

# The Durham HPC Hardware Lab

Institute for Computational Cosmology

Durham University

Alastair Basden, Peter Draper, Mark Lovell,  
Richard Regan, Paul Walker

# Contents

- COSMA
- DiRAC
- The HPC Hardware Lab
  - Overview
  - Location
- Hardware access
- Hardware details

# Alastair Basden

- Dept of Physics
  - Institute for Computational Cosmology
- HPC Manager - COSMA
  - DiRAC national facility
- DiRAC Technical Directorate



- The COSmology MACHine
  - Est. 2001
- Run by the ICC on behalf of DiRAC
  - Primarily covering STFC science areas
- Largest HPC system in the country
  - (by some metrics)
- 3 generations in operation
- Newest being COSMA8
  - ~70k cores, 0.5PB RAM, 20PB storage, 528 nodes



# DiRAC

- Established 2009
- Provides HPC to the STFC theory community
  - Particle Physics, Astrophysics, Cosmology, Solar System and Planetary Science and Nuclear Physics
- Three services:
  - Extreme Scaling: Edinburgh (TURSA)
  - Data Intensive: Cambridge and Leicester (CSD3 and DiaL)
  - Memory Intensive: Durham (COSMA)
- Co-designed and tailored for specific workloads
  - Bespoke systems for the science being carried out

# HPC Hardware Lab

- Mission: Provide access to the latest HPC hardware to users from across the UK
  - For code testing, performance tuning and debugging
  - To advise on purchase of future technologies
  - To allow informed decisions to be made whenever funding appears

# History

- Came together almost by accident
  - 2019: Intel provide a 56-core 6TB Cascade Lake system
    - For testing Non-volatile DIMM performance (Apache Pass)
  - 2019: University funding for the DINE cluster
    - BlueField DPU test system, 16 nodes
  - 2019: ExCALIBUR announced, with a H&ES component
    - Hardware and Enabling Software: \$4.5m
    - The UK £45m “preparation-for-exascale” fund

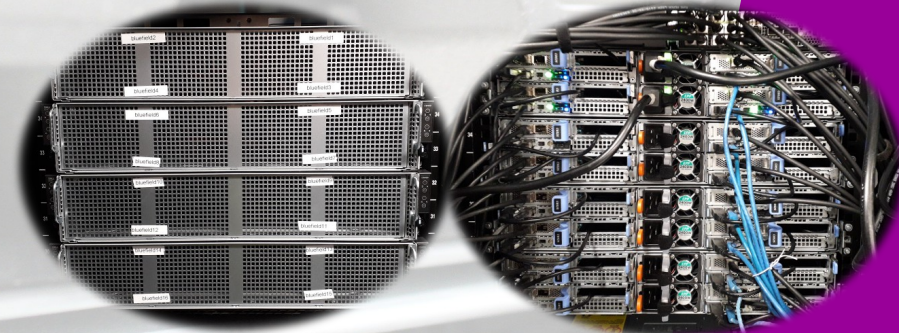
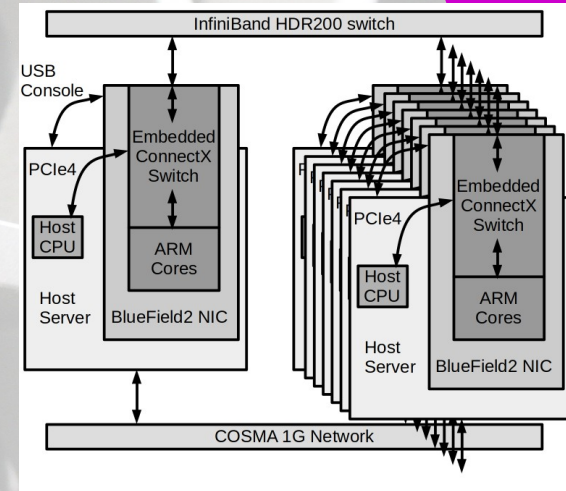
# Hardware Lab Components

