

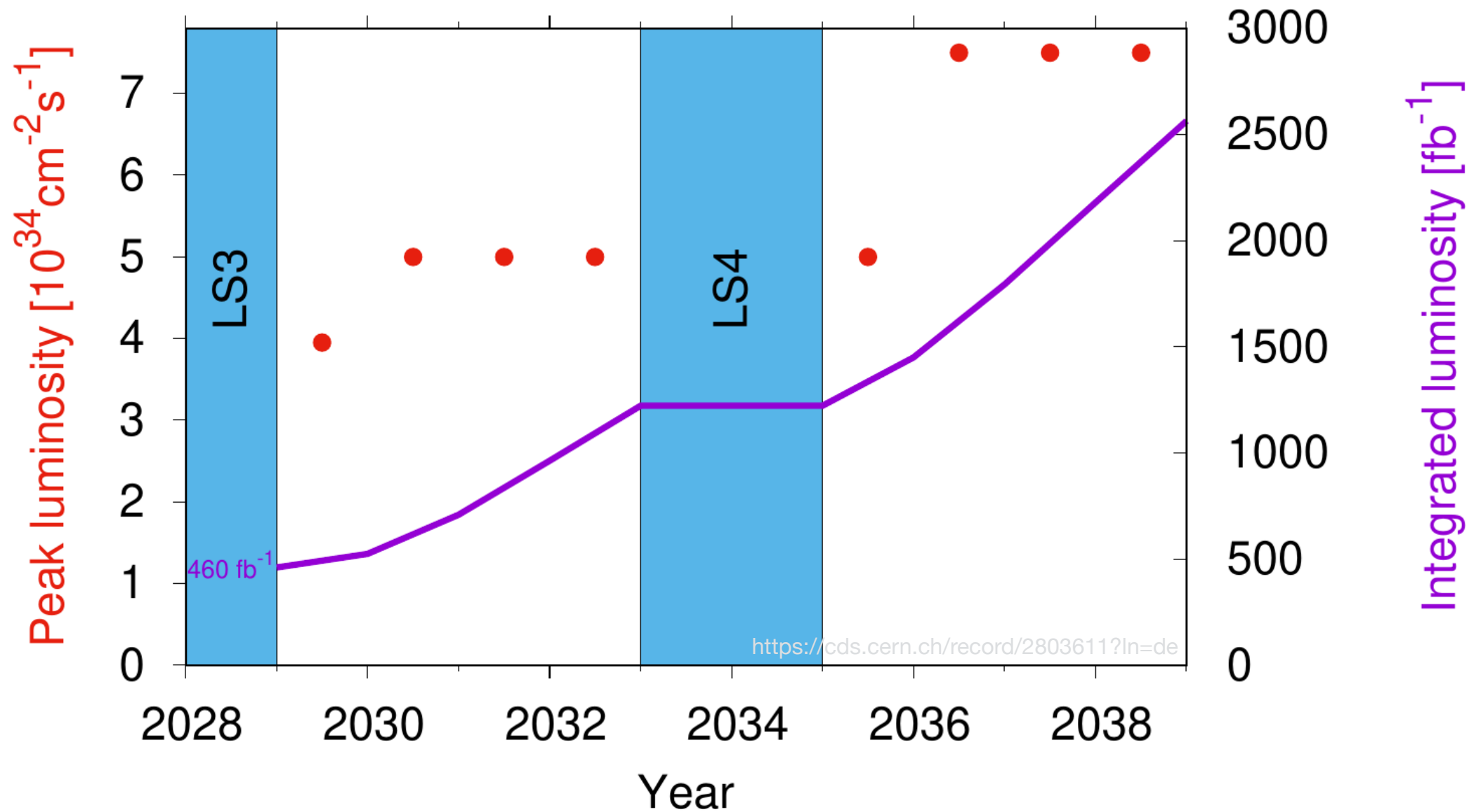
Transparent extension of the WLCG using the non-HEP side JURECA as an example

