



# ExCALIBUR: Hardware & Enabling Software Testbeds

February 2021

<https://excalibur.ac.uk>

[researchinfrastructure@epsrc.ukri.org](mailto:researchinfrastructure@epsrc.ukri.org)

## Contents

Novel hardware/software architecture testbed - University of Birmingham .....	3
Graphcore testbed – University of Bristol.....	4
Exascale Data Testbed for Simulation, Data Analysis & Visualization, University of Cambridge .....	5
AMD GPU testbed – Durham University .....	6
Storage and RAM as a service – Durham University .....	7
Wafer scale testbed – University of Edinburgh.....	8
FPGA testbed – University of Edinburgh, UCL, University of Warwick.....	9
ARM+GPU Demonstrator, University of Leicester .....	10
The UCL Adaptable Cluster Project .....	11

## Novel hardware/software architecture testbed - University of Birmingham

This project will create a testbed featuring novel accelerator technology from [NextSilicon](#) in collaboration with University of Birmingham HPC systems partner Lenovo, as part of a co-design partnership. The project will evaluate the performance of the main codes used by UKRI researchers, with a particular emphasis on evaluating some of the major algorithm classes used in supercomputing and data science, as well as other HPC apps.

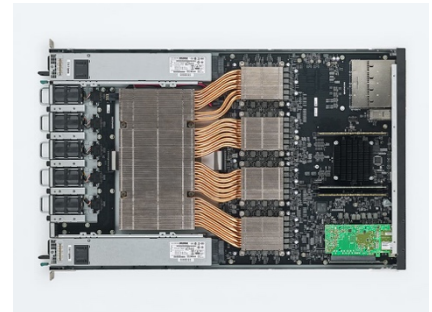


Should the technology fulfil its promised potential in the evaluation, it will be made available to the UK HPC community as part of the recently awarded [Baskerville EPSRC Tier 2 service](#) which has recently been awarded to the University of Birmingham. This facility, which will be installed in Q1 2021, will include 184 Nvidia A100 GPU accelerator cards. Other novel accelerator technologies are also planned for introduction over the lifetime of the facility, thus enabling us to assess the new technology in a heterogeneous accelerated live HPC environment.

## Graphcore testbed – University of Bristol

This testbed will evaluate the [Graphcore IPU-M2000](#) system for high performance and scientific computing applications and provide a novel architecture for the community to test and develop AI compatible codes on. The IPU (Intelligent Processing Unit) is a completely new kind of massively parallel processor, co-designed from the ground up to accelerate machine intelligence.

Each MK2 GC200 IPU in the IPU-M2000 unit has 1472 processor cores, running nearly 9,000 independent parallel program threads with 900MB in processor memory and 250 TeraFlops of AI compute at FP16.16 and FP16.SR (stochastic rounding). The IPU-M2000 system has four IPUs, delivering approximately 1 PetaFlop of AI compute, and supporting ultra-low latency IPU-Fabric interconnect.

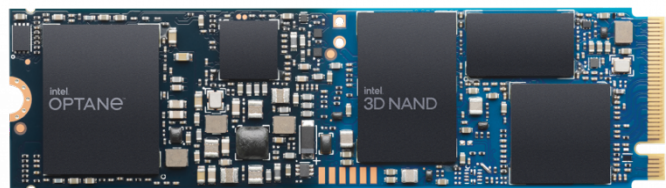


The testbed includes four IPU-M2000 systems, which will enable the interconnect to be tested and characterised. The project will evaluate the Graphcore system's intended use cases around AI training and inference, and also look at a subset of HPC codes that may be suitable for this platform. The Graphcore system will also be made available to the ExCALIBUR and wider UK research community with support and a training programme from the Bristol team. It should be noted that codes will need to fit into small memories and must be single or half precision due to IPU requirements.

