

# First Deployment of XCache for Workflow and Efficiency Optimizations on Opportunistic HPC Resources in Germany

Robin Hofsaess<sup>1,\*</sup>, Manuel Giffels<sup>1</sup>, Artur Gottmann<sup>1</sup>, Günter Quast<sup>1</sup>, Andreas Petzold<sup>1</sup>, Matthias Schnepf<sup>1</sup>, and Achim Streit<sup>1</sup>

<sup>1</sup>Karlsruhe Institute of Technology, Karlsruhe, 76131, Germany

**Abstract.** The future German HEP computing strategy will partially rely on national scientific HPC centers for providing the CPU pledges to the WLCG. The pilot phase for this endeavor is starting in 2025 from which on the share provided by HPC resources will gradually increase. To make this courageous step a success, it is essential that the integrated HPC sites can be utilized just as reliable and efficient for HEP workflows as the traditional, dedicated Grid sites. Motivated by I/O limitations observed at HoreKa, the scientific HPC cluster at KIT, the integration was optimized and an XRootD-based approach for data access bottleneck mitigation was developed and deployed, resulting in a comparable performance of the cluster. This is an important result supporting the future HEP computing strategy in Germany and shows that national scientific HPC centers are capable of delivering the performance and reliability required to be used for providing pledges to the WLCG. Moreover, the findings underline the potential of integrating HPC clusters into the Grid, paving the way for a scalable and sustainable computing infrastructure in the HL-LHC era.

## 1 Introduction

In 2022, a novel High-Energy Physics (HEP) computing strategy in preparation for the HL-LHC era was decided by the Committee for Elementary Particle Physics [1] in Germany. This strategy focuses on a more efficient and sustainable computing environment and therefore foresees a gradual consolidation of storage and compute resources. The university Tier-2 centers currently contributing to the Worldwide LHC Computing Grid (WLCG) are planned to be replaced in this process. While storage will be taken over by the Helmholtz Centers, Karlsruhe Institute of Technology (KIT) and Deutsches Elektronen-Synchrotron (DESY), for the contribution of compute power, the strategy includes a complete novelty: Official WLCG pledges will be provided with shares on scientific High-Performance Computing (HPC) centers within the *National High Performance Computing Alliance* [2]. While the opportunistic integration of HPC centers is already today making an important contribution to the WLCG, the 'pledged' integration of such resources is entirely new – as a fixed pledge will be guaranteed with resources that are neither fully controlled nor solely dedicated to the WLCG.

---

\*e-mail: Robin.Hofsaess@cern.ch

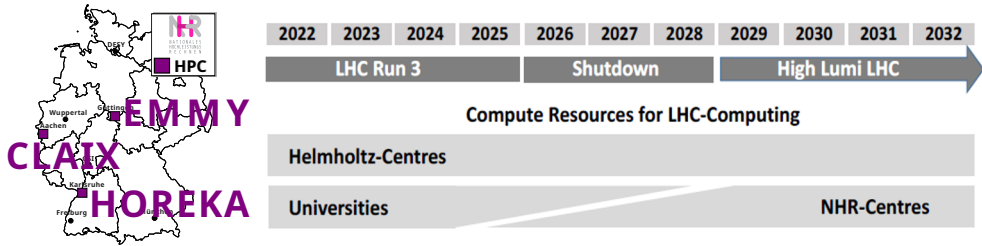


Figure 1: The German HEP computing strategy foresees a gradual transition of compute resources from the university Tier-2 centers to the Helmholtz Centers, KIT and DESY, and scientific HPC clusters within the NHR Alliance – namely CLAIX, EMMY, and HoreKa. *Adapted from:* Markus Schuhmacher/[1].

For the beginning, three clusters of the NHR Alliance – CLAIX (Aachen), EMMY (Göttingen), and HoreKa (Karlsruhe for Freiburg and Wuppertal) – were selected that are collocated to the university Tier-2 centers that will be replaced in the long term, as indicated in Fig. 1. To increase cost effectiveness, the gradual transition allows to still use the available Tier-2 resources during the transition process, but instead of buying new hardware, the HPC share is increased over time. It is important to mention that this transition away from dedicated university WLCG sites does not mean that the Tier-2 groups are replaced. On the contrary, these groups will take responsibility for the reliable and efficient integration of the HPC clusters and will be the direct link between the centers and the WLCG. They will therefore undergo a shift from the pure provision of resources to a research and development focus for efficient integration as well as experiment and user support [1].

