

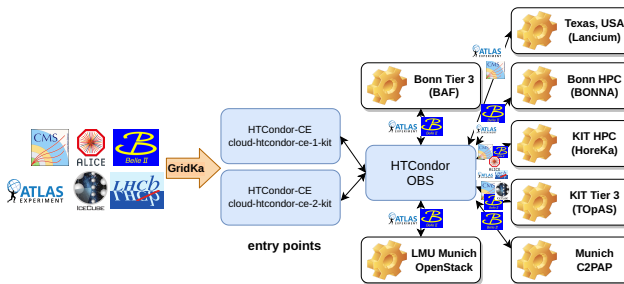
Data-Aware Scheduling with Hash-Based Data Distribution

M. Giffels, A. Gottmann, **R. Hofsaess**, M. Horzela, G. Quast, M. Schnepf | 30.03.2023

FTS and XRootD Workshop @ JSI (27-31 March 2023)

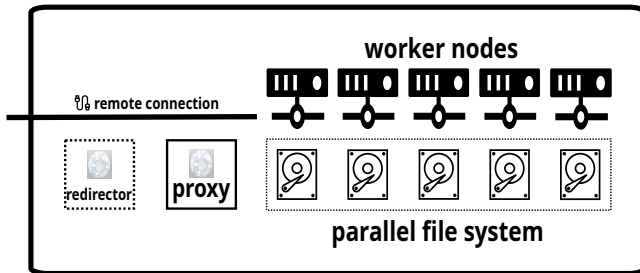


Integration of Opportunistic Resources



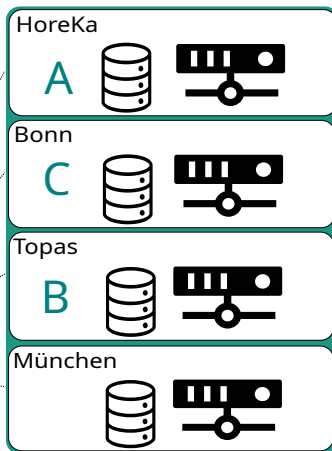
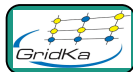
- Opportunistic resources are integrated dynamically with COBaID/TARDIS
- Very heterogeneous infrastructure
- No permanent, managed storage at integrated sites – only caches

