
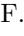














# Characterizing GPU Energy Usage in Exascale-Ready Portable Science Applications

William F. Godoy<sup>1</sup>  , Oscar Hernandez<sup>1</sup> , Paul R. C. Kent<sup>1</sup> ,  
Maria Patrou<sup>1</sup> , Kazi Asifuzzaman<sup>1</sup> , Narasinga Rao Miniskar<sup>1</sup> ,  
Pedro Valero-Lara<sup>1</sup> , Jeffrey S. Vetter<sup>1</sup> , Matthew D. Sinclair<sup>2</sup> ,  
Jason Lowe-Power<sup>3</sup> , and Bobby R. Bruce<sup>3</sup> 

<sup>1</sup> Oak Ridge National Laboratory, Oak Ridge, TN, USA  
{godoywf,oscar,kentpr,patroum,asifuzzamank,miniskarnr,valerolarap,  
vetter}@ornl.gov

<sup>2</sup> University of Wisconsin-Madison, Madison, WI, USA  
sinclair@cs.wisc.edu

<sup>3</sup> University of California, Davis, Davis, CA, USA  
{jlowepower,bbruce}@ucdavis.edu

**Abstract.** We characterize the GPU energy usage of two widely adopted exascale-ready applications representing two classes of particle and mesh solvers: (i) QMCPACK, a quantum Monte Carlo package, and (ii) AMReX-Castro, an adaptive mesh astrophysical code. We analyze power, temperature, utilization, and energy traces from double-/single (mixed)-precision benchmarks on NVIDIA's A100 and H100 and AMD's MI250X GPUs using queries in NVML and rocm\_smi.lib, respectively. We explore application-specific metrics to provide insights on energy vs. performance trade-offs. Our results suggest that mixed-precision energy savings range between 6–25% on QMCPACK and 45% on AMReX-Castro. Also, we found gaps in the AMD tooling used on Frontier GPUs that need to be understood, while query resolutions on NVML have little variability between 1 ms-1 s. Overall, application level knowledge is crucial to define energy-cost/science-benefit opportunities for the codesign of future supercomputer architectures in the post-Moore era.

**Keywords:** Energy efficiency · HPC Applications · GPU Power

## 1 Introduction

As energy consumption and costs have grown exponentially in the post-Moore exascale era, high-performance computing (HPC) faces new challenges [23]. Interest in the energy-cost/science-benefit trade-offs is again gaining traction<sup>1</sup> as HPC systems become more heterogeneous [32]. Since HPC traditionally focused on optimizing time-to-solution, it is crucial to understand applications characteristics to design future energy-efficient hardware. Here, we provide insights

<sup>1</sup> <https://www.ora.gov/2024EECWorkshop>.

on the GPU energy characteristics of two applications developed under the US Department of Energy’s (DOE’s) Exascale Computing Project (ECP, 2016–2023) [16] that are widely used at HPC production facilities across the world: (i) QMCPACK [14, 15, 19], a quantum Monte Carlo (QMC) package, and (ii) the mesh-based AMReX-Castro’s [1] compressible astrophysics code. We capture power, utilization, temperature, and calculate energy traces on NVIDIA’s A100 and H100 GPUs and AMD’s MI250X GPU. To capture these measurements, we designed an open-source `HWEnergyTracer.jl`<sup>2</sup> tool that runs side-by-side with an application and captures queries from NVIDIA’s Management Library (NVML), and AMD’s ROCm System Management Interface Library (`rocm_smi_lib`).

The paper is organized as follows: Sect. 2 provides information for QMCPACK, AMReX-Castro and selected benchmarks. Section 3 describes our methodology and the targeted GPU systems. Section 4 presents our results and analysis of the applications’ energy characteristics. Related work in HPC is shown in Sect. 5. Finally, Sect. 6 provides our conclusions and future directions. To the best of our knowledge, our contributions on quantifying science-per-energy has not previously been an integral part of the applications’ development process.

## 2 Background

*QMCPACK and the NiO benchmark:* QMCPACK is an open-source, many-body, ab-initio QMC framework solving the Schrödinger equation for atoms, molecules, 2D nanomaterials, and solids. QMC methods lead to far greater accuracy, but at a much greater computational cost. Key recent QMCPACK improvements made during the ECP included (i) a redesigned diffusion Monte Carlo (DMC) solver [19] – the focus on this study – using OpenMP offload capabilities on GPUs, and (ii) software engineering improvements for CPU/GPUs [9]. We use the nickel oxide (NiO) supercell benchmark<sup>3</sup> which is characterized by the total number of electrons, for which its required memory grows quadratically, while its computation increases at a cubic rate.

*AMReX-Castro and the Sedov case:* AMReX [36] is a widely used adaptive mesh refinement framework that powers several HPC applications running at scale. AMReX decomposes a problem into levels of adaptive resolution and rectangular patches. During the ECP, AMReX’s solver capabilities were advanced for new CPU/GPU architectures, while energy-efficiency is in the product roadmap [22]. We use the Castro astrophysical radiation-magneto-hydrodynamics code that builds on AMReX to run the 2D Sedov spherical blast wave standard problem on a rectangular AMR mesh. Its computational costs are typically driven by the time evolution of the total number of AMR cells [8].

