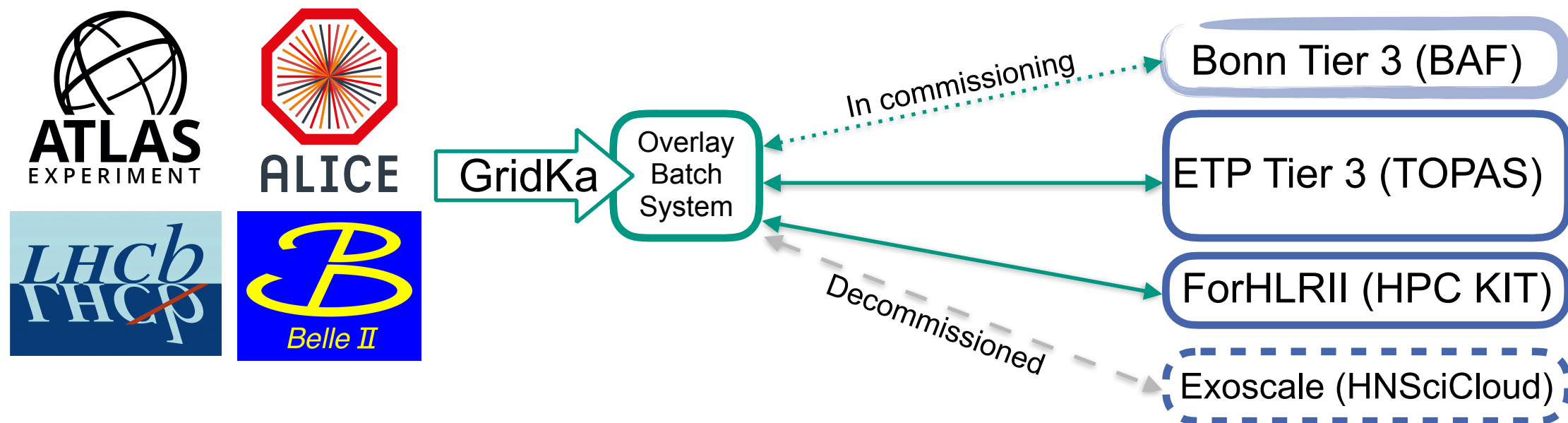# COBalD Resource Pool Model

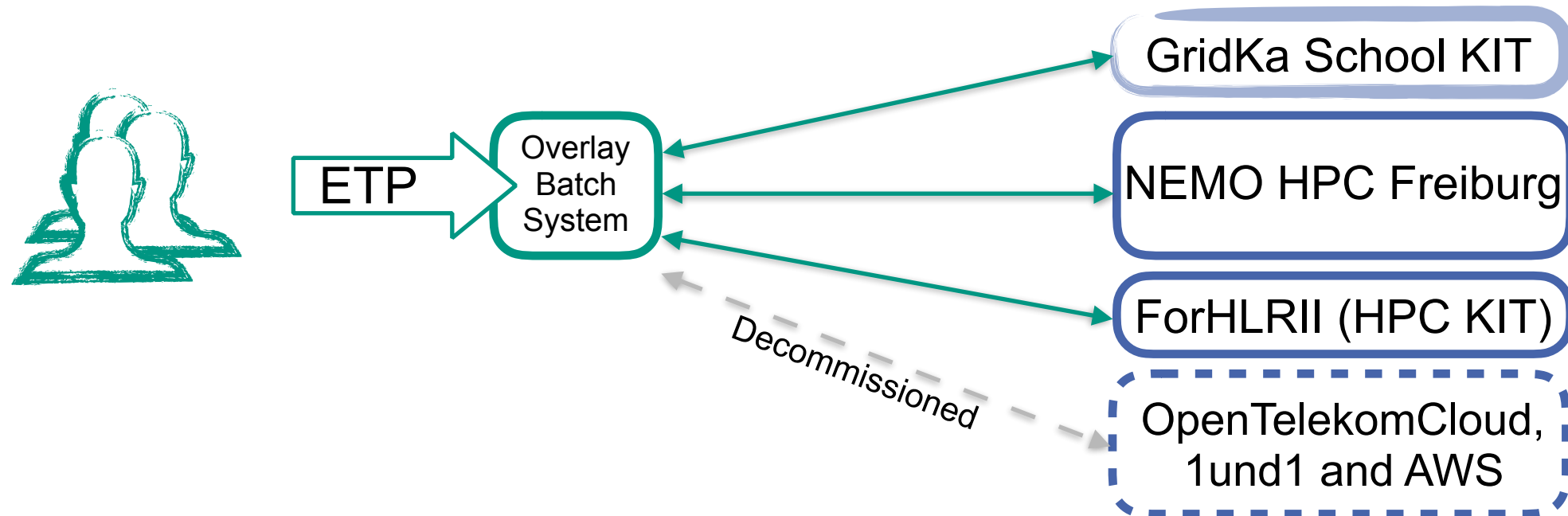# COBalD Resource Pool Model

```python
if utilisation < self.low_utilisation:
    return supply * self.low_scale
elif allocation > self.high_allocation:
    return supply * self.high_scale
```