In general, this strategy is motivated by the expected benefits of these resources for the HEP community. The NHR centers provide an enormous performance potential with modern, scalable hardware, offering enough compute power to meet the growing demand of HEP research in the future. Together with highly efficient cooling systems (see e.g. Ref. [3]), this means a significant increase in performance and sustainability compared to operating several smaller, dedicated WLCG computing centers with typically longer hardware life cycles. Furthermore, the joint utilization of the NHR resources not only by different LHC experiments but also between users from various scientific fields fosters collaboration between different science communities and allows to share innovations.

The utilization of HPC resources therefore can pave the way for a more sustainable and collaborative computing infrastructure in Germany – provided they can be used efficiently. This is, of course, inevitably necessary for the success of the future HEP computing strategy. The experience from the last years with the utilization of HoreKa has shown, this is not trivial to accomplish. The NHR center in Karlsruhe is one of the pilot centers for the future HPC strategy and is used as an opportunistic resource since 2021. From the beginning, it served as a testbed and helped to explore and improve the utilization of HPC centers for HEP workflows. Based on the experience with HoreKa, challenges and limitations were identified – as discussed in the next section – and possible improvements for a more efficient utilization were explored – described afterwards.

## 2 Challenges and Limitations

Using the NHR centers for official WLCG compute tasks and HEP workflows in general introduces various challenges. First of all, the transition to the HPC centers goes hand in hand with a transition from administrative access on the dedicated resources to simple, non-privileged user access. This needs to be considered when integrating the centers into the WLCG and introduces additional prerequisites, such as availability of container technologies, to realize the integration and to provide the standardized Grid and experiment software. In the German HEP computing environment, this is often realized with COBaD/TARDIS [4], a software tool that allows a dynamic and transparent integration of heterogeneous resources into one overlay batch system, which then can be filled with WLCG computing jobs. A more detailed explanation can, e.g., be found in Ref. [5]. Thanks to the tool, the utilization of the NHR centers is possible, but not only the mode of operation as a user is different, but also the resources themselves. While the dedicated Grid sites forming the WLCG are especially tailored to fit the needs of High-Energy Physics research and provide a standardized environment, HPC resources typically have a different focus and may not be inherently suited for HEP workloads. The computing model of multi-purpose HPC clusters aims at high-parallelized, low-latency calculations, often spanning many individual compute nodes. However, high data throughput – as required by many particle physics workflows – is not usually a priority. Should more data be required for an HPC compute task, data transfer nodes are often provided with faster connectivity and the data is transferred in advance. Such a prefetching, however, is not always feasible for HEP workflows, particularly with the CMS computing model that heavily relies on streaming of (partial) datasets for more efficient data transfers. As a consequence, data intensive HEP jobs can, e.g., suffer from I/O limitations and bottlenecks on such centers.

For the opportunistic integration of HoreKa, for example, a comparably lower CPU efficiency was observed during the initial phase of the integration. Fig. 2 shows the comparison to GridKa (the German Tier-1), RWTH (Tier-2 in Aachen), TOPAS (Tier-3 at KIT [6]), and RWTH-HPC (CLAIX, the NHR center in Aachen). While CLAIX/RWTH-HPC shows that a comparably efficient utilization of NHR centers is in principle possible, HoreKa is clearly lacking behind.

But when it comes to actually identify the reasons causing the degrading in CPU efficiency, two other challenges arise: First, the security policies on HPC centers often follow a zero trust policy, from which also the HEP community is not excluded. As a result, the monitoring capabilities are often limited, making the bottleneck identification more complicated. Nevertheless, clear indications for a data access bottleneck via the external 1 Gb/s link of the worker nodes (indicated in brown in Fig. 3) were found<sup>1</sup>. And second, in contrary to traditional Grid resources which can be upgraded when such a data access bottleneck is observed, the hardware setup of the NHR centers cannot be modified – it can only be influenced in the future at best. As a consequence, the mitigation of any kind of limitations has to be solely software-based and must not require root access, making the optimization of the HPC center utilization significantly more complex.