<sup>2</sup> <https://github.com/JuliaORNL/HWEnergyTracer.jl>.

<sup>3</sup> [https://www.olcf.ornl.gov/wp-content/uploads/OLCF-6\\_QMCPACK\\_description-1.pdf](https://www.olcf.ornl.gov/wp-content/uploads/OLCF-6_QMCPACK_description-1.pdf).

### 3 Methodology

We use the Julia language [3] to call queries listed in Table 1 in NVIDIA’s NVML by using the CUDA.jl package [2], and AMD’s `rocm_smi_lib` C API using Julia’s foreign function interface. The result is the `HWEnergyTracer.jl` tool that runs 15s before the start and 10s after the end of the application run to guarantee a steady static power state. Power values are integrated over time to obtain energy traces, and to minimize observed false positives (over-/under-shoots). We use three different GPU systems, 2 NVIDIA and 1 AMD, as described in Table 2. No power capping was applied in this work.

**Table 1.** Queries used for NVIDIA’s NVML and AMD’s `rocm_smi_lib`

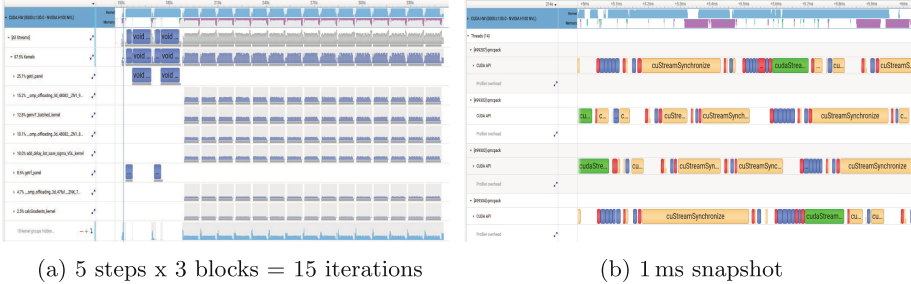
Metric	Relevant Query	Description
<b>NVIDIA</b>	<b><code>nvmlDeviceGet*</code></b>	
Power (W)	<code>PowerUsage</code>	Power usage of the GPU and its associated circuitry (e.g., memory) averaged over a 1 s interval [34]
Utilization (%)	<code>UtilizationRates</code>	Percent of time over the past sample period, between 1 and $\frac{1}{6}$ s, during which kernels were executing
Temperature (°C)	<code>Temperature</code>	Current temperature readings for the device
<b>AMD</b>	<b><code>rsmi_dev_*</code></b>	
Power (W)	<code>power_ave_get</code>	device energy counter average for a short time (1 ms)
Utilization (%)	<code>busy_percent_get</code>	Percentage of time busy processing
Temperature (°C)	<code>temp_metric_get</code>	Retrieved from the temperature sensor for the device

**Table 2.** System hardware and software used in this study

System	Milan0	Hudson	Frontier
<b>Hardware</b>			
GPU-per-node	2 NVIDIA A100	2 NVIDIA H100	8 GCD AMD MI250X
Memory(GB)/Bandwidth(GB/s)	HBM2E 80/1,940	HBM3 94/1,940	HBM2E 64/3,276
Thermal Design Power (W)	300	400	500
<b>Software</b>			
GPU Tool Chain	NVHPC 24.9	NVHPC 24.9	ROCm 6.2
QMCPACK	v3.17.1	v3.17.1	v3.17.1
Compiler	Clang 19.1	Clang 19.1	AMDClang 6.2
Programming Model	OpenMP-offload	OpenMP-offload	OpenMP-offload
AMReX-Castro	v24.12	v24.12	v24.12
Compiler	GCC 13.2	GCC 13.2	GCC 12.3
Programming Model	CUDA 12.4	CUDA 12.4	HIP 6.2

## 4 Results

Results are presented for the application energy characteristics traces. We discuss: (i) regions of interests based on power consumption, (ii) measurement resolutions from 1 ms to 1 s to capture difference of the averages values given by the vendor tools, (iii) impact of double/single precision, and (iv) exploring application-specific energy-efficiency metrics (e.g. science/Joule).