COBalD Pool

Allocation

Resource Usage

Utilisation

# Currently Available Opportunistic Resources

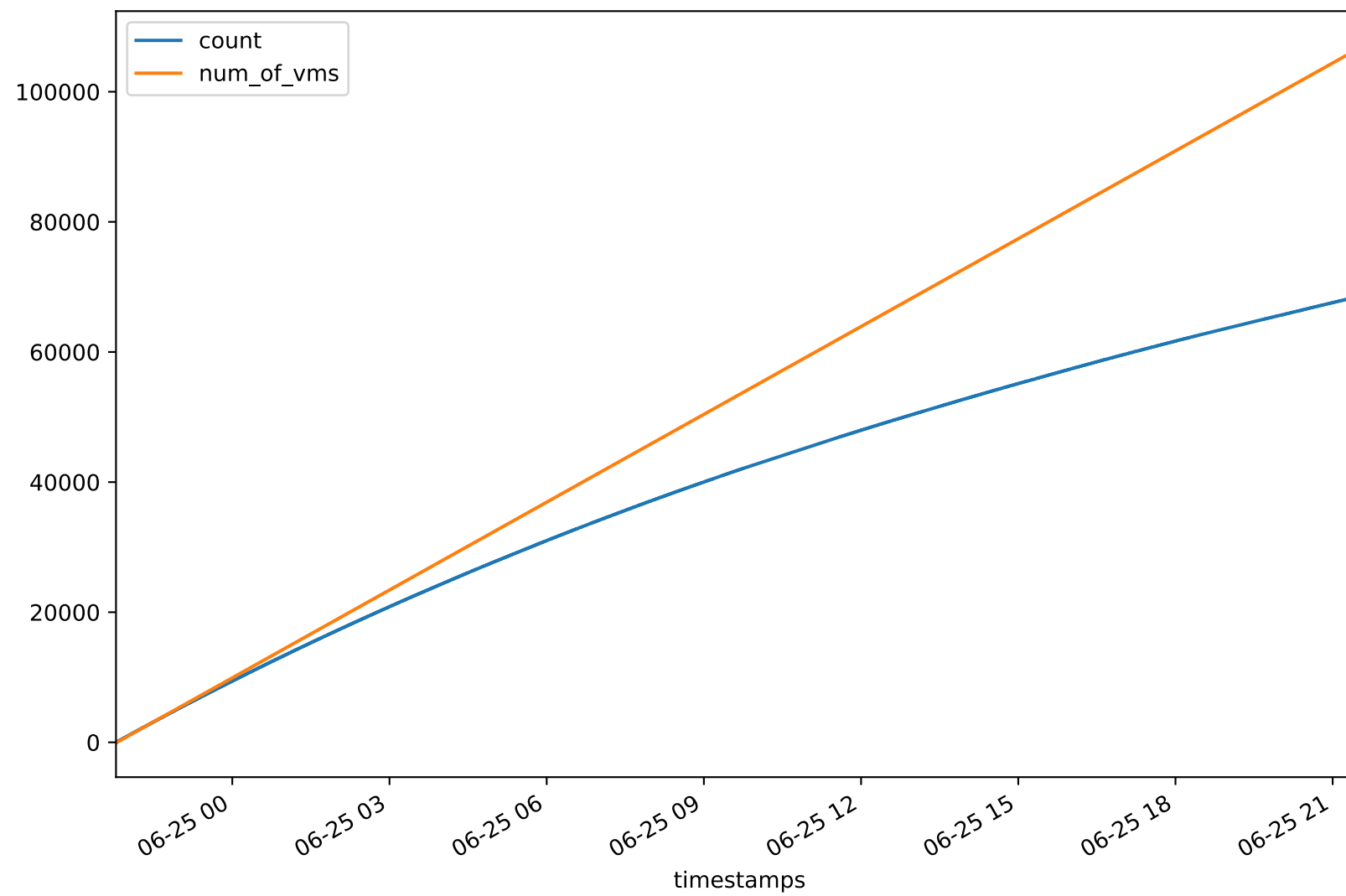# Currently Available Opportunistic Resources