Opportunistic Resource

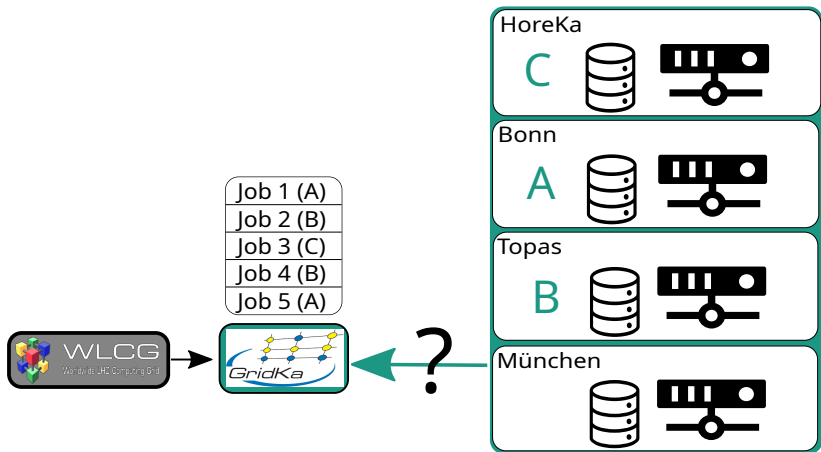


Data-Aware Scheduling for ORs

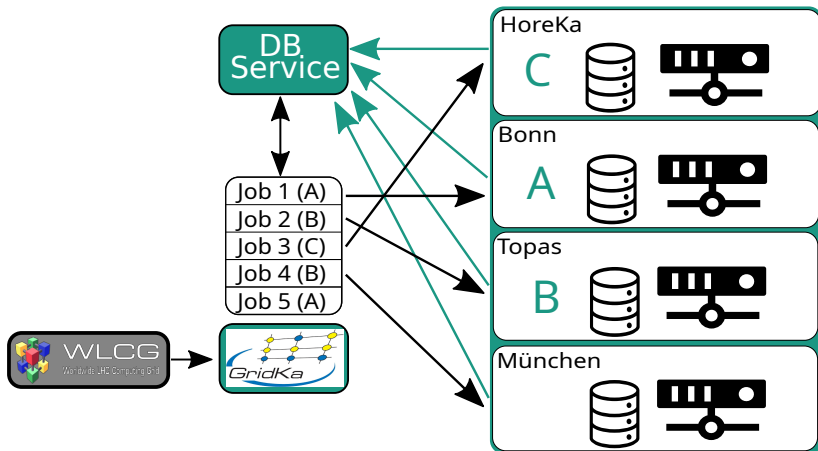
Sites are integrated
dynamically with
COBaID/TARDIS



Data-Aware Scheduling for ORs



Data-Aware Scheduling for ORs

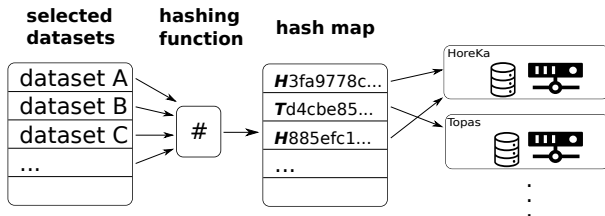


Hash-Based Approach

- Hash-based distribution of files inspired by [Ceph/CRUSH](#)
- based on the file name, a hashing function calculates the data placement

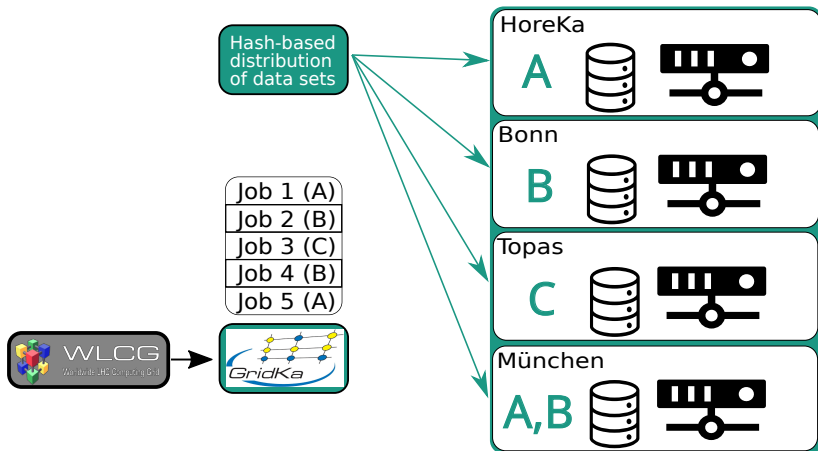
Hash-Based Approach

- Hash-based distribution of files inspired by [Ceph/CRUSH](#)
- based on the file name, a hashing function calculates the data placement