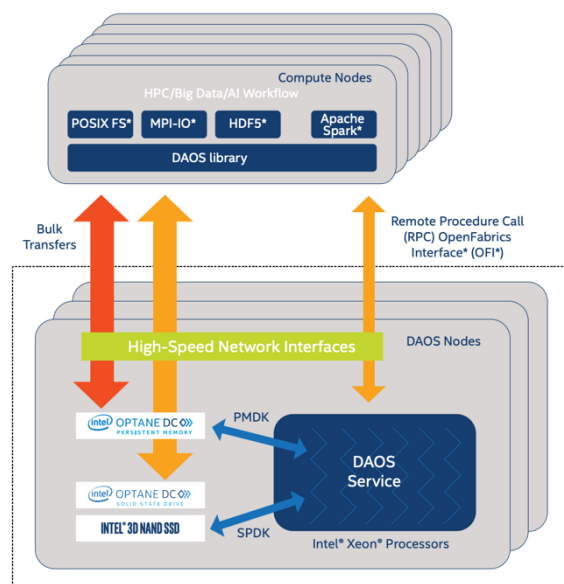
## Exascale Data Testbed for Simulation, Data Analysis & Visualization, University of Cambridge

This testbed utilises world leading HPC systems development, deployment and operational skills housed within the Cambridge Research Computing Service to build a next generation high performance PCI-Gen-4 solid state I/O testbed utilising a range of file systems including Lustre, Intel DAOS, BeeGFS and HDF5 on state-of-the-art solid state storage hardware.

The system utilises the latest Intel PCI Gen-4 NVMe drives and the new Intel PCI-Gen-4 Optane Data Centre Persistent Memory. The project will see the deployment in April 2021 of the UK's fastest HPC storage testbed delivering over 500GB/s bandwidth and over 20 million IOPS of raw I/O performance which can be deployed across applications via a range of leading HPC file systems such as Lustre, Intel DAOS or BeeGFS as well as other more low level direct I/O protocols.



It is expected that the solution will also be one of the fastest HPC storage solutions in the world, being ranked high in the worldwide I/O 500 listing. Intel [DAOS](#) is of particular interest since it has been developed from the ground up to provide a persistent file system utilising both NVMe drives and Optane DCP memory. DAOS is still at proof of concept stage but is shown to deliver far higher performance than traditional parallel file systems. This project is supported directly by Intel in terms of hardware, staff effort and strong co-design work in collaboration with Intel engineers developing the DAOS file system. The system will represent Intel's largest DAOS testbed.



In addition to the I/O hardware and various file system technologies the testbed is configured with comprehensive system level telemetry monitoring capability provided by the UKRI funded Scientific OpenStack middleware layer combined with a range of other more specialised application I/O profiling tools. The UK [Scientific OpenStack](#) is a world leading HPC middleware layer developed at Cambridge and funded by over 4 years investment from STFC, EPSRC and MRC. System I/O telemetry combined with application level I/O profiling is vital if we are to fully exploit emerging I/O and file system technologies by helping application developers understand how to implement the most efficient I/O mechanisms within the application code. Without such tools developers will be blind in terms of how to best utilise the new I/O platforms.

## AMD GPU testbed – Durham University

The Durham AMD GPU testbed provides researchers with the opportunity to test their code on the [AMD MI50 GPU](#). The testbed consists of a Gigabyte server with:

- 2 x AMD EPYC 7282 16 core 2.8GHz CPUs
- 1TB RAM
- 6 x AMD MI50 GPUs
- AMD software, ROCM, AOCC, AOMP, GCC with offload support installed



The GPU testbed extends the existing AMD cluster at Durham and is already in use by researchers at various sites including the University of Bristol and the Hartree Centre, with the first journal submission for research using this testbed already submitted.

## Storage and RAM as a service – Durham University

The Durham Adaptable Memory System has distinct components that are required to investigate adaptable memory technologies that will function as a testbed and demonstrator for the ExCALIBUR HES programme. These components, namely a [BlueField-2](#) cluster and a [Gen-Z](#) equipped cluster, their functions and resource requests are given in turn below.

