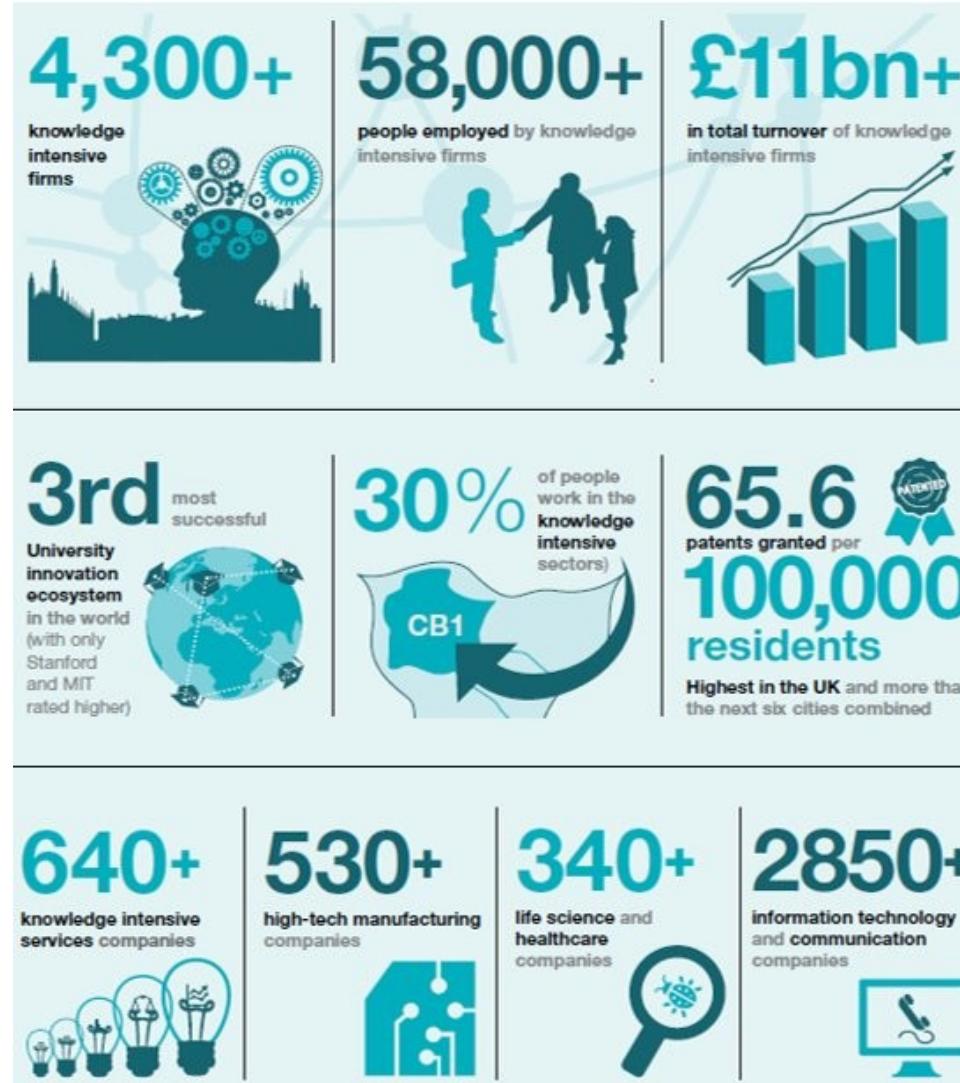


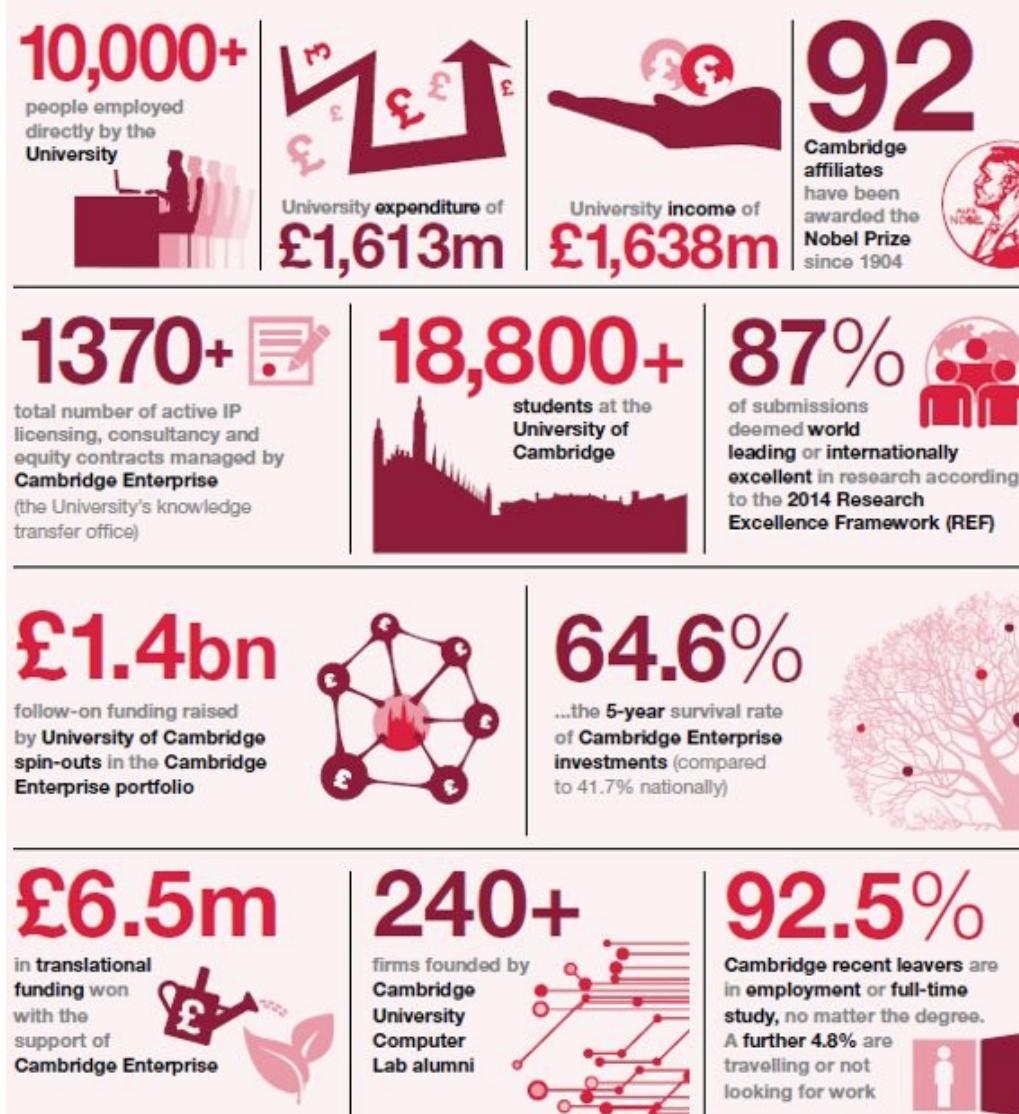
MVAPICH at the Cambridge Open Zettascale Lab

Chris Edsall
MVAPICH User Group August 2022

Cambridge Silicon Fen



University of Cambridge



Research Computing Services

Provide leadership class research computing services - hardware, software and knowledge resources, to:-

- The university of Cambridge research community
- The UK national academic research community
- UK industrial users
- Wider international organisations via consultancy and system integration services via the Dell-Cambridge HPC solution centre & Cambridge spinout “Cambridge Research Computing Ltd”

