

Evaluation of the NVIDIA Grace Superchip in the HPE/Cray XD Isambard 3 supercomputer

Thomas Green, Sadaf Alam, Richard Gilham,
Simon McIntosh-Smith, Will Wishart

Bristol Centre for Supercomputing (BriCS)

Introduction

- Brief history
 - Isambard service
- Details on Isambard 3
 - Nvidia Grace CPU Superchip
- Motivation
 - Self-service approach
- Methodology
 - Tools used
- Results
- Future



Brief history

- **2018** – First Isambard system in production with Arm-based ThunderX2 processors
- **2020** – Isambard 2 launched with increased capacity of ThunderX2, included a small-scale A64FX cluster.
 - Hosted at the **Met Office** in Exeter, UK.
 - **328** nodes or 20,992 ThunderX2 cores (featured 4-way SMT)
- **2024** – Isambard 2 decommissioned
 - Funded by **EPSRC** (UK research agency)
 - In collaboration with **GW4** universities.



What is Isambard 3?

- A new general purpose air-cooled CPU HPC machine, ~300kW
- Based on **NVIDIA Grace CPU Superchip**
- Delivered by HPE
- 384 nodes, 55,236 cores
- 2 PBytes storage, Slingshot 11 network
- Production from Jan 25
- Funded by UKRI
- Collaboration with GW4 universities.
- Hosted at Bristol (alongside **Isambard-AI**)



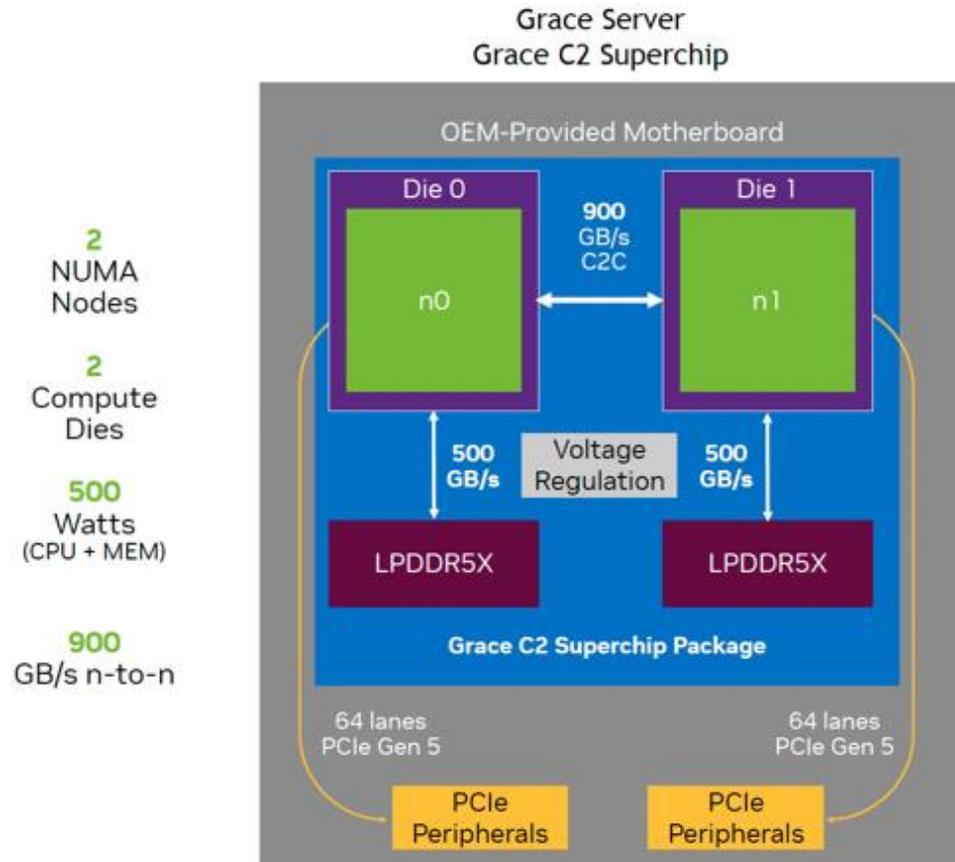
UK Research
and Innovation



Isambard 3 technical summary

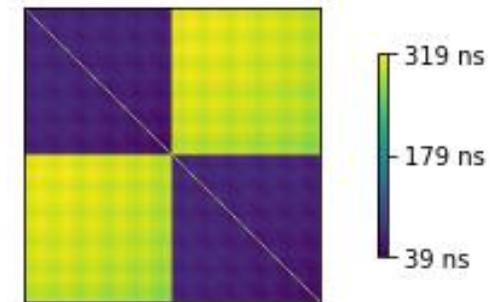
- **55,296** Armv9 cores, **384** nodes, 6 racks of NVIDIA Grace CPU Superchip (2x72 core @3.2GHz)
 - ~2 PFLOP/s HPL, just outside the Top500
- **Slingshot 11** dragonfly network @ 200Gbps
- 2.0 PiByte of Cray/HPE ClusterStor E1000 **Lustre** storage
 - Mix of NVME/HDD, ~50GB/s
- Multi Architecture Comparison System (**MACS**)
 - AMD Milan, Genoa, Bergamo, Intel Sapphire Rapids + HBM, AMD/Nvidia GPUs...
 - Limited number of nodes
- In a new, dedicated **MDC** equivalent size of 15 racks total (also hosts Isambard-AI phase 1).

