# Maximising Sustainability of Isambard AI Exascale Supercomputing Platform, from Data Centre to Compute Nodes

Sadaf R Alam
University of Bristol, UK
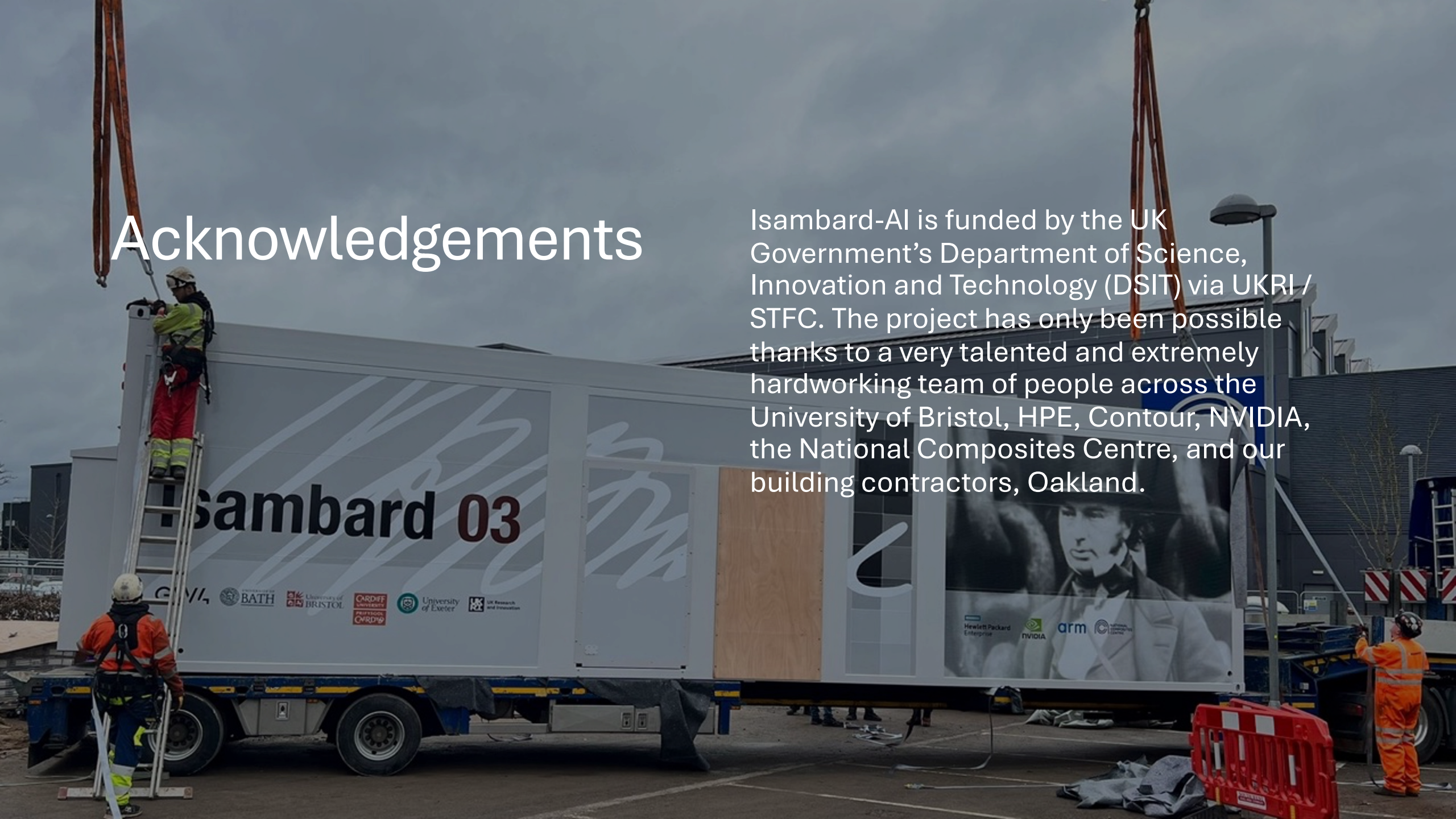ExaComm 2024
May 16, 2024
Hamburg, Germany

# Outline

- Background and timeline

- Design specifications for an AI Research Resource (AI RR)
  - Sustainability
  - Performance
  - Accessibility

- Next steps

University of BRISTOL

# Acknowledgements

# Bristol set to host UK's most powerful supercomputer to turbocharge AI innovation

A new supercomputer set is to be built in Bristol, in a move to drive pioneering AI research and innovation in the UK.

From: **Department for Science, Innovation and Technology** and **The Rt Hon Michelle Donelan MP**
Published 13 September 2023



- UK AI Research Resource dubbed Isambard-AI will be one of Europe's most powerful supercomputers
- new facility will serve as national resource for researchers and industry experts spearheading AI innovation and scientific discovery
- plans for the supercomputer backed by £900 million investment announced in March to transform UK's computing capacity

A new supercomputer set to be one the most powerful in Europe is to be built in Bristol, in a move to drive pioneering AI research and innovation in the UK.

The UK government has confirmed the University of Bristol will host the new AI Research Resource (AIRR), which will serve as a national facility to help researchers maximise the potential of AI and support critical work into the potential and safe use of the technology.

The world-class AIRR cluster will vastly increase the UK's compute capacity – essential to achieving the UK's AI ambitions and securing its place as a world-leader in harnessing the rapidly developing technology. The cluster, which will be made up of thousands of state-of-the-art graphics processing units, or GPUs, will be able to train the large language models that are at the forefront of AI research and development today.

