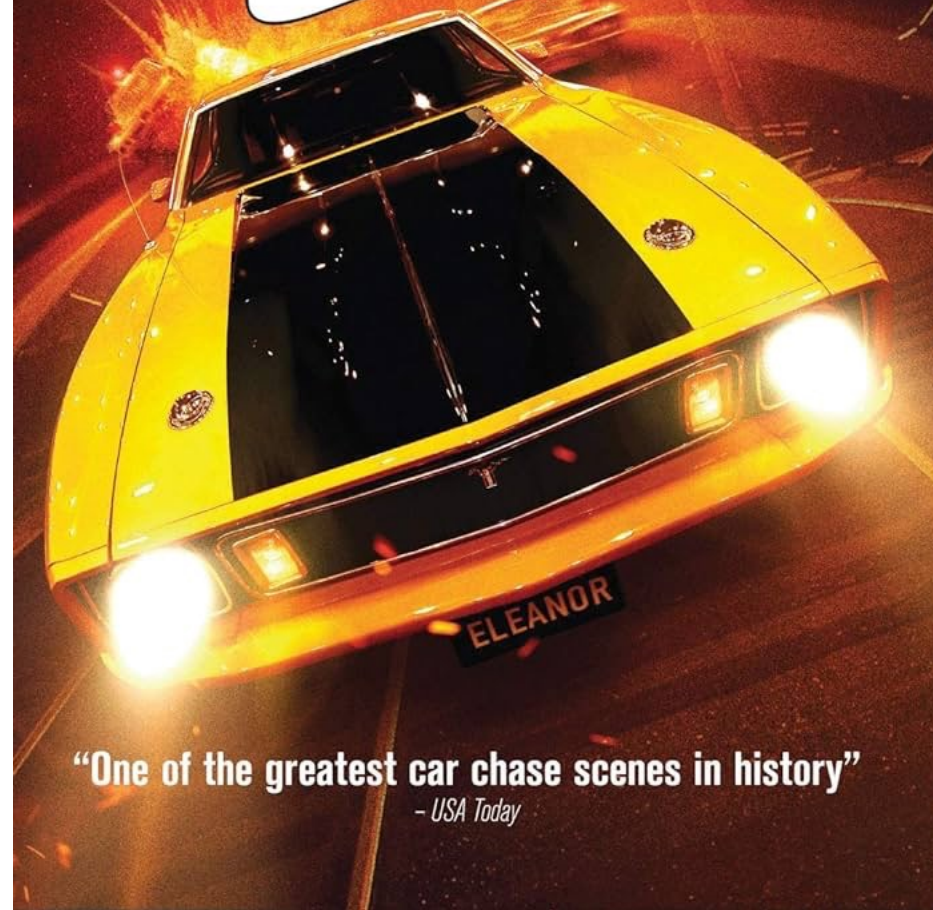


Isambard-AI: Extremely rapid deployment of leadership-class AI



Prof Simon McIntosh-Smith (PI)
Director, Bristol Centre for Supercomputing (BriCS)

HERE IN **60** WEEKS
THE ORIGINAL CLASSIC



“One of the greatest car chase scenes in history”

- USA Today

What is Isambard-AI?



- **>£300M** investment by UK Government in AI capability
- Funding **~5,500 NVIDIA Grace-Hopper GPUs** in a new, 5MW HPE modular data centre (**MDC**) facility in Bristol, UK
 - ~21 ExaFLOP/s of 8-bit for AI, ~250 PFLOP/s 64-bit
- **Extremely rapid deployment** a key requirement:
 - **First conversation** with UK Government on **Aug 18th 2023**
 - **>£200M procurement** written in **1 week**, run in just **2 weeks**
 - Contract signed and ground broken in November 2023
 - Site chosen for power availability and rapid planning permission procedures (8 weeks)