➜ Perfectly suited to transparently integrate various resources into WLCG computing

Manuel Giffels

# Scalability

Manuel Giffels

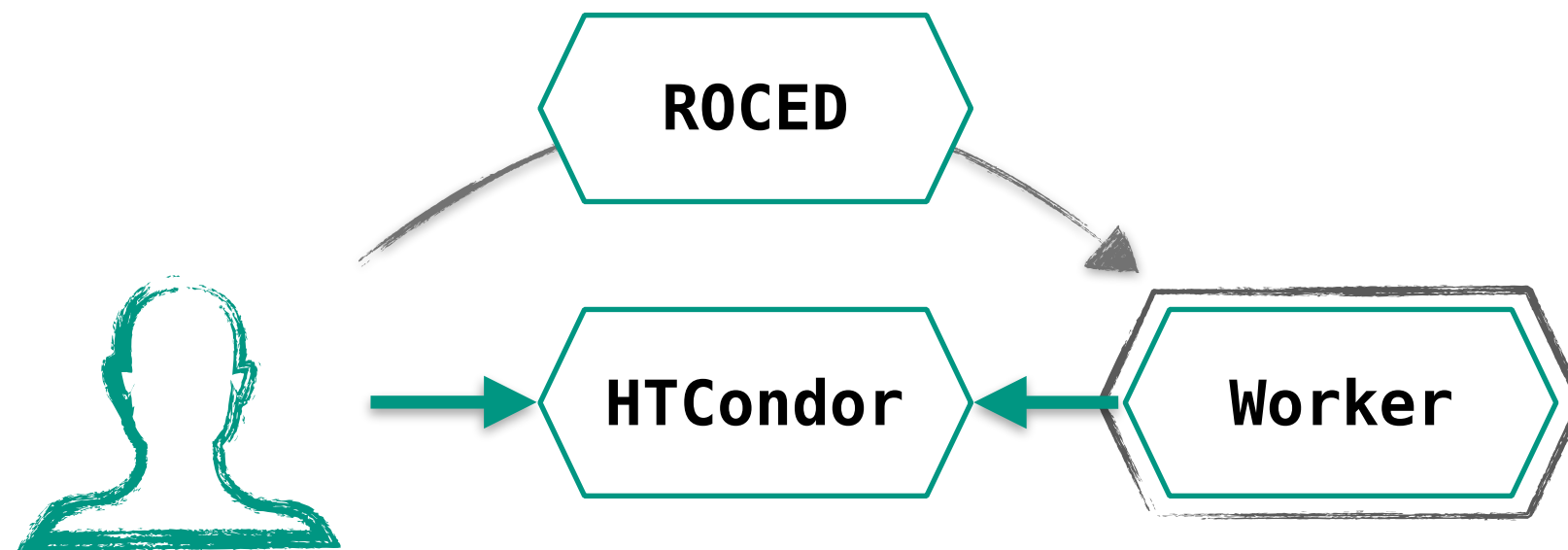ETP / SCC

# Opportunistic Resources @ KIT (2007-present)

- Longterm experience with opportunistic resources, virtualization and containers
- Software development (Cloud Resource Manager) ViBatch, ROCED (➟ COBalD/TARDIS)[1]
- Opportunistic Resources
  - Institute resources
    - Desktop cluster (Docker)
    - GridKa School Virtual Infrastructure (OpenStack)
  - HPC Cluster
    - ~~IC1@Uni Karlsruhe (ViBatch)~~
    - ForHLR II @ KIT (Singularity)
    - bwFORcluster NEMO @ Uni Freiburg (OpenStack)
  - Commercial cloud providers
    - AWS, 1&1 Cloud Services, OTC, ExoScale