To achieve an adequate and comparable utilization of HoreKa, various improvements were tested during its opportunistic integration and an XRootD-based concept for an improved usage of HPC centers as Grid resources was developed, which is described below.

---

<sup>1</sup>One strong indication is the very comparable performance of CLAIX, which provides 10 Gb/s per worker node.

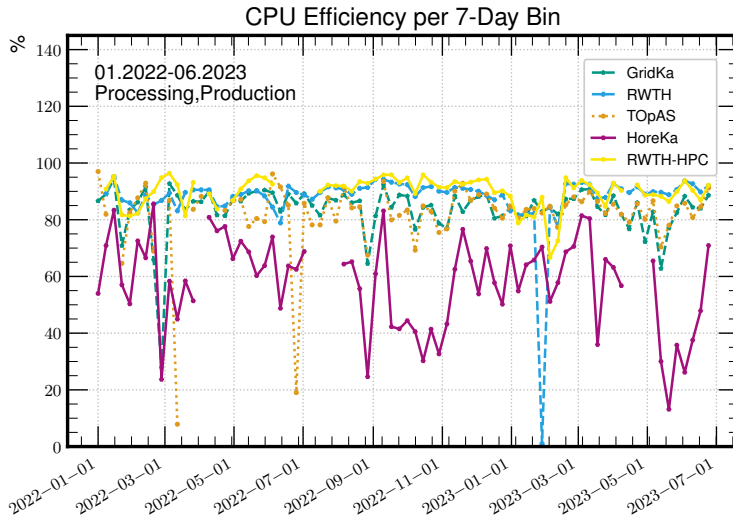


Figure 2: CPU efficiency comparison for official WLCG sites (GridKa, RWTH) and opportunistic sites in Germany. While RWTH-HPC (CLAIX, the NHR center in Aachen) shows a comparable efficiency during the pilot phase, HoreKa is clearly lacking behind the WLCG sites. [7]

### 3 Workflow and Efficiency Optimizations with XRootD

For the identification of bottlenecks and limitations, the period from January 2022 until June 2023 was investigated, where HoreKa was integrated as is – without any optimizations. After an in-depth investigation of the reasons for the less efficient utilization of the NHR center<sup>2</sup>, different optimization steps were developed based on the observations. Alongside the optimization of the site configuration<sup>3</sup>, an XRootD-based concept for data access bottleneck mitigation of such centers with limited I/O capabilities per worker node was developed.

#### 3.1 Concept

The operational mode of the HoreKa HPC cluster does not foresee large-scale data transfers directly over the worker nodes. Instead, login nodes with a faster external connectivity (indicated in green in Fig. 3) are provided. To better align the utilization for HEP workflows to the operational concept of HPC workflows, a fully containerized, XRootD-based concept with minimal requirements was developed to overcome the observed limitations. For this, an XRootD Caching Proxy (XCache) is deployed on a login node of the cluster with a 50 Gb/s external link and access to the parallel filesystem, which is accessible by every node and can be used as cache space. HEP jobs running on the worker nodes are then configured to direct their request for data to the Proxy (dashed, blue arrow) by simply prefixing the local proxy

<sup>2</sup>More details on the process and results for the investigation can be found in Ref. [7].

<sup>3</sup>An improvement was achieved by disabling the legacy *lazy download* feature and an intelligent adaption of the workflow mix. More details are provided in [7].