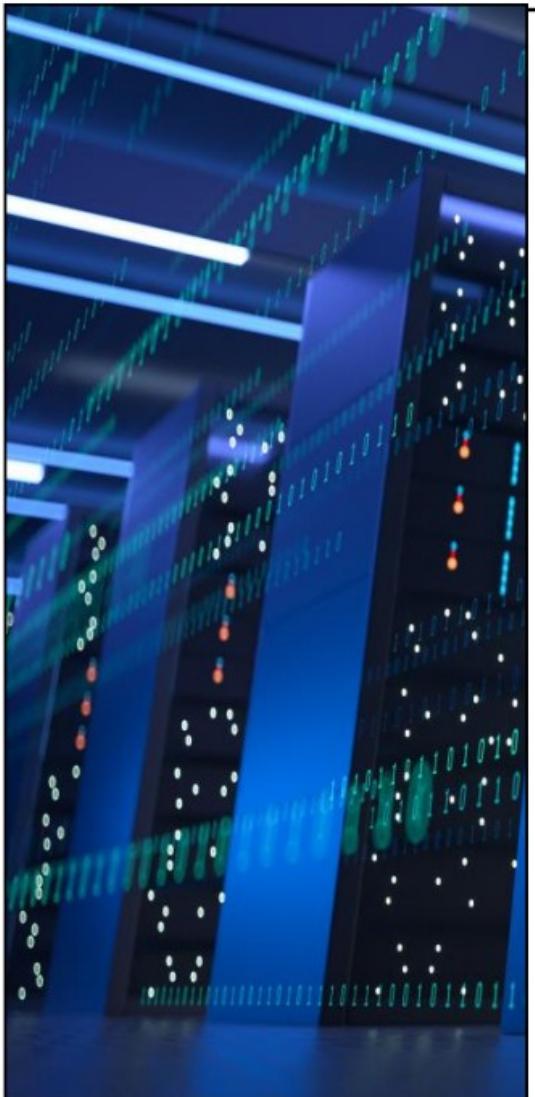
CSD3

- The Cambridge System for Data Driven Discovery
- Top100 Supercomputer
- Dell Servers
- Partitions
 - Cascade Lake
 - 762 nodes
 - HDR100
 - Ice Lake
 - 544 nodes
 - HDR200
 - "ampere"
 - 80 Nodes
 - 4x A100 80 GiB
 - Dual rail ConnectX6
 - Bluefield2 DPU





UNIVERSITY OF
CAMBRIDGE



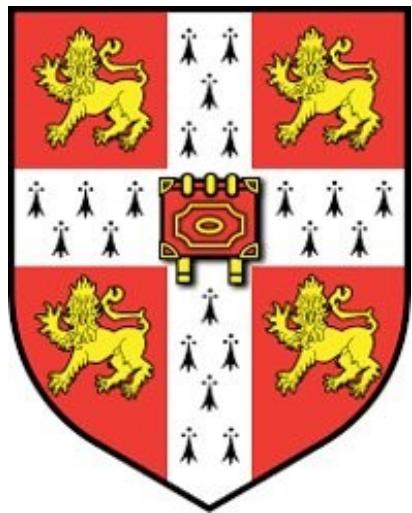
- \$5M public/industry funding - technology scanning Lab, with the goal of making exascale practical and ZettaScale possible
- Academic / industrial partnership for the testing & co-design of leading edge HPC, AI and DA solutions -20 engineers
- Pushing the boundaries of performance but more importantly making large scale HPC system more accessible thereby democratising HPC/AI and DA technologies for everyone at every-scale



Cambridge Open Zettascale Lab - Themes

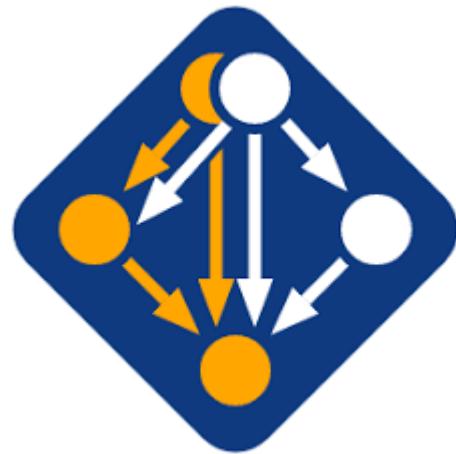
- Energy-Efficient Computing
- Extreme-Scale Visualisation
- High-Performance Interconnects for Zettascale
- oneAPI Centre of Excellence
- Scientific OpenStack: Creating a Cloud-Native, Software-Defined Supercomputer
- Zettascale Systems for Urgent Computing
- Zettascale-Class Storage Solutions

Cambridge <-> OSU Collaboration



- Long standing cooperation, all the way back to Wilkes-1
- Increasing collaboration across a range of activities