**All dynamically and transparently integrated in a single HTCondor instance using ROCED and CoBalD/TARDIS**



Legend:
- Middleware
- Provider
- Infrastructure technology
- Virtualisation technology

[1] COBalD - the Opportunistic Balancing Daemon (http://cobald.readthedocs.io/) by M. Fischer & E. Kühn
TARDIS - Transparent Adaptive Resource Dynamic Integration System (https://github.com/giffels/tardis) Illustration by E. Kuehn (KIT)

# ROCED Resumé



Slide by Max Fischer (KIT)

Manuel Giffels
SCC / ETP

# ROCED Resumé

Manuel Giffels

SCC / ETP

# ROCED Resumé

- Dynamic resources matching user demand
  - Trivial to support new providers for many users
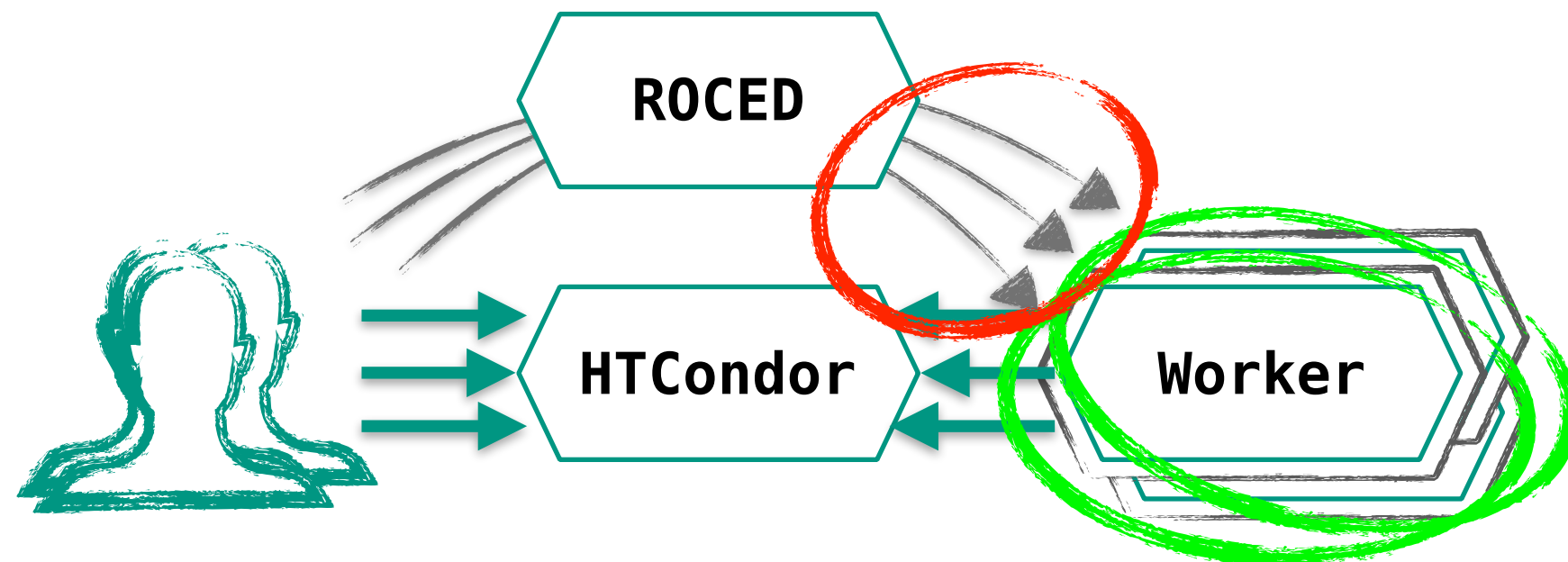  - Difficult to manage several providers for many users