University of BRISTOL

https://www.gov.uk/government/news/bristol-set-to-host-uks-most-powerful-supercomputer-to-turbocharge-ai-innovation

(ARM based Isambard 1, 2) Isambard 3 in 2023/4 (no data centre and Isambard PI Simon MS with a non-dedicated GW4 team)



Isambard-AI procurement



Isambard hiring started



Service configuration and hardening for AI users



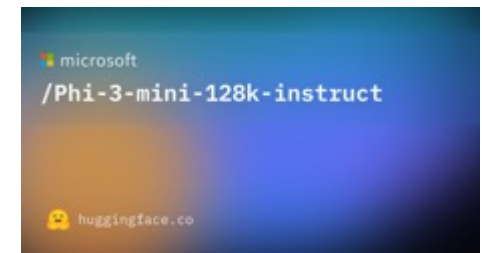| 2016–July 2023 | Aug. 2023 | Sep.–Oct. 2023 | Nov. 2023 | Dec. 2023 | Mar. 2024 | Apr. 2024 | May 2024 |

UK govt feedback on Bristol's Isambard-AI proposal

AI Safety Summit

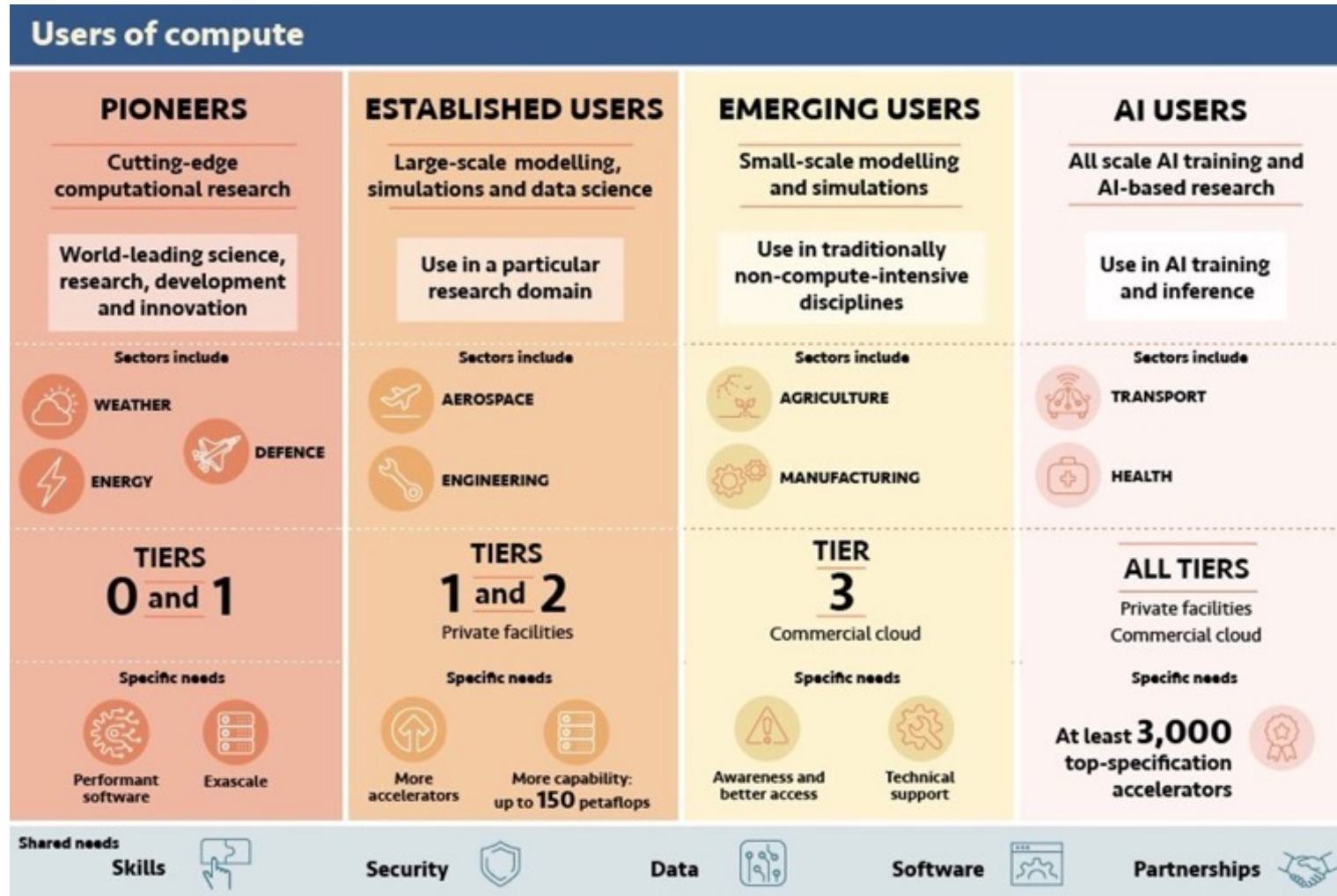Modular data centre (POD) and Isambard-AI phase 1 installed + a new team in ~3 months

Running AI and HPC applications









"Who is Isambard?" running on Isambard-AI


University of BRISTOL

# Design Specifications for AI Research Resource (RR)



| Users of compute | | | |
|---|---|---|---|
| **PIONEERS** | **ESTABLISHED USERS** | **EMERGING USERS** | **AI USERS** |
| Cutting-edge computational research | Large-scale modelling, simulations and data science | Small-scale modelling and simulations | All scale AI training and AI-based research |
| World-leading science, research, development and innovation | Use in a particular research domain | Use in traditionally non-compute-intensive disciplines | Use in AI training and inference |
| Sectors include: WEATHER, ENERGY, DEFENCE | Sectors include: AEROSPACE, ENGINEERING | Sectors include: AGRICULTURE, MANUFACTURING | Sectors include: TRANSPORT, HEALTH |
| **TIERS 0 and 1** | **TIERS 1 and 2** Private facilities | **TIER 3** Commercial cloud | **ALL TIERS** Private facilities, Commercial cloud |
| Specific needs: Performant software, Exascale | Specific needs: More accelerators, More capability: up to 150 petaflops | Specific needs: Awareness and better access, Technical support | Specific needs: At least 3,000 top-specification accelerators |

