

How to Measure, Visualise and Interpret Performance Portability



Dr Tom Deakin

Senior Research Associate
University of Bristol

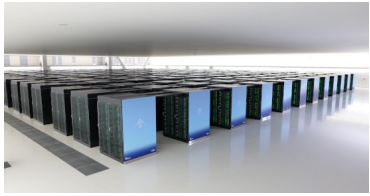
tom.deakin@bristol.ac.uk



@tjdeakin

<https://hpc.tomdeakin.com>

Processor diversity at (pre-)Exascale



At RIKEN: Fujitsu A64fx CPUs.



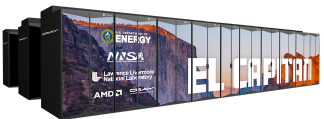
At NERSC: AMD EPYC Milan CPUs and NVIDIA A100 GPUs.



At ORNL: AMD EPYC custom CPUs and Radeon Instinct GPUs (4 per node).



At ALCF: Intel Xeon Sapphire Rapids CPUs and Xe Ponte Vecchio GPUs (6 per node).



At LLNL: AMD EPYC Genoa CPUs and Radeon Instinct GPUs (4 per node).

Processor diversity in the Cloud

- AWS: Amazon Graviton 2 CPUs, Intel Xeon, AMD EPYC.
- Oracle Cloud: Ampere Altra CPUs, Intel Xeon, AMD EPYC.
- Plus GPU offerings.

<https://www.arm.com/solutions/infrastructure/cloud-computing>

Recent architectural trends



CPU

- Many “complex” cores (80 per socket).
- Wide vectors (AVX-512, SVE 128-2048 bits).
- Chiplet manufacturing.
- Deep cache hierarchy. NUMA.
- Mainly DRAM, but...
 - Intel Xeon Phi MCDRM
 - Fujitsu A64FX HBM2
 - NVIDIA Grace LPDDR5x

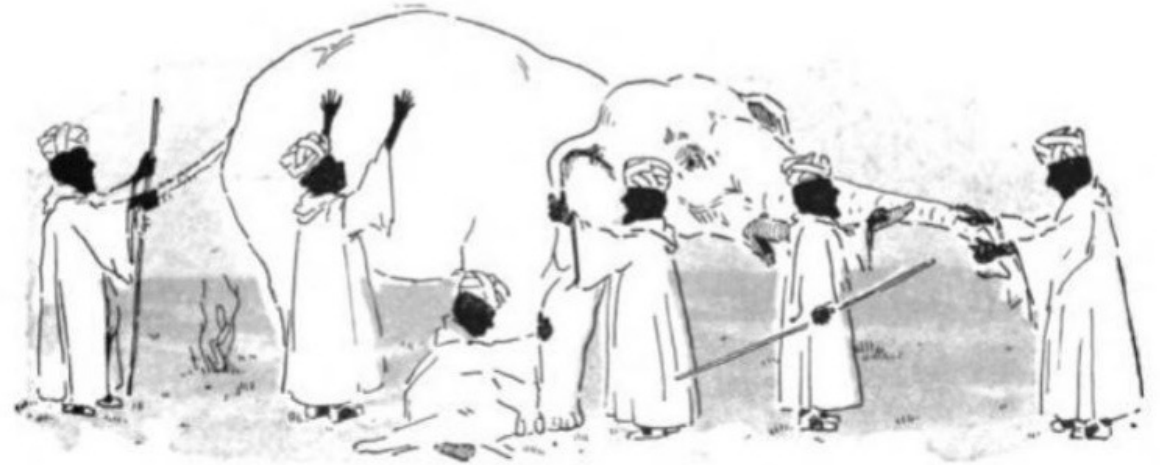
GPU

- Lots of “lightweight” cores.
- Very wide vector units (warp).
- Cores becoming more complex:
 - Specialised in-core accelerators.
- Interconnects (NVLink).
- Latest (specialised) memory technology:
 - GDDR
 - HBM
- Deepening memory hierarchy:
 - Caches, scratchpad (shared), ...

What is performance portability?

“A code is performance portable if it can achieve a similar fraction of peak hardware performance on a range of different target architectures”

- Needs to be a good fraction of best achievable (i.e. hand optimised).
- Range of architectures depends on your goal, but important to allow for future developments.



