

# Dynamic Transparent Integration and Management of Opportunistic Resources with COBaID/TARDIS

ErUM Data Cloud Workshop, Karlsruhe 27.06.2019

Manuel Giffels, René Caspart Max Fischer, Eileen Kühn, Matthias Schnepf, Florian v. Cube

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



# Opportunistic Resources and their Challenges

- Very different to the traditional HEP environment
  - ➔ Virtualisation and containerisation techniques
- Temporary availability of opportunistic resources
  - ➔ Dynamic integration and workflow management
- Varying demand for opportunistic resources
  - ➔ Dynamic integration and workflow management
- Availability and suitability for given tasks is not predictable
  - ➔ Dynamic reaction on availability and utilisation
- Each opportunistic resource is different (very heterogeneous system)
  - ➔ Hard to manage for users and computing operations of experiments
  - ➔ Transparent integration of resources needed

# Transparent Integration: Overlay Batch System (OBS)

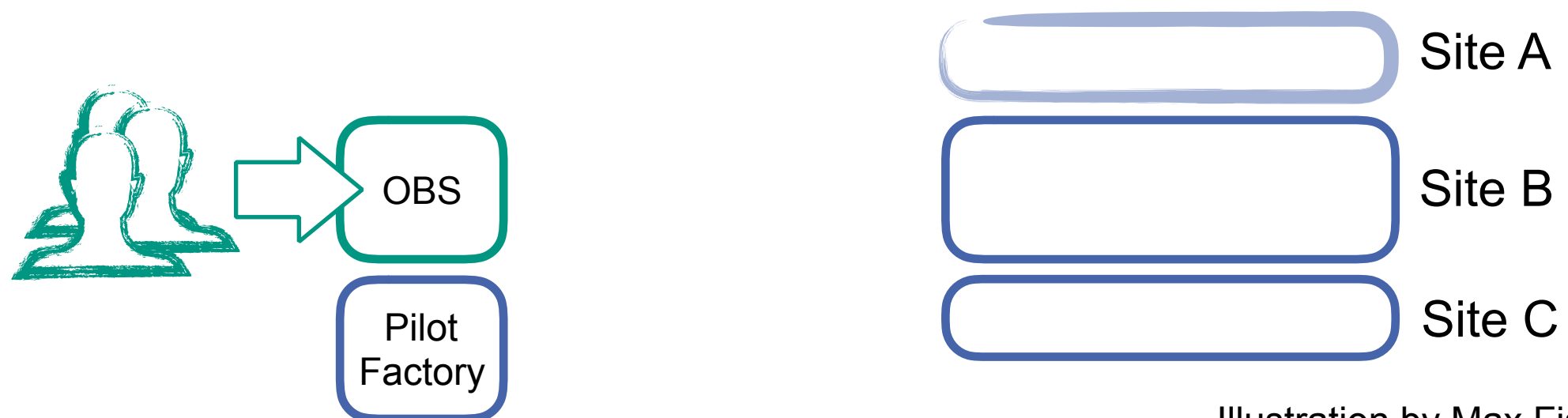


Illustration by Max Fischer (KIT)

# Transparent Integration: Overlay Batch System (OBS)

Or how the Grid is used today:

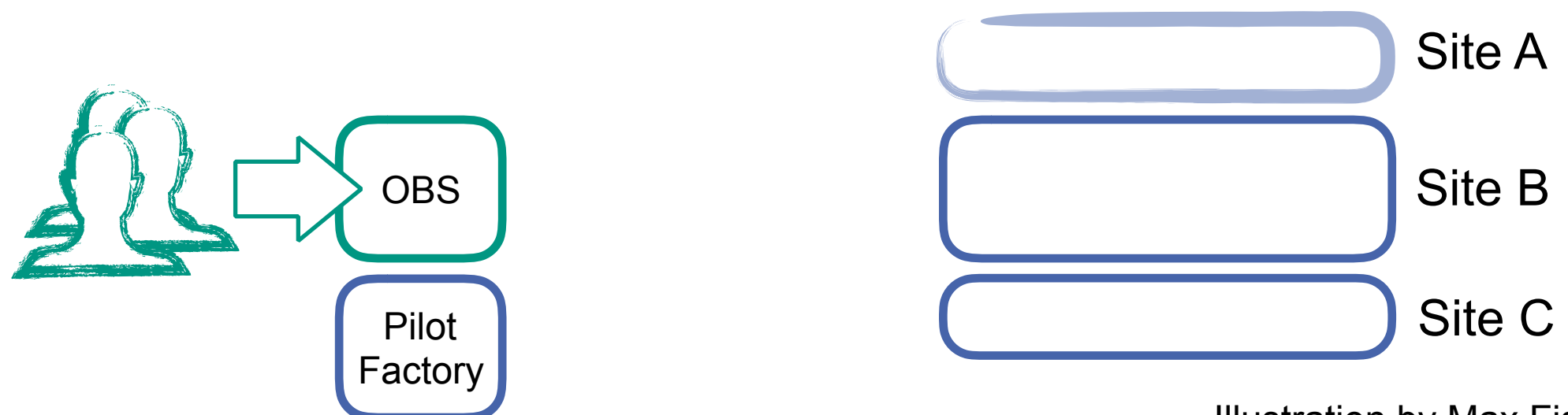


Illustration by Max Fischer (KIT)



# Transparent Integration: Overlay Batch System (OBS)

Or how the Grid is used today:

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

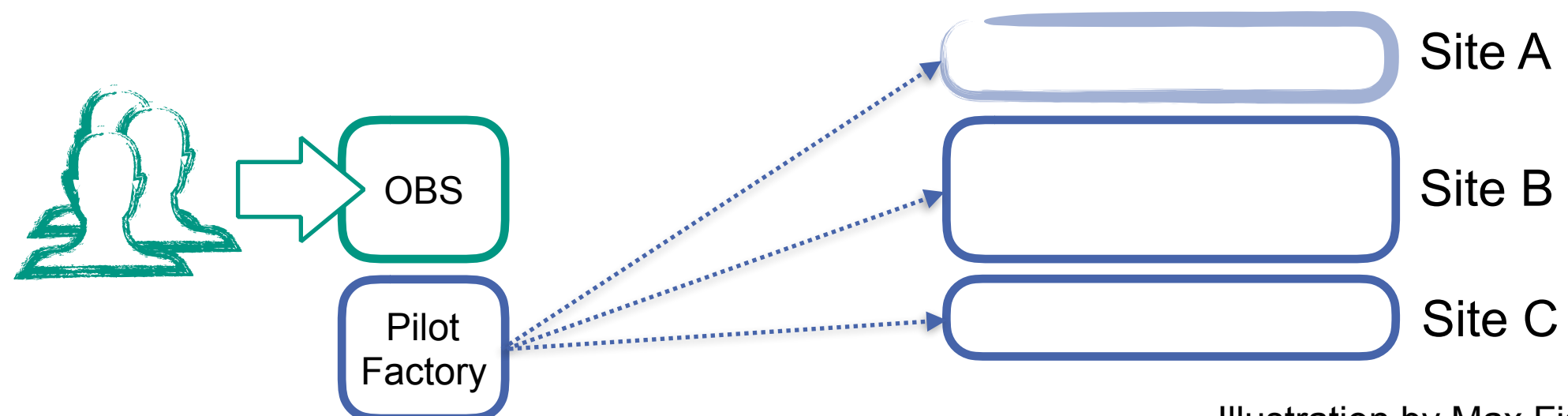


Illustration by Max Fischer (KIT)

# 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

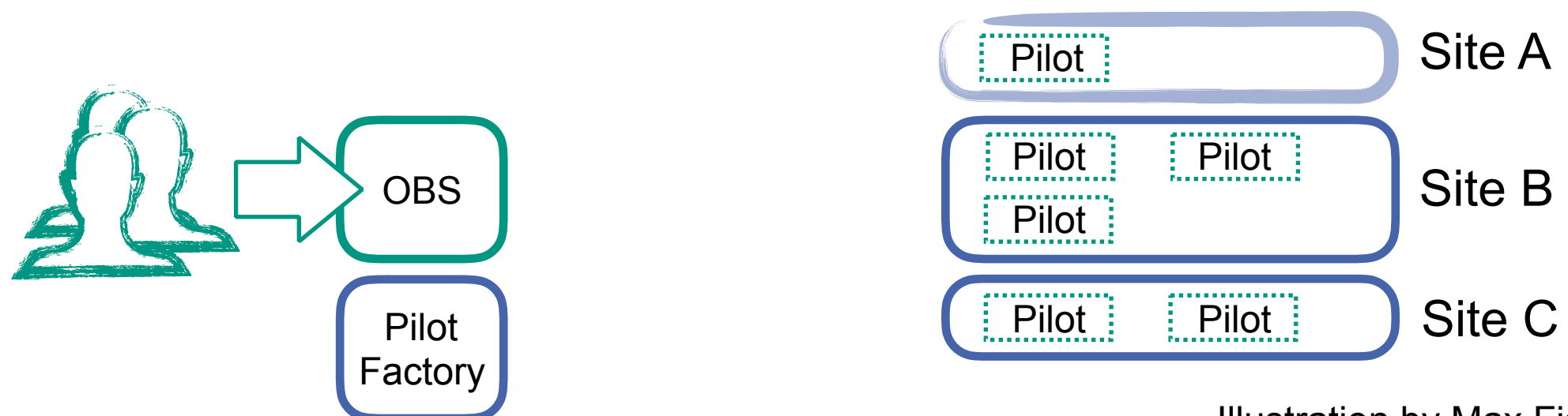


Illustration by Max Fischer (KIT)