Slide by Max Fischer (KIT)

Manuel Giffels

# ROCED Resumé
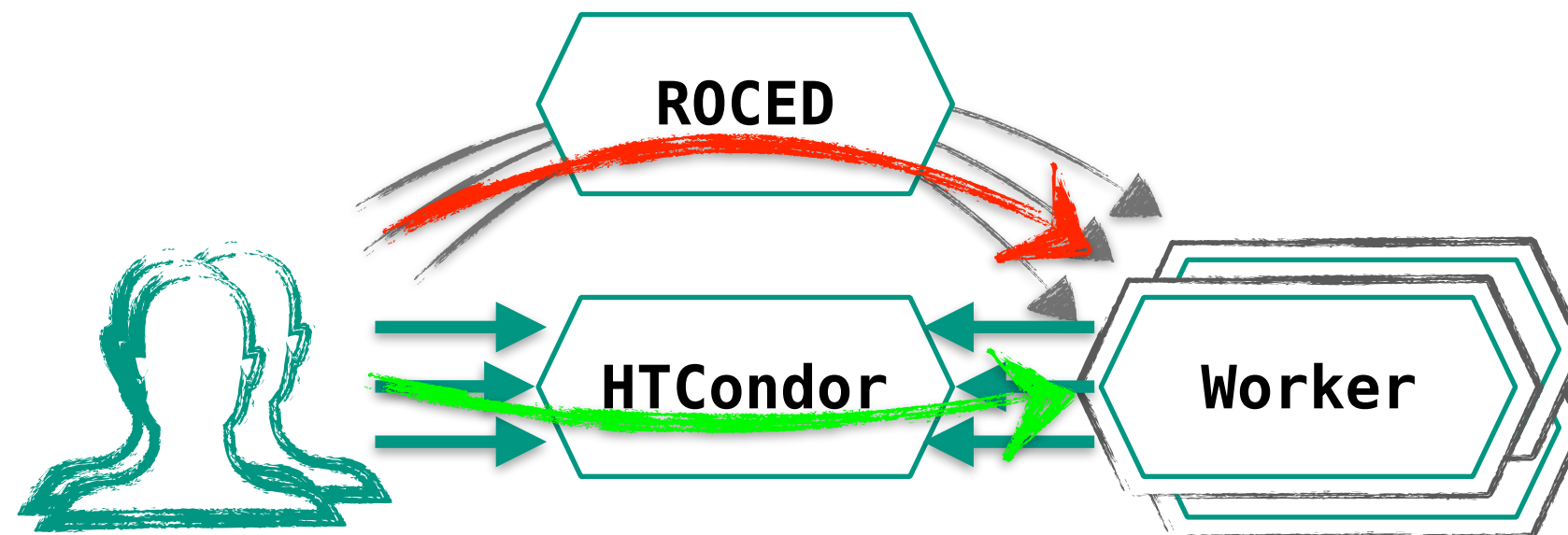
- Dynamic resources matching user demand
  - Trivial to support new providers for many users
  - Difficult to manage several providers for many users
- Resource aggregation in overlay batch system
  - Unreliable to predict resources required for jobs
  - Efficient to integrate resources, then match jobs