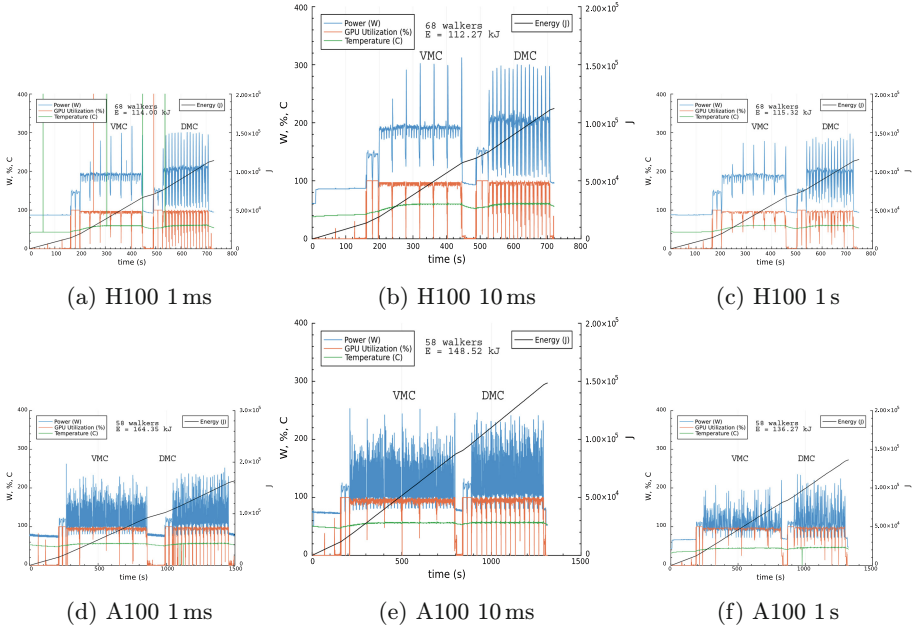
**Fig. 1.** QMCPACK NiO Benchmark DMC GPU traces on an NVIDIA H100.

### 4.1 QMCPACK NiO Benchmark Results

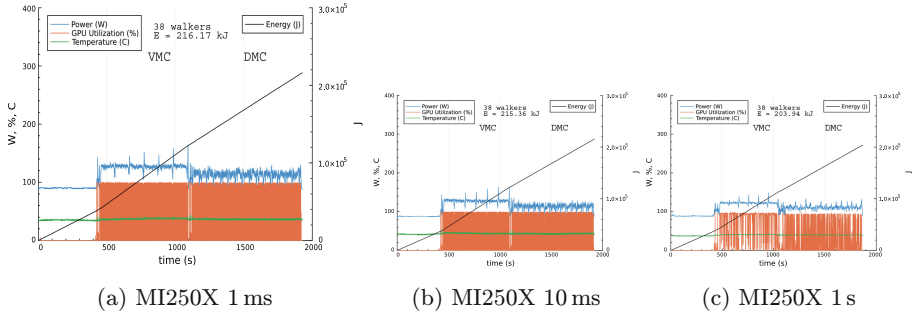
We run the NiO benchmark for a representative supercell of 512 atoms and 6,144 electrons. Figure 1 shows the DMC code profile showing the asynchronous nature of kernel launching for each MC “walker” to maximize GPU usage (blue top row) [19]. The number of walkers is maximized to fit on GPU memory [6].

Figures 2 and 3 show the power, temperature, utilization and energy traces for the maximum number of walkers on the NVIDIA H100 and A100, and AMD MI250X GPUs, respectively, for different query time resolutions. In all cases, we use double-precision configurations. The query time resolution used in subsequent runs is highlighted as the larger figure. Also, the four run stages are captured by the traces for spline data initialization (showing low GPU activity), Variational Monte Carlo (VMC), matrix inversions, and DMC.

On the H100 and A100, the DMC power and utilization patterns match the kernel behavior shown in Fig. 1 a. As for the MI250X, the DMC stage power shows similar characteristics for minimum values, but peaks are not signaled as those in NVIDIA GPUs. This might be explained by the different query methodologies used by NVIDIA NVML, which has a wider time average frame of 1 s. The MI250X utilization shows fast variations between 100 and 0% highlighting the need for more robust energy tooling for further investigation. Energy consumption varies nearly linearly for all cases allowing for simple modeling of this benchmark. On the H100, the overall integrated energy usage ( $\approx 112\text{--}115\text{ kJ}$ ) shows little sensitivity,  $\approx 2.6\%$  up to 1 s, to the different resolutions. Importantly, the 1 ms resolution yields false positive results, as shown in the spikes in Fig. 2 a,