# 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

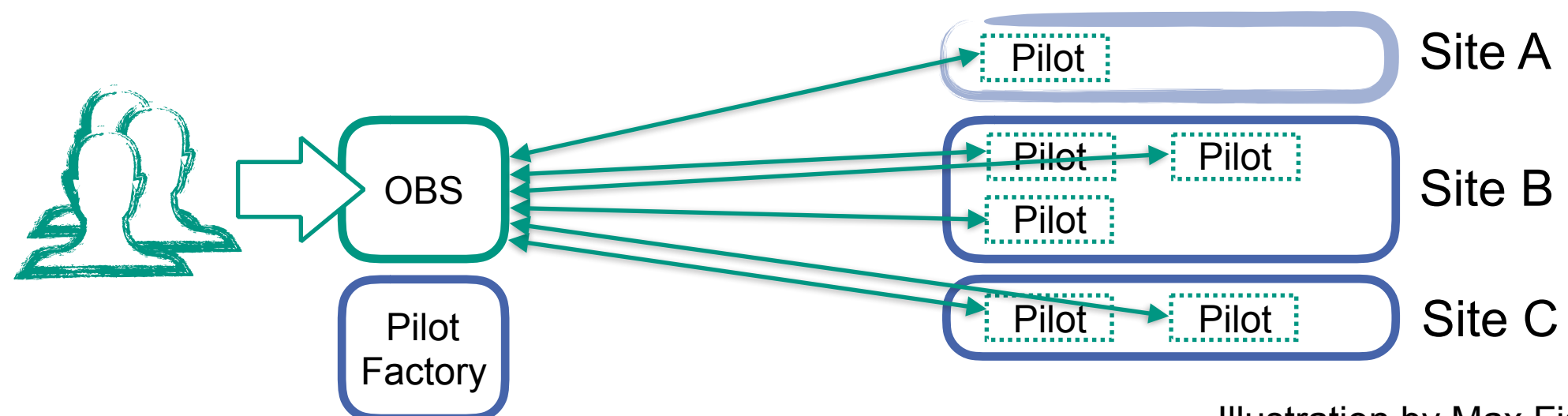
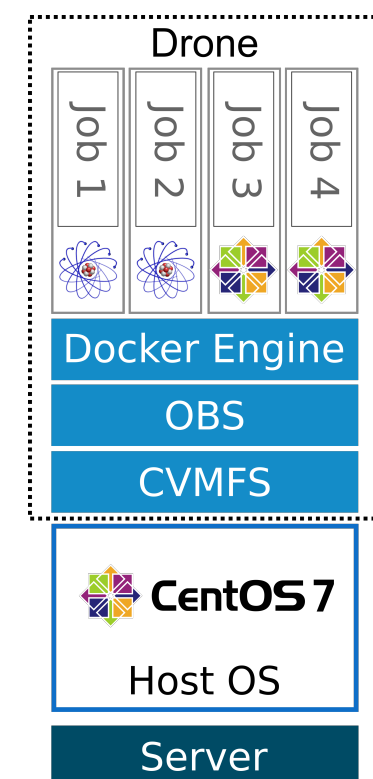
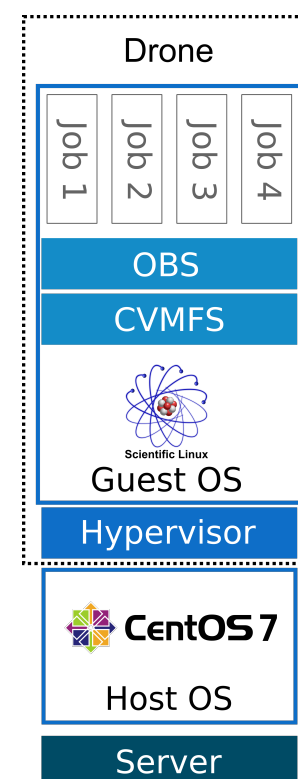
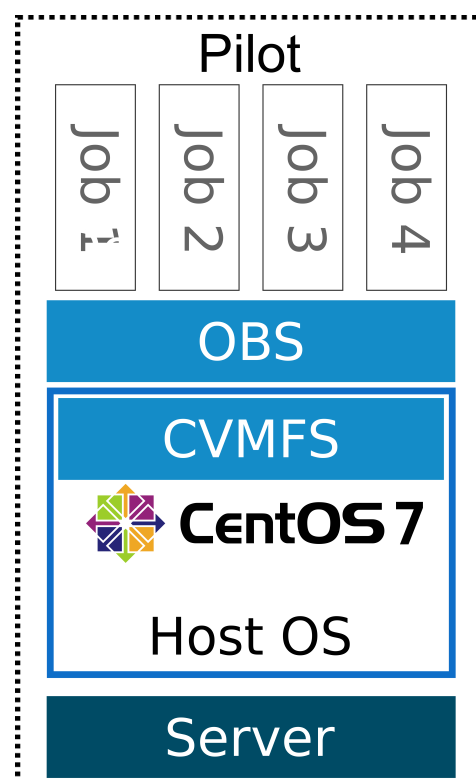


Illustration by Max Fischer (KIT)

# Drone Concept

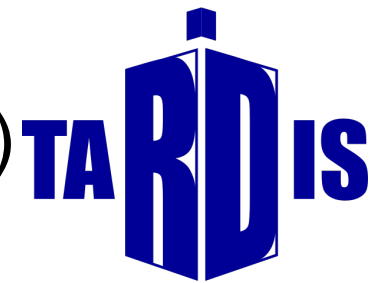
- Pilot concept
  - Allocate and integrate resources (via OBS)
  - Usually just a batch system daemon plus wrapper scripts
- Drone is a generalised pilot
  - Allocate and integrate resources (via OBS same as a pilot)
  - Provide dedicated environment (OS/Software) via virtualisation or containerisation techniques



# Dynamic Resource Integration with TARDIS

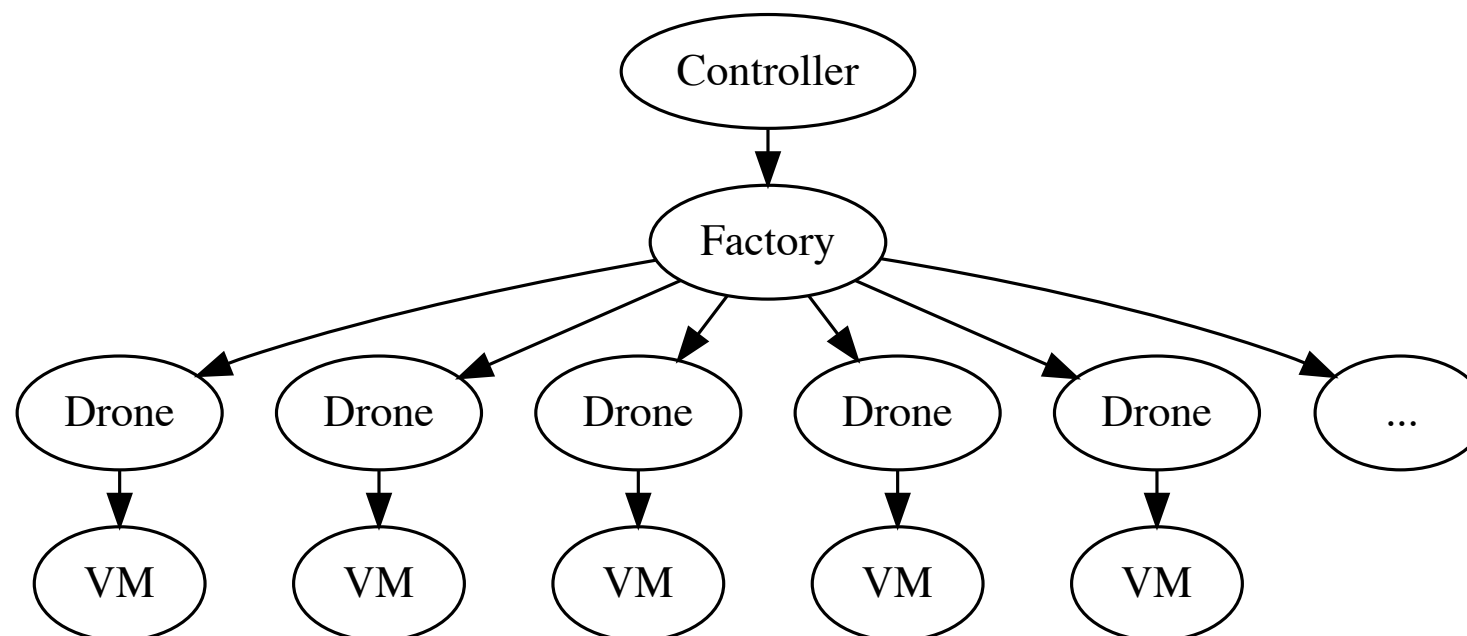
Transparent Adaptive Resource Dynamic Integration System (TARDIS):

- Provides **backends** to **various cloud providers** and **batch systems**
- Allows **dynamic orchestration** of pre-built VMs and containers
- **Asynchronous** multi-agent COBaID service (ensures scalability)



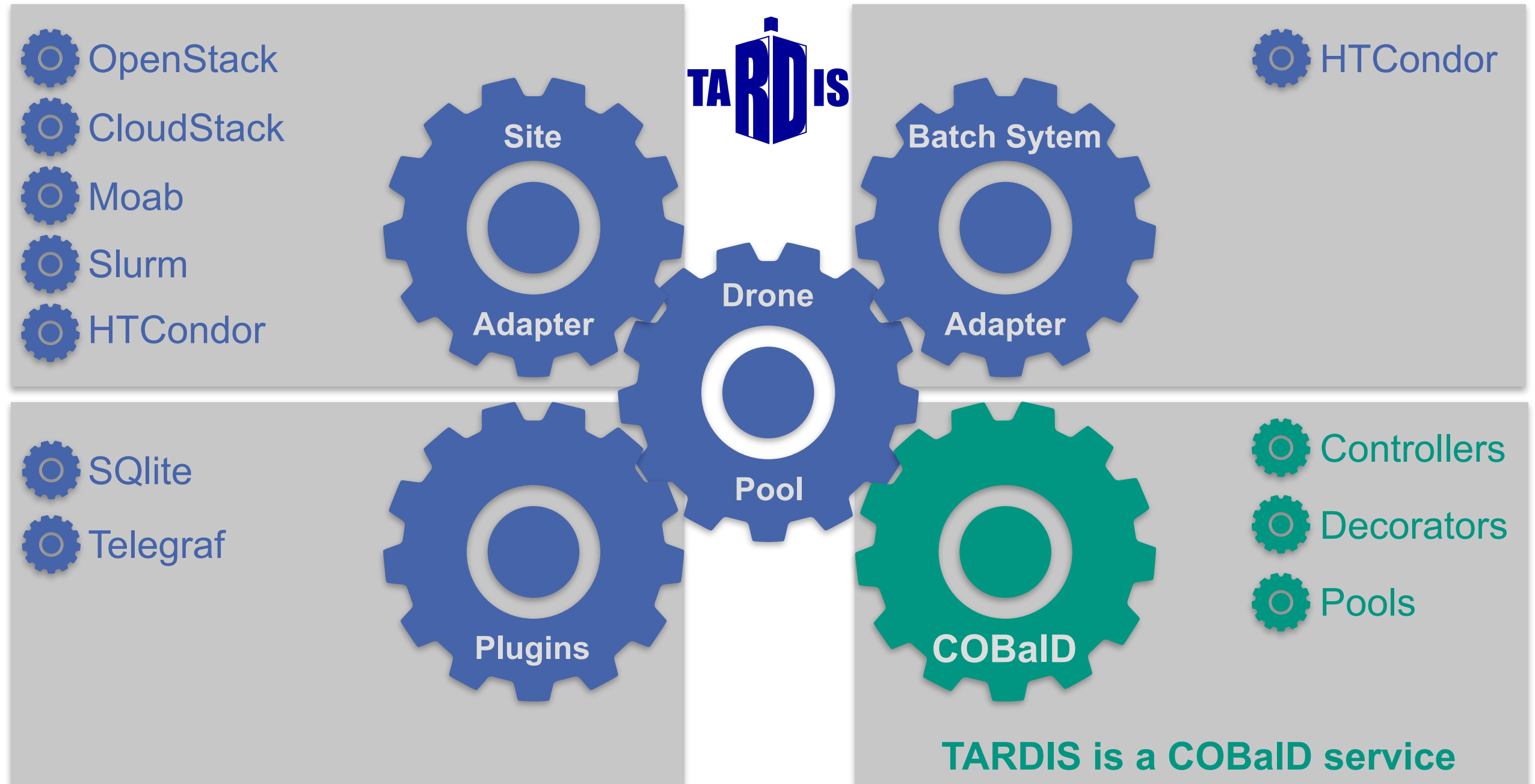
Structure:

- A resource is represented by a Drone implementing the pool interface
- Drones are aggregated in dynamically sized composite pools
- Factory pool acquires and releases resources as desired by the controller



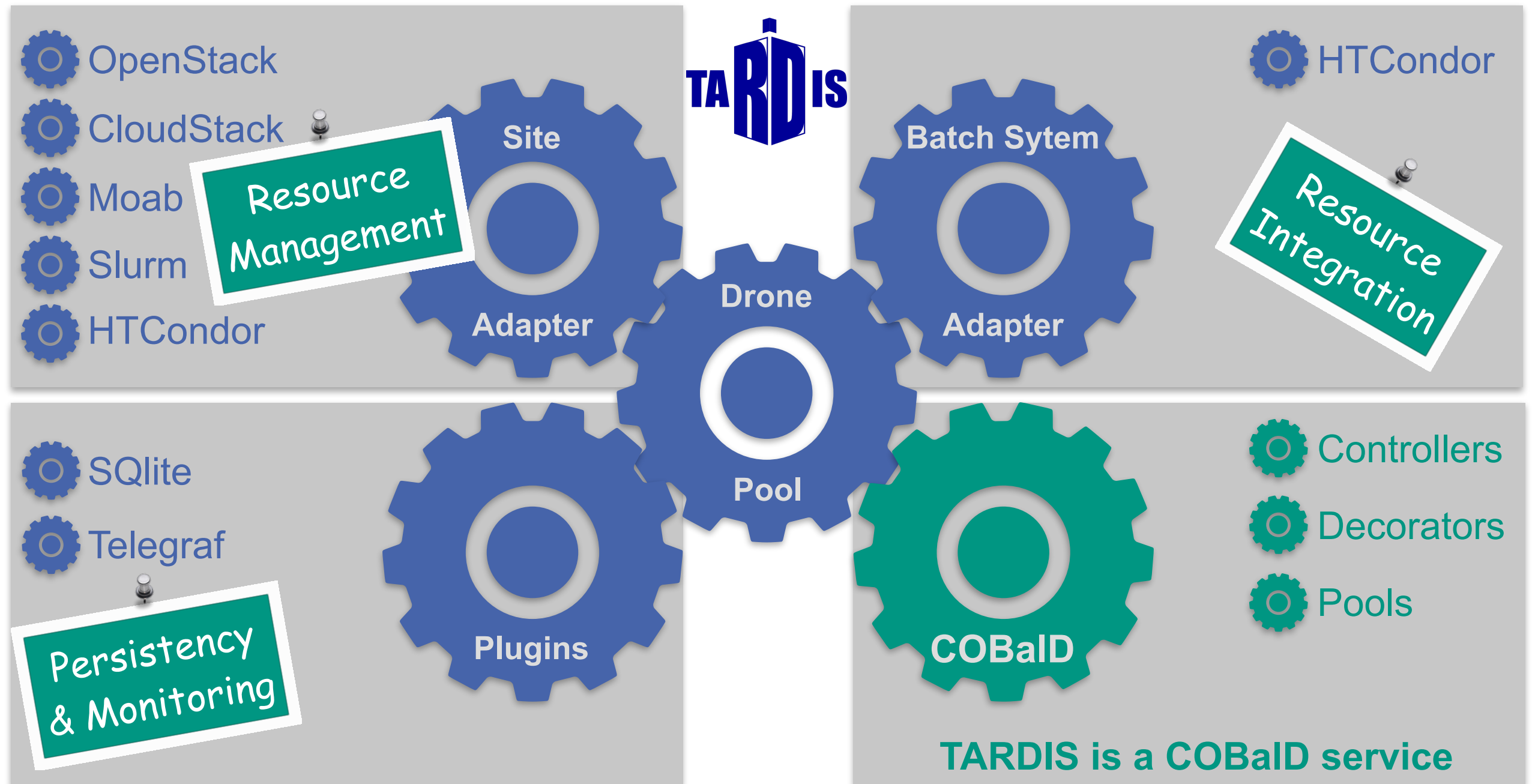


# Components of COBaID-TARDIS



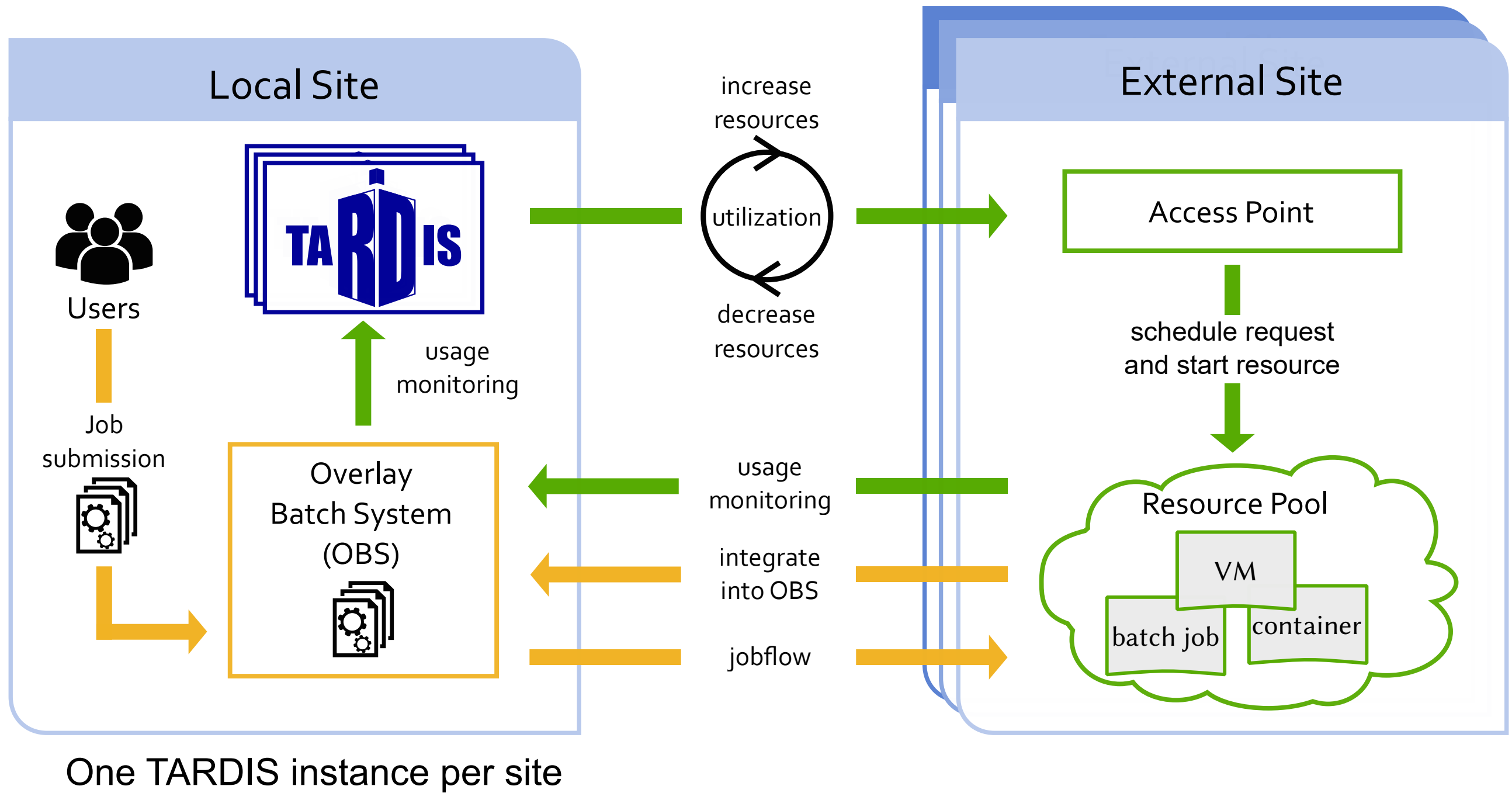
Well defined abstract base class interfaces to all components available!

# Components of COBaID-TARDIS

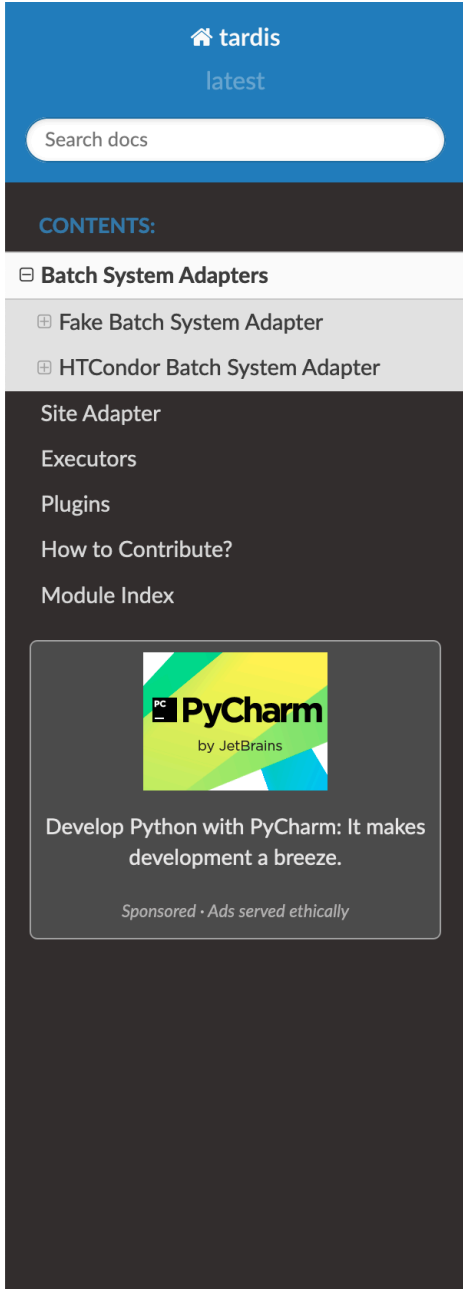


Well defined abstract base class interfaces to all components available!