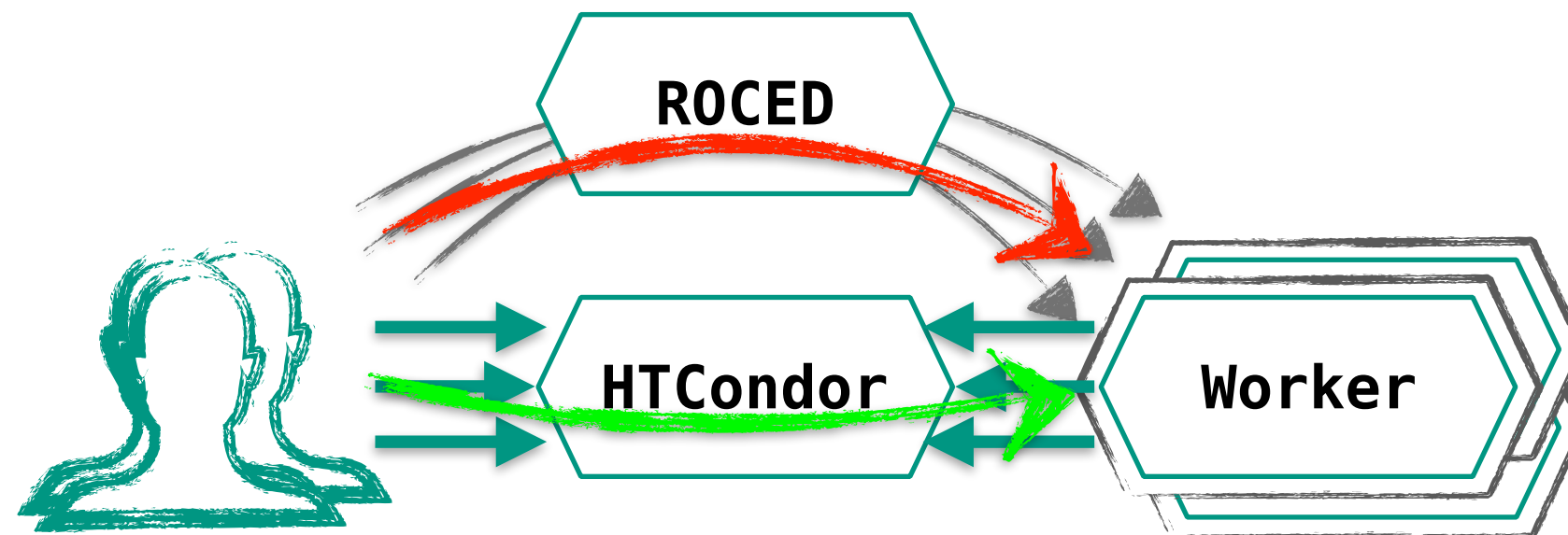Slide by Max Fischer (KIT)

Manuel Giffels

# ROCED Resumé

- Dynamic resources matching user demand
    - Trivial to support new providers for many users
    - Difficult to manage several providers for many users
- Resource aggregation in overlay batch system
    - Unreliable to predict resources required for jobs
    - Efficient to integrate resources, then match jobs
- Yet it really works!



Slide by Max Fischer (KIT)

# Innovative Digital Technologies for Exploring Universe and Matter

- Joint proposal by HEP, Physics of Hadrons and Nuclei, Astroparticle Physics
- <u>Covered Topics:</u>
  - Development of technologies to utilize heterogeneous computing resources (Integration of Opportunistic Resources, Caching Technologies, Workflow Management)
  - Application and testing of those technologies in heterogenous computing resources
  - Deep Learning - Achieving knowledge through profound data-driven methods (Hardware-related Data Processing, Object Reconstruction, Simulation, Quality of Network Predictions)
  - Event reconstruction: Cost- and energy efficient utilization of computing resources (Alternative Algorithms and Architectures like GPUs)
- Funded in the scope of Digital Agenda programme (BMBF)