Shared needs: Skills, Security, Data, Software, Partnerships

Accessible to all users
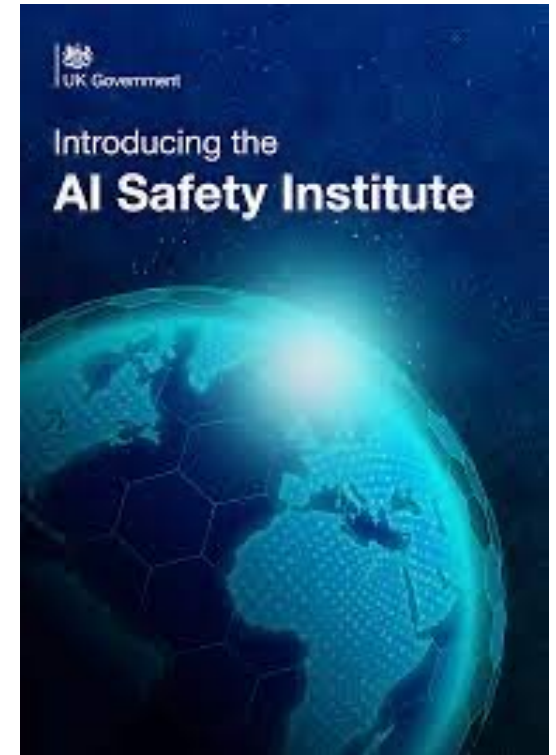
Sustainable AI supercomputing

Accessible to all sectors

Performance for all tiers

Sustainable AI RR program

References: Independent Review of The Future of Compute: Final report and recommendations, March 2023; National AI Strategy - AI Action Plan, July 2022; £300 million to launch first phase of new AI Research Resource

University of BRISTOL

# What is Isambard-AI for?

- UKRI-funded AI research in the UK, e.g.:
  - Training large language models
  - Large-scale inference
  - Foundational AI research
  - AI safety and understanding
  - Hybrid AI + simulation workflows
  - Machine learning
- Research on Isambard-AI must have a strong AI component
- Accommodate GPU jobs at any scale
  - Interactivity via JupyterHub—single to 100s of GPUs
  - Long running jobs for large-scale training—10s to 1000s of GPUs

# Sustainability as a Key Design Principle

- Optimisation targets
  - PUE = Power Usage Effectiveness
    - Target <1.1
  - CUE = Carbon Usage Effectiveness
    - Non-fossil fuel sources
  - Plan for heat reuse for nearby buildings and local district heat circuit in future
- Aligning with university of Bristol Net Zero and sustainability targets for 2030
  - Categorising emissions
    - Scope 1 (~0%), 2 (90%) and 3 (10%)—based on an average UK data of 0.2123 kg $CO_2$/kWh (IEA 2022 data)
  - Recycling 90% of components at the end of life in the UK

# Isambard 1, 2 & 3 – Leading ARM for HPC since 2016 as a UK national tier-2 resource

- Isambard 1 and 2 hosted at the UK Met Office data centre

- Options considered:
  - Renting space in a DC—£££ plus not available for hundreds of KW DLC cabinets like Cray HPE XE
  - Building new-–time and ££££

- Solution— containersied data centre or MDC



**HPC**wire
Since 1987 - Covering the Fastest Computers in the World and the People Who Run Them

- Home
- Topics
- Sectors
- Exascale
- Specials
- Resource Library
- Podcast
- Events
- Job Bank
- About

**Behind the Met Office's Procurement of a Billion-Dollar Microsoft System**
By Oliver Peckham

May 13, 2021

The UK's national weather service, the Met Office, caused shockwaves of curiosity a few weeks ago when it formally announced that its forthcoming billion-dollar supercomputer – expected to be the most powerful weather and climate-focused supercomputer in the world when it launches in 2022 – would come from an unlikely source: Microsoft. At the HPC User Forum yesterday, Richard Lawrence, an IT fellow for supercomputing at the Met Office, detailed the service's hunt for its next generation of supercomputing.

**Out with the old, in with the new**

University of BRISTOL

# What is an MDC

- Modular, self-contained & agile
  - Described as Lego blocks—designed and tuned for functional and performance specifications e.g. high availability, security, etc.
  - Everything from all IT (compute, network & storage), power (UPS, batteries) and cooling can be included as self-contained units

- Efficient and flexible deployments
  - Typically built and commissioned in months
  - Accommodates different environmental conditions
    - On-site integration options
    - Offsite integration options

- Sustainable
  - Holistic, fine-grain telemetry via DCIM
  - Upgradable, refreshable, recyclable with a lifespan on 10-15 years

The Isambard Site (Isambard AI phase 1 & 2, and Isambard 3)



University of BRISTOL

# Isambard Site – National Composite Centre (NCC) Facility in Bristol

- NCC—UK's Centre of Excellence for Composites Research and Development
  - Availability of power (~10 MW), networking and cooling
  - Heat reuse options
  - Co-location with industrial user community that has a digitalisation first approach

# Physical Space Constraints at NCC



**FIGURE 1.** HPE Cray EX cabinet exploded view

~5,300 GPUs in **12 EX4000 cabinets**

New Class of AI Supercomputer Connects 256 Grace Hopper Superchips Into Massive, 1-Exaflop, 144TB GPU for Giant Models Powering Generative AI, Recommender Systems, Data Processing

May 28, 2023



DGX SuperPOD with **256 GPUs in 16 cabinets**

University of BRISTOL

# HPE EX Series DLC and Nvidia GH 200

- **HPE EX solution**
  - Direct liquid cooling for high performance computing and networking

- **4-way Nvidia GHr superchip**
  - NVLink-C2C also only uses 1.3 picojoules/bit transferred—5x more energy efficient than PCIe Gen 5

Figure 5. Memory Accesses across NVLink-connected Grace Hopper Superchips

Source: NVIDIA Grace Hopper Superchip Architecture Whitepaper

**FIGURE 1.** HPE Cray EX cabinet exploded view

Source: HPE CRAY EX Liquid-Cooled Cabinet for Large Scale Systems brochure

University of BRISTOL

# Grace-Hopper Superchip & HPE EX Compute Blade



Source: NVIDIA Grace Hopper Superchip Architecture Whitepaper



Source: HPE EX4000 Grace-Hopper blade

4 x Grace ARM CPUs
288 cores
512 GB Memory

4 x Hopper GPUs
~260 64-bit Tflops, ~16k 8-bit Tflops
384 GB High Bandwidth Memory