NVIDIA Grace Superchip



Source: NVIDIA

- Provides an upgrade route to previous ARM-based Isambard 2 (**ThunderX2**).
- Strength in its memory bandwidth.
- Use of LPDDR5X makes power of CPU+MEM efficient.
- C2C latency
- 4x 128b SVE2 per core



Available systems within Isambard 3

- Shared infrastructure

- 2 PB ClusterStor (Lustre)
- Slingshot **FMN 2.2**
- Slingshot **SHS 11.1** (**11.0** on MACS)
- SLES15sp5
- Cray Programming Environment
- HPCM 1.11

- Comparison given various CPU and memory configurations

Processor	#	Mem [GB]	Cores	Base Clock Speed [GHz]	FP64 peak [TFLOP/s]	Default TDP [W]	Bandwidth [GB/s]
NVIDIA Grace CPU Superchip	384	240	2 x 72	3.1	7.1	1 x 500 (including memory)	1024.0
AMD EPYC 7713 (Milan)	12	256	2 x 64	2.0	4.0	2 x 225	409.6
AMD EPYC 9354 (Genoa)	2	384	2 x 32	3.25	3.3	2 x 280	921.6
AMD EPYC 9754 (Bergamo)	2	192	1 x 128	2.25	4.6	1 x 360	460.8
Intel Xeon Gold 6430 (Sapphire Rapids)	2	256	2 x 32	2.1	4.3	2 x 270	614.4
Intel Xeon CPU Max 9462 (Sapphire Rapids)	2	120	2 x 32	2.7	5.5	2 x 350	3276.8

Motivation

- **UKRI** launch regular Access calls.
 - Evaluation using codes from recent project applications.
 - Currently 20 projects from different fields of science.

- Evaluation will cover:
 - Performance **characteristics**
 - Suitability for a user software **self-service** approach
- Characteristics include:
 - Nvidia Grace Superchip
 - System design with **Slingshot 11**
- Self-service:
 - to provide the tools and information to allow users to build their own software.

Benchmarks*

Question from researchers

“How does X perform on Isambard compared to system Y”

- Synthetic
 - **STREAM** [[link](#)]
 - Arm-kernels [[link](#)]
 - CloverLeaf [[link](#)]
 - TeaLeaf [[link](#)]
 - **SNAP** [[link](#)]
 - **Neutral** [[link](#)]
 - **OSU Micro-benchmarks** [[link](#)]
- Applications
 - **CASTEP** [[link](#)]
 - CP2K [[link](#)]
 - **GROMACS** [[link](#)]
 - **NAMD** [[link](#)]
 - **OpenFOAM** [[link](#)]

Based on previous studies on Isambard 2

* N.B. focus on effortless **science!**

Method to run benchmarks

Reframe

- Previous **scripting** method required changes to support Isambard 3.
- Reframe provided method to run across **all clusters**.
- Supported **Spack** as the self-service approach.

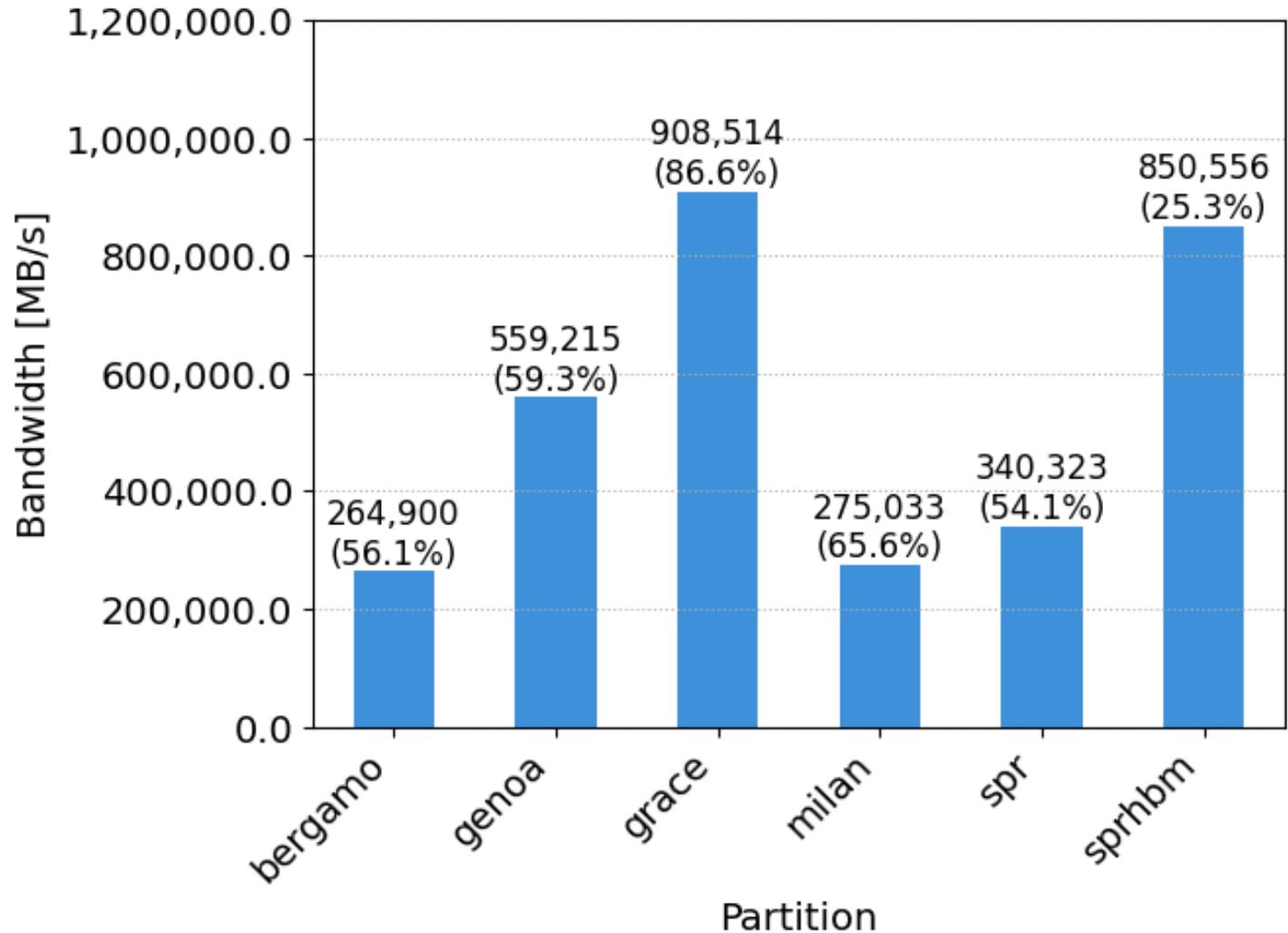
Spack

- Supported **packages** being explored.
- Provided mechanism to try **different** compilers.
- Experienced issues with compilers except with GCC.
 - Documented HPE approach results in CCE not in archspec
 - Mixing compilers seems to have issues in 0.23.1

