

Joshua H. Davis, Pranav Sivaraman, Isaac Minn, Abhinav Bhatele  
Department of Computer Science, University of Maryland

## Abstract

Maintaining a single codebase that can achieve good performance on a range of accelerator-based supercomputing platforms is of extremely high value for productive scientific application development. However, the large quantity of programming models available which claim to provide performance portability leaves developers with a complex choice when picking a model to use, potentially requiring an intensive effort to test each available model with kernels from their app. In order to better understand the current state of performance portable programming models, this project evaluates seven of the most popular programming models using two memory-bound mini-applications on two leadership-class supercomputers, Summit and Perlmutter. These results provide a useful evaluation of how well each programming model provides true performance portability in real-world usage for memory-bound applications.

## Background

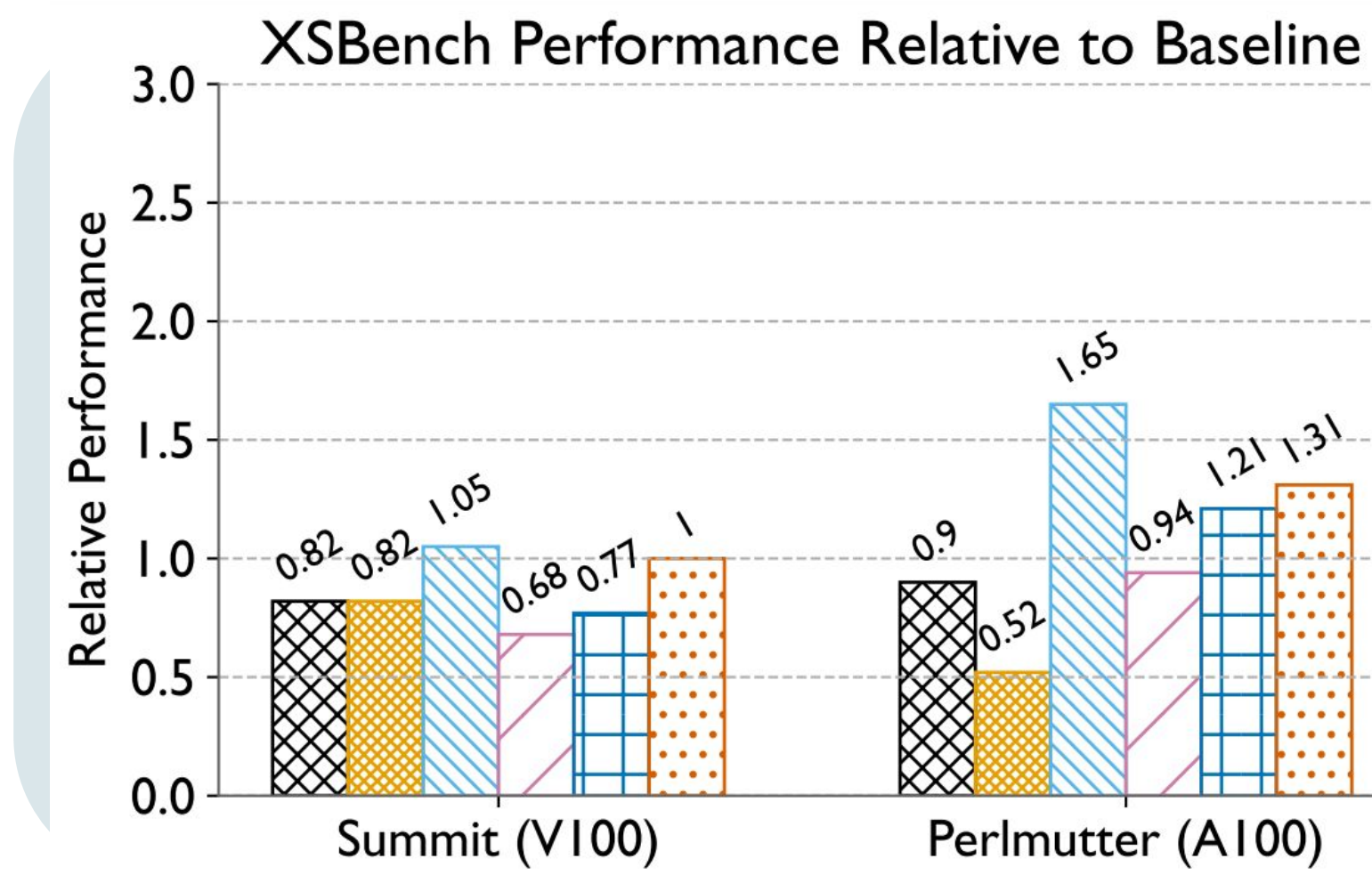
- **Performance Portability:** the ability for one single-source application to run on a range of hardware platforms with good performance
- OpenMP target offload (OMPT), OpenACC (ACC), Kokkos, RAJA, SYCL and HIP are programming models providing portable abstractions