896 GB Memory Total
NVLink-C2C = 900 GB/s

Isambard AI node = 4 x GH200
Injection bandwidth = 4 x 200 Gbps

University of BRISTOL

# Memory Architecture of GH 200

## A boon for developers and users



Figure 7.     NVIDIA Hopper System with Disjoint Page Tables

Figure 8.     ATS in an NVIDIA Grace Hopper Superchip System

Source:  NVIDIA Grace Hopper Superchip Architecture Whitepaper

In PCIe-connected x86+Hopper systems, the CPU and the GPU have independent per process page tables, and system allocated memory is not directly accessible from the GPU

Address Translation Service (ATS) enables the CPU and GPU to share a single per-process page table, enabling all CPU and GPU threads to access all system-allocated memory

University of BRISTOL

# HPE SlingShot High Speed Interconnect for AI & HPC

## Ethernet
- Standards based / interoperable
- Commodity technology
- Converged network
- Limited HPC features
- High latency
- Efficient for large payloads only
- Limited scalability for HPC

## Slingshot
- Standards based / interoperable
- Commodity technology
- Converged network
- Full set of HPC features
- Low latency
- Efficient for small to large payloads
- Very scalable for HPC & Big Data

## HPC Networks
- Proprietary (single vendor)
- Non-commodity technology
- HPC interconnect only
- Full set of HPC features
- Low latency
- Efficient for small to large payloads
- Very scalable for HPC & Big Data

Source: https://www.nextplatform.com/2019/08/16/how-cray-makes-ethernet-suited-for-hpc-and-ai-with-slingshot/

**Liquid cooled interconnect (sustainability & scalability)**
Example with 16-switch group
2 switches per chassis for single injection to 32 compute nodes (8 compute blades)
6 switches per cabinet for single injection to 256 compute nodes (64 compute blades)



Optical connections (inter-switch group)

Electrical connections (inter-switch group)

Up to 100s of cabinets

Cableless connections (direct switch to compute blade)

16 endpoints

Example: 16 endpoints per cabinet

Example: 256 endpoints per cabinet

16 endpoints

**FIGURE 7.** Example of Dragonfly topology in HPE Slingshot switches

Source: HPE CRAY EX Liquid-Cooled Cabinet for Large Scale Systems brochure

University of BRISTOL

| | |
|---|---|
| **Application** | **AI and ML Applications and Frameworks** |
| **Environment** | **NVIDIA Containers**<br>**Standard conda / pip environments**<br>**Custom conda / pip environments**<br>**Install / compile your own software** |
| **Interface** | **Notebooks and Dashboards** / **Job Scripts and Graphical Interfaces** |
| **Platform** | **JupyterHub** · **Kubeflow** · **Custom Platforms** · **Batch Jobs** · **Container Runtimes** · **VSCode**<br>**Kubernetes** · **Shell access (slurm)** |
| **Tenancy** | **Multi-tenant Partitions** |
| **Infrastructure** | **CSM – Cloud Native Supercomputing** |

# Isambard-AI >> Who is Isambard?



Isambard Kingdom Brunel was a renowned British engineer and architect who lived from April 9, 1806, to September 15, 1859. He is best known for his significant contributions to the development of the United Kingdom's infrastructure during the Industrial Revolution. Brunel designed and built numerous important structures, including the Great Western Railway, which connected London to the west of England and Wales. He also designed several iconic bridges, tunnels, and ships, such as the SS Great Britain, the first iron-hulled, screw propelled ship. Brunel's innovative designs and engineering feats have left a lasting legacy in the field of engineering.

Phi-3 Mini installed on pytorch through, pip, using cuda 11.8 on GH200 GPU, using ~7GB of HBM3.