FIDIUM Collaboration Meeting 2022

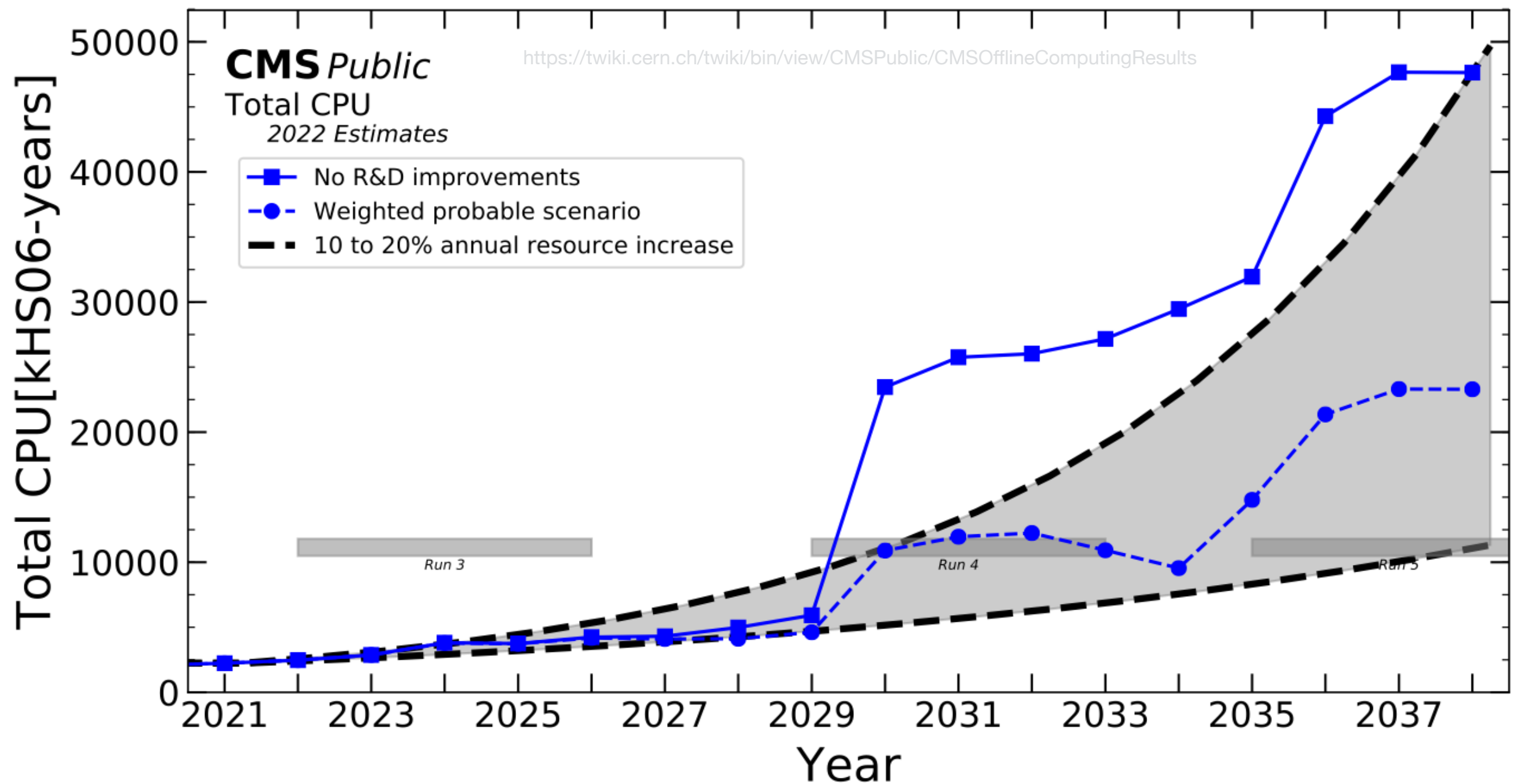
20.10.22

Manuel Giffels, Alexander Jung, Thomas Kreß, Thomas Madlener,
Andreas Nowack, Alexander Schmidt, Christoph Wissing





- Estimated luminosity for Run 4 and beyond
- HL-LHC: challenges in the areas of data acquisition, processing, simulation, and analysis