This will be the first UK install of both BlueField-2 and Gen-Z for HPC within the UK. Expertise in BlueField-1 exists at Durham, with a 16 node system already in operation. Both of these systems will be integrated with [COSMA](#), allowing the existing login nodes, LDAP servers and administration consoles to be used. Users will be able to request an account through the [SAFE](#) system managed by EPCC.



BlueField-2 technology is not yet available on the market, and so we will be getting pre-release access for this proof of concept (PoC) cluster. BlueField-2 has significant advantages over the original BlueField-1 cards, namely increased processing power and clock rate of the embedded Arm cores. This is therefore an ideal time to realise this test cluster environment, placing the UK on the leading edge of this novel technology. The expertise exists within Durham based on experience with BlueField-1 systems.

Gen-Z technology is also not yet available on the market. However, Durham University have joined the Gen-Z consortium, and will have access to pre-market proof of concept equipment which will be obtained as part of this proposal. This will place the UK in an advantageous position to test and evaluate this new technology.

## Wafer scale testbed – University of Edinburgh

This project brings a [Cerebras CS-1 Wafer Scale Engine](#) system to the UK –the first such system in Europe. This enables performance and usability exploration for UK academic and industrial users. The majority of the system has been funded by the University of Edinburgh, however the support from ExCALIBUR HE&S has allowed for a more general access service to be provided to researchers from across ExCALIBUR and the wider computational science and AI community in the UK.

Cerebras Systems have developed the world’s largest processor, the Wafer Scale Engine (WSE), at over 46,000 square millimetres, with 1.2 trillion transistors, 400,000 processor cores, 18 gigabytes of SRAM, and an interconnect between processors capable of moving 100 million billion bits per second. With the WSE at its core, the Cerebras CS-1 system is firmly focussed on neural network training and according to Cerebras the CS-1 provides 3000x more capacity and 10,000x greater bandwidth than the leading competitor.

From a software perspective, Cerebras Systems have integrated their hardware into common machine learning frameworks such as TensorFlow and PyTorch2, opening up the potential for easy porting of existing application to the system. They also provide a graph compiler (CGC) and optimised library kernels, to efficiently map applications to the many processors on the WSE and ensure optimal use of the resource.



With potential for extreme performance for a wide range of machine learning training tasks, the Cerebras CS-1 is a very exciting new technology. However, there is currently a lack of user experience and application performance data to assess the suitability of the hardware for actual applications, and the requirements/costs for porting codes to the system. With a software environment that partially resembles standard CPU- and GPU-based systems, and partially resembles FPGA-based systems, with associated placement and routing requirements, it is important to be able evaluate both performance and usability of the CS-1 for end user applications. Such end user applications may also include more traditional numerical applications and this will be an area of exploration on the system.



## FPGA testbed – University of Edinburgh, UCL, University of Warwick

This testbed system, and associated effort for enabling software, is aimed at allowing researchers to port their scientific and data-science applications to Field Programmable Gate Arrays (FPGAs) and explore performance and power advantages such technology provides. Composed of next-generation hardware and software, this will form an important UK resource for exploring the future role of FPGA technology in science, engineering, and the broader computational science communities.

In addition to the testbed hardware itself the project is supported by Research Software Engineer (RSE) effort to develop the software stack to enable easier usage of FPGAs, which will be driven by specific use cases from the Excalibur Design and Development Working Groups and other interested application communities.

Project partners EPCC, UCL, and Warwick will work in collaboration with FPGA vendor Xilinx, Inc. the leader in adaptive and intelligent computing, to deliver and operate the testbed.

It is their intention that this will be a first step towards building a future community and ecosystem around the role of FPGAs in HPC, data science, AI, and machine learning workloads in the UK. The project will also be running a series of training events and workshops, and developing training material to ensure the system is accessible and usable.