Stack
HPC



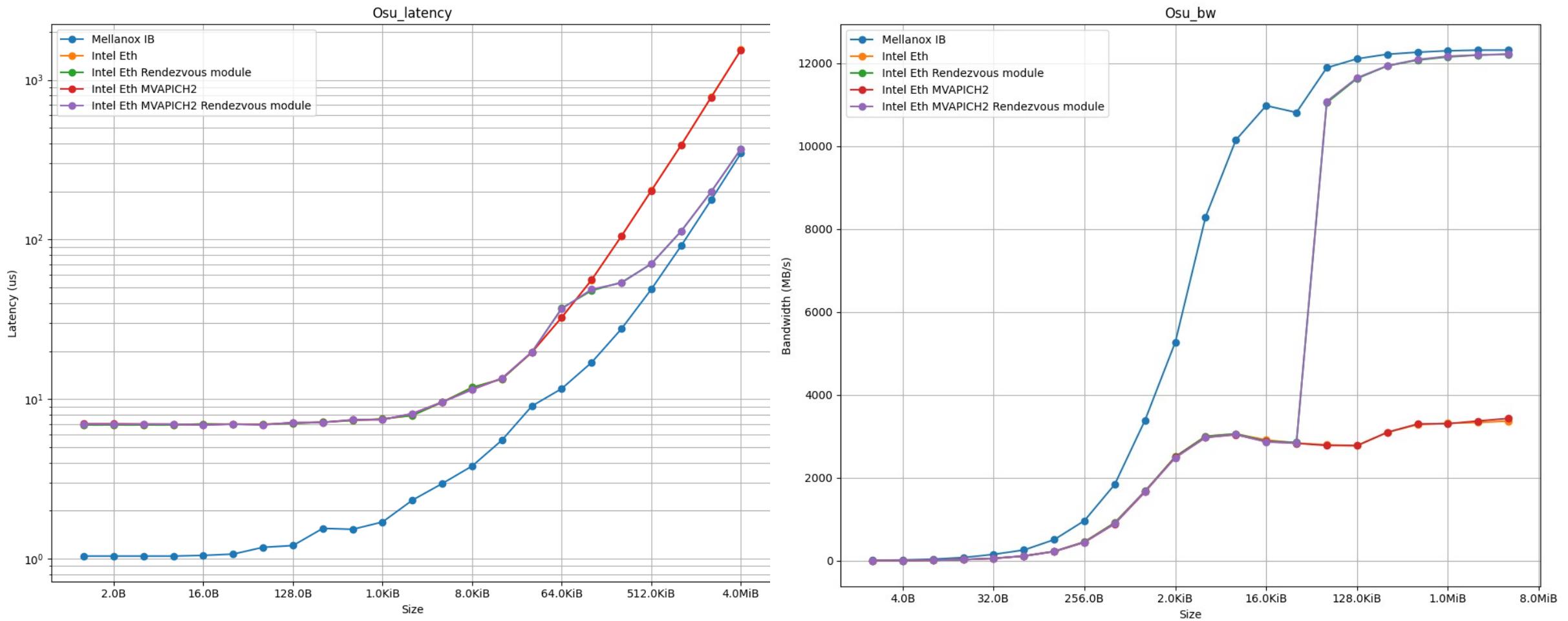
Re³Frame

- Open, reproducible
- Spack, ReFrame
- Suite of Benchmarks
 - Synthetic (OSU Micro-Benchmarks)
 - Applications
 - OpenFOAM
 - GROMACS
 - WRF 4.x
 - CP2K
- Using to, e.g.
 - Compare ethernet NICs
 - MPI libraries

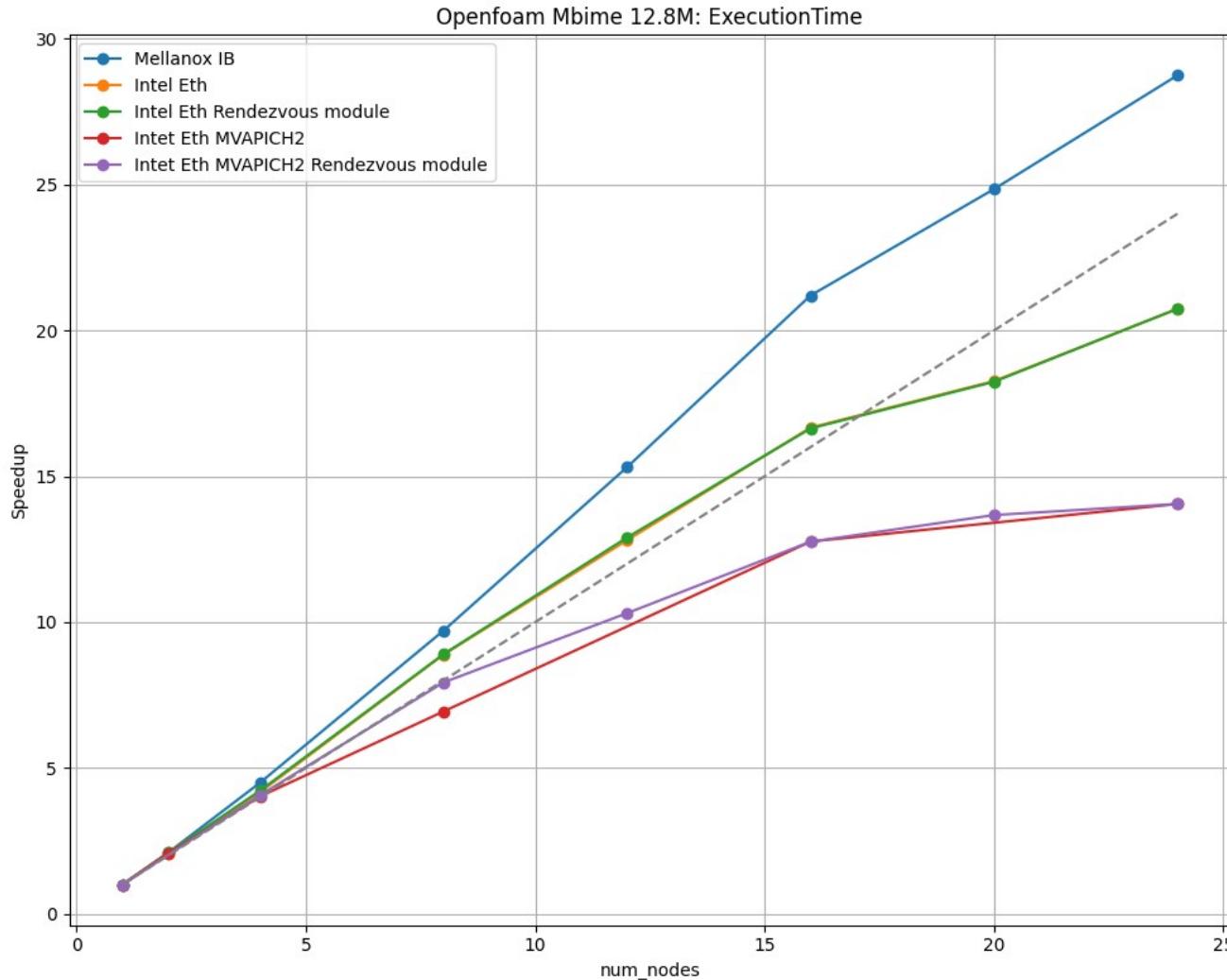
The screenshot shows a GitLab interface with the following details:

- Repository Path:** Files · rocky8 · Research Computing Services / Research Software Engineering / Ethernet Benchmarks / hpc-t...
- File List:**
 - systems: add alaska HPL runs on IB (1 year ago)
 - tools: add gromacs message pars... (1 year ago)
 - .gitignore: gitignore gromacs downlo... (1 year ago)
 - LICENSE: Create LICENSE (2 years ago)
 - README.md: Fix used reframe version i... (1 year ago)
 - environment.yml: add scaling ratio to cp2k + ... (1 year ago)
 - reframe_config.py: Configure omb test on icel... (2 months ago)
- Section:** hpc-tests
- NB This is a work in progress and READMEs etc may not be up to date.

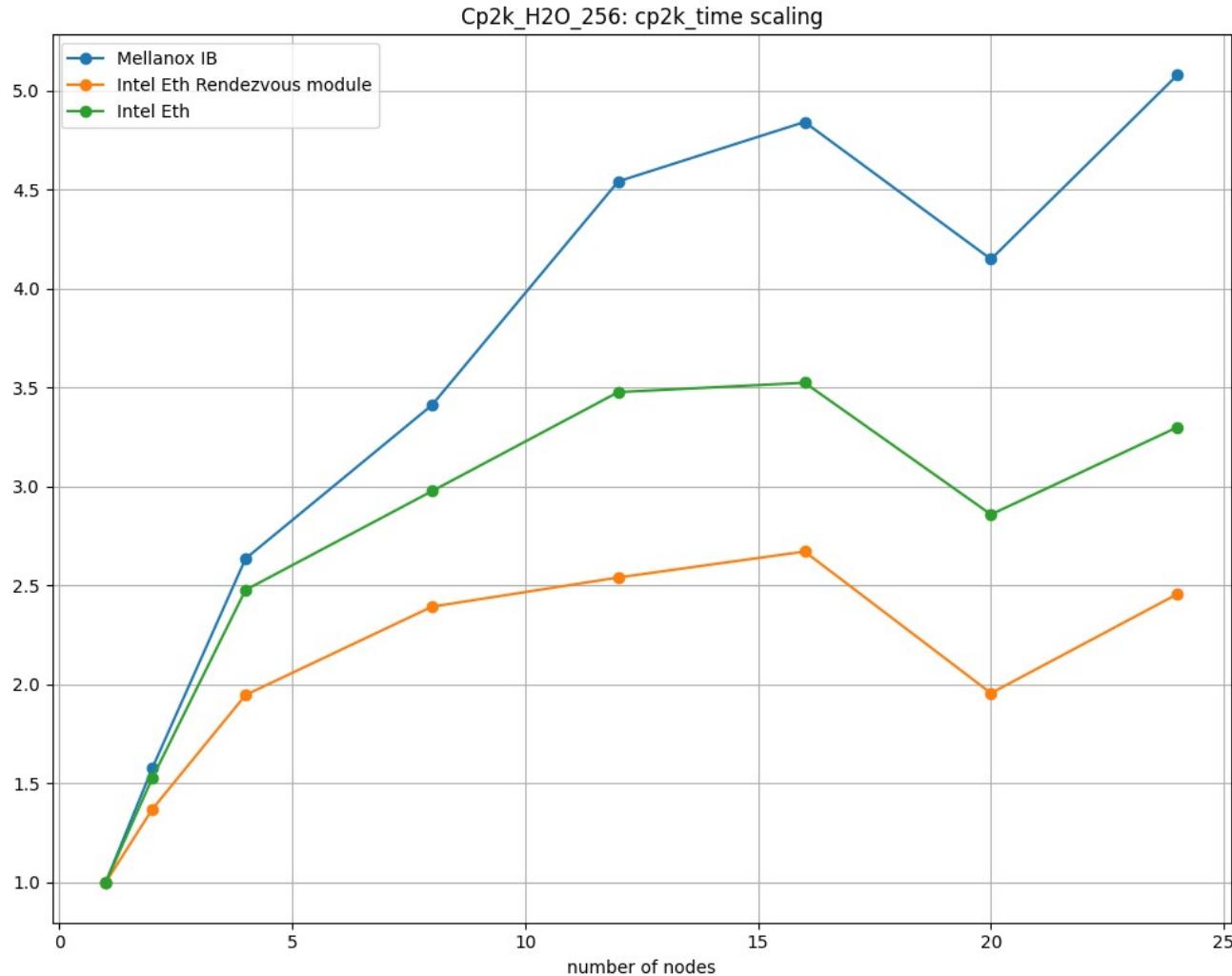
- Work by **Kacper Kornet** (Zettascale) and **Moustafa Abduljabbar** (OSU)
- Experimental partition on CSD3
 - Intel® Ethernet 800 Series (code named Columbiaville) (100Gb/s)
 - CascadeLake 56 nodes available, currently 28 each configured Intel Ethernet Fabric Suite and MOFED (because of conflict)
 - Dell ethernet switch (model Z9332, 400 Gb/s ports with 4-way split cables, 4X100 Gb/s)
- Morocco Cluster (same nodes and software as CSD3)
 - Mellanox ConnectX-6 (limited in our case 50 Gb/s because of split cables)
- Requires MVAPICH ≥ 3 for PSM3 libfabric provider
 - Two choices for PSM3:
 - Intel Ethernet Fabric Suite <- this is the one we used in this work
 - Latest Intel MPI (2021.6) also ships PSM3