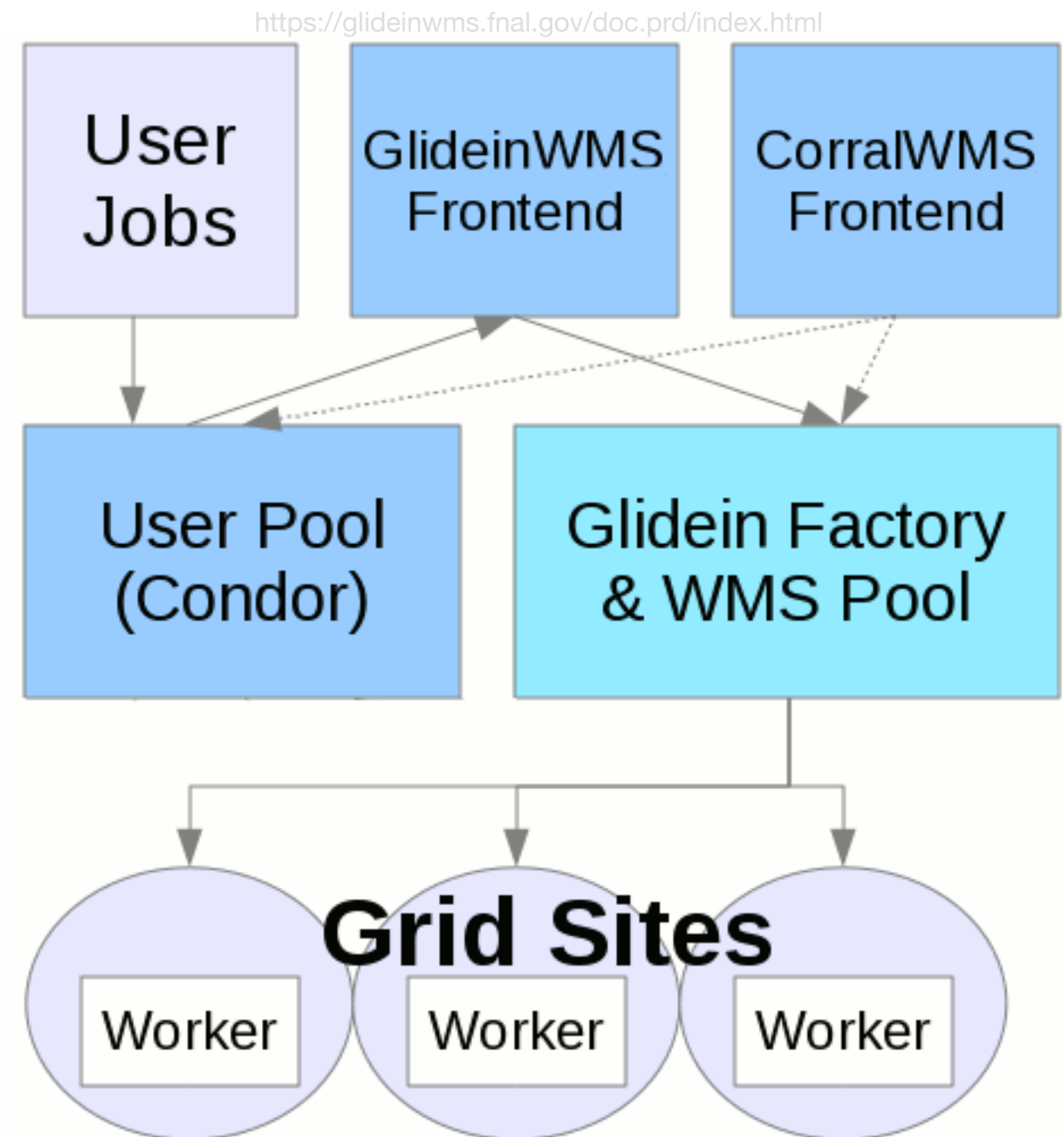
- Grid Resources may no longer be sufficient as of Run 4
- We need additional (opportunistic) resources

- Transparent for end users (integration into global CMS collector)
- Dynamic, without knowing/predicting future demand
- Expected WLCG environment is provided