Contact Wahab Kawafi a.kawafi@bristol.ac.uk for details

University of BRISTOL

# IsamBot running in a Jupyter notebook on Isambard-AI (spawned by JupyterHub on a rCN)

Contact James Womack
j.c.womack@bristol.ac.uk
for details



University of BRISTOL

# MLPerf Training Early Results (Bert-Large)

https://mlcommons.org/benchmarks/training/

Contact Wahab Kawafi for details:
a.kawafi@bristol.ac.uk



With Slingshot 2.1, libfabric 1.15.2.0 and GPU RDMA enabled. Promising single node results, but more fine tuning required to scale.
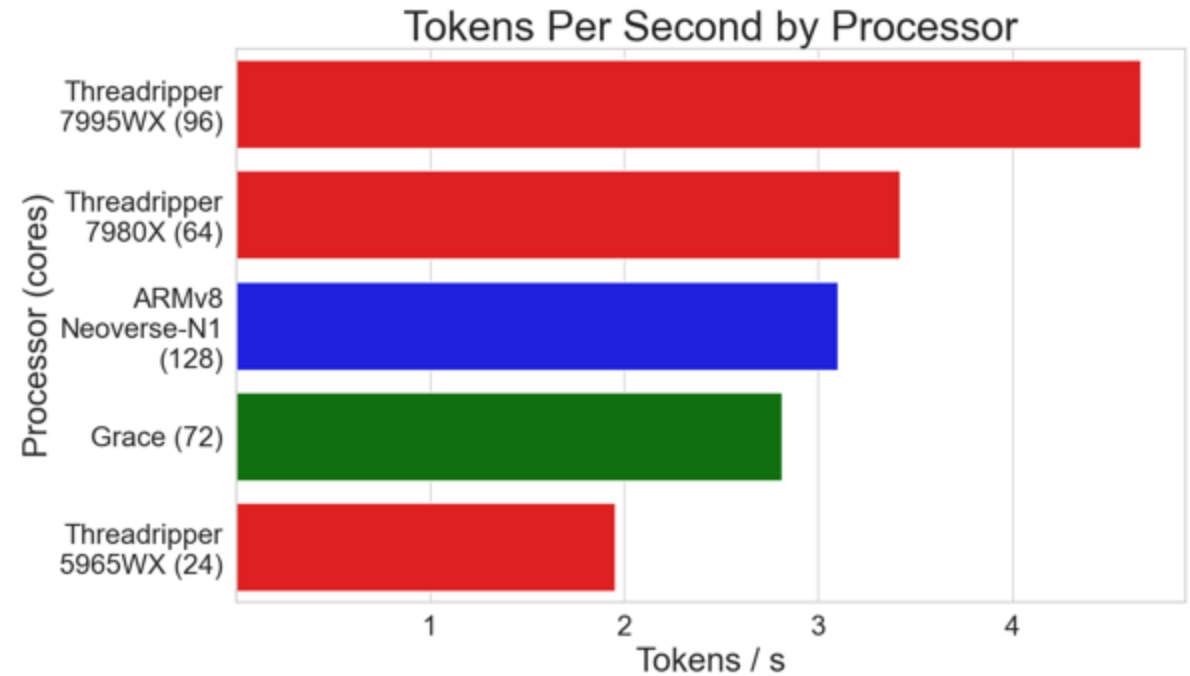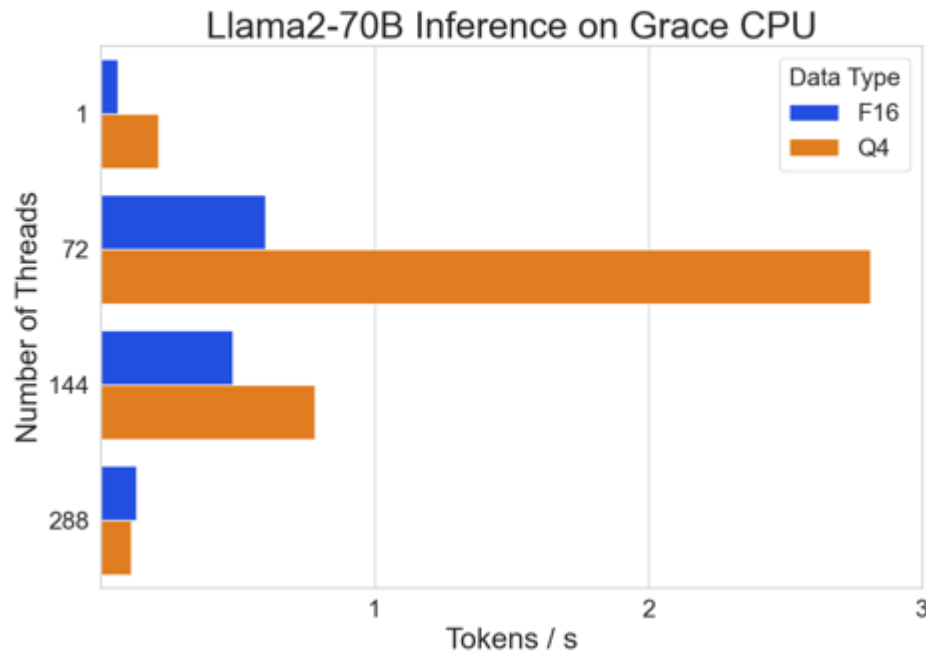
University of BRISTOL

# llama.cpp Benchmark Early Results

LLM inference on CPU
https://github.com/ggerganov/llama.cpp
OpenBenchmarking
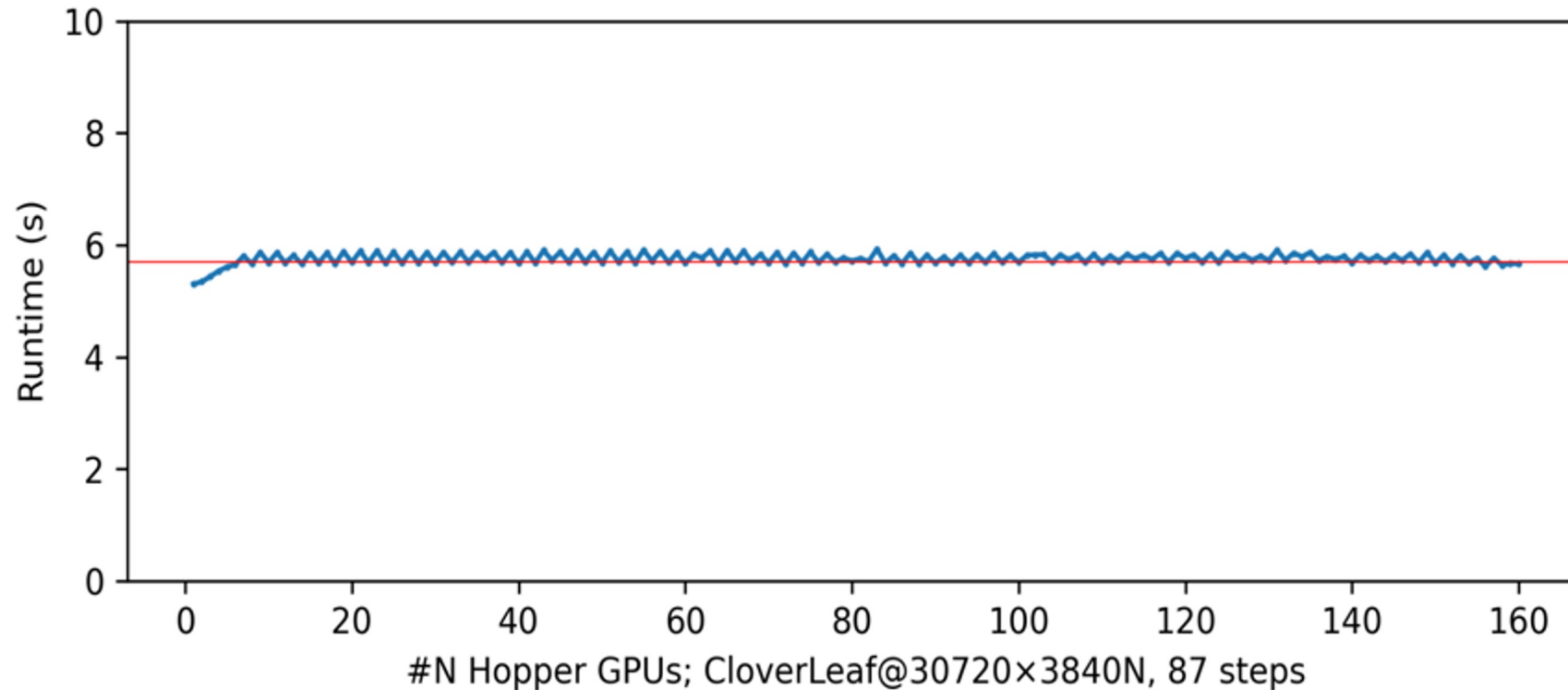
Contact Wahab Kawafi for details:
a.kawafi@bristol.ac.uk



Early results. Grace GH200 CPUs have promising results to complement inference on LLMs (70B) even with a relatively low thread count per socket.