Enabling performance portability

Open standard parallel programming models



Open-source programming abstractions



Your favourite
DSL and its
compiler

BabelStream

- Benchmarks achievable (main) memory bandwidth.
- Based on McCalpin STREAM, except:
 - Arrays allocated on the heap.
 - Problem size known only at runtime.
- Written in many programming models.
- Constructed of simple vector operations, e.g.:
 - Triad: $a[i] = b[i] + \text{scalar} * c[i]$

<https://github.com/UoB-HPC/BabelStream>



Measuring efficiency

- Compare relative application performance on different processors.
- Processors have different performance characteristics.
- Architectural efficiency:
 - Percentage of peak hardware performance.
 - E.g. achieved GB/s or FLOP/s vs theoretical tech sheet.
- Application efficiency:
 - Performance relative to specialised, hand-tuned, unportable, “best” version.
 - I.e. vs “World record”.

BabelStream heatmaps

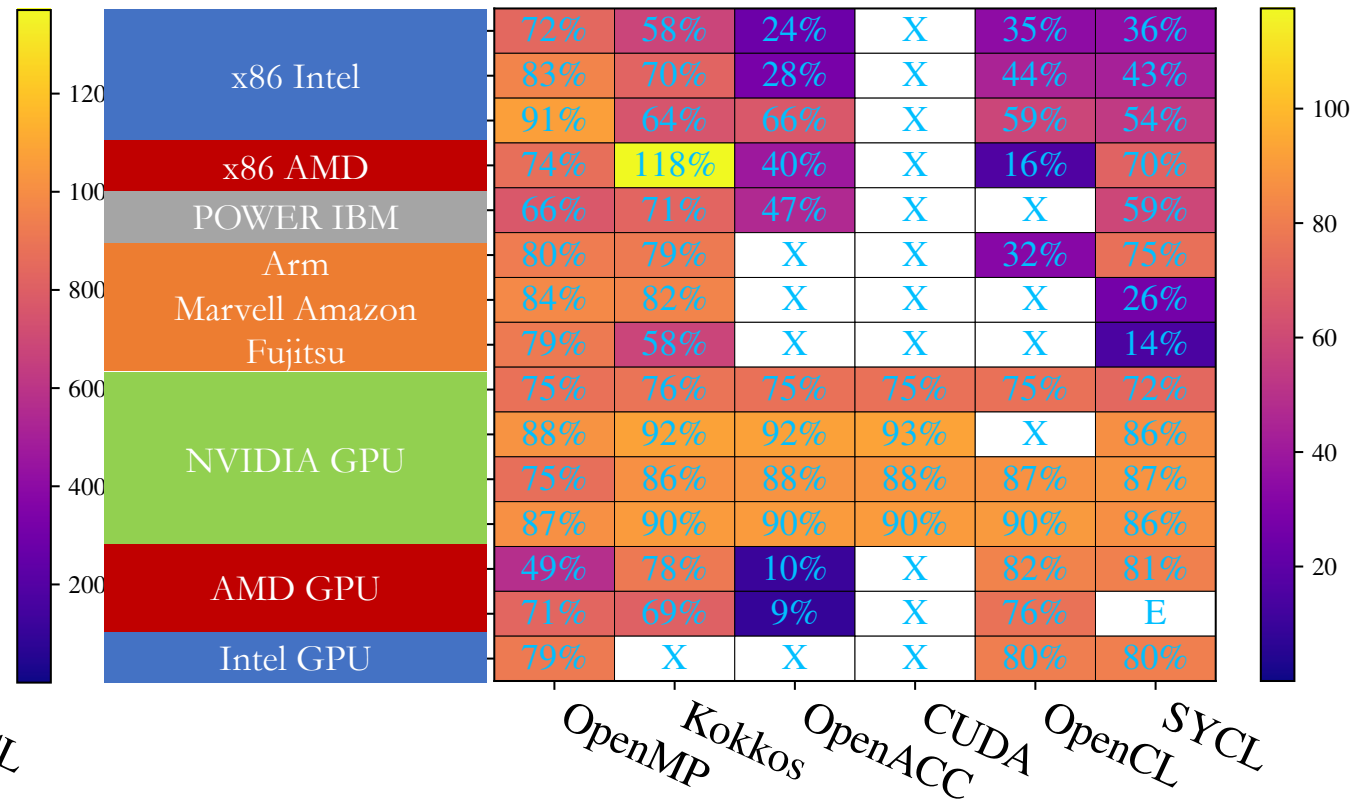
Peak performance

BabelStream Triad array size=2**25

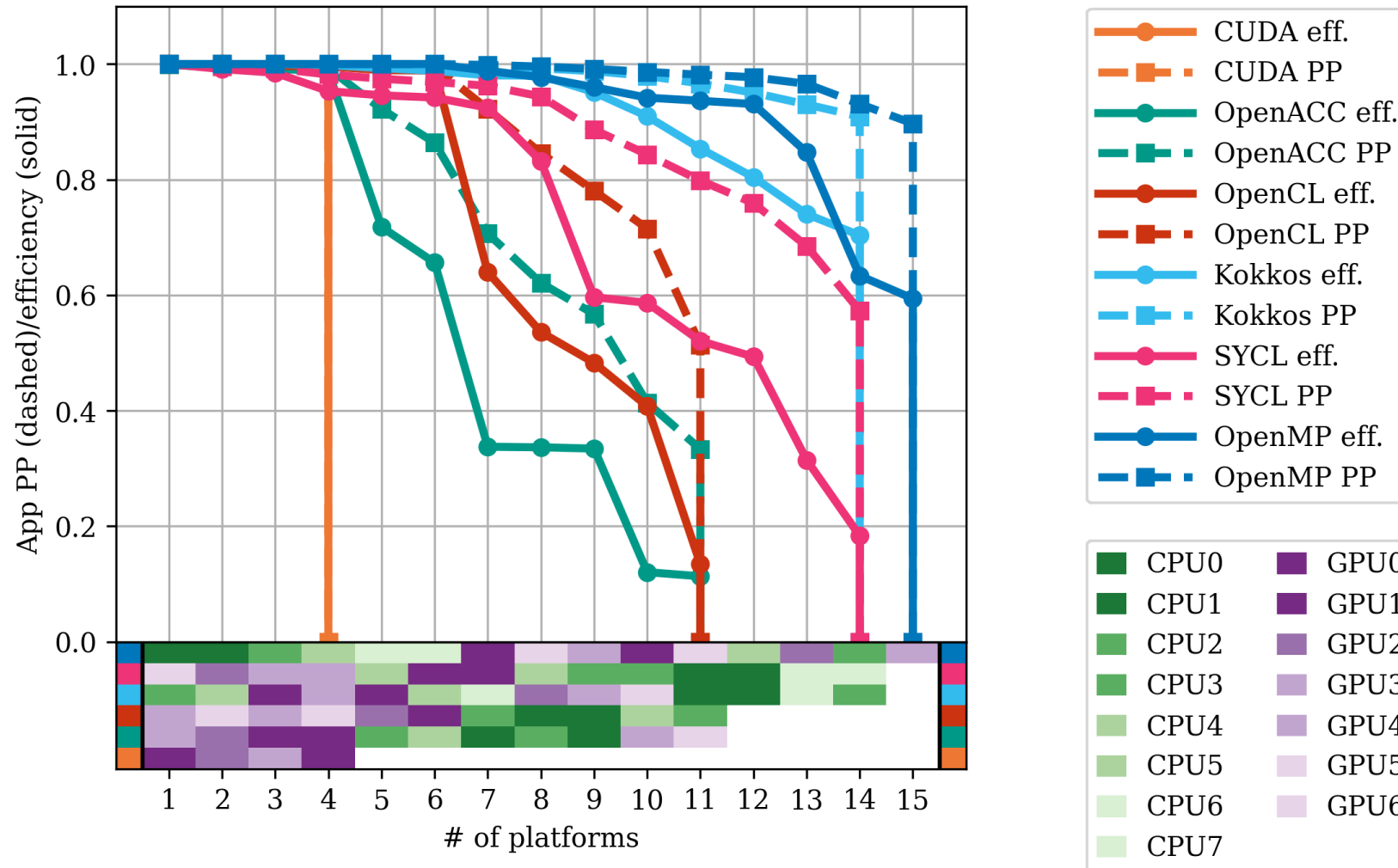
Cascade Lake	203	163	68.5	X	98.0	100
Skylake	211	180	70.6	X	113	110
Knights Landing	448	315	322	X	287	263
Rome	305	481	162	X	64.9	287
Power 9	224	241	158	X	X	200
ThunderX2	229	226	X	X	93.5	217
Graviton 2	173	169	X	X	X	54.3
A64FX	804	595	X	X	X	147
P100	552	559	551	551	552	526
V100	788	831	830	837	X	774
A100	1161	1343	1361	1370	1356	1358
Turing	533	555	555	556	554	530
Radeon VII	488	781	99.0	X	821	808
MI50	729	708	88.3	X	778	E
IrisPro Gen9	26.8	X	X	X	27.3	27.4
	OpenMP	Kokkos	OpenACC	CUDA	OpenCL	SYCL