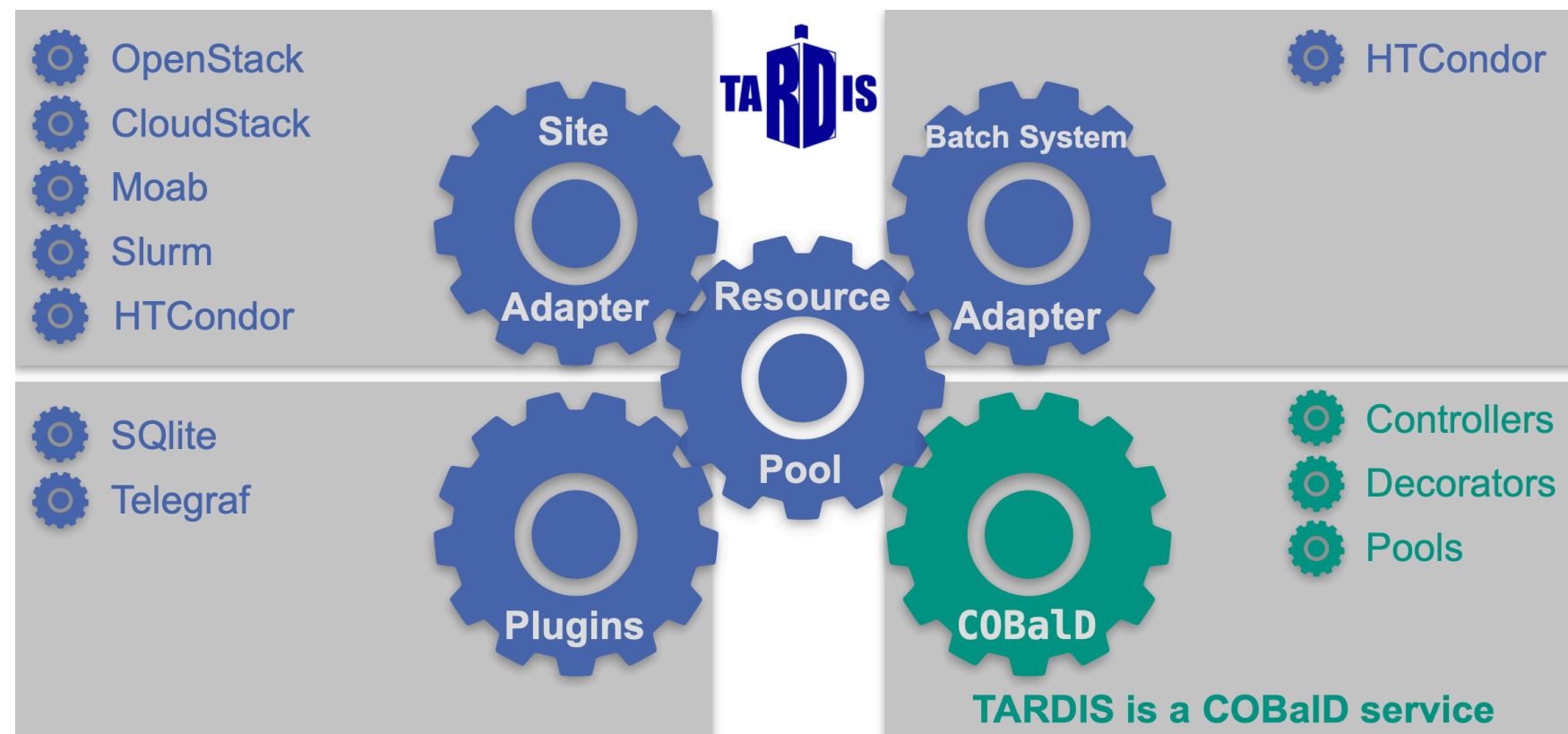
5 Expected setup by grid jobs

- Connection to external site to store output
- CVMFS present to grant access to CMSSW
- Access to conditions database
- Potentially access to input data