National Composites Centre

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

- 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

March 13th 2024



March 13th 2024



March 13th 2024



Winds gusting to 20mph

20 tonnes

Millions of Pounds worth
of MDC 30ft off the ground

March 14th 2024

168 GPUs installed

HPE EX2500

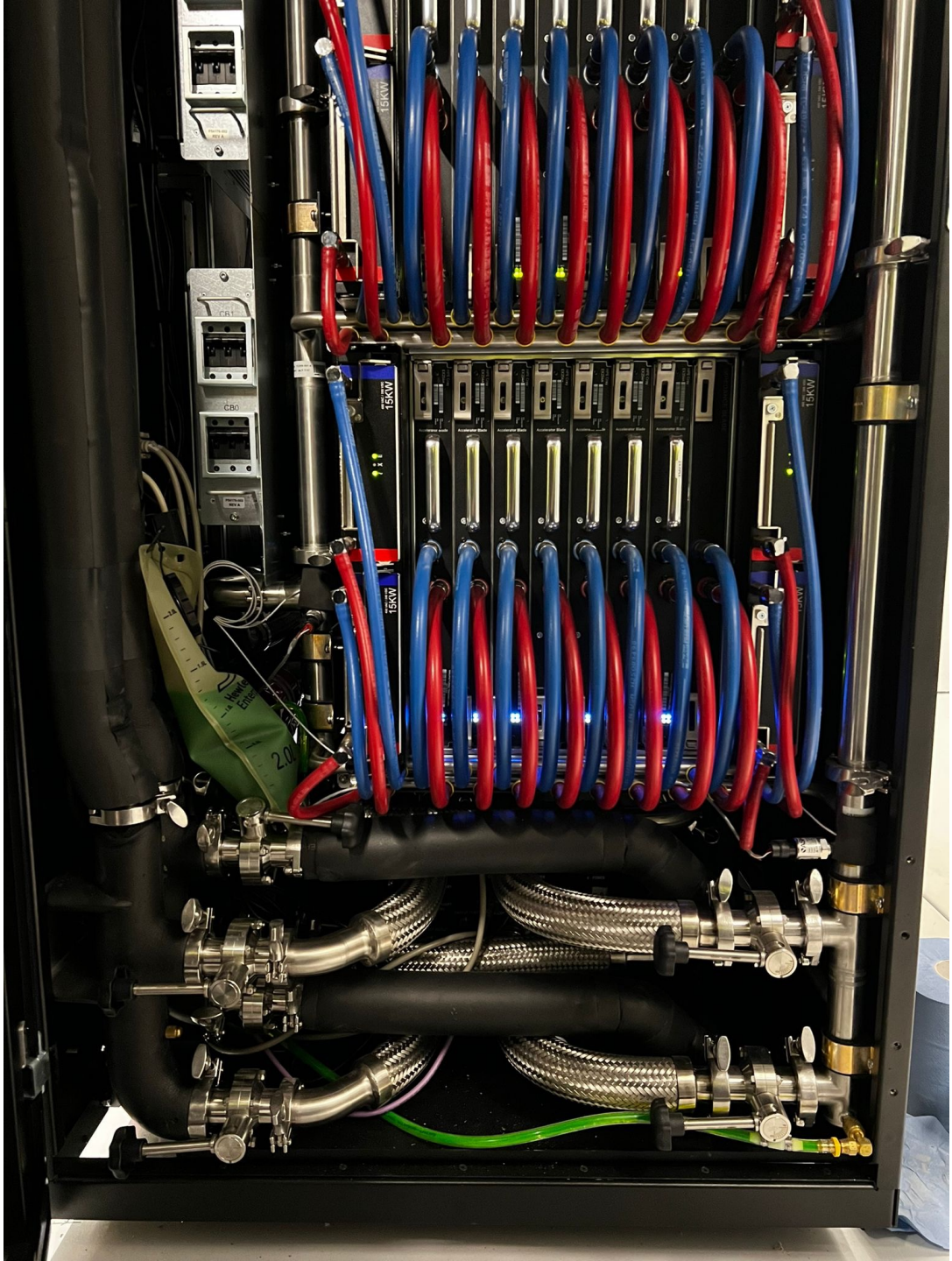
NVIDIA Grace-Hoppers



March 27th 2024





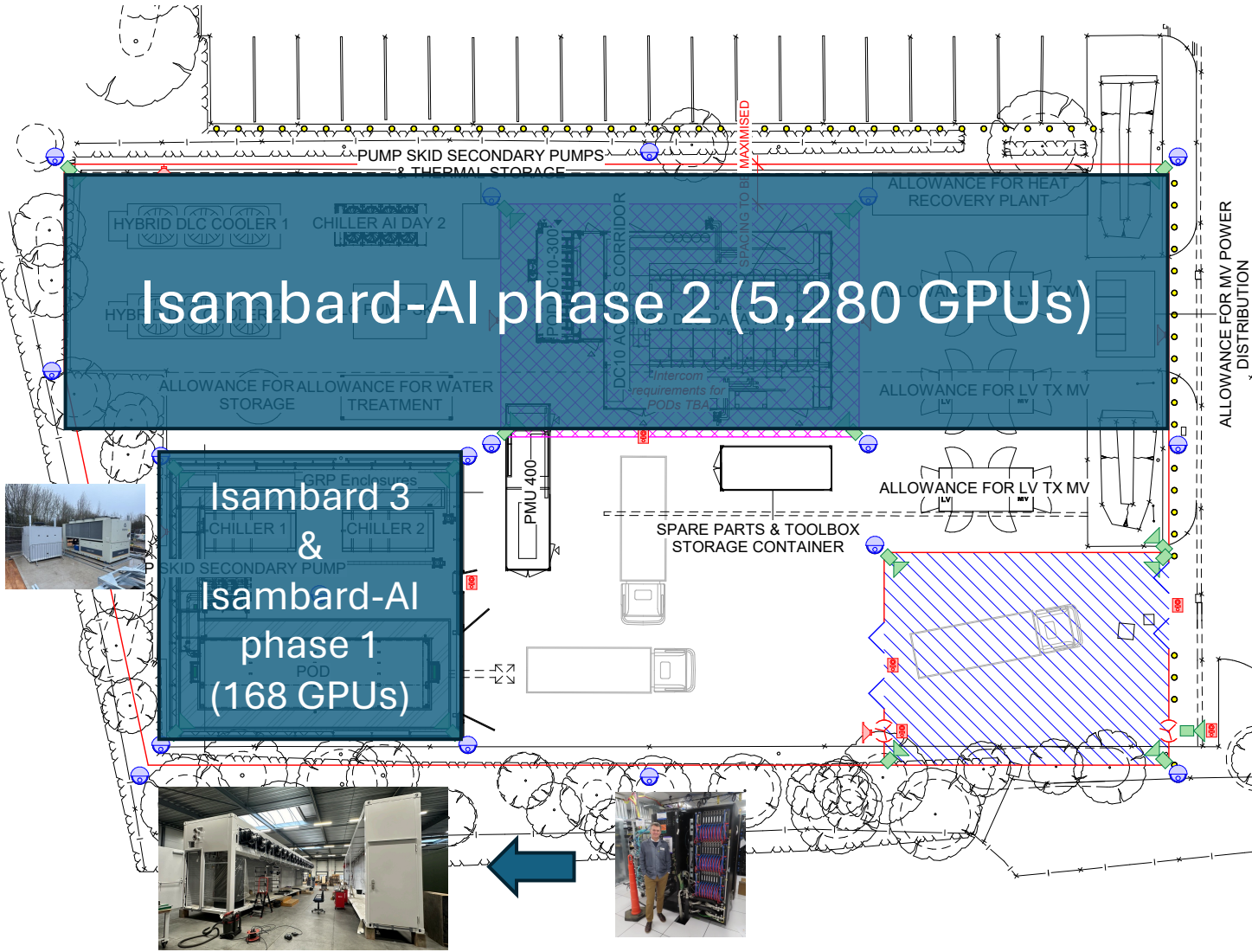






**WHAT
HAPPENS
NEXT?**

The Isambard Site



Isambard-AI AIRR:

- Phase 1 arrived in March 2024
- Phase 2 arrives August 2024
- >£200M including Modular Data Centre, all cooling etc.
- >£300M total investment over 5 years
- 5MW
- Extremely energy efficient, PUE <1.1
- Almost entirely direct liquid cooled, plans to reuse the waste heat

MDC before assembly

Isambard-AI phase 1 168 Grace-Hopper cabinet

HPE EX Series DLC and GH200

- HPE EX4000 solution
 - Direct liquid cooling for high performance computing and networking
- 4-way Grace-Hopper nodes

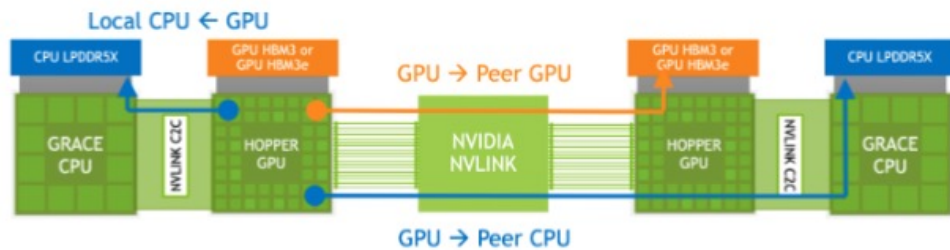


Figure 5. Memory Accesses across NVLink-connected Grace Hopper Superchips

Source: NVIDIA Grace Hopper Superchip Architecture Whitepaper

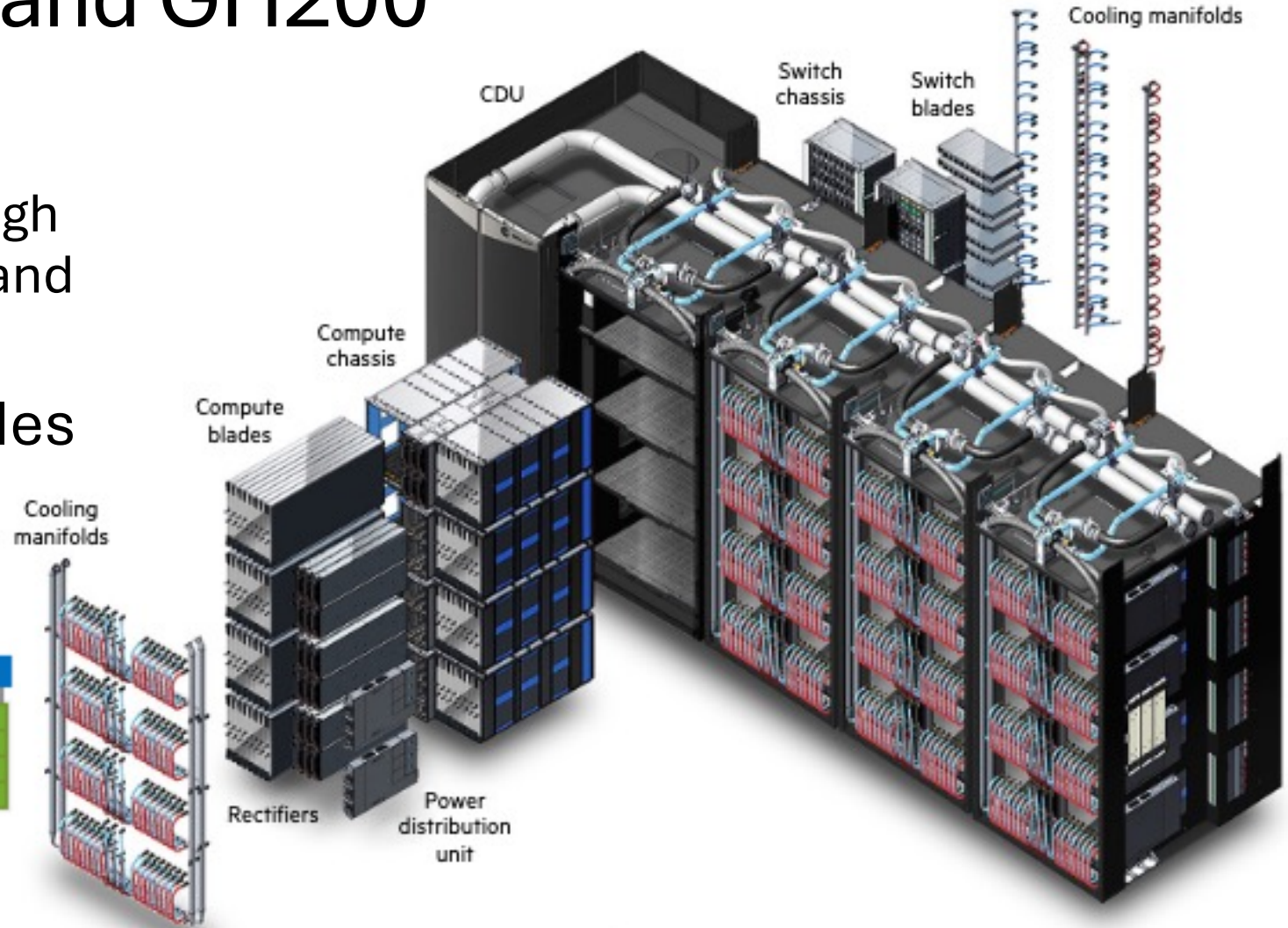


FIGURE 1. HPE Cray EX cabinet exploded view

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

What is Isambard-AI for?

- 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

AISI | AI SAFETY
INSTITUTE

**UK
RI** UK Research
and Innovation



Department for
Science, Innovation,
& Technology



For more detail, see our upcoming CUG paper!

Isambard-AI: a leadership class supercomputer optimised specifically for Artificial Intelligence

Simon McIntosh-Smith
University of Bristol
Bristol, United Kingdom
S.McIntosh-Smith@bristol.ac.uk

Sadaf R Alam
University of Bristol
Bristol, United Kingdom
Sadaf.Alam@bristol.ac.uk

Christopher Woods
University of Bristol
Bristol, United Kingdom
Christopher.Woods@bristol.ac.uk