- DINE: Durham Intelligent Networking Environment
- DINE2
- CPU compute
- GPU compute
  - NVIDIA, AMD, Intel
- Composability
- Rockport 6D torus network fabric
- Storage laboratory
- Environmental
  - Solar panel installation
  - Immersion cooling
  - Heat storage
  - Logging and awareness
- Quantum
- Leading to bespoke system design



# DINE

- Durham Intelligent Network Environment
  - A 24 node (initially 16) system for investigation of networking technologies
    - And other things
  - 32 cores, 512GB RAM per node
- Has hosted:
  - BlueField-1
  - BlueField-2
  - Rockport Ethernet
- UK's first production AMD EPYC HPC system
- Funded by Durham



# DINE-2

- **Durham Investigatory Node Environment**
  - 8 node Intel Sapphire Rapids system
    - 64 cores, 2TB RAM per node
  - Currently hosts a CerIO composable PCIe fabric
    - 8x A30 GPUs, assignable in any number to any host
- **Funded by DiRAC, IRIS and SKA**

# CPU compute

- Providing users with access to cutting edge CPU technologies:
- Coming soon: AMD Turin
- AMD Genoa and Bergamo
- NVIDIA Grace
- Intel Sapphire Rapids
- AMD Milan-X (extreme cache version)
- AMD Milan, Rome
- Intel Cascade Lake (with Apache Pass RAM, 6TB)
- Funded by OEMs, DiRAC, ExCALIBUR

# GPU Compute

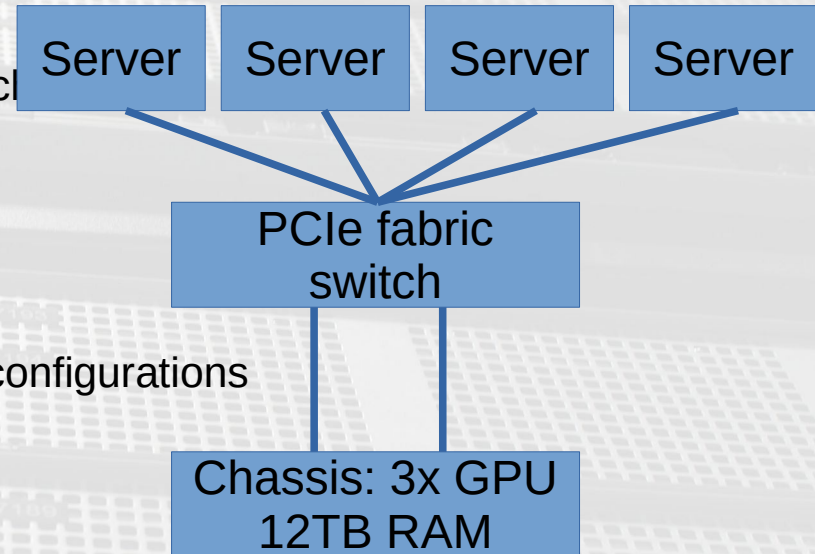
- Access to small numbers of latest GPUs
- Coming soon: AMD MI300X
  - And hopefully MI300A
- AMD MI210, MI100, MI50
- NVIDIA H100 (Grace-hopper system)
- NVIDIA A100, A30, V100
- Intel Ponte Vecchio
- Direct and queue-based access
- Funded by Dell, AMD, Intel, IRIS

# Composability

- Infrastructure-as-a-service
  - The ultimate goal for cloud-type systems
  - Is it relevant for HPC?
    - How does performance suffer?
    - Is it stable?
    - What are the use cases?