The testbed will be physically based in EPCC's Advanced Compute Facility, and will be made publicly available. It will form a unique resource within UK academic computing, as a single system that provides access to next-generation Versal Adaptive Compute Acceleration Platform (ACAP) technology from Xilinx, which includes their revolutionary AI engines; hierarchical memory hardware provision, with high bandwidth (HBM2) and Non-Volatile (NVRAM) memory on some of the hosted hardware, providing a unique resource for software developers and algorithm designers to investigate this emerging field in computing hardware; multiple networking options including a high performance node-level network and direct FPGA to FPGA networking to enable system designers and applications developers to assess the relative merits of both approaches; and multiple families of FPGA, allowing evaluation of a range of technologies by users.

The system will be hosted within an existing, established and modern HPC system which provides sufficient resources to enable developers to quickly and efficiently develop application kernels, synthesise their FPGA bitstreams and test their codes in emulation. Finally, the RSE effort will provide an enabling software stack that should significantly reduce the barrier to entry in utilising FPGAs for scientific and data-science applications.

## ARM+GPU Demonstrator, University of Leicester

The aim of this testbed is to ensure that:

- ARM servers work harmoniously with accelerators (such as GPUs)
- Any shortcomings are understood, documented and reported to vendors
- ExCALIBUR and the wider UK research community has access to an ARM-GPU testbed

The ARM+GPU testbed system is based on HPE's Apollo 70 platform and Marvell's ThunderX2 processors, with 2x NVidia V100 GPU per server, incorporated into the University's existing [ARM Catalyst](#) system.



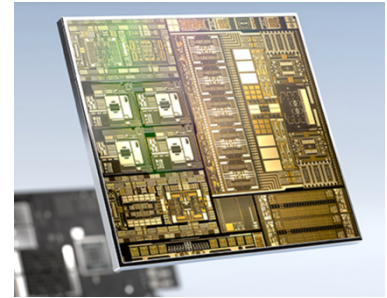
The project includes Research Software Engineering (RSE) effort for the porting, benchmarking and development of existing codes to support the programme of work described above. The RSE will also contribute to the creation of digital assets including progress reports, whitepapers and how-to documents, as well as software enhancements and modification.

The ARM+GPU testbed builds on the University of Leicester's existing experience managing ARM based HPC systems, our Software engineering team's CUDA expertise, and our existing relationship with the technical teams at ARM.

## The UCL Adaptable Cluster Project

The ExCALIBUR Interconnect Demonstrator consists of two non-blocking interconnect fabrics supporting up to 60 attached nodes in a dual fabric configuration.

One fabric is 200 Gbps HDR Mellanox Infiniband configured so that it is possible to construct multi-hop routes between nodes. The second fabric is 100Gbps Mellanox Ethernet, with [BlueField](#) adaptors on each node. This allows us to measure the impacts of a variety of in-network technologies – doing computation at the switch level (requiring multiple hops) and looking at the possibility of using acceleration on the adaptor to off-load some of the work of the host machine (the BlueField cards). We also aim to compare “state of the art” in using Ethernet as an Interconnect with Infiniband to measure whether on RDMA on Converged Ethernet has reached the point where it is a performant, cost effective interconnect.



In order to understand system and application performance the Adaptable Cluster collects metrics from several sources in the system and dashboards to visualise them, which then allow focus on how to improve system design and resource usage. Alerts can be set up to draw attention to performance issues as well. The testbed uses components such as Elasticsearch, Kibana, Logstash and Prometheus to provide insights into both breadth and depth of system and application performance.



UCL is the location of the ExCALIBUR instance of the [ARM FORGE](#) Application. This is an application that supports the debugging, profiling and optimisation of codes that use distributed resources, such as a cluster. It is both CPU and GPU enabled. UCL will support ARM FORGE for key centres in the ExCALIBUR project. It will also be available to UCL projects that are not associated with ExCALIBUR. This package enables jobs that use up to 2048 cores to be analysed in terms of code efficiency. One outcome of this project will be methodologies that enable results from Prometheus and ARM Forge to be used to improve system design, architecture performance and application performance.