## Methodology

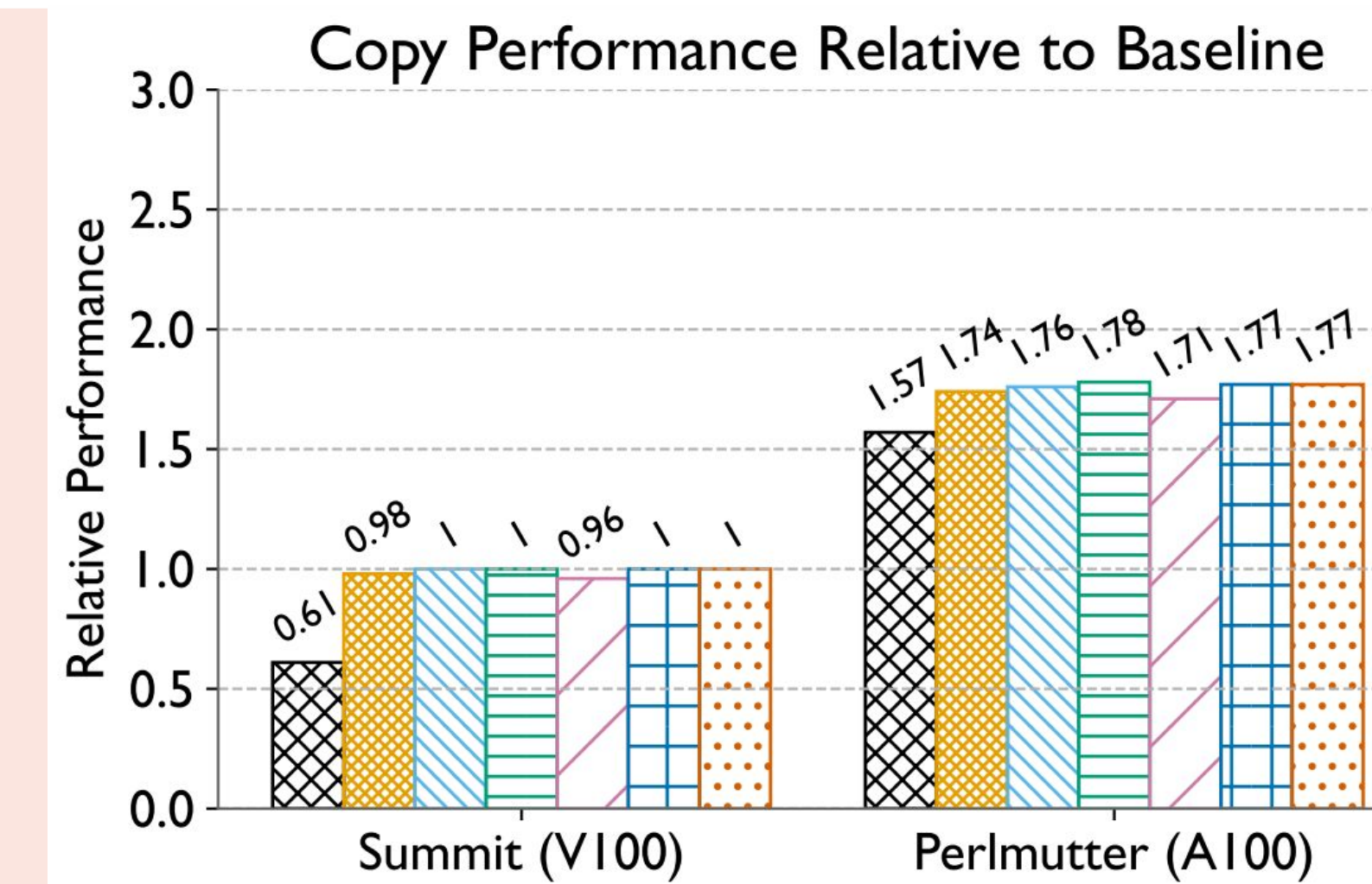
- We surveyed available proxy applications and benchmarks, and selected those with the most available implementations. This poster focuses on two memory-bound codes.
- **XSbench [1]:** memory-bound proxy app from OpenMC (Monte Carlo), evaluated with the “large” problem size (355 isotopes, 11303 grid points)
  - We implemented a new Kokkos port of XSbench for this effort
- **BabelStream [2]:** a memory bandwidth benchmark. We evaluate the dot, triad, and copy kernels for 800 iterations each:
- Evaluation platforms:
  - OLCF Summit: IBM Power 9 CPU and NVIDIA V100 GPU
  - NERSC Perlmutter: AMD EPYC CPU and NVIDIA A100 GPU