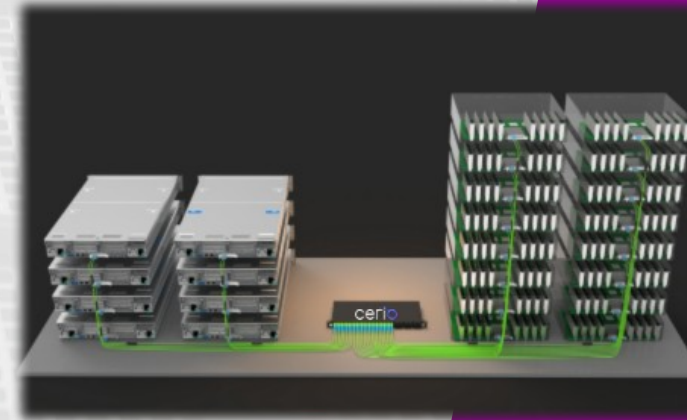
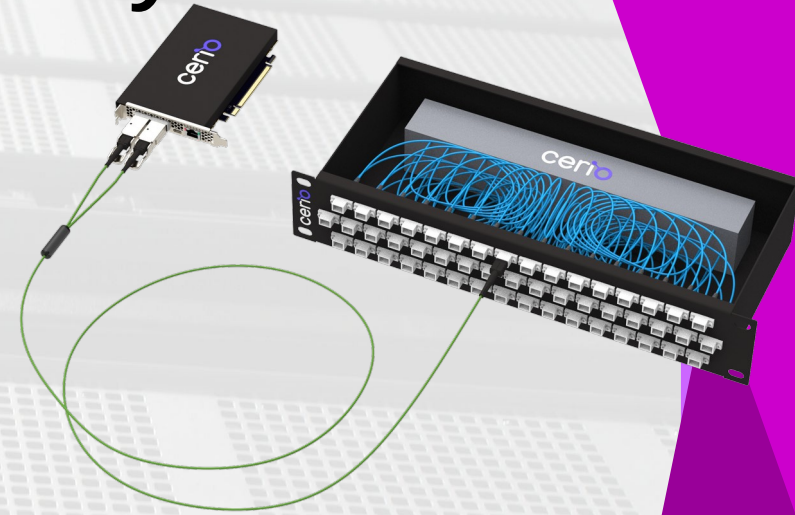
# Liquid composable system

- Installed 2021
- A PCIe4-based composable fabric: 4 lanes to each
- 3 A100 GPUs shared between 4 servers
  - Including a login node
  - GPUs per node can be changed in a few clicks
- 12TB RAM shared between these servers
  - Can be changed and reconfigured
- Positives: It works, can allow high RAM or GPU configurations
- Negatives:
  - Bandwidth is shared
  - GPUs sometimes fail requiring a full stack reboot
  - RAM/kernel issues (and no Rocky9 support yet)
  - Rack-scale limitations
  - Bottlenecks
- ExCALIBUR funded



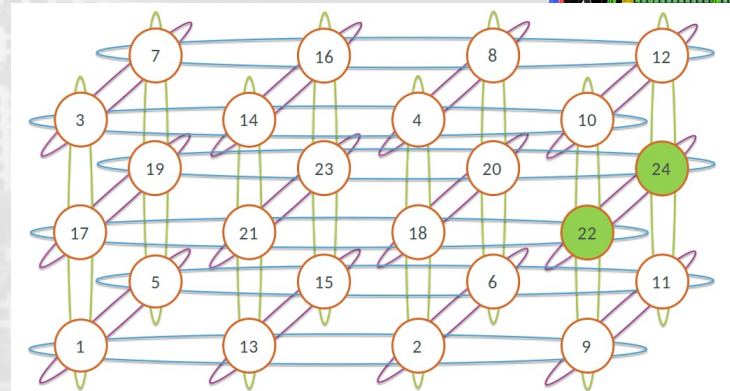
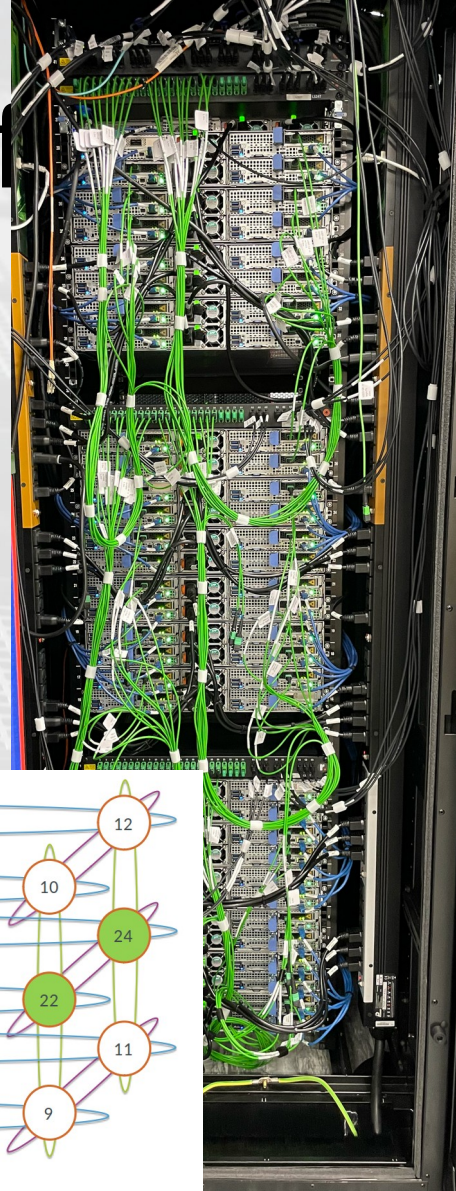
# CerIO composable system