# The Entire Picture



# First Documentation Available



[Docs » Batch System Adapters](#) [Edit on GitHub](#)

## Batch System Adapters

### Fake Batch System Adapter

The `FakeBatchSystemAdapter` implements a batch system adapter that mocks the response of hypothetical batch system. It can be used for testing purposes and as a demonstrator in workshops and tutorials.

The mocked response to the `get_allocation()`, `get_utilization()` and `get_machine_status()` API calls is configurable statically in the adapter configuration.

#### Available configuration options

Option	Short Description	Requirement
adapter	Name of the adapter (FakeBatchSystem)	Required
allocation	Mocked response to <code>get_allocation()</code> call	Required
utilization	Mocked response to <code>get_utilization()</code> call	Required
machine_status	Mocked response to <code>get_machine_status()</code> call	Required

#### Example configuration

```
BatchSystem:
  adapter: FakeBatchSystem
  allocation: 1.0
  utilization: 1.0
  machine_status: Available
```

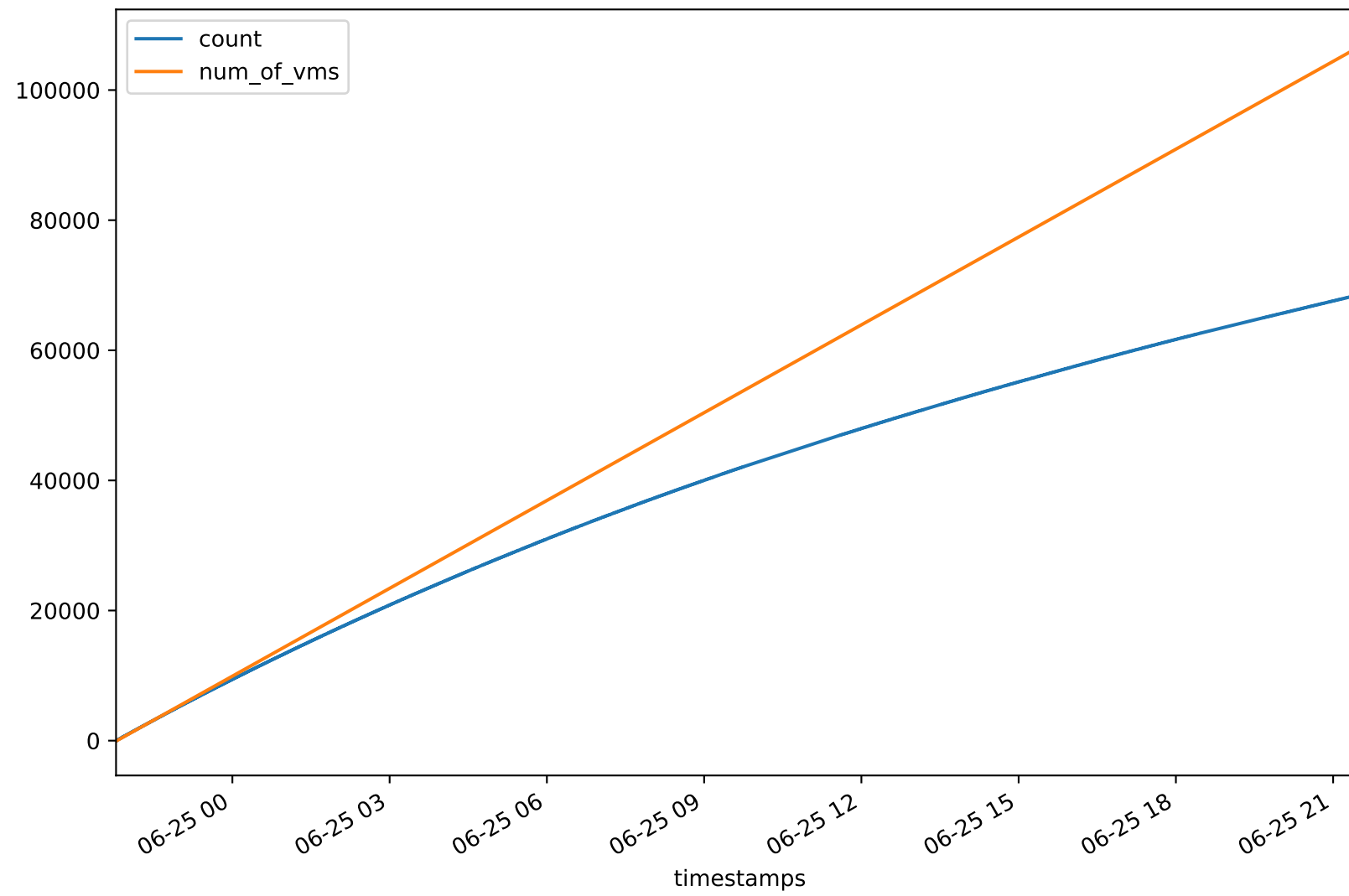
### HTCondor Batch System Adapter

<https://cobald-tardis.readthedocs.io/en/latest/index.html>

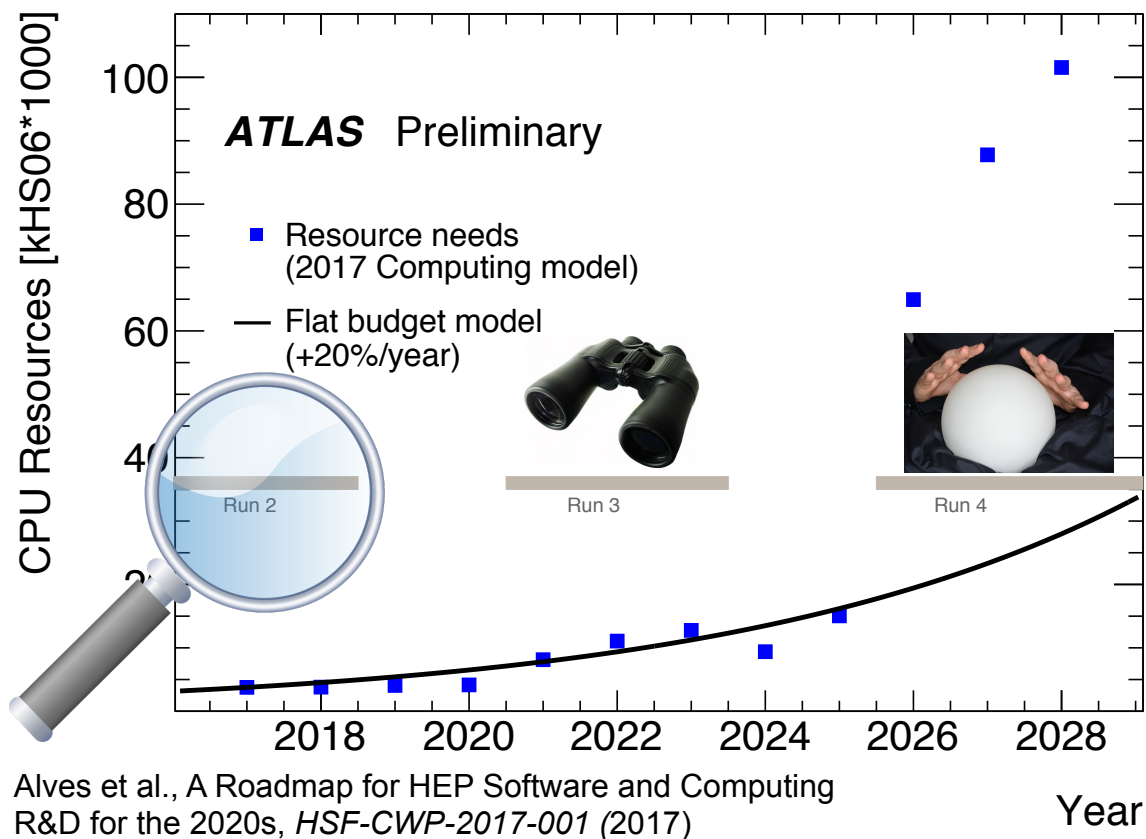
# BACKUP



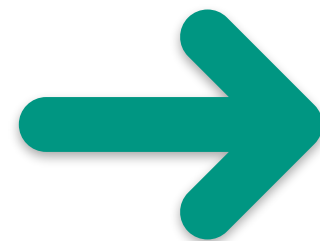
# Scalability



# HEP Computing Challenge Ahead



- ATLAS/CMS CPU resource estimates
- Assuming flat budget and 20% technology advance per year
- CPU shortfall between needs and technology gains about factor 4 in 2027
- Situation for disk storage is even worse (Shortfall is about factor 7 in 2027)
- ➔ Critical conditions for HL-LHC (Run 4)