## Results

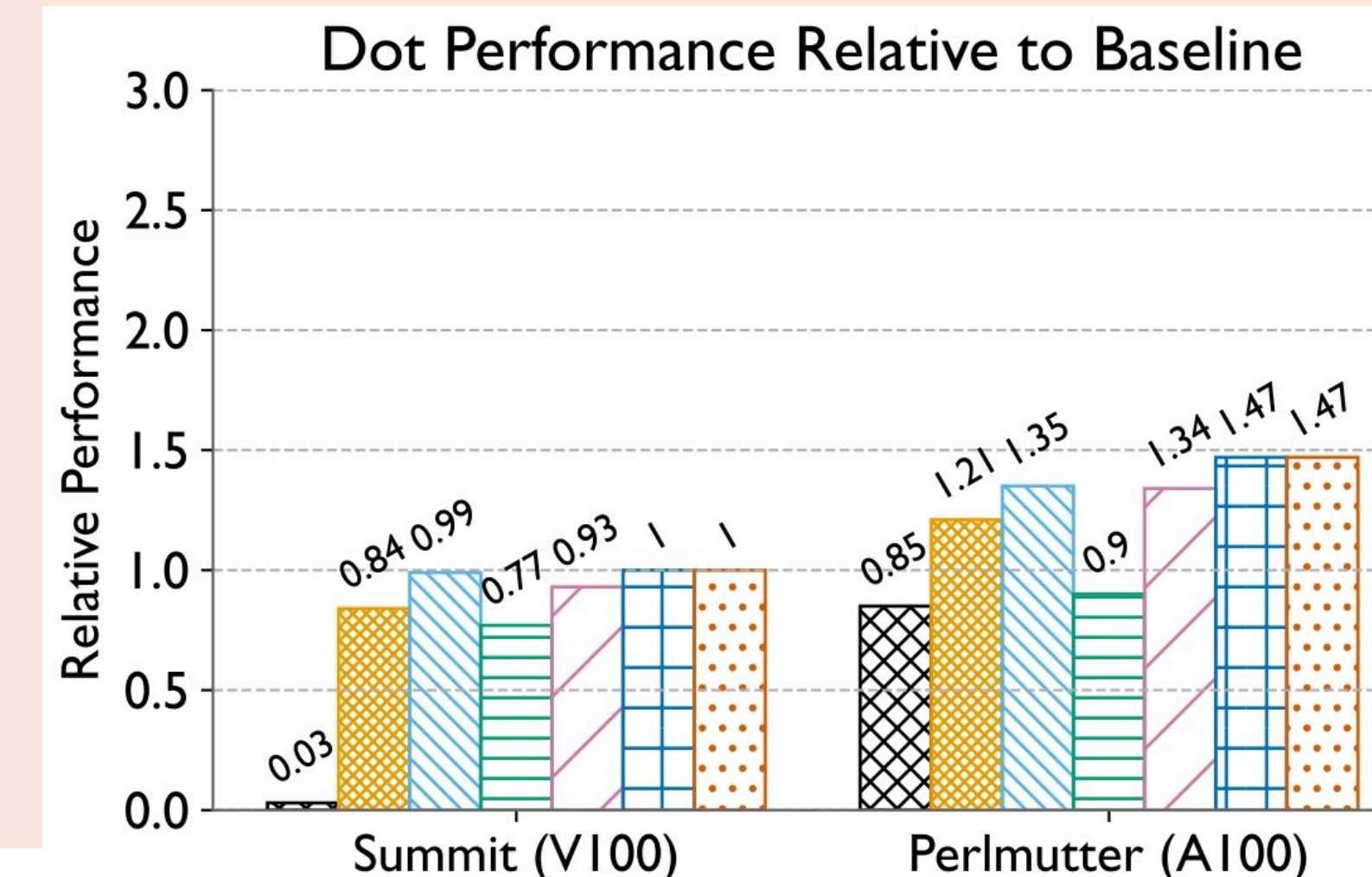