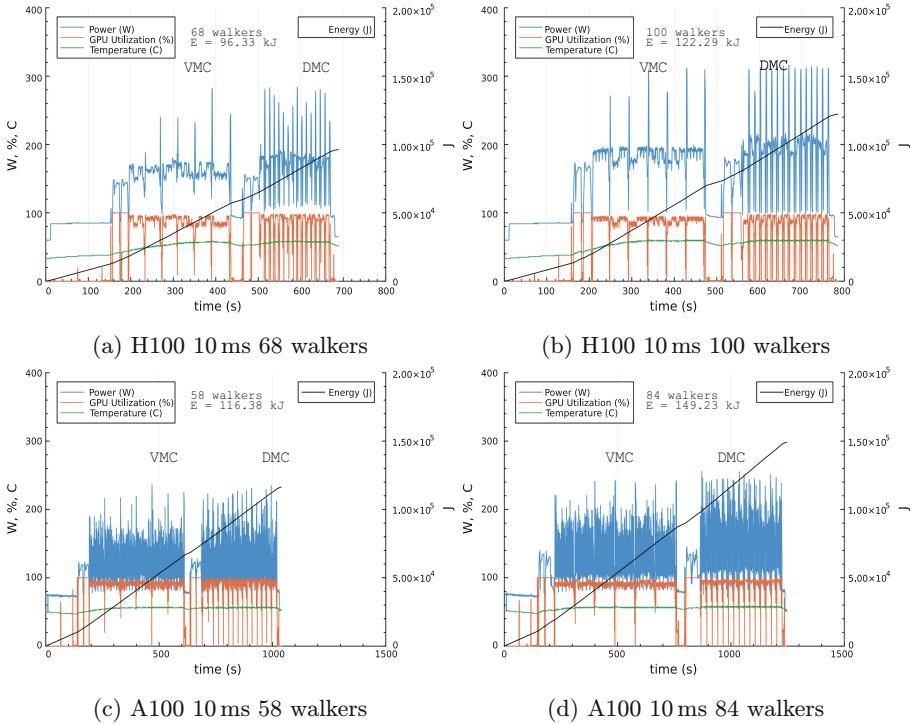


**Fig. 2.** Energy characteristics of the QMCPACK NiO benchmark on NVIDIA H100 and A100 for different query time resolutions.



**Fig. 3.** Energy characteristics of the QMCPACK NiO benchmark on an AMD MI250X for different query time resolutions.

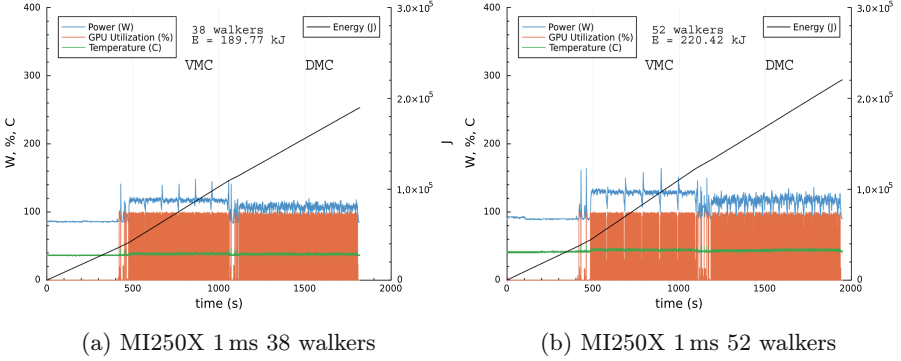
which are a repeatable pattern in our experiments. Additionally, the 1 s resolution is not enough to capture variations reinforcing the need for experimenting with finer resolutions. On the A100, the power traces have more variability than for the H100. It is unclear if this is due to changes in the methodology for measuring power in more recent NVIDIA GPUs or if it is an actual hardware characteristic. Importantly, querying at 1 ms introduces overhead in the time-to-solution on the A100, we ensured this is a repeatable pattern. Overall, both the A100 and H100 runs have similar energy characteristics, although as expected



**Fig. 4.** Mixed-precision traces on NVIDIA H100 and A100 for (a) max double-precision walkers (68 and 58) and (b) max mixed-precision walkers (100 and 84).

the H100 is more performance and energy efficient. As for AMD, energy characteristics are slightly different. As for the cause, either the time-to-solution (and therefore the energy consumption) is higher than it is on the A100 (and the H100). Two things stand out: (i)  $\leq 10$  ms measurements do not introduce noticeable overhead, but 1 s resolutions might not be sufficient to capture variability for this workload. (ii) The DMC average power is lower for the AMD MI250X than for the NVIDIA GPUs at the expense of longer time-to-solution for a smaller number of walkers. In all cases, temperature values remain nearly constant, thereby indicating that thermal solutions for cooling the devices run efficiently for these benchmarks.

*Mixed Precision Runs:* QMCPACK uses a mixed-precision (single/double) option that reduces memory consumption at the expense of less accurate results. This is applied in two ways: (i) a shorter time-to-solution for the same number of walkers, or (ii) more walkers to fill the GPU’s memory and increase utilization. These two options are illustrated in Figs. 4 and 5 for the NVIDIA H100 and A100, and AMD MI250X GPUs, respectively.



**Fig. 5.** AMD MI250X mixed-precision traces for (a) max double-precision walkers (38) and (b) max mixed-precision walkers (52).

For these cases, (a) and (c) contain the same number of walkers as the maximum double-precision runs, and (b) contains the maximum number of walkers for a mixed-precision run. In all cases, mixed-precision runs result in energy savings, thanks to faster times-to-solution than the double-precision runs. This is more evident on the runs with NVIDIA’s H100 and A100 GPUs than with AMD’s MI250X GPU. Because some components keep double-precision representations, the savings are not expected to be close to 50%. QMCPACK’s mixed-precision has not been thoroughly studied for the latest GPUs. For example, this work identified different default behaviors (e.g., DMC mixed precision was running extra calculations) that has been corrected in QMCPACK. Thus showing the importance of monitoring energy characteristics during development. A simple power and utilization trace can showcase hardware usage correlations with expected algorithm behavior.