Configuration available at: <https://github.com/isambard-sc/buildit>

Stream

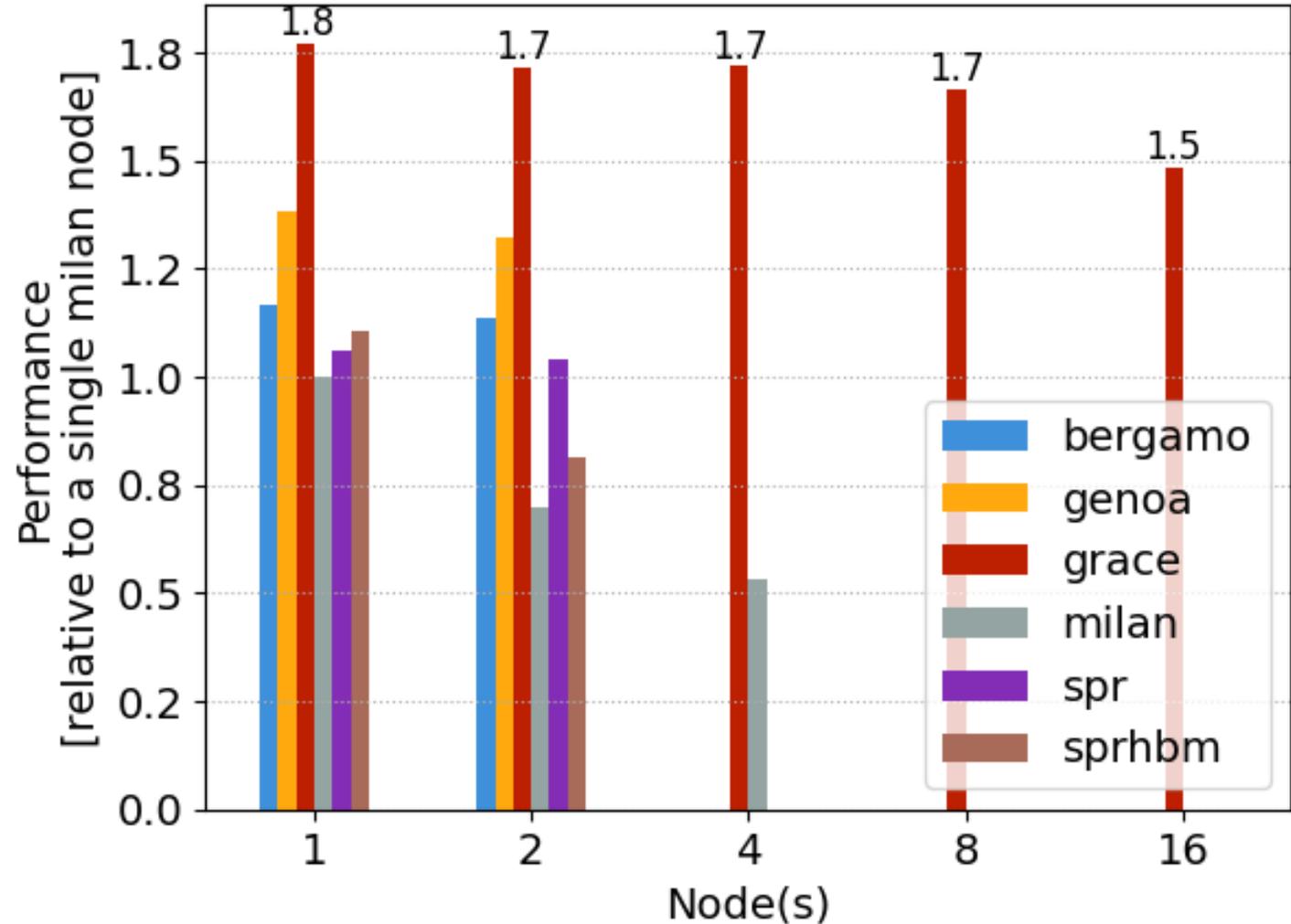
- This synthetic benchmark provides a measure of bandwidth in MB/s
- Grace and Sapphire Rapids HBM are the expected top performers with their memory design.
- Compiler choice influences results.



System: Isambard 3
Build: GCC 12.3, OpenMP,
Source: core Spack package, 5.10

SNAP – uob-hpc

- The SNAP mini-app solves the linear radiation pseudo-transport problem on a structured mesh
- Influenced by **cache behaviour**
- The weak scaling behaviour shown by Grace.
- Some show reduction in scaling even at modest core that is worth exploring further.