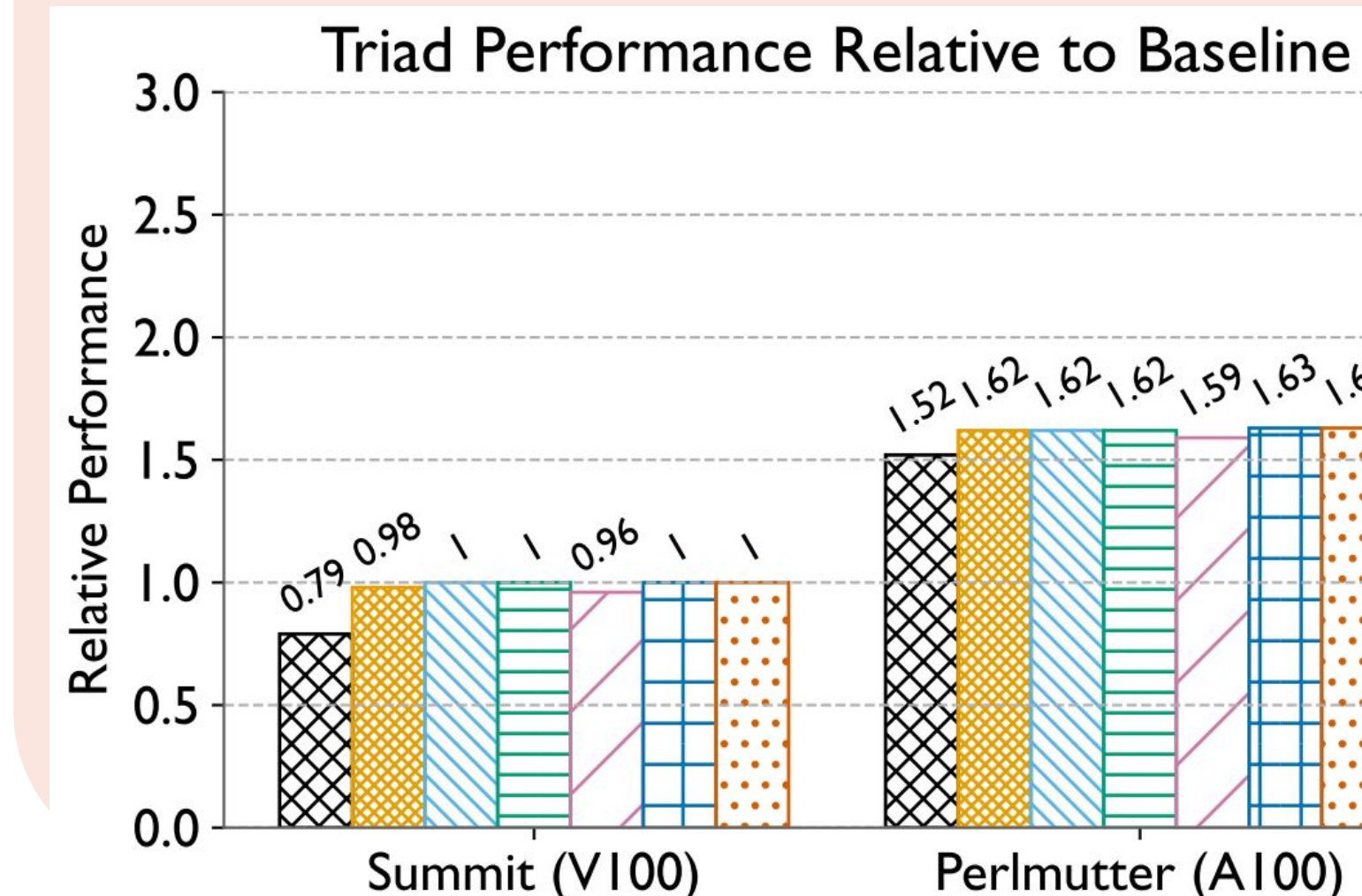