*Energy Metric.* The DMC performance metric (throughput) is measured as the computational cost of advancing walkers at each step and block [6]. We adapt this metric for energy efficiency (science/energy) by changing the denominator from time-to-solution to energy consumption in:

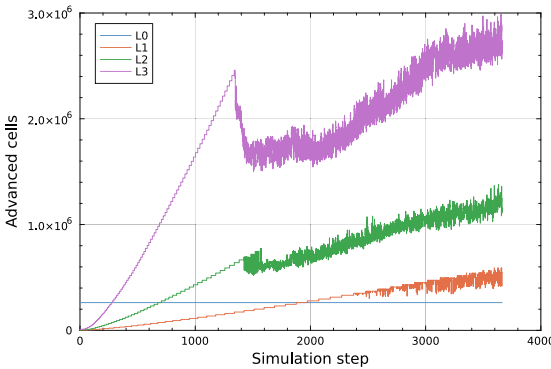
$$\text{Throughput}_{\text{Energy}} = \frac{\text{walkers} \times \text{blocks} \times \text{steps}}{\text{DMC energy}}, \quad (1)$$

where  $\text{DMC energy}$  is the energy needed by the DMC region.

Metric results are presented in Table 3 along with the average GPU power and utilization. Using the maximum number of walkers for mixed precision on H100 leads to the greatest science/energy ratio, and running in double precision draws higher energy rates. Energy savings from mixed-precision for the same number of walkers are on the order of 6%–25%. We also added cases to compare H100 versus A100, and A100 versus MI250X runs for the same number of walkers. Energy improvements between A100 and H100 cases are roughly  $1.5\times$ , whereas

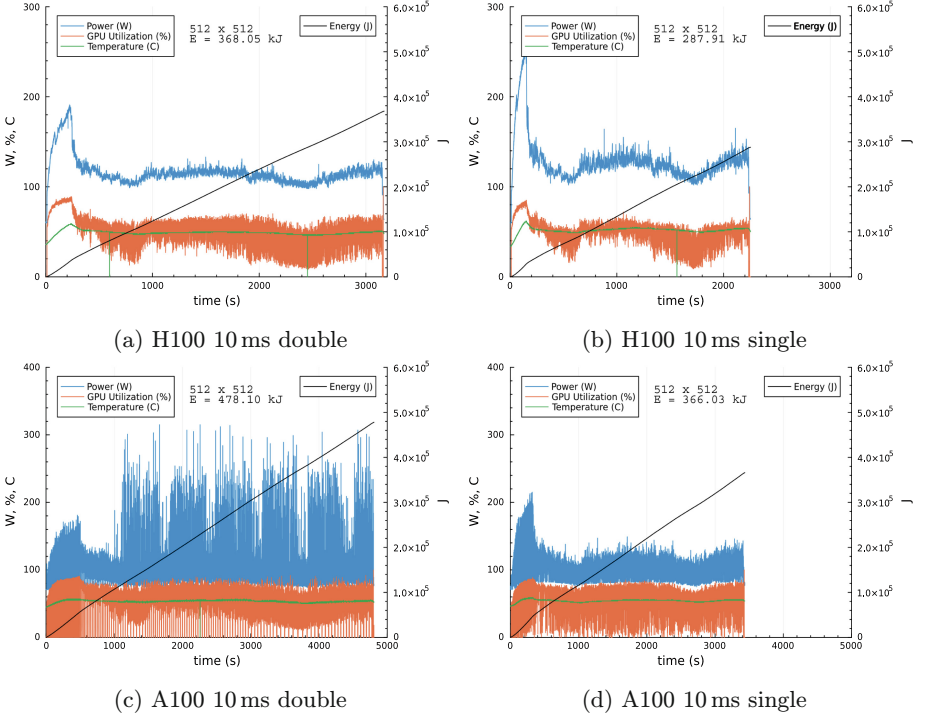
**Table 3.** Energy metrics for QMCPACK’s DMC on the NiO a512 benchmark for multiple GPU configurations

Configuration	Walkers *Max	Throughput Energy (1/kJ)	Power (W)	GPU (%)
<b>NVIDIA</b>				
H100-mixed	*100	38.69	190.02	72.54
H100-mixed	*68	33.20	172.60	74.60
H100-double	*68	27.26	191.87	83.92
H100-mixed	84	37.56	182.31	74.69
H100-mixed	58	32.42	174.24	75.93
H100-double	58	26.28	181.89	78.58
A100-mixed	*84	25.25	136.10	85.76
A100-mixed	*58	20.86	121.59	85.08
A100-double	*58	17.31	124.97	88.03
A100-mixed	52	21.41	122.22	83.84
A100-mixed	38	16.13	109.01	86.92
A100-double	38	15.58	119.11	86.69
<b>AMD</b>				
MI250X-mixed	*52	9.12	115.97	39.33
MI250X-mixed	*38	7.57	106.36	40.12
MI250X-double	*38	6.19	112.61	39.27



**Fig. 6.** AMReX-Castro Sedov problem advanced cells per AMR mesh level as a function of simulation step.

A100 runs are 2× more efficient than MI250X runs. Results for AMD’s MI250X suggest that improvements are needed on either the software or vendor tool stack to achieve results comparable to NVIDIA’s A100.