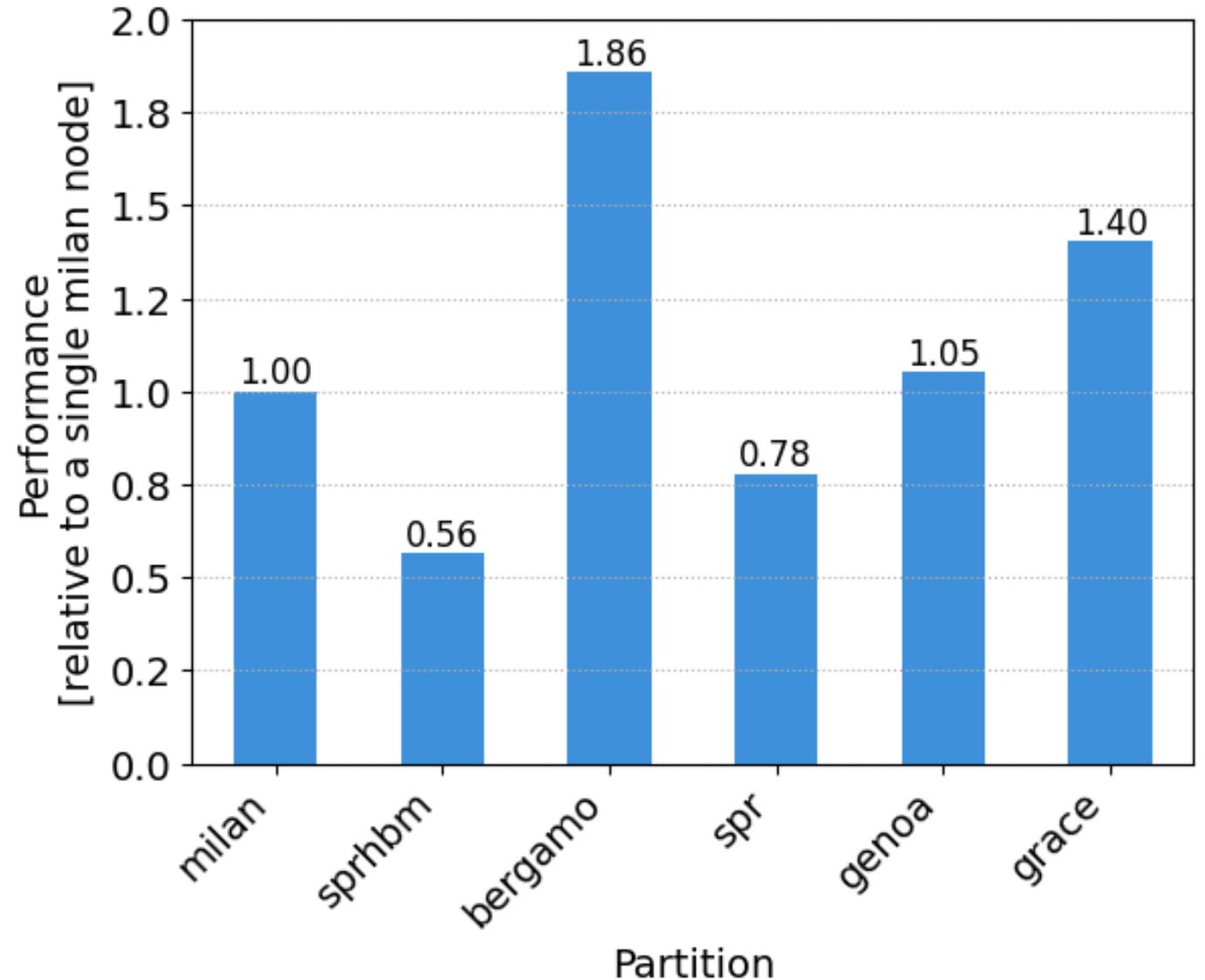
System: Isambard 3

Build: GCC 12.3, cray-mpich 8.1.30, OpenMP

Source: created Spack package, e7ab43d

Neutral

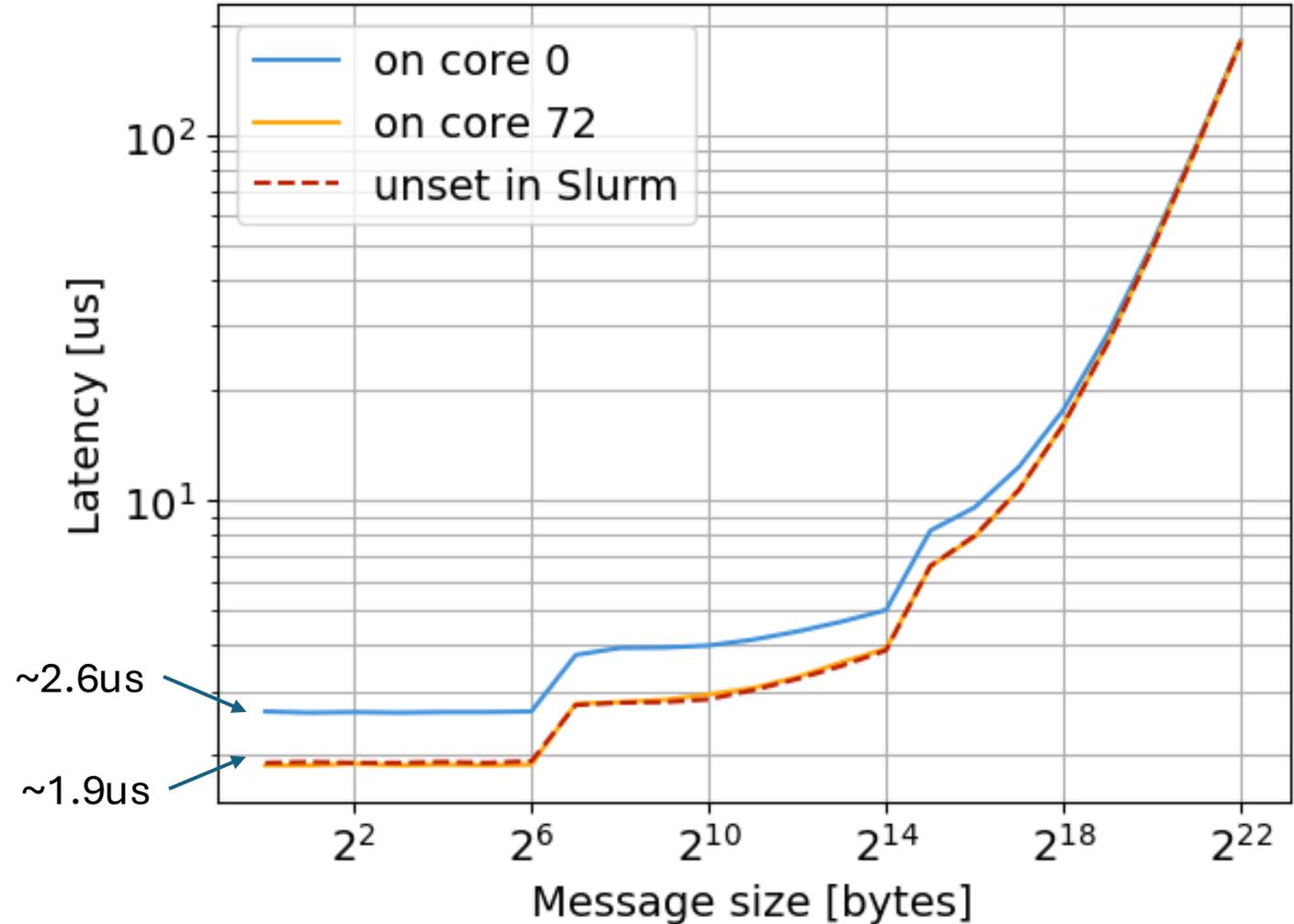
- Neutral is a Monte Carlo neutron transport mini-app (single node).
- Influenced by **cache behaviour**
- Random memory access benefits the Bergamo configuration.
- Grace performs well within the range of other configurations.