Architectural efficiency

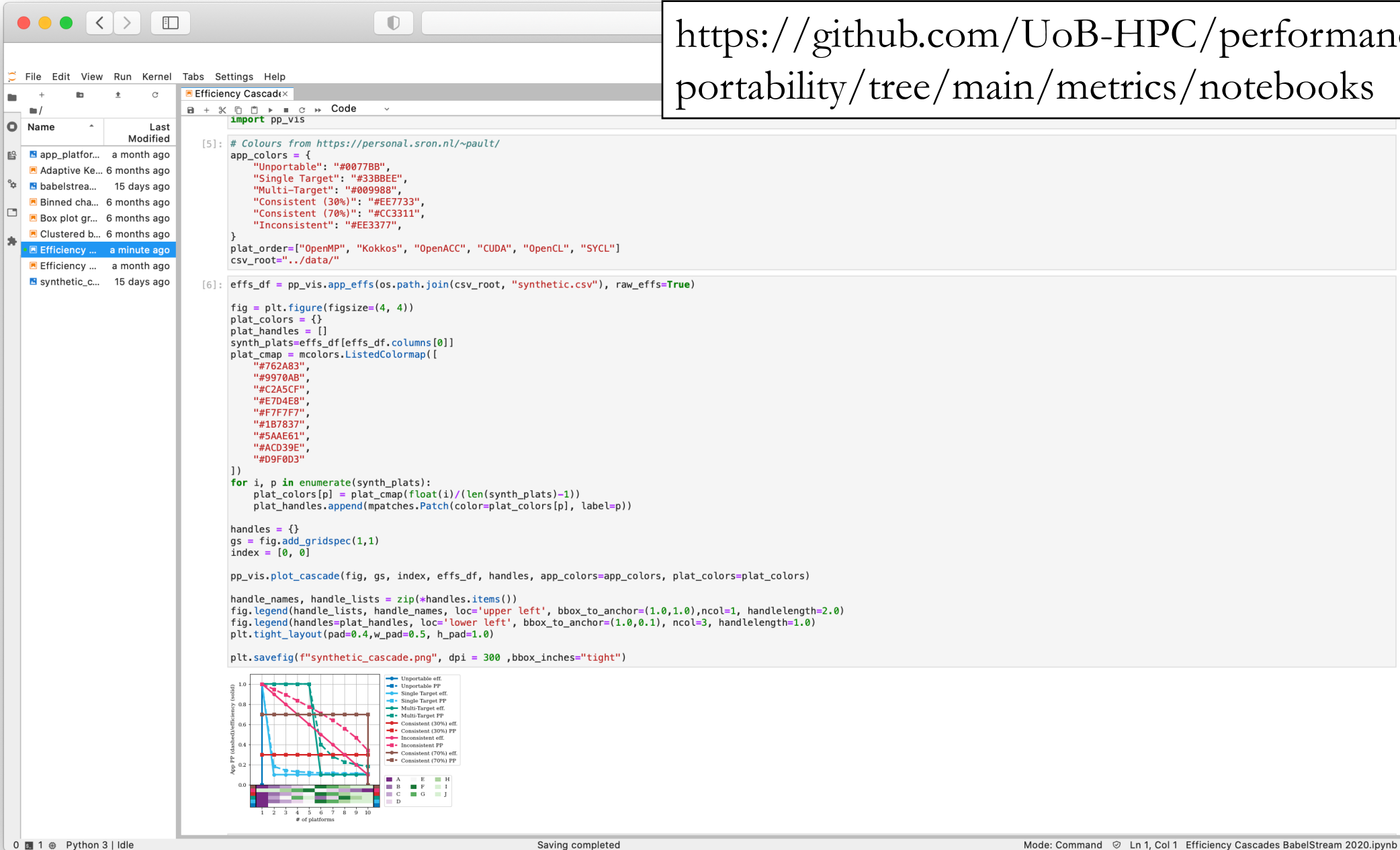
BabelStream Triad array size=2**25



BabelStream Cascade plot

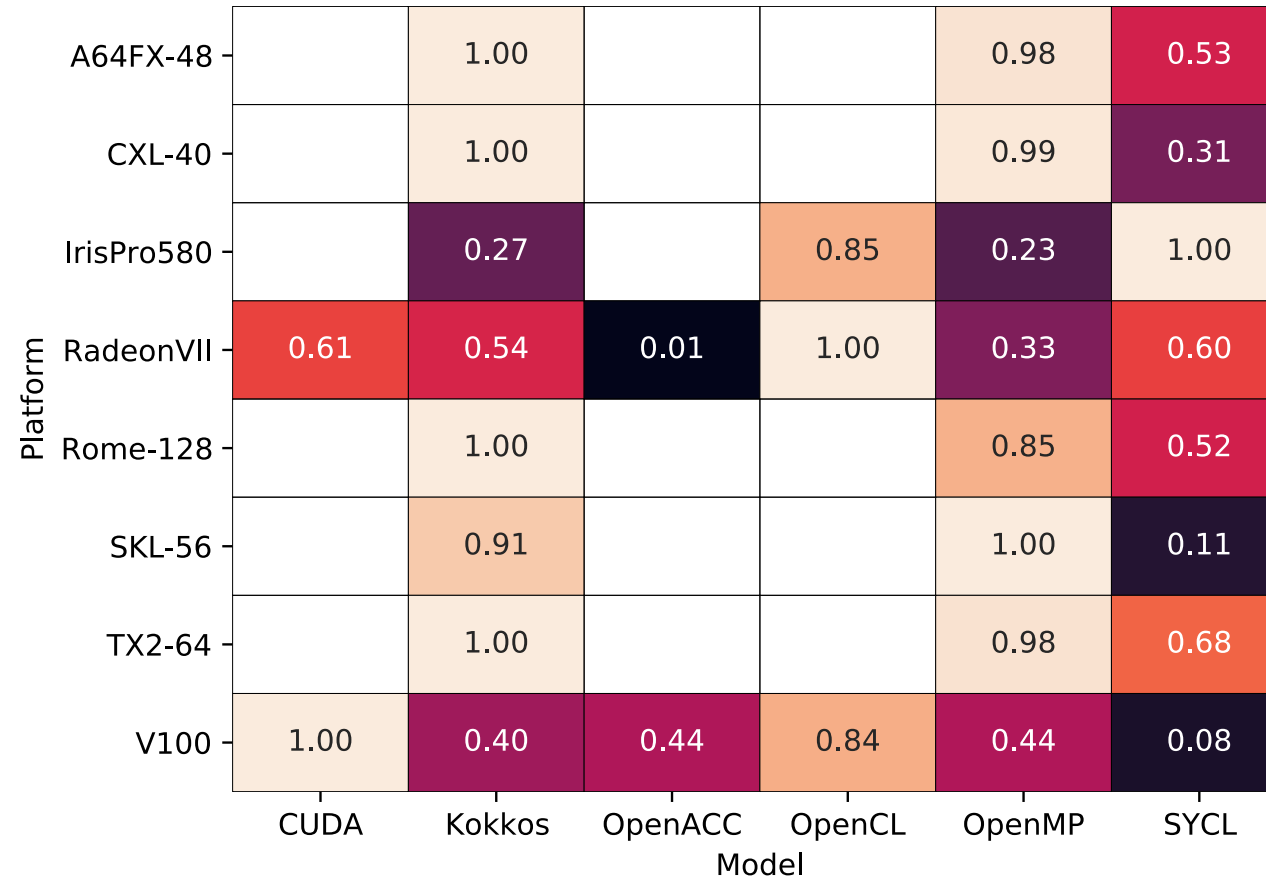


<https://github.com/UoB-HPC/performance-portability/tree/main/metrics/notebooks>



Compute bound: miniBUDE

>56% of peak compute performance on CXL and Rome



Measuring Productivity

- “Ideal” application has one version that is Performant, Portable and Productive.
- Significant specialisation for Performance and/or Portability can impact Productivity.
- Intel Code Base Investigator measures code divergence.
 - Specialisation using C pre-processor.

