# Performance Portability for Next-Generation Heterogeneous Systems

## Dr Tom Deakin

| Rank | System | Accelerator |
|------|--------|-------------|
| 1 | Frontier | ✅ |
| 2 | Supercomputer Fugaku | ❌ |
| 3 | LUMI | ✅ |
| 4 | Leonardo | ✅ |
| 5 | Summit | ✅ |
| 6 | Sierra | ✅ |
| 7 | Sunway TaihuLight | ❌ |
| 8 | Perlmutter | ✅ |
| 9 | Selene | ✅ |
| 10 | Tianhe-2A | ✅ |

**Latency** ←→ **Throughput**

| Latency | Throughput |
|---|---|
| "Complex" cores | More "simple" cores |
| Instruction Level Parallelism | Very wide SIMD |
| Deep cache hierarchy | Fast context switching |
| NUMA | Programable memory hierarchy |
| Wide SIMD | Latest memory technology |

# NVIDIA Grace-Hopper

# Apple M1

None · 332
NVIDIA GPU · 154
AMD GPU · 8
Other · 6

Data: TOP500 June 2022
Graph: doi.org/10.1109/P3HPC56579.2022.00006

5

Tension between migrating to next system (which may be GPUs), and keeping running on current system
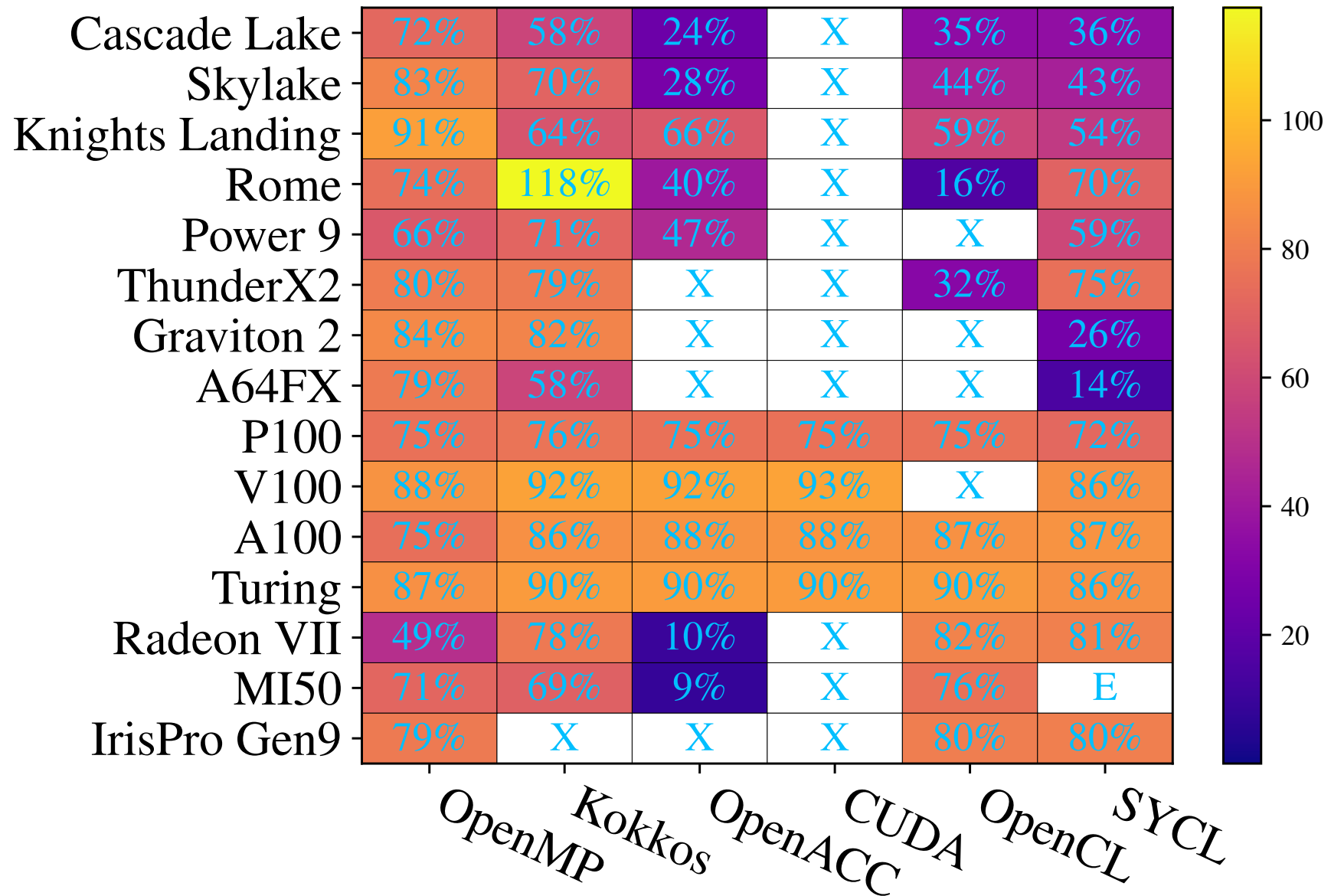
# Performance, Portability, and Productivity

