Comparative Overview of HPC Frameworks for CPU/GPU Programming

Ivan Smirnov, Vladislav Veselov, Maryana Smirnova, Markus Rampp





Introduction

Results & Discussion

Modern High-Performance Computing (HPC) relies on a wide range of hardware, including CPUs, GPUs, and accelerators. While the Message Passing Interface (MPI) remains the standard for distributed-memory computing, developers must choose from an increasing number of node-level programming models, such as OpenMP, CUDA, HIP, OpenACC, oneAPI (SYCL), Kokkos , and RAJA . These frameworks vary in their ability to deliver portability, performance, and ease of use, making it essential to carefully evaluate their features.	Frame- work	Primary Model	Target Hardware	Functional portability*	Performance Portability	Ecosystem Maturity	Use Cases
This poster compares these programming models based on key factors such as <i>functional portability</i> (the ability to run code across different architectures), <i>performance portability</i> (maintaining efficiency across platforms), <i>ecosystem maturity</i> (tooling, libraries, and community support), and <i>use cases</i> . Drawing from published studies, benchmarks, and real-world applications, we highlight each framework's strengths, limitations, and trade-offs. This analysis aims to provide developers and researchers with practical insidets to		Cross-platform, kernel- based, host-device(with OpenCL C use)	CPU, GPU, FPGA, DSP	(🔀) Was created as cross- vendor [1]	Varies significantly across platforms [3], and even in convenient single-node cases is 1.3 times slower than CUDA [2]	Never gained much traction in the HPC-GPU space, mostly due to the lukewarm support by NVIDIA [22]	Cross-vendor HPC, embedded systems, AI/ML and scientific computing
guide the selection of the best framework for optimizing HPC workloads on increasingly diverse hardware systems. Method This evaluation systematically analyzes eight prominent parallel programming frameworks in HPC: OpenMP, OpenACC, CUDA, RAJA, Kokkos, HIP, SYCL, and OpenCL. The analysis is based on established	SYCL	Cross-platform, single- source C++	CPU, GPU	implementations are available from an increasing number of vendors, including adding support for diverse acceleration API back-ends in addition to OpenCL: Intel oneAPI, AdaptiveCpp, triSYCL, neoSYCL, SimSYCL [4]	It is high on NVIDIA and Intel GPU, but limited on CPU [5]	(?) Growing tooling and libraries through Intel oneAPI; still developing maturity compared to CUDA [6]	HPC, scientific computing, AI/ML and data-parallel tasks
 criteria to ensure consistency and scientific rigor. Evaluation Criteria Primary Model — Frameworks are classified by their parallelization approach (e.g., shared memory, host-device, or abstraction-based). Target Hardware — Compatibility with CPUs, GPUs, FPGAs, and hybrid systems is assessed. Functional Portability — The ability to execute code across vendor platforms with minimal modification is evaluated. 	RAJA	Abstraction layer, loop- level parallelism (multi- backend)	CPU, GPU	Vendor interactions to support new hardware from IBM, NVIDIA, AMD, Intel, and Cray [7]	It is high on NVIDIA GPU, but limited on AMD GPU [8]	(?) Well-supported within DOE but slightly less comprehensive than Kokkos [9]	Scientific simulations, multi-backend HPC and loop management, also performance-portable HPC applications at LLNL
 Performance Portability — The capability to maintain efficient performance across diverse hardware with varying levels of tuning is analyzed. Ecosystem Maturity — Tool availability, community activity, and quality of documentation are considered. Use Cases — Frameworks are examined for their applicability to specific HPC domains such as scientific simulations and AI/ML. 	Kokkos	Abstraction layer, parallel execution and memory management (multi- backend)	CPU, GPU (NVIDIA, AMD, Intel)	(] Provides backend switching between OpenMP, CUDA, and HIP for portability across vendors [12]	Achieves close-to-native performance with tuning [12] [13]	Strong DOE backing, integrated with major HPC libraries like Trilinos [14]	HPC simulations, computational science, fine-grained parallelism and performance-portable C++ applications [13]
 Evaluation Approach Functional Portability: Measures how easily code runs across platforms. High (green): Supports 3+ vendors with minimal code changes. Medium (yellow): Supports 2 vendors, moderate adaptations required. Low (red): Vendor-specific, significant rewrites needed. 	Open ACC	Directive-based, host- device (focused on GPU offloading)	NVIDIA and AMD GPUs	(?) Supports multi-vendor systems but favors NVIDIA GPUs due to more mature implementations	(🔀) Performance depends heavily on compiler quality and vendor support [10], [11]	(?) Limited tools and libraries, mostly focused on legacy projects [11]	Climate modeling, GPU- accelerated legacy applications [10]
 Performance Portability: Assesses how consistently frameworks achieve high performance. High (green): Strong performance across CPUs and GPUs with little tuning. Medium (yellow): Good performance on one platform, acceptable on others with moderate tuning. Low (red): Optimized for one platform only, requiring extensive reimplementation. 	OpenMP	Directive-based, shared memory (with GPU offloading support)	CPU, GPU	Vendor-neutral [17], [18]	Tuning required for GPUs [19], [8]	Robust tools, broad adoption, and active vendor/community support [17]	Shared-memory HPC, engineering simulations, hybrid Al/ML [20]
 <i>boostion Constant of the second states and the second states and the second state of the second state and the second state an</i>	CUDA	Hardware-specific, kernel- based, host-device	NVIDIA GPUs	Only NVIDIA hardware [15]	Performance portability across vendors is non- existent, but high within NVIDIA GPUs [15], [16]	Extensive libraries (cuBLAS, cuDNN), industry-standard tools (Nsight), strong NVIDIA support [16]	GPU-accelerated AI/ML, scientific simulations, rendering[16]
 Conclusion This work reveals the following key observations regarding GPU-focused programming models and complementary CPU-based paradigms: Hardware-Specific Approaches (for example, CUDA for NVIDIA, HIP for AMD, and oneAPI for Intel) typically achieve excellent performance on their target architectures but may increase 	HIP	Hardware-specific, kernel-based, host-device (CUDA-like)	AMD GPUs	Portable for AMD and convertible CUDA applications with HIPIFY [21]	Optimized for AMD, tuning required for other vendors [8], [21]	AMD-focused tools and libraries, still maturing	AMD-targeted HPC, AI/ML and engineering simulations