University of BRISTOL

# CloverLeaf Benchmark Early Results

https://github.com/UoB-HPC/CloverLeaf

Part of SPEChpc2021, primarily mem-BW, structured grid, stencil pattern, we use the CUDA port
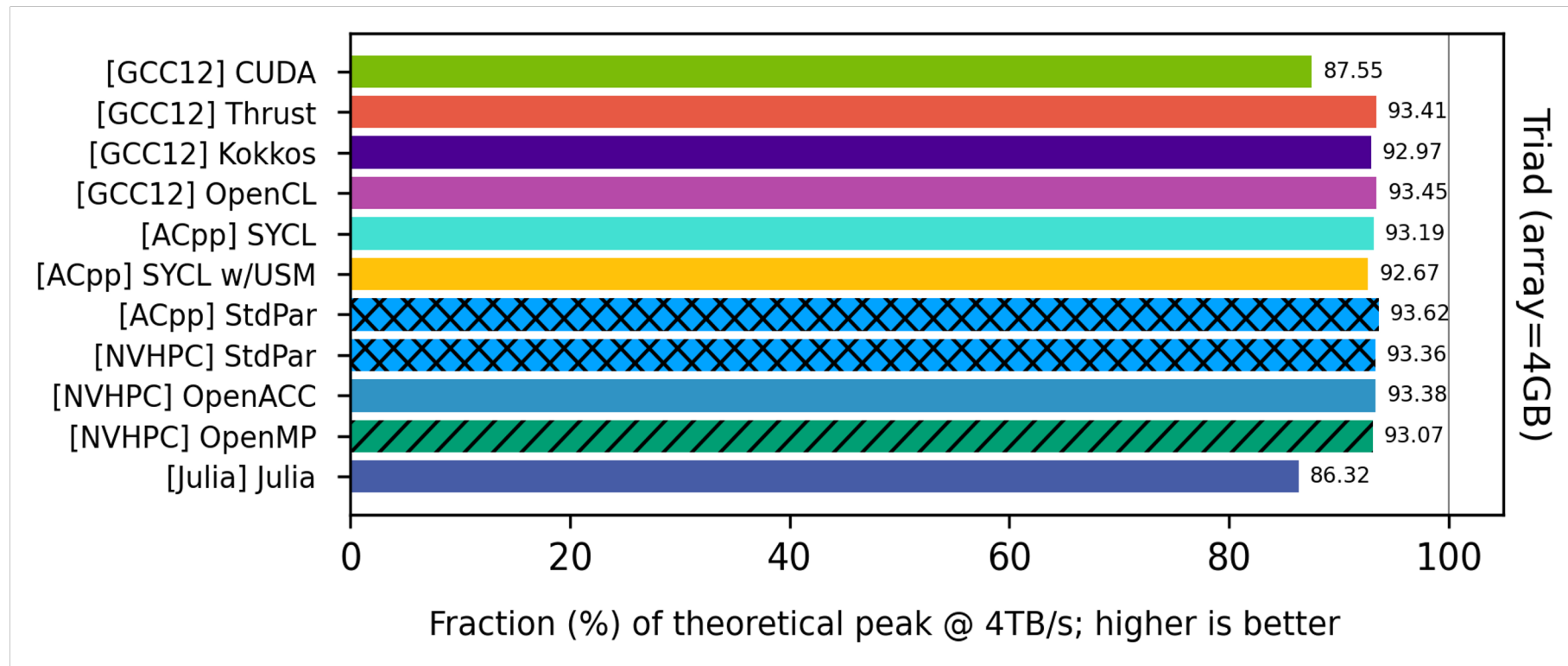
Contact Tom Lin for details: wl14928@bristol.ac.uk



Weak scaling up to almost the entire machine of 160 GPUs (2x4 GPUs are on login node)
Similar results for TeaLeaf (part of SPEChpc2021) primarily mem-BW, MPI collectives, SpMV

University of
BRISTOL

# PE and Memory Bandwidth Benchmarking

https://github.com/UoB-HPC/BabelStream

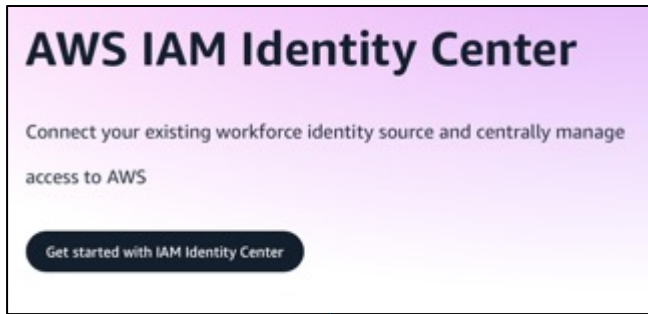Contact Tom Lin for details: wl14928@bristol.ac.uk



CUDA, Thrust, Kokkos, OpenCL, SYCL (via AdaptiveCpp), StdPar, OpenACC, Julia, OpenMP targets
Some of the compilers (GCC, Acpp, Clang) are built from source and everything worked as expected