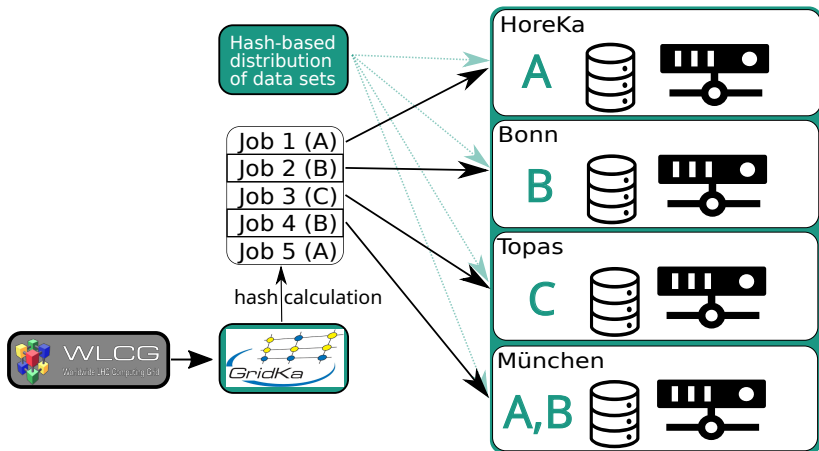


e.g.:
/SingleMuon/Run2022C-PromptNanoAODv10-v1/NANOAOB

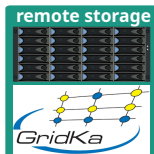
Hash-Based Distribution of Datasets



Hash-Based Distribution of Datasets

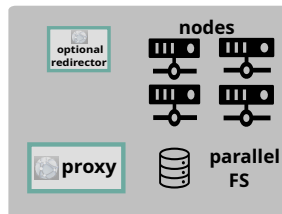
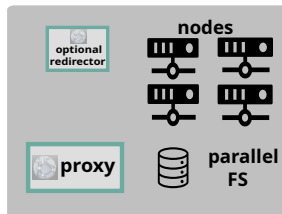
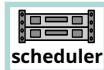


Concept and Implementation (I)



Remote storage
at GridKa initially
provides the data

HTCondor OBS



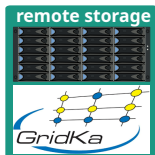
Introduction
○○

Data-Aware Scheduling
○○○

Setup
●○○○

Other Projects
○○○○○○

Concept and Implementation (I)

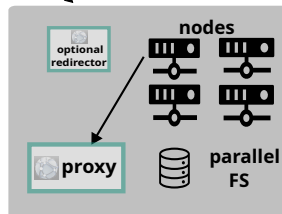
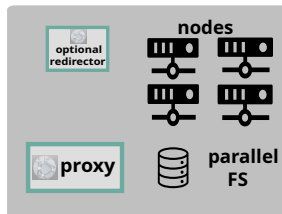


Remote storage
at GridKa initially
provides the data

HTCondor OBS



calculates hash and
distributes jobs



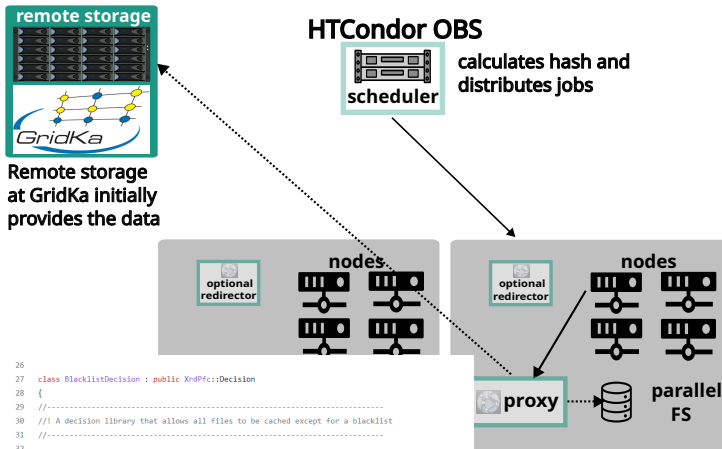
Introduction
○○

Data-Aware Scheduling
○○○

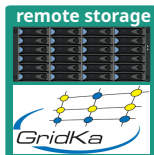
Setup
●○○○

Other Projects
○○○○○○

Concept and Implementation (I)

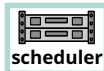


Concept and Implementation (II)