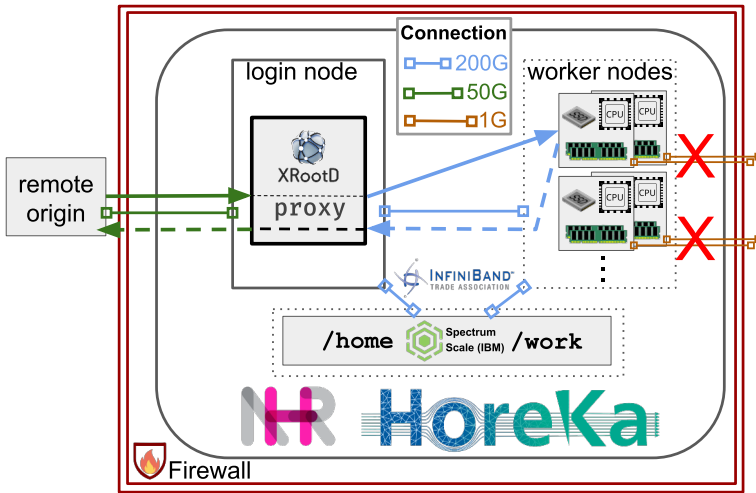


Figure 3: Overview of the cluster setup including the connectivity and the XRootD (Caching) Proxy deployed on the login node. Instead of using the slow, external links of the worker nodes directly (brown, 'X'), HEP jobs are configured to use the XCache deployed on the login node, which transfers the data with its faster external connection to the Internet (green) and provides the data internally over the fast Infiniband (blue). [7]

address in the site configuration<sup>4</sup>. The data request is then sent via the XRootD Proxy to the remote origin (dashed, green arrow) and the transfer is executed from the login node with its 50 Gb/s WAN connection (green arrow) instead of the worker's 1 Gb/s link (brown, red 'X'). Fetched data is finally provided by the XCache server over the fast internal Infiniband to the initial client (blue arrow), resulting in an accelerated transfer and a significant increase in available bandwidth – efficiently mitigating the data access bottleneck.

For a further improvement of the data rates, the caching feature of XRootD can be employed. In this case, the XCache server caches the transferred data after fetching it from remote on-the-fly on the parallel filesystem. If a requested data block is locally available, it can be provided directly via the internal Infiniband, reaching several GB/s and therefore a significant acceleration. Whether caching is beneficial, however, strongly depends on the circumstances and different aspects need to be considered. It is important that a sufficiently large cache (or share on the cluster's filesystem) is available, matching the expected number of parallel running jobs and the processed remote data. Only if these aspects are in an adequate relation, a decent cache hit rate can be expected, resulting in an actual improvement of the utilization. As an example, in case of HoreKa, 250 TB are usable for the prototype setup. Since the NHR center was mainly used for Monte Carlo simulation tasks for CMS, caching was observed to be rather inefficient due to the sheer size of the CMS pileup-mixing datasets that are used. However, under different circumstances, caching can definitely be beneficial, as for example shown in Ref.[8].

Even if the caching is not leading to significant improvements, another feature of XCache can be utilized to further optimize the data access. If enough bandwidth is available, prefetching can be enabled with the presented setup. XCache then fetches additional data blocks in advance, which are consequently provided via the fast Infiniband when actually requested.

<sup>4</sup>Implementation details are available in Ref. [7].

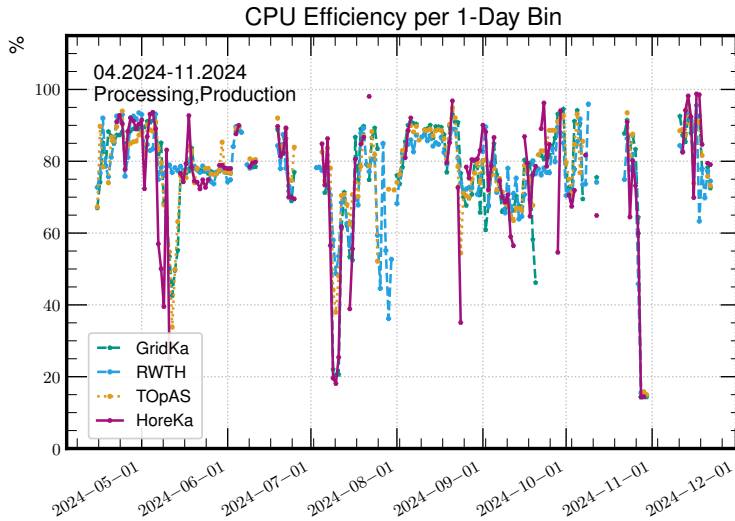


Figure 4: CPU efficiency comparison between HoreKa and the other investigated sites. With the optimizations in place, HoreKa is now clearly in line with the traditional WLCG sites most of the time, which is a great improvement in comparison to Fig. 2.