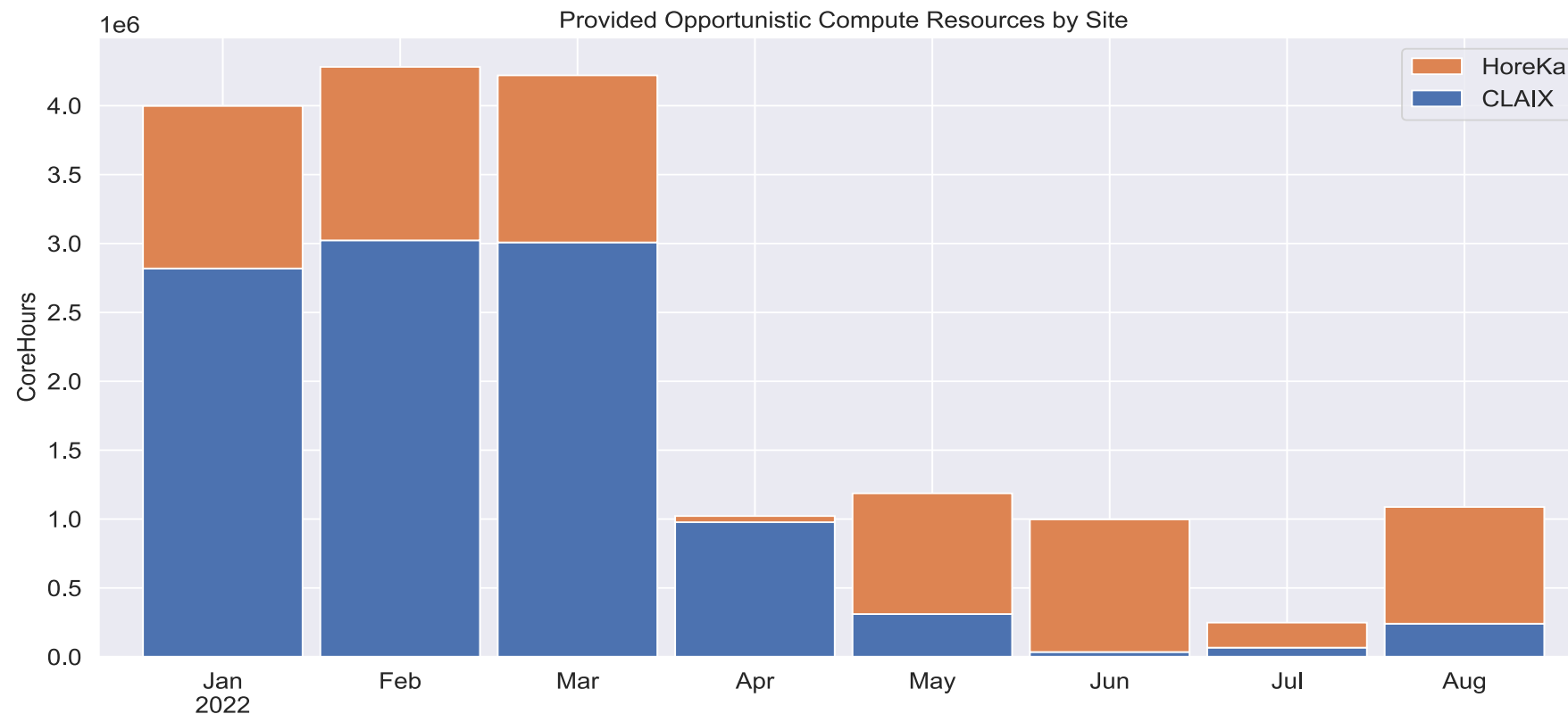
- Placeholder jobs bind resources
- Environment is provided by a Singularity container
- Integration into overlaying batch system (OBS)
- OBS: central management of the actual jobs



- Automate the startup of glideins (manual still possible)
- COBald = Opportunistic Balancing Daemon
 - Monitors allocation and utilization of resources (e.g. RAM, CPUs...)
 - Forms abstract metrics
 - Requests more glideins or lets existing ones expire, based on metrics
- TARDIS = Transparent Adaptive Resource Dynamic Integration System
 - Interface between jobs, external resources, batch systems and COBald



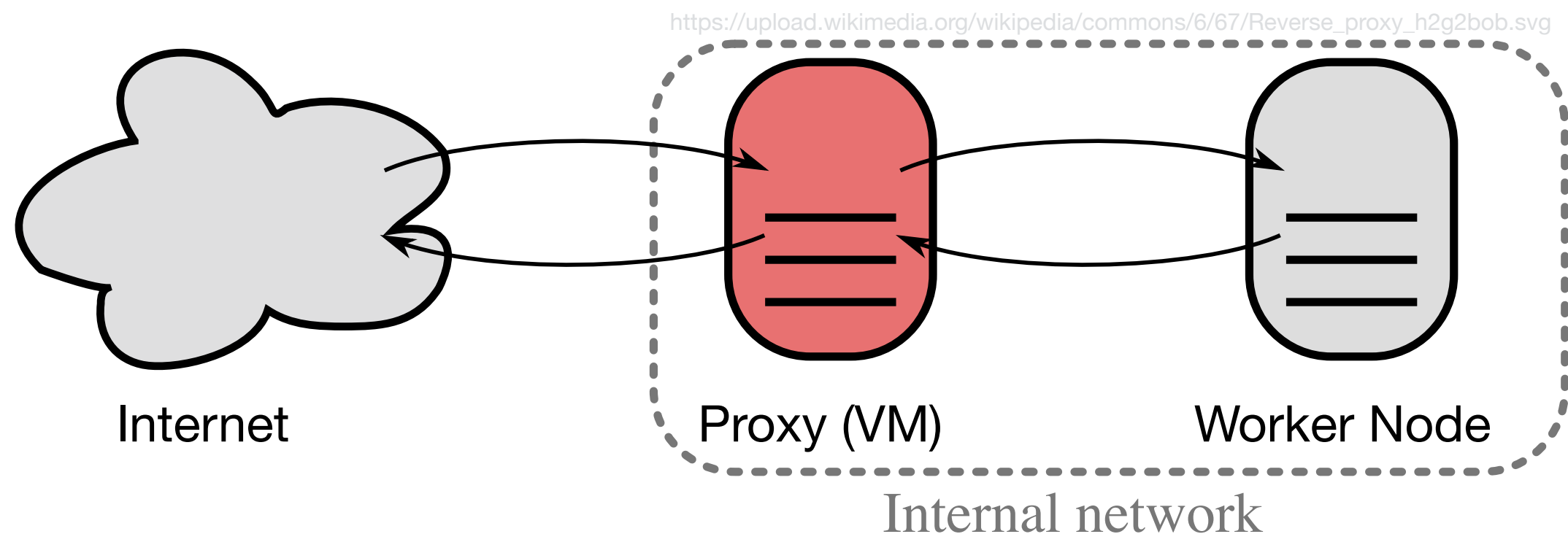
- Already in use to extend German Tier 1 and Tier 2 to local HPC Centres
 - KIT's Tier 1 extended to HoreKa
 - RWTH's Tier 2 extended to CLAIX
- } ~18 million corehours provided to CMS this year so far