Abstract — Isambard-AI is a new, leadership-class supercomputer, designed to support AI-related research. Based on the HPE Cray EX4000 system, and housed in a new, energy efficient Modular Data Centre in Bristol, UK, Isambard-AI employs 5,448 NVIDIA Grace-Hopper GPUs to deliver over 21 ExaFLOP/s of 8-bit floating point performance for LLM training, and over 250 PetaFLOP/s of 64-bit performance, for under 5MW. Isambard-AI integrates two, all-flash storage systems: a 20 PiByte Cray ClusterStor and a 3.5 PiByte VAST solution. Combined these give Isambard-AI flexibility for training, inference and secure data accesses and sharing. But it is the software stack where Isambard-AI will be most different from traditional HPC systems.

signature, design, build, order, integration, delivery to operational lights on inside 4 months. The early results from the Isambard-AI phase 1 system validate our design choices, from the data centre to direct liquid cooled cabinets to the AI optimised compute ecosystem.

At the time of writing this manuscript in the first week of April'24, we have demonstrated that a modular data centre can be ready in 2-3 months. The site work started in December'23 for preparing the concrete base for the placement of the HPE Modular Data Centre (MDC) technology, called Performance Optimised Datacentre (POD) [21]. The model that was initially



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

Thank you

Stay tuned!