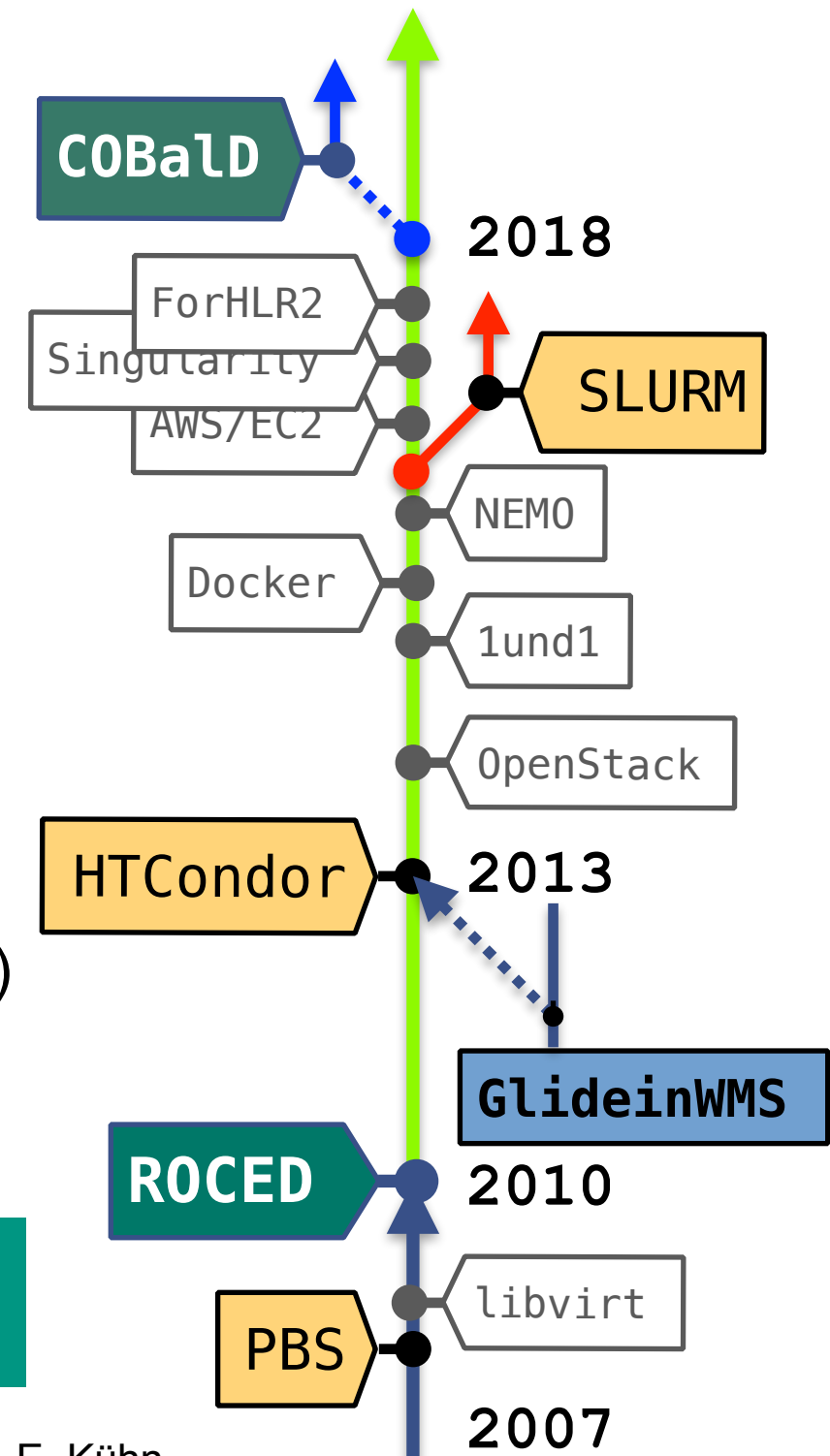
## New interesting research topics:

- 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

# Opportunistic Resources @ KIT (2007-present)

- Longterm experience with opportunistic resources, virtualization and containers
- Software development (Cloud Resource Manager) ViBatch, ROCED (→ COBaID/TARDIS)<sup>1</sup>
- 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

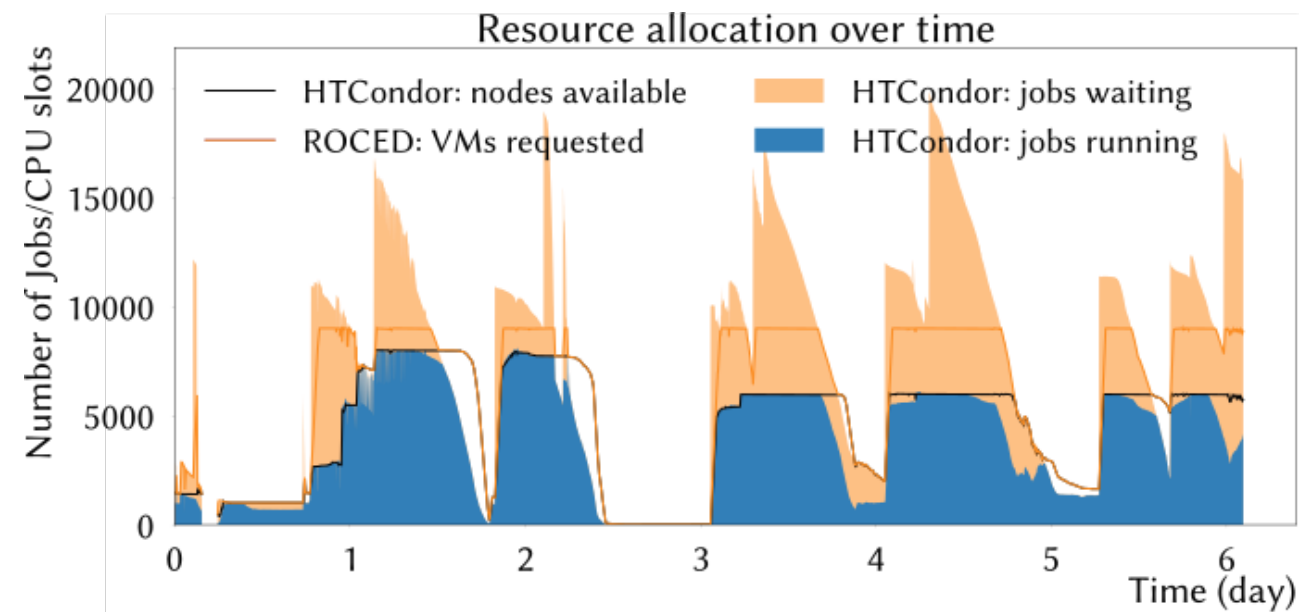
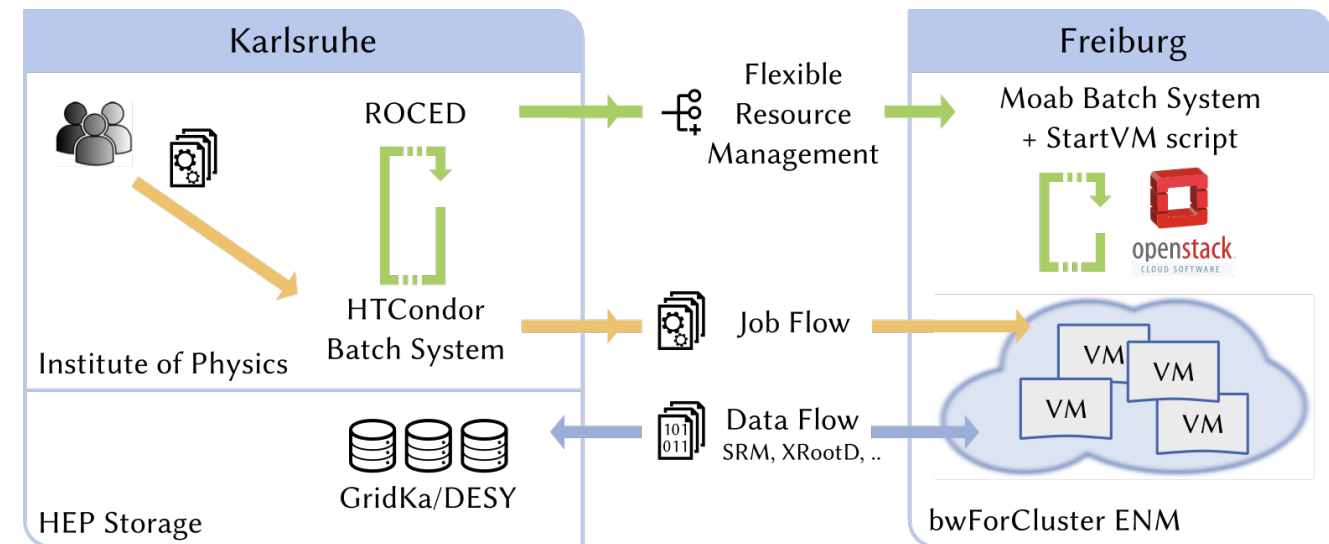
All dynamically and transparently integrated in a single HTCondor instance using ROCED and CoBaID/TARDIS



<sup>1</sup>COBaID - 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 Max Fischer (KIT)

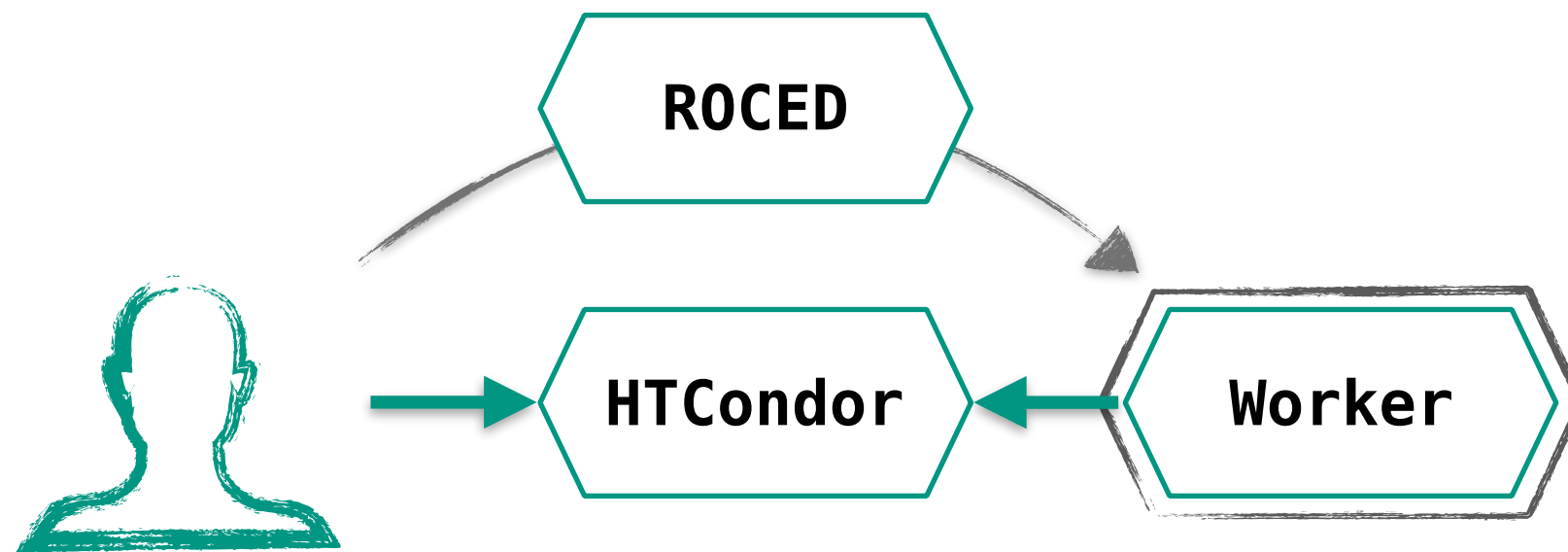
# Success Story - Opportunistic “Tier 1“ for a Day

- Dynamically shared HPC Centre at Freiburg (three diverse communities)
- Virtualization is key component to:
  - Allow dynamic resource partitioning
  - Meet OS & software requirements
- ROCED cloud scheduler developed at KIT
  - On-demand resource provisioning
  - Transparent resource integration
- Suitable for CPU-intensive workflows
- I/O-intensive workflows need data caches and fast networks (more later!)