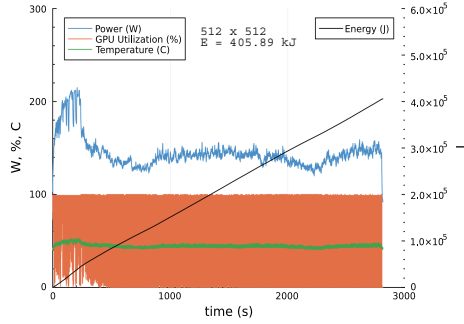


**Fig. 7.** AMReX-Castro Sedov energy characteristics on NVIDIA H100 and A100 for a  $512 \times 512$  base mesh with  $\text{CFL} = 0.25$  using (a) double and (b) single precision.

## 4.2 AMReX-Castro Sedov

We ran the 2D Sedov case to capture an AMR simulation’s influence on energy characteristics. Compute activity comes from advancing a variable evolution (governed by PDEs) on cells placed at different mesh levels. Figure 6 shows the evolution of four mesh levels (L0–L3) in a coarser-to-finer order as a function of simulation steps. As seen, the finer mesh (L3) dominates the computation, and all of them except the base mesh (L0) evolve similarly. We expect that energy consumption would be dominated by this evolution and the use of double- or single-precision representations.

*Double/single-precision* In all cases, we use the previously set query resolution for NVIDIA (10 ms) and AMD (1 ms) GPUs, and fixed the Courant-Friedrichs-Lewy (CFL) condition to a value of 0.25. Figure 7 shows the traces on NVIDIA H100 and A100 runs when using the largest possible base  $512 \times 512$  mesh and double- and single-precision representations. As expected, power and utilization show a strong correlation with the finer mesh evolution in Fig. 6. Power variability is high for the rapid evolution of the mesh levels (1,200–1,400 simulation steps), reaching an absolute peak and stabilizing until the end of the



**Fig. 8.** AMReX-Castro Sedov energy characteristics on an AMD MI250X for a  $512 \times 512$  base mesh using double precision, 1 ms resolution, and  $\text{CFL} = 0.25$ .

run. The use of single precision (Figs. 7 band 7 d) represents a faster code execution ( $\approx 2200/3200 = 68\%$ ) when compared with double precision (Figs. 7 a and 7 c) but a larger fraction in terms of energy ( $\approx 287/368 = 78\%$ ). In fact, the initial peak power is higher than for double precision. Nevertheless, the H100 runs demonstrate the improvements in energy consumption over the previous A100 by using only 60% and 78% as much energy for the double- and single-precision cases, respectively. For AMD, the MI250X is more energy and time efficient than the A100 for the Sedov case when using double precision (Fig. 8). Power characteristics are very similar to those observed on NVIDIA’s H100.

We did not observe differences between double- and single-precision runs on the MI250X, so the latter is not provided. Energy-efficiency metrics still need to be studied for this application as the influence of more complex performance and energy drivers (e.g., non-linear evolution in the AMR mesh sizes and CFL number) is highly problem- and hardware-dependent.

## 5 Related Work

Energy efficiency in HPC has garnered significant attention in recent years at all levels (including applications, facilities, and tools). Muriedas, et al. and Yang, et al. measured HPC application’s energy consumption on Intel and NVIDIA GPUs [11, 35]. Schieffer et al. characterized energy on AMD matrix cores using `rocm-smi-lib` [25]. Other work has designed a variety of tools for predicting GPU power consumption. This includes AccelWattch [12], machine learning models to predict power consumption [17, 26, 33], and flexible simulator interfaces [31]. However, most of these approaches focus on modeling and simulation tools, unlike our work. Simsek et al. [29, 30] studied the energy efficiency of the SPH-EXA astrophysics application on CPUs and GPUs as well as the impact of dynamic frequency scaling using the open-source Power Measurement Toolkit [5]. Govind et al. [10] investigated the power characteristics of scientific and machine learning applications on the Perlmutter supercomputer. Zhao et al. studied power traces of the popular Vienna Ab initio Simulation Package on

NVIDIA A100 GPUs including power capping techniques [37]. Foster et al. [7] studied the energy efficiency of machine learning benchmarks. Mantovani et al. studied the performance and energy consumption of HPC workloads on Arm ThunderX2 CPU cores [20]. Mittal and Vetter presented a survey of methods for improving GPU energy efficiency [21]. Bridges et al. provided an overview of techniques for obtaining GPU power consumption data [4]. At the facility level, Karimi et al. presented a system-wide HPC monitoring framework on the Summit supercomputer [13]. Shin et al. introduced a comprehensive strategy for the sustainability of future HPC systems [27]. Silva et al. presented a comprehensive survey on the energy of state-of-the-art supercomputers [28]. Rrapaj et al. quantified the long-term energy consumption in systems of the National Energy Research Scientific Computing Center [24].