<u>Proposal of:</u>

<u>Associated Partners:</u>

# National Research Data Infrastructure (NFDI)

Particle, Astroparticle and Hadron & Nuclear Physics Accelerates the NFDI:

- Task areas:
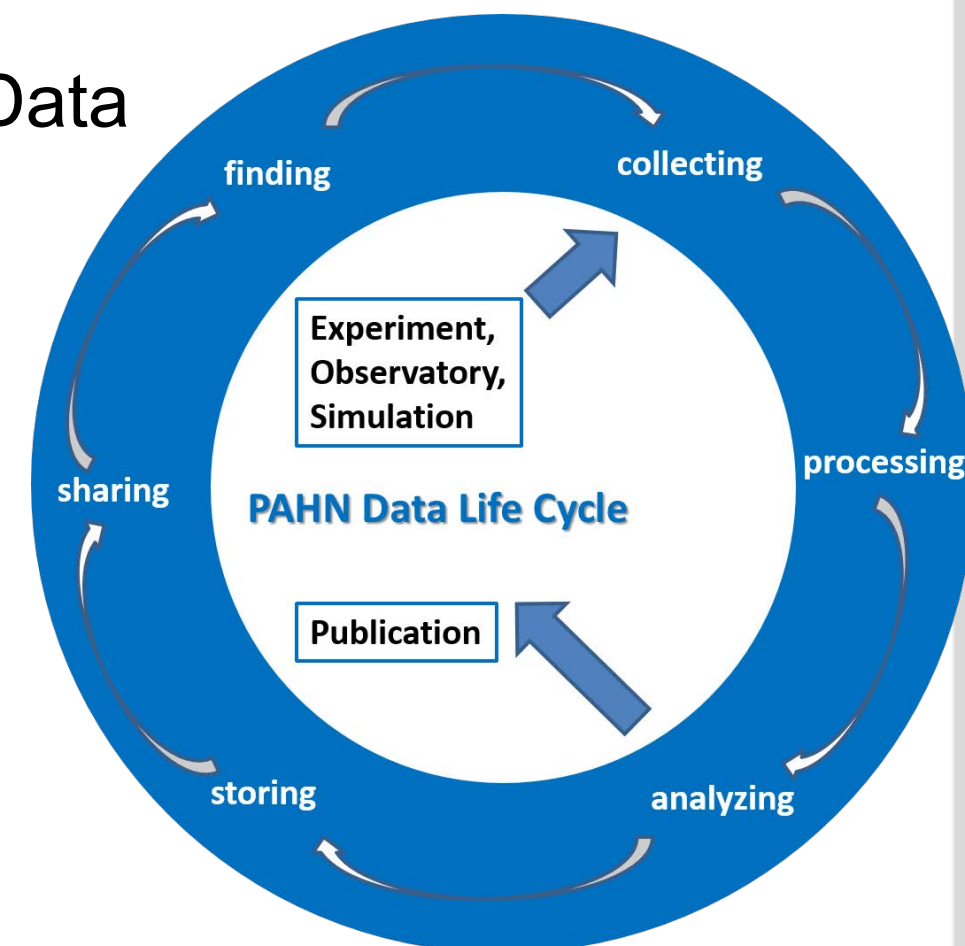  - Developing workflows and tools for data management
    - Tools to utilize heterogenous and HPC resources
  - …
- FAIR Data Life Cycle Concepts and Open Data
- Data analysis procedures and services
- Real-time data analysis and selection



Proposal deadline in October!

# Resources

- COBalD: http://cobald.readthedocs.io/
- ROCED: https://github.com/roced-scheduler/ROCED
- TARDIS: http://cobald-tardis.readthedocs.io/
- COBalD Simulation: https://git.scc.kit.edu/fq8360/cobalt_sim
- COBalD Demo: https://github.com/MaineKuehn/cobald_demo

Manuel Giffels