- Installed 2024
- PCIe5-based fabric
  - No central switch: Uses a flit-based torus topology
  - Full data centre scalability
  - 200Gb/s to nodes
  - 300Gb/s inter-card bandwidth
- 8 compute nodes, 8x A30 GPUs
- Will allow networking and composability within a single fabric
- IRIS/SKA/ExCALIBUR/DiRAC



# Rockport 6D Torus Ethernet fabric

- A “switchless” fabric for 100G Ethernet
- Trailed on DINE in 2021
- Installed on COSMA7 in 2022
  - 224 nodes (half the cluster) replaced IB
  - Allows direct comparison of fabrics
  - At full HPC problem-size scale
- Works well
  - Performance comparable to InfiniBand
    - For real workloads
  - Handles congestion well
- ExCALIBUR/DiRAC funded





# Storage sub-lab

- Various different storage technologies
  - Most make it into production
- High-performance scratch Lustre (NVMe)
- DAOS (NVMe)
- Ceph
- StorJ private cloud
- VAST (NVMe)
- Globus (data transfer)
- Lustre (efficient bulk storage)
- Tape
- Funded by DiRAC/IRIS/SKA



# The snap file systems

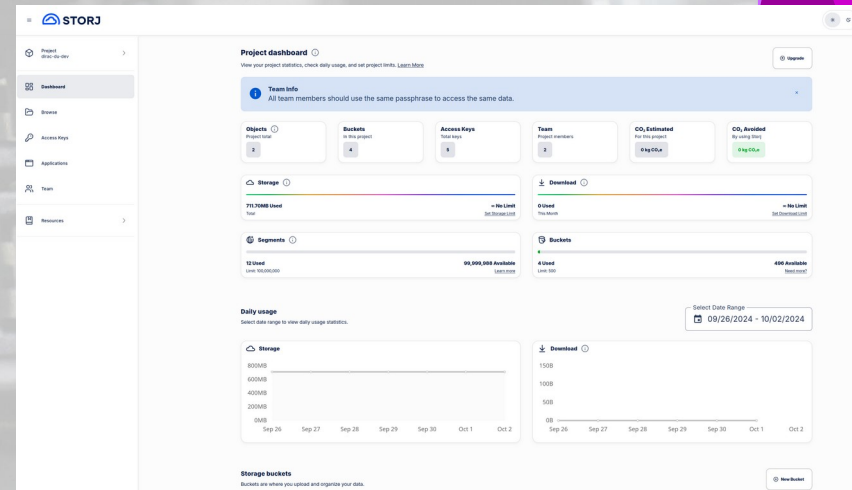
- 1.2PB fast parallel file system
  - Lustre
  - ~400 GBytes/s performance
- Use case: Dumping simulation snapshots
- 25 server nodes, each with 8x 6.4TB NVMe drives

# DAOS

- A new(ish) open-source file system from Intel
  - 4-node system
  - Providing bulk storage to DINE
  - Fast performance at low latency