- maintenance complexity when porting to alternative hardware.
- Directive-Based Methods (for example, OpenMP and OpenACC) offer convenient multi-vendor support, but performance optimization for each backend may lag behind native solutions.
- Abstraction Layers (for example, Kokkos and RAJA) provide singlesource development for multiple platforms, helping manage code complexity. Nonetheless, consistent performance across different architectures depends on the maturity of underlying compilers and runtimes.

Selecting an optimal framework involves balancing immediate performance needs against longer-term sustainability. Although CUDA remains dominant in many NVIDIA-based environments, advanced solutions from AMD, Intel, and high-level abstractions like Kokkos continue to expand the possibilities for portable HPC development. Future studies could further explore the role of emerging frameworks in large-scale applications and evaluate their performance across a wider range of accelerators. This poster supports HPC researchers and developers in navigating this complex landscape by highlighting key criteria for achieving both functional and performance portability.

Table. Confidence indicators:

- Question mark (?) — uncertainty due to limited data or conflicting studies

- Hourglass (👗) — information older than 10 years, potential obsolescence

*more on portability [22]:

	CUDA		HIP		SYCL		OpenACC		OpenMP		Kokkos		RAJA		OpenCL	
	C++	Fortran	C++	Fortran	C++	Fortran	C++	Fortran	C++	Fortran	C++	Fortran	C++	Fortran	C++	Fortrai
NVIDIA	•1	2	— ³	★/4	5	16	•7	8	9	10	1 3	$+^{14}$		1		1
AMD	— ¹⁸	* 19	●20	★/4	A ²¹	16	A ²²	2 3	2 4	2 5	A ²⁸	$+^{14}$		1	•	1
Intel	31	/32	▲ ³³	/34	● ³⁵	_ 6	★36	*37	● ³⁸	9 39	4 2	★14		1	٠	1

Full vendor support

Indirect, but comprehensive support, by vendor Vendor support, but not (yet) entirely comprehensive

Comprehensive support, but not by vendor Limited, probably indirect support - but at least some