System: Isambard 3
Build: GCC 12.3, OpenMP
Source: created Spack package, d983598

OSU Microbenchmarks

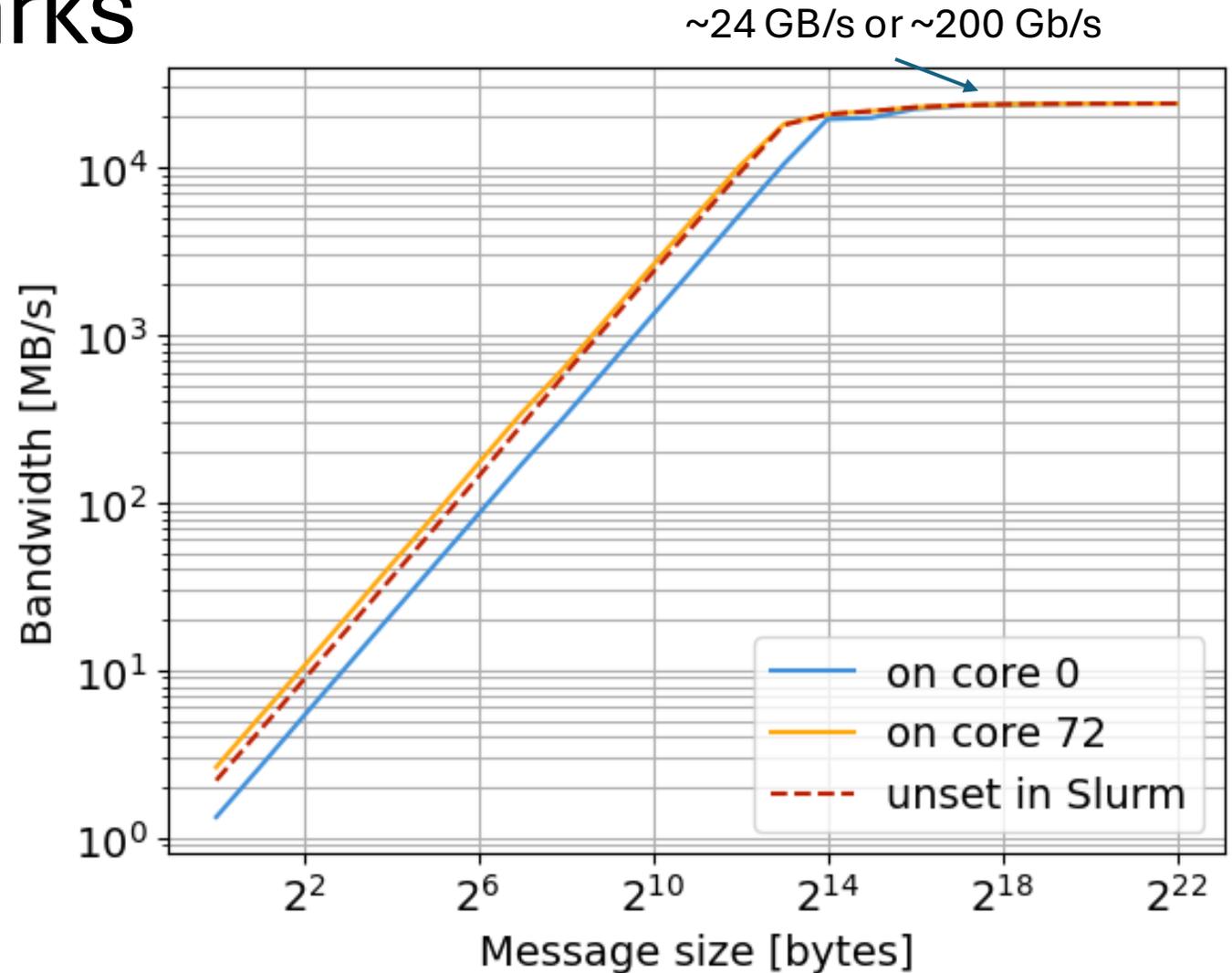
- **Latency** within Grace clearly depends on MPI task placement
- Slingshot card connected to 2nd “socket”. Default behaviour performs sensibly with MPI being placed near Slingshot.