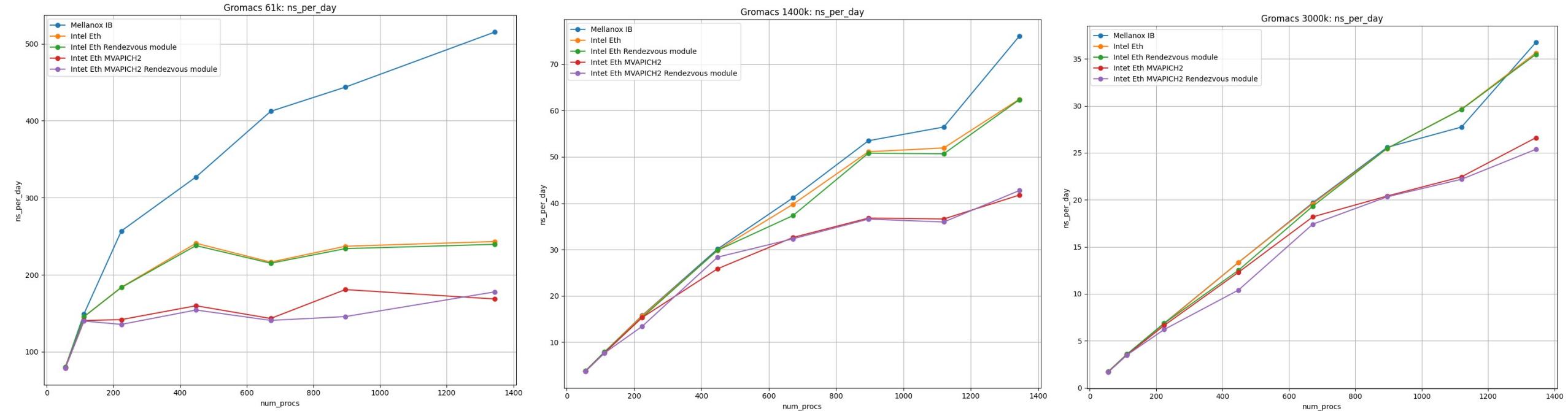
CFD - OpenFOAM



Materials CP2K



Molecular Dynamics - GROMACS



- MVAPICH2-DPU
- Accelerate MPI collective operations
- Applications?
 - CASTEP – UK based materials science code
 - MPI/OpenMP + accelerator offload
 - Scales to O(10k) processors
 - Top5 code on ARCHER2
 - Heavily FFT bound (30% wall time waiting on alltoall comms)
 - Aim: communication hiding by overlapping computation and communication
 - Preliminary work (**Arjen Tumerus**, Zettascale Lab): would require significant re-architecting of the code.





