

# Bioinformatics and Astrophysics Cluster (BinAC)

Jens Krüger<sup>1</sup>, Volker Lutz<sup>1</sup>, Felix Bartusch<sup>1</sup>, Werner Dilling<sup>1</sup>, Anna Gorska<sup>2</sup>, Christoph Schäfer<sup>3</sup>, and Thomas Walter<sup>1</sup>

<sup>1</sup>Zentrum für Datenverarbeitung, Eberhard Karls Universität Tübingen

<sup>2</sup>Algorithms in Bioinformatics, Eberhard Karls Universität Tübingen

<sup>3</sup>Computational Physics, Institut für Astronomie und Astrophysik,  
Eberhard Karls Universität Tübingen

BinAC provides central high performance computing capacities for bioinformaticians and astrophysicists from the state of Baden-Württemberg. The bwForCluster BinAC is part of the implementation concept for scientific computing for the universities in Baden-Württemberg. Community specific support is offered through the bwHPC-C5 project.

## 1 Introduction

The bwForCluster BinAC offers cutting edge computing capabilities for users from the scientific domains of bioinformatics, astrophysics and related fields [1]. In both domains, scientific challenges are addressed through highly demanding computational approaches. As thus, the need for efficient and highly performant compute resources was immanent during the procurement phase for BinAC. Anyhow, the precise requirements from individual user groups dealing with particular use cases led to a broad spectrum of desired characteristics.

## 2 Resources and Architecture

Maximum performance on single cores, low latency network connections for improved scaling and the most recent GPU technology were the most prominent requests. All of BinAC's 300 individual nodes have a common base configuration. Each node holds two Intel Xeon E5-2680v4 (Broadwell) making a total of 28 cores available on each node. Aggregated over all 11,184 cores this leads to a nominal peak performance of more than 534 Tera-Flop/s, which put BinAC briefly on the TOP500 list of the world fastest supercomputers [2]. Additionally, each node provides 128 GB



DDR4-RAM of memory and a 256 GB local SSD harddisk. Also, 60 nodes of BinAC are equipped with two Nvidia Tesla K80 cards each. As each of these GPU accelerators has two Kepler GK210 Chips a total of 4 GPUs are available on each of the GPU nodes. Furthermore, to cope with high memory jobs, four SMP nodes with 40 cores each were equipped with 1 TB DDR4-RAM of memory. In order to provide interactive visualization capabilities of simulation data four dedicated nodes are available. They are equipped with NVIDIA Quadro M4000 graphic cards.

As depicted in Figure 1, all compute nodes are part of a low-latency high-bandwidth FDR Infiniband fabric and connected with the outside world through a 10 GBit/s uplink. The global workspace file system is based on the scalable parallel file system BeeGFS providing up to 720 TB raw storage space.

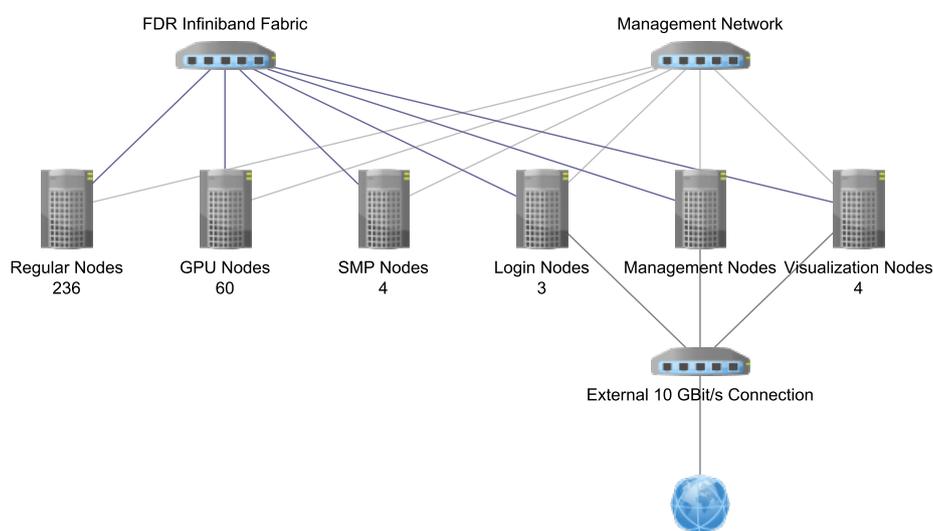


Figure 1: Architecture overview of the BinAC Cluster environment. All nodes are part of a low-latency high-bandwidth FDR Infiniband fabric and connected to the outside world through a 10 GBit/s connection.