System: Isambard 3 (2 nodes)
Build: GCC 12.3 , cray-mpich 8.1.30
Source: core Spack package, 7.5

OSU Microbenchmarks

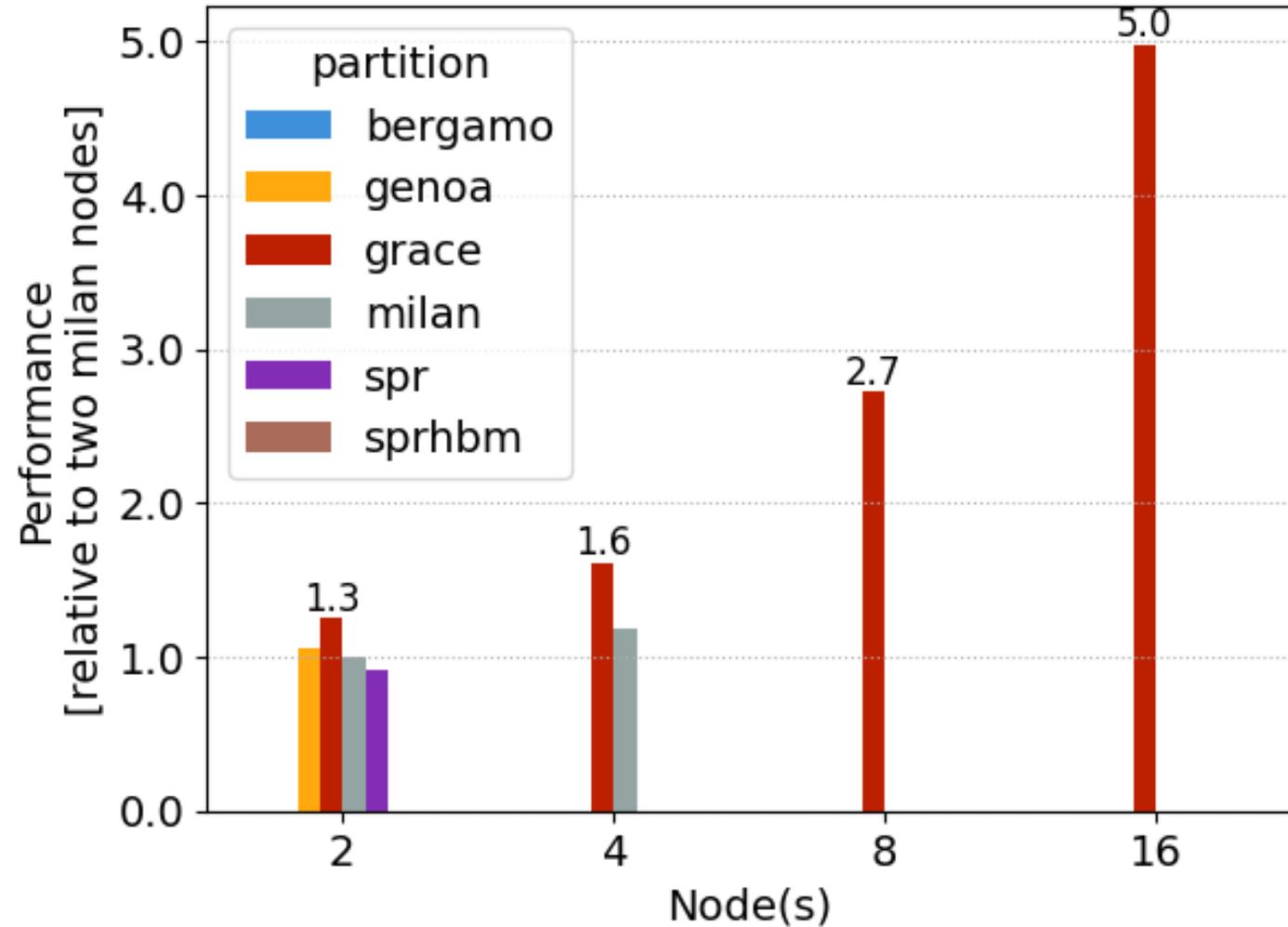
- **Bandwith** penalty when not near the Slingshot card.
- Slingshot card connected to 2nd “socket”. Default behaviour performs sensibly with MPI being placed near Slingshot.



System: Isambard 3 Grace (2 nodes)
Build: GCC 12.3, cray-mpich 8.1.30
Source: core Spack package, 7.5

CASTEP - crambin

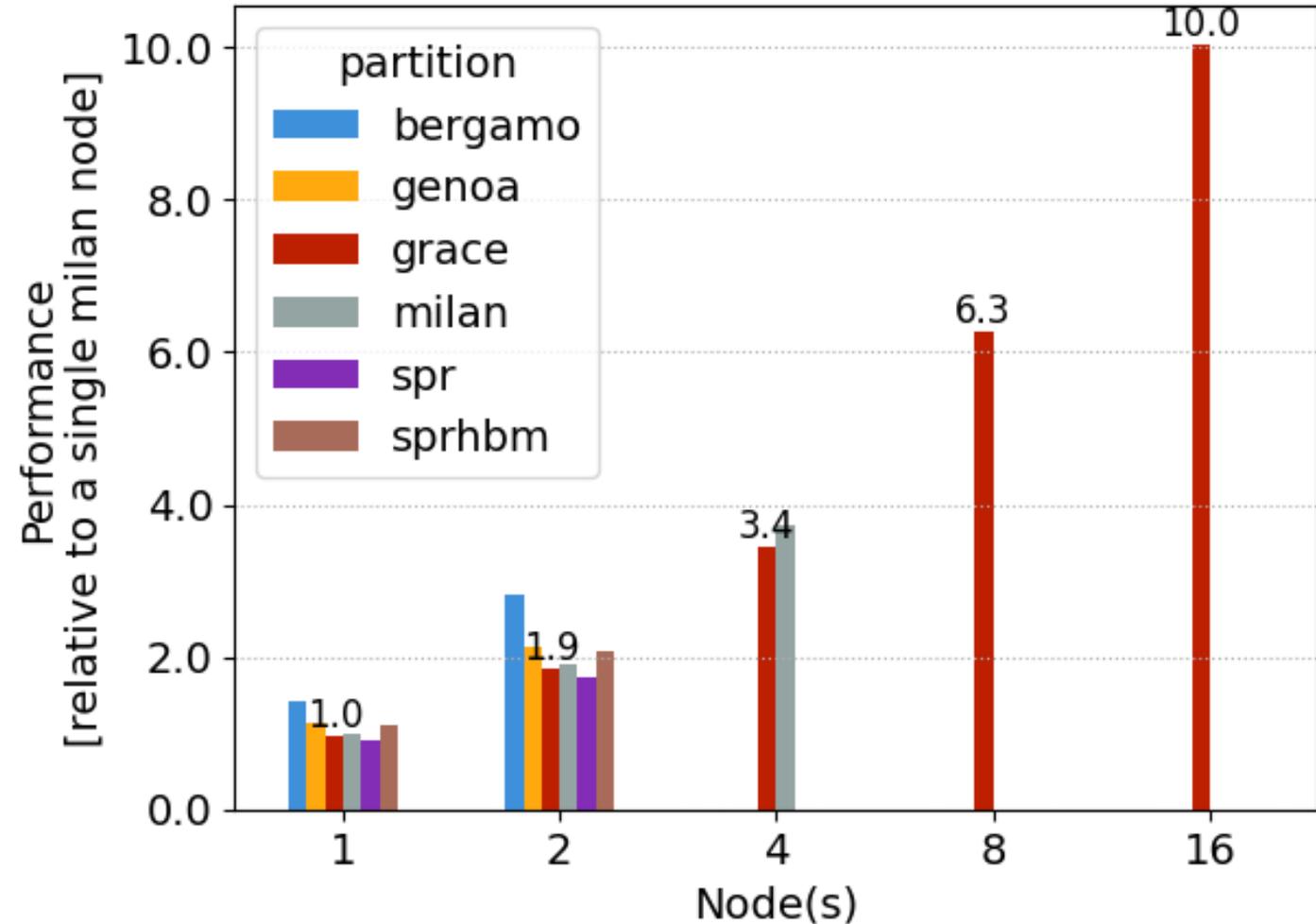
- Popular code used to calculate material properties from first principles
- Considered mainly **memory bandwidth bound**.
- Memory requirements resulted in minimum of 2 nodes
- For Bergamo and Sapphire Rapids HBM there was not enough memory available.
- Strong scaling seems good.