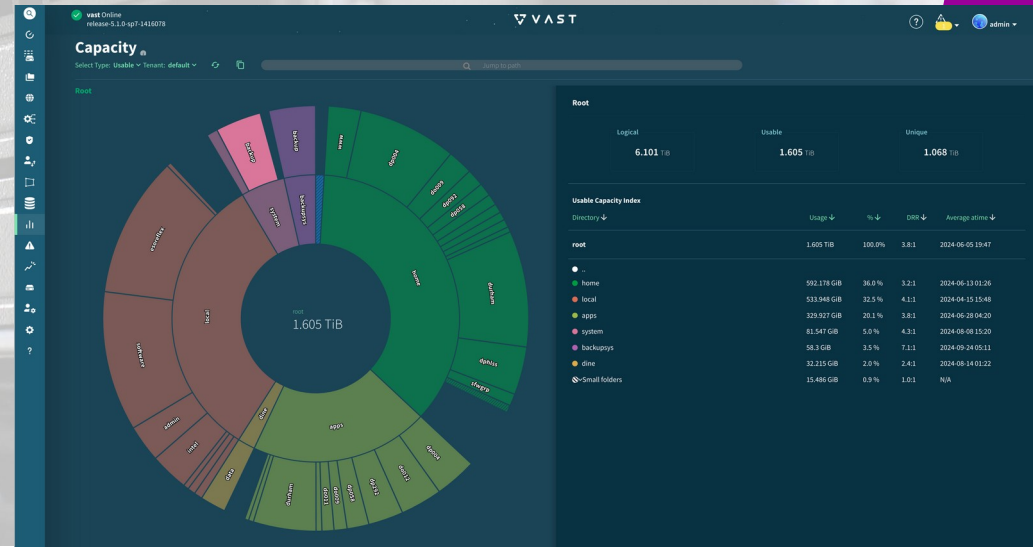
# StorJ

- A distributed storage object cloud provider
- We are working closely with them to have a private cloud instance spanning the 4 DiRAC sites
  - Good performance for hosting buckets
  - Ideal for data sharing with collaborators



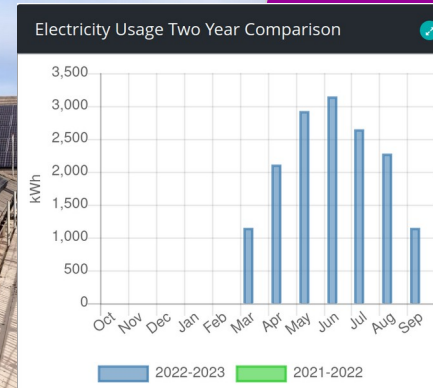
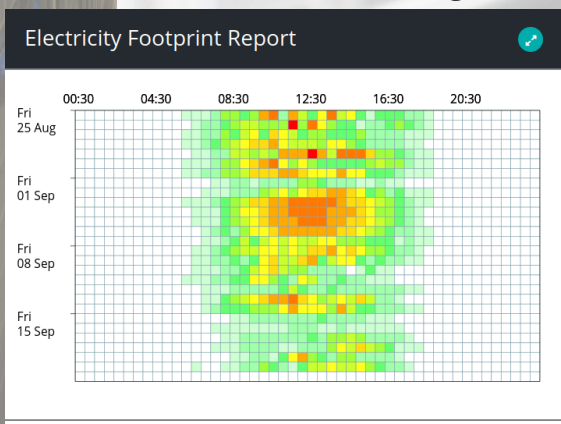
# VAST

- A new NVMe-based flash file system
  - Offering high performance with deduplication
    - Thus good data compression ratios (currently ~4:1)
- COSMA homespace
- IRIS/SKA



# Environmental-related

- HPC is a huge energy user
  - COSMA ~1MW at peak
  - Responsibility to keep this as low as possible
- 2023: Installation of ~£1m solar panels
  - Funded by DiRAC
  - Investigation into the interplay between supply and demand