- Work in progress: Making the HPC Centre JURECA located at Forschungszentrum Jülich (FZJ) available
- 3 major challenges arise

Strict firewall (almost)

- Worker nodes have no connection to the outside world
- Worker nodes route all traffic via a virtual machine (VM) in the HDF Cloud
- Bandwidth Worker Nodes \leftrightarrow VM: $\sim 10\text{Gbit/s}$
- Usage of ssh and proxychains from Worker Node (client) to VM (host)

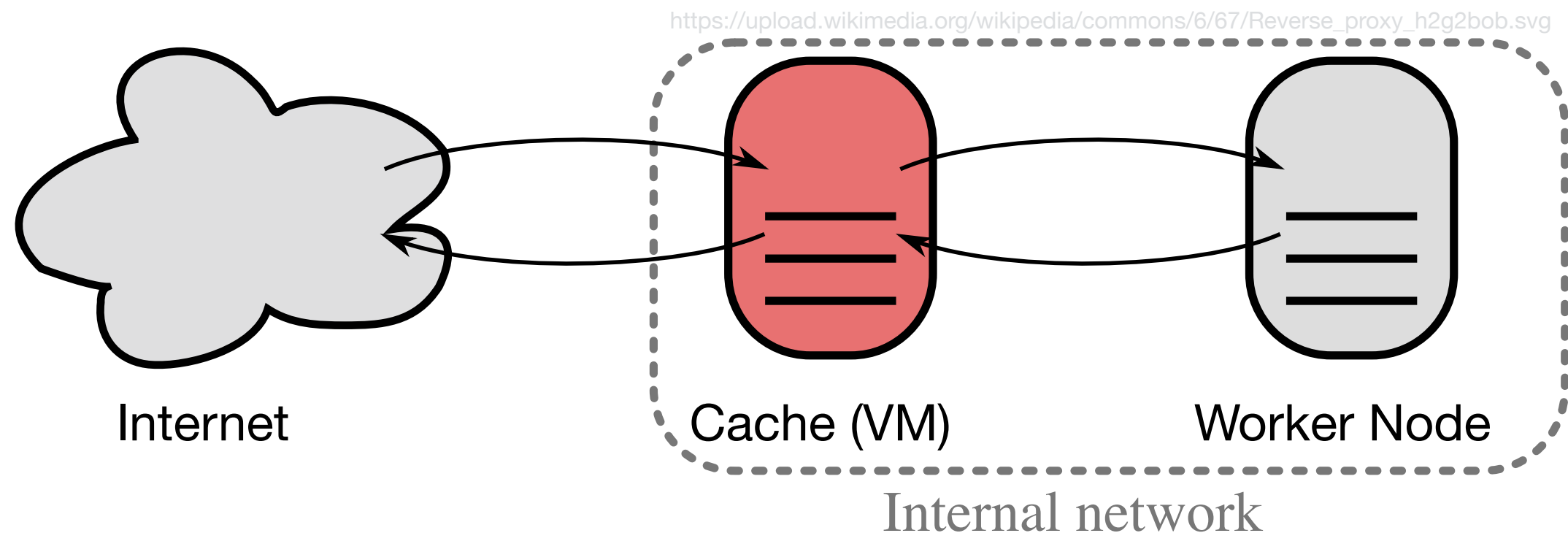


Environment (✓)