System: Isambard 3
Build: GCC 12.3, , cray-mpich 8.1.30
Source: created Spack package, 25.1.1

Gromacs - TestCaseB

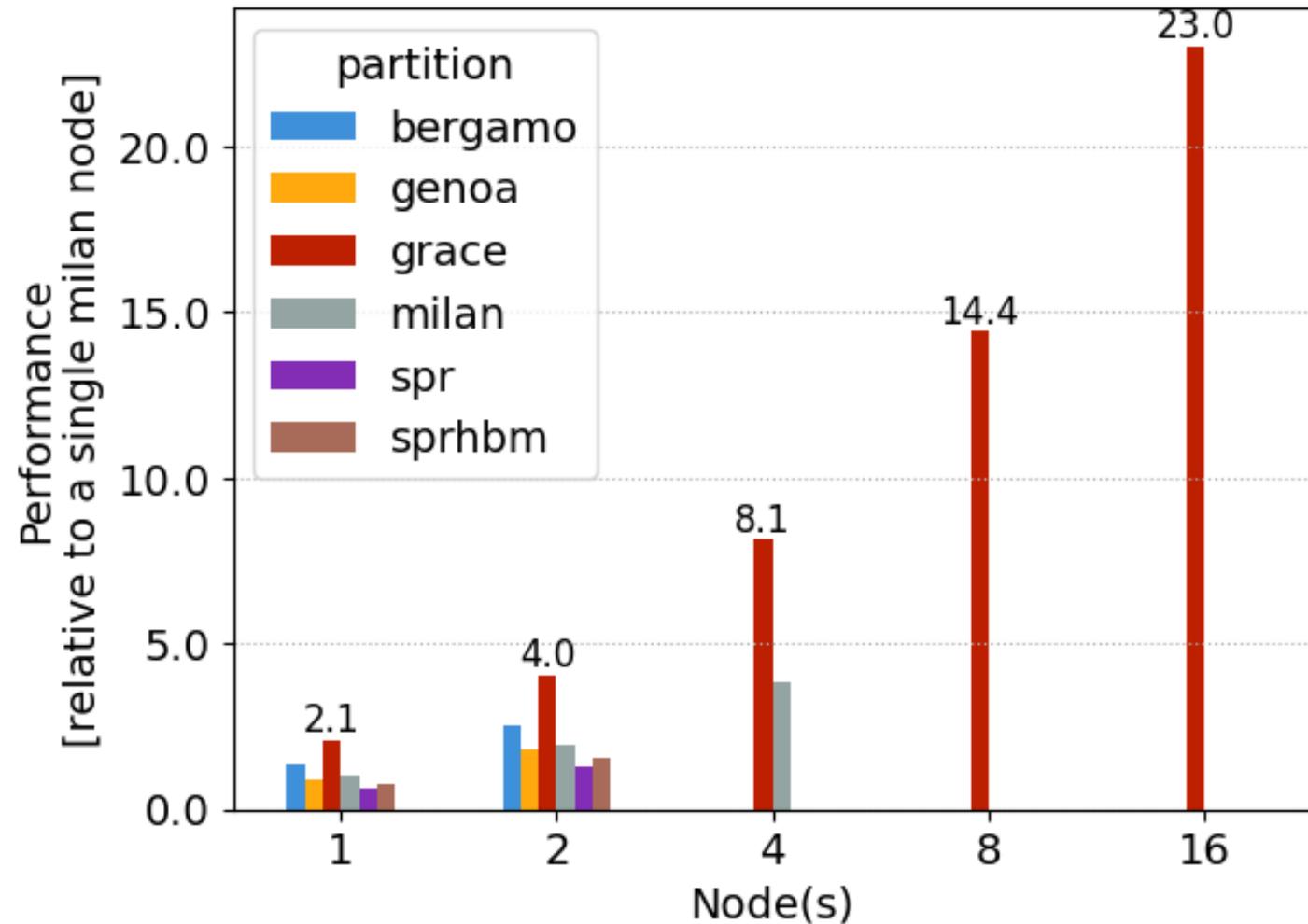
- A molecular dynamics package that solves Newton's equations of motion
- Considered **compute bound** problem.
- Grace is within the range of performances.