## 6 Discussion, Conclusions and Future Work

We quantified and analyzed the energy consumption characteristics of two science applications that are widely used in HPC systems—QMCPACK and AMReX-Castro—on recent NVIDIA A100 and H100, and AMD MI250X GPUs. Kernel characteristics were mapped to the power, temperature, utilization, and energy traces for different configurations, including reduced floating-point precision. Application-specific metrics were discussed compare the science-per-energy. Some key observations:

- **Observation 1:** small variability for the power, utilization, temperature traces and integrated energy calculations were observed for vendor APIs, NVML and rocm\_smi\_lib, in our application measurements.
- **Observation 2:** reduced floating-point precision resulted in more energy-efficient runs, while not at an ideal 50% rate these were on the order of 6%–25% on QMCPACK (on NVIDIA and AMD) and 45% for AMReX-Castro (on NVIDIA only).
- **Observation 3:** Energy-efficiency improvements (on the order of  $1.5\times$ ) were shown for NVIDIA’s H100 over their A100. Room for improvement exists for AMD’s GPU tools and applications as the ecosystem matures.
- **Observation 4:** the proposed science-per-energy metric for QMCPACK allows for comparison across GPU vendors and generations. It factors in computational time aspects into the total energy required. A AMReX-Castro science-per-energy metric requires factoring-in the AMR variability in the number of cells advanced at every level.

Next steps include expanding to more HPC-relevant scientific workloads, understanding energy at the different GPU component levels, and integrating modeling for design space exploration (e.g., GEM5 [18]). As HPC facilities energy costs increase, we conclude that this type of analysis at the application level is crucial in the co-design of future supercomputing architectures.

**Acknowledgements.** This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the US Department of Energy (DOE). The US government retains and the publisher, by accepting the article for publication, acknowledges that the US government retains a nonexclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for US government purposes. DOE will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<https://www.energy.gov/doe-public-access-plan>). This material is based on work supported by the DOE's Office of Science, Office of Advanced Scientific Computing Research through EXPRESS: 2023 Exploratory Research for Extreme Scale Science. PRCK was supported by the DOE's Office of Science, Basic Energy Sciences, Materials Sciences and Engineering Division as part of the Computational Materials Sciences Program and the Center for Predictive Simulation of Functional Materials. This research used resources of the Oak Ridge Leadership Computing Facility and the Experimental Computing Laboratory at the Oak Ridge National Laboratory, which is supported by the DOE's Office of Science under Contract No. DE-AC05-00OR22725. WG would like to acknowledge Brandon Tran from the University of Wisconsin for the valuable discussion on NVML.

## References

1. Almgren, A., et al.: Castro: a massively parallel compressible astrophysics simulation code. *JOSS* **5**(54) (2020). <https://doi.org/10.21105/joss.02513>
2. Besard, T., et al.: Effective extensible programming: unleashing julia on GPUs. *IEEE TPDS* **30**(4) (2019). <https://doi.org/10.1109/TPDS.2018.2872064>
3. Bezanson, J., et al.: Julia: a fresh approach to numerical computing. *SIAM Rev.* **59**(1), 65–98 (2017). <https://doi.org/10.1137/141000671>
4. Bridges, R.A., Imam, N., Mintz, T.M.: Understanding GPU power: a survey of profiling, modeling, and simulation methods. *ACM Comput. Surv.* **49**(3) (2016)
5. Corda, S., et al.: Pmt: power measurement toolkit. In: *IEEE/ACM HUST Workshop* (2022). <https://doi.org/10.1109/HUST56722.2022.00011>
6. Elwasif, W., et al.: Application experiences on a GPU-accelerated arm-based HPC testbed. In: *Proceedings of the HPC Asia 2023 Workshops*. p. 35–49. *HPCAsia '23 Workshops*, Association for Computing Machinery, New York (2023). <https://doi.org/10.1145/3581576.3581621>
7. Foster, B., et al.: Evaluating energy efficiency of GPUs using ML benchmarks. In: *IPDPSW* (2023). <https://doi.org/10.1109/IPDPSW59300.2023.00019>
8. Godoy, W.F., et al.: Modeling pre-exascale AMR Parallel I/O workloads via proxy applications. In: *IPDPSW* (2022). <https://doi.org/10.1109/IPDPSW55747.2022.00153>
9. Godoy, W.F., et al.: Software stewardship and advancement of a high-performance computing scientific application: QMCPACK. *FGCS* **163**, 107502 (2025)
10. Govind, A., et al.: Comparing power signatures of HPC workloads: machine learning vs simulation. In: *SC-W* (2023). <https://doi.org/10.1145/3624062.3624274>
11. Gutiérrez Hermosillo Muriedas, J.P., et al.: Perun: benchmarking energy consumption of high-performance computing applications. In: *Euro-Par* (2023). [https://doi.org/10.1007/978-3-031-39698-4\\_2](https://doi.org/10.1007/978-3-031-39698-4_2)
12. Kandiah, V., et al.: AccelWattch: a power modeling framework for modern GPUs. In: *MICRO* (October 2021)