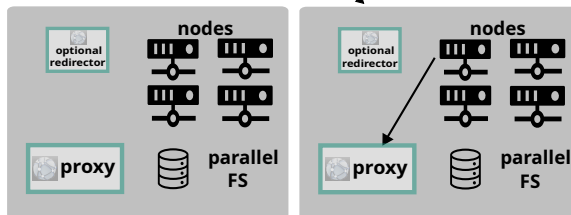


Remote storage
at GridKa initially
provides the data

HTCondor OBS



calculates hash and
distributes jobs



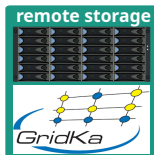
Introduction
○○

Data-Aware Scheduling
○○○

Setup
○○●○

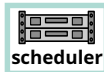
Other Projects
○○○○○○

Concept and Implementation (II)



Remote storage
at GridKa initially
provides the data

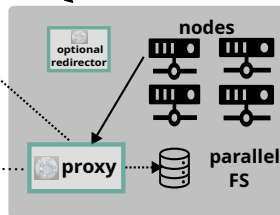
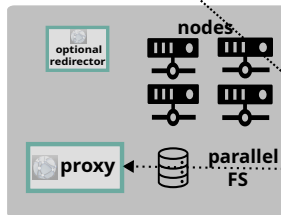
HTCondor OBS



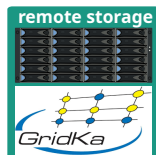
calculates hash and
distributes jobs

job gets scheduled

XRootD Caching
Proxy:
Implements the
hashing and
provides cached
data if available



Concept and Implementation (II)



Remote storage
at GridKa initially
provides the data

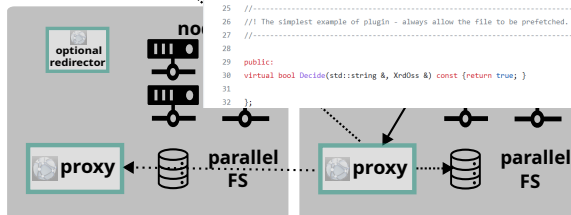
HTCondor OBS



calculates hash and
distributes jobs

job gets sched

XRootD Caching
Proxy:
Implements the
hashing and
provides cached
data if available



Next Steps

- Proof of Concept:
 - Setup a test system at our institute
 - There, we also integrate some external resources (NEMO@Freiburg and our own T3) into our HTCondor batch system
 - First: fixed file lists.
 - Later: hash-based data distribution
- Finally: adapt system to the opportunistic resources

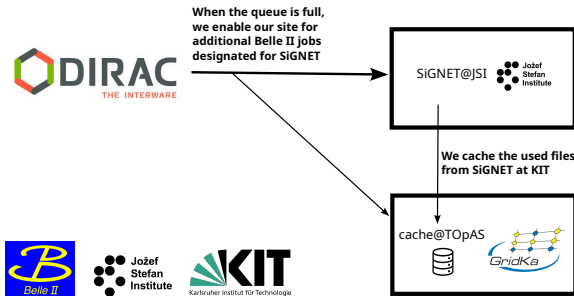
Summary

With our hash-based approach:

- We avoid actually keeping track of thousands of files on volatile caches
 - The caching decision is determined locally (or centrally)
→ **horizontally scalable**
 - With local on-the-fly caching, no prefetching is necessary (but possible)
 - No database service needed (or at least more lightweight, tbd)
 - Highly automatizable → **no active data management necessary**
 - Allows inclusion of HPC sites without remote connection of WNs
- reduction of data transfers and more efficient utilization of resources achievable

Belle II Caching@GridKa

Currently, we have a XRootD caching project with SiNET/JSI ongoing!
(I am not directly involved, credits go to: [Moritz Bauer](#) and [Matthias Schnepf](#))