### 3 Application Show Cases

The following examples illustrate the effectiveness of BinAC to address current scientific questions from different fields.

A classical approach for the metagenomic analysis of the human gut microbiome is to compare reads against the chosen database and then analyze the counts of reads assigned to the certain taxa [3]. In the case of phages this approach is not feasible since there are not enough known phage genomes to compare against. Therefore, one first needs to assemble the reads to be able to analyze long stretches of the phage genomes instead of the short reads. Assembly, especially of big datasets, is computationally expensive. For the case at hand, the sequencing reads were pulled together, resulting in two datasets, one per participant, comprising 134 million and 110 million reads. The assembly program Ray was run on 112 CPUs on BinAC cluster, 13 times (different k-mer size) for each participant (see Table 1). The authors of Ray have tested it on a

k-mer	A	B	k-mer	A	B
19	2h 24m 3s	3h 34m 59s	33	2h 19m 34s	4h 9m 3s
21	2h 7m 27s	3h 44m 35s	35	2h 21m 6s	4h 8m 24s
23	2h 8m 31s	4h 16m 23s	37	2h 19m 15s	3h 22m 32s
25	3h 40m 12s	3h 23m 53s	39	2h 18m 35s	4h 9m 39s
27	3h 40m 31s	4h 9m 8s	55	2h 20m 4s	4h 7m 54s
29	2h 17m 17s	4h 7m 52s	77	4h 3m 15s	3h 21m 34s
31	2h 17m 41s	3h 22m 16s			

Table 1: Runtimes for the assembly of reads using Ray.

roughly four times bigger dataset. Using 128 cores it took them 13 hours to complete [4]. But the datasets are not directly comparable, since they used a simulated metagenome dataset with entire bacterial genomes, therefore the resulting contigs were much larger.

Ray is not the only program applicable for this use case. Thanks to BinAC it is possible to construct a pipeline comprising several steps, (1) Ray (assembly), (2) Bowtie2 (read mapping back to the assembled scaffolds), (3) Blast (comparison to the Viral and CARD databases), (4) Prodigal (gene prediction) and (5) Aragorn (tRNA gene prediction), and run it in one go, with various parameters. Such an approach would be impossible using other resources, not only because of the CPU and memory usage, but also because it would not be possible to occupy a lab server completely for longer periods of time.

Gromacs is a molecular dynamics package intended for the simulation of biochemical molecules like proteins, lipids and nucleic acids [5]. It supports CUDA-based GPU acceleration. Benchmarking the alcohol dehydrogenase protein (ADH) system with 130,000 atoms using a rectangular box and PME ([http://www.gromacs.org/GPU\\_acceleration](http://www.gromacs.org/GPU_acceleration)) already results in a simulation speed of 23 ns/day just using the 28 CPU cores available on a single compute node. When adding one GPU the speed is amplified by factor 2 (47 ns/day) without any further optimization. When simulating larger systems the effect is even more pronounced. A membrane system with two ion channels consisting of 430,000 atoms (rectangular box, PME) can be simulated with 6.0 ns/day just using CPU cores. Using all four GPUs available a speed of 24.3 ns/day can be reached which corresponds to a speedup of 4. Hence, the GPU nodes of BinAC enable researchers to simulate their large biomolecular systems on biological relevant time scales which had been much more challenging before.

The smooth particle hydrodynamics (SPH) equations allow to model gas, liquids and elastic, and plastic solid bodies. The approach is used for the simulation of the collision between Ceres-sized objects. Schäfer et. al. implemented a GPU version of this method also considering self-gravity of the simulated objects [6]. Comparing the runtimes of the new implementation to an existing OpenMP version they can report a speedup of at least two orders of magnitude for all relevant substeps of the SPH method. Given that the runtime evaluation was 'only' carried out on a single Nvidia GTX Titan even larger speedup and further parallelization options arise from the four K80 GPU cores available on each of BinAC's GPU nodes.