MVAPICH



dask



UNIVERSITY OF
CAMBRIDGE

Data Analytics - Methods

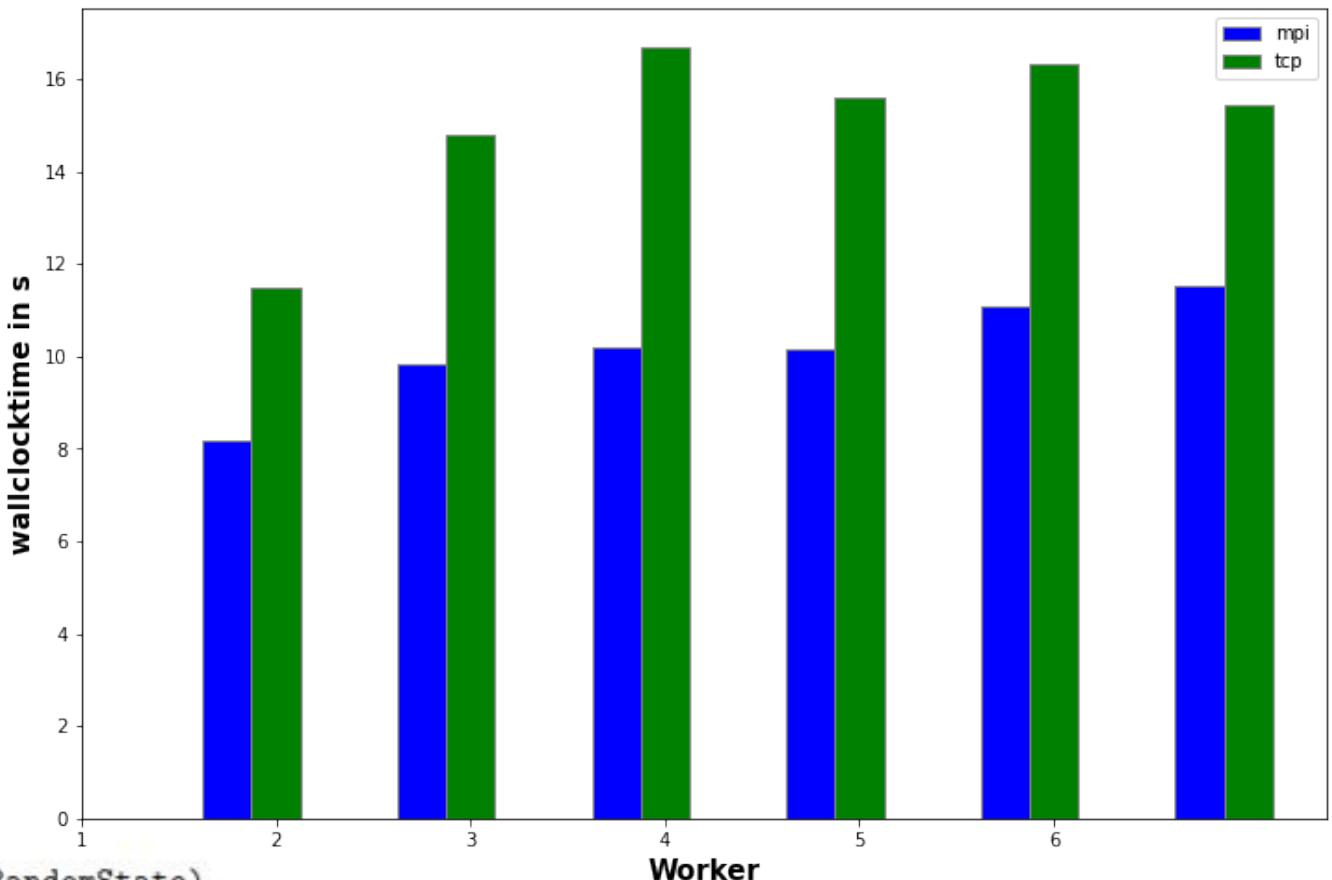
- Work by **Stefanie Reuter** (Zettascale Lab) with assistance from **Amir Shafi** (OSU)
- Installed MVAPICH2 2.3.7
 - (Next try MVAPICH2-GDR, MVAPICH2-X)
- Installed the coda environment (based on python 3.8, dask 2021.1.1, distribute 2021.1.1)
- mpi4dask-0.2
- Ran test cases
 - numpy_sum_mpi.py on CSD3 Icelake
 - cupy_sum.py on CSD3 ampere partition (Wilkes-3)