- CVMFS not present, but provided by prepared cvmfsexec tarball
- User jobs run inside singularity container
- Bind-mount unpacked tarball to /cvmfs/ inside container
- Bind-mount additional needed stuff inside container

Caching (✓)

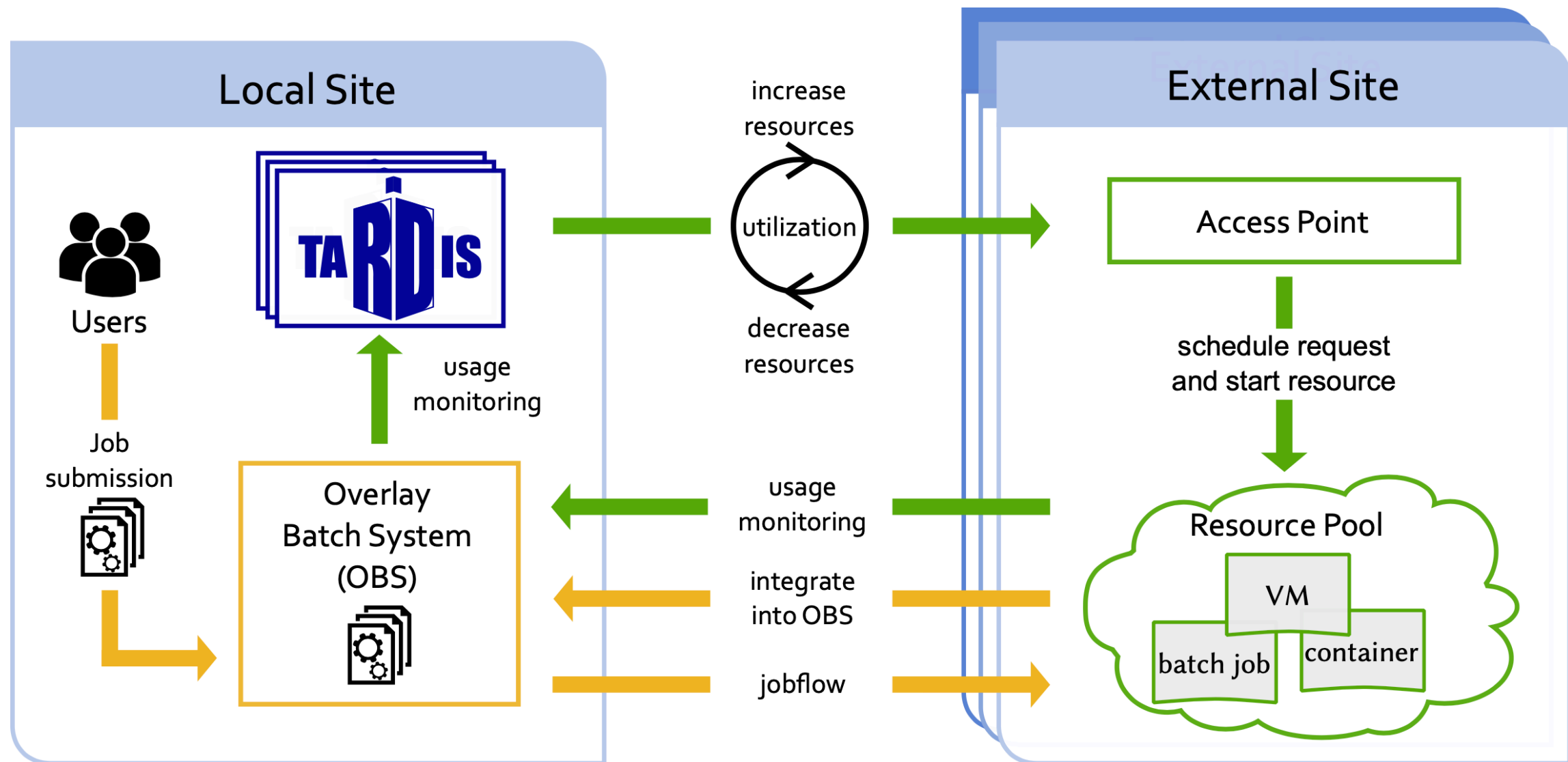
- CVMFS traffic handled by Squid Proxy running on the VM
- Each glidein sets up its own cache on the worker node and can serve up to 256 payloads with it
- Caching of conditions database by Squid Proxy on VM
- Once cache is established: little traffic to the outside, much internal traffic



- It works so far, but...
- ...connection glidein <-> HTCondor pool unstable
- Seemingly random outages between a few minutes and ~1h, connection comes back again
- Payload not effected

- Fixing the connection issues
- Ask the responsible admins for an exception in the fire wall
- Could maybe fix the connection problems
- Scale test with fully working setup/stress test

Backup



- Set I contains resource types (RAM, CPU,...)