<https://github.com/intel/code-base-investigator>

PP-CC plane

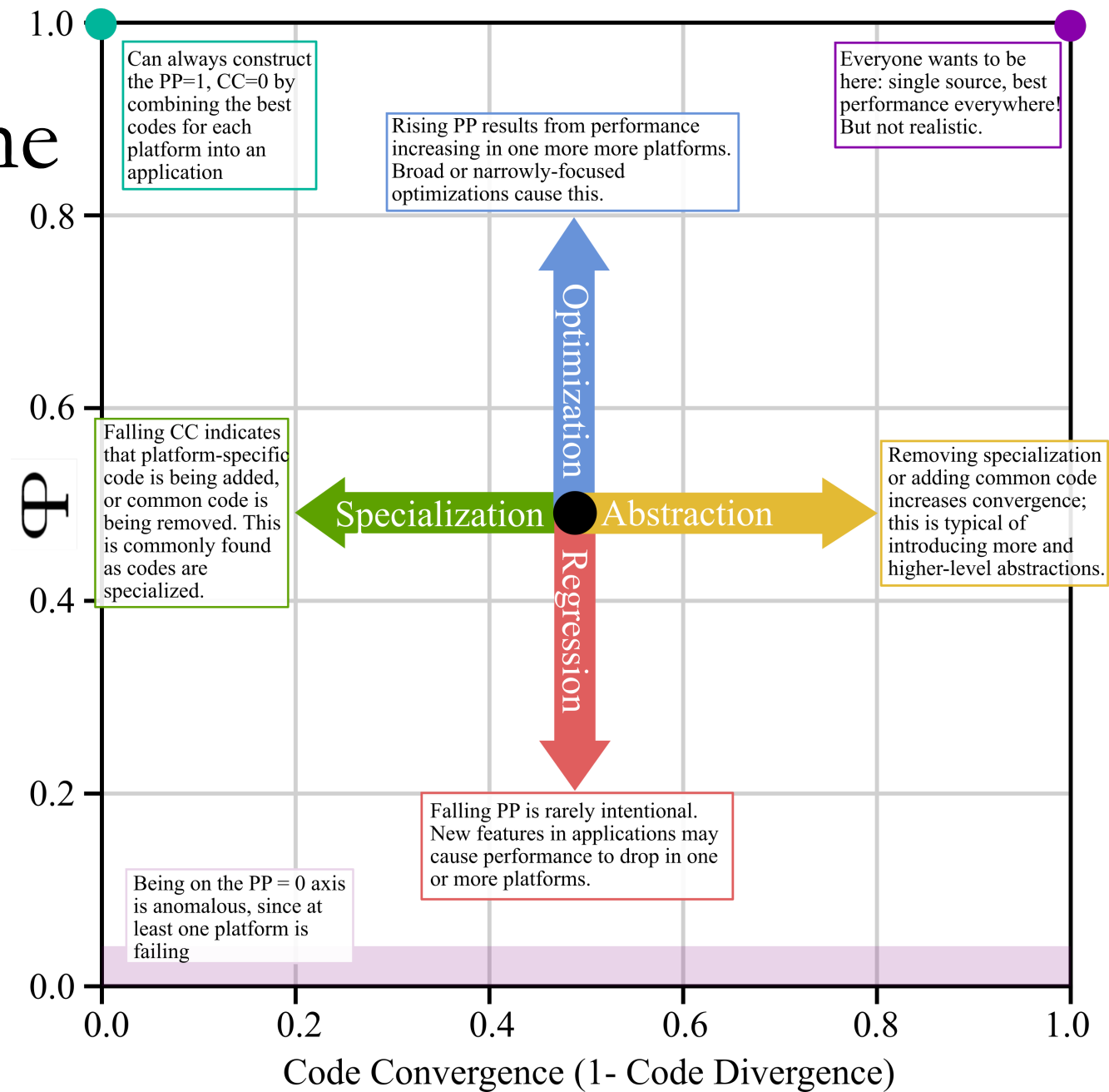


Figure by J. Sewall and J. Pennycook from upcoming article from Sewall, Pennycook, Jacobsen, Deakin and McIntosh-Smith

Top Tips for Performance Portability

- Measure and track performance and portability.
- Use Open Standards.
- Expose maximum concurrency/parallelism at all levels.
- Optimise for the diverse landscape - avoid over-optimising for one platform.
- Build in paths for specialisation and auto-tuning.

<http://uob-hpc.github.io/publications/>