Dask – CPU Sum

```
for i in range(RUNS):
    start = time.time()
    rs = da.random.RandomState(RandomState=numpy.random.RandomState)
    a = rs.normal(100, 1, (int(3e4), int(3e4)), chunks=(int(5e3), int(5e3)))

    x = a + a.T
    xx = await client.compute(x)

    duration = time.time() - start
```



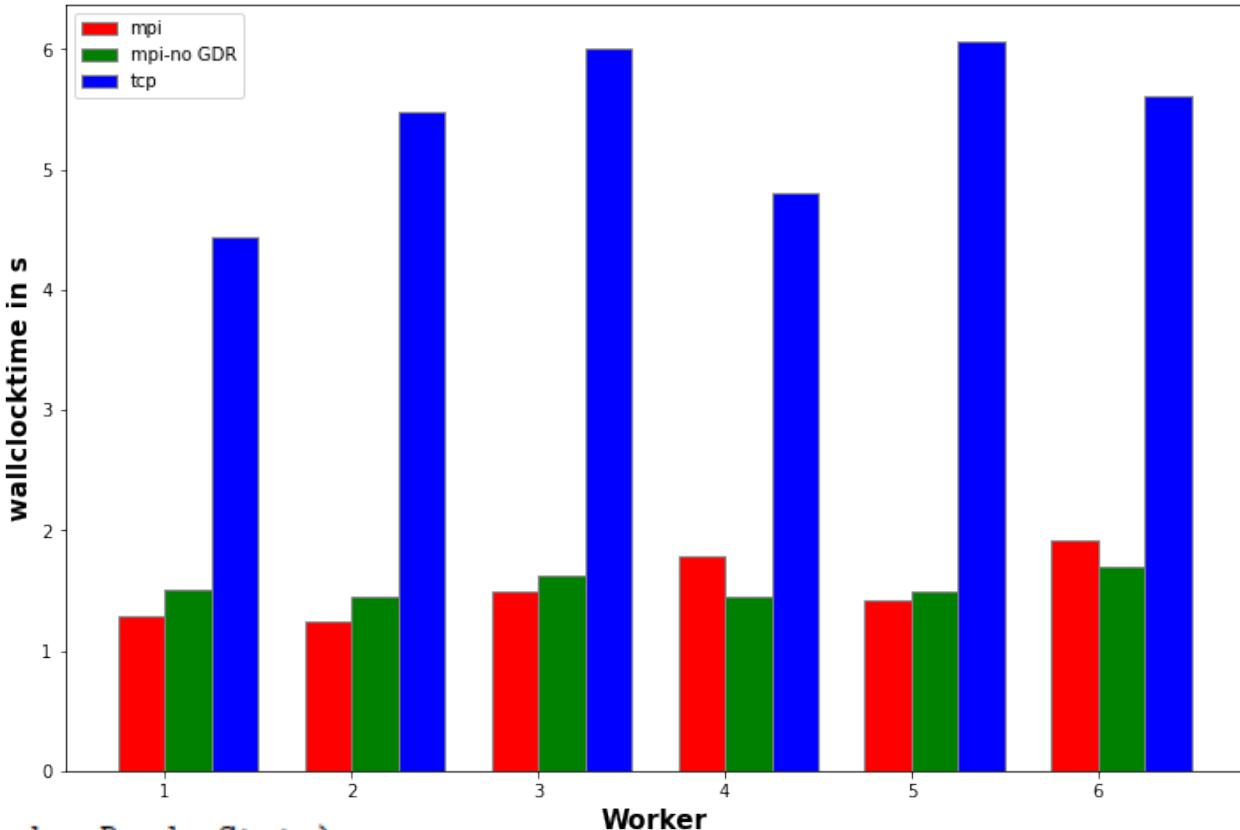
Dask – GPU Cupy Sum

```
for i in range(RUNS):
    start = time.time()

    rs = da.random.RandomState(RandomState=cupy.random.RandomState)
    a = rs.normal(100, 1, (int(20000), int(20000)), \
                 chunks=(int(4000), int(4000)))
    x = a + a.T

    xx = await client.compute(x)

    duration = time.time() - start
```



DASK - Applications

Once we have got the software stack running and validated the next step is to try it on real world applications:

- Climate
- Square Kilometer Array

Summary

- We've set up the Zettascale Lab to look at challenges around the next generation of supercomputers
- Cambridge and OSU have begun a deeper collaboration

Contacts

- <https://www.zettascale.hpc.cam.ac.uk/>
- info-cozl@hpc.cam.ac.uk
- <https://twitter.com/ZettascaleLab>
- cje57@cam.ac.uk
- <https://twitter.com/hpcchris>



About Research Media Events Careers

Get in touch



The power to
change our
world

Our mission is to make zettascale possible, exascale practical.