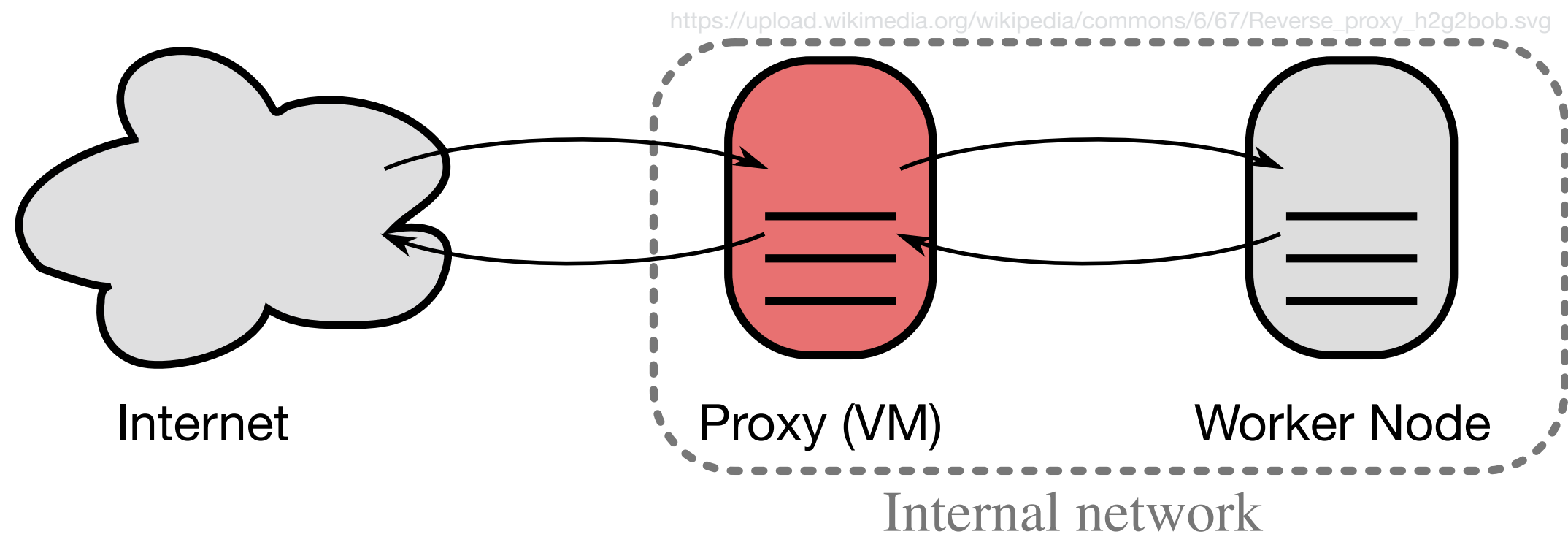
- Allocation: measures how full a resource is/whether jobs still fit on it

$$\text{allocation} = \max_{i \in I} \left(\frac{\text{used}(i)}{\text{requested}(i)} \right)$$

- Utilization: represents the suitability of the resource for the jobs (high utilization hints a good fit between jobs and requested resource)

$$\text{utilization} = \min_{i \in I} \left(\frac{\text{used}(i)}{\text{requested}(i)} \right)$$

- Cloud computing platform offered by Jülich Supercomputing Centre (JSC)
- Resources produced as part of the Helmholtz Data Federation
- Provides access to the HPC file system among other things



- Connection made from port 5555 on Worker Node is forwarded to port 5432 on VM (inside container)
- Access to and caching of conditions database and CVFMS by Squid Proxy on VM via port 3306 (outside container)

- Based on standard Squid http proxy cache software
- Contains configuration defaults and bug fixes that are known to work well with the applications used on the grid
- Pre-configured for use by the Frontier distributed database caching system
- Used by ATLAS and CMS to access conditions data

- Setup ssh tunnel „Setup“ on the worker node
- Mount CVMFS repositories
- Launch singularity container and bind-mount CVMFS repositories to /cvmfs

- Enter the now running container

- Launch glidein Work that is done inside the container
- ...
- Profit

- Code available at <https://gitlab.desy.de/thomas.madlener/cms-drp-jsc>