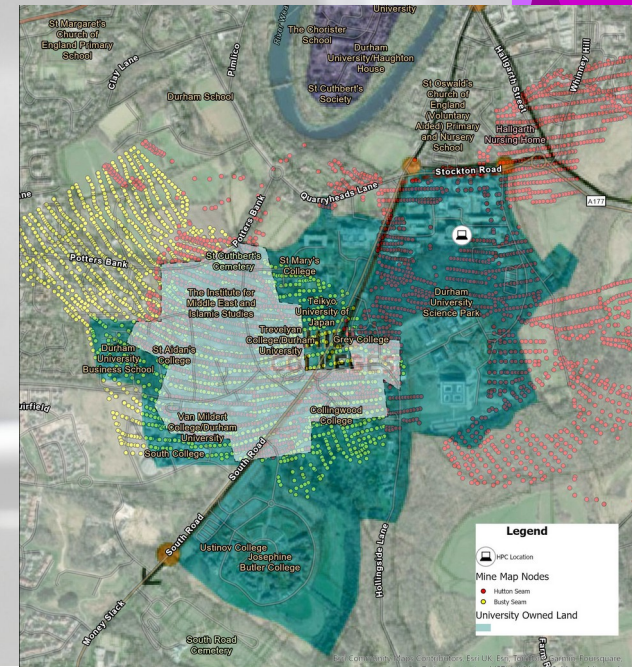
# Cooling sub-lab

- A sequence of technologies:
  - Passive cooled rear doors (pre 2010)
  - Active cooled rear doors (2018)
  - Direct liquid cooling (2020)
  - Immersion cooling (2024-5)
    - As a national object-of-study
    - Support for visits to Durham for operators to learn this technology
    - Reduced operational and embodied CO2
    - EPSRC DRI funding



# Mine water heat storage

- We are sitting on old, flooded mine workings
- HPC produces a lot of heat
  - We can heat buildings with this in the winter
    - (why don't we!?)
  - What can we do with the heat in the summer?
    - Store it underground for extraction in the winter
- This project will investigate feasibility
  - In particular, how fast does the water flow?
- EPSRC funded



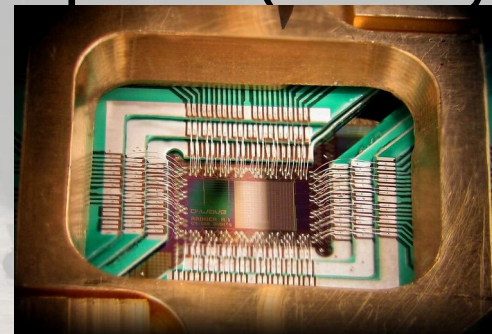


# Energy awareness

- Quarterly reporting to users around energy usage
  - Ability to query energy usage for each submitted job
- Monitoring of power on a node, rack, room and system scale
- Power-down of unused nodes

# Quantum computing

- No quantum compute in Durham (yet)
- However, the hardware lab controls national access to:
  - DWAVE quantum annealer
  - QuEra neutral atom quantum computer (shortly)
- ExCALIBUR funded



# Accessing the hardware lab

- Sign up on SAFE:
  - [safe.epcc.ed.ac.uk/dirac](http://safe.epcc.ed.ac.uk/dirac)
- Apply to join an appropriate project code:
  - do009: General purpose
  - do015: Cerio compasable system
  - do016: NVIDIA GPUs
  - do017: Intel GPUs
  - do018: AMD GPUs
- And feel free to arrange a visit to the data centre

# Leading to bespoke system design

- Key outputs from the hardware lab are:
  - Up to date knowledge of performance on new technologies
  - Experience profiling and optimising codes
  - Code preparation for future systems
  - Training on new technologies and tools
  - User awareness
  - Input into future system design

# Future plans

- MI300X and possibly MI300A systems
- UntetherAI card
- Turin CPU
- CXL composable systems
- Ultra-Ethernet fabric
- Funding:
  - DiRAC, IRIS, UKRI calls, ExCALIBUR-2?, Computer Science?, Physics

# Conclusion

- The Durham HPC Hardware Laboratory
  - Accessible for UK researchers
    - Single login
  - Cutting edge technologies
  - Let us know if there is something of particular interest