Very good experience, however HEP was involved in the project from the early beginning

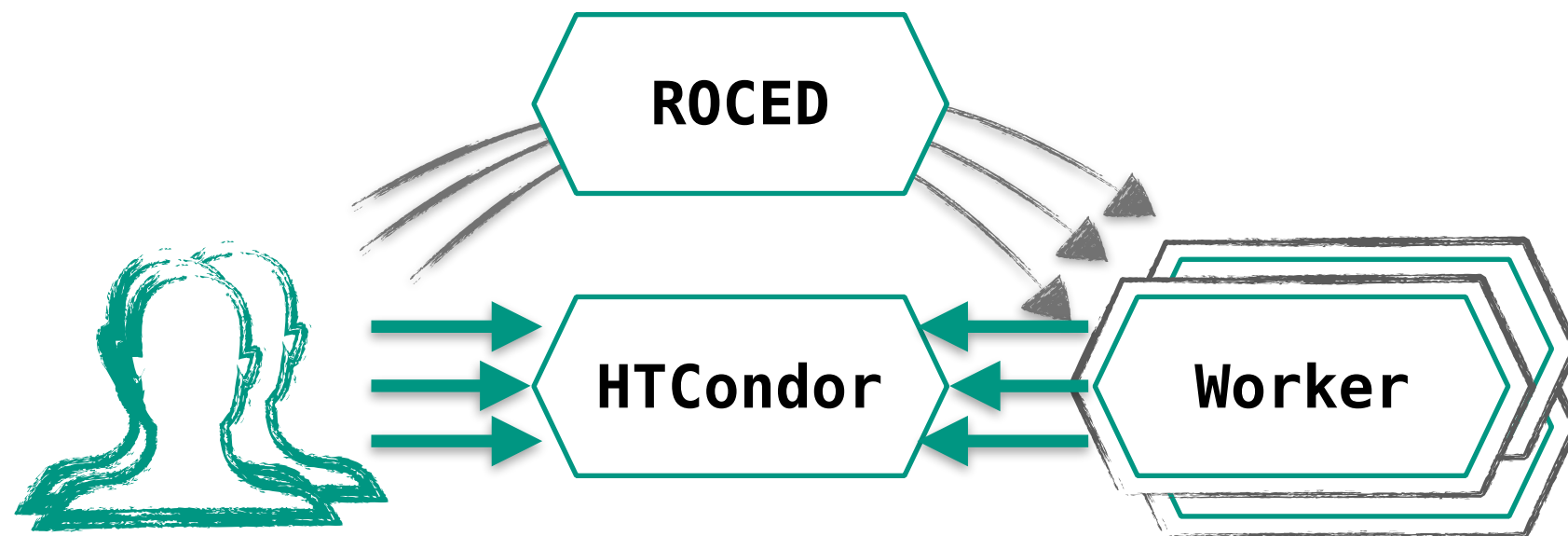
# ROCED Resumé



Slide by Max Fischer (KIT)



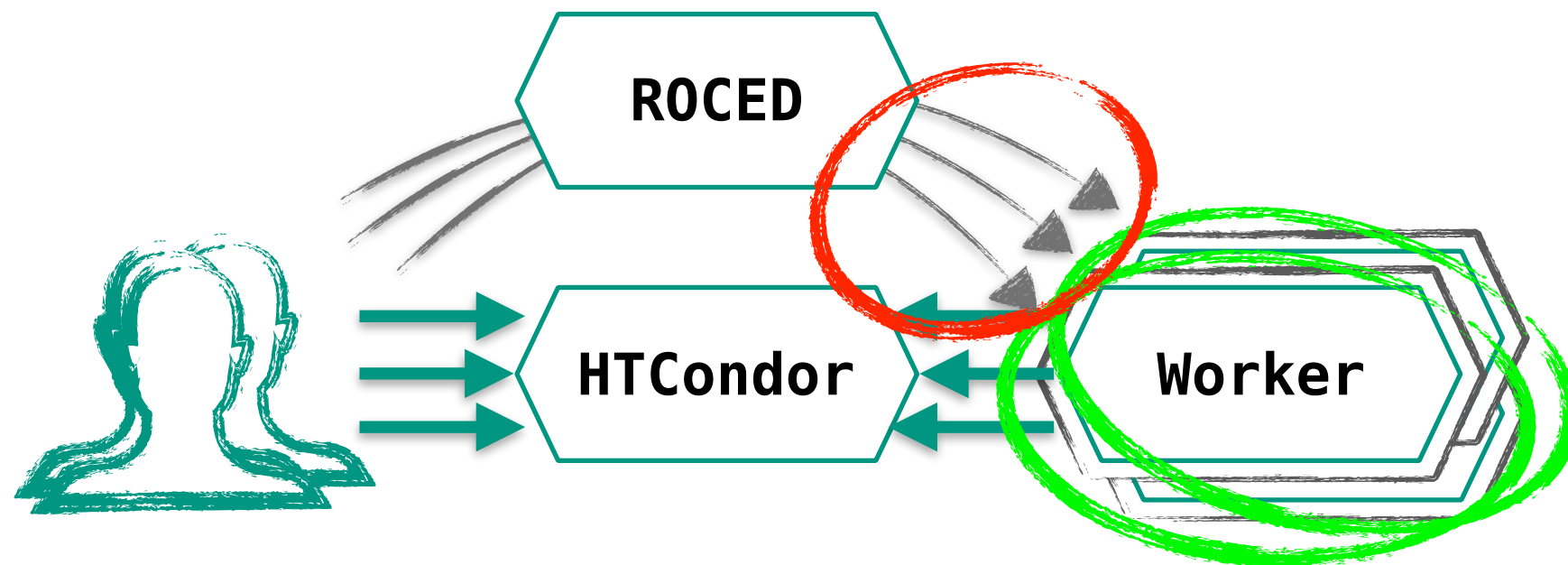
# ROCED Resumé



Slide by Max Fischer (KIT)

# ROCED Resumé

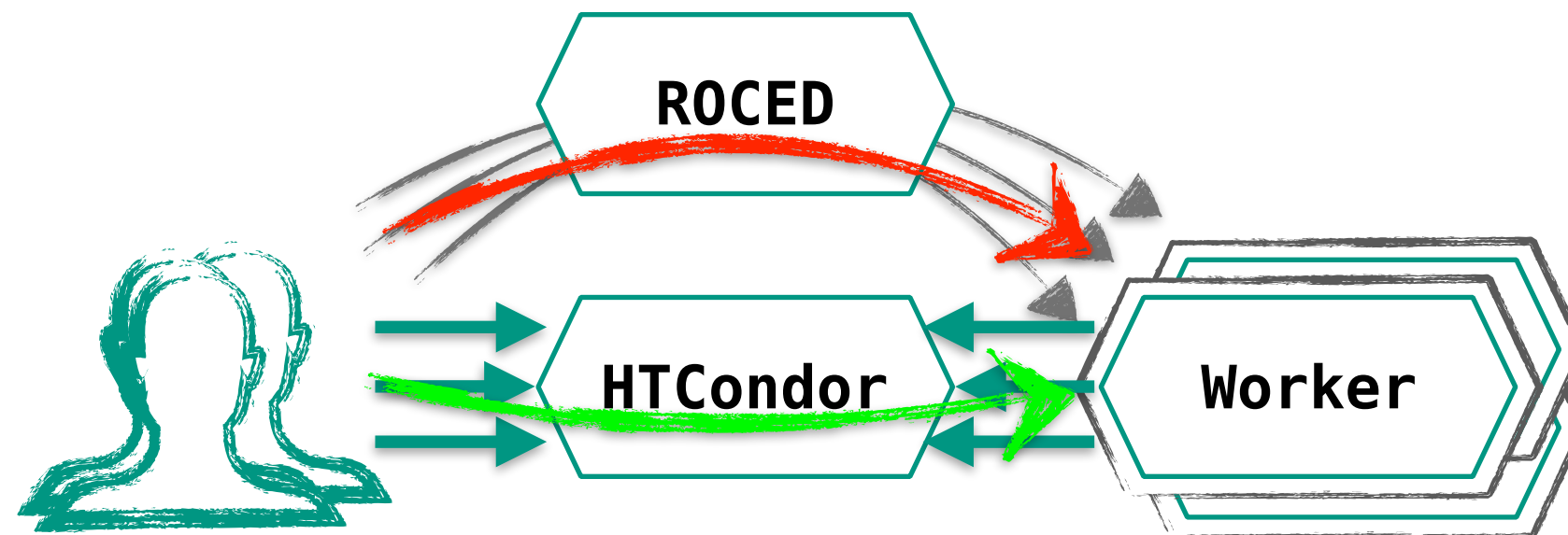
- Dynamic resources matching user demand
  - Trivial to support **new providers** for many users
  - Difficult to manage **several providers** for many users



Slide by Max Fischer (KIT)