"A code is performance portable if it can achieve a similar fraction of peak hardware performance on a range of different target architectures".
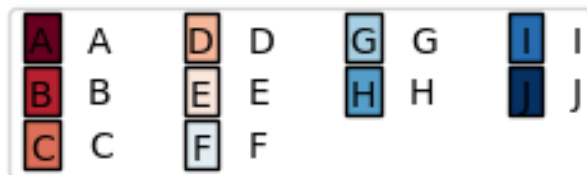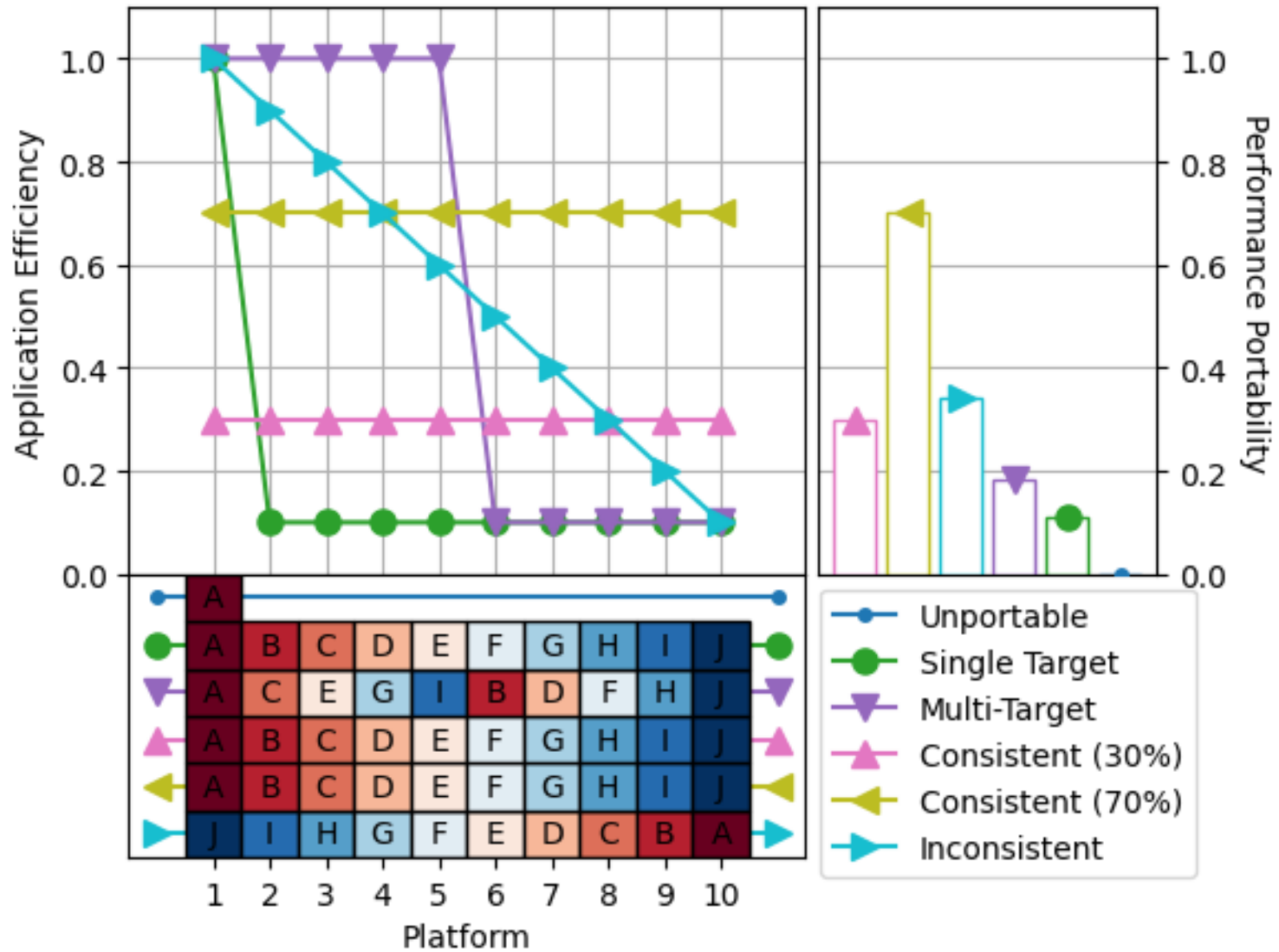
Problem
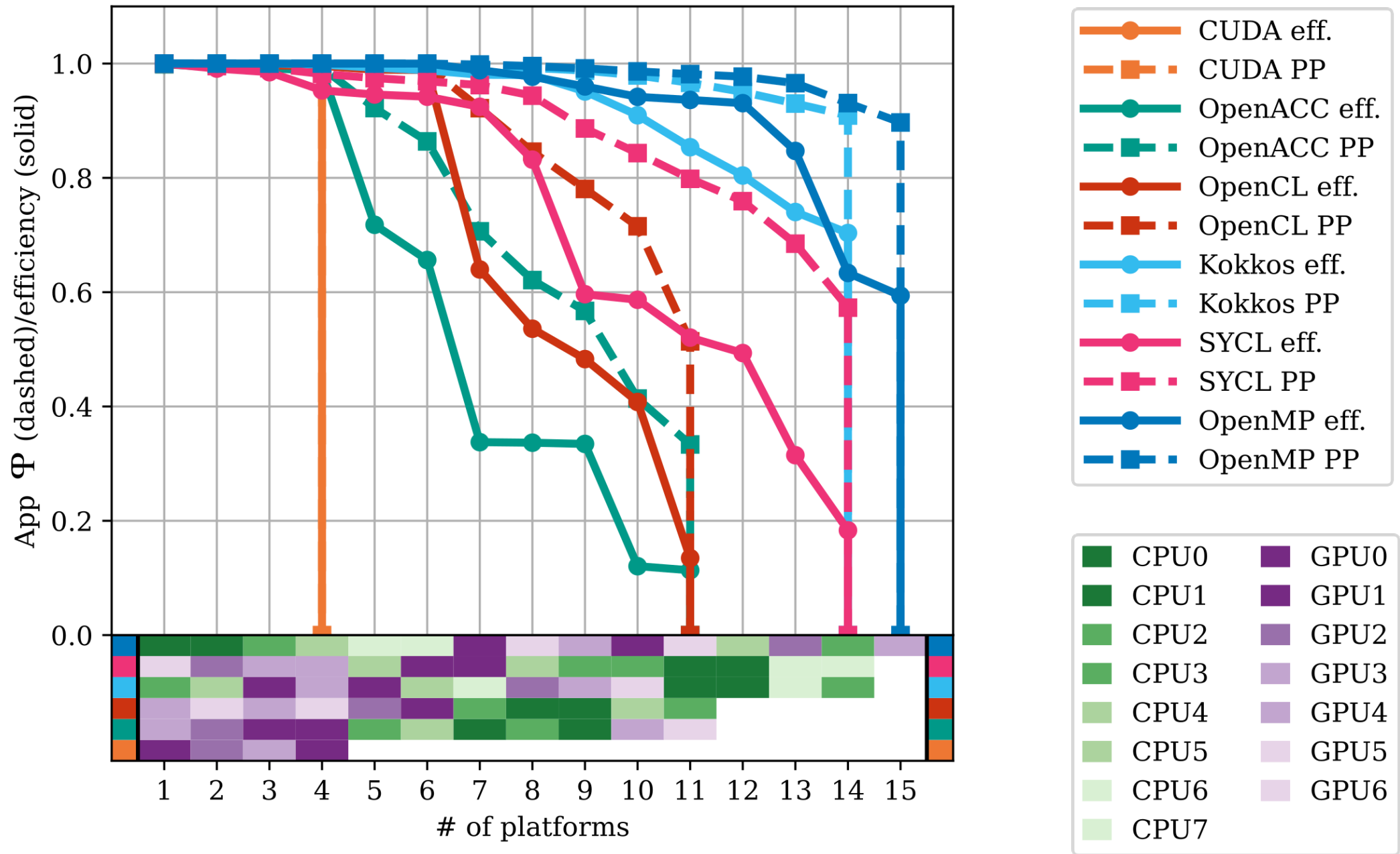
Application