# Lowering Access Barriers via IAM and Single Sign-On (SSO)

- OIDC single-sign on – bring your own <u>high-level trust</u> identity with federated academic & research credentials
  - Security via multi-factor authentication (MFA) for web & ssh
  - Okta a preferred option for govt public cloud AI users
- Self-service, cloud-native user and project management portal (single pane of glass for accessing all services plus accounting, reporting and audit trails)
- Waldur: single source of user truth
  - Provides Authorisation via OIDC
  - Manages projects, groups and roles
  - HTTPS (Connects to Waldur via standard OIDC)
  - An SSH key signing CA gets authorisation from Waldur (via OIDC)— Signs a short-lived SSH certificate
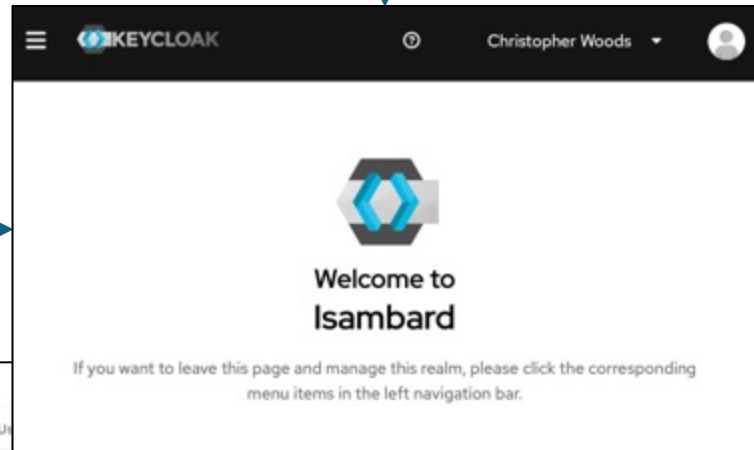
Administrator Identities

Academic Identities

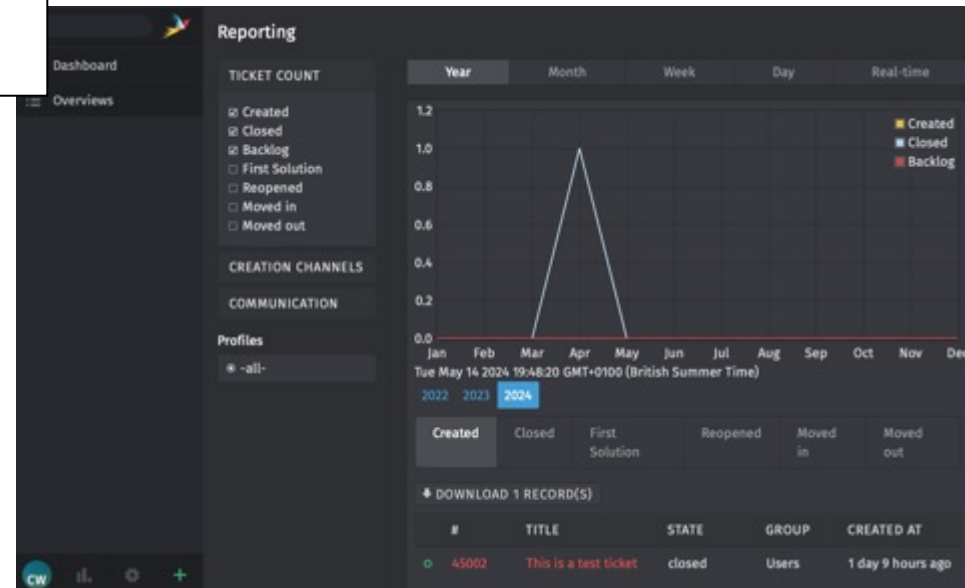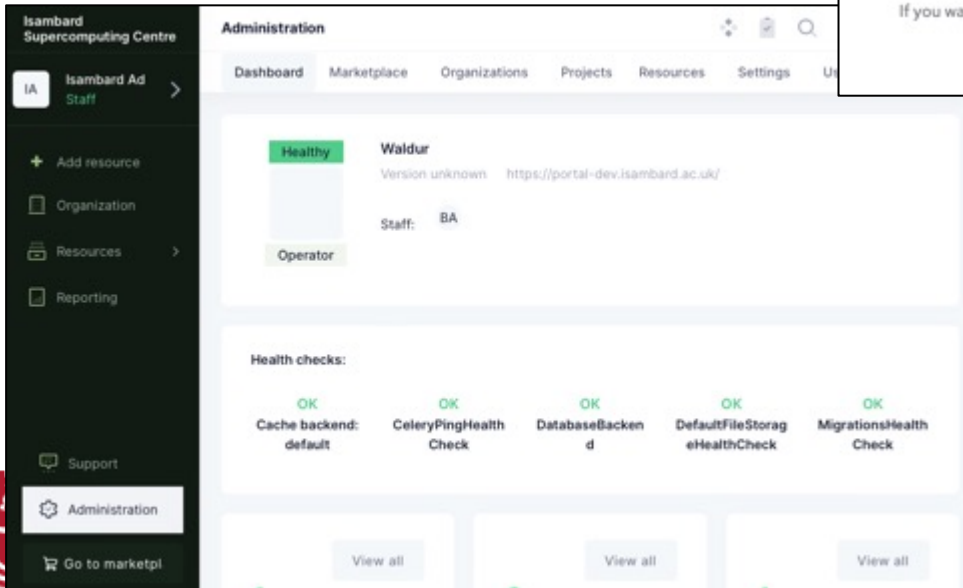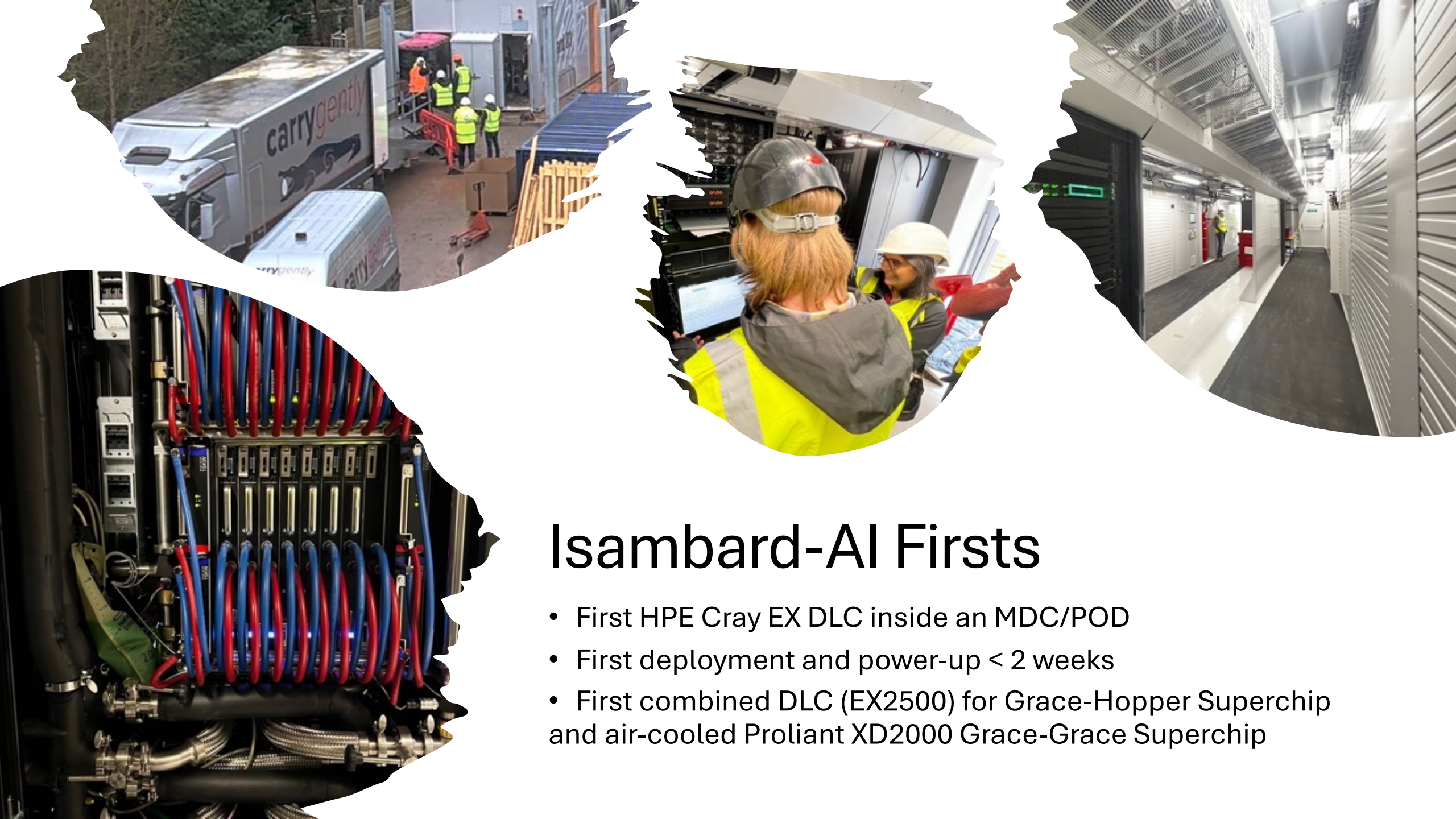Non-academic Identities