System: Isambard 3
Build: GCC 12.3, , cray-mpich 8.1.30
Source: core Spack package, 2024.3

NAMD - STMV

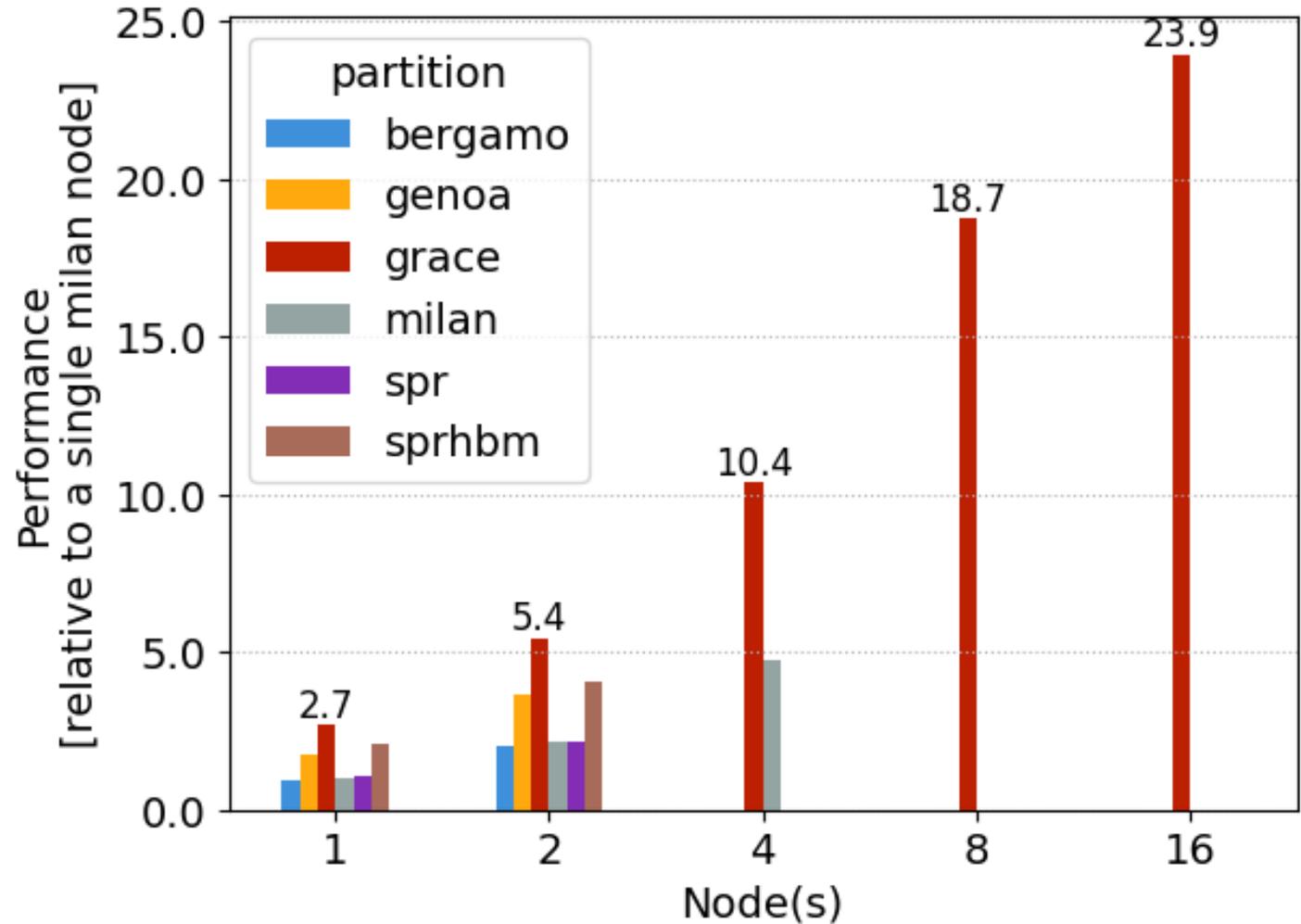
- A molecular dynamics simulation program designed to scale up to millions of atoms
- **No single bound** and Grace is suited to workload.
- Strong scaling is clearly shown without signs of reaching limit.



System: Isambard 3
Build: GCC 12.3, , cray-mpich 8.1.30
Source: created Spack package, 3.0

OpenFOAM – HPC Motorbike

- OpenFOAM is a modular C++ framework aiming to simplify writing custom computational fluid dynamics (CFD) solvers
- A **memory bandwidth bound** code will perform well on Grace and Sapphire Rapids HBM
- OpenFoam performs well in strong scaling as expected.



System: Isambard 3
Build: GCC 12.3, , cray-mpich 8.1.30
Source: core Spack package, 2312

“Supercomputer for scientific applications!”

- NVIDIA Grace Superchip
 - Performs well across **range of software**
 - **Competitive** against a range of other offerings.
 - **Memory bandwidth** clear advantage (e.g. OpenFOAM)
- System design
 - Well suited but Slingshot-11 requires careful consideration of application likely due to Grace C2C.
- Suitability of **self-service approach** with Spack
 - Worked well with GCC.
 - Other compilers need further work with packages.
 - Approach in HPE documentation to use CCE requires consideration of compiler flags from archspec.
- Early user feedback
 - Positive experience
 - Fix signed vs unsigned char

Future plans

- Welcome **further projects** onto Isambard 3!
- Try the soon to be released **Spack 1.0** where compiler configuration has major changes.
- Further investigate **Slingshot 11** behaviour.
- Continue to improve the **build configuration** for applications.



Acknowledgements

The authors acknowledge the use of resources provided by the Isambard 3 Tier-2 HPC Facility.

Isambard 3 is hosted by the University of Bristol and operated by the GW4 Alliance (<https://gw4.ac.uk>) and is funded by UK Research and Innovation; and the Engineering and Physical Sciences Research Council [EP/X039137/1].

brics-enquiries@bristol.ac.uk

<https://www.bristol.ac.uk/supercomputing/>