Platform

Efficiency

$$\mathbb{P}(a, p, H) = \begin{cases} \dfrac{|H|}{\displaystyle\sum_{i \in H} \dfrac{1}{e_i(a, p)}} & \text{if, } \forall i \in H \\ & e_i(a, p) \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

BabelStream Triad array size=2**25

| | OpenMP | Kokkos | OpenACC | CUDA | OpenCL | SYCL |
|---|---|---|---|---|---|---|
| Cascade Lake | 72% | 58% | 24% | X | 35% | 36% |
| Skylake | 83% | 70% | 28% | X | 44% | 43% |
| Knights Landing | 91% | 64% | 66% | X | 59% | 54% |
| Rome | 74% | 118% | 40% | X | 16% | 70% |
| Power 9 | 66% | 71% | 47% | X | X | 59% |
| ThunderX2 | 80% | 79% | X | X | 32% | 75% |
| Graviton 2 | 84% | 82% | X | X | X | 26% |
| A64FX | 79% | 58% | X | X | X | 14% |
| P100 | 75% | 76% | 75% | 75% | 75% | 72% |
| V100 | 88% | 92% | 92% | 93% | X | 86% |
| A100 | 75% | 86% | 88% | 88% | 87% | 87% |
| Turing | 87% | 90% | 90% | 90% | 90% | 86% |
| Radeon VII | 49% | 78% | 10% | X | 82% | 81% |
| MI50 | 71% | 69% | 9% | X | 76% | E |
| IrisPro Gen9 | 79% | X | X | X | 80% | 80% |

From doi.org/10.1109/P3HPC51967.2020.00006

11

12

# BabelStream

**https://github.com/uob-hpc/babelstream**

Device discovery and control

Data location and movement in discrete memory spaces

Expressing concurrent and parallel work

Field Summary

From doi.org/10.1109/P3HPC54578.2021.00007

19

# x86 CPU

## Architectural efficiency

20

# NVIDIA GPUs

The University of Bristol is an Intel oneAPI Center of Excellence helped support this work.

Architectural efficiency

| | DoConcurrent | OpenMPTarget | OpenMPTargetLoop | OpenACC | OpenACCArray | CUDA | CUDAKernel |
|---|---|---|---|---|---|---|---|
| **A100 40GB Cray** Copy | | | 87 | 87 | 87 | | |
| Mul | | | 86 | 86 | 86 | | |
| Add | | | 89 | 89 | 89 | | |
| Triad | | | 89 | 89 | 89 | | |
| Dot | | | 86 | 86 | 87 | | |
| **A100 40GB NVHPC** Copy | | | | | | | |
| Mul | | | | | | | |
| Add | | | | | | | |
| Triad | | | | | | | |
| Dot | | | | | | 93 | |

# Specialisation?

Can always construct the PP=1, CC=0 by combining the best codes for each platform into an application

Rising PP results from performance increasing in one more more platforms. Broad or narrowly-focused optimizations cause this.

Everyone wants to be here: single source, best performance everywhere! But not realistic.

Optimization

Falling CC indicates that platform-specific code is being added, or common code is being removed. This is commonly found as codes are specialized.

Specialization

Abstraction

Removing specialization or adding common code increases convergence; this is typical of introducing more and higher-level abstractions.

Regression

Falling PP is rarely intentional. New features in applications may cause performance to drop in one or more platforms.

Being on the PP = 0 axis is anomalous, since at least one platform is failing

Code Convergence (1- Code Divergence)

# Which performance portable programming model should I use?

Use open standard parallel programming models

Express all concurrent work asynchronously

Build in tuning parameters

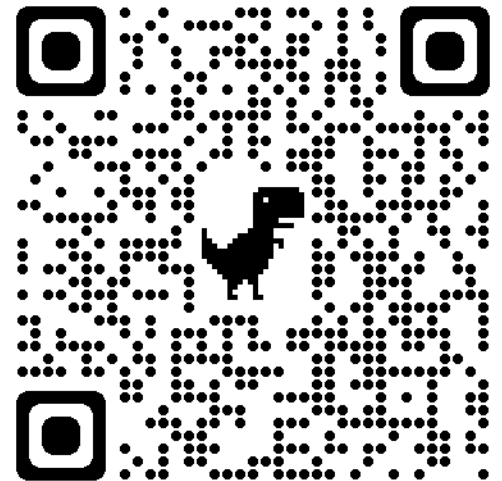Test all compilers & runtimes, on all systems

Tell your vendor

University of
BRISTOL

https://hpc.tomdeakin.com

@tjdeakin

tom.deakin@bristol.ac.uk