How to Measure, Visualise and Interpret Performance Portability



Dr Tom Deakin

Senior Research Associate University of Bristol tom.deakin@bristol.ac.uk



The Next Platform, Jan 13th 2020: "HPC in 2020: compute engine diversity gets real" https://www.nextplatform.com/2020/01/ 13/hpc-in-2020-compute-engine-diversitygets-real/

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)

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

All images copyright of respective owners.

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] + 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".

https://doi.org/10.1109/P3HPC51967.2020.00006

BabelStream heatmaps

Peak performance



Architectural efficiency



https://doi.org/10.1016/j.future.2017.08.007

PP metric

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

1

|TT|

```
from statistics import harmonic_mean
def pp(n):
    if 0 in n:
        return 0
    return harmonic_mean(n)
```

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

https://doi.org/10.1109/P3HPC51967.2020.00007

Cascade plots



💭 File Edit View Run Kernel Tabs Settings Help

https://github.com/UoB-HPC/performanceportability/tree/main/metrics/notebooks

1 m.	+ 83	± C	■ Efficiency Cascade×	portability/tiee/main/metrics/hotebooks
_	— /		■ + % □ □ ▶ ■ C → Code ∨	
0	Name ^	Last	import pp_vis	
		Modified	[5]: # Colours from https://personal.srop.pl/~pault/	
6	🖪 app_platfor	. a month ago	app_colors = {	
	Adaptive Ke.	6 months ago	"Unportable": "#0077BB",	
°0	N habelstrea	15 days ago	"Single Target": "#33BBEE",	
		C months ago	"Multi-Target": "#009988", "Consistent (2004)", "#EET732"	
	Binned cha	Innee cra o months ago Cunsistent (30%)": "#tt//33", "Consistent (70%)": "#tt//33",		
	Box plot gr 6 months ago "Inconsistent": "#EE3377",			
1.	Clustered b	. 6 months ago	}	
	Efficiency	a minute ago	<pre>plat_order=["OpenMP", "Kokkos", "OpenACC", "CUDA", "OpenCL", "SYCL"]</pre>	
	Efficiency	a month ago	csv_root="/data/"	
	synthetic c	. 15 davs ago	[6]: effs df = np vis app effs(os path ioin(csv root "synthetic csv") raw effs	
		<pre>plat_colors = {} plat_colors = {} plat_colors = {} plat_colors = {} plat_colors = {} plat_cmap = mcolors.ListedColormap(["#752A83", "#9970A8",</pre>		
	<pre>"#CASSCF", "#F7F7F", "#F7F7F", "#SAE61", "#ACD39E", "#ACD39E", "#OPF0D3"]) for i, p in enumerate(synth_plats): plat_colors[p] = plat_cmap(float(i)/(len(synth_plats)-1)) plat_handles.append(mpatches.Patch(color=plat_colors[p], label=p)) handles = {} gs = fig.add_gridspec(1,1) index = [0, 0]</pre>			
			pp_vis.plot_cascade(fig, gs, index, effs_df, handles, app_colors=app_colors, plat_colors=plat_colors)	
			<pre>handle_names, handle_lists = zip(+handles.items()) fig.legend(handle_lists, handle_names, loc='upper left', bbox_to_anchor=(1.0,1.0),ncol=1, handlelength=2.0) fig.legend(handles=plat_handles, loc='lower left', bbox_to_anchor=(1.0,0.1), ncol=3, handlelength=1.0) plt.tight_layout(pad=0.4,w_pad=0.5, h_pad=1.0) plt.savefig(f"synthetic_cascade.png", dpi = 300 ,bbox_inches="tight")</pre>	
	$ \begin{array}{c} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\$			
0	🛯 1 @ Python	3 Idle	Saving completed	Mode: Command 💿 Ln 1, Col 1 Efficiency Cascades BabelStream 2020.ipynb

Data from https://doi.org/10.1109/P3HPC51967.2020.00006

BabelStream Cascade plot





CPU7

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



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

Summary

- Three Ps: Performance Portability and Productivity.
- Open standard programming models enable the three Ps.
- Measure the three Ps and use PP-CC plane to guide changes.
- Try out the methodology for yourself:
 - Scripts and tools: <u>https://github.com/UoB-HPC/performance-portability</u>
 - Code Base Investigator: <u>https://github.com/intel/code-base-investigator</u>

Thanks: Jason Sewall and John Pennycook at Intel





Deadline for submissions: Monday 14 June 2021 AoE

https://openmpcon.org

Opportunities for online participation