## **4 Access and Support**

Generally all academic users who have some kind of affiliation with a research institution within the state of Baden-Württemberg are eligible to apply for access to BinAC. In order to provide optimal assistance to researchers looking for high performance compute resources, a lightweight application system was introduced by the bwHPC initiative ([https://www.bwhpc-c5.de/zas\\_info\\_bwforcluster.php](https://www.bwhpc-c5.de/zas_info_bwforcluster.php)). The researchers are asked to provide a short description of the planned research, desired compute resources and software environments. This information serves as basis for an appropriate resource allocation so the task can be distributed to the compute center which is best suited to handle it. The bwHPC-C5 partners, including the team at the compute center of the University of Tübingen, provide assistance and consulting for all questions related to the usage of BinAC ([hpcmaster@uni-tuebingen.de](mailto:hpcmaster@uni-tuebingen.de)). Some scientific communities might require more effort to enable an appropriate research environment, like installing special software packages or enabling access to external data sources. For this purpose the bwHPC-C5 partners offer the opportunity to form Tigerteams offering extended support through experts from the associated compute centers and experts from the affected scientific fields.

## **5 User Perspective**

The queues on BinAC are filling up fast. Users are encouraged to apply quickly for an appropriate entitlement. The highly motivated HPC team at the ZDV is supporting users from the related communities, trying to provide the necessary software environment, stable operation of the underlying hardware and consequently a positive user experience. Recent additions to the module environment used for handling the simulation software covers machine learning applications such as Tensorflow and Caffe.

## **6 Conclusions**

BinAC offers state of the art compute capabilities accompanied by result oriented support through the bwHPC-C5 project. A powerful research tool is offered to the users from the research areas of bioinformatics and astrophysics, enabling them to carry out their cutting edge research.

## **Acknowledgments**

The bwHPC-C5 project has been funded by the Ministry of Science, Research and the Arts of the state of Baden-Württemberg, Germany. The Universities of Freiburg, Heidelberg, Hohenheim, Konstanz, Mannheim, Stuttgart, Tübingen and Ulm, the Karlsruhe Institute of Technology, as well as the Universities of Applied Sciences in Stuttgart and Esslingen are the partners carrying out the project.

## **References**

- [1] H. Hartenstein, T. Walter, and P. Castellaz. “Aktuelle Umsetzungskonzepte der Universitäten des Landes Baden-Württemberg für Hochleistungsrechnen und datenintensive Dien-

- ste.” *Praxis der Informationsverarbeitung und Kommunikation*, Band 36, Heft 2 (2013): 99-108.
- [2] <https://www.top500.org/system/178843> (accessed 2017-02-02)
- [3] M. Willmann, M. El-Hadidi, D. H. Huson, M., Schütz, C. Weidenmaier, I. B. Autenrieth, and S. Peter. “Antibiotic selection pressure determination through sequence-based metagenomics.” *Antimicrobial Agents and Chemotherapy*, 59(12) (2015): 7335-7345. <http://doi.org/10.1128/AAC.01504-15>
- [4] S. Boisvert, F. Raymond, E. Godzaridis, F. Laviolette, and J. Corbeil. “Ray Meta: scalable de novo metagenome assembly and profiling.” *Genome Biology*, 13(12), (2012) R122. <http://doi.org/10.1186/gb-2012-13-12-r122>
- [5] M. J. Abraham, T. Murtolad, R. Schulz, S. Páll, J. C. Smith, Berk Hess, and E. Lindahl. “GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers” *SoftwareX*, Volumes 1–2 (2015): 19–25.
- [6] C. Schäfer, S. Riecker, T. I. Maindl, R. Speith, S. Scherrer, and W. Kley. “A smooth particle hydrodynamics code to model collisions between solid, self-gravitating objects” *Astronomy and Astrophysics* 590, A19 (2016).