No direct support available C++ C++ (sometimes also C) Fortran Fortran

Added by us

References

- 1. Stone, J. E., Gohara, D., & Shi, G. (2010). OpenCL: A parallel programming standard for heterogeneous computing systems. Computational Science and Engineering, 12(3), 66–72.
- 2. Pennycook, S. J., Hammond, S. D., Wright, S. A., Herdman, J. A., Miller, I., & Jarvis, S. A. (2013). An investigation of the performance portability of OpenCL. Journal of Parallel and Distributed Computing, 73(11), 1439–1450. https://doi.org/10.1016/j.jpdc.2012.07.005
- 3. Bertoni, C., Kwack, J., Applencourt, T., Ghadar, Y., Homerding, B., Knight, C., Videau, B., Zheng, H., Morozov, V., & Parker, S. (2020). Performance portability evaluation of OpenCL benchmarks across Intel and NVIDIA platforms. Argonne National Laboratory. Retrieved from https://www.anl.gov
- 4. The Khronos Group. (n.d.). SYCL. Retrieved from https://www.khronos.org/sycl/
- 5. Reguly, I. Z. (2023). Evaluating the performance portability of SYCL across CPUs and GPUs on bandwidth-bound applications. Workshops of the International Conference on High Performance Computing, Network, Storage, and Analysis (SC-W 2023), Denver, CO, USA, November 12–17, 2023.
- 6. HPCwire. (2023, February 28). State of SYCL: ECP BoF showcases progress and performance. Retrieved from https://www.hpcwire.com/2023/02/28/state-of-syclecp-bof-showcases-progress-and-performance/
- 7. Lawrence Livermore National Laboratory. (n.d.). RAJA Portability Suite: Enabling performance portable CPU and GPU HPC applications. Retrieved from <u>https://computing.llnl.gov/projects/raja-managing-application-</u> portability-next-generation-platforms
- 8. Davis, J. H., Sivaraman, P., Minn, I., Parasyris, K., Menon, H., Georgakoudis, G., & Bhatele, A. (2023). An evaluative comparison of performance portability across GPU programming models. Department of Computer Science, University of Maryland, & Lawrence Livermore National Laboratory.
- 9. RAJA Documentation. (n.d.). Retrieved from <u>https://raja.readthedocs.io/en/develop/</u>
- 10.Sabne, A., Sakdhnagool, P., Lee, S., & Vetter, J. S. (2014). Evaluating performance portability of OpenACC. In Languages and Compilers for Parallel Computing (pp. 63–77). Springer.
- 11.Deakin, T., et al. (2019). Performance portability across diverse computer architectures. 2019 IEEE/ACM International Workshop on Performance, Portability and *Productivity in HPC (P3HPC), Denver, CO, USA, 2019, pp. 1–13. https://doi.org/10.1109/P3HPC49587.2019.00006*
- 12. Edwards, H. C., Trott, C. R., & Sunderland, D. (2014). Kokkos: Enabling manycore performance portability through polymorphic memory access patterns. *Journal of* Parallel and Distributed Computing, 74(12), 3202–3216.
- 13. Kokkos Documentation. (n.d.). Available at <u>https://kokkos.org/kokkos-core-wiki/</u>
- 14. Kokkos Abstract. (2023). Available at https://kokkos.org/about/abstract/

*

- 15. NVIDIA Developer. (n.d.). CUDA Toolkit Documentation. Retrieved from https://docs.nvidia.com/cuda/
- 16. NVIDIA Developer. (n.d.). CUDA Zone. Retrieved from <u>https://developer.nvidia.com/cuda-zone</u>
- 17. OpenMP. (n.d.). OpenMP (Open Multi-Processing). Retrieved from https://en.wikipedia.org/wiki/OpenMP
- 18. OpenMP Architecture Review Board. (n.d.). *OpenMP Compilers & Tools*. Retrieved from <u>https://www.openmp.org/resources/openmp-compilers-tools/</u>
- 19. Malik, D. (2022). Performance Portability of OpenMP. Technical University of Munich. Retrieved from https://events.gwdg.de/event/243/contributions/503/attachments/139/174/OpenMP.Pd
- 20. van Waveren, M. (2020). *OpenMP Use Cases*. OpenMP ARB & CS GROUP. Retrieved from <u>https://www.openmp.org/wp-content/uploads/OpenMP-Use-Cases</u>vanWaveren.pdf
- 21. AMD. (n.d.). HIP Documentation: Performance Portability for Heterogeneous Systems. Retrieved from https://rocm.docs.amd.com/projects/HIP/en/latest/index.html
- 22. Herten, A. (2023). Many cores, many models: GPU programming model vs. vendor compatibility overview. Proceedings of the P3HPC Workshop, hosted at SC23 (International Conference for High Performance Computing, Networking, Storage, and Analysis). Retrieved from <u>https://doi.org/10.48550/arXiv.2309.05445</u>