# ROCED Resumé

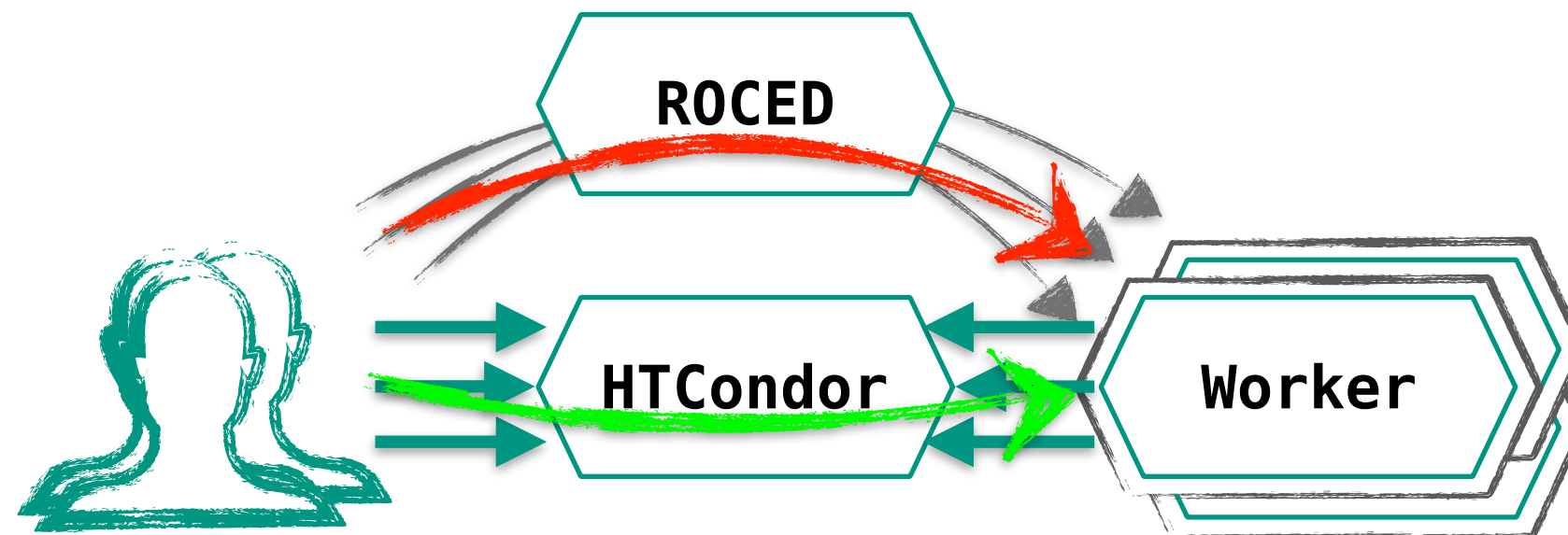
- Dynamic resources matching user demand
  - Trivial to support **new providers** for many users
  - Difficult to manage **several providers** for many users
- Resource aggregation in overlay batch system
  - Unreliable to **predict** resources required for jobs
  - Efficient to **integrate** resources,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



Slide by Max Fischer (KIT)

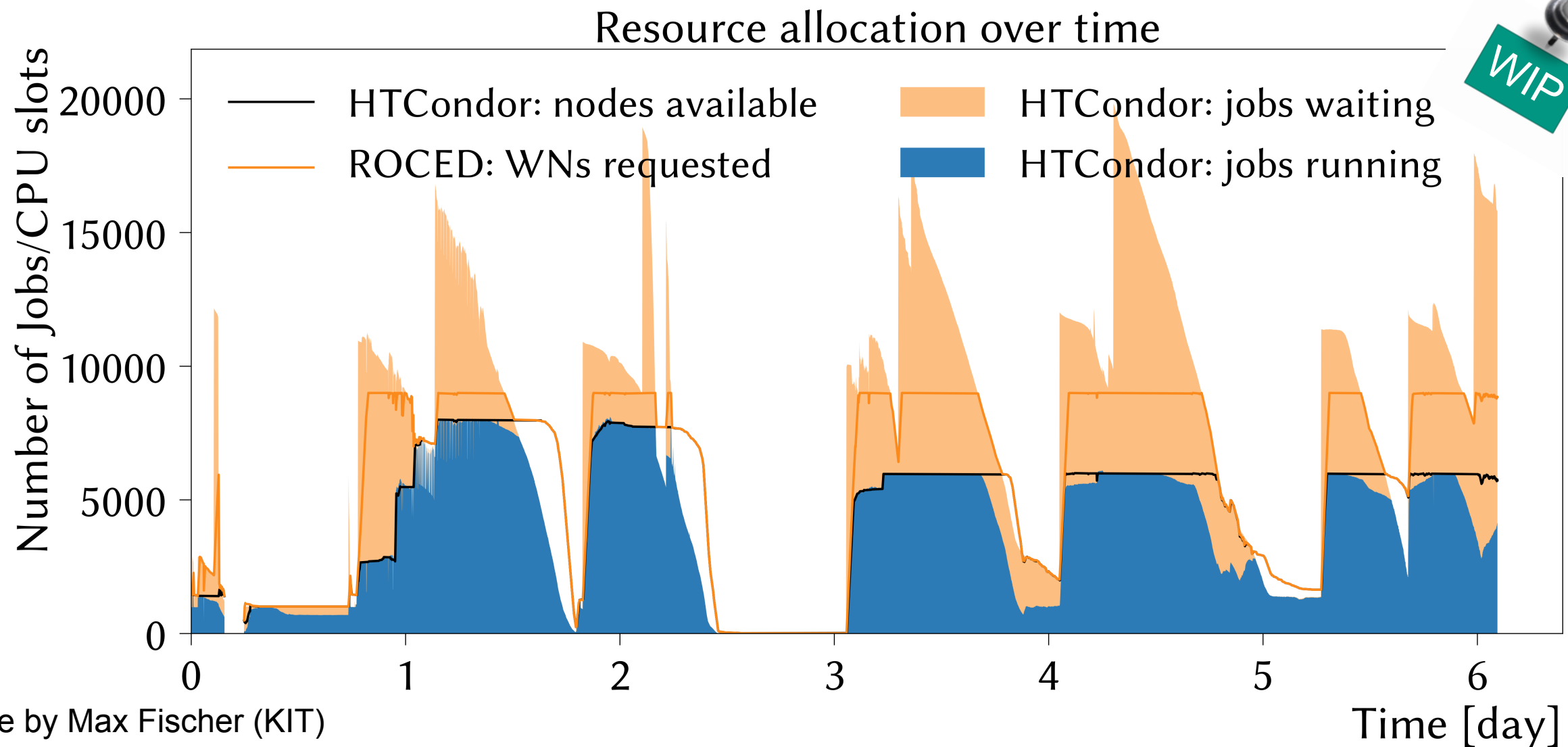
# ROCED Resumé

- 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)

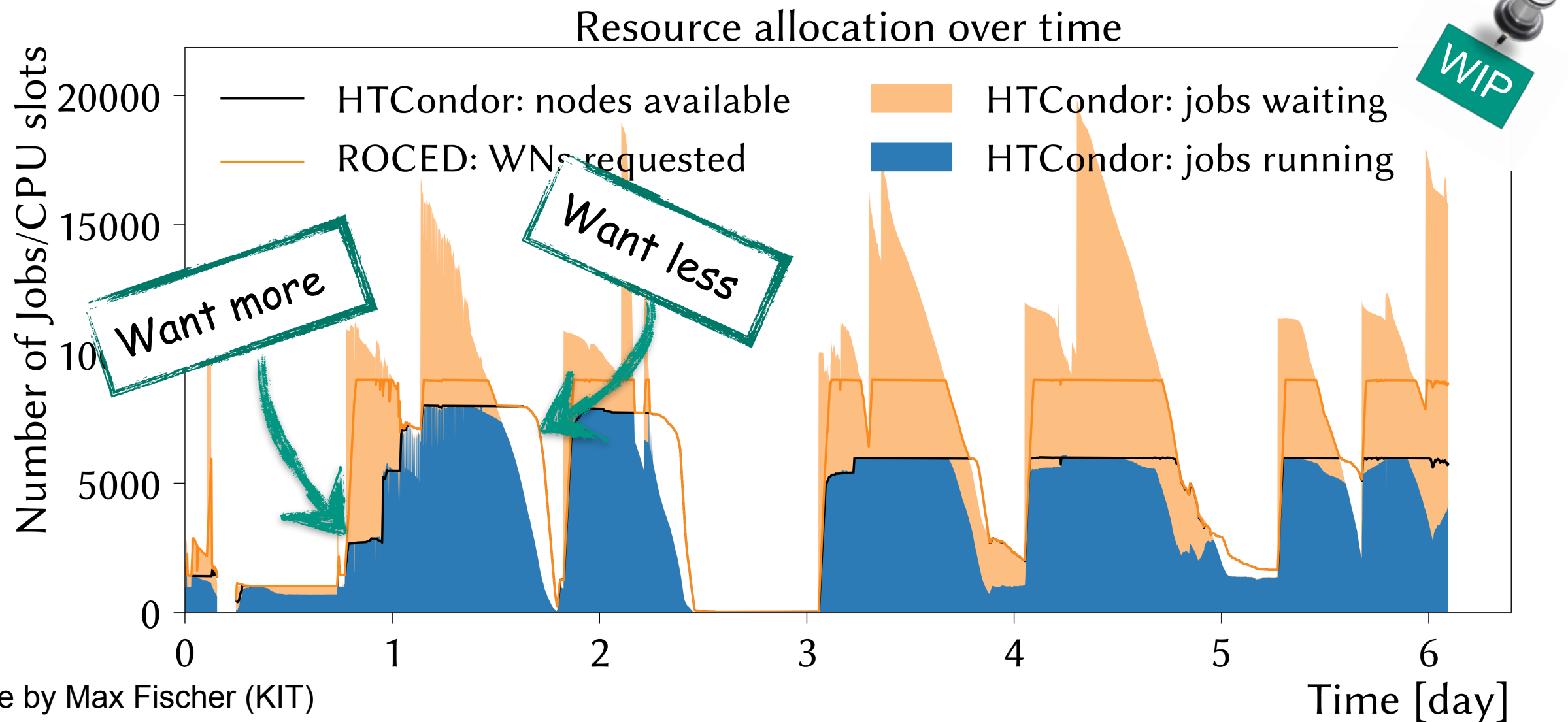
# Pragmatic View to Resource Management



Slide by Max Fischer (KIT)