Belle II Caching@GridKa

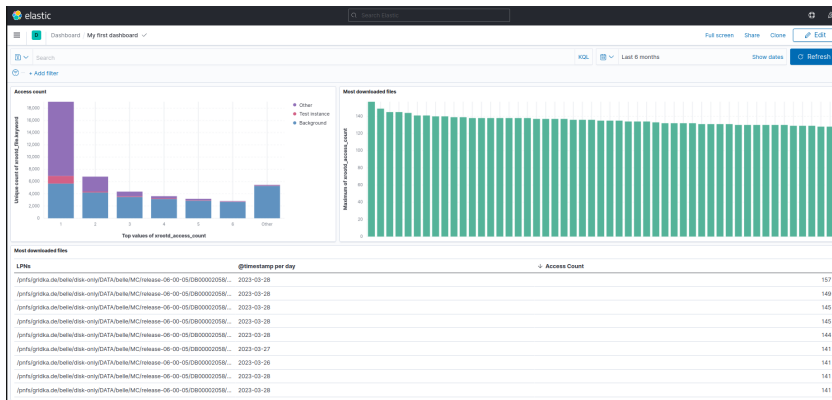
```
xrd.port 1094
all.export /xroot:/ nostage
all.export /root:/ nostage
all.export /https:/ nostage
all.export /davs:/ nostage
all.export /http:/ nostage

xrootd.chksum max 4 Adler32
ofs.osslib libXrdPss.so
pss.cachelib default
ofs.ckslib * libXrdPss.so
pss.origin =
oss.localroot /ceph/belle2-xrootd-cache/
pfc.ram 32g
pfc.diskusage 0.9 0.95
pss.debug

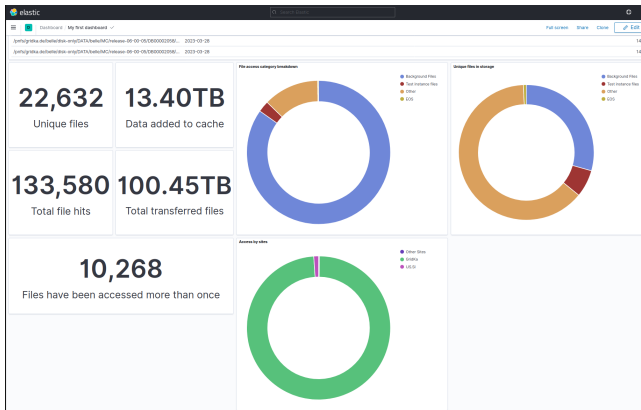
xrootd.monitor all flush 10s window 5s fstat 5 lfn gbuff 8k dest fstat info user pfc 129.13.101.187:9921
xrootd.mongostream pfc use flush 10s

continue /etc/opt/belle2-xrootd/config.d/
~
~
~
~
~
~
~
"b2-xrootd-cache-server.cfg" 21L, 537C 1,1 ALL
```

Belle II Caching@GridKa



Belle II Caching@GridKa



XRootD Interactive

Available at GitHub: [xrd-interactive](#)

- set of convenience functions to use xrootd interactively (as python [questionary](#))
- Intended to make the usage more easy for non-powerusers (bachelor/master students etc)
- It is not really matured, nor generalized... (mainly adapted to GridKa)
- No plans on further development, just a small gimmick

XRootD Interactive

```
rhofsaess@portal: ~/development x + -  
rhofsaess@portal1:~/development/xrd-interactive$ python3 xrootd_interactive.py -u rhofsaess -b /store/user/rhofsaess/  
? Please select a redirector: root://cmsxrootd-kit.gridka.de:1094/, (RW) [default]  
Redirector selected: root://cmsxrootd-kit.gridka.de:1094/  
Selected base path: /store/user/rhofsaess/  
Current base path: /store/user/rhofsaess/  
? What do you want to do? (Use arrow keys)  
» exit  
ls  
interactive ls  
stat  
stat directory  
dir size  
dir content  
rm file  
interactive file rm  
rm dir  
interactive dir rm  
mv  
mkdir  
copy file to  
copy file from  
create file list  
change base path  
change redirector  
help
```