13. Karimi, A.M., et al.: Power profile monitoring and tracking evolution of system-wide HPC workloads. In: ICDCS (2024). <https://doi.org/10.1109/ICDCS60910.2024.00018>
14. Kent, P.R.C., et al.: QMCPACK: advances in the development, efficiency, and application of auxiliary field and real-space variational and diffusion quantum Monte Carlo. *J. Chem. Phys.* **152**(17) (2020). <https://doi.org/10.1063/5.0004860>
15. Kim, J., et al.: QMCPACK: an open source ab initio quantum Monte Carlo package for the electronic structure of atoms, molecules and solids. *J. Phys. Cond Matter* **30**(19), 195901 (2018)
16. Kothe, D., et al.: Exascale computing in the united states. *CiSE* **21**(1) (2019)
17. Lee, W., et al.: PowerTrain: a learning-based calibration of McPAT power models. In: ISLPED, pp. 189–194 (2015)
18. Lowe-Power, J., et al.: The gem5 simulator: V20.0 (2020). <https://arxiv.org/abs/2007.03152>
19. Luo, Y., Doak, P., Kent, P.: A high-performance design for hierarchical parallelism in the QMCPACK Monte Carlo code. In: SC-W HiPar (2022). <https://doi.org/10.1109/HiPar56574.2022.00008>
20. Mantovani, F., et al.: Performance and energy consumption of HPC workloads on a cluster based on Arm ThunderX2 CPU. *FGCS* **112** (2020)
21. Mittal, S., Vetter, J.S.: A survey of methods for analyzing and improving GPU energy efficiency. *ACM Comput. Surv.* **47**(2) (2014)
22. Myers, A., et al.: AMReX and pyAMReX: looking beyond the exascale computing project. *IJHPCA* **38**(6) (2024). <https://doi.org/10.1177/10943420241271017>
23. Reed, D., Gannon, D., Dongarra, J.: Reinventing high performance computing: challenges and opportunities (2022). <https://arxiv.org/abs/2203.02544>
24. Rrapaj, E., et al.: Power consumption trends in supercomputers: a study of NERSC’s Cori and perlmutter machines. In: ISC (2024). <https://doi.org/10.23919/ISC.2024.10528943>
25. Schieffer, G., et al.: On the rise of AMD matrix cores: performance, power efficiency, and programmability. In: ISPASS (2024). <https://doi.org/10.1109/ISPASS61541.2024.00022>
26. Shim, J.S., et al.: DeepPM: transformer-based power and performance prediction for energy-aware software. In: DATE, pp. 1491–1496 (2022)
27. Shin, W., et al.: Towards sustainable post-exascale leadership computing. In: SC-W (2024). <https://doi.org/10.1109/SCW63240.2024.00225>
28. Silva, C., et al.: A review on the decarbonization of high-performance computing centers. *Renew. Sustain. Energy Rev.* **189**, 114019 (2024)
29. Simsek, O.S., et al.: Accurate measurement of application-level energy consumption for energy-aware large-scale simulations. In: SC-W (2023). <https://doi.org/10.1145/3624062.3624272>
30. Simsek, O.S., et al.: Increasing energy efficiency of astrophysics simulations through GPU frequency scaling. In: SC-W (2024). <https://doi.org/10.1109/SCW63240.2024.00229>
31. Smith, A., et al.: Designing generalizable power models for open-source architecture simulators. In: OSCAR (2024)
32. Vetter, J.S., et al.: Productive computational science in the era of extreme heterogeneity. Tech. rep., USDOE Office of Science (SC), Washington, D.C. (United States) (2018). <https://doi.org/10.2172/1473756>
33. Wu, G., et al.: GPGPU performance and power estimation using machine learning. In: HPCA, pp. 564–576 (2015). <https://doi.org/10.1109/HPCA.2015.7056063>

34. Yang, Z., et al.: Accurate and convenient energy measurements for GPUs: a detailed study of NVIDIA GPU's built-in power sensor. In: SC24 (2024). <https://doi.org/10.1109/SC41406.2024.00028>
35. Yang, Z., et al.: Accurate and convenient energy measurements for GPUs: a detailed study of NVIDIA GPU's built-in power sensor. In: SC (2024). <https://doi.org/10.1109/SC41406.2024.00028>
36. Zhang, W., et al.: AMReX: a framework for block-structured adaptive mesh refinement. JOSS 4(37) (2019). <https://doi.org/10.21105/joss.01370>
37. Zhao, Z., et al.: Understanding VASP Power Profiles on NVIDIA A100 GPUs. In: SC-W (2024). <https://doi.org/10.1109/SCW63240.2024.00189>