- Performance is measured in terms of runtime in XSbench
- Presented relative to the **Summit CUDA baseline**
- Kokkos outperforms even CUDA on both systems
- OpenACC lags far behind on Perlmutter
- SYCL and HIP perform poorly on Summit
- OpenMP falls in the middle on both systems
- Higher variability across models on Perlmutter

- BabelStream performance is measured in terms of memory transfer bandwidth, relative to Summit CUDA baseline
- Performance is relatively consistent, except for OpenMP and for the Dot kernel
- RAJA and OpenACC also fall behind on the BabelStream Dot kernel
- With Copy kernel in CUDA we observe a 77% performance benefit from moving to A100 GPU
- Higher variability on Perlmutter for Dot



Higher is better



## Acknowledgements

This material is based upon work supported in part by the National Science Foundation Graduate Research Fellowship under Grant No. 1650114. This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. This research used resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility located at Lawrence Berkeley National Laboratory, operated under Contract No. DE-AC02-05CH11231 using NERSC award m2404 for 2023.

## Analysis & Discussion

Performance Portability Metric (Pennycook et al., [3])

Application	OMPT	ACC	Kokkos	RAJA	SYCL	HIP	CUDA
XSbench	0.64	0.45	1.00	N/A	0.61	0.73	0.87
BS-Copy	0.72	0.98	0.99	1.00	0.96	0.99	0.99
BS-Triad	0.85	0.99	1.00	1.00	0.97	1.00	1.00
BS-Dot	0.06	0.83	0.95	0.68	0.92	1.00	1.00

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

- Performance portability metric from Pennycook et al. is defined as the harmonic mean of performance efficiency
- We define performance efficiency as the implementation performance divided by peak performance achieved on the platform across all implementations
- Kokkos, CUDA, and HIP achieve the best performance portability, OpenMP and OpenACC the worst
- The much larger and more complex kernel in XSbench and the reduction operation in BabelStream-dot lead to worse performance portability for most models

## Conclusion and Future Work

- Results set expectations for developers looking to select a programming model for a memory-bound application, and for those porting their application from Summit V100s to Perlmutter A100s
- Summit and Perlmutter both use NVIDIA GPUs – moving to Frontier (AMD) and Aurora (Intel) will provide even greater challenge
- Continuing to analyze the performance of additional applications and programming models using this methodology, including adding more missing programming model implementations

## References

[1] John R. Tramm, Andrew R. Siegel, Tanzima Islam, and Martin Schulz. 2014. XSbench—the development and verification of a performance abstraction for Monte Carlo reactor analysis. *PHYSOR*, (2014).  
 [2] Tom Deakin, James Price, Marc Martineau, and Simon McIntosh-Smith. 2018. Evaluating Attainable Memory Bandwidth of Parallel Programming Models via BabelStream. *Int. J. Comput. Sci. Eng.* 17, 3 (Jan 2018), 247–262.  
 [3] Simon J. Pennycook, Jason D. Sewall, and Victor W. Lee. 2016. A metric for performance portability. In *Proceedings of the 7th International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems*.