Waldur Portal

Zammad Helpdesk

Keycloak brokers identities, enabling users to bring their own identity to the service. Administrator identities are strictly separated and closely monitored.

# Isambard-AI Firsts

- First HPE Cray EX DLC inside an MDC/POD

- First deployment and power-up < 2 weeks

- First combined DLC (EX2500) for Grace-Hopper Superchip and air-cooled Proliant XD2000 Grace-Grace Superchip

# scaling out in two phases

| Phase 1 (~0.7 8-bit AI Exaflops) | Phase 2 (~21 8-bit AI Exaflops) |
|---|---|

**Phase 1 (~0.7 8-bit AI Exaflops)**

Arrived in March 2024 – in Isambard 3 MDC
Piloting, on-boarding and staging services

**1 x DLC EX2500 cabinet**

21 blades (4-way Grace-Hopper)

42 nodes

168 GH superchips

12,096 Neoverse V2 Armv9 CPU cores

168 Hopper GPUs

21.5 TB CPU memory

16.1 TB high bandwidth GPU memory

37.6 TB total memory

**AI high performance storage**

~1 PB all-flash ClusterStor Lustre

**Phase 2 (~21 8-bit AI Exaflops)**

Arriving Summer 2024 – new Isambard-AI MDC
Delivery of AI services

**12 x DLC EX4000 cabinets**

660 blades (4-way Grace-Hopper)

1,320 nodes

5,280 GH superchips

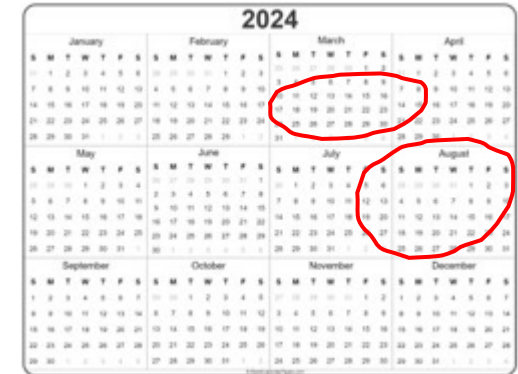380,160 Neoverse V2 Armv9 CPU cores

5,280 Hopper GPUs

675 TB CPU memory

506 TB high bandwidth GPU memory

1.18 PB total memory

**AI high performance storage**

~27 PB all-flash storage!

(~20 PB Lustre, ~7 PB software defined VAST)

University of
BRISTOL

# Thank you

Stay tuned!

University of BRISTOL



UK Government

AI SAFETY SUMMIT

**THE BLETCHLEY DECLARATION**

WORLD FIRST AGREEMENT ON SAFE AND RESPONSIBLE DEVELOPMENT OF FRONTIER AI

- 28 COUNTRIES FROM ACROSS THE GLOBE, AND THE EU

- IDENTIFYING AI OPPORTUNITIES AND RISKS

- BUILDING A SHARED UNDERSTANDING OF THESE RISKS

- INTERNATIONAL COLLABORATION ON SCIENCE AND RESEARCH