The setup can therefore better be described as 'XBuffer' instead of XCache, essentially serving as a buffer for the cluster. Here, it is important to mention that the configuration should be carefully adapted to avoid introducing bandwidth limitations again, when too much data is transferred in advance. The results for the prototype deployment of such an XBuffer at HoreKa are presented in the next section.

### 3.2 Results for the First Prototype Deployment of XBuffer at HoreKa

The ultimate goal for the utilization of HPC centers within the German HEP computing strategy is to being able to process HEP workflows as reliable and efficient as on classic WLCG resources. While CLAIX was already able to compete in terms of efficiency without any improvements, HoreKa was clearly lacking behind in the initial phase of the integration, as shown in Fig. 2. With the optimized configuration and the deployment of an XBuffer prototype, however, it was possible to improve the utilization of HoreKa significantly. This is depicted in Fig. 4. With the optimizations in place, the cluster is now competing well with the other sites, indicating that a comparable utilization of the NHR centers for HEP workflows is definitely possible.

When comparing the utilization of HoreKa directly with the university Tier-2 in Aachen (RWTH, Fig. 5) and the KIT Tier-3 center (TOpAS, Fig. 6), a great improvement can be observed in terms of efficiency. The NHR center performs very comparable to the other sites and reaches the target region (90 % to 100 % of the WLCG site in comparison), as indicated by the ratios in the bottom part of both figures. In terms of reliability, HoreKa is comparable to the Tier-2 center, but with an overall higher failure rate slightly lacking behind the Tier-3 at KIT.

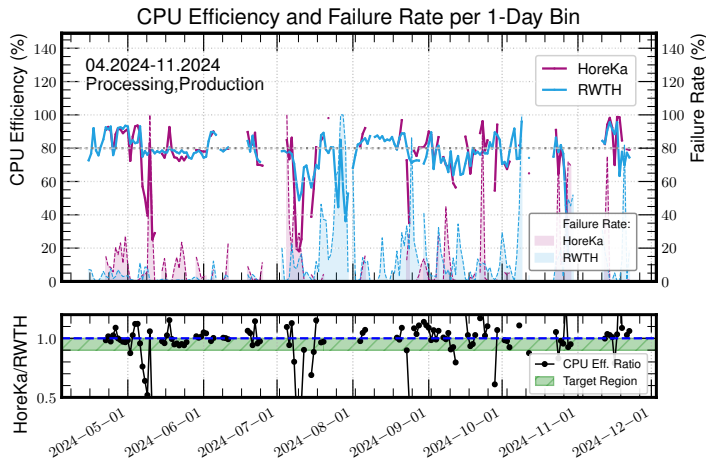


Figure 5: With optimizations, HoreKa is most of the time performing very similar to the Tier-2 center in Aachen. CPU efficiency and failure rate (top) are overall very comparable. The efficiency ratio per bin even shows that the NHR center is not only inside the target region of a comparable efficiency, but even above in several bins, reflecting a considerable improvement compared to the initial phase (Fig. 2).

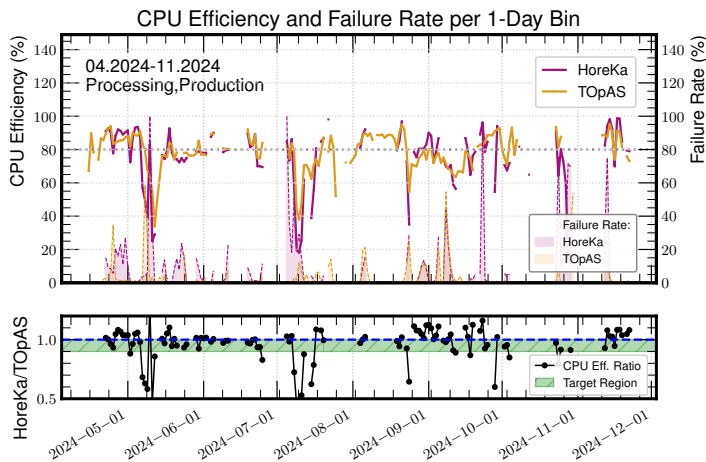


Figure 6: In this figure, the utilization of HoreKa is compared to TOpAS, the KIT Tier-3 center. While the Tier-3 is with an overall lower failure rate slightly better in terms of reliability, the HPC cluster now can compete well in terms of efficiency. Just as for the Tier-2, the target region is reached in most bins.