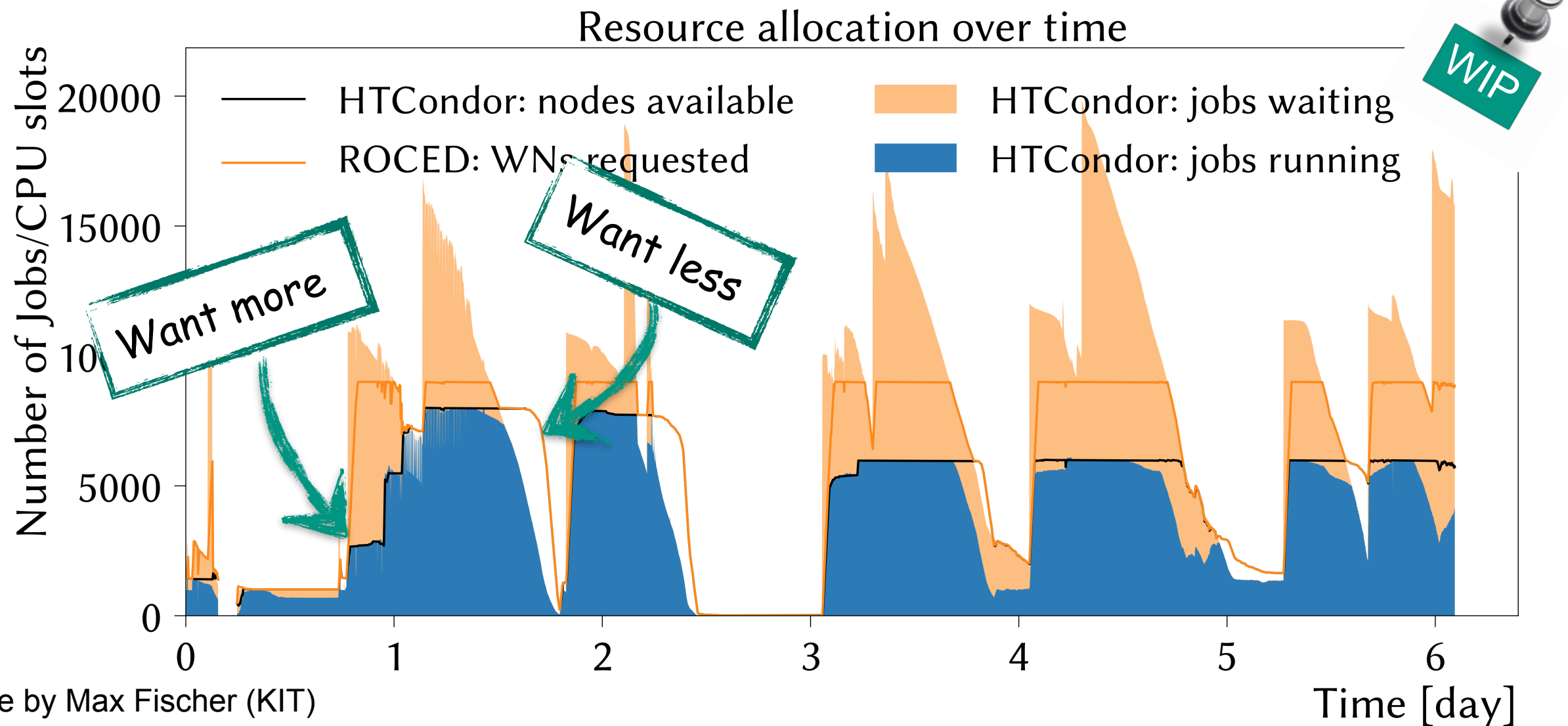
# Pragmatic View to Resource Management



Slide by Max Fischer (KIT)

# Pragmatic View to Resource Management

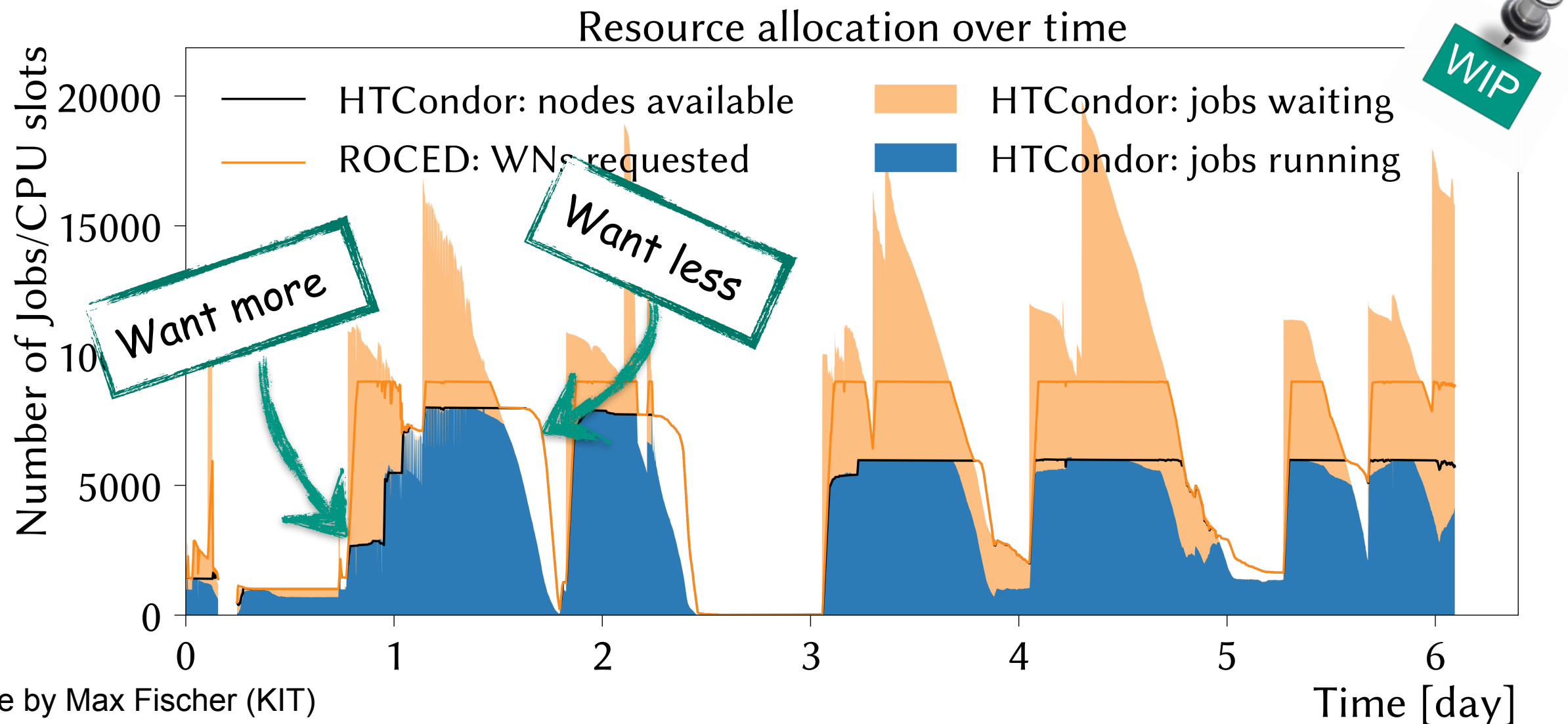
- New development for scalability and maintainability in HNSciCloud



Slide by Max Fischer (KIT)

# Pragmatic View to Resource Management

- New development for scalability and maintainability in HNSciCloud
- Simple logic: **more used** resources, **less unused** resources
  - COBalD only watches, creates and disables resources
  - Batch system scheduler selects appropriate resources



Slide by Max Fischer (KIT)

# Innovative Digital Technologies for Exploring Universe and Matter

- Joint proposal by HEP, Physics of Hadrons and Nuclei, Astroparticle Physics
- Covered Topics:
  - Development of technologies to utilize heterogeneous computing resources (Integration of Opportunistic Resources, Caching Technologies, Workflow Management)
  - Application and testing of those technologies in heterogeneous 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)

Proposal of:



FRIEDRICH-ALEXANDER  
UNIVERSITÄT  
ERLANGEN-NÜRNBERG



GEORG-AUGUST-UNIVERSITÄT  
GÖTTINGEN



Universität Hamburg  
DER FORSCHUNG | DER LEHRE | DER BILDUNG



Associated Partners:

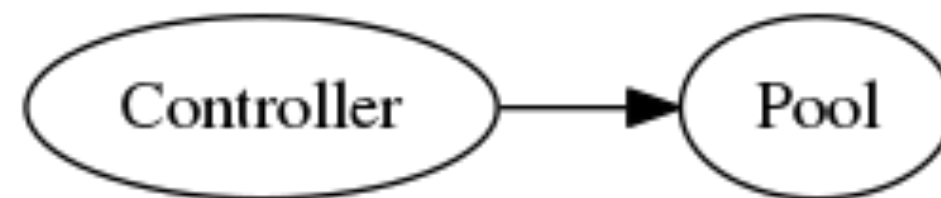


WWU  
MÜNSTER

# Resources

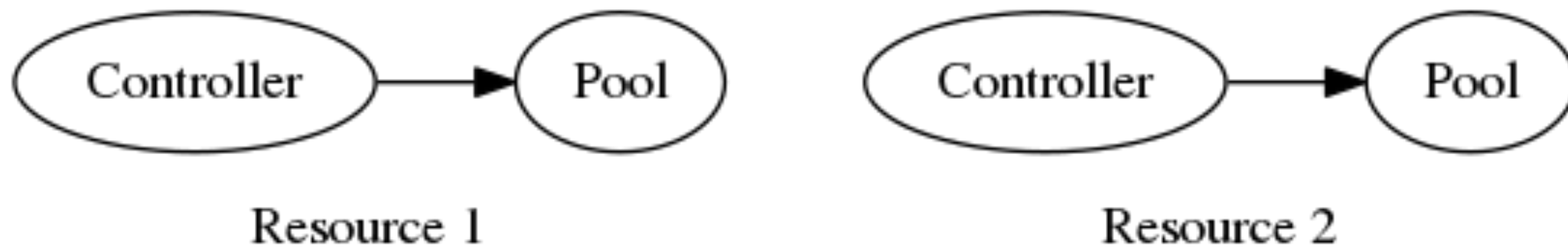
- COBaID: <http://cobald.readthedocs.io/>
- ROCED: <https://github.com/roced-scheduler/ROCED>
- TARDIS: <http://cobald-tardis.readthedocs.io/>
- COBaID Simulation: [https://git.scc.kit.edu/fq8360/cobalt\\_sim](https://git.scc.kit.edu/fq8360/cobalt_sim)
- COBaID Demo: [https://github.com/MaineKuehn/cobald\\_demo](https://github.com/MaineKuehn/cobald_demo)

# COBaID Resource Pool Model





# COBaID Resource Pool Model



# COBaID Resource Pool Model