## 4 Conclusion and Outlook

With this contribution, the first deployment of an XRootD-based prototype for the optimized utilization of German scientific HPC clusters for HEP workflows was presented. The deployed XBuffer concept – together with site specific configuration optimizations – lead to a

significant improvement of HoreKa's integration. With all optimizations in place, it was possible to achieve a comparable – or even better – performance in most of the investigated time periods, clearly showing that a comparable utilization of HPC centers as part of the NHR strategy will be possible in the future. Furthermore, the presented concept can be flexibly extended in the future. With the current HEP job allocation at the cluster, the described setup is sufficient to achieve the desired goals. When scaling up the shares on the HPC site over the next years, however, the prototype setup can also run into limitations with only a single node utilized for the proxying. These can be hardware limitations when too many jobs are using the XCache/XBuffer, or again I/O limitations, when too much data is transferred. For this case, two possible solutions were developed. The first one is a simple scaling of the presented prototype setup. Instead of only using one proxy node, the setup can easily be scaled horizontally by adding further proxies on multiple login nodes managed by a local XRootD Redirector. With this, the available bandwidth can be upgraded, and as an additional benefit, an implicit load balancing mechanism is introduced when clients contact the Redirector, which selects the best matching Proxy for the transfer. As an alternative, one or more dedicated transfer nodes could be deployed at the collocated WLCG site (GridKa, in case of HoreKa), if the NHR center allows a direct connection to the internal Infiniband network. This approach even has additional benefits, as on the one hand, the hardware is again under WLCG control, simplifying the deployment and maintenance of the setup, and on the other, this setup would allow a transfer over the dedicated Grid network (LHCONE) instead of the public Internet, reducing the efficiency loss through several firewalls.

In summary, this work not only demonstrates the feasibility of the HPC integration for the future German HEP computing strategy but also provides a robust foundation for further optimizations in the future.

## References

- [1] Komitee für Elementarteilchenphysik, Perspektivpapier der Teilchenphysiker:innen in Deutschland. (2022). [https://www.ketweb.de/sites/site\\_ketweb/content/e199639/e312771/KET-Computing-Strategie-HL-LHC-final.pdf?preview=preview](https://www.ketweb.de/sites/site_ketweb/content/e199639/e312771/KET-Computing-Strategie-HL-LHC-final.pdf?preview=preview)
- [2] NHR-Verein e.V., National High Performance Computing (accessed: 17.02.2025). <https://www.nhr-verein.de/en>
- [3] Könemann, C. and Wiebe, S., KIT Supercomputer One of the World's Most Energy-Efficient. KIT Press Release 038/2024 (2024). [https://www.kit.edu/kit/english/pi\\_2024\\_038\\_kit-supercomputer-one-of-the-worlds-most-energy-efficient.php](https://www.kit.edu/kit/english/pi_2024_038_kit-supercomputer-one-of-the-worlds-most-energy-efficient.php)
- [4] Giffels, M. and others, MatterMiners/tardis, Zenodo, (2024). <https://doi.org/10.5281/zenodo.2240605>
- [5] von Cube, R. F., *Dynamic Integration of Heterogeneous Computing Resources and Jet Energy Calibration for the CMS Experiment* (PhD Thesis, Karlsruhe Institute of Technology, 2022). <https://doi.org/10.5445/IR/1000151252>
- [6] Caspart, R. and others, Setup and commissioning of a high-throughput analysis cluster. EPJ Web of Conferences **245**, 07007 (2020). <https://doi.org/10.1051/epjconf/202024507007>
- [7] Hofsaess, R., *Optimized Utilization of HPC Centers for High-Energy Physics Workflows and Jet Energy Corrections for the CMS Experiment* (PhD Thesis, Karlsruhe Institute of Technology, 2025), to be submitted
- [8] Flix Molina, J. and others, Optimal XCache service for the CMS experiment in Spain. **22nd** conference in the ACAT series, submitted (2024). <https://indico.cern.ch/event/1